

بناام خدا

هفتمین کنفرانس آمار ایران

مجموعه مقالات

جلد دوّم (فارسی)

۲ تا ۴ شهریور ماه ۱۳۸۳

دانشگاه علامه طباطبایی

تهران - ایران

مجموعه مقالات هفتمین کنفرانس آمار ایران

آماده‌سازی و صفحه‌آرایی : TEX-پاکی ، یاشار اکبری
طراحی جلد : نیما دانش‌پرور
تیراژ : ۸۰۰ نسخه
چاپ و صحافی : مرکز آمار ایران
تاریخ انتشار : بهمن ماه ۱۳۸۳

بنام خدا

پیشگفتار

خداوند بزرگ را سپاسگزاریم که ما را یاری نمود تا هفتمین کنفرانس آمار ایران را با همکاری انجمن آمار ایران، کوشش همکاران، دانشجویان، کارکنان و حمایت‌های مسئولین محترم دانشگاه علامه طباطبایی و سازمانها و دستگاه‌های مختلف کشور در روزهای ۲ تا ۴ شهریور ۱۳۸۳ در دانشگاه علامه طباطبایی برگزار نمائیم.

تلاش کمیته علمی کنفرانس برای برنامه‌ریزی در جهت ارتقای سطح علمی کنفرانس از ابتدای سال ۱۳۸۲ آغاز گردید و استقبال قابل توجه اساتید، دانشجویان و محققین کشور در ارائه مقاله در زمینه‌های نظری و کاربردی آمار چشمگیر بود، به طوری که ناچار شدیم به دلیل محدودیت زمانی تعدادی از مقالات را که از کیفیت برتری برخوردار بودند انتخاب نماییم.

متقاضیان شرکت در کنفرانس تعداد ۳۹۲ مقاله به دبیرخانه ارسال نمودند و کمیته‌های تخصصی یازده گانه کمیته علمی کنفرانس تعداد ۱۹۷ مقاله را به عنوان سخنرانی، تعداد ۱۰۱ مقاله را برای ارائه به صورت پوستر و ۶ مورد را به صورت کارگاه آموزشی پذیرش نمودند. به علاوه از دانشمندان صاحب‌نام علم آمار داخل و خارج کشور برای ارائه سخنرانی‌های عمومی و تخصصی دعوت به عمل آمد که تعدادی از آنها به دلایلی از شرکت در کنفرانس عذرخواهی نمودند و نهایتاً تعداد ۵۷ مقاله به صورت سخنرانی و ۴ کارگاه آموزشی توسط مدعوین داخلی و خارجی در کنفرانس ارائه گردید. مجموعه حاضر شامل ۳۲ مقاله پذیرفته شده توسط داوران کمیته‌های تخصصی یازده گانه کمیته علمی کنفرانس می‌باشد.

برگزاری هفتمین کنفرانس آمار ایران بدون مساعدت‌های بی‌دریغ همکاران گرامی در دانشکده اقتصاد، تعدادی از اساتید دانشگاه‌های سراسر کشور به عنوان هیأت داوران، دانشجویان دوره کارشناسی و کارشناسی ارشد آمار و حمایت‌های هیأت رئیسه محترم و تلاش‌های کارمندان و کارکنان دانشگاه علامه طباطبایی امکان‌پذیر نبود و جا دارد از تمامی عزیزانی که در برگزاری کنفرانس زحماتی را متحمل شدند قدردانی و تشکر شود.

هماهنگی جلسات کمیته علمی با تلاش آقای مجید رجبی و انجام کلیه امور کامپیوتری و تهیه و تنظیم مجموعه مقالات با همکاری و پشتکار بی‌دریغ و خستگی‌ناپذیر خانم تی تی ماندگاریان و آقای یاشار اکبری صورت پذیرفت که از این عزیزان نهایت تشکر و امتنان را دارد.

بهمن ماه ۱۳۸۳

کمیته‌های علمی و برگزاری هفتمین کنفرانس آمار ایران

حمایت کنندگان

هفتمین کنفرانس آمار ایران توسط دانشگاه علامه طباطبایی و با همکاری انجمن آمار ایران و حمایت‌های علمی و مالی مؤسسات زیر برگزار گردید که بدینوسیله مراتب قدردانی خود را از آنان ابراز می‌داریم.

پژوهشکده آمار

مرکز آمار ایران

سازمان مدیریت و برنامه‌ریزی کشور

بانک مرکزی جمهوری اسلامی ایران

بانک کشاورزی

پژوهشکده هواشناسی و علوم جو

اسامی اعضای کمیته‌های علمی و برگزاری هفتمین کنفرانس آمار ایران

اعضای کمیته برگزاری

۱. دکتر فرزاد اسکندری (مسئول دبیرخانه کمیته علمی)
۲. آقای محمدرضا اصغری اسکویی (مسئول امور رایانه)
۳. آقای عبدالرحیم بادامچی زاده (دبیر پنجمین مسابقه دانشجویی آمار و مسئول نمایشگاه کتاب)
۴. دکتر محمد با منی مقدم (مسئول امور میهمانان خارجی)
۵. آقای محمدجواد زالی (دبیر کمیته اجرایی)
۶. دکتر شهرام سلیمی (مسئول پذیرش)
۷. دکتر سید علی میریان (نماینده دانشگاه در کمیته برگزاری)
۸. دکتر نادر نعمت‌الهی (مسئول کمیته و دبیر کنفرانس)
۹. دکتر حمیدرضا نواب‌پور (دبیر کمیته علمی)
۱۰. دکتر عبدالساده نیسی (مدیر امور اجرایی)

اعضای کمیته علمی

- | | |
|-------------------------------|--|
| دانشگاه علامه طباطبایی | ۱. دکتر فرزاد اسکندری |
| دانشگاه صنعتی اصفهان | ۲. دکتر احمد پارسیان |
| دانشگاه علامه طباطبایی | ۳. دکتر محمد جلوداری ممقانی |
| دانشگاه الزهرا | ۴. دکتر صدیقه شمس |
| پژوهشکده آمار | ۵. دکتر عباس گرامی |
| دانشگاه علوم پزشکی شهید بهشتی | ۶. دکتر یدالله محرابی |
| دانشگاه تربیت مدرس | ۷. دکتر محسن محمدزاده |
| دانشگاه شهید بهشتی | ۸. دکتر محمدرضا مشکانی |
| دانشگاه علامه طباطبایی | ۹. دکتر نادر نعمت‌الهی (دبیر کنفرانس) |
| دانشگاه علامه طباطبایی | ۱۰. دکتر حمیدرضا نواب‌پور (دبیر و مسئول کمیته) |
| دانشگاه شهید بهشتی | ۱۱. دکتر محمدقاسم وحیدی اصل |

فهرست مطالب

عنوان صفحه

مقالات جلد اول

- تحلیل داده‌های گسسته با روشهای متفاوت در رده آمار بیزی و انتخاب بهترین روش ۱
علی‌آذر بر، فرزاد اسکندری
- کاربرد آزمون‌های آماری جهت کشف نفوذهای خرابکارانه در شبکه‌های کامپیوتری ۲۸
افشین آشفته، سید مهدی امیرجهانشاهی
- بیز و ناریبی تحت تابع زیان لاینکس ۳۸
منصور آقابائنی جزی، احمد پارسیان
- برآورد توابع چگالی با روش ماکسیمم آنترابی زمانی که چهارگشتاور اولیه توزیع معلوم است ۶۳
علی‌آقا محمدی
- تحلیل چند متغییره شاخص چاقی کودکان زیر دو سال و رابطه ساختاری آن با چاقی والدین آنها در شیراز ۷۶
سید محمدتقی آیت‌اللهی، سید تقی حیدری
- برآورد پارامترهای توزیع توانی بر اساس داده‌های رکوردی با در نظر گرفتن زمان رخداد رکورد ۹۱
جعفر احمدی، محمد آرشی
- اطلاع فیشور در مشاهدات فرین توزیعهای متقارن ۱۰۳
جعفر احمدی، بهاره خطیب آستانه، مصطفی رزمخواه
- مقایسه اطلاع فیشور نهفته در رکوردهای ضعیف و قوی با مشاهدات مستقل و هم‌توزیع ۱۱۶
جعفر احمدی، مصطفی رزمخواه، بهاره خطیب آستانه
- تحلیل چند بعدی فقر توسط مجموعه‌های فازی ۱۳۰
حمید اردهه، صدیقه علیمردانی

ب هفتمین کنفرانس آمار ایران

عنوان صفحه

۱۵۰ تحلیل مدل‌های عاملی پویای بیزی در روش‌های چند متغیره غیر خطی
فرزاد اسکندری، پیمان آقابیگی

۱۷۶ تحلیل توزیع درآمد به کمک توزیع پارتو (رگرسیون پارتو)
فرزاد اسکندری، رقیه قبادی

۱۸۷ ابعاد و ملاکهای کیفیت و روش‌های ارزیابی کیفیت آمارهای رسمی
محمد رضا اناری

۱۹۹ الگوریتم بوت استرپ در آمار فضایی
نصراله ایران‌پناه، محسن محمدزاده

۲۱۳ آن-تروپی-رکدها
سیمیندخت براتپور، جعفر احمدی

مدل تحلیل عاملی بررسی ارتباط بین فن‌آوری اطلاعات و سبک رهبری در
سازمانها
۲۲۴ رضا برادران کاظم‌زاده، مهدی شریفی زمانی

۲۳۷ معرفی معادلات دیفرانسیل تصادفی و روشهای عددی حل آنها
قاسم برید لقمانی

۲۵۱ پیش‌بینی قابلیت اعتماد سیستم‌های پیچیده بکمک آمارهای ترتیبی
حسین بیورانی، ویکتور یوریویچ کارلیوف

۲۵۷ استفاده از آن‌تروپی در اندازه‌گیری نابرابری توزیع درآمد
عین‌الله پاشا، احمد زنده‌دل

۲۷۳ روش بوت استرپ برای برآوردگرهای فازی
عین‌الله پاشا، اشکان شباک

۲۸۷ پیش‌بینی دینامیکی-آماري دمای ماکزیمم و می‌نیمم شهر تهران
مژده پدرام، نجمه ابوفاضلی

مجموعه مقالات ج

عنوان صفحه

- ۳۰۰ هموارسازی هسته‌ای تابع دوره نگار با استفاده از فاصله کولبک - لیبلر
غلامعلی پرهام، سحر درنیانی
- ۳۱۰ الگوهای نقطه‌ای
رضا پورطاهری، محمدقاسم وحیدی اصل
- ۳۲۳ روش کمترین قدر مطلق خطا در مدل‌های خطی
حمزه ترابی
- ۳۴۱ ارزیابی توزیع درآمد در ایران
محسن جلالی
- ۳۵۷ شعاع طیفی قدم زدن تصادفی ساده بر گروه‌های کاکستر مثلثی
محمد جلوداری ممقانی
- ۳۷۲ روش مناسب برآورد متوسط هزینه سالانه خانوار
مریم جوادی، فاطمه هرندی، زهره فالاح محسن‌خانی
- ۳۸۰ سازمانهای غیر دولتی (NGOs)، خالاهای آماری، راهکارها
جواد حسین‌زاده
- ۳۹۳ پیش‌بینی فضایی به روش کریگینگ گوسی تبدیل یافته
فاطمه حسینی، محسن محمدزاده
- ۴۰۰ برآوردگرهای بیز و مینیماکس بر پایه‌ی تابع زیان فازی
ایرج حقیقی‌نژاد، سید محمود طاهری
- ۴۱۴ شرح قضیه آگاتیکه و کاربرد آن در رگرسیون
سمیه حیدری، ناصر رضا ارقامی
- تعمین متوسط اندازه نمونه لازم برای انجام آزمون نسبت احتمال دنباله‌ای
پارامتر مقیاس در یک زیر خانواده نمایی از توزیع‌ها
ناصر داورزنی، محمد جعفری جوزانی
- ۴۴۲ مدل‌بندی بارش با استفاده از نظریه فرآیندهای نقطه‌ای
باقر ذهبیون، رضا پورطاهری

د هفتمین کنفرانس آمار ایران

عنوان صفحه

مقالات جلد دوم

- برآورد ضریب جینی و بررسی چگونگی توزیع درآمد نقاط شهری و روستایی
استان های کشور ۴۵۵
جعفر رحمانی شمسی
- شرایط آغازین برای مدل های تصادفی - انتقالی در داده های طولی با پاسخ
دودویی کامل و ناقص ۴۶۲
زهرا رضایی قهرودی، مجتبی گنجعلی
- تحلیل ممیزی با استفاده از آمیخته های نرمال ۴۷۲
قاسم رکابدار، رحیم چینی پرداز
- خطای مطلق یا نسبی در تعیین اندازه نمونه ۴۸۴
علیرضا زاهدیان
- مطالعه آنتروپی باقیمانده توزیع های طول عمر ۴۹۷
یونس زهرهوند، مجید اسدی
- یک الگوریتم مقدماتی برای شبیه سازی عددی معادلات دیفرانسیل تصادفی . ۵۰۷
پرویز سرگلزایی، محمد امینی، محمود دادخواه
- رگرسیون ژرفا در حالت چند گانه ۵۲۶
حمید شریف
- استفاده از معیارهای کولبک - لیبیلر و فاصله چرنوف برای آنالیز ممیزی سری های
زمانی چند متغیره ۵۴۰
سارا شفیع بابایی، رحیم چینی پرداز
- مدل تاثیر تصادفی با داده هایی با پاسخ های آمیخته ۵۵۸
علی صادقی، مجتبی گنجعلی
- آزمون تصادفی شده برای میانگین توزیع نرمال بر پایه ی داده های نادقیق . . . ۵۶۷
سید محمود طاهری، تکتم بزرگوار

- تعمین جایگاه اقتصادی استان‌های کشور - ۱۳۷۹ از نظر ارزش افزوده و درصد نیروی شاغل در بخش‌های اقتصادی به روش تحلیل عاملی ۵۸۱
وحید طیفوری
- برآورد فاصله‌ای برای خانواده توزیع‌های نمایی طبیعی ۶۰۸
محسن عارفی، غلامرضا محتشمی برزادران
- روش درست‌نمایی وزنی در برآوردیابی برای ناحیه کوچک ۶۲۲
ملیحه عباس نژاد مشهدی
- رگرسیون چندگانه ۶۳۱
فرهاد فتاحی، عباس گرامی
- میانگین‌گیری بیزی مدل‌های رگرسیونی ۶۴۶
افشین فلاح، محسن محمدزاده
- نمونه‌های گردان، جایگزینی برای سرشماری‌های جمعیتی ۶۵۶
زهره فلاح محسن‌خانی، فاطمه هرندی، فرشید جمشیدی
- فواصل پیش‌بینی برای مدل‌های اتورگرسیو خطی با استفاده از تکنیک‌های بوت استرپ ۶۶۶
ملیحه قابل، صادق رضایی
- تعمین تعداد خوشه‌ها در تحلیل‌های آمیخته با استفاده از آنتروپی نرمال شده ۶۷۶
محمد قربانی، محسن محمدزاده
- شناخت اشیاء (شکل و ساختار) با استفاده از آمار ۶۸۵
هادی گنجی، محمود صفارزاده
- مقایسه روش‌های مختلف نمونه‌گیری با استفاده از الگوریتم‌های مونت کارلو و بوت استرپ ۶۹۶
عباس محمدخانی، نصراله ایران‌پناه
- پیشگویی فضایی بیزی با استفاده از روش‌های مونت کارلو ۷۰۸
محسن محمدزاده، مجید جعفری خالدي

عنوان صفحه

- ۷۱۵ پیش‌بینی بیزی طول عمر برای مدل پارتو با حجم نمونه تصادفی
محسن محمدزاده، منصور زرگر
- ۷۳۱ شبیه‌سازی احتمالات ورشکستگی در فرآیندهای مخاطره بیمه
نادر مظاهری
- ۷۴۵ تعیین دبی‌ها اوج با استفاده از داده‌های بالاتر از یک آستانه معین
سید سعید موسوی ندوشنی
- ۷۶۰ روش‌های نوین اقتصادسنجی در تحلیل داده‌های مکانی
رزیتا مویدفر، آذر ابراهیمی
- ۷۸۰ هم‌انباشتگی کسری در سری‌های زمانی و بررسی نرخ تورم در ایران
سید محمود میرجلایی، قاسم تارمست
- معیارهای انتخاب متغیر در داده‌های چندمتغیره بر اساس ساختار کوواریانس
و همبستگی
نادر نعمت‌الهی، قدرت‌الله رحمتی
- ۷۹۸ پیشگویی در نمونه‌گیری از جامعه متناهی تحت مدل‌های خطای اندازه‌گیری
و کاربرد آن در برآورد هزینه و درآمد خانوار
نادر نعمت‌الهی، پوریا رضاسلطانی
- استفاده از برآوردگرهای عادی و رگرسیونی نمونه‌گیری مجموعه رتبه‌دار در
برآورد مقدار کل محصول گندم ایران
نادر نعمت‌الهی، لطیف سعادت‌ی
- ۸۱۹ انتشار داده‌های آماری و کنترل افشای اطلاعات فردی
حمیدرضا نواب‌پور، محمد بردبار عشرت‌آبادی
- ۸۵۹ برآورد واریانس به روش جک‌نایف در آمارگیریه‌های چندچارچوبی
مرجان نورینی
- ۸۷۷ آزمون برازش توزیع لجستیک در تعیین توانایی و رتبه‌بندی در امتحان‌ها
سید مقتدی هاشمی پرست، سید حسن مرتضوی نصیری، کمال عقیق

برآورد ضریب جینی و بررسی چگونگی توزیع درآمد نقاط شهری و روستایی استان‌های کشور

جعفر رحمانی شمسی

سازمان مدیریت و برنامه‌ریزی استان یزد

چکیده: در این مقاله روشی برای برآورد مقدار ضریب جینی^۱ به عنوان مهمترین شاخص بررسی توزیع درآمد^۲ ارائه شده است. این روش بر خلاف روش‌های معمول محاسبه ضریب جینی متکی بر داده‌های خام نبوده و در آن می‌توان از داده‌های گروه بندی شده مرکز آمار ایران استفاده کرد. کاربرد اصلی این روش در محاسبه ضریب جینی استان‌ها می‌باشد زیرا دسترسی به داده‌های خام استان‌ها بسادگی امکان پذیر نمی‌باشد. همچنین بر اساس این روش برای اولین بار ضرایب جینی کلیه استان‌های کشور به تفکیک نقاط شهری و روستایی محاسبه گردیده تا چگونگی توزیع درآمد در استان‌ها مورد بررسی قرار گیرد همچنین با استفاده از روش آنالیز خوشه‌ای استان‌های کشور از نظر چگونگی توزیع درآمد مقایسه و طبقه‌بندی شده‌اند.

واژه‌های کلیدی: توزیع درآمد، ضریب جینی، آنالیز خوشه‌ای

۱ مقدمه

در سالهای اخیر مسئله توزیع درآمد و توزیع امکانات زندگی بین اقشار مختلف جامعه مورد توجه بیشتری قرار گرفته است. مطالعات مربوط به توزیع درآمد سابقه‌ای کمتر از یک قرن دارد و در این سالها شاخصهای گوناگونی برای محاسبه میزان نابرابری توزیع درآمد مطرح شده است که یکی از مهمترین آنها ضریب جینی است.

۲ ضریب جینی

رایج‌ترین شاخص نابرابری توزیع درآمد که در اکثر مطالعات و بررسی‌های مربوط به توزیع درآمد از آن استفاده می‌شود ضریب جینی است. این معیار در سال ۱۹۱۲ به وسیله آماردان ایتالیایی کورودن سی جینی^۳ ارائه شده است و به روش زیر محاسبه می‌شود:

$$G = \frac{\Delta}{2\mu}$$

1) Gini Coefficient 2) Income Distribution 3) Gini

و

$$\Delta = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1}^n |x_i - x_j| \quad (1)$$

که در آن:

x_i : بیان کننده درآمد^۲ امین فرد جامعه

n : تعداد کل افراد جامعه (نمونه)

μ : میانگین درآمد افراد جامعه (نمونه)

بر اساس رابطه فوق اگر Δ در حداقل مقدار ممکن (یعنی صفر یا حالتی که توزیع درآمد کاملاً عادلانه است) باشد، اندازه ضریب جینی برابر صفر و چنانچه Δ در حداکثر مقدار ممکن یعنی 2μ باشد (در حالت توزیع کاملاً نا عادلانه درآمد) اندازه ضریب جینی برابر یک خواهد بود.

به طور کلی اندازه ضریب جینی بین صفر و یک تغییر می‌کند که مقدار صفر نشان دهنده برابری کامل توزیع درآمد و مقدار یک نشان دهنده نابرابری کامل توزیع درآمد است، مقادیر بین

صفر و یک معمولاً بدین صورت تفسیر می‌گردد:

۱- $G \leq 0.3$ توزیع نسبتاً عادلانه درآمد

۲- $0.3 < G < 0.5$ توزیع نیمه عادلانه درآمد

۳- $G \geq 0.5$ توزیع نا عادلانه درآمد

۳ فرمول برآوردی ضریب جینی

با توجه به فرمول محاسباتی ضریب جینی در خواهیم یافت که محاسبه دقیق ضریب جینی نیاز به درآمد (هزینه) تک تک افراد جامعه (نمونه) دارد. این مقدار برای نقاط شهری و روستایی کشور در اکثر سالها توسط مرکز آمار ایران و یا برخی از محققان محاسبه شده است، اطلاعات مورد نیاز برای محاسبه این معیار از داده‌های خام طرح هزینه و درآمد خانوارهای شهری و روستایی مرکز آمار ایران بدست آمده است. همچنین این معیار در برخی از استان‌های کشور در بعضی از سالها بر اساس داده‌های خام این طرح محاسبه شده است ولی تاکنون تحلیل جامعی از چگونگی توزیع درآمد کلیه استان‌های کشور انجام نگرفته است و دلیل اصلی آن عدم دسترسی ساده به این داده‌ها به تفکیک استان‌های کشور و مشکلاتی است که محاسبه این ضریب از روی داده‌های خام دارد. با این وجود اطلاعات گروه بندی شده هزینه‌های خانوارهای شهری و روستایی به تفکیک استان‌های کشور هر ساله توسط مرکز آمار ایران منتشر می‌شود. در نشریات این مرکز خانوارهای شهری و روستایی در ۱۰ گروه هزینه‌ای گروه بندی شده‌اند و متوسط هزینه هر گروه و نیز فراوانی خانوارهای درون این گروهها در اختیار قرار گرفته است. در واقع بجای هزینه هر

خانوار درون گروههای هزینه‌ای تنها متوسط هزینه آن خانوارها برای هر استان در اختیار است. بر اساس ایده داده‌های گمشده (missing values) می‌توان هزینه هر خانوار درون گروهها را به عنوان داده گمشده در نظر گرفت و متوسط هزینه آن گروه را جایگزین آنها کرد. در واقع جداولی که در اختیار است را می‌توان به عنوان جداول فراوانی داده‌ها در نظر گرفت که داده‌های خام توسط آنها در ۱۰ گروه طبقه‌بندی شده‌اند و متوسط هر طبقه (نماینده طبقه) و نیز فراوانی هر طبقه در اختیار است. در این حالت فرمول برآوردی ضریب جینی را می‌توان به صورت زیر نوشت:

$$\Delta = \frac{1}{n(n-1)} \sum_{i=1}^{10} \sum_{j=1}^{10} f_i f_j | \bar{x}_i - \bar{x}_j |$$

که در آن:

\bar{x}_i : متوسط هزینه خانوارهای گروه هزینه‌ای i ام
 n : تعداد کل افراد جامعه (نمونه)
 f_i : تعداد خانوارهای درون گروه هزینه‌ای i ام

۴ برآورد ضریب جینی نقاط شهری و روستایی استان‌های کشور بین سالهای ۱۳۷۴ تا ۱۳۸۰

در این بخش بر اساس نتایج طرح هزینه و درآمد خانوارهای شهری و روستایی مرکز آمار ایران و با استفاده از فرمول برآوردی ارائه شده، برآورد ضریب جینی نقاط شهری و روستایی تمام استان‌های کشور در سالهای ۱۳۷۴ تا ۱۳۸۰ بدست آمده است، همچنین با استفاده از میانگین هندسی متوسط ضریب جینی هر استان بین سالهای ۱۳۷۴ تا ۱۳۸۰ نیز محاسبه شده است. برای سادگی محاسبات شبه نرم‌افزاری در محیط Excel طراحی شده و تنها با وارد کردن اطلاعات مورد نیاز یعنی متوسط هزینه هر گروه‌های هزینه‌ای و تعداد خانوارهای درون گروههای هزینه‌ای می‌توان ضریب جینی را مشاهده کرد. نتایج بدست آمده به تفکیک نقاط شهری و روستایی در جدول ۱ و ۲ آمده است.

لازم به ذکر است استان قم در سال ۱۳۷۶ و استان‌های قزوین و گلستان در سال ۱۳۷۷ به جمع استان‌های کشور اضافه شده‌اند و به همین دلیل نتایج سالهای قبل در مورد آنها وجود ندارد. جدول ۲ ضرایب جینی نقاط روستایی استان‌های کشور را نمایش می‌دهد.

جدول ۱: ضریب جینی نقاط شهری استان های کشور

| متوسط سالهای ۷۴ تا ۸۰ | سال ۸۰ | سال ۷۹ | سال ۷۸ | سال ۷۷ | سال ۷۶ | سال ۷۵ | سال ۷۴ | استان |
|-----------------------|--------|--------|--------|--------|--------|--------|--------|---------------------|
| ۰/۳۹۶ | ۰/۴۰ | ۰/۳۹ | ۰/۳۹ | ۰/۳۹ | ۰/۳۹ | ۰/۴۰ | ۰/۴۰ | کل کشور |
| ۰/۳۸۷ | ۰/۴۲ | ۰/۴۲ | ۰/۳۹ | ۰/۳۸ | ۰/۳۷ | ۰/۳۷ | ۰/۳۸ | آذربایجان شرقی |
| ۰/۳۶۶ | ۰/۳۸ | ۰/۳۵ | ۰/۳۷ | ۰/۳۴ | ۰/۳۷ | ۰/۳۶ | ۰/۴۰ | آذربایجان غربی |
| ۰/۳۷۶ | ۰/۳۸ | ۰/۳۹ | ۰/۳۸ | ۰/۳۷ | ۰/۳۷ | ۰/۳۷ | ۰/۳۸ | اردبیل |
| ۰/۴۰۶ | ۰/۴۲ | ۰/۴۰ | ۰/۳۹ | ۰/۴۲ | ۰/۴۰ | ۰/۴۲ | ۰/۴۰ | اصفهان |
| ۰/۳۵۲ | ۰/۳۴ | ۰/۳۳ | ۰/۳۹ | ۰/۳۶ | ۰/۳۶ | ۰/۳۳ | ۰/۳۶ | ایلام |
| ۰/۴۱۷ | ۰/۴۴ | ۰/۴۵ | ۰/۳۸ | ۰/۴۶ | ۰/۴۵ | ۰/۴۳ | ۰/۴۴ | بوشهر |
| ۰/۳۶۹ | ۰/۳۸ | ۰/۳۷ | ۰/۳۷ | ۰/۳۵ | ۰/۳۸ | ۰/۳۷ | ۰/۳۷ | تهران |
| ۰/۳۶۰ | ۰/۳۴ | ۰/۳۳ | ۰/۳۴ | ۰/۳۶ | ۰/۳۹ | ۰/۳۷ | ۰/۴۰ | چهارمحال و بختیاری |
| ۰/۴۰۰ | ۰/۴۰ | ۰/۴۳ | ۰/۳۷ | ۰/۳۸ | ۰/۳۹ | ۰/۴۰ | ۰/۴۳ | خراسان |
| ۰/۳۴۰ | ۰/۳۰ | ۰/۳۳ | ۰/۳۶ | ۰/۳۵ | ۰/۳۲ | ۰/۳۲ | ۰/۳۳ | خوزستان |
| ۰/۳۷۳ | ۰/۳۴ | ۰/۳۰ | ۰/۳۸ | ۰/۳۴ | ۰/۳۹ | ۰/۳۹ | ۰/۳۹ | زنجان |
| ۰/۳۵۰ | ۰/۳۳ | ۰/۳۲ | ۰/۳۳ | ۰/۳۳ | ۰/۴۰ | ۰/۳۶ | ۰/۳۹ | سمنان |
| ۰/۲۸۶ | ۰/۴۳ | ۰/۴۲ | ۰/۴۰ | ۰/۳۷ | ۰/۳۴ | ۰/۳۵ | ۰/۳۹ | سیستان و بلوچستان |
| ۰/۳۵۸ | ۰/۳۴ | ۰/۳۶ | ۰/۳۶ | ۰/۳۶ | ۰/۳۸ | ۰/۳۳ | ۰/۳۹ | فارس |
| ۰/۲۶۸ | ۰/۳۷ | ۰/۳۸ | ۰/۳۸ | ۰/۳۵ | -- | -- | -- | قزوین |
| ۰/۳۳۵ | ۰/۴۵ | ۰/۲۸ | ۰/۳۳ | ۰/۳۲ | ۰/۳۱ | -- | -- | قم |
| ۰/۳۴۳ | ۰/۳۱ | ۰/۳۷ | ۰/۳۶ | ۰/۳۴ | ۰/۳۱ | ۰/۳۲ | ۰/۳۳ | کردستان |
| ۰/۳۷۱ | ۰/۳۶ | ۰/۳۹ | ۰/۳۹ | ۰/۳۴ | ۰/۳۷ | ۰/۳۷ | ۰/۳۷ | کرمان |
| ۰/۳۴۷ | ۰/۴۱ | ۰/۳۶ | ۰/۳۴ | ۰/۳۰ | ۰/۳۱ | ۰/۳۵ | ۰/۳۷ | کرمانشاه |
| ۰/۲۵۹ | ۰/۳۶ | ۰/۳۷ | ۰/۳۷ | ۰/۳۸ | ۰/۳۲ | ۰/۳۶ | ۰/۳۷ | کهگیلویه و بویراحمد |
| ۰/۴۰۵ | ۰/۴۴ | ۰/۴۵ | ۰/۳۸ | ۰/۳۶ | -- | -- | -- | گلستان |
| ۰/۳۷۹ | ۰/۳۸ | ۰/۳۹ | ۰/۴۰ | ۰/۴۰ | ۰/۳۴ | ۰/۳۹ | ۰/۳۷ | گیلان |
| ۰/۳۳۵ | ۰/۳۵ | ۰/۳۴ | ۰/۳۳ | ۰/۳۴ | ۰/۳۲ | ۰/۳۵ | ۰/۳۲ | لرستان |
| ۰/۳۷۹ | ۰/۳۸ | ۰/۳۸ | ۰/۳۶ | ۰/۳۸ | ۰/۴۰ | ۰/۴۰ | ۰/۳۷ | مازندران |
| ۰/۳۵۵ | ۰/۳۷ | ۰/۳۷ | ۰/۳۷ | ۰/۳۶ | ۰/۳۰ | ۰/۳۵ | ۰/۳۷ | مرکزی |
| ۰/۳۲۵ | ۰/۳۶ | ۰/۲۷ | ۰/۳۱ | ۰/۳۵ | ۰/۳۰ | ۰/۳۶ | ۰/۳۳ | هرمزگان |
| ۰/۴۲۰ | ۰/۴۰ | ۰/۴۳ | ۰/۴۶ | ۰/۴۴ | ۰/۴۱ | ۰/۳۸ | ۰/۴۳ | همدان |
| ۰/۴۲۰ | ۰/۳۸ | ۰/۳۷ | ۰/۴۱ | ۰/۴۴ | ۰/۴۸ | ۰/۴۵ | ۰/۴۱ | یزد |

جدول ۲: ضریب جینی نقاط روستایی استان های کشور

| استان | سال ۷۴ | سال ۷۵ | سال ۷۶ | سال ۷۷ | سال ۷۸ | سال ۷۹ | سال ۸۰ | متوسط سالهای ۷۴ تا ۸۰ |
|--------------------|--------|--------|--------|--------|--------|--------|--------|-----------------------|
| کل کشور | ۰/۴۳ | ۰/۳۹ | ۰/۴۱ | ۰/۴۲ | ۰/۴۲ | ۰/۴۲ | ۰/۴۱ | ۰/۴۱ |
| آذربایجان شرقی | ۰/۲۸ | ۰/۳۶ | ۰/۴۱ | ۰/۴۰ | ۰/۳۶ | ۰/۴۳ | ۰/۴۶ | ۰/۳۹۹ |
| آذربایجان غربی | ۰/۲۷ | ۰/۳۷ | ۰/۳۸ | ۰/۴۱ | ۰/۳۶ | ۰/۳۹ | ۰/۳۹ | ۰/۳۸۰ |
| اردبیل | ۰/۳۵ | ۰/۳۴ | ۰/۴۱ | ۰/۳۵ | ۰/۴۲ | ۰/۴۲ | ۰/۴۲ | ۰/۳۸۵ |
| اصفهان | ۰/۴۳ | ۰/۴۴ | ۰/۴۳ | ۰/۴۳ | ۰/۴۴ | ۰/۴۲ | ۰/۴۲ | ۰/۴۲۶ |
| ایلام | ۰/۳۲ | ۰/۳۱ | ۰/۳۰ | ۰/۳۶ | ۰/۳۲ | ۰/۳۴ | ۰/۳۴ | ۰/۳۲۶ |
| بوشهر | ۰/۴۰ | ۰/۳۶ | ۰/۴۲ | ۰/۴۱ | ۰/۴۰ | ۰/۳۵ | ۰/۳۷ | ۰/۳۷۲ |
| تهران | ۰/۳۶ | ۰/۳۲ | ۰/۳۵ | ۰/۳۰ | ۰/۳۳ | ۰/۳۸ | ۰/۳۴ | ۰/۳۳۹ |
| چهارمحال و بختیاری | ۰/۴۱ | ۰/۳۸ | ۰/۴۱ | ۰/۴۲ | ۰/۳۸ | ۰/۳۹ | ۰/۳۴ | ۰/۳۹۰ |
| خراسان | ۰/۴۶ | ۰/۴۲ | ۰/۴۲ | ۰/۴۵ | ۰/۴۲ | ۰/۴۳ | ۰/۴۱ | ۰/۴۲۹ |
| خوزستان | ۰/۳۴ | ۰/۲۹ | ۰/۳۷ | ۰/۳۵ | ۰/۳۶ | ۰/۳۱ | ۰/۲۴ | ۰/۳۱۸ |
| زنجان | ۰/۴۲ | ۰/۳۸ | ۰/۴۳ | ۰/۴۴ | ۰/۴۶ | ۰/۴۴ | ۰/۴۲ | ۰/۴۲۶ |
| سمنان | ۰/۴۶ | ۰/۴۲ | ۰/۴۳ | ۰/۳ | ۰/۳۹ | ۰/۴ | ۰/۴۱ | ۰/۳۹۷ |
| سیستان و بلوچستان | ۰/۴۷ | ۰/۴۳ | ۰/۳۸ | ۰/۳۸ | ۰/۴۴ | ۰/۳۷ | ۰/۴۲ | ۰/۴۱۱ |
| فارس | ۰/۳۹ | ۰/۳۶ | ۰/۳۹ | ۰/۴۲ | ۰/۴۲ | ۰/۳۸ | ۰/۳۴ | ۰/۳۸۵ |
| فروین | -- | -- | -- | ۰/۴۳ | ۰/۴۲ | ۰/۴۲ | ۰/۴۲ | ۰/۴۲۸ |
| قم | -- | -- | ۰/۴۶ | ۰/۳۶ | ۰/۳۶ | ۰/۳۴ | ۰/۳۸ | ۰/۳۷۷ |
| کردستان | ۰/۳۳ | ۰/۳۲ | ۰/۴۲ | ۰/۴۲ | ۰/۳۷ | ۰/۳۸ | ۰/۳۷ | ۰/۳۶۶ |
| کرمان | ۰/۴۴ | ۰/۴۲ | ۰/۳۸ | ۰/۴۴ | ۰/۴۵ | ۰/۴۳ | ۰/۴۱ | ۰/۴۲۲ |
| کرمانشاه | ۰/۴۲ | ۰/۲۹ | ۰/۲۹ | ۰/۳۲ | ۰/۳۷ | ۰/۳۳ | ۰/۳۷ | ۰/۳۲۸ |
| کهگیویه و بویراحمد | ۰/۴۵ | ۰/۴۵ | ۰/۴۷ | ۰/۳۴ | ۰/۴۳ | ۰/۴۱ | ۰/۳۹ | ۰/۴۱۸ |
| گلستان | -- | -- | -- | ۰/۴۴ | ۰/۴۸ | ۰/۴۹ | ۰/۴۶ | ۰/۴۶۵ |
| گیلان | ۰/۴۲ | ۰/۴۱ | ۰/۴۴ | ۰/۴۷ | ۰/۴۴ | ۰/۴۳ | ۰/۴۱ | ۰/۴۳۲ |
| لرستان | ۰/۳۷ | ۰/۳۱ | ۰/۳۵ | ۰/۳۶ | ۰/۳۴ | ۰/۳۶ | ۰/۳۴ | ۰/۳۴۷ |
| مازندران | ۰/۴۵ | ۰/۴۱ | ۰/۴۳ | ۰/۴۳ | ۰/۴۲ | ۰/۴۱ | ۰/۳۹ | ۰/۴۱۸ |
| مرکزی | ۰/۴۱ | ۰/۳۹ | ۰/۳۹ | ۰/۴۱ | ۰/۳۹ | ۰/۴۴ | ۰/۴۴ | ۰/۴۰۶ |
| هرمزگان | ۰/۳۹ | ۰/۳۵ | ۰/۴۲ | ۰/۳۸ | ۰/۴۰ | ۰/۳۴ | ۰/۳۶ | ۰/۳۷۵ |
| همدان | ۰/۴۰ | ۰/۴۸ | ۰/۳۸ | ۰/۳۹ | ۰/۴۳ | ۰/۴۹ | ۰/۴۱ | ۰/۴۳۳ |
| یزد | ۰/۴۲ | ۰/۴۰ | ۰/۵۰ | ۰/۵۵ | ۰/۴۵ | ۰/۴۵ | ۰/۴۶ | ۰/۴۶۲ |

جدول ۳: خوشه‌بندی استان‌های کشور با استفاده از ضریب جینی نقاط شهری

| | |
|----------|--|
| خوشه اول | بوشهر، خوزستان، قم، کردستان، کرمانشاه، لرستان، هرمزگان |
| خوشه دوم | آذربایجان شرقی، آذربایجان غربی، اردبیل، ایلام، تهران، چهارمحال و بختیاری، زنجان، سمنان، سیستان و بلوچستان، فارس، قزوین، کرمان، کهگیلویه و بویراحمد، گیلان، مازندران، مرکزی |
| خوشه سوم | اصفهان، خراسان، گلستان، همدان، یزد |

۵ خوشه‌بندی استان‌های کشور بر اساس برآورد ضرایب جینی و تحلیل نتایج بدست آمده

در این بخش با استفاده از روش آنالیز خوشه‌ای^۴ استان‌های کشور با استفاده از متوسط ضریب جینی نقاط شهری و روستایی طی سالهای ۷۴ تا ۸۰ در ۳ گروه قرار گرفته‌اند تا استان‌های مشابه از نظر چگونگی توزیع درآمد خانوارهای شهری و روستایی مشخص گردند. نتایج به تفکیک نقاط شهری و روستایی در جداول ۳ و ۴ آمده است.

بر اساس ضرایب جینی بدست آمده برای نقاط شهری استان‌های کشور می‌توان گفت توزیع درآمد کل کشور (نقاط شهری) در حالت نیمه متعادل است و این مورد برای استان‌ها نیز به همین گونه می‌باشد، با این وجود بین برخی از استان‌های کشور از نظر چگونگی توزیع درآمد خانوارهای شهری تفاوتی وجود دارد که این تفاوت را می‌توان با خوشه‌بندی آنها بیشتر مشخص کرد. براساس این خوشه‌بندی استان‌های کشور در ۳ خوشه واقع شده‌اند، متوسط ضریب جینی استان‌های خوشه اول برابر ۰/۳۳، خوشه دوم ۰/۳۷ و خوشه سوم ۰/۴۱ می‌باشد. در این بین استان بوشهر عادلانه‌ترین توزیع درآمد نقاط شهری را داشته و استان‌های همدان و یزد نا عادلانه‌ترین وضعیت را دارا بوده‌اند.

بر اساس ضرایب جینی بدست آمده برای نقاط روستایی استان‌های کشور می‌توان گفت توزیع درآمد کل کشور (نقاط روستایی) در حالت نیمه متعادل است و این مورد برای استان‌ها نیز به همین گونه می‌باشد، با این وجود بین برخی از استان‌های کشور از نظر چگونگی توزیع درآمد خانوارهای روستایی تفاوتی وجود دارد که این تفاوت را می‌توان با خوشه‌بندی آنها بیشتر مشخص کرد. براساس این خوشه‌بندی استان‌های کشور در ۳ خوشه واقع شده‌اند متوسط ضریب جینی استان‌های خوشه اول برابر ۰/۳۴، خوشه دوم ۰/۴۰ و خوشه سوم ۰/۴۶ می‌باشد.

4) Cluster Analysis

جدول ۴: خوشه‌بندی استان‌های کشور با استفاده از ضریب جینی نقاط روستایی

| | |
|----------|--|
| خوشه اول | ایلام، تهران، خوزستان، کردستان، کرمانشاه، لرستان |
| خوشه دوم | آذربایجان شرقی، آذربایجان غربی، اردبیل، اصفهان، چهارمحال و بختیاری، بوشهر، خراسان، زنجان، سمنان، سیستان و بلوچستان، فارس، قزوین، قم، کرمان، کهگیلویه و بویر احمد، گیلان، مازندران، مرکزی، هرمزگان، همدان |
| خوشه سوم | گلستان، یزد |

در این بین استان خوزستان عادلانه‌ترین توزیع درآمد نقاط روستایی را داشته و استان‌های گلستان ناعادلانه‌ترین وضعیت را دارا بوده‌اند.

از نتایج دیگری که می‌توان از برآورد ضریب جینی گرفت این است که به‌طور کلی توزیع درآمد در نقاط روستایی کشور ناعادلانه‌تر از توزیع درآمد در نقاط شهری است و این موضوع در اکثر استان‌های کشور نیز دیده می‌شود و تنها در استان‌های ایلام، تهران، خوزستان و کرمانشاه توزیع درآمد خانوارهای روستایی عادلانه‌تر از نقاط شهری است.

مراجع

- [۱] نتایج آمارگیری از هزینه و درآمد خانوارهای شهری، (۱۳۷۴ تا ۱۳۸۰)، مرکز آمار ایران
- [۲] نتایج آمارگیری از هزینه و درآمد خانوارهای روستایی، (۱۳۷۴ تا ۱۳۸۰)، مرکز آمار ایران
- [3] Kokwani, N.C (1980), Incomm Equality and poverty, A World Bank Research Publication.
- [4] Khan, Azizur Rahman (1993), Structural Adjustment and Income distribution, ILO.

شرایط آغازین برای مدل‌های تصادفی - انتقالی در داده‌های طولی با پاسخ دودویی کامل و ناقص

زهرا رضایی قهرودی^۱، مجتبی گنجعلی^۲

^۱ کارشناس ارشد آمار دانشگاه شهید بهشتی

^۲ عضو هیات علمی دانشگاه شهید بهشتی

چکیده: فرض کنید برای داده‌های طولی با پاسخ دودویی مدلی را استفاده کنیم که در آن پاسخ در زمان t به پاسخهای زمانهای گذشته و نیز به عوامل طبیعی غیرقابل مشاهده شده به عنوان اثرهای تصادفی وابسته است. چنین مدلی مدل انتقالی - تصادفی نامیده می‌شود. همچنین فرض کنید فرایند در زمان صفر شروع ولی در زمان $J, J \geq 1$ مشاهده شده باشد. بنابراین برای فرد i ام فرض کنید مشاهدات در زمانهای $J, J+1, \dots, T$ مشاهده شده است. از آنجا که برای مثال در مدل مارکوف مرتبه اول پاسخ در زمان J توسط مدل انتقالی - تصادفی به پاسخ در زمان $J-1$ وابسته است ولی پاسخ $J-1$ مشاهده نمی‌شود، مسئله شرایط آغازین ایجاد می‌شود. در این مقاله با استفاده از مفهوم متغیر پنهان برای پاسخ دودویی یک فرایند مارکوف مرتبه اول در نظر گرفته و فرض می‌کنیم که اختلالی که فرایند را ایجاد می‌کند به دو قسمت خطای اندازه‌گیری و تاثیر تصادفی تجزیه می‌شود. روشهای مختلف مواجهه با شرایط آغازین را در دو مورد (i) پاسخهای دودویی کامل و (ii) پاسخهای دودویی ناقص مورد بحث و بررسی قرار می‌دهیم.

واژه‌های کلیدی: داده‌های طولی، مدل تصادفی - انتقالی، فرایند مارکوف مرتبه اول، شرایط آغازین

مقدمه

اخیرا مطالعات طولی از اهمیت به سزایی برخوردار شده است. از مشخصه‌های بارز مطالعات طولی اندازه‌های مکرر برای افراد مختلف در طول زمان است. در مطالعات طولی همبستگی بین اندازه‌ها در طول زمان مورد بررسی قرار می‌گیرد. در داده‌های طولی امکان گم شدن داده‌ها نیز وجود دارد. در این مقاله هدف بررسی مسئله شرایط آغازین در مدل‌های اثرات تصادفی - انتقالی با پاسخ دودویی کامل و ناقص است. فرض می‌کنیم فرایند در زمان صفر شروع ولی در زمان $J, J \geq 1$ مشاهده شده باشد. بنابراین برای فرد i ام فرض کنید مشاهدات در زمانهای $J, J+1, \dots, T$ مشاهده شده است. از آنجا که برای مثال در مدل مارکوف مرتبه اول پاسخ در زمان J توسط مدل انتقالی - تصادفی به پاسخ در زمان $J-1$ وابسته است ولی پاسخ

۱ - J مشاهده نمی‌شود، مسئله شرایط آغازین ایجاد می‌شود. با توجه به اینکه بحث ما جنبه‌های ضروری مسئله شرایط آغازین و راه‌حل‌های آن می‌باشد، در این مقاله روش‌های مختلف مواجهه با مقادیر مشاهده نشده (در داده‌های کامل و ناقص) که ممکن است دنباله‌ای از y_{it} برای $t = 0, \dots, J-1$ باشد و یا در حالت خاص تنها y_{i0} ($J=1$) باشد را مورد بحث و بررسی قرار می‌دهیم.

۱ گم شدن داده‌ها در مطالعات طولی

در داده‌های طولی گم شدن داده‌ها به دو صورت الگوی عمومی و الگوی یکنواکه همان انصراف است طبقه‌بندی می‌شود. در الگوی عمومی هیچ روند خاصی مشاهده نمی‌شود و مقادیر گمشده در هر زمانی می‌تواند رخ دهد. بالاخص آزمودنی خاصی می‌تواند از مطالعه در زمان خاصی خارج و در زمان دیگری دوباره به مطالعه بازگشت کند. ولی انصراف در داده‌های طولی به این معناست که وقتی آزمودنی در زمان t پاسخ نمی‌دهد هیچوقت به مطالعه بازگشت ندارد و ما پاسخهای زمانهای $t+1, \dots$ را نیز از دست می‌دهیم. به منظور مدل‌بندی همزمان انصراف و پاسخ از دیدگاه تابع درست‌نمایی تعاریف اساسی را که اولین بار در رابین (۱۹۷۶) مطرح شده است، می‌آوریم. فرض کنید در صورت عدم داشتن داده گمشده برای پاسخهای آزمودنی i ام پاسخهای دودویی y_{ij} برای $j = 1, 2, \dots, T$ مشاهده شود. برای نشان دادن آن که متغیر y گمشده است از تابع نشانگر R استفاده می‌کنیم که به صورت زیر تعریف می‌شود.

$$y_{it} = \begin{cases} 1 & \text{اگر پاسخ فرد } i \text{ ام در زمان } t \text{ مشاهده شده باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

همچنین پاسخ دودویی y_{it} را برآمده از یک متغیر پنهان y_{it}^* فرض می‌کنیم، که به صورت زیر تعریف می‌شود.

$$y_{it} = \begin{cases} 1 & y_{it}^* > c \\ 0 & \text{در غیر این صورت} \end{cases}$$

برای درک مفهوم متغیر پنهان فرض کنید که تنها مثبت شده است که آزمودنی چاق یا لاغر است. اما اگر فردی چاق ($y_{it} = 1$) است این به خاطر وزن زیاد ($y_{it}^* > c$) او بوده برای بردار تصادفی $Y_i = (Y_{i1}, \dots, Y_{in})^T$ بردار نشانگر پاسخ $R_i = (R_{i1}, \dots, R_{in})^T$ را استفاده می‌کنیم.

۲ مدل اثرهای تصادفی و انتقالی

این مدل ترکیبی از مدل اثرهای تصادفی و مدل انتقالی است به این صورت که پاسخ در زمان t به پاسخ در زمانهای گذشته و نیز به عوامل طبیعی غیر قابل مشاهده شده به عنوان اثرهای تصادفی وابسته است. به عنوان مثال برای بررسی علت تصادف برای آزمودنیها باید این مسئله را بررسی کنیم که آیا علت تصادف به دلیل شتابزدگی فرد (اثر تصادفی) و یا به دلیل عادت فرد به تصادف (پاسخهای گذشته) بوده است. برای بیان مدلی مناسب به صورت زیر عمل می‌کنیم. فرض کنید که $y_i = (y_{i1}, \dots, y_{iT_i})^T$ پاسخهای دودویی برای n آزمودنی ($i = 1, \dots, n$) در T_i موقعیت متفاوت ($t = 1, \dots, T_i$) با بردار p بعدی متغیرهای تبیینی $x_{it} = (x_{it1}, \dots, x_{itp})^T$ که به زمان وابسته است فرض شود. مدل زیر مدلی مناسب برای بررسی اثرهای تصادفی و پاسخهای گذشته است.

$$Y_{it}^* = x_{it}'\beta + z_{it}'\tau_i + \sum_{j=1}^q \gamma_j y_{i,t-j} + \varepsilon_{it}$$

که در آن Y_{it}^* متغیر پنهانی متناظر با Y_{it} است که قبلاً معرفی شده است و $y_{i,t-j}$ پاسخ آزمودنی i ام در زمان $t - j$ است که تاثیر آن بر y_{it} توسط پارامتر γ_j سنجیده می‌شود. β بردار p بعدی ضرایب متغیرهای تبیینی و $\tau_i = (\tau_{i1}, \dots, \tau_{is})^T$ اثرهای تصادفی برای $i = 1, \dots, n$ ، با تابع چگالی چندمتغیره معلوم $g(\cdot; \phi)$ فرض می‌شوند. همچنین z_{it} زیر مجموعه‌ای از x_{it} است. در برخورد با رویه انصراف، برای مدل‌بندی مکانیسم گم شدن، به روشهای مختلف می‌توان عمل کرد. مدل مکانیسم گم شدن که الگوی گم شدن انصراف و تناوبی را در نظر می‌گیرد به صورت زیر بیان می‌شود.

$$R_{it}^* = v_{it}'\alpha_l + \theta_2 \tau_i + \varepsilon_{it}$$

که R_{it}^* متغیر پنهان برای مدل مربوط به مکانیسم گم شدن است. α_l بردار پارامتری حاصل از تاثیر متغیرهای تبیینی بر مکانیسم گم شدن و θ_2 پارامتر منعکس کننده تاثیر تصادفی بر مکانیسم گم شدن است. این مدل به مدل تاثیر تصادفی مشترک (۱۹۹۵) موسوم است.

۳ مسئله شرایط آغازین و بعضی راه‌حل‌های معمول

با توجه به اینکه بحث ما جنبه‌های ضروری مسئله شرایط آغازین و راه‌حل‌های آن می‌باشد، بنابراین یک فرآیند مارکوف مرتبه اول را در نظر می‌گیریم و متغیرهای برون‌زا (متغیرهای تبیینی) را فعلاً در نظر نمی‌گیریم. به علاوه فرض شده است که احتمالی که فرآیند را ایجاد می‌کند به دو قسمت

2) Follmann D. and Wu M.

خطای اندازه‌گیری و تاثیر تصادفی تجزیه می‌شود. فرآیند به وسیله یک متغیر پنهان دودویی Y_{it}^* تعریف شده است که از مدل زیرین تبعیت می‌کند:

$$Y_{it}^* = \beta_0 + \gamma y_{i,t-1} + \varepsilon_{it} \quad i = 1, \dots, n, \quad t = 1, \dots, T \quad (1.3)$$

که در آن

$$\varepsilon_{it} = \tau_i + U_{it}$$

و

$$E(U_{it}) = 0 = E(\tau_i), \quad E(\tau_i^2) = \sigma_\tau^2, \quad E(U_{it}^2) = \sigma_u^2 = 1$$

$$E(U_{it}U_{i't''}) = 0, \quad t \neq t'', \quad E(\tau_i U_{i't''}) = 0 \quad \forall i, i', t''$$

که در آن τ_i اثر تصادفی و U_{it} خطای مدل است و احتمالهای تغییر وضعیت برای فرد i ام در زمان t ام به شرط τ_i به صورت زیر است.

$$f(y_{it}|y_{i,t-1}, \tau_i) = \Phi\{\beta_0 + \gamma y_{i,t-1} + \tau_i \cdot [2y_{it} - 1]\}$$

که در آن Φ تابع توزیع تجمعی نرمال است. در اینجا از یک مدل پروبیت استفاده می‌شود، اما بدون از دست دادن هیچ کلیتی مدل لوژستیک نیز می‌تواند مورد بررسی قرار گیرد. اگر شروع مشاهده فرآیند در زمان J باشد به این معنی که شروع اولیه فرآیند زمان 0 بوده و محقق مشاهدات خود را از زمان J ثبت کرده است، مسئله شرایط آغازین ایجاد می‌شود. در حالتی که $J > 1$ باشد، $f(y_{iJ}|\tau_i)$ نشان می‌دهد که چگونه مسئله مقادیر اولیه به وجود می‌آید. توجه کنید که $f(y_{iJ}|\tau_i)$ به تمام دنباله‌های ممکن پیشامدهای قبل از زمان J ام که مشاهده نمی‌شود و منجر به یک مقدار خاص برای $y_{i,J}$ می‌شود وابسته است. لازم به ذکر است که عموماً محققین مقادیر متغیرهای تمیینی مناسب را در دوره‌های زمانی $t = 1, \dots, J-1$ نیز نمی‌دانند. اکنون پنج روش مختلف مواجهه با مقادیر مشاهده نشده که ممکن است دنباله‌ای از y_{it} برای $t = 0, \dots, J-1$ باشد و یا در حالت خاص ($J = 1$) (y_{i0} مشاهده نشده است) را که ما در اینجا در نظر گرفته‌ایم، بکار می‌بریم. در هر حالت سعی شده است که تابع درستنمایی برای بهینه‌سازی داده شود.

روش اول:

اگر $J = 1$ باشد، y_{i1} را یک مقدار ثابت غیرتصادفی و معلوم در نظر می‌گیریم و تابع درستنمایی را به شرط y_{i1} به صورت زیر می‌نویسیم یعنی:

$$\prod_{i=1}^n f(y_{i2}, \dots, y_{iT_i}|y_{i1}) = \prod_{i=1}^n \int \prod_{t=2}^{T_i} f(y_{it}|\tau_i, y_{i,t-1}) g(\tau_i) d\tau_i \quad (2.3)$$

که در این روش عملاً اثر پاسخهای مشاهده شده زمان 1 به عنوان متغیر تمیینی در نظر گرفته می‌شود و از یک مدل مارکوف مرتبه اول استفاده شده است.

روش دوم:

حالتی است که ایتکن و آلفو^۳ (۲۰۰۰) معرفی کرده‌اند که آنها نیز مانند روش اول عمل می‌کنند با این تفاوت که پاسخ مشاهده شده $y_{i\lambda}$ اثر تصادفی را نیز تحت تاثیر قرار می‌دهد. تابع درستیابی ارائه شده توسط ایتکن و آلفو (۲۰۰۰) به صورت زیر است.

$$\prod_{i=1}^n f(y_{i2}, \dots, y_{iT_i} | y_{i\lambda}) = \prod_{i=1}^n \int \prod_{t=2}^{T_i} f(y_{it} | \tau_i, y_{i,t-1}) g(\tau_i | y_{i\lambda}) d\tau_i \quad (3.3)$$

که در آن همانطور که مشاهده می‌کنید توزیع شرطی اثر تصادفی به شرط $y_{i\lambda}$ بایستی در نظر گرفته شود.

روش سوم:

روش سوم است که برای پاسخ در زمان ۱ تابع توزیع مجزایی ارائه دهیم، که در این روش به صورت زیر عمل می‌شود.

$$\prod_{i=1}^n f(y_{i\lambda}, \dots, y_{iT_i}) = \prod_{i=1}^n \int f(y_{i\lambda} | \tau_i) \prod_{t=2}^{T_i} f(y_{it} | \tau_i, y_{i,t-1}) g(\tau_i) d\tau_i \quad (4.3)$$

توجه کنید که در رابطه فوق $f(y_{i\lambda} | \tau_i)$ نیز مدل بندی می‌شود.

روش چهارم:

اگر فرآیند تعادلی باشد، احتمالهای حاشیه‌ای برای حالت $y_{it} = 1$ برای همه t ها (با فرض گذشته بی‌نهایت) به صورت زیر است.

$$\Pi_1(\tau_i) = \frac{\Phi(\beta_0 + \tau_i)}{1 - \Phi(\beta_0 + \gamma + \tau_i) + \Phi(\beta_0 + \tau_i)} \quad (5.3)$$

و احتمال حاشیه‌ای برای حالت $y_{it} = 0$ ، $\Pi_0(\tau_i) = 1 - \Pi_1(\tau_i)$ است. که Π_0 و Π_1 به ترتیب احتمالهای حاشیه‌ای در حالت پاسخ ۱ و ۰ است. لازم به ذکر است که احتمالهای تعادلی به صورت زیر به دست می‌آیند:

$$P(Y_{it} = 1 | \tau_i) = \sum_{y_{i,t-1}} P(Y_{it} = 1 | Y_{i,t-1} = y_{i,t-1}, \tau_i) P(Y_{i,t-1} = y_{i,t-1} | \tau_i)$$

که چون در شرایط تعادلی هستیم، خواهیم داشت

$$P(Y_{i,t-1} = y_{i,t-1} | \tau_i) = P(Y_{it} = y_{it} | \tau_i)$$

با استفاده از این واقعیت که $\Pi_0(\tau_i) = 1 - \Pi_1(\tau_i)$ است، خواهیم داشت:

$$P(Y_{it} = 1 | \tau_i) = P(Y_{it} = 1 | Y_{i,t-1} = 0, \tau_i) [1 - \Pi_1(\tau_i)] + P(Y_{it} = 1 | Y_{i,t-1} = 1, \tau_i) \Pi_1(\tau_i)$$

یعنی

$$\Pi_1(\tau_i) = \Phi(\beta_0 + \tau_i) \Pi_0(\tau_i) + \Phi(\beta_0 + \gamma + \tau_i) \Pi_1(\tau_i)$$

و بنابراین اگر فرآیند در وضعیت تعادل باشد،

$$P(Y_{iJ} = y_{iJ} | \tau_i) = \Pi_1(\tau_i)^{y_{iJ}} \Pi_0(\tau_i)^{1-y_{iJ}}$$

است. روش دیگر مواجهه با شرایط آغازین استفاده از احتمالات تعادلی است. روش پنجم:

روش آخر که هکمن^۴ (۱۹۸۱) پیشنهاد کرده است حالتی است که در آن برای پاسخها در زمان $t = 1$ اثر تصادفی ای متفاوت با اثر تصادفی در زمانهای $t = 2, \dots, T_i$ در نظر می‌گیریم. مدل ارائه شده توسط هکمن (بدون حضور متغیرهای تبیینی) به صورت زیر است:

$$Y_{it}^* = \beta_0 + \gamma y_{i,t-1} + \varepsilon_{it} \quad i = 1, \dots, n, t = 2, \dots, T$$

که در آن

$$\varepsilon_{it} = \tau_{i1} + U_{it}, \quad t = 2, \dots, T$$

$$Y_{i1}^* = \beta_0 + \varepsilon_{i1}$$

و

$$\varepsilon_{i1} = \tau_{i2} + U_{i1}$$

و τ_{i2} و τ_{i1} اثرات تصادفی با میانگین صفر هستند که به طور نرمال توزیع شده است. همچنین τ_{i2} با τ_{i1} همبسته‌اند. بنابراین تابع درستنمایی ارائه شده توسط هکمن به صورت زیر است:

$$\prod_{i=1}^n f(y_{i1}, \dots, y_{iT_i}) = \prod_{i=1}^n \int \int f(y_{i1} | \tau_{i2}) \times \prod_{t=2}^{T_i} f(y_{it} | \tau_{i1}, y_{i,t-1}) g(\tau_{i1}, \tau_{i2}) d\tau_{i1} d\tau_{i2} \quad (6 \cdot 3)$$

که برآورد پارامترهای مدل به روش ماکسیمم درستنمایی با در نظر گرفتن همبستگی بین τ_{i2} و τ_{i1} به روش مربع بندی گاوسی به دست می‌آید. هکمن (۱۹۸۱) در یک مطالعه شبیه‌سازی برتری این روش را نسبت به برخی روشهای دیگر مورد بررسی قرار داده است. در بخش بعد

4) Heckman

مسئله شرایط آغازین را با در نظر گرفتن پاسخ گمشده در پنج روش فوق مورد بحث و بررسی قرار می‌دهیم.

مسئله دیگری که در شرایط آغازین مطرح است، مسئله پارامترهای فرعی است. در این حالت اثرهای ثابت τ_i را به عنوان پارامتر در نظر می‌گیریم و تابع درستیابی شرطی y_{i2}, \dots, y_{iT} به شرط y_{i1} که خود y_{i1} تابعی از پاسخهای قبلی مشاهده نشده است، را به صورت زیر بیان می‌کنیم:

$$f(y_{i2}, \dots, y_{iT} | y_{i1}) = \prod_{i=1}^n \prod_{t=2}^T \frac{\Phi\{\beta_0 + \gamma y_{i,t-1} + \tau_i [2y_{it} - 1]\} \times P(y_{i1} | \tau_i)}{P(y_{i1} | \tau_i)}$$

$$= \prod_{i=1}^n \prod_{t=2}^T \Phi\{\beta_0 + \gamma y_{i,t-1} + \tau_i [2y_{it} - 1]\} \quad (7.3)$$

همانگونه که ملاحظه می‌شود توزیع $f(y_{i1} | \tau_i)$ از معادله بالا حذف می‌شود و این نشان می‌دهد که استفاده از روش پارامترهای فرعی، راه حلی برای رفع شرایط آغازین است. این مدل یک مدل رگرسیون خطی با اثرات ثابت است و محاسبه آن بسیار راحت است. با ماکسیم کردن تابع (7.3) نسبت به τ_i برای $i = 1, \dots, n$ و $\hat{\beta}$ و $\hat{\gamma}$ زمانی که n و T بزرگند برآوردهای سازگاری هستند (هکمن، ۱۹۸۱). در حالتی که T کوچک است زیاد بودن پارامترها مشکلی است که این روش با آن مواجه است.

۴ شرایط آغازین برای مدل‌های با پاسخ دودویی و ناقص

همانطور که در بخش قبل مسئله شرایط آغازین را بیان کردیم، در این بخش نیز به این مسئله می‌پردازیم با این تفاوت که در این بخش پاسخهای گمشده در داده‌های ما وجود دارد. مدل در نظر گرفته شده برای پاسخ مدل مارکوف مرتبه اول همراه با اثر تصادفی است و مدلی که برای پاسخ‌های گمشده در نظر گرفته می‌شود مدل پیشنهادی دیگل و کنورد (۱۹۹۴) است که مکانیسم گم شدن در زمان t به پاسخهای زمانهای t و $t-1$ بستگی دارد با این تفاوت که در مدل دیگل و کنورد اثرهای تصادفی را نیز وارد مدل می‌کنیم. مدل فوق به صورت زیر بیان می‌شود.

$$Y_{it}^* = \beta + \gamma y_{i,t-1} + \varepsilon_{it1} \quad i = 1, \dots, n, \quad t = 1, \dots, T \quad (10.4)$$

که در آن

$$\varepsilon_{it1} = \theta_1 \tau_i + U_{it1}$$

برای مکانیسم گم شدن داریم:

$$R_{it}^* = \alpha + \psi_{12} y_{i,t-1} + \psi_{22} y_{it} + \varepsilon_{it2} \quad i = 1, \dots, n, \quad t = 2, \dots, T$$

که در آن فرض شده برای پاسخ در زمان ۱ مشاهده گمشده نداریم و ψ_{12} و ψ_{22} اثر پاسخهای قبلی و جاری را برگمشدن پاسخ در زمان جاری منعکس می‌کنند. همچنین

$$\varepsilon_{it2} = \theta_2 \tau_i + U_{it2}$$

و داریم:

$$E(U_{it1}) = E(U_{it2}) = 0 = E(\tau_i), \quad E(\tau_i^2) = 1$$

$$E(U_{it1}^2) = E(U_{it2}^2) = \sigma_u^2 = 1, \quad E(U_{it1} U_{it'1}) = E(U_{it2} U_{it'2}) = 0, \quad t \neq t''$$

$$E(\tau_i U_{it'1}) = E(\tau_i U_{it'2}) = 0 \quad \forall i, t', t''$$

که در آن τ_i اثرهای تصادفی و U_{it1} و U_{it2} به ترتیب خطاهای مدل پاسخ و مکانیسم گم شدن هستند. توجه داشته باشید که مکانیسم گم شدن در صورتی کاملاً تصادفی است که پارامترهای $\psi_{12}, \psi_{22}, \theta_2$ معنی‌دار نباشند. احتمالهای تغییر وضعیت برای فرد i ام در زمان t ام به شرط τ_i به صورت زیر است.

$$f(y_{it}|y_{i,t-1}, \tau_i) = \Phi\{\beta + \gamma_1 y_{i,t-1} + \tau_i [2y_{it} - 1]\}$$

و

$$P(R_{it}|y_{i,t-1}, y_{it}, \tau_i) = \Phi\{\alpha + \psi_{12} y_{i,t-1} + \psi_{22} y_{it} + \theta_2 \tau_i [2R_{it} - 1]\}$$

که در آن Φ تابع توزیع تجمعی نرمال است.

در مواجهه با مسئله شرایط آغازین وقتی داده‌های گمشده داریم ممکن است به پنج روش برخورد کنیم. در زیر این روشها برای حالتی که فرآیند در زمان ۱ شروع به مشاهده شدن می‌کند و مشاهدات قبل از این در دسترس نیست، مورد بررسی قرار گرفته‌اند. در هر حالت تابع درستنمایی برای بهینه‌سازی داده شده است.

روش اول:

در این روش پاسخ y_{i1} را به عنوان یک پاسخ ثابت و معلوم در نظر می‌گیریم، تابع درستنمایی به شکل زیر است.

$$\begin{aligned} & \prod_{i=1}^n f(y_{i2}, \dots, y_{iT}, R_{i2}, \dots, R_{iT} | y_{i1}) \\ &= \prod_{i=1}^n \int \prod_{t=2}^T f(y_{it} | \tau_i, y_{i,t-1}) P(R_{it} | y_{it}, y_{i,t-1}, \tau_i) g(\tau_i) d\tau_i \quad (2.4) \end{aligned}$$

که با توجه به تابع چگالی‌های ارائه شده برای پاسخ و مکانیسم گم شدن، تابع درستنمایی فوق را می‌توان ماکسیمم نمود.

روش دوم:

روش دوم مانند روش اول است با این تفاوت که در آن اثر تصادفی به y_{i1} وابسته است که تابع درستیابی در این روش به شکل زیر است.

$$\prod_{i=1}^n f(y_{i2}, \dots, y_{iT}, R_{i2}, \dots, R_{iT} | y_{i1})$$

$$= \prod_{i=1}^n \int \prod_{t=2}^T f(y_{it} | \tau_i, y_{i,t-1}) P(R_{it} | y_{it}, y_{i,t-1}, \tau_i) g(\tau_i | y_{i1}) d\tau_i \quad (3.4)$$

روش سوم:

روشی است که در آن توزیع مجزایی را برای پاسخ مشاهده شده y_{i1} در نظر گیریم. در این حالت تابع درستیابی به شکل زیر است.

$$\prod_{i=1}^n f(y_{i1}, \dots, y_{iT}, R_{i1}, \dots, R_{iT})$$

$$= \prod_{i=1}^n \int f(y_{i1} | \tau_i) \prod_{t=2}^T f(y_{it} | \tau_i, y_{i,t-1}) P(R_{it} | y_{it}, y_{i,t-1}, \tau_i) g(\tau_i) d\tau_i \quad (4.4)$$

که در آن $f(y_{i1} | \tau_i)$ می‌تواند جداگانه مدل‌بندی شود.

روش چهارم:

اگر فرآیند در وضعیت تعادلی باشد آنگاه خواهیم داشت:

$$f(y_{i1} | \tau_i) = \Pi_1(\tau_i)^{y_{i1}} \Pi_0(\tau_i)^{1-y_{i1}} \quad (5.4)$$

که در آن Π_1 و Π_0 به ترتیب احتمالهای حاشیه‌ای برای $y_{it} = 1$ و $y_{it} = 0$ است و می‌توان از این توزیع در (۴.۴) استفاده کرد.

روش پنجم:

نهایتاً روش آخر استفاده از مدل همگن است که با کمی تغییر و با در نظر گرفتن مکانیسم گم شدن در تابع درستیابی، طبق مدل در نظر گرفته شده زیر خواهیم داشت:

$$Y_{it}^* = \beta + \gamma_1 y_{i,t-1} + \varepsilon_{it} \quad i = 1, \dots, n, t = 2, \dots, T$$

$$\varepsilon_{it} = \tau_{i1} + U_{it}$$

$$Y_{i1}^* = \beta + \varepsilon_{i1}$$

که در آن

$$\varepsilon_{i1} = \tau_{i2} + U_{i2}$$

و

$$R_{it}^* = \alpha + \psi_{12}y_{i,t-1} + \psi_{22}y_{it} + \varepsilon_{it2} \quad t = 2, \dots, T$$

اگر فرض کنیم ε_{it3} از ε_{it2} و ε_{it1} مستقل است، تابع درستنمایی مدل فوق به صورت زیر است.

$$\prod_{i=1}^n f(y_{i1}, \dots, y_{iT}, R_{i2}, \dots, R_{iT}) = \prod_{i=1}^n \int \int f(y_{i1} | \tau_{i2}) \prod_{t=2}^{T_i} f(y_{it} | \tau_{i1}, y_{i,t-1}) \times P(R_{it} | y_{it}, y_{i,t-1}) g(\tau_{i1}, \tau_{i2}) d\tau_{i1} d\tau_{i2} \quad (6.4)$$

به عنوان کار بعدی مقایسه‌ای عمومی بین این روشها پیشنهاد می‌شود. می‌توان همچون هکمن (۱۹۸۱) که مدل با استفاده از روشهای (۱)، (۲)، (۳) و (۵) را در داده‌های کامل با استفاده از شبیه‌سازی مونت کارلو، مقایسه کرده و به این نتیجه رسیده است که روش (۵) نتایج بهتری ارائه می‌دهد، روش (۴) را نیز با این روشها مقایسه نمود. همچنین مقایسه‌ای بین روشهای ارائه شده برای داده‌های با پاسخ‌های گمشده، پیشنهاد می‌شود. برای ملاحظه چگونگی کاربرد برخی از روشهای ذکر شده در داده‌های واقعی به رضایی (۱۳۸۲) مراجعه کنید.

مراجع

- [1] Aitkin, M. and Alfo, M. (2000). Random coefficient models for binary longitudinal responses with attrition. *Statistics and Computing*. 10, 279-287.
- [2] Diggle, P. and Kenward, M.G. (1994). Informative drop-out in longitudinal data analysis (with discussion). *Appl. Statist.*, 43, 49-93.
- [3] Follmann D. and Wu M. (1995). An approximate generalized linear model with random effects for informative missing data. *Biometrics* 51: 151-168.
- [4] Heckman, J., (1981). Statistical models for discrete panel data, in Manski, C. & McFadden D., *Structural analysis of discrete data with econometric applications*. 114-195, Cambridge, Mass: MIT press.
- [5] Rubin D. B. (1976). Inference and missing data. *Biometrika* 63, 581-592.

[۶] رضایی قهرودی، ۱۳۸۲، تحلیل داده‌های طولی دودویی ناشی از متغیرهای پنهان با پاسخهای گمشده، دانشگاه شهید بهشتی.

تحلیل ممیزی با استفاده از آمیخته‌های نرمال

قاسم رکابدار^۱، رحیم چینی‌پرداز^۲

^۱ دانشگاه آزاد اسلامی واحد آبادان، خرمشهر

^۲ دانشگاه شهید چمران، گروه آمار

چکیده: تحلیل ممیزی خطی فیشر^۱ (LDA) ابزار مهمی برای رده‌بندی بین چندگروه می‌باشد. هنگامیکه کلاسها دارای توزیع نرمال چند متغیره با ماتریس کواریانس مشترک باشند طبق لم نیمن-پیرسن LDA بهینه بوده و معادل با قاعده رده‌بندی ماکزیمم درست‌نمایی است. LDA دارای معیابی است از جمله فرض نرمال بودن کلاسها اغلب بندرت پیش می‌آید و کرانه‌های تصمیم خطی، کلاسها را به قدر کافی از هم جدا نمی‌کنند. برای تعمیم LDA می‌توان فرض کرد که کلاسها از زیر کلاسهایی تشکیل شده‌اند که دارای توزیع نرمال با میانگینهای متفاوت و ماتریس کواریانس مشترک هستند. ممیزی در این حالت آمیخته^۲ (MDA) می‌باشد و کرانه‌های تصمیم بین کلاسها غیر خطی و پیچیده بوده و برآورد پارامترها تنها از طریق الگوریتم EM امکان‌پذیر است. کاهش ابعاد در ممیزی یکی از اهداف مهمی است که به وسیله ممیزی خطی فیشر امکان‌پذیر می‌گردد و بوسیله متغیرهای ممیزی فیشر می‌توان کلاسها را در نمودارهای دو بعدی مشاهده نمود. ممیزی آمیخته نیز این خاصیت را در بردارد و کاهش ابعاد داده‌ها در فضای برازش یافته توسط میانگین زیر کلاسها امکان‌پذیر است. در این مطالعه ممیزی بوسیله آمیخته‌های نرمال توسعه داده می‌شود و سپس داده‌های گزارش توسعه انسانی در سال ۲۰۰۲ را برای مقایسه روشهای ممیزی برای رده‌بندی کشورها به سه کلاس کشورها با توسعه انسانی بالا، کشورها با توسعه انسانی متوسط و کشورها توسعه انسانی پایین در نظر گرفته شدند و نشان داده شد که ممیزی آمیخته دارای نرخ خطای کمتری نسبت به دیگر روشها برای رده‌بندی کشورها است.

واژه‌های کلیدی: رده‌بندی، تحلیل خوشه‌ای، الگوریتم EM، تحلیل ممیزی انعطاف پذیر، شاخص توسعه انسانی

۱ مقدمه

رده‌بندی دویا چند کلاس یا ممیزی از موضوعات مهم در آمار چند متغیره است که دارای کاربردهای مختلف در حوزه‌های گوناگون است. در حالت کلی ممیزی اختصاص یک یا چند مشاهده x با

1) Linear Discriminant Analysis 2) Mixture Discriminant Analysis

کلاس نامعلوم به جوامعی معلوم G_1, G_2, \dots, G_J است. برای این کار فرض می‌شود مجموعه راهنما $\{x_i, g_i\}_{i=1}^n$ است که در آن $x_i \in \mathbb{R}^p$ موجود g_i مشخص کننده کلاس نمونه نام است. با استفاده از مجموعه راهنما پارامترها در مدل برآورد و سپس با استفاده از پارامترهای برآورد شده هدف بدست آوردن یک قاعده برای ممیزی بصورت یک تابع جهت پیش‌بینی مشاهده جدید x به یکی از کلاسها است.

روش سنتی و کلاسیک آماری ممیزی خطی و رگرسیون لجستیک است این روش‌ها تا حدود زیادی توسعه یافته و معمولاً در متون آماری موجود هستند. به عنوان مثال (ماردیا و دیگران ۱۹۷۹ و آندرسن ۱۹۸۴) برای ممیزی خطی و ممیزی لجستیک مراجع خوبی محسوب می‌شوند. هرگاه جوامع نرمال باشند تابع ممیزی خطی نیز نرمال و بنابراین خطای ممیزی نیز براحتی قابل محاسبه خواهد بود. روش درجه دوم^۳ در ممیزی بدلیل پیچیده بودن توزیع ممیزی بوسیله آماردانان هنوز در حال توسعه است. وب (۱۹۹۹) روش‌های ممیزی غیر خطی از جمله درجه دوم، ناپارامتری و . . . را مورد بررسی قرار داده است.

متأسفانه فرض نرمال بودن کلاسها بندرت پیش می‌آید. در چنین صورتی تابع ممیزی داده‌ها را بصورت موثری از هم جدا نمی‌کند و خطای ممیزی افزایش خواهد یافت. یک طریق برای حل این مساله فرض تعمیم LDA است به اینکه فرض شود کلاسها شامل زیرکلاسهای نرمال اند یعنی کلاسها دارای توزیع آمیخته‌اند و بنابراین ممیزی بکار رفته نیز آمیخته خواهد بود. این موضوع در نوشتجات آماری بوسیله مک لچلان (۱۹۹۲) و تکست و دیگران (۱۹۹۱) پیشنهاد شده است. هاستی و تیمشیرانی (۱۹۹۶) این مطلب را گسترش داده‌اند. آنها فرض کرده‌اند هرکلاس از زیر کلاسهایی با میانگین متفاوت و ماتریس کواریانس درون و بین کلاس‌های مشترک هستند. LDA فقط برای رده‌بندی استفاده نمی‌شود بلکه ابزار سودمندی برای کاهش ابعاد داده‌هاست با LDA می‌توان چندین کلاس از داده‌ها را در غالب نمودارهای دو بعدی نمایش داد بگونه‌ای که بیشترین جداسازی بین میانگین کلاسها را داشته باشند MDA نیز این خاصیت را دارد و می‌توان کاهش ابعاد را روی میانگین کلاسها (زیر کلاسها) مشاهده نمود. در LDA کاهش ابعاد در حالتی که ممیزی بین دو کلاس باشد امکان‌پذیر نیست اما در MDA این امر امکان‌پذیر است.

در این مقاله تحلیل ممیزی آمیخته بررسی و الگوریتم EM برای برآورد پارامترها و ممیزی پیشنهاد شده است. این الگوریتم بوسیله هاستی و تیمشیرانی پیشنهاد و بوسیله نویسندگان مقاله انطباق و سپس در داده‌های توسعه انسانی بکار گرفته شده است و در بخش دوم مقاله‌های ممیزی خطی مرور و بخش سوم مقاله توسعه LDA به MDA را شامل می‌شود. بخش چهارم به انطباق الگوریتم EM برای استفاده در ممیزی MDA می‌پردازد و در بخش پنجم داده‌های توسعه انسانی را برای مقایسه بین MDA با LDA و QDA در نظر گرفته شده‌اند سرانجام بخش‌های نهایی مقاله یک نتیجه‌گیری در خصوص MDA خواهد بود.

۲ تحلیل ممیزی خطی LDA

فرض کنید در مساله ممیزی مجموعه راهنما $\{x_i, g_i\}_{i=1}^n$ در اینصورت چگالی شرطی X در کلاس j ام به صورت زیر خواهد بود:

$$p_j(x) = \phi(\mu_j, \Sigma) = P(X = x | G = j) \\ = |\sqrt{\pi} \Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} D(x_i, \mu_j) \right\} \quad (1)$$

در اینصورت لگاریتم درستتمایی مشاهدات به صورت زیر خواهد بود:

$$\ell(\mu_j, \Sigma) = - \sum_{j=1}^J \sum_{g_i=j} D(x_i, \mu_j) - n \log |\Sigma| \quad (2)$$

در اینجا $D(x, \mu) = (x - \mu)' \Sigma^{-1} (x - \mu)$ فاصله مایلانوبیس بین x و μ و $\Sigma_{g_i=j}$ به معنی مجموع مشاهدات مجموعه راهنما در کلاس j ام است. گیریم احتمال پیشین کلاس j ام بصورت $P(G = j) = \Pi_j$ باشد که معمولاً از قبل معلوم است و یا از مجموعه راهنما برآورد می شود اگر هیچ دلیلی بر ترجیح کلاسها بر یکدیگر نباشد $\frac{1}{J} = \Pi_1 = \Pi_2 = \dots = \Pi_J$ در نظر گرفته می شود سپس مشاهده x_0 به کلاس j ام رده بندی می شود $C(x_0) = j$ اگر احتمال پسین اینکه از کلاس j ام باشد در میان بقیه کلاسها ماکزیمم شود:

$$P(G = j | X = x_0) = \max_{\ell} P(G = \ell | X = x_0) \quad (3)$$

هرگاه $P(G = j) = \frac{1}{J}$ باشد این روش با روش ممیزی ماکزیمم درستتمایی معادل می شود. با توجه به (۱) اگر کلاسها دارای توزیع نرمال با ماتریس کواریانس مشترک باشند چگالی پسین کلاس j ام به صورت زیر خواهد بود:

$$P(G = j | X = x) = \frac{\Pi_j \phi(\mu_j, \Sigma)}{\sum_{i=1}^J \Pi_i \phi(\mu_i, \Sigma)} \\ \propto \exp \left\{ x' \Sigma^{-1} \mu_j - \frac{1}{2} \mu_j' \Sigma^{-1} \mu_j + \log \Pi_j \right\} \\ \propto \exp \{ x' \beta_j + \alpha_j \}$$

توجه کنید ثابت بودن مخرج کسر نسبت به j تاثیری در ماکزیمم کردن کلاسها ندارد. کران تصمیم بین دو کلاس مانند i و j ام به صورت مجموعه ای از نقاط تعریف می شود که دارای احتمالات پسین مساوی هستند یعنی:

$$b_{i,j} = \{x \in \mathbb{R}^p; P(G = j | X = x) = P(G = i | X = x)\}$$

و بنابراین کران تصمیم بین دو کلاس بصورت زیر است:

$$(\beta_j - \beta_i)'x + (\alpha_j - \alpha_i) = 0$$

که تابعی خطی از x است. تابع ممیزی برای کلاس j ام به صورت زیر تعریف می‌شود:

$$\delta_j(x) = -(x'\beta_j + \alpha_j) \quad (4)$$

و قاعده رده‌بندی انتساب x به کلاس j ام است اگر رابطه زیر برقرار باشد:

$$\delta_j(x) = -\min_{\ell} (x'\beta_{\ell} + \alpha_{\ell})$$

باید توجه شود که در حالت مساوی بودن احتمال پیشین کلاسها مشاهده جدید x به کلاسی رده‌بندی می‌شود که کمترین فاصله ماکزیمم را با مرکز آن داشته باشد. α_{ℓ} و β_{ℓ} با استفاده از مجموعه راهنما برآورد می‌شوند بنابراین ماکزیمم (۲) بصورت زیر حاصل می‌شوند.

$$\hat{\mu}_j = \frac{1}{n_j} \sum_{g_i=j} x_i \quad \hat{\Sigma}_j = \frac{1}{n} \sum_{g_i=j} \sum_{g_i=j} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)' \quad (5)$$

لازم بذکر است اگر کلاسهای دارای ماتریس کواریانس مختلف باشند یعنی:

$$p_j(x) = \phi(\mu_j, \Sigma_j)$$

کران تصمیم بین دو کلاس i و j غیر خطی به صورت تابع درجه دوم از x خواهد بود در این حالت تعداد پارامترهایی که باید برآورد شوند $JP + \frac{P(P+1)}{2}$ خواهد بود چون اندازه نمونه‌های مجموعه راهنما ثابت است بنابراین برآورد پارامترها در ممیزی درجه دوم نیرومند نیست و برای حالت $P > 2$ پیشنهاد نمی‌شود.

۱.۲ کاهش ابعاد در LDA

کاهش ابعاد در ممیزی را می‌توان با تجزیه فیشر (۱۹۳۶) بدست آورد در روش پیشنهادی فیشر هدف بدست آوردن $\alpha \in \mathbb{R}^p$ بگونه‌ای که کسر زیر ماکزیمم شود:

$$\frac{\alpha' B \alpha}{\alpha' W \alpha}$$

در اینجا B ماتریس کواریانس بین گروهها و W ماتریس کواریانس درون گروهی است. بهترین جواب برای α بردار ویژه مطابق با بزرگترین مقدار ویژه ماتریس $W^{-1}B$ است در حالت کلی این ماتریس حداکثر دارای $\min(J-1, p)$ مقدار ویژه غیر صفر است بنابراین متغیرهای ممیزی

فیشر به صورت $y_k = \alpha'_k x$ و $k = 1, 2, \dots, K$ می‌باشند که $K \leq \min(J - 1, p)$ یعنی عملاً از K ممیز اول فیشر در LDA به عنوان متغیر ممیزی استفاده می‌شود دو متغیر ممیزی اول فیشر بیشترین جداسازی را بین میانگین کلاسها دارند و برای نمایش کلاسها در غالب نمودار استفاده می‌شوند. در ممیزی خطی فیشر مشاهده جدید x_0 به کلاس j ام رده‌بندی می‌شود اگر فاصله اقلیدسی زیر کمترین در میان همه کلاسها باشد:

$$\delta_j^2(x_0) = \sum_{k=1}^K (y_k(x_0) - \bar{y}_k^j)^2$$

که \bar{y}_k^j میانگین متغیر ممیزی k ام در کلاس j ام است. در حالت نرمال بودن کلاسها هاستی و تیبشیرانی (۱۹۹۶) نشان داده‌اند که رده‌بندی در فضای ممیزی خطی فیشر معادل با LDA است. همچنین هاستی و دیگران (۱۹۹۴) نشان دادند که فضای ممیزی خطی فیشر (LDA) معادل با رگرسیون خطی چند متغیره است که از مقادیر بهینه برای نشان دادن کلاسها استفاده شده است در حالت کلی LDA را می‌توان بصورت دنباله‌ای از رگرسیونهای چندگانه خطی $\eta_k(X) = X' \beta_k$ و $k = 1, \dots, K$ نمایش داد این شکل از ممیزی پس از تعیین مقادیر بهینه $\theta_1, \dots, \theta_K$ و انتساب آنها به کد کلاسها بدست آورده می‌شود. روش کار بصورت زیر است فرض شود که Y ماتریسی است که اگر x_i متعلق به کلاس j ام باشد آنگاه در ردیف j ام و ستون j ام مقدار آن یک است و بقیه ستونها در ردیف j ام مقدار صفر را دارند:

(۱) رگرسیون چند متغیره خطی: رگرسیون خطی چند متغیره از Y در مقابل X (ماتریس پیش بینی‌گرها) برازش می‌شود گیریم \hat{Y} ماتریس مقادیر برازش شده باشد و $\eta(X)$ بردار توابع رگرسیونی باشد.

(۲) مقادیر بهینه: K تا از بزرگترین بردارهای ویژه $Y'Y$ با شرط $\Theta' D_{\Pi} \Theta = I_K$ که $D_{\Pi} = \frac{Y'Y}{n}$ ماتریس قطری از احتمال پیشین کلاسهاست فرض شود Θ ماتریس مقادیر بهینه برآورد شده باشد.

(۳) با استفاده از ماتریس مقادیر بهینه در مرحله (۲) و بردار توابع رگرسیونی در مرحله (۱) $\eta(X) \leftarrow \Theta' \eta(X)$

در این حالت مشاهده جدید x_0 به کلاس j ام رده‌بندی می‌شود اگر فاصله اقلیدسی زیر در میان همه کلاسها کمترین باشد:

$$\delta_j^2(x_0) = \sum_{k=1}^K w_k (\eta_k(x_0) - \bar{\eta}_k^j(x))^2$$

که $\bar{\eta}_k^j(x)$ میانگین مقادیر برازش شده تابع رگرسیونی $\eta_k(x)$ در کلاس j ام است و $w_k = \frac{1}{r_k^2(1-r_k^2)}$ و r_k وزنهایی هستند که فضای برازش شده رگرسیونی را به فضای ممیزی فیشر تبدیل می‌کنند و k ام مقدار ویژه‌ای است که در مرحله (۲) الگوریتم محاسبه می‌شود. اگر در مرحله (۱) از الگوریتم

از رگرسیون ناپارامتری استفاده شود در این صورت ممیزی غیر خطی خواهد بود. هاستی و دیگران (۱۹۹۴) هر روشی را که از یک پیش‌پردازش رگرسیونی برای مسائل ممیزی استفاده شود را ممیزی انعطاف‌پذیر^۴ (FDA) نامیده‌اند بنابراین اگر پیش‌پردازش رگرسیون خطی باشد ممیزی نیز خطی خواهد بود.

۳ تعمیم LDA با استفاده از آمیخته‌های نرمال

فرض کنید در بخش قبلی تابع چگالی تابع چگالی در هر کلاس آمیخته نرمال باشد که در آن هر مولفه دارای میانگین مربوط به خود و ماتریس کواریانس مشترک با بقیه مولفه‌ها باشد بنابراین در کلاس j ام به صورت زیر خواهد بود:

$$p_j(x) = \sum_{r=1}^{R_j} \pi_{jr} \phi(x_i, \mu_{jr}, \Sigma)$$

$$= |\sqrt{\pi} \Sigma|^{-1} \exp \left\{ -\frac{1}{\sqrt{\pi}} D(\mu_{jr}, \Sigma) \right\} \quad (۶)$$

در اینجا R_j تعداد مولفه‌های کلاس j ام و π_{jr} نسبت‌های آمیخته‌گی هر یک از مولفه‌های نرمال است که $\sum_{r=1}^{R_j} \hat{\pi}_{jr} = 1$ می‌باشد. در مدل ممیزی آمیخته احتمال پسین کلاس j ام در صورتی که x مشاهده شده باشد به صورت زیر است:

$$P(G = j | X = x) = \frac{\sum_{r=1}^{R_j} \pi_{jr} \phi(x_i, \mu_{jr}, \Sigma)}{\sum_{j=1}^J \sum_{r=1}^{R_j} \pi_{jr} \phi(x_i, \mu_{jr}, \Sigma)} \quad (۷)$$

لگاریتم درست‌نمایی شرطی در اینجا بصورت زیر است:

$$\ell_m(\mu_{jr}, \Sigma, \pi_{jr}) \propto \sum_{j=1}^J \log \left(\sum_{r=1}^{R_j} \pi_{jr} \exp \left\{ -\frac{1}{\sqrt{\pi}} D(\mu_{jr}, \Sigma) \right\} \right) - \frac{n}{\sqrt{\pi}} \log |\Sigma| \quad (۸)$$

در MDA مشاهده جدید x_0 به کلاس j ام نسبت داده می‌شود، $C(x_0) = j$ ، اگر رابطه (۷) در میان تمام کلاس‌های دیگر ماکزیمم شود:

$$P(G = j | X = x_0) = \max_{\ell} P(G = \ell | X = x_0)$$

در اینجا پارامترهای مدل $R_j, \pi_{jr}, \mu_{jr}, \Sigma$ که $r = 1, \dots, R_j, j = 1, \dots, J$ می‌باشد این پارامترها از طریق مجموعه راهنما با ماکزیمم کردن رابطه (۸) بدست می‌آیند بدیهی است که در MDA ماکزیمم کردن رابطه (۸) از طریق مشتق‌گیری امکان‌پذیر نیست بنابراین باید از طریق روشهای عددی ماکزیمم شود.

۴ برآورد پارامتر مدل ممیزی آمیخته

دمپستر و دیگران (۱۹۷۷) برای برآورد پارامترها مطرح کرده‌اند که در هنگام ناکامل بودن مشاهدات کاربرد دارد. الگوریتم EM روشی تکراری تا رسیدن به همگرایی در لگاریتم درستنمایی داده‌هاست که در هر تکرار مرحله E (امید ریاضی) و مرحله M (برآورد ماکزیمم درستنمایی) وجود دارد. در مدل ممیزی آمیخته بطور دقیق مشخص نیست که هر نمونه در مجموعه راهنما به کدام زیرکلاس تعلق دارد بنابراین برای ماکزیمم کردن (۲) باید از الگوریتم عددی EM استفاده نمود.
مرحله E:

مرحله E از الگوریتم EM شامل امید ریاضی لگاریتم درستنمایی داده‌های کامل است به شرط اینکه مقدار مشاهده شده x و مقدار صحیح پارامترها در هر تکرار در اختیار باشد در داده‌های آمیخته نرمال می‌توان مشاهدات را به صورت داده‌های ناکامل در نظر گرفت که مقادیر گمشده در آن کد زیر کلاس هر مولفه می‌باشد در اینجا با فرض اینکه پارامترهای ابتدایی مشخص هستند احتمال یا وزن اینکه مشاهده x_i متعلق به زیر کلاس r ام از کلاس j ام باشد برای همه نمونه‌های مجموعه راهنما در کلاس j ام محاسبه می‌شود که برای اثبات می‌توان به مک لاجلان (۱۹۹۲) مراجعه کرد:

$$P(C_{jr}(x_i) | x_i, j) = \frac{\pi_{jr} \phi(x_i, \mu_{jr}, \Sigma)}{\sum_{l=1}^{R_j} \pi_{jl} \phi(x_i, \mu_{jl}, \Sigma)} \quad (9)$$

مرحله M:

با استفاده از وزنهای (احتمالات) برآورد شده در مرحله E رابطه (۹) برآوردهای ماکزیمم درستنمایی برای پارامترهای هر مولفه نرمال در هر کلاس محاسبه می‌شوند:

$$\hat{\pi}_{jr} \propto \sum_{g_i=j} P(C_{jr}(x_i) | x_i, j) \quad \sum_{r=1}^{R_j} \hat{\pi}_{jr} = 1 \quad (10)$$

$$\hat{\mu}_{jr} = \frac{\sum_{g_i=j} P(C_{jr}(x_i) | x_i, j) x_i}{\sum_{g_i=j} P(C_{jr}(x_i) | x_i, j)} \quad (11)$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^J \sum_{g_i=j}^{R_j} \sum_{r=1}^{R_j} P(C_{jr}(x_i) | x_i, j) (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)' \quad (12)$$

همانطور که از مقایسه رابطه (۵) با رابطه‌های (۱۱) و (۱۲) دیده می‌شود برآوردهای ماکزیمم درست‌نمایی مانند حالت نرمال کامل هستند با این تفاوت که وزنه‌های $P(C_{jr}(x_i) | x_i, j)$ جایگزین تابع نشانگر کلاسها شده‌اند. با توجه به اینکه وزنه‌های $P(C_{jr}(x_i) | x_i, j)$ خود تابعی از μ_{jr} و Σ هستند رابطه‌های (۱۰) تا (۱۲) در رابطه (۹) جایگزین شده تا وزنه‌های بعدی محاسبه شوند این مراحل آنقدر تکرار می‌شوند تا همگرایی در لگاریتم درست‌نمایی شرطی (۸) روی دهد یکی از خواص مناسب الگوریتم EM این است که در هر تکرار لگاریتم درست‌نمایی داده‌ها نسبت به تکرار قبلی افزایش می‌یابد بنابراین اگر در مرحله‌ای رابطه (۸) نسبت به تکرار قبل از آن تغییر کمی داشته باشد الگوریتم متوقف می‌شود و برآوردهای نهایی بدست آورده می‌شوند.

تکرارهای الگوریتم EM به اندازه خوشه‌ها R_j و برآوردهای ابتدایی برای پارامترهای μ_{jr} و Σ و وزنه‌های $P(C_{jr}(x_i) | x_i, j)$ دارند برای این منظور می‌توان از الگوریتم خوشه‌بندی میانگین K ام^۵ یک مقدار ثابت برای خوشه‌ها (تعداد زیر کلاسها) تعیین کرد و سپس از این الگوریتم برای برآورد مراکز زیر کلاسها μ_{jr} در هر کلاس استفاده شود سپس برای همه مشاهدات مجموعه راهنما در کلاس j ام $P(C_{jr}(x_i) | x_i, j)$ برابر با یک است اگر μ_{jr} نزدیکترین مرکز و در غیر اینصورت صفر است. پس از برآورد پارامترها مشاهده جدید x مطابق با رابطه (۷) به خوشه مورد نظر نسبت داده می‌شود. در MDA با مساوی قرار دادن احتمال پسین کلاسها به تابعی پیچیده و غیر خطی خواهیم رسید بنابراین در MDA حتی با فرض مساوی بودن ماتریس کواریانس کلاسها (زیر کلاسها) کرانه‌های تصمیم بین کلاسها غیر خطی است که این یکی از تعمیمهای مهمی است که هنگام فرض نرمال آمیخته برای کلاسها در LDA بدست می‌آید.

فرض مساوی بودن ماتریس کواریانس زیر کلاسها این امکان را بوجود می‌آورد تا کاهش ابعاد داده‌ها در MDA بدست آورده شود به عبارت دیگر می‌توان فرض نمود که با نسخه وزنی LDA روبرو هستیم که زیر کلاسها جایگزین کلاسها شده‌اند اگر $R = \sum_{j=1}^J R_j$ تعداد کل زیر کلاسها باشد که هر کدام از مشاهدات با وزن خود به هر یک از زیر کلاسها تعلق می‌گیرد در این حالت کلاس j ام شامل $n_j R_j$ مشاهده است بنابراین مجموعه راهنما در این حالت مجموعه‌ای وزنی و افزوده شده به تعداد $n' = \sum_{j=1}^J n_j R_j$ از مشاهدات است. دوباره از الگوریتم EM کمک گرفته می‌شود و پس از محاسبه وزنه‌ها در مرحله E در مرحله M مسئله LDA وزنی حل می‌شود نرم افزار مورد استفاده این امکان را می‌دهد که در مرحله M از مقادیر بهینه برای محاسبه متغیرهای ممیزی استفاده شود اگر $Z_{n \times R}$ ماتریس پاسخ ظاهری باشد که ردیفهای این ماتریس وزن مشاهدات است بگونه‌ای که اگر مشاهده i ام متعلق به کلاس j ام باشد Z_{ij} باشد Z_{ij} باشد آنگاه زیر کلاسهای مربوط به کلاس j ام با وزن این مشاهده پر می‌شوند و بقیه زیر کلاسها در ردیف i ام مقدار صفر را می‌گیرند باید توجه داشت که مجموع وزنه‌ها برابر n است. $(n = \sum_{j=1}^J \sum_{q_i=j} \sum_{r=1}^{R_j} P(C_{jr}(x_i) | x_i, j))$

5) Kmeans

جدول ۱: نرخ خطای رده‌بندی نادرست برای داده‌های توسعه انسانی

| روش | مجموعه راهنما | مجموعه آزمون |
|-----|---------------|--------------|
| LDA | ۰/۱۰۳ | ۰/۱۰۳ |
| QDA | ۰/۹۰۶ | ۰/۱۲۱ |
| MDA | ۰/۰۲۶ | ۰/۰۵۲ |

پس از محاسبه ماتریس ظاهری Z با استفاده از مراحل (۳-۱) گفته شده در بخش ۱.۲ مدل کاهش یافته MDA در هر تکرار محاسبه می‌شود. در LDA کاهش ابعاد داده‌ها توسط تعداد کلاسها محدود شده است بنابراین در حالتی که ممیزی بین دو کلاس باشد ($\bar{J} = 2$) کاهش ابعاد امکان‌پذیر نیست اما در MDA چون زیر کلاسها جایگزین کلاسها می‌شوند بنابراین همواره متغیرهای ممیزی قابل محاسبه بوده و کاهش ابعاد امکان‌پذیر است.

۵ مثال کاربردی: رده‌بندی کشورها

از سال ۱۹۹۰ سازمان ملل متحد هر ساله گزارشی تحت عنوان گزارش توسعه انسانی را منتشر می‌کند در این گزارش برای رده‌بندی کشورها از نظر توسعه انسانی آنها متغیرهای زیر در نظر گرفته می‌شود:

۱: امید به زندگی در بدو تولد.

۲: نرخ باسوادی بزرگسالان.

۳: نرخ رشد نام نویسی در دوره‌های تحصیلی.

۴: سرانه محصول ناخالص داخلی.

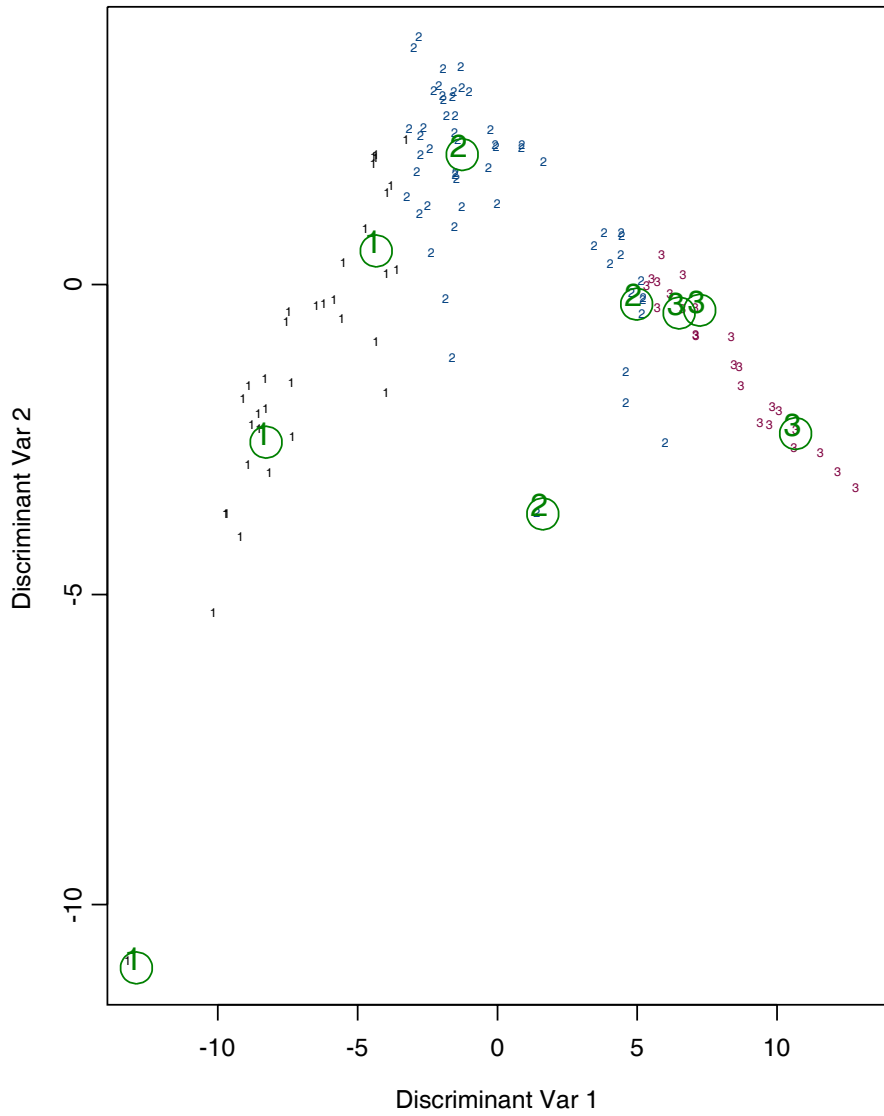
جمع‌آوری اطلاعات معمولاً با استفاده از اطلاعات رسمی منتشره کشورها و در صورت غیر قابل اعتماد بودن اطلاعات بوسیله برنامه عمران سازمان ملل برآورد می‌شوند. در این گزارش از داده‌های بدست آمده شاخص توسعه انسانی^۶ (HDI) محاسبه می‌شود و با استفاده از آن کشورها به سه رده کشورها یا توسعه انسانی بالا، کشورها با توسعه انسانی متوسط و کشورها با توسعه انسانی پایین رده‌بندی می‌شوند.

۱.۵ رده‌بندی کشورها با تحلیل ممیزی

در گزارش توسعه انسانی سال ۲۰۰۲ آمارها برای ۱۷۳ کشور گزارش شده است برای مقایسه روشهای ممیزی خطی LDA، درجه دوم QDA و ممیزی آمیخته MDA، ۱۱۵ کشور بطور تصادفی برای مجموعه راهنما انتخاب شدند و نرخ خطای رده‌بندی نادرست برای مجموعه راهنما

6) Human Development Index

Discriminant Plot for true classes



شکل ۱: متغیرهای ممیزی: متغیرهای ممیزی اول و دوم بهتر کشورها را از هم جدا کرده‌اند

محاسبه شد سپس برای بررسی اعتبار مدل‌های ارائه شده ۵۸ کشور دیگر به عنوان مجموعه آزمون در نظر گرفته شده‌اند و نرخ خطای رده‌بندی نادرست برای مشاهدات مجموعه آزمون با استفاده از توابع ممیزی برآورد شده توسط مجموعه راهنما محاسبه شده است همانطور که از جدول (۱) مشاهده می‌شود MDA با سه زیرکلاس در هر کلاس دارای نرخ خطای کمتری نسبت به روشهای ممیزی خطی و درجه دوم است بنابراین از ممیزی آمیخته برای رده‌بندی کشورها استفاده شده است.

شکل (۱) نمودار مربوط متغیرهای ممیزی برای ممیزی خطی و ممیزی آمیخته را برای مجموعه راهنما نشان می‌دهد همانطور که مشاهده می‌شود در MDA چون زیر کلاسها جایگزین کلاسها شده‌اند بنابراین تعداد متغیرهای ممیزی بیشتر از حالت خطی است اعداد درون دایره‌ها مشخص کننده مرکز زیر کلاسها می‌باشند و متغیرهای ممیزی اول و دوم بهتر از متغیرهای ممیزی دوم و سوم کلاسها (زیر کلاسها) را از هم جدا کرده‌اند.

کشورها با توسعه انسانی بالا:

آرژانتین، آنتیگوآ و باربودا، آلمان، اتریش، اروگوئه، اسپانیا، استرالیا، استونی، اسلواکی، اسلوانی، ایتالیا، امارات متحده عربی، انگلستان، ایالات متحده آمریکا، ایرلند، ایسلند، باربادوس، باهاماس، بحرین، برونئی، بلژیک، پرتغال، ترنیداد و توباگو، جمهوری چک، دانمارک، ژاپن، سنگاپور، سوئد، سوئیس، سیشل، فرانسه، رژیم اشغالگر قدس، فنلاند، قبرس، قطر، کاستاریکا، کانادا، کرواسی، کره جنوبی، کویت، لوکزامبورگ، لهستان، ماریتوس، مالت، مالزی، مجارستان، مکزیک، نروژ، نیوزلند، هلند، هنگ کنگ، یونان.

کشورها با توسعه انسانی متوسط:

آذربایجان، آفریقای جنوبی، آلبانی، اردن، ارمنستان، ازبکستان، اکراین، اکوادور، اکوتوریال، الجزایر، السالوادور، اندونزی، ایران، برزیل، بلغارستان، بلین، بولیوی، بوتسوانا، پاپوآ، پاراگوئه، پاناما، پرو، تاجیکستان، تایلند، ترکمنستان، ترکیه، تونس، جامائیکا، جزایر سلیمان، دمنیکا، جمهوری دمنیکن، چین، روسیه، روسیه سفید، رومانی، زیمبابوه، ساموآ، سائوتومو، سری لانکا، سنت اوسیا، سنت وینت، سوازیلند، سورینام، سوریه، عربستان، عمان، غنا، فیلیپین، فیجی، قرقیزستان، قزاقستان، کاپورد، کامبوج، کلمبیا، کنگو، کنیا، کوبا، گرجستان، گرانادا، گواتمالا، گوانا جدید، گویان، مالدیو، مراکش، مصر، مغولستان، مقدونیه، ملداوی، میانمار، لاتویا، لیتوانی، لبنان، لسوتو، لیبی، لیتونی، نامیبیا، نیکاراگوئه، ونزوئلا، ویتنام، هند، هندوراس.

کشورها با توسعه انسانی پایین:

آفریقای مرکزی، آنگولا، اتیوپی، اریتره، اوگاندا، بنگلادش، بنین، بوتان، بورکینافاسو، برونزی، پاکستان، تانزانیا، توگو، چاد، جمهوری مردمی لئو، جیبوتی، رواندا، زامبیا، سنگال، سودان، سیرالئون، کامرون، کنگو، کوته‌دوبری، کوموروس، گابن، گامبیا، گینه، گینه بیسائو، ماداگاسکار، مالاوی، مالی، موریتانی، موزامبیک، نپال، نیجر، نیجریه، وائوتوآ، هائیتی، یمن.

۶ نتیجه‌گیری

همانطور که گفته شد بدلیل غیر نرمال بودن کلاسها ممیزی متداول کلاسیک توان کافی در جداسازی کلاسها را ندارند. یک روش پارامتری را در ممیزی میتوان با فرض داشتن توزیع نرمال آمیخته در کلاسها بدست آورد با این فرض دو تعمیم مهم در ممیزی خطی بدست می‌آید:

- (۱) کرانه‌های تصمیم بین کلاسها غیر خطی و پیچیده است.
 - (۲) هنگامیکه ممیزی بین دو کلاس باشد کاهش ابعاد بر خلاف LDA امکان‌پذیر است.
- روش ممیزی آمیخته اگر چه منجر به محاسبات پیچیده و استفاده از روشهای عددی بجای تحلیلی می‌گردد ولی با توجه به دسترسی به رایانه‌ها و نرم‌افزارهای دقیق می‌توان از روش MDA استفاده کرد.

مراجع

- [1] Anderson, T. W.(1984), An Introduction to Multivariate Statistical Analysis, John Wiley, New York.
- [2] Dempster, A. P., Laird, N. M. and Rubin, D. B.(1977), Maximum likelihood from incomplete data via EM algorithm (with discussion), R. Statis. Sac. B, 39:1-38.
- [3] Hastie, T. Tibshirani, R. and Buja, A.(1994), Flexible Discriminant Analysis by Optimal Scoring, J. Amer. Statis. Assoc. 89, 1225-1270.
- [4] Hastie, T. Tibshirani, R. (1996), Discriminant Analysis by Gaussian Mixtures, J. R. Statis. Soc. B, 58:155-176.
- [5] Mardia, K. Kent, J. and Bibby, J.(1979), Multivariate Analysis, Academic Press, London.
- [6] McLachlan, G. J. (1992), Discriminant Analysis and Statistical Pattern Recognition, John Wiley, New York.
- [7] Taxt, T. Hjot, N. and Eikvil, L. (1991), Statistical classification using a linear mixture of multinormal probability densities, Pattern Recognition Letters, 12:731-737.
- [8] UNDP. (2002), Human Development Report 2002, New York.
- [9] Webb, A. (1999), Statistical Pattern Recognition, Arnold, London.

خطای مطلق یا نسبی در تعیین اندازه نمونه

علیرضا زاهدیان

مرکز آمار ایران

چکیده: یکی از مهمترین مسائلی که در شروع طراحی یک طرح نمونه‌گیری با آن روبرو هستیم این است که اندازه نمونه مورد نیاز باید چقدر باشد تا با اطمینان کافی انتظار داشته باشیم خطای نمونه‌گیری در برآورد پارامترهای مورد نظر در حد قابل قبولی باشد. در تعیین اندازه نمونه، عواملی مانند پراکندگی واحدهای جامعه آماری از نظر صفت مورد بررسی، اندازه جامعه، توزیع صفت مورد بررسی و حداکثر خطای نمونه‌گیری قابل قبول، نقش تعیین کننده‌ای دارند. در بین عوامل فوق، خطای نمونه‌گیری تنها عاملی است که کنترل آن در اختیار طراح آمارگیری قرار دارد و با تغییر آن می‌تواند حجم کار و هزینه طرح را تا جایی که به اهداف آن صدمه وارد نشود، کاهش دهد. در این مقاله روشهای مختلف تعیین حداکثر خطای نمونه‌گیری قابل پذیرش برای تعیین اندازه نمونه در طرحهای آمارگیری و اثرات آنها در فرمولهای تعیین اندازه نمونه مورد بحث قرار می‌گیرد.

واژه‌های کلیدی: اندازه نمونه، خطای نسبی نمونه‌گیری، خطای مطلق نمونه‌گیری

۱ مقدمه

یکی از عواملی که در دقت نتایج یک آمارگیری نمونه‌ای تاثیر مستقیم دارد، اندازه نمونه است. بگونه‌ای که با افزایش اندازه نمونه خطای نمونه‌گیری کاهش پیدا کرده و به تعبیری نتایج بدست آمده از اعتبار آماری بیشتری برخوردار خواهند بود. ولی در عمل با افزایش اندازه نمونه، هزینه آمارگیری و میزان خطای غیر نمونه‌گیری خصوصا در آمارگیریهای بزرگ افزایش پیدا می‌کنند. علاوه بر این در برخی از زمینه‌ها مانند تحقیقات پزشکی که واحدهای آماری انسانها هستند، برای افزایش اندازه نمونه محدودیتهای بسیاری وجود دارد. بنابراین آمارشناس باید به کمک محقق طرح اندازه نمونه را بگونه‌ای تعیین کند که بین دقت آمارگیری و هزینه آن تعادل لازم برقرار شود. مسئله تعیین اندازه نمونه از جنبه‌های مختلف در کتب و مقالات آماری مورد بحث و بررسی قرار گرفته است. برای مثال می‌توان به مراجعی مانند دمینگ^۱ (۱۹۶۶)، کوکران^۲ (۱۹۷۷)، کرامر و تیمن^۳ (۱۹۸۷)، کوهن^۴ (۱۹۸۸ و ۱۹۹۷)، شیفر^۵ و همکاران (۱۹۹۰)، کیش^۶ (۱۹۹۵)، لنت^۷ (۲۰۰۱) و تامپسون^۸ (۲۰۰۲) اشاره کرد. در این مراجع عواملی مانند تعیین یا برآوردهای تخمین

1) Deming 2) Cochran 3) Kramer and Thiemann 4) Cohen 5) Schaeffer
6) Kish 7) Lenth 8) Thompson

میزان پراکندگی صفت مورد بررسی، توزیع آماری صفت مورد بررسی، توان آزمونهای آماری مربوط به میانگین، مجموع و نسبت، در ارتباط با اندازه نمونه، مورد تجزیه و تحلیل قرار گرفته‌اند و در مورد مسئله تعیین حداکثر خطای نمونه‌گیری قابل پذیرش و رابطه آن با اندازه نمونه نیز اشاراتی شده است ولی این مسئله به طور خاص مورد بررسی قرار نگرفته است. این در حالی است که حداکثر خطای نمونه‌گیری قابل پذیرش (یا به اختصار «خطا») در بین عوامل فوق تنها عاملی است که تعیین آن در اختیار طراح طرح بوده و با تغییر آن در حدی که به کیفیت و اهداف طرح صدمه‌ای وارد نشود، می‌توان از اندازه نمونه و به تبع آن از هزینه طرح کاست. در این مقاله سعی شده است فرمولهای تعیین اندازه نمونه مورد نیاز برای پارامترهایی مانند میانگین، مجموع و نسبت، با تمرکز بر نقش خطای مذکور در تعیین اندازه نمونه مورد بررسی قرار گیرد.

تعبیر ریاضی حداکثر خطای نمونه‌گیری قابل پذیرش را می‌توان در قالب رابطه زیر بیان کرد:

$$Pr(|\hat{\Theta} - \Theta| \leq d) = 1 - \alpha \quad (۱)$$

که بر اساس آن با احتمال $1 - \alpha$ می‌توانیم انتظار داشته باشیم که خطای نمونه‌گیری در برآورد پارامتر مورد نظر (Θ) حداکثر برابر با d یا همان حداکثر خطای نمونه‌گیری قابل پذیرش باشد. رابطه بالا خطای نمونه‌گیری به صورت مطلق و برابر با تفاوت بین مقدار واقعی پارامتر و برآورد آن بیان شده است. حال اگر بخواهیم خطای نمونه‌گیری را به صورت نسبی در نظر بگیریم، رابطه (۱) به صورت زیر تبدیل می‌شود:

$$Pr\left(\left|\frac{\hat{\Theta} - \Theta}{\Theta}\right| \leq r\right) = 1 - \alpha$$

و یا به عبارت دیگر

$$Pr\left(\left|\hat{\Theta} - \Theta\right| \leq r\Theta\right) = 1 - \alpha \quad (۲)$$

که در آن r حداکثر خطای نسبی قابل پذیرش است. با مقایسه این رابطه با رابطه (۱) به این نتیجه می‌رسیم که اندازه نمونه با حداکثر خطای مطلق d معادل با اندازه نمونه با حداکثر خطای مطلق $r\Theta$ است. بنابراین در تعیین اندازه نمونه، خطای نمونه‌گیری را به دو صورت می‌توان بیان کرد:

۱- خطای مطلق ۲- خطای نسبی.

استفاده از خطای نسبی در بین آمارشناسان از عمومیت بیشتری برخوردار است ولی باید توجه داشت حداکثر خطای قابل پذیرش، توسط محقق طرح تعیین می‌شود و باید به گونه‌ای بیان شود که برای وی به راحتی قابل درک باشد. به عبارت دیگر در تعیین حداکثر خطای قابل پذیرش، شیوه بیان خطا نیز از اهمیت خاصی برخوردار است.

همانطور که اشاره شد، حداکثر خطای قابل پذیرش باید توسط محقق طرح تعیین شود ولی معمولاً

وی نمی‌تواند تعابیر آماری نظرات خود را بطور دقیق بیان نماید. از اینرو باید آمارشناس با توجه به اهداف طرح سوالاتی را از محقق طرح بپرسد تا بتواند حداکثر خطای قابل پذیرش را با تعابیر آماری بیان کند. برخی از این سوالات را می‌توان به صورت زیر بیان نمود:

- چقدر خطا در برآورد پارامتر مورد بررسی قابل چشم‌پوشی است؟
 - چه میزان تغییر در پارامتر مورد بررسی را می‌توان نادیده گرفت؟
 - پارامتر مورد بررسی را در چه فاصله‌ای تخمین می‌زنید؟
 - آیا از نظر شما تفاوت بین $\Theta \pm d$ و Θ معنی‌دار است؟
 - آیا از نظر شما $r\Theta \pm \Theta$ با Θ تفاوت معنی‌داری دارد؟ (منظور مقادیر عددی r و Θ است که توسط آمارشناس پیشنهاد می‌شود)
 - مقدار پارامتر مورد بررسی را تا چند رقم می‌توان به بالا یا پایین گرد کرد؟
- برای مثالهای بیشتری در این زمینه می‌توانید مقاله لنت (۱۰۲۰) را ببینید. علاوه بر این در تعیین حداکثر خطای قابل پذیرش باید به نکات زیر توجه داشت:
- حداکثر خطای قابل پذیرش نباید از خطای گرد کردن قابل پذیرش بیشتر باشد.
 - در بعضی موارد بیان خطا به صورت نسبی منطقی نیست. برای مثال، در برآورد متوسط دمای بدن بیماران مورد آزمایش نمی‌توان خطا را به صورت درصدی از دمای بدن بیان کرد یا در برآورد متوسط قطر قطعات تراش داده شده توسط یک دستگاه نمی‌توان خطا را به صورت درصدی از متوسط قطر قطعات تعیین نمود.
 - در مواردی که تعاریف یا استانداردهای خاصی برای اندازه‌گیری صفت مورد بررسی وجود دارد باید حداکثر خطا را بر اساس آنها بیان کرد. برای مثال، حداکثر خطا در مورد برآورد دمای بدن انسان باید ۰/۵ درجه باشد.
 - استفاده از خطای نسبی مستلزم تخمینی از پارامتر مورد بررسی است و در مواردی که این پارامتر دارای نوسانات زیادی باشد استفاده از خطای نسبی توصیه نمی‌شود. بنابراین معمولاً از خطای نسبی در مواردی که آمارگیری دارای سابقه تکرار (در سال گذشته) است، استفاده می‌شود.
 - در مواردی که پارامتر مورد بررسی طی زمان دارای روند صعودی است (برای مثال، مجموع درآمد سالیانه خانوارهای کشور هر سال افزایش پیدا می‌کند)، استفاده از خطای نسبی محتاطانه است و در مواردی که پارامتر مورد بررسی روند نزولی دارد (برای مثال، نسبت جمعیت روستایی به کل جمعیت کشور) غیر محتاطانه است.
 - در مواردی که جامعه آماری به تعدادی زیرجامعه (مانند استان، مرد و زن، شهری و روستایی و . . .) شکسته شده و اندازه نمونه برای هر زیر جامعه به صورت جداگانه تعیین می‌شود، استفاده از خطای نسبی موجب می‌شود علاوه بر عوامل دیگر، عامل بزرگی زیر جامعه از

نظر صفت مورد بررسی نیز در توزیع اندازه نمونه بین زیر جامعه‌ها دخالت داده شود. به عبارت دیگر در این حالت از زیرجامعه بزرگتر نمونه بیشتری انتخاب می‌شود. حال رابطه (۱) را در نظر بگیرید. اگر $\hat{\Theta}$ برآورد نا اریبی از Θ بوده و دارای توزیع نرمال باشد، دارای توزیع نرمال استاندارد است. بنابراین داریم:

$$Pr \left(\frac{|\hat{\Theta} - \Theta|}{\sqrt{V(\hat{\Theta})}} \leq Z_{(1-\frac{\alpha}{2})} \right) = Pr \left(|\hat{\Theta} - \Theta| \leq Z_{(1-\frac{\alpha}{2})} \sqrt{V(\hat{\Theta})} \right) = 1 - \alpha$$

که در آن $Z_{(1-\frac{\alpha}{2})}$ صدک $(1 - \frac{\alpha}{2})^{100}$ ام توزیع نرمال استاندارد است. با توجه به اینکه واریانس $\hat{\Theta}$ با افزایش n کاهش پیدا می‌کند، اندازه نمونه (n) را بگونه‌ای باید انتخاب کرد که رابطه $Z_{(1-\frac{\alpha}{2})} \sqrt{V(\hat{\Theta})} \leq d$ برقرار باشد. در قسمتهای بعد چگونگی تعیین n با توجه به رابطه فوق، برای پارامترهایی مانند میانگین، مجموع و نسبت در روشهای نمونه‌گیری که از عمومیت بیشتری برخوردار هستند، مورد بحث قرار خواهد گرفت.

۲ برآورد میانگین با نمونه‌گیری تصادفی ساده

برآورد هر پارامتری معمولاً با دو هدف انجام می‌گیرد:
 ۱- ساخت فاصله اطمینان ۲- آزمون فرض. در قسمتهای بعد روش تعیین اندازه نمونه برای هر یک از این دو هدف، مورد بررسی قرار می‌گیرد.

۱.۲ برآورد میانگین با هدف ساخت فاصله اطمینان (برآورد فاصله‌ای)

در این حالت n را طوری انتخاب می‌کنیم که رابطه زیر برقرار باشد:

$$Pr \left(|\bar{x} - \mu| \leq Z_{(1-\frac{\alpha}{2})} \sqrt{V(\bar{x})} \right) = 1 - \alpha$$

با توجه به اینکه $V(\bar{x}) = \frac{N-n}{N} \frac{\sigma^2}{n}$ می‌توانیم بنویسیم:

$$d = Z_{(1-\frac{\alpha}{2})} \sqrt{\frac{N-n}{N} \frac{\sigma^2}{n}}$$

و بر این اساس اندازه نمونه مورد نیاز از رابطه زیر بدست می‌آید:

$$n = \frac{n_0}{1 + \frac{n_0}{N}}$$

که

$$n_o = \frac{Z_{(1-\frac{\alpha}{4})}^2 \sigma^2}{d^2} \quad (3)$$

در روابط بالا اگر از ضریب تصحیح جامعه (fpc) که برابر $\frac{N}{N-n}$ است، صرف نظر کنیم، n_o همان اندازه نمونه مورد نیاز خواهد بود. $Z_{(1-\frac{\alpha}{4})}$ تعیین کننده اطمینان مورد نظر است و معمولاً به ازای $\alpha = 0.05$ برابر با ۱.۹۶ یا تقریباً ۲ در نظر گرفته می‌شود ولی اگر سطح اطمینان بالایی مدنظر باشد، (مثلاً در آزمایشات پزشکی، پذیرش محموله‌های دارو یا کنترل تولید قطعات خاص در صنعت) α برابر ۰.۰۱ و در نتیجه $Z_{(1-\frac{\alpha}{4})} = 2.58$ یا تقریباً ۳ در نظر گرفته می‌شود. روشهای تعیین σ^2 عبارتند از:

۱- آمارگیریهایی قبلی ۲- آمارگیری اولیه ۳- تقسیم آمارگیری به دو مرحله، یکی برای تخمین σ^2 و دیگری برای برآورد میانگین که اطلاعات مرحله اول را نیز شامل می‌شود. در حالتی که هیچ اطلاعی در مورد σ^2 در دست نباشد، با فرض نرمال بودن توزیع صفت مورد بررسی، می‌توان σ^2 را بر اساس دامنه تغییرات تخمین زد:

$$\sigma^2 = \frac{range}{6}$$

که تخمین محافظه‌کارتر $\sigma^2 = \frac{range}{4}$ است. زیرا ۹۵ درصد از مقادیر یک جامعه نرمال در فاصله ۲ انحراف معیار و ۹۹.۷ درصد از مقادیر آن در فاصله ۳ انحراف معیار از میانگین قرار دارند. شيفر و همکارانش (۱۹۹۰) و لور^۹ (۱۹۹۹) مثالهایی را در مورد استفاده از این روش ارائه کرده‌اند.

حال به سراغ d می‌رویم همانطور که قبلاً اشاره شد، می‌توان d را به صورت مطلق یا نسبی تعیین کرد. اگر d را به صورت نسبی تعیین کنیم با صرف نظر از ضریب جامعه محدود (fpc) داریم:

$$n = \frac{Z_{(1-\frac{\alpha}{4})}^2 \sigma^2}{r^2 \mu^2} = \frac{Z_{(1-\frac{\alpha}{4})}^2}{r^2} CV$$

که در این حالت به جای واریانس (σ^2) از ضریب تغییرات (CV) برای تعیین اندازه نمونه استفاده می‌شود. با توجه به اینکه ضریب تغییرات یک متغیر در طول زمان نسبت به واریانس آن از پایایی بیشتری برخوردار است، استفاده از خطای نسبی در آمارگیریهایی که در دوره‌های زمانی تکرار می‌شوند، مناسب‌تر است. شيفر و همکارانش (۱۹۹۰) از رابطه زیر برای تعیین اندازه نمونه استفاده کرده‌اند:

$$n = \frac{CV}{CV_o(\bar{x})}$$

که در آن $CV_o(\bar{x})$ ضریب تغییرات مطلوب در برآورد میانگین است.

9) Lohr

۲.۲ برآورد میانگین با هدف آزمون فرض

فرض کنید برآورد میانگین با هدف انجام آزمون فرض $H_0: \mu = \mu_0$ در مقابل $H_1: \mu \neq \mu_0$ انجام شود. در این حالت ناحیه رد فرض H_0 عبارتست از:

$$R = \{\bar{x} > \mu_0 + Z_{(1-\frac{\alpha}{2})}\sigma_{\bar{x}}\} \cup \{\bar{x} < \mu_0 - Z_{(1-\frac{\alpha}{2})}\sigma_{\bar{x}}\}$$

و ناحیه قبول عبارتست از:

$$A = \{|\bar{x} - \mu_0| < Z_{(1-\frac{\alpha}{2})}\sigma_{\bar{x}}\}$$

در این صورت توان آزمون (احتمال رد یک فرض صفر نادرست) از رابطه زیر بدست می‌آید:

$$\begin{aligned} Power &= 1 - \beta = Pr(\bar{x} > \mu_0 + Z_{(1-\frac{\alpha}{2})}\sigma_{\bar{x}}) + Pr(\bar{x} < \mu_0 - Z_{(1-\frac{\alpha}{2})}\sigma_{\bar{x}}) \\ &= Pr\left(\frac{\bar{x} - \mu}{\sigma_{\bar{x}}} > \frac{\mu_0 - \mu}{\sigma_{\bar{x}}} + Z_{(1-\frac{\alpha}{2})}\right) + Pr\left(\frac{\bar{x} - \mu}{\sigma_{\bar{x}}} < \frac{\mu_0 - \mu}{\sigma_{\bar{x}}} - Z_{(1-\frac{\alpha}{2})}\right) \\ &= \Phi\left(-\sqrt{n}\frac{\mu - \mu_0}{\sigma} - Z_{(1-\frac{\alpha}{2})}\right) + \Phi\left(\sqrt{n}\frac{\mu - \mu_0}{\sigma} - Z_{(1-\frac{\alpha}{2})}\right) \end{aligned}$$

که در آن، α مقدار خطای نوع اول قابل پذیرش (معمولاً برابر با 0.05) و $\Phi(\cdot)$ تابع احتمال توزیع نرمال استاندارد است. اگر $\sqrt{n}\frac{\mu - \mu_0}{\sigma} \geq 1$ باشد، می‌توانیم $\Phi\left(-\sqrt{n}\frac{\mu - \mu_0}{\sigma} - Z_{(1-\frac{\alpha}{2})}\right)$ را تقریباً مساوی صفر در نظر بگیریم. در این صورت داریم:

$$Power = \Phi\left(\sqrt{n}\frac{\mu - \mu_0}{\sigma} - Z_{(1-\frac{\alpha}{2})}\right)$$

و اندازه نمونه مورد نیاز از رابطه زیر بدست می‌آید:

$$n = (Z_{(1-\frac{\alpha}{2})} + Z_{(1-\beta)})^2 \frac{\sigma^2}{(\mu - \mu_0)^2}$$

معمولاً $\alpha = 0.05$ و $\beta = 0.2$ در نظر گرفته می‌شود. فرمول توان آزمون نشان می‌دهد با افزایش n توان آزمون در سطح معنی‌داری α افزایش پیدا می‌کند. معمولاً فرضهای H_0 و H_1 به گونه‌ای تعیین می‌شوند که اهمیت α بیشتر از β باشد. به عبارت دیگر در بسیاری از موارد محقق فقط α را تحت کنترل قرار می‌دهد و نیاز چندانی به کنترل β احساس نمی‌کند و در نتیجه فرضهای H_0 و H_1 را نیز بر همین اساس تعریف می‌کند. در این موارد خصوصاً در نمونه‌گیری از جوامع بزرگ می‌توان از $Z_{1-\beta}$ در فرمول اندازه نمونه چشم‌پوشی کرد که در این صورت داریم:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 \sigma^2}{(\mu - \mu_0)^2} \quad (4)$$

از مقایسه روابط (۳) و (۴) می‌توانیم نتیجه بگیریم در صورتی که کنترل β مد نظر نباشد، در فرمول تعیین اندازه نمونه، رابطه $d \leq \mu - \mu_0$ باید برقرار باشد. در این صورت اگر μ_0 پس از انجام آمارگیری خارج از فاصله اطمینان قرار گیرد، فرض H_0 در سطح معنی‌داری α رد می‌شود. برای مثال فرض کنید بخواهیم آزمون کنیم افزایش متوسط درآمد قشر خاصی از جامعه نسبت به سال گذشته بیشتر از نرخ تورم (برای مثال ۲۰ درصد) بوده است یا خیر. در این حالت باید $d \leq 0.2\mu'$ باشد که μ' متوسط درآمد در سال گذشته است.

۳ برآورد مجموع با نمونه‌گیری تصادفی ساده

در این حالت $\hat{X} = N\bar{x}$ که در آن N تعداد واحدهای جامعه است، برآوردی ناریب از مجموع صفت مورد بررسی می‌باشد. بنابراین مشابه قسمت قبل، اندازه نمونه مورد نیاز برای برآورد مجموع بر اساس خطای مطلق از رابطه زیر بدست می‌آید:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 N^2 \sigma^2}{d^2} \quad (5)$$

حال اگر $d = r \times X$ در نظر گرفته شود، اندازه نمونه بر اساس خطای نسبی عبارتست از:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 N^2 \sigma^2}{r^2 X^2} = \frac{Z_{(1-\frac{\alpha}{2})}^2 N^2 \sigma^2}{r^2 N^2 \mu^2} = \frac{Z_{(1-\frac{\alpha}{2})}^2 \sigma^2}{r^2 \mu^2}$$

بنابراین فرمول تعیین اندازه نمونه برای برآورد مجموع بر اساس خطای نسبی با فرمول تعیین اندازه نمونه برای برآورد میانگین یکسان است. حال فرض کنید آزمون فرض $H_0: X = X_0$ در مقابل $H_1: X \neq X_0$ مد نظر باشد. در این حالت، اندازه نمونه مورد نیاز از رابطه زیر بدست می‌آید:

$$n = (Z_{(1-\frac{\alpha}{2})} + Z_{(1-\beta)})^2 \frac{N^2 \sigma^2}{(X - X_0)^2}$$

۴ برآورد نسبت با نمونه‌گیری تصادفی ساده

می‌دانیم \hat{P} که برابر با تعداد واحدهای دارای صفت مورد نظر در نمونه تقسیم بر تعداد واحدهای نمونه است، برآوردگری ناریب برای نسبت در جامعه (P) است و واریانس آن از رابطه

برآورد P بر اساس خطای مطلق را می‌توان از رابطه زیر بدست آورد:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 P(1-P)}{d^2} \quad (۶)$$

در رابطه بالا P با یکی از روشهای زیر تعیین می‌شود:

- استفاده از اطلاعات آمارگیریهای قبلی
 - استفاده از آمارگیری اولیه
 - تقسیم آمارگیری به دو مرحله، یکی برای برآورد اولیه P و دیگری برای برآورد نهایی P
 - مساوی گرفتن P با $\frac{1}{2}$ که موجب ماکزیمم شدن n به ازای Z و d ثابت می‌شود.
- اگر به جای d از خطای نسبی استفاده کنیم فرمول تعیین اندازه نمونه به صورت زیر تبدیل می‌شود:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 P(1-P)}{r^2 P^2} = \frac{Z_{(1-\frac{\alpha}{2})}^2 (1-P)}{r^2 P} \quad (۷)$$

موضوع مهمی که باید به آن توجه داشت این است که اندازه نمونه در رابطه (۶) به ازای $P = \frac{1}{2}$ ماکزیمم می‌شود ولی اندازه نمونه در رابطه (۷) با کاهش P بزرگتر می‌شود. برای مثال نرخ بیکاری را به عنوان یک نسبت در نظر بگیرید. اگر اطلاعات قبلی حاکی از این باشد که نرخ بیکاری در یک کشور بین 10° تا 15° درصد در نوسان است، بر اساس رابطه (۶)، n به ازای $P = 0.1$ و بر اساس رابطه (۷) به ازای $P = 0.15$ ماکزیمم خواهد بود. بنابراین در تعیین اندازه نمونه برای برآورد P استفاده از خطای نسبی در مقایسه با خطای مطلق، موجب بروز تغییرات مهمی می‌شود و به این نتیجه می‌رسیم که استفاده از خطای نسبی در مواردی قابل توصیه است که واحدهای دارای صفت مورد بررسی، نادر باشند. کوکران (۱۹۷۷) واحدهایی را که نسبت آنها کمتر از 10° درصد باشد، نادر معرفی کرده است. کوکران (۱۹۷۷) و زاکولا^{۱۰} (۱۹۹۹) استفاده از خطای نسبی را برای تعیین اندازه نمونه مورد نیاز برای برآورد تعداد $(N \times P)$ مناسب دانسته‌اند ولی دلیل خاصی برای آن ذکر نکرده‌اند. شاید دلیل آن این باشد که اندازه نمونه مورد نیاز برای برآورد

$N \times P$ بر اساس خطای مطلق و خطای نسبی به ترتیب از دو رابطه $n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 N^2 P(1-P)}{d^2}$ و $n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 (1-P)}{r^2 P}$ بدست می‌آیند و در فرمول دوم نیازی به N نیست و اندازه نمونه برای برآورد P و $N \times P$ یکسان است.

در تعیین حداکثر خطای قابل پذیرش در برآورد P باید توجه داشت که معمولاً نسبتها به صورت درصد و تا دو رقم اعشار گزارش می‌شوند (در اینگونه موارد حداکثر خطای مطلق نباید بیشتر

10) Zakkula

از ۰/۰۰۵ یا همان خطای گرد کردن باشد). بنابراین در تعیین اندازه نمونه برای برآورد نسبت راحتتر می‌توان حداکثر خطای مطلق قابل پذیرش را تعیین نمود و استفاده از خطای نسبی فقط در مواردی که واحدهای دارای صفت مورد نظر، نادر هستند، قابل توجیه است. بدیهی است در اینگونه موارد نیز برای آنکه حداکثر خطای مطلق کمتر از خطای گرد کردن (معمولاً ۰/۰۰۵) باشد، باید حداکثر خطای نسبی را ۰/۰۵ در نظر بگیریم.

۵ نمونه‌گیری با طبقه‌بندی و با انتساب متناسب

در این حالت کوکران (۱۹۷۷) و شیفر (۱۹۹۰) از رابطه زیر برای تعیین اندازه نمونه مورد نیاز برای برآورد میانگین استفاده کرده‌اند:

$$n_0 = \frac{1}{V_0} \sum_{h=1}^H \frac{N_h}{N} \sigma_h^2$$

که در آن N تعداد واحدهای جامعه، N_h تعداد واحدهای جامعه در طبقه h ام، σ_h^2 واریانس صفت مورد بررسی در طبقه h ام، H تعداد طبقه و V_0 واریانس مطلوب در برآورد میانگین است که می‌توان آن را از رابطه $V_0 = \frac{d^2}{Z_{(1-\frac{\alpha}{2})}^2}$ بدست آورد. V_0 را می‌توان بر اساس ضریب تغییرات مطلوب در برآورد میانگین (CV_0) نیز بیان کرد که در این صورت $V_0 = CV_0^2 \mu^2$ خواهد بود. $\mu Z_{(1-\frac{\alpha}{2})}$ و حداکثر خطای نسبی برابر $0/1 Z_{(1-\frac{\alpha}{2})}$ خواهد بود. با این توضیحات می‌توان فرمول اندازه نمونه بر اساس خطای مطلق را به صورت زیر بیان کرد:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 \sum N_h \sigma_h^2}{d^2 N}$$

در این فرمول d از همان تعابیر مربوط به نمونه‌گیری تصادفی ساده برخوردار است. همچنین اندازه نمونه بر اساس خطای نسبی از رابطه زیر بدست می‌آید:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 \sum N_h \sigma_h^2}{r^2 \mu^2 N}$$

برای برآورد مجموع نیز اندازه نمونه مورد نیاز بر اساس خطای مطلق و خطای نسبی از روابط زیر بدست می‌آید:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{2})}^2 N \sum N_h \sigma_h^2}{d^2}$$

$$n_o = \frac{Z_{(1-\frac{\alpha}{r})}^2 N \sum N_h \sigma_h^2}{r^2 X^2} = \frac{Z_{(1-\frac{\alpha}{r})}^2 N \sum N_h \sigma_h^2}{r^2 N \mu^2}$$

و بطور مشابه برای برآورد P اندازه نمونه بر اساس خطای مطلق و خطای نسبی از روابط زیر بدست می‌آید:

$$n_o = \frac{Z_{(1-\frac{\alpha}{r})}^2 \sum N_h P_h (1 - P_h)}{d^2 N}$$

$$n_o = \frac{Z_{(1-\frac{\alpha}{r})}^2 \sum N_h P_h (1 - P_h)}{r^2 P^2 N}$$

بنابراین تغییری که به دلیل استفاده از خطای نسبی در فرمول تعیین اندازه نمونه در نمونه‌گیری تصادفی ساده بوجود می‌آید، در اینجا مشاهده نمی‌شود.

۶ نمونه‌گیری با طبقه‌بندی و با انتساب متناسب

در این حالت فرمولهای زیر توسط کوکران (۱۹۷۷) و شیفر (۱۹۹۰) پیشنهاد شده است:

$$n = \frac{(\sum N_h \sigma_h)^2}{V_o N^2 + \sum N_h \sigma_h^2} \quad n_o = \frac{(\sum N_h \sigma_h)^2}{V_o N^2}$$

بر این اساس می‌توان اندازه نمونه مبتنی بر خطای مطلق را از روابط زیر بدست آورد:

$$n = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{d^2 N^2 + Z_{(1-\frac{\alpha}{r})}^2 \sum N_h \sigma_h^2} \quad n_o = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{d^2 N^2}$$

و برای اندازه نمونه مبتنی بر خطای نسبی داریم:

$$n = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{r^2 N^2 \mu^2 + Z_{(1-\frac{\alpha}{r})}^2 \sum N_h \sigma_h^2} \quad n_o = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{r^2 N^2 \mu^2}$$

برای برآورد مجموع، اندازه نمونه مورد نیاز بر اساس خطای مطلق از روابط زیر بدست می‌آید:

$$n = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{d^2 + Z_{(1-\frac{\alpha}{r})}^2 \sum N_h \sigma_h^2} \quad n_o = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{d^2}$$

این اندازه نمونه بر اساس خطای نسبی عبارتست از:

$$n = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{r^2 X^2 + Z_{(1-\frac{\alpha}{r})}^2 \sum N_h \sigma_h^2} \quad n_0 = \frac{Z_{(1-\frac{\alpha}{r})}^2 (\sum N_h \sigma_h)^2}{r^2 X^2}$$

که اگر در روابط بالا X را برابر $N\mu$ بگیریم، اندازه نمونه بر اساس خطای نسبی برای برآورد مجموع برابر با اندازه نمونه بر اساس خطای نسبی برای برآورد میانگین خواهد بود. به همین ترتیب اندازه نمونه مورد نیاز بر اساس خطای مطلق برای برآورد P از روابط زیر بدست می‌آید:

$$n = \frac{Z_{(1-\frac{\alpha}{r})}^2 \left(\sum N_h \sqrt{P_h(1-P_h)} \right)^2}{d^2 N^2 + Z_{(1-\frac{\alpha}{r})}^2 \sum N_h P_h(1-P_h)}$$

$$n_0 = \frac{Z_{(1-\frac{\alpha}{r})}^2 \left(\sum N_h \sqrt{P_h(1-P_h)} \right)^2}{d^2 N^2}$$

و اندازه نمونه مورد نیاز بر اساس خطای نسبی برای برآورد P نیز عبارتست از:

$$n = \frac{Z_{(1-\frac{\alpha}{r})}^2 \left(\sum N_h \sqrt{P_h(1-P_h)} \right)^2}{r^2 N^2 P^2 + Z_{(1-\frac{\alpha}{r})}^2 \sum N_h P_h(1-P_h)}$$

$$n_0 = \frac{Z_{(1-\frac{\alpha}{r})}^2 \left(\sum N_h \sqrt{P_h(1-P_h)} \right)^2}{r^2 N^2 P^2}$$

اگر برای برآورد پارامتر مورد نظر Θ از نمونه‌گیری خوشه‌ای استفاده کنیم، تعبیر حداکثر خطای قابل پذیرش همانند نمونه‌گیری تصادفی ساده است زیرا در نمونه‌گیری خوشه‌ای ابتدا اندازه نمونه از طریق فرمولهای مربوط به نمونه‌گیری تصادفی ساده محاسبه می‌شود و سپس بر اساس متوسط اندازه خوشه و ضریب همبستگی درون خوشه‌ای تعدیل می‌گردد. در برآورد نسبی نیز اگر برآورد نسبت $R = \frac{\bar{Y}}{X}$ مد نظر باشد، اندازه نمونه مورد نیاز بر اساس خطای مطلق از رابطه زیر [شیفر و همکاران (۱۹۹۰)] بدست می‌آید:

$$n_0 = \frac{Z_{(1-\frac{\alpha}{r})}^2 \sigma^2 \mu_x}{d^2}$$

که d همان تعبیرهای مربوط به نمونه‌گیری تصادفی ساده را دارا است. همچنین اگر صفت مورد بررسی دارای توزیع غیر نرمال باشد، می‌توان بر اساس فرمول میانگین و واریانس این توزیعها فرمول تعیین اندازه نمونه را برای توزیعهای خاص به صورت تقریبی بدست آورد. برای مثال، در توزیع پواسون رابطه $\sigma^2 = \mu$ برقرار است. بنابراین اندازه نمونه بر اساس خطای نسبی از رابطه $n = \frac{Z(CV)}{r}$ بدست می‌آید. بنابراین می‌توان اظهار داشت با تغییر توزیع، فرمول تعیین اندازه نمونه تغییر می‌کند ولی تعبیر حداکثر خطای قابل پذیرش به قوت خود باقی است. در مورد آزمون فرضهای پیچیده‌تر مانند $H_0: \mu_1 = \mu_2$ یا $H_0: P_1 = P_2$ نیز این امر صادق است. برای جزئیات بیشتر در این زمینه می‌توانید به کتاب رن بیل^{۱۱} (۲۰۰۲) مراجعه کنید.

۷ نتیجه‌گیری

- برای برآورد میانگین و مجموع اندازه مورد نیاز بر اساس خطای مطلق تفاوتی با اندازه مورد نیاز بر اساس خطای نسبی ندارد ولی برای برآورد P استفاده از خطای نسبی موجب تغییر اساسی در فرمول اندازه نمونه می‌شود.
- استفاده از خطای نسبی در تعیین اندازه نمونه برای برآورد P فقط برای صفات نادر قابل توجیه است.
- در آمارگیریهایی که به صورت دوره‌ای (سالانه) اجرا می‌شوند، استفاده از خطای نسبی مناسبتر است.
- در مواردی که جامعه آماری به چند زیرجامعه شکسته می‌شود و برای هر زیر جامعه برآورد جداگانه‌ای مورد نیاز است، استفاده از خطای نسبی مناسبتر است.
- استفاده از خطای نسبی موجب می‌شود در تعیین اندازه نمونه، به جای واریانس از ضریب تغییرات استفاده شود که نسبت به واریانس از پایایی بیشتری برخوردار است.
- در استفاده از خطای نسبی باید توجه داشت حداکثر خطای نسبی باید در حدی باشد که خطای مطلق متناظر با آن قابل پذیرش باشد.

مراجع

- [1] Cochran, William G. (1977), Sampling Techniques. New York: John Wiley.
- [2] Cohen, J. (1988), Statistical Power Analysis For The Behavioural Sciences. New York: Academic Press.

- [3] Deming, W. E. (1966), Some Theory of Sampling. New York: John Wiley.
- [4] Hansen, M. H. Horwitz, W. N. and Madow, W. G. (1953), Sample Survey Methods and Theory. New York: John Wiley.
- [5] Kish, L. (1995), Survey Sampling. New York: John Wiley.
- [6] Lenth, R. V. (2001), Some Practical Guidelines for Effective Sample Size Determination. The American Statistician, Vol. 55, No. 3, 187-193.
- [7] Lohr, S. L. (1999), Sampling Design and Analysis. Duxbury Press.
- [8] Ran Belle, G. (2002), Statistical Rules of Thump. New York: John Wiley.
- [9] Scheaffer, R. L. Mendenhall, W. and Ott, L. (1990), Elementary Survey Sampling. Boston: PWS-KENT.
- [10] Thompson, S. K. (2002), Sampling. New York: John Wiley.
- [11] Zakkula, G. (1999), Elements of Sampling Theory and Methods. Prentice Hall.

مطالعه آنتروپی باقیمانده توزیع‌های طول عمر

یونس زهره‌وند^۱، مجید اسدی^۲

^۱ دانشجوی کارشناسی ارشد آمار کاربردی دانشگاه اصفهان

^۲ عضو هیئت علمی دانشگاه اصفهان

چکیده: در مبحث قابلیت اعتماد برای مطالعه طول عمر یک سیستم، تاکنون معیارهای مختلفی از قبیل؛ تابع نرخ شکست، میانگین باقیمانده عمر و ... توسط محققین معرفی شده‌اند و بر اساس هر یک از این معیارها کلاسهای متفاوتی از توزیعها ارائه شده و مطالعات فراوانی روی هر کدام از این کلاسها صورت گرفته است. در این مقاله با استفاده از معیار میزان عدم حتمیت یک متغیر تصادفی طول عمر و مباحث موجود در نظریه اطلاع دو کلاس جدید از توزیع‌های طول عمر را معرفی می‌کنیم و خواص آنها را مورد بررسی قرار می‌دهیم.

واژه‌های کلیدی: نظریه اطلاع، آنتروپی شانون، توزیع‌های طول عمر، میانگین باقیمانده عمر، تابع نرخ شکست، تابع نرخ شکست معکوس

۱ مقدمه

فرض کنید T یک متغیر تصادفی مطلقاً پیوسته و نامنفی با تابع توزیع $F(t)$ و تابع چگالی احتمال $f(t) = F'(t)$ باشد. در مباحث نظریه قابلیت اعتماد اندازه‌های متفاوتی برای مدل‌سازی و تجزیه و تحلیل داده‌های مورد مطالعه معرفی شده‌اند. مهمترین آنها عبارت است از تابع قابلیت اعتماد که آن را با $\bar{F}(t)$ نمایش می‌دهیم و عبارت است از: $\bar{F}(t) = P(T > t)$ یکی دیگر از اندازه‌های مطرح در قابلیت اعتماد تابع نرخ شکست می‌باشد که آن را با $\lambda_F(t)$ نمایش می‌دهیم و به صورت زیر تعریف می‌شود:

$$\lambda_F(t) = \lim_{\varepsilon \rightarrow 0} \frac{P(t < T \leq t + \varepsilon | T > t)}{\varepsilon}$$

$$= \frac{f(t)}{\bar{F}(t)}$$

احتمال موجود در فرمول فوق در واقع احتمال از کار افتادن بلافاصله بعد از زمان t است با این فرض که تا زمان t کار کرده است. همچنین میانگین باقیمانده عمر T نیز از اندازه‌های مهم در

مطالعات طول عمر است که آن را با $\mu_{F(t)}$ نمایش می‌دهیم و بصورت زیر تعریف می‌شود:

$$\mu_{F(t)} = E(T - t | T > t) = \frac{\int_t^{\infty} \bar{F}(x) dx}{\bar{F}(t)}$$

تابع بقا توسط $\lambda_{F(t)}$ و $\mu_{F(t)}$ به ترتیب بصورت زیر نمایش داده می‌شود:

$$\bar{F}(t) = \exp\left(-\int_0^t \lambda_{F(x)} dx\right)$$

$$\bar{F}(t) = \frac{\mu_{F(0)}}{\mu_{F(t)}} \exp\left(-\int_0^t \frac{1}{\mu_{F(x)}} dx\right)$$

همچنین روابط بین $\mu_{F(t)}$ و $\lambda_{F(t)}$ بصورت زیر است:

$$\lambda_{F(t)} = \frac{\mu'_{F(t)} + 1}{\mu_{F(t)}}$$

میزان عدم حتمیت T بر مبنای تعریف شانون (۱۹۴۸) بصورت زیر تعریف می‌شود:

$$H(f) = -\int_0^{\infty} f(t) \log f(t) dt$$

$H(f)$ معیاری برای اندازه‌گیری میزان عدم حتمیت و یا قابلیت پیش‌بینی متغیر تصادفی T می‌باشد. هر چه مقدار $H(f)$ برای یک متغیر تصادفی بیشتر باشد، توزیع آن متغیر یکنواخت‌تر می‌باشد و پیش‌بینی برآمدهای متغیر تصادفی T سخت‌تر است و بالعکس هر چه میزان $H(f)$ کمتر باشد توزیع T متمرکز و کشیده‌تر است و پیش‌بینی مقادیر برآمدهای آنی T ساده‌تر خواهد بود. آنتروپی شانون در حالت گسسته همواره کمیتی مثبت می‌باشد ولی در حالت پیوسته می‌تواند هر مقدار حقیقی را اختیار کند.

اغلب در مطالعات طول عمر و تحلیل بقا با حالاتی روبرو هستیم که در آنها واحد مورد مطالعه طول عمری بیشتر از مقدار معین t دارد. به عبارت دیگر فرض کنید T طول عمر یک قطعه باشد. حال اگر این قطعه تا زمان t عمر کرده باشد، یعنی $T > t$ آنگاه باقیمانده عمر آن $T - t$ است. لذا در چنین حالتی $H(f)$ معیار مناسبی برای اندازه‌گیری عدم حتمیت باقیمانده عمر سیستم نمی‌باشد. در این مقاله سعی می‌کنیم مفهوم آنتروپی باقیمانده و آنتروپی گذشته متغیر تصادفی T را که به ترتیب توسط ابراهیمی (۱۹۹۵) کرسنزو و لانگوباردی (۲۰۰۲) ارایه شده است را معرفی کرده و بعضی از خواص آنها را مورد بررسی قرار دهیم. به عبارت دیگر آنتروپی متغیرهای $T - t | T > t$ و $t - T | T < t$ را مورد مطالعه قرار می‌دهیم و بعضی از خواص آنها را استخراج می‌کنیم. به عنوان مثال نشان می‌دهیم آنتروپی باقیمانده و آنتروپی گذشته تحت بعضی شرایط تابع توزیع متغیر تصادفی T را مشخص می‌کنند. همچنین رفتار آنتروپی باقیمانده و آنتروپی گذشته را در سیستم‌های موازی و متوالی مورد مطالعه قرار می‌دهیم.

۲ آنترویی باقیمانده

قبل از معرفی آنترویی باقیمانده مثال زیر را در نظر بگیرید.
 مثال: فرض کنید متغیر تصادفی T دارای توزیع احتمالی بصورت زیر باشد:

$$p(T = t) = \begin{cases} 1/5 & t = 1, 2, 3, 4, 5 \\ 0 & o.w \end{cases}$$

با توجه به فرم توزیع T پیش‌بینی برآمد متغیر تصادفی T دشوارترین حالت ممکن در بین تمام توزیع‌های گسسته ۵ مقداری است. حال فرض کنید به طریقی مطلع می‌شویم که $T \geq 3$. در این صورت:

$$p(T = t | T \geq 3) = \begin{cases} 1/3 & t = 3, 4, 5 \\ 0 & o.w \end{cases}$$

ملاحظه می‌شود که در توزیع جدید پیش‌بینی مقدار T از بین سه مقدار ممکن ۳، ۴، ۵ خواهد بود که بوضوح ساده‌تر از قبل است. یعنی آنترویی متغیر تصادفی $(T | T \geq 3)$ کمتر از آنترویی T می‌باشد.

فرض کنید $f_t(x)$ تابع چگالی احتمال متغیر تصادفی $T | T > t$ باشد آنگاه

$$f_t(x) = \frac{f(x)}{F(t)}, \quad x > t$$

ابراهیمی (۱۹۹۵) آنترویی $f_t(x)$ را بصورت زیر تعریف کرد:

$$H(f; t) = - \int_t^\infty \left(\frac{f(x)}{F(t)} \right) \log \left(\frac{f(x)}{F(t)} \right) dx \\ = 1 - \frac{1}{F(t)} \int_t^\infty (\log \lambda_F(x)) f(x) dx \quad (2-1)$$

$H(f; t)$ به عنوان یک معیار پویا برای اندازه‌گیری عدم حتمیت T ، تمام خواص $H(f)$ را داراست. واضح است که

$$H(f; 0) = H(f)$$

لازم به ذکر است که: $-\infty < H(f; t) < +\infty$.

به راحتی می‌توان نشان داد که: $H(f; t) = c$ (c یک عدد ثابت است) اگر و تنها اگر T دارای توزیع نمایی باشد. ابراهیمی (۱۹۹۵) بر مبنای معیار $H(f; t)$ دو کلاس از توزیع‌های احتمال به صورت زیر تعریف کرد.

تعریف ۲.۱: متغیر تصادفی مطلقاً پیوسته و نامنفی T دارای آنترویی باقیمانده عمر نزولی (صعودی) $DURL$ ($IURL$) است هرگاه $H(f; t)$ تابعی نزولی (صعودی) از t باشد.

با توجه به تعریف توزیع نمایی با تابع نرخ شکست $\lambda_F(t) = c$ مرز بین توزیع‌های $IURL, DURL$ می‌باشد و این دو کلاس از توزیعها را از هم تفکیک می‌کند.
 تبصره ۲.۱: فرض کنید توزیع متغیر تصادفی T در کلاس توزیع‌های $DURL$ قرارگیرد بنابراین با گذشت زمان ($t > 0$)، توزیع متغیر تصادفی جدید ($T|T > t$) اطلاع بخش‌تر از توزیع T خواهد بود زیرا با گذشت زمان عدم حتمیت این متغیر کمتر می‌شود.
 ابراهیمی (۱۹۹۵) نتایج زیر را اثبات کرد:

نتیجه ۲.۱: اگر در توزیع F , $\mu_F(t) < \infty$ آنگاه $H(f; t) \leq 1 + \mu_F(t)$
 نتیجه ۲.۲: اگر توزیع F , $IURL$ ($DURL$) باشد آنگاه $H(f; t) \leq (\geq) 1 - \log \lambda_F(0)$
 نتیجه ۲.۳: اگر توزیع F , $DURL$ باشد آنگاه $H(f; t) \leq 1 + \log \mu$

یکی از خواص مهم توابعی مانند $\lambda_F(t)$ و $\mu_F(t)$ آن است که توزیع احتمال را به صورت منحصر بفرد مشخص می‌کنند. اکنون این سوال مطرح است که آیا رابطه بین F و $H(f; t)$ یک به یک است. در ادامه نشان می‌دهیم که در کلاس وسیعی از توزیعها، $H(f; t)$ به طور منحصر به فرد f را مشخص می‌کند.

قضیه ۲.۱: فرض کنید X و Y دو متغیر تصادفی نامنفی و مطلقاً پیوسته به ترتیب با توابع توزیع F و G باشند و F و G از خانواده توزیع‌های $IURL$ باشند و به ازای هر t داشته باشیم

$$H(f; t) = H(g; t) \quad (2-2)$$

آنگاه: $G(t) = F(t)$

اثبات: با مشتق‌گیری از طرفین تساوی (۲-۲) داریم:

$$H'(f; t) = H'(g; t)$$

که این نتیجه می‌دهد:

$$\lambda_F(t)[H(f; t) - 1 + \log \lambda_F(t)] = \lambda_G(t)[H(g; t) - 1 + \log \lambda_G(t)] \quad (2-3)$$

اگر به ازای هر t , $\lambda_F(t) = \lambda_G(t)$ آنگاه $G(t) = F(t)$ و قضیه اثبات می‌شود فرض کنید:

$$A = \{t; t > 0, \lambda_F(t) \neq \lambda_G(t)\}$$

و A تهی نباشد، یعنی

$$\exists t_0 \in A: \lambda_F(t_0) \neq \lambda_G(t_0)$$

بدون خلل به کلیت مسئله فرض کنید $\lambda_F(t_0) > \lambda_G(t_0)$ با توجه به صعودی بودن $H(f; t)$ و $H(g; t)$ به ازای هر t و به خصوص برای $t = t_0$ داریم: $0 \leq H(f; t_0) - 1 + \log \lambda_F(t_0)$ و $0 \leq H(g; t_0) - 1 + \log \lambda_G(t_0)$ و لذا از (۲-۳) نتیجه می‌گیریم:

$$H(f; t_0) - 1 + \log \lambda_F(t_0) \leq H(g; t_0) - 1 + \log \lambda_G(t_0)$$

در نتیجه $\lambda_F(t_0) \leq \lambda_G(t_0)$ و این با فرض $\lambda_F(t_0) > \lambda_G(t_0)$ تناقض دارد. بنابراین A شامل هیچ عضوی نیست و لذا به ازای هر t

$$\lambda_F(t) = \lambda_G(t)$$

که معادل است با

$$F(t) = G(t)$$

بلاک و همکارانش (۱۹۸۵) نشان دادند که اگر F و G توابع توزیع مربوط به دو متغیر تصادفی مطلقاً پیوسته X و Y با توابع نرخ خطر λ_F و λ_G باشند و به ازای هر t داشته باشیم $\lambda_G(t) = \theta(t)\lambda_F(t)$ که در آن $0 \leq \theta(x) \leq 1$ یک تابع صعودی (نزولی) از x است، در اینصورت اگر F ، IFR ، (DFR) ، $(DMRL)$ ، $(IMRL)$ ، $IFRA$ ، $(DFRA)$ ، NBU ، (NWU) باشد آنگاه G نیز آن خاصیت را خواهد داشت. اسدی و ابراهیمی (۲۰۰۰) نتیجه زیر را برای $H(f; t)$ نشان دادند.

لم ۲.۱: فرض کنید X و Y دو متغیر تصادفی مطلقاً پیوسته با توابع توزیع F و G و توابع نرخ شکست λ_F و λ_G و آنتروپی باقیمانده عمر $H(f; t)$ و $H(g; t)$ باشند و به ازای هر $t > 0$

$$\lambda_G(t) = \theta(t)\lambda_F(t) \quad (۲ - ۴)$$

که در آن $0 \leq \theta(t) \leq 1$ تابعی صعودی از t است. در اینصورت اگر F ، $DURL$ باشد آنگاه

$$G \text{ نیز } DURL \text{ است مشروط بر آنکه } \lim_{t \rightarrow \infty} \frac{\bar{G}(t)}{\bar{F}(t)} < \infty$$

اثبات: داریم

$$H(f; t) = 1 - \frac{1}{\bar{F}(t)} \int_t^{\infty} \log \lambda_F(x) f(x) dx = 1 - E_X(\lambda_F(X) | X > t)$$

$$H(g; t) = 1 - \frac{1}{\bar{G}(t)} \int_t^{\infty} \log \lambda_G(x) g(x) dx = 1 - E_Y(\lambda_G(Y) | Y > t)$$

بنابراین برای اثبات نزولی بودن $H(g; t)$ کافی است نشان دهیم که $E_Y(\lambda_G(Y) | Y > t)$ تابعی صعودی از t است از طرفی از (۲ - ۴) داریم:

$$\log \lambda_G(y) = \log \theta(y) + \log \lambda_F(y)$$

در نتیجه

$$E_Y(\lambda_G(Y) | Y > t) = E_Y(\theta(Y) | Y > t) + E_Y(\lambda_F(Y) | Y > t)$$

از طرفی چون $\theta(t)$ تابعی صعودی از t است بنابراین $E_Y(\theta(Y) | Y > t)$ نیز تابعی صعودی از t خواهد بود. لذا برای صعودی بودن $E_Y(\lambda_G(Y) | Y > t)$ کافی است نشان دهیم

می‌کنیم: $E_Y(\lambda_F(Y) | Y > t)$ صعودی است. بدین منظور m_1 و m_2 را بصورت زیر تعریف

$$m_1(t) = E_X(\lambda_F(X) | X > t)$$

$$m_2(t) = E_Y(\lambda_F(Y) | Y > t)$$

و قرار می‌دهیم

$$\beta(t) = \bar{G}(t)[m_1(t) - m_2(t)]$$

داریم:

$$\beta'(t) = \bar{G}(t)\{[\log \lambda_F(t) - m_1(t)]\lambda_G(t) + m_1'(t)\}$$

از طرفی $m_1'(t) = \lambda_F(t)[m_1(t) - \log \lambda_F(t)]$ که نتیجه می‌دهد:

$$\beta'(t) = \bar{G}(t)m_1'(t)\left[1 - \frac{\lambda_G(t)}{\lambda_F(t)}\right] > 0$$

در نتیجه $\beta(t)$ تابعی صعودی است. حال با فرض آنکه $\lim_{t \rightarrow \infty} \frac{\bar{G}(t)}{\bar{F}(t)} < \infty$ داریم:

$$\lim_{t \rightarrow \infty} \beta(t) = \lim_{t \rightarrow \infty} \left[\frac{\bar{G}(t)}{\bar{F}(t)} \int_t^\infty \log \lambda_F(x) f(x) dx - \int_t^\infty \log \lambda_F(x) g(x) dx \right] = 0$$

و بنابراین به ازای هر t ، $\beta(t) \leq 0$ و $m_1(t) \leq m_2(t)$ و چون $\lambda_F(t) \leq m_1(t)$ در نتیجه به ازای هر t ، $\lambda_F(t) \leq m_2(t)$ و بنابراین $m_2(t)$ تابعی صعودی است و بنابراین $H(g; t)$ نیز نزولی است.

در مباحث قابلیت اعتماد یک روش برای افزایش قابلیت سیستم، اضافه کردن مؤلفه بیشتر به سیستم است. یکی از مهمترین ساختارهای سیستم‌هایی که دارای مؤلفه اضافی هستند، سیستم‌های k از n می‌باشد.

سیستم k از n : سیستمی که عملکرد آن مستلزم عملکرد حداقل k تا از کل n مؤلفه‌اش می‌باشد. در حالت خاص اگر $k = n$ سیستم را سری گوئیم. و در حالت $k = n$ سیستم را موازی گوئیم.

فرض کنید $T_1, T_2, \dots, T_n \stackrel{iid}{\sim} F(\cdot)$ طول عمر n مؤلفه یک سیستم را نشان دهند و $T_{(1)}, \dots, T_{(n)}$ آماره‌های مرتب نظیر این نمونه باشند. در این صورت $T_{(k)}$ توزیع طول عمر یک سیستم $(n - k + 1)$ از n را نشان می‌دهد و توابع چگالی و توزیع و نرخ شکست طول

عمر این سیستم به ترتیب بصورت زیر می باشد.

$$f_{(k)}(t) = \frac{n!}{(k-1)!(n-k)!} F(t)^{k-1} f(t) \bar{F}(t)^{n-k}$$

$$F_{(k)}(t) = \sum_{i=k}^n \binom{n}{i} F(t)^i \bar{F}(t)^{n-i}$$

$$\lambda_{F_{(k)}}(t) = \frac{f_{(k)}(t)}{\bar{F}_{(k)}(t)}$$

$$= \frac{n!}{(k-1)!(n-k)!} \cdot \lambda_F(t) \cdot \frac{(F(t)/\bar{F}(t))^{k-1}}{\sum_{i=0}^{k-1} \binom{n}{i} \left(\frac{F(t)}{\bar{F}(t)}\right)^i}$$

از لم (۱-۲) قضیه زیر نتیجه می شود.

قضیه ۲.۲: اگر $F \stackrel{iid}{\sim} T_1, \dots, T_n$ با تابع چگالی f و تابع نرخ شکست $\lambda_F(t)$ و آنتروپی باقیمانده عمر نزولی $H(f; t)$ باشند، آنگاه $H(f_n; t)$ (آنتروپی باقیمانده عمر آماره مرتب $m|n$) نیز نزولی است.
اثبات: داریم

$$\lambda_{F_{(n)}}(t) = \frac{n \left(\frac{F(t)}{\bar{F}(t)}\right)^{n-1}}{\sum_{i=0}^{n-1} \binom{n}{i} \left(\frac{F(t)}{\bar{F}(t)}\right)^i} \lambda_F(t)$$

$$\theta(t) = n \frac{(F(t)/\bar{F}(t))^{n-1}}{\sum_{i=0}^{n-1} \binom{n}{i} (F(t)/\bar{F}(t))^i}$$

اگر قرار دهیم

$$\theta(t) = n \frac{(F(t)/\bar{F}(t))^{n-1}}{\sum_{i=0}^{n-1} \binom{n}{i} (F(t)/\bar{F}(t))^i}$$

آنگاه به راحتی ملاحظه می شود که $0 \leq \theta(t) \leq 1$ و $\theta(t)$ تابعی صعودی از t است. بنابراین با استفاده از لم قبل $H(f_n; t)$ هم تابعی نزولی از t است یعنی $DURL, F_{(n)}(t)$ است.

در واقع قضیه فوق بیان می‌کند که اگر یک سیستم موازی با n مولفه با توزیع طول عمر مشترک F و آنتروپی باقیمانده نزولی داشته باشیم آنگاه آنتروپی باقیمانده عمر کل سیستم نیز نزولی خواهد بود.

۳ آنتروپی گذشته

در بسیاری از حالات سیستم‌های مورد استفاده بصورت پیوسته مورد بازرسی قرار نمی‌گیرند. اکنون فرض کنید قطعه‌ای با طول عمر T را در زمان $t = 0$ بکار گرفته‌ایم و در زمان $t > 0$ مورد بازرسی قرار داده‌ایم و مشاهده نموده‌ایم که سیستم از کار افتاده است. به عبارت دیگر طول عمر سیستم کمتر از t است یعنی $T < t$. اکنون این سوال مطرح است که آنتروپی $t - T$ (یعنی زمان از کار افتادگی قطعه تا زمان t) $(t - T | T < t)$ را بصورت زیر تعریف کردند:

$$\begin{aligned} \bar{H}(f; t) &= - \int_0^t \left(\frac{f(x)}{F(t)} \right) \log \left(\frac{f(x)}{F(t)} \right) dx \\ &= 1 - \frac{1}{F(t)} \int_0^t (\log r_F(x)) f(x) dx \quad (3-1) \end{aligned}$$

که در آن $y < t$ ، $\frac{f(y)}{F(t)}$ تابع چگالی $(T | T < t)$ است و $r_F(t) = \frac{f(t)}{F(t)}$ تابع نرخ شکست معکوس می‌باشد. در زیر نشان می‌دهیم که $\bar{H}(f; t)$ تحت شرایطی F را بصورت منحصر بفرد مشخص می‌کند. به عبارت دیگر یک رابطه یک به یک بین F و $\bar{H}(f; \cdot)$ وجود دارد. **قضیه ۳.۱:** فرض کنید X و Y دو متغیر تصادفی نامنفی و مطلقاً پیوسته به ترتیب با توابع توزیع F و G باشند و به ازای هر t ، $\bar{H}(f; t) = \bar{H}(g; t)$ آنگاه در کلاس همه توزیع‌هایی که \bar{H} آنها تابعی نزولی از t است، داریم:

$$F(t) = G(t)$$

اثبات: با مشتق‌گیری از طرفین تساوی فوق داریم:

$$\bar{H}'(f; t) = \bar{H}'(g; t)$$

$$r_F(t)[1 - \bar{H}(f; t) - \log r_F(t)] = r_G(t)[1 - \bar{H}(g; t) - \log r_G(t)] \quad (3-1)$$

از آنجایی که r_F از اندازه‌های مشخصه‌ساز توزیع احتمال است اگر به ازای هر t ، $r_F(t) = r_G(t)$ آنگاه $F(t) = G(t)$ و قضیه اثبات می‌شود فرض کنید

$$\exists t_0 > 0 \quad r_F(t_0) \neq r_G(t_0)$$

بدون خلل به کلیت مسئله فرض کنید $r_F(t_0) > r_G(t_0)$ بنابراین با توجه به رابطه (۲ - ۳) و نزولی بودن $\bar{H}(f; t)$ به ازای هر t و به خصوص $t = t_0$ داریم:

$$1 - \bar{H}(f; t_0) - \log r_F(t_0) > 1 - \bar{H}(g; t_0) - \log r_G(t_0)$$

و یا

$$1 - \bar{H}(f; t_0) - \log r_F(t_0) = 1 - \bar{H}(g; t_0) - \log r_G(t_0)$$

که از دو رابطه فوق به راحتی نتیجه می شود $r_F(t_0) \leq r_G(t_0)$ و این با فرض $r_F(t_0) > r_G(t_0)$ تناقض دارد. بنابراین به ازای هر t

$$r_F(t) = r_G(t)$$

که معادل است با

$$F(t) = G(t)$$

نتیجه مهم زیر را برای آنتروپی گذشته نشان می دهیم.

لم ۳.۱: فرض کنید X و Y دو متغیر تصادفی مطلقاً پیوسته با توابع توزیع F و G و توابع نرخ شکست معکوس r_F و r_G باشند و داشته باشیم:

$$\forall t \geq 0; \quad r_G(t) = \theta(t)r_F(t)$$

که در آن $0 \leq \theta(t) \leq 1$ تابعی نزولی از t باشد. در این صورت اگر $\bar{H}(f; t)$ تابعی صعودی از

$$t \text{ باشد آنگاه } \bar{H}(g; t) \text{ نیز چنین است مشروط بر آنکه } \lim_{t \rightarrow 0} \frac{G(t)}{F(t)} < \infty$$

اثبات: مشابه اثبات لم (۳.۱) می باشد و از آوردن آن صرفنظر می شود.

در قضیه زیر رفتار $\bar{H}(f; t)$ یک سیستم موازی را بر اساس آنتروپی گذشته مولفه هایش تحت شرایطی خاص بررسی می کنیم.

قضیه ۳.۲: اگر $T_1, \dots, T_n \stackrel{iid}{\sim} F$ با تابع چگالی مشترک f و تابع نرخ شکست معکوس $r_F(t)$ و آنتروپی گذشته صعودی $\bar{H}(f; t)$ باشند. آنگاه $\bar{H}(f_{(1)}; t)$ آنتروپی گذشته آماره مرتب اول، نیز صعودی است.

اثبات: اگر $T_{(1)} = \min\{T_1, \dots, T_n\}$ آنگاه داریم

$$\begin{aligned} r_{T_{(1)}}(t) &= \frac{f_{T_{(1)}}(t)}{F_{T_{(1)}}(t)} \\ &= \frac{f(t)}{F(t)} \cdot \frac{n\bar{F}(t)^{n-1}}{\sum_{i=0}^{n-1} \bar{F}(t)^i} \end{aligned}$$

با قرار دادن

$$\theta(t) = \frac{n\bar{F}(t)^{n-1}}{\sum_{i=0}^{n-1} \bar{F}(t)^i}$$

بوضوح $0 \leq \theta(t) \leq 1$ و $\theta(t)$ تابعی نزولی از t است بنابراین با استفاده از لم قبل $\bar{H}(f_{(1)}; t)$ تابعی صعودی از t است و آنتروپی گذشته $F_{(1)}(t)$ هم صعودی است. در واقع قضیه فوق بیان می‌کند که اگر یک سیستم سری با n مولفه با توزیع طول عمر مشترک F و آنتروپی گذشته صعودی داشته باشیم آنگاه آنتروپی گذشته کل سیستم نیز صعودی خواهد بود.

مراجع

- [1] Ebrahimi, N.(1996). How to measure uncertainty in the residual life time distribution. *Sankhya*, vol 58 p;48-56.
- [2] Asadi, M. Ebrahimi, N.(2000).Residual entropy and its characterizaations in terms of hazard function and mean residual life function. *Statistics and Probability letters*, vol 49 P;263-269.
- [3] DI Cresenzo,A. Longobardi,M.(2002) Entropy-based measure of uncertainty in past lifetime distributions.*J. Appl. Prob.* vol 39 P;434-440.
- [4] Ebrahimi, N. Pelleray, F.(1995) New partial ordering of survival functions based on the notion of uncertainty.*Appl. Prob.* vol 32 p;202-211.

یک الگوریتم مقدماتی برای شبیه‌سازی عددی معادلات دیفرانسیل تصادفی

پرویز سرگلزایی^۱، محمد امینی^۱، محمود دادخواه^۲

^۱ دانشگاه سیستان و بلوچستان، گروه ریاضی

^۲ دانشجوی کارشناسی ارشد ریاضی دانشگاه سیستان و بلوچستان

چکیده: در این مقاله چند روش عددی سودمند برای حل معادلات دیفرانسیل تصادفی ارائه می‌شود. انتگرالهای تصادفی، روشهای اولر - مارویاما و ملشتاین برای حل عددی معادلات دیفرانسیل تصادفی و همگرایی قوی و ضعیف مورد توجه قرار گرفته است. در ادامه چند مثال عددی هم ارائه می‌گردد.

واژه‌های کلیدی: روش اولر - مارویاما، روش ملشتاین، روش مونت کارلو، همگرایی ضعیف و قوی، بسط تصادفی تیلور

۱ مقدمه

امروزه مدل‌های معادلات دیفرانسیل تصادفی نقشی بسیار مهم و برجسته در حوزه‌های مختلف علوم ایفا می‌کنند. از جمله‌ی این علوم می‌توان به زیست‌شناسی، شیمی، مکانیک، اقتصاد، پزشکی، الکترونیک، ریاضی و ... اشاره کرد. تفهیم و آشنایی کامل با نظریه‌ی معادلات دیفرانسیل تصادفی منوط به آشنایی پیشرفته با احتمال و فرایندهای تصادفی است اما فهم شبیه‌سازی عددی مورد بحث در این مقاله برای معادلات دیفرانسیل تصادفی، نیازمند آشنایی و داشتن زمینه‌ای در مورد روش اولر برای حل عددی معادلات دیفرانسیل معمولی و فهم شهودی متغیرهای تصادفی است. بعلاوه آشنایی با روشهای عددی حل معادلات دیفرانسیل معمولی، گام مفید و موثر اولیه برای کار با معادلات دیفرانسیل تصادفی را در اختیار ما قرار می‌دهد. بهر حال ما در این مقاله فرض کرده‌ایم که خواننده حداقل به طور سطحی با مفاهیمی از قبیل متغیرهای تصادفی، استقلال، مقادیر مورد انتظار (امیدها)، واریانس و بخصوص مفهوم متغیرهای تصادفی با توزیع نرمال آشناست. آزمایشها و نتایج عددی ما از روش مونت کارلو اقتباس شده است (که بیان می‌دارد: متغیرهای تصادفی به کمک تولید اعداد تصادفی و مقادیر مورد انتظار با میانگینهای محاسبه شده، قابل حصول‌اند). برای اجرای هر چه سریع‌تر و بهتر برنامه‌ها در *Matlab*، یک تکنیک جدید ارائه کرده‌ایم و سعی کرده‌ایم کدهای برنامه را به صورت برداری در آورده و اجرا نماییم.

۲ حرکت براونی (فرایند وینر)

حرکت براونی استاندارد یا فرایند وینر استاندارد روی $[0, T]$ ، یک متغیر تصادفی $W(t)$ است که بطور پیوسته به $t \in [0, T]$ وابسته بوده و در سه شرط ذیل هم صدق می نماید:

$$(1) \quad W(0) = 0 \quad (w.p.1)$$

(۲) برای هر $0 \leq s < t \leq T$ ، متغیرهای تصادفی داده شده با نمو $W(t) - W(s)$ توزیع نرمال با میانگین صفر و واریانس $t - s$ داشته باشند به عبارت دیگر داشته باشیم:

$$W(t) - W(s) \sim N(0, t - s)$$

(۳) برای هر $0 \leq s < t < u < v \leq T$ ، نموهای $W(t) - W(s)$ و $W(v) - W(u)$ مستقل از هم باشند.

برای کارهای محاسباتی، در نظر گرفتن حرکت براونی گسسته سازی شده که در آن $W(t)$ در مقادیر گسسته سازی شده t معین شده است، مناسب است. بنابراین فرض می کنیم $\delta t = \frac{T}{N}$ (عدد صحیح) و نیز $t_j = j\delta t$ ، $W_j = W(t_j)$. بنا بر شرط (۲) و (۳) داریم:

$$(1) \quad W_j = W_{j-1} + dW_j, \quad j = 1, 2, \dots, N$$

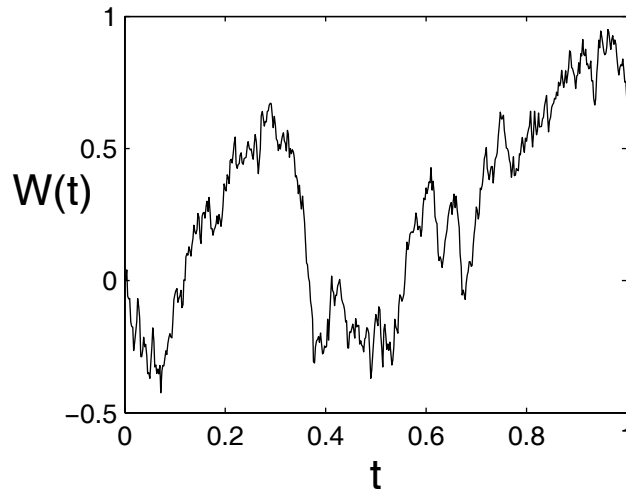
که هر dW_j یک متغیر تصادفی مستقل به فرم $N(0, \delta t)$ می باشد. برنامه‌ی *Matlab* ذیل شبیه سازی ساده و متفاوتی را با توجه به مقادیر اخذ شده از کاربر برای فرایند وینر روی $[0, T]$ می دهد. از تابع مولد اعداد تصادفی *randn* استفاده کرده ایم تا در هر بار بازخوانی آن عددی تصادفی با توزیع $N(0, 1)$ داشته باشیم. از وضعیت *state* هم برای دریافت نتایج عددی مشابه در هر تکرار استفاده شده است. برای تکرار پذیر بودن نتایج عددی، مقدار اولیه را هم از کاربر دریافت کرده ایم اما معمول آنست که آنرا صفر در نظر بگیریم. ما نمودار را برای $T = 1$ و $N = 500$ رسم کرده ایم.

```

randn('state',100)
T=input(' final time:');
N=input(' number of divisions:');
dt=T/N; dW=zeros(1,N);
W=zeros(1,N);
w(1)=dW(1);
for j=2:N
dW(j)=sqrt(dt)*randn;
W(j)=W(j-1)+dW(j);
end
subplot(2,2,1);
plot([0:dt:T],[0,W], 'r-')
xlabel('t','FontSize',16) , ylabel('W(t)','FontSize',16,'Rotation',0)

```

***** win.m : Brownian Path Simulation *****



حالا می‌خواهیم با یک تکنیک خوب، معادل برنامه‌ی قبلی را به صورتی کارتر نوشته و از مزاحمت‌های حلقه‌ی *For* رهایی پیدا کنیم. این کار را با جانشین سازی حلقه‌ی *For* با دستور سطح بالاتری بنام *Cumsum* انجام می‌دهیم. این کاریکی از تکنیک‌های بالا بردن کارایی برنامه است. زامین مولفه‌ی این دستور بصورت زیر محاسبه می‌شود:

$$dW(1) + dW(2) + \dots + dW(j)$$

در برنامه‌ی *win2.m* این کار را انجام داده‌ایم. اگر شکل مربوط به این برنامه را هم رسم کنیم، دقیقاً همانند شکلی است که در حالت قبلی برای حلقه‌ی *For* رسم کرده بودیم، با این تفاوت که از مزاحمت حلقه‌ی *For* در اینجا خبری نیست.

```
randn('state',100)
T=input('Input final time:');
N=input('Input number of divisions:');
dt=T/N; dW=sqrt(dt)*randn(1,N);
W=cumsum(dW); subplot(2,2,1);
plot([0:dt:T],[0,W], 'r-')
xlabel('t','FontSize',16)
ylabel('W(t)','FontSize',16,'Rotation',0)
```

*** *win2.m* ***

حالا می‌خواهیم یک تابع ساده‌ی تصادفی را شبیه‌سازی کنیم. در برنامه‌ی *bpath.m* ذیل، تابع $U(W(t)) = e^{t + \frac{1}{2}W(t)}$ را حول ۱۰۰۰ مسیر براونی گسسته‌سازی شده رسم و ارزیابی کرده‌ایم.

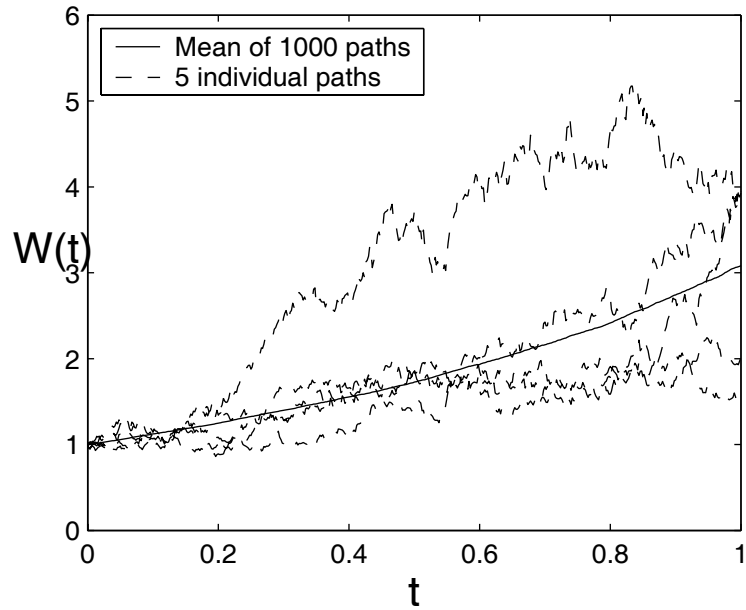
```
randn('state',100)
T=1; N=500;
```

```

dt=T/N; t=[dt:dt:T]; M=1000;
dW=sqrt(dt)*randn(M,N);
W=cumsum(dW,2);
U=exp(repmat(t,[M 1])+0.5*W);
Umean=mean(U);
subplot(1.5,1.5,1);
plot([0,t],[T,Umean],'b-'),
hold on
plot([0,t],[ones(5,T),U(T:5,:)],...
'r-'), hold off
xlabel('t','FontSize',16)
ylabel('W(t)','FontSize',16,...
'Rotation',0)
legend('Mean of 1000 paths','5 individual paths',2)
avver=norm((Umean-exp(9*t/8)),'inf')

```

*** *bpath.m* : Function along a Brownian Path ***



این تابع به عنوان حالت خاصی از جوابهای معادلات دیفرانسیل خطی مطرح است. با فرض $T = 1$ و $N = 500$ ، ۵ مسیر مجزا و میانگین این مسیرها را هم رسم کرده‌ایم. برای افزایش کارایی برنامه سعی کرده‌ایم برنامه را برداری نماییم، بطوریکه dW یک آرایه $M \times N$ است که $dW(i, j)$ نمو dW_j برای i امین مسیر را می‌دهد. دستور *repmat* برای معرفی یک آرایه‌ی $M \times N$ بکار می‌رود بطوریکه سطرهاى آن تماماً نسخه‌ای از t می‌باشد. با توجه به

شکل ترسیمی در می‌یابیم که علی‌رغم ناهموار بودن $U(W(t))$ در طی مسیرهای مجرد، میانگین آن هموار به نظر می‌رسد. امید یا مقدار مورد انتظار برای $U(W(t))$ برابر $e^{\frac{\lambda t}{\lambda}}$ است. مقدار $avver = 0,0504$ بدست آمده است. هرگاه جواب یک معادله‌ی دیفرانسیل برای مسیر معینی خواسته شده باشد، گوییم جواب قوی مورد نیاز است و این منجر به مفهوم همگرایی قوی می‌شود. در مقابل هرگاه فقط اطلاعات نوع مقدار مورد انتظار (امید) در مورد جواب مورد مطالعه قرار گیرد، منجر به مفهوم همگرایی ضعیف خواهد شد.

۳ انتگرالهای تصادفی

اگر تابع مناسبی همانند f داده شده باشد، انتگرال $\int_0^T f(t)dt$ را ممکن است با جمع سری ریمن ذیل تقریب بزنیم:

$$\sum_{j=0}^{N-1} f(t_j)(t_{j+1} - t_j) \quad (2)$$

که نقاط گسسته‌سازی شده $t_j = j\delta t$ در بخش (۲) معرفی شدند. هرگاه از سری (۲) وقتی $\delta t \rightarrow 0$ حد گرفته شود، این عبارت به شکل انتگرال تعریف خواهد شد. به طریقی مشابه می‌توان مجموعی به شکل ذیل را هم در نظر گرفت:

$$\sum_{j=0}^{N-1} f(t_j)(W(t_{j+1}) - W(t_j)) \quad (3)$$

که در مقایسه با سری (۲) در بالا، می‌توان آنرا به عنوان تقریبی برای انتگرال تصادفی $\int_0^T f(t)dW(t)$ در نظر گرفت (انتگرال ایتو). از طرفی می‌توان یک تقریب دیگر نیز برای سری (۲) به فرم زیر، برای بدست آوردن انتگرال، در نظر گرفت:

$$\sum_{j=0}^{N-1} f\left(\frac{t_j + t_{j+1}}{2}\right)(t_{j+1} - t_j) \quad (4)$$

فرم متناظر تصادفی برای تقریب (۴) می‌تواند چنین بیان شود (انتگرال استراتونویچ):

$$\sum_{j=0}^{N-1} f\left(\frac{t_j + t_{j+1}}{2}\right)(W(t_{j+1}) - W(t_j)) \quad (5)$$

در حالتی که $f(t) = W(t)$ باشد، نیاز به محاسبه $W(t)$ در نقطه $t = \frac{t_j+t_{j+1}}{2}$ داریم. می‌توان نشان داد که به ازای یک نمو مستقل $N(\circ, \frac{\Delta t}{2})$ و با $\frac{W(t_{j+1})+W(t_j)}{2}$ ، مقداری برای $W(\frac{t_j+t_{j+1}}{2})$ بدست می‌آید که در سه شرط بیان شده در ابتدای بخش (۲)، مطابقت دارد. به کمک همین روش، در برنامه‌ی *stint.m*، با تقریب انتگرالهای ایتو و نیز استراتونوویچ، مقدار خطا را در مورد انتگرال ایتو برابر $itoerr = \circ.158$ و خطا در مورد انتگرال استراتونوویچ برابر $straerr = \circ.186$ بدست آورده‌ایم.

```

randn('state',100)
T=input('Input final time:');
N=input('Input number of divisions:');
dt=T/N; dW=sqrt(dt)*randn(1,N);
W=cumsum(dW);
ito=sum([0,W(1:end-1)].*dW)
strat=sum((0.5*[0,W(1:end-1)]+W)+0.5*sqrt(dt)*randn(1,N)).*dW)
itoerr=abs(ito-0.5*(W(end)^2-T))
straterr=abs(strat-0.5*W(end)^2)

```

*** stint.m : Approximate Stochastic Integrals ***

در حالت انتگرال ایتو داریم:

$$\begin{aligned}
 \sum_{j=\circ}^{N-1} W(t_j)(W(t_{j+1}) - W(t_j)) &= \frac{1}{2} \sum_{j=\circ}^{N-1} \{ (W(t_{j+1}))^2 - W(t_j)^2 \\
 &- (W(t_{j+1}) - W(t_j))^2 \} \\
 &= \frac{1}{2} (W(T)^2 - W(\circ)^2) \\
 &- \frac{1}{2} \sum_{j=\circ}^{N-1} \{ (W(t_{j+1}) - W(t_j))^2 \} \quad (6)
 \end{aligned}$$

حال جمله‌ی $\sum_{j=\circ}^{N-1} \{ (W(t_{j+1}) - W(t_j))^2 \}$ در (۶) دارای امید ریاضی یا مقدار مورد انتظار T و واریانس $O(\delta t)$ است. بنابراین برای δt ی کوچک، این متغیر تصادفی به ثابت T نزدیک خواهد شد و بنابراین:

$$\int_{\circ}^T W(t)dW(t) = \frac{1}{2} \{ (W^2(T) - T) \} \quad (7)$$

و معادل (۶) و (۷) برای انتگرال استراتونوویچ چنین می شود $(\Delta Z_j \sim N(0, \frac{\Delta t}{\nu}))$:

$$\sum_{j=0}^{N-1} \left(\frac{W(t_{j+1}) + W(t_j)}{2} + \Delta Z_j \right) (W(t_{j+1}) - W(t_j)) = \frac{1}{2} (W(T)^2 - W(0)^2) + \sum_{j=0}^{N-1} \Delta Z_j \{ (W(t_{j+1}) - W(t_j)) \} \quad (8)$$

جمله‌ی دوم در (۸)، دارای امید ریاضی ۰ و واریانس $O(\delta t)$ است. بنابراین داریم:

$$(Stra..) \int_0^T W(t) dW(t) = \frac{1}{2} W^2(T)$$

هر کدام از انتگرالهای ایتو و استراتونوویچ کاربرد خاص خود را در علوم مختلف دارند و بسته به مواقع نیاز به کار گرفته می شوند، با این وجود، این انتگرالها قابل تبدیل به یکدیگر هستند و می توان با داشتن یکی از آنها، فرم دیگری را بدست آورد.

۴ روش اولر - مارویاما

معادله‌ی دیفرانسیل تصادفی زیر به فرم دیفرانسیلی آن را در نظر بگیرید:

$$dX(t) = f(X(t))dt + g(X(t))dW(t) \quad X(0) = X_0, \quad 0 \leq t \leq T \quad (9)$$

که می توان آنرا به فرم انتگرالی زیر هم نوشت:

$$X(t) = X(0) + \int_0^t f(X(s))ds + \int_0^t g(X(s))dW(s) \quad 0 \leq t \leq T \quad (10)$$

که در آن f و g توابع اسکالری و شرط اولیه‌ی X_0 یک متغیر تصادفی می باشد. هنگام بررسی (10) ، انتگرال دوم در سمت راست آن، باید با توجه به فرایند براونی (وینر) که در بخش قبلی توضیح دادیم، برآورد گردد. در این مقاله ما از شکل ایتو استفاده خواهیم کرد. جواب $X(t)$ برای هر t ، یک متغیر تصادفی خواهد بود. از توضیح بیشتر در مورد تمام خواص جواب $X(t)$ برای (10) خودداری می نماییم و بجای آن روشهای عددی حل (10) را معرفی می نماییم. بنابراین جواب $X(t)$ را متغیر تصادفی حاصل از میل دادن حد طول گام به صفر، در نظر می گیریم. برای راحتی کار، بیشتر با فرم (9) کار خواهیم کرد. (توجه کنید چون فرایند وینر با احتمال ۱، هیچ جا مشتق پذیر نیست، اجازه نداریم از $\frac{dW}{dt}$ استفاده کنیم.)

حال اگر $g \equiv 0$ و X_0 ثابت باشد، آنگاه معادله از حالت تصادفی به حالت جبری تغییر فرم می‌دهد. اکنون برای اعمال یک روش عددی بر (۹) روی بازه $[0, T]$ ، ابتدا بازه را گسسته‌سازی می‌کنیم. فرض کنید $\Delta t = \frac{T}{N}$ (N عدد صحیح) و $t'_j = j\Delta t$ باشد. تقریب عددی برای $X(t'_j)$ را با X_j نمایش داده و روش اولر - مارویاما را چنین معرفی می‌نماییم:

$$X_j = X_{j-1} + f(X_{j-1})\Delta t + g(X_{j-1})(W(t'_j) - W(t'_{j-1})) \quad (11)$$

در این مقاله، مسیرهای براونی گسسته‌سازی شده را محاسبه کرده و از آنها برای حصول نمودار $W(t'_j) - W(t'_{j-1})$ که در (۱۱) لازم داریم، استفاده خواهیم کرد. برای راحتی کار در محاسبات کامپیوتری، تقریباً همه جا طول گام Δt برای نتایج عددی را، مضرب $R \geq 1$ از نمودار δt مسیر براونی فرض خواهیم کرد. این کار به ما اطمینان می‌دهد که مجموعه نقاط $\{t_j\}$ که مبنای گسسته‌سازی مسیر براونی است، شامل مجموعه نقاط $\{t'_j\}$ که جواب روش اولر - مارویاما در آن محاسبه شده است، می‌باشد. به هر حال گاهی مسیر براونی به عنوان قسمتی از داده‌ی مسئله بیان شده است اما اگر مسیر تحلیلی بکار گرفته شود، Δt به اندازه‌ی دلخواه کوچک می‌تواند بکار گرفته شود. حالا روش اولر - مارویاما را برای حل معادله‌ی اسکالری ذیل بکار می‌بندیم:

$$dX(t) = \lambda X(t)dt + \mu X(t)dW(t), \quad X(0) = X_0. \quad (12)$$

فرض می‌کنیم که λ و μ ثابت‌های حقیقی‌اند. این نوع معادله‌های اسکالری در بسیاری از علوم کاربرد دارند و من جمله می‌توان از آن به عنوان مدل ارزش سرمایه در ریاضیات مالی یا مدل رشد سهام در بازار بورس استفاده کرد [۴]. جواب واقعی این نوع معادلات را می‌توان به صورت زیر بیان کرد [۸]:

$$X(t) = X_0 e^{(\lambda - \frac{1}{2}\mu^2)t + \mu W(t)} \quad (13)$$

در برنامه‌ی ذیل، $\lambda = 2$ و $\mu = 1$ و $X_0 = 1$ در نظر گرفته‌ایم. مسیر براونی گسسته‌سازی شده حول $[0, 1]$ با $\delta t = 2^{-8}$ محاسبه شده است. بعلاوه جواب واقعی [یعنی (۱۳)] را در X_{true} قرار داده‌ایم. ضریب نمودار یعنی $R = 4$ قرار داده شده است.

```

randn('state',100)
lambda=2; mu=1; X0=1; T=1;
N=2^8; dt=T/N;
dW=sqrt(dt)*randn(1,N);
W=cumsum(dW);
Xtrue=X0*exp((lambda-0.5*mu^2)...
*(dt:dt:T)+mu*W);
subplot(1,4,1,4,1);
plot([0:dt:T],[X0,Xtrue], 'm-'),...
hold on
R=4; Dt=R*dt; l=N/R;

```

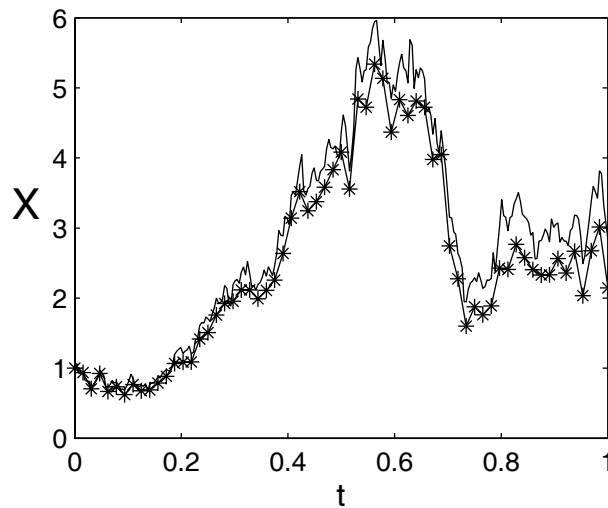


```

Xem=zeros(1,1); Xtemp=X0;
for j=1:l
Winc=sum(dW(R*(j-1)+1:R*j));
Xtemp=Xtemp+Dt*lambda*Xtemp+mu*Xtemp*Winc;
Xem(j)=Xtemp;
end
subplot(1,4,1,4,1);
plot([0:Dt:T],[X0,Xem], 'r-*'), hold off
xlabel('t','FontSize',12), ylabel('X','FontSize',16,'Rotation',0)
emerr=abs(Xem(end)-Xtrue(end))

```

*** em.m : Euler – Maroyama method for Scalar SDE ***



با تابع $emerr$ مقدار خطا را در نقطه‌ی انتهایی ارزیابی کرده‌ایم و در اینجا مقدار آنرا $۰٫۶۹۰۷$ بدست آورده‌ایم. با اخذ مقادیر کوچکتری مثل $R = ۲$ و $R = ۱$ خطا را بصورت $۰٫۱۵۹۵$ و $۰٫۸۲۱$ بدست می‌آوریم. توجه کنید که در گام کلی، روش اولر- مارویاما به نمو $W(t'_j) - W(t'_{j-1})$ نیاز دارد که با عبارت زیر محاسبه شده است:

$$W(t'_j) - W(t'_{j-1}) = W(jR\delta t) - W((j-1)R\delta t) = \sum_{k=jR-R+1}^{jR} dW_k$$

در برنامه‌ی نوشته شده دستور $Winc = sum(dW(R*(j-1)+1:R*j))$ این کار را انجام می‌دهد.

۵ همگرایی قوی و ضعیف روش اولر - مارویاما

اگر برنامه‌ی روش اولر - مارویاما (یعنی $em.m$) را با Δt های کوچکتری بیازماییم، در می‌یابیم که هر قدر Δt کوچکتر می‌شود، جواب روش اولر - مارویاما سازگاری و نزدیکی بیشتری با جواب واقعی پیدا می‌کند. با توجه به اینکه، X_n و $X(t_n)$ متغیرهای تصادفی هستند، برای بیان دقیق و واضح مفهوم همگرایی باید تصمیم بگیریم که چگونه تفاضل بین آنها را اندازه‌گیری کنیم. اگر از $E|X_n - X(t_n)|$ ، که E امید ریاضی را می‌نمایاند، استفاده شود به مفهوم همگرایی قوی منجر می‌شود.

گوییم یک روش مرتبه‌ی قوی همگرایی γ دارد هرگاه ثابتی همانند C چنان باشد که برای Δt های بقدر کافی کوچک داشته باشیم:

$$E|X_n - X(t')| \leq C\Delta t^\gamma, \quad (t' = n\delta t \in [0, T]) \quad (14)$$

به عنوان مثال اگر توابع f و g توابعی هموار بوده و در شرایط خاصی صدق کنند، ثابت می‌شود که روش اولر - مارویاما همگرایی قوی مرتبه‌ی $\frac{1}{2}$ دارد (این مطلب را با توجه به شکل برنامه $em.m$ می‌توان دید). در بررسیهای عددی که انجام خواهیم داد، بر روی خطا در آخرین نقطه تمرکز می‌نماییم. پس فرض کنیم:

$$e_{\Delta t}^S = E|X_N - X(T)|, \quad N\Delta t = T \quad (15)$$

نمایانگر خطای نقطه‌ی انتهایی روش اولر - مارویاما (به مفهوم قوی همگرایی) باشد. اگر کران $C\Delta t^\gamma$ در (۱۴) با $\frac{1}{2}$ $\gamma =$ ، برای هر نقطه‌ی ثابت در $[0, T]$ برقرار باشد، قطعاً در نقطه‌ی انتهایی هم برقرار است. لذا برای Δt های بقدر کافی کوچک می‌توان داشت:

$$e_{\Delta t}^S \leq C\Delta t^{\frac{1}{2}} \quad (16)$$

در برنامه‌ی ذیل بدنبال همگرایی روش اولر - مارویاما برای معادله‌ای به شکل (۱۲) هستیم. به ازای $\lambda = 2$ و $\mu = 1$ و $X_0 = 1$ و $T = 1$ و $N = 2^9$ ، 1000 مسیر براونی گسسته‌سازی شده را بکار برده‌ایم. در هر مسیری روش اولر - مارویاما، به ازای 5 طول گام متفاوت $\Delta t = 2^{p-1}\delta t$ ($1 \leq p \leq 5$) محاسبه شده است. مقدار خطا در s امین مسیر برای p امین طول گام، با عبارت $Xerr(s, p)$ بررسی شده است. تابع $mean$ میانگین مسیرها را محاسبه می‌کند طوری که p امین مولفه از $mean(Xerr)$ تقریبی از $e_{\Delta t}^S$ به ازای $\Delta t = 2^{p-1}\delta t$ است.

```
randn('state',100)
lambda=input(' drift coefficient:');
mu=input(' diffusion coefficient:');
X0=input(' initial value:');
```

```

T=input(' final time:');
N=input(' number of divisions:');
dt=T/N; M=1000; Xerr=zeros(M,5);
for s=1:M
dW=sqrt(dt)*randn(1,N);
W=cumsum(dW);
Xtrue=X0*exp((lambda-0.5*mu^2)+mu*W(end));
for p=1:5
R=2^(p-1); Dt=R*dt; l=N/R;
Xtemp=X0;
for j=1:l
Winc=sum(dW(r*(j-1)+1:r*j));
Xtemp=Xtemp+Dt*lambda*Xtemp+mu*Xtemp*Winc;
end
Xerr(s,p)=abs(Xtemp-Xtrue);
end
end
Dtvals=dt*(2.^([0:4])); subplot(221);
loglog(Dtvals,mean(Xerr),'b*-'), hold on
loglog(Dtvals,(Dtvals.^(.5)),'r-'), hold off
axis([1e-3 1e-1 1e-4 1])
xlabel('\Delta t'), ylabel('Sample average of |X(T) - X_t| ')
title('emstrong.m','FontSize',10)
•• Least squares fit of error=C*Dt^q••
A=[ones(5,1), log(Dtvals)']; rhs=log(mean(Xerr));
sol=A\rhs, q=sol(2), resid=norm(A*sol-rhs)

*** ems.m : Strong convergence of Euler - Maroyama method ***

```

حالا اگر فرض کنیم نامساوی (۱۶) بطور تقریبی به تساوی برقرار باشد، داریم: $e_{\Delta t}^S \approx C \Delta t^{\frac{1}{4}}$ که اگر از طرفین آن لگاریتم بگیریم داریم:

$$\log e_{\Delta t}^S \approx \log C + \frac{1}{4} \log \Delta t \quad (17)$$

برای رسم لگاریتمی تقریبی از $e_{\Delta t}^S$ از دستور $\log\log(Dtvals, mean(Xerr))$ استفاده کرده‌ایم. اگر در شکل رسم شده (شکل ۱ بالای گوشه‌ی سمت چپ) دقت شود، می‌توان دید که دو منحنی با هم به خوبی سازگارند که بیانگر دقیق و خوب بودن فرمول (۱۷) است. برای بررسی بیشتر کارایی فرمول (۱۷)، فرض کنیم که قاعده‌ی توانی $e_{\Delta t}^S = C \Delta t^q$ برای ثابت‌هایی همانند C و q برقرار باشد طوری که:

$$\log e_{\Delta t}^S = \log C + q \log \Delta t$$

به کمک تقریب کمترین مربعات برای $\log C$ و q در انتهای برنامه‌ی *ems.m*، مقدار $q = ۰.۵۳۸۴$ و باقیمانده‌ی ۰.۲۶۶ بدست می‌آید که خود بیانگر اعتبار نتایج حاصله برای مرتبه‌ی قوی $\frac{1}{4}$

است. ذکر این نکته لازم است که ما در بررسی خطای $e_{\Delta t}^S$ ، از منابع دیگر خطا که قابل چشم پوشی بودند، گذشته‌ایم. این منابع عبارتند از:

- (۱) خطای نمونه: یعنی خطایی که موقع تقریب میانگین نمونه با امید ریاضی حاصل می‌شود.
- (۲) خطای اریب اعداد تصادفی: یعنی خطای ذاتی مولد اعداد تصادفی
- (۳) خطای گرد کردن: که همان خطای گرد کردن اعداد ممیز شناور می‌باشد.

در محاسبات اصلی از میان سه مرجع خطای ذکر شده‌ی فوق، خطای نمونه از بقیه‌ی خطاها مهمتر است. از اینرو در بعضی از برنامه‌های این مقاله برای حصول مرتبه‌ی همگرایی پیش‌بینی شده مجبور شده‌ایم حجم نمونه‌ها را بسیار بزرگ و طول گام را تا حد ممکن کوچک اختیار کنیم (اگر M حجم نمونه باشد، خطای نمونه متناسب با $\frac{1}{\sqrt{M}}$ کاهش می‌یابد). البته یک بررسی در [۷] نشان داده است که مراجع دیگر ذکر شده‌ی خطا قبل از اینکه در مسئله اهمیت پیدا کنند از بین می‌روند.

همگرایی قوی ذکر شده در (۱۴) میزان سرعت کاهش میانگین خطا را وقتی $\Delta t \rightarrow 0$ ، اندازه گیری می‌کند اما ممکن است یک الگوی کم کاربرد دیگر هم گاهی مد نظر باشد و آن عبارتست از اندازه‌گیری سرعت کاهش خطای میانگین‌ها و این منجر به مفهوم همگرایی ضعیف می‌شود. گویم روشی دارای مرتبه‌ی ضعیف همگرایی γ است هرگاه ثابت C چنان باشد که برای Δt های بقدر کافی کوچک و هر تابع g داشته باشیم:

$$|Eg(X_n) - Eg(X(t'))| \leq C\Delta t^\gamma, \quad (t' = n\delta t \in [0, T]) \quad (18)$$

البته تابع g باید تابعی بطور مناسب هموار بوده و در شرایط رشد چند جمله‌ایها صدق کند. در این مقاله توجه ما به حالتی معطوف خواهد بود که در آن g یک تابع یکه (همانی) است. در [۱۱] ثابت شده که برای توابع f و g بطور مناسب هموار، روش اولر - مارویاما مرتبه‌ی همگرایی ضعیف ۱ دارد. به پیروی از آنچه در مورد همگرایی قوی گفتیم، فرض کنیم:

$$e_{\Delta t}^W = |EX_N - EX(T)|, \quad N\Delta t = T$$

نمایانگر خطای نقطه‌ی انتهایی روش اولر - مارویاما (به مفهوم ضعیف همگرایی) باشد. حال اگر $g(x) = x$ و $\gamma = 1$ باشد، از (۱۸) نتیجه می‌شود که برای Δt های بقدر کافی کوچک، داریم:

$$e_{\Delta t}^W \leq C\Delta t$$

در برنامه‌ی $emw.m$ ذیل همگرایی مرتبه‌ی ضعیف روش اولر - مارویاما را اجرا کرده‌ایم. در این برنامه معادله‌ی (۱۲) را به ازای $\lambda = 2$ و $\mu = 1$ و $\sigma = 1$ اجرا کرده‌ایم. مطابق مطالبی که قبلاً بیان شد، حدود 50000 مسیر براونی گسسته شده را آزموده‌ایم تا نتیجه‌ی لازم حاصل

شود. $\Delta t = 2^{p-1}$, $(1 \leq p \leq 5)$ فرض شده است. این برنامه در مقایسه با $ems.m$ یک سطح اضافی برداری شدن در بر دارد یعنی ما شبیه‌سازی را با تمام 50000 مسیر انجام داده‌ایم. این کار زمان اجرای انجام برنامه را بهبود می‌بخشد اما حافظه‌ی زیادتری اشغال می‌کند. برای جبران این اشکال مسیرهای مختلفی را استفاده کرده‌ایم، طوری که بجای تمامی مسیرها، فقط کفیسست نمونه‌های جاری ذخیره شوند. علاوه بر این برای افزایش کارایی برنامه $\Delta t = \delta t$ فرض شده است. تقریبهای میانگین نمونه برای EX_N در Xem و خطای روش در $Xerr$ ذخیره شده‌اند. با توجه به (۱۲) داریم: $E(X(t)) = e^{\lambda t}$ و از این مطلب جواب واقعی را بدست آورده‌ایم. در اینجا، مطابق قانون توانی که قبلاً گفتیم، مقدار $q = 0.9858$ و باقیمانده 0.0142 بدست آمده است که موید همگرایی ضعیف مرتبه‌ی ۱ روش است. (شکل ۱ بالای گوشه‌ی سمت راست)

```

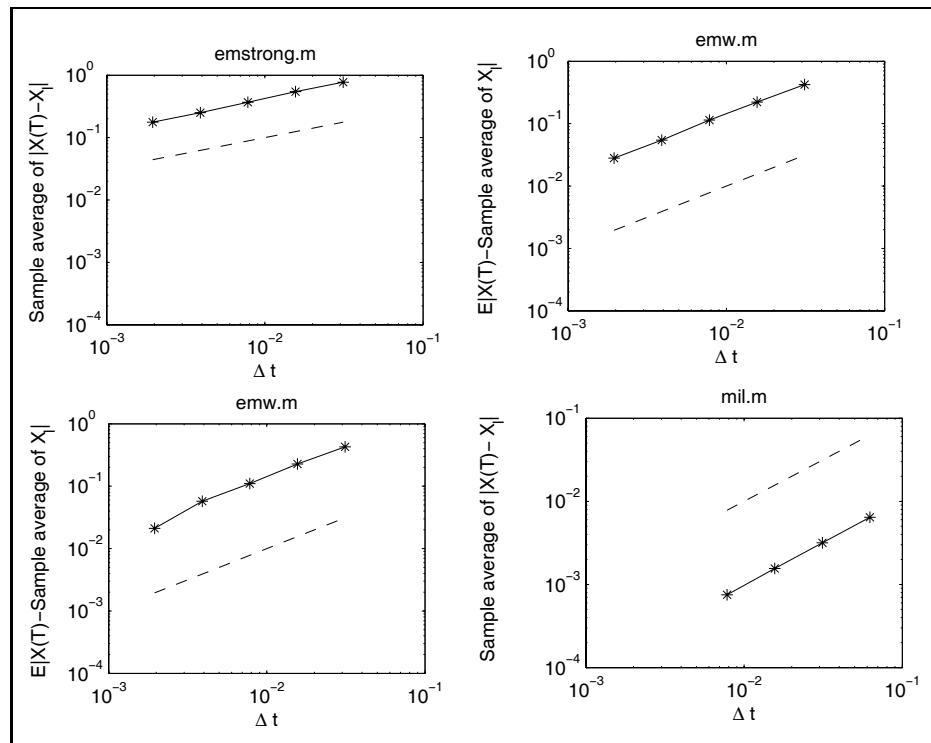
randn('state',100)
lambda=input(' drift coefficient:');
mu=input(' diffusion coefficient:');
X0=input(' initial value:');
T=input(' final time:');
M=50000; Xem=zeros(5,1);
for p=1:5
Dt=2^(p-10); L=T/Dt;
Xtemp=X0*ones(M,1);
for j=1:L
Winc=sqrt(Dt)*randn(M,1);
••Winc=sqrt(Dt)*sign(randn(M,1)); Weak E-M method ••
Xtemp=Xtemp+Dt*lambda*Xtemp+mu*Xtemp.*Winc;
end
Xem(p)=mean(Xtemp);
end
Xerr=abs(Xem-exp(lambda));
Dtvals=2.^([1:5]-10);
subplot(222);
loglog(Dtvals,Xerr,'b*-'), hold on
loglog(Dtvals,Dtvals,'r-'), hold off
axis([1e-3 1e-1 1e-4 1])
xlabel('\ Delta t'), ylabel('E|X(T) - Sample average of X_t| ')
title('emw.m', 'FontSize',10)
A=[ones(p,1), log(Dtvals)']; rhs=log(Xerr);
sol=A\rhs; q=sol(2)
resid=norm(A*sol-rhs)

```

***** emw.m : Weak convergence of Euler – Maroyama method *****

لازم است تأکید کنیم که ما در $emw.m$ برای هر طول گام Δt ، مسیرهای گوناگونی را استفاده کرده‌ایم و این کار هم معقولانه است زیرا همگرایی ضعیف تنها با میانگین جوابها سروکار دارد و ما مختار هستیم که هر نمونه‌ی $(0, 1)$ $\sqrt{\Delta t}N(0, 1)$ را برای نمو $W(t'_j) - W(t'_{j-1})$ در (۱۱) در هر

گامی قرار دهیم و در واقع حتی اگر نمورا با متغیر تصادفی دو نقطه‌ای مستقل $\sqrt{\Delta t}V_j$ (مقادیر $+1$ و -1 را با احتمال مساوی می‌پذیرد). جانشین نماییم، مرتبه‌ی همگرایی ضعیف همچنان ثابت باقی می‌ماند. (توجه کنید که V_j دارای میانگین و واریانس همانند $N(0, 1)$ است). اگر در این مرحله نمودهای وینر را با $\sqrt{\Delta t}V_j$ عوض کنیم، روش منسوب به روش ضعیف اولر - مارویاما بدست می‌آید که مرتبه‌ی ضعیف ۱ دارد اما چون از اطلاعات مسیر استفاده نمی‌کند، همگرایی قوی در بر ندارد. نکته‌ی اصلی در استفاده از روش ضعیف اولر - مارویاما در این است که اعداد تصادفی که توسط مولد اعداد تصادفی از V_j تولید می‌شوند، نسبت به آنهایی که از $N(0, 1)$ تولید می‌شوند، معمولاً کاراتر واقع می‌شوند. ما این موضوع را در برنامه‌ی *emw.m* در خطی که بین دو • قرار دارد، بیان کرده‌ایم. اگر این خط هم در اجرای برنامه قرار گیرد، مقدار $q = 1.671$ و باقیمانده برابر 0.2096 بدست می‌آید. (شکل ۱ پایین سمت چپ)



شکل ۱

۶ روش ملشتاین با مرتبه‌ی بالاتر

در بخش گذشته دیدیم که روش اولر - مارویاما مرتبه‌ی قوی همگرایی $\frac{1}{4}$ و ضعیف ۱ دارد. این در حالیست که روش اولر برای معادلات جبری، مرتبه‌ی قوی ۱ دارد. امکان بالا بردن مرتبه‌ی قوی روش اولر - مارویاما با اضافه کردن یک تصحیح در نمو آن وجود دارد. این کار منجر به روش منسوب به ملشتاین^۱ می‌گردد. امکان این افزایش با توجه به اینکه بسط روش قدیمی تیلور باید با توجه به فرمول ایتو تصحیح شود، وجود دارد. بسط تصادفی ایتو - تیلور یاریگر ما در این زمینه خواهد بود. هرگاه این بسط را از نقطه‌ی معینی برش بزنیم و باقی جملات را بعنوان خطا فرض کنیم، روش ملشتاین چنین بدست می‌آید:

$$X_j = X_{j-1} + f(X_{j-1})\Delta t + g(X_{j-1})(W(t'_j) - W(t'_{j-1})) \\ + \frac{1}{4}g'(X_{j-1})g(X_{j-1})\{(W(t'_j) - W(t'_{j-1}))^2 - \Delta t\}$$

در برنامه‌ی *mil.m* ذیل، ما روش ملشتاین را برای معادله‌ی زیر اجرا کرده‌ایم:

$$dX(t) = rX(t)(k - X(t))dt + \beta X(t)dW(t), \quad X(0) = X_0. \quad (19)$$

این معادله که در آن β, k, r مقادیر ثابتی هستند، یکی از معادلات خاص رشد جمعیتی است [۹].

در اجرای برنامه $r = 3, k = 1.5, \beta = 0.25, X_0 = 0.5$ فرض شده است. مسیر حول $[0, 1]$ با $\delta t = 2^{-11}$ گسسته‌سازی شده است.

جواب معادله‌ی (۱۹) را می‌توان به شکل انتگرالی نوشت. برای راحتی، جواب ملشتاین را با $\Delta t = \delta t$ در نظر گرفته‌ایم تا تقریب خوبی برای جواب واقعی باشد و و آنرا با جواب ملشتاین در ازای $\Delta t = 128\delta t, \Delta t = 64\delta t, \Delta t = 32\delta t, \Delta t = 16\delta t$ با 50° مسیر مقایسه کرده‌ایم. در مقایسه با *ems.m*، این برنامه یک سطح اضافی دارد یعنی بجای استفاده از حلقه‌ی *for* برای تغییر حول مسیرهای نمونه، آنرا با تمامی مسیرها بطور همزمان محاسبه کرده‌ایم. در این برنامه $q = 10^3 16$ و باقیمانده برابر 0.0095 شده است. (شکل ۱ پایین سمت راست)

```
randn('state',100)
r=2; k=1; beta=0.25; X0=0.5;
T=1; N=2^(11); dt=T/N;
M=500;
R=[1; 16; 32; 64; 128];
dW=sqrt(dt)*randn(M,N);
Xmil=zeros(M,5);
for p=1:5
Dt=R(p)*dt; L=N/R(p);
Xtemp=X0*ones(M,1);
```

1) Milstein

```

for j=1:L
Winc=sum(dW(:,R(p)*(j-1)+1:R(p)*j),2);
Xtemp=Xtemp+Dt*r*Xtemp.*(k-Xtemp)+beta*Xtemp.*Winc...
+0.5*beta^2*Xtemp.*(Winc.^2-Dt);
end
Xmil(:,p)=Xtemp;
end
Xref=Xmil(:,1);
Xerr=abs(Xmil(:,2:5)-repmat(Xref,1,4));
mean(Xerr);
Dtvals=dt*R(2:5);
subplot(2.3,2.3,4);
loglog(Dtvals,mean(Xerr),'b*-'), hold on
loglog(Dtvals,Dtvals,'r-'), hold off
axis([1e-3 1e-1 1e-4 1])
xlabel('\Delta t'), ylabel('Sample average of |X(T) - X_l| ')
title('mil.m', 'FontSize',10)
A=[ones(4,1), log(Dtvals)]; rhs=log(mean(Xerr)');
sol=A\rhs; q=sol(2)
resid=norm(A*sol-rhs)

```

*** mil.m : Milstien method with Strong convergence ***

۷ قانون زنجیره‌ای تصادفی

در بخش ۳ دیدیم که بیش از یک راه برای تعریف انتگرالها به مفهوم تصادفی وجود دارد. در این بخش به طور خلاصه به یکی دیگر از تفاوت‌های محاسبات جبری و تصادفی اشاره می‌نماییم. در حالت جبری، اگر $dX/dt = f(X)$ ، آنگاه برای هر تابع هموار U بر طبق قانون زنجیره‌ای می‌توان نوشت:

$$\frac{dU(X(t))}{dt} = \frac{dU(X(t))}{d(X(t))} \cdot \frac{d(X(t))}{dt} = \frac{dU(X(t))}{d(X(t))} \cdot f(X(t)) \quad (۲۰)$$

حالا فرض کنید که $X(t)$ در معادله‌ی ای‌توی (۹) صدق می‌کند. در مقایسه با (۲۰)، در مورد $U(X)$ در این وضع چه می‌توان گفت؟ یک حدس معقول و مستدل این است که: $dU = \left(\frac{dU}{dX}\right)dX$ طوریکه به کمک (۹)؛

$$dU(X(t)) = \frac{dU(X(t))}{d(X)}(f(X(t))dt + g(X(t))dW(t))$$

تجزیه و تحلیلی کمی مشکل، به کمک محاسبات ایتو، وجود یک جمله‌ی اضافی را آشکار ساخته و فرمول درست چنین بدست می‌آید:

$$dU(X(t)) = \frac{dU(X(t))}{dX}dX + \frac{1}{2}g(X(t))^2 \frac{d^2U(X(t))}{dX^2}dt$$

که با کمک (۹) از آن نتیجه می‌شود:

$$dU(X(t)) = (f(X(t))\frac{dU(X(t))}{dX} + \frac{1}{2}g(X(t))^2 \frac{d^2U(X(t))}{dX^2})dt + g(X(t))\frac{dU(X(t))}{dX}dW(t) \quad (21)$$

با توجه به اینکه جزئیات کامل این مطالب در [۹] و [۱۱] و [۶] وجود دارد، ما در پی اثبات یا توجیه روابط فوق نیستیم، اما یک آزمایش عددی را با آن شکل می‌دهیم. معادله‌ی دیفرانسیل تصادفی ذیل را در نظر می‌گیریم:

$$dX(t) = (\alpha - X(t))dt + \beta\sqrt{X(t)}dW(t), \quad X(0) = X_0 \quad (22)$$

که α, β پارامترهای مثبت و ثابت‌اند. این یک معادله‌ی دیفرانسیل ریشه‌ی مربعی است که ارزش سرمایه‌ی گذاشته شده را مدل‌سازی می‌کند [۸]. می‌توان نشان داد که اگر $X(0) \geq 0$ (با احتمال ۱)، آنگاه این مثبت بودن برای تمام $t > 0$ ، باز هم باقیست. با فرض $U(X) = \sqrt{X}$ ، اجرای روش (۲۱) نتیجه می‌دهد:

$$dU(t) = \left(\frac{\alpha - \beta^2}{4U(t)} - \frac{1}{4}U(t) \right)dt + \frac{1}{4}\beta W(t) \quad (23)$$

در برنامه‌ی *chain.m* ذیل، روش اولر - مارویاما را برای معادله‌ی (۲۲) با $\alpha = 2, \beta = 1$ اجرا کرده‌ایم. مسیر براونی حول $[0, 1]$ با $\delta t = \frac{1}{100}$ و $\Delta t = \delta t$ را اجرا کرده‌ایم. همچنین (۲۳) را برای U حل کرده‌ایم و آنرا هم رسم کرده‌ایم. با توجه به شکل در می‌یابیم که تطابق خوبی بین این دو جواب وجود دارد، زیرا تفاضل بین آنها $Xdiff = 0.151$ شده است.

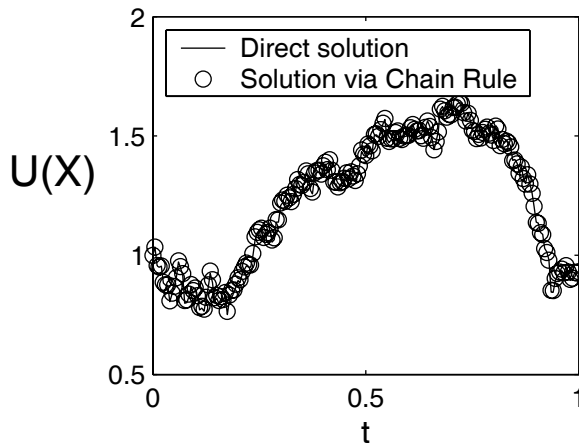
```
randn('state',100)
alpha=2; beta=1; X0=1; T=1;
N=200; dt=T/N; X01=1/sqrt(X0);
Dt=dt;
Xem1=zeros(1,N); Xem2=zeros(1,N);
Xtemp1=X0; Xtemp2=X01;
for j=1:N
Winc=sqrt(dt)*randn;
f1=(alpha-Xtemp1);
```

```

g1=beta*sqrt(abs(Xtemp1));
Xtemp1=Xtemp1+Dt*f1+Winc*g1;
Xem1(j)=Xtemp1;
f2=(4*alpha-beta^2)/(8*Xtemp2)-Xtemp2/2;
g2=beta/2;
Xtemp2=Xtemp2+Dt*f2+Winc*g2;
Xem2(j)=Xtemp2;
end
subplot(1,2,1,2,1)
plot([0:Dt:T],[sqrt([X0,abs(Xem1)])], 'b-',[0:Dt:T],[X0,Xem2], 'ro')
legend('Direct solution','Solution via Chain Rule',2)
xlabel('t','FontSize',12)
ylabel('U(X)','FontSize',16,'Rotation',0,'HorizontalAlignment','right')
Xdiff=norm(sqrt(Xem1)-Xem2,'inf')

```

**** chain.m : Test Stochastic Chain Rule ****



۸ نتیجه گیری

در این مقاله ضمن آنکه از کمترین جزییات بهره گرفته شده است، سعی شده یک تکنیک کاربردی و عملی برای شبیه سازی عددی معادلات دیفرانسیل تصادفی در اختیار خواننده قرار گیرد. بنابراین از بسیاری از روشهای پیشرفته تر صرف نظر شده و اثبات مطالب اصلی ارجاع گردید. به عنوان مثال بیان نکردیم که چه شرایطی باید بر f, g اعمال شود تا معادله دیفرانسیل ایتوی (۹) جواب تضمینی و منحصر بفرد داشته باشد یا در جای دیگر، ما بیشتر توجه خود را به معادلات دیفرانسیل اسکالری معطوف کردیم. روش اولر- مارویاما تفاوت چندانانی در حالت اسکالری یا غیر اسکالری

ندارد اما روش ملشتاین در حالت غیر اسکالری بسیار پیچیده می‌شود. به هر حال خواننده می‌تواند بدنبال روشهای عددی کاراتر و با مرتبه‌ی همگرایی بالاتر هم باشد.

مراجع

- [1] T. C. GARD, Introduction to Stochastic Differential Equations, Marcel Dekker, New York (1988).
- [2] D. J. HIGHAM, Mean-Square and asymptotic stability of the stochastic theta method, SIAM J. Numer. Anal., 38 (2000), pp. 753-769.
- [3] D. J. HIGHAM AND N. J. HIGHAM, MATLAB Guide, SIAM, Philadelphia, 2000.
- [4] J. C. HULL, Options, Futures, and Other Derivatives, 4th ed., Prentice-Hall, Upper Saddle River, NJ, 2000.
- [5] I. KARATZAS AND S. E. SHREVE, Brownian Motion and Stochastic Calculus, 2th ed., Springer-Verlag, Berlin, 1991.
- [6] P. E. KLOEDEN AND E. PLATEN, Numerical Solution of Stochastic Differential Equations, Springer-Verlag, Berlin, 1999.
- [7] Y. KOMORI, Y. SAITO, AND T. MITSUI, Some issues in discrete approximate solution for stochastic differential equations, Comput. Math. Appl., 28 (1994), pp. 269-276.
- [8] X. MAO, Stochastic Differential Equations, Hoewood, Chichester, 1977.
- [9] B. ØKSENDAL, Stochastic Differential Equations, 5th ed., Springer-Verlag, Berlin, 1998.
- [10] H. C. ÖTTINGER, Stochasti Processes in polymeric Fluids, Springer-Verlag, Berlin, 1996.
- [11] E. PLATEN, An introduction to numerical methods for stochastic differential equations, Acta Numer., 8 (1999), pp. 197-246.
- [12] Y. SAITO, AND T. MITSUI, Stability analysis of numerical schemes for stochastic differential equations, SIAM J. Numer. Anal., 33(1996), pp. 2254-2267.
- [13] THE MATH WORKS, INC, MATLAB User's Guide, Natick, Massachusetts, 1992.
- [14] L. N. TREFETHEN, Spectral Methods in MATLAB, SIAM, Philadelphia, 2000.

رگرسیون ژرفا در حالت چند گانه

حمید شریف

بانک مرکزی جمهوری اسلامی ایران

چکیده: در این مقاله به بیان نظریه ژرفا در رگرسیون می‌پردازیم. رگرسیون ژرفا^۱ که عدد صحیحی بین 0 و n است به عنوان خاصیتی از برازش رگرسیونی در نظر گرفته می‌شود که رتبه برازش رگرسیونی تعبیر می‌شود و معیاری برای مقایسه برازشهای مختلف رگرسیونی می‌باشد. در این مقاله با بیان رگرسیون ژرفا در حالت چندگانه به بررسی خواص آن پرداخته و روش ژرفترین رگرسیون را که نسبت به تبدیلات یکنوا بر متغیر پاسخ هم‌ورد است، ارائه می‌دهیم. با معرفی معیار مقدار فروریزش نشان می‌دهیم که این روش، روشی استوار می‌باشد. در پایان برآوردکننده ژرفترین رگرسیون را با دیگر برآوردکننده‌های رگرسیونی با محاسبه رگرسیون ژرفای آنها برای داده‌های نمونه‌ای موجود مقایسه می‌کنیم.

واژه‌های کلیدی: ژرفا، نابرازا^۲، ژرفترین رگرسیون، مقدار فروریزش، رگرسیون ژرفای ماکسیمال، هم‌وردایی

۱ مقدمه

وقتی رگرسیون کمترین توانهای دوم را با استفاده از n مشاهده برای یک مدل پارامتری $Y = X\beta + \varepsilon$ به کار می‌بریم فرضهایی را در مورد بردار خطاها ε در نظر می‌گیریم. یکی از این فرضها نرمال بودن توزیع خطاهاست. در عمل انحرافات از این فرضها رخ می‌دهد. مثلاً ممکن است توزیع زیربنایی خطاها متقارن اما غیر نرمال باشد یا تیزتر از نرمال بوده و دمهای کوتاهتری داشته باشد یا کشیدگی کمتر از نرمال با دمهای پهن‌تر داشته باشد و یا ممکن است توزیع به صورت نرمال باشد ولی دارای دور افتاده‌هایی باشد.

در این موارد به جای روش کمترین توانهای دوم از روشهای رگرسیون استوار استفاده می‌شود که در مقایسه با روش کمترین توانهای دوم نسبت به این انحرافات حساسیت کمتری دارند. یک روش رگرسیونی استوار که اخیراً توسط روسیف و هوبرت^۳ (۱۹۹۹) مطرح شده است روش ژرفترین رگرسیون نام دارد که بر پایه نظریه رگرسیون ژرفا بنا شده است. رگرسیون ژرفا کیفیت برازش را اندازه‌گیری کرده و میزان دوری آن را از هر نابرازا اندازه می‌گیرد و بیان می‌کند که ابرصفحه^۴ برازشی در توصیف داده‌ها به چه اندازه خوب عمل می‌کند. بنابراین برازشی با ژرفای بزرگ نسبت به

1) regression depth 2) nonfit 3) Rousseeuw and Hubert

داده‌ها متعادلتر است و لذا یک برازش خوب، ژرفای بزرگتری نسبت به یک برازش بد دارد. در این مقاله با مرور رگرسیون ژرفا در حالت چندگانه به بررسی خواص آن پرداخته و روش ژرفترین رگرسیون را بیان می‌کنیم.

در بخش دوم این مقاله، بطور مختصر رگرسیون ژرفا را در حالت ساده مرور می‌کنیم. در بخش سوم رگرسیون ژرفا را در حالت چندگانه معرفی می‌کنیم. رگرسیون ژرفای ماکسیمال و ژرفترین رگرسیون به ترتیب در بخشهای چهارم و پنجم بیان می‌شود و در بخش ششم خواص ژرفترین رگرسیون در حالت چندگانه بررسی می‌گردد. این خواص شامل هم‌وردایی و استواری ژرفترین رگرسیون است و در انتها با ارائه مثال برآوردکننده ژرفترین رگرسیون را با دیگر برآوردکننده‌های رگرسیونی مقایسه می‌کنیم.

۲ مروری بر رگرسیون ژرفا در حالت ساده

هدف در رگرسیون ساده برازش یک خط راست $y = \theta_1 x + \theta_2$ به یک مجموعه داده $Z_n = \{(x_i, y_i); i = 1, \dots, n\} \subseteq R^2$ است. یک برازش را به صورت $\theta = (\theta_1, \theta_2)$ نشان می‌دهیم که مؤلفه اول آن برآورد شیب و مؤلفه دوم آن جمله عرض از مبدا است. مانده‌های مجموعه داده Z_n متناسب با برازش θ را به صورت $r_i(\theta) = y_i - \theta_1 x_i - \theta_2$ نشان می‌دهیم. برای معرفی ژرفای یک برازش ابتدا یک نابرازا را تعریف می‌کنیم.

تعریف ۱: برازش $\theta = (\theta_1, \theta_2)$ برای مجموعه داده Z_n نابرازا نامیده می‌شود اگر و فقط اگر یک عدد حقیقی $v = v_\theta$ مخالف همه x_i ها وجود داشته باشد به طوری که

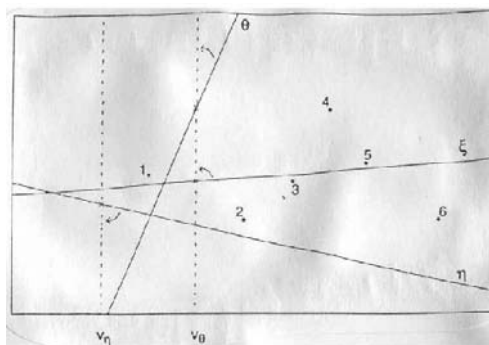
$$\forall x_i < v, \quad r_i(\theta) < 0, \quad \forall x_i > v, \quad r_i(\theta) > 0$$

یا

$$\forall x_i < v, \quad r_i(\theta) > 0, \quad \forall x_i > v, \quad r_i(\theta) < 0$$

خطی که بالا یا پایین همه مشاهدات واقع می‌شود همیشه یک نابرازا است. بعداً خواهیم دید که نابرازاها در واقع برازش‌هایی با ژرفای صفر هستند.

شکل ۱ یک مجموعه داده را با ۶ مشاهده θ و η نابرازا نشان می‌دهد. مقادیر مربوطه v_θ و v_η نشان داده شده‌اند. با توجه به شکل، وجود v متناظر با وجود یک ضربدر روی خط عمودی که از نقطه v رسم شده، است. چون $\forall x_i < v_\theta, r_i(\theta) > 0$ و همچنین $\forall x_i > v_\theta, r_i(\theta) < 0$ لذا با توجه به تعریف فوق θ ، نابرازا است. همچنین چون خط η در پایین همه مشاهدات قرار دارد این خط نیز نابرازا است. از آن جا که تعریف فوق برای خط ξ صادق نیست این خط نابرازا نیست. همان طور که ذکر شد با مشخص شدن v ، نقطه متناظر با آن که محل تلاقی خط عمودی ترسیم شده از نقطه v با خط رگرسیونی است و با ضربدر نشان داده شده است، مشخص می‌گردد. به عبارت دیگر محل تلاقی خط برازش داده شده با خط عمودی که از نقطه v متناظر با آن خط ترسیم می‌گردد با ضربدر مشخص شده است. از این



شکل ۱: مجموعه داده دو متغیره با نابرازهای θ و η و یک برازش ξ با رگرسیون ژرفای ۲

رو می‌توان خطوط نابرازا را از روی شکل ۱ مشخص کرد. به این ترتیب که خط برازش داده شده را حول نقطه‌ای که با ضربدر بر روی آن خط مشخص شده است می‌گردانیم تا عمودی شود. اگر در حین چرخاندن، خط از هیچ مشاهده‌ای عبور نکند نابرازا است. بنابراین خطوط θ و η ، نابرازا هستند ولی خط ξ ، نابرازا نیست چون وقتی آن را حول نقطه مذکور (مشخص شده با ضربدر) می‌گردانیم تا عمودی شود از مشاهدات ۴ و ۵ عبور می‌کند.

به طور کلی ژرفای یک برازش θ برای یک مجموعه داده Z_n به حجم n به صورت زیر بیان می‌شود:

تعریف ۲: رگرسیون ژرفای ($rdepth$) یک برازش θ متناسب با مجموعه داده Z_n عبارت است از کمترین تعداد مشاهداتی که باید برداشته شود تا θ نابرازا شود. به طور معادل، $rdepth(\theta, Z_n)$ عبارت است از کوچکترین تعداد مانده‌هایی که می‌بایست علامتشان تغییر کند تا θ نابرازا شود.

برای مثال خط ξ را در شکل ۱ در نظر می‌گیریم. این خط با حذف مشاهدات ۴ و ۵ نابرازا می‌شود (زیرا با قرار دادن v_ξ مساوی با v_θ و حذف مشاهدات ۴ و ۵ می‌توان خط ξ را به صورت عمودی درآورد بدون آن که از مشاهده‌ای عبور کند). چون خط ξ با حذف حداقل ۲ مشاهده نابرازا می‌شود لذا $rdepth(\xi, Z_n) = 2$.

توجه ۱: تعاریف ۱ و ۲ در مواقعی که در x_i ها تکرار وجود داشته باشد (x_i) ها با هم مساوی باشند نیز برقرار است. و x_i ها به هیچ فرض توزیعی نیاز ندارند.

مرتبه زمانی: گوئیم $f(n)$ از مرتبه زمانی $g(n)$ است و آن را با نماد $f(n) = O(g(n))$ نشان می‌دهیم اگر و تنها اگر اعداد صحیح و مثبتی مانند n_0 و c وجود داشته باشد بطوریکه $f(n) \leq c g(n)$ برای تمام n های $n \geq n_0$ برقرار باشد.

برای محاسبه $rdepth(\theta, Z_n)$ ابتدا مشاهدات را به صورت $x_1 \leq x_2 \leq \dots \leq x_n$ در

مرتبه زمانی $O(n \log n)$ مرتب می‌کنیم. تعریف می‌کنیم:

$$L^+(v) = \#\{j; x_j \leq v, r_j \geq 0\}$$

و

$$R^-(v) = \#\{j; x_j > v, r_j \leq 0\}$$

L^+ و R^- به صورت مشابه تعریف می‌شوند. سپس $L^+(x_i), L^-(x_i), R^+(x_i)$ و $R^-(x_i)$ برای هر $i = 1, \dots, n$ مشابه تعاریف فوق، محاسبه می‌شوند. $rdepth(\theta, Z_n)$ در $O(n)$ عمل به صورت زیر محاسبه می‌شود.

$$rdepth(\theta, Z_n) = \min_{1 \leq i \leq n} (\min\{L^+(x_i) + R^-(x_i), R^+(x_i) + L^-(x_i)\}) \quad (1)$$

مثال زیر نحوه محاسبه رگرسیون ژرفای یک خط را با استفاده از رابطه (۱) نشان می‌دهد.
مثال ۱: شکل ۲، ۷ مشاهده با خط ξ را نشان می‌دهد همانطور که دیده می‌شود این خط نابراز نیست. با استفاده از رابطه (۱) رگرسیون ژرفای این خط به صورت زیر محاسبه می‌شود.

$$i = 1: \min\{L^+(x_1) + R^-(x_1), R^+(x_1) + L^-(x_1)\} = \min\{1 + 3, 3 + 1\} = 4$$

$$i = 2: \min\{L^+(x_2) + R^-(x_2), R^+(x_2) + L^-(x_2)\} = \min\{1 + 2, 3 + 2\} = 5$$

$$i = 3: \min\{L^+(x_3) + R^-(x_3), R^+(x_3) + L^-(x_3)\} = \min\{2 + 2, 2 + 2\} = 4$$

$$i = 4: \min\{L^+(x_4) + R^-(x_4), R^+(x_4) + L^-(x_4)\} = \min\{2 + 1, 2 + 3\} = 3$$

$$i = 5: \min\{L^+(x_5) + R^-(x_5), R^+(x_5) + L^-(x_5)\} = \min\{3 + 1, 1 + 3\} = 4$$

$$i = 6: \min\{L^+(x_6) + R^-(x_6), R^+(x_6) + L^-(x_6)\} = \min\{3 + 0, 1 + 4\} = 3$$

$$i = 7: \min\{L^+(x_7) + R^-(x_7), R^+(x_7) + L^-(x_7)\} = \min\{4 + 0, 0 + 4\} = 4$$

بنابراین

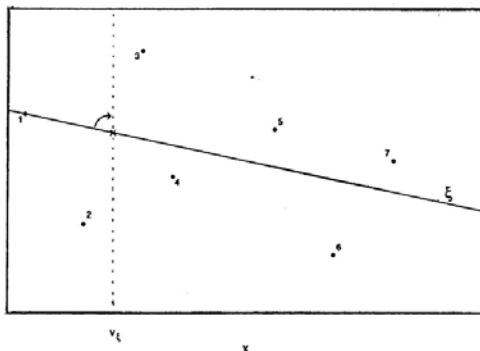
$$rdepth(\xi, Z_n) = \min(4, 5, 4, 3, 4, 3, 4) = 3$$

از روی شکل ۲ نیز دیده می‌شود که با برداشتن مشاهدات ۱ و ۴ و ۶ می‌توان خط ξ را بدون این که از مشاهدات دیگر عبور کند حول نقطه متناظر با v که با ضربدر روی آن خط مشخص شده است دوران داد تا به صورت عمودی در آید. به عبارت دیگر با برداشتن مشاهدات ۱ و ۴ و ۶ خط ξ ، نابراز می‌شود. رگرسیون ژرفای خط ξ برابر با کران پایین رگرسیون ژرفای ماکسیمال یعنی $\lceil \frac{7}{3} \rceil = \lceil \frac{7}{3} \rceil = 3$ است که در بخش‌های بعد بیان می‌شود.
توجه ۲: $\lceil \lambda \rceil$ کوچکترین عدد صحیح بزرگتر یا مساوی λ است.

۳ رگرسیون ژرفا در حالت چندگانه

در رگرسیون چندگانه مجموعه داده Z_n به صورت

$$Z_n = \{(x_{i1}, \dots, x_{i,p-1}, y_i); i = 1, \dots, n\} \subset \mathbb{R}^p$$



شکل ۲: مجموعه داده دو متغیره و برازش ξ با رگرسیون ژرفای ۳.

است. X را به عنوان قسمتی از مختصات هر نقطه به صورت

$$X_i = (x_{i,1}, \dots, x_{i,p-1}) \in \mathbb{R}^{p-1}$$

در نظر می‌گیریم. اکنون می‌خواهیم y_i را به وسیله

$$\theta_1 x_{i,1} + \dots + \theta_{p-1} x_{i,p-1} + \theta_p = (X_i, 1) \theta'$$

یعنی به وسیله یک ابرصفحه آفین در \mathbb{R}^p ، که $\theta = (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$ است، برازش دهیم. در اینجا هیچ فرض توزیعی وجود ندارد.

قبل از بیان تعاریف و مطالب این بخش ابتدا ابرصفحه و ابرصفحه آفین را تعریف می‌کنیم.

تعریف ۳: فرض کنید $[f = \alpha]$ نشان دهنده مجموعه سطح $\{x; f(x) = \alpha\}$ باشد. یک ابرصفحه مجموعه‌ای به شکل $[f = \alpha]$ است که f تابع خطی غیرصفر روی فضای برداری X و α عددی حقیقی است. به عبارت دیگر مجموعه $H \subseteq \mathbb{R}^n$ یک ابرصفحه است اگر اعداد حقیقی a_1, \dots, a_n و C (با $a_i \neq 0$ برای حداقل یک i) وجود داشته باشند به طوری که H متشکل از همه نقاط $x = (x_1, \dots, x_n)$ ایی باشد که در رابطه $\sum_i a_i x_i = C$ صدق می‌کنند.

تعریف ۴: یک ابرصفحه آفین در فضای برداری X مجموعه‌ای مانند \mathcal{M} است به طوری که برای هر x در \mathcal{M} ، $\mathcal{M} - x$ یک ابرصفحه باشد. به عبارت دیگر $\mathcal{M} \subseteq X$ یک ابرصفحه آفین است اگر یک تابع خطی غیرصفر $f: X \rightarrow \mathbb{F}$ و یک α در \mathbb{F} وجود داشته باشد به طوری که $\mathcal{M} = \{x \in X; f(x) = \alpha\}$ (میدان حقیقی \mathbb{R} ، یا میدان مختلط \mathbb{C} است).

مثال ۲: فضای برداری $X = \mathbb{R}^2$ و تابع $f(x, y) = y$ را در نظر می‌گیریم. در این صورت مجموعه $\{f = 0\} = \{(x, y); f(x, y) = 0\} = \{(x, 0); x \in \mathbb{R}\}$ ک ابرصفحه است که

فضای \mathbb{R}^2 را به دو نیم‌فضای بسته $\{(x, y); f(x, y) \geq 0\}$ و $\{(x, y); f(x, y) \leq 0\}$ افزایش می‌کند. مجموعه

$$C = \{(x, y); f(x, y) = 1\} = \{(x, 1); x \in \mathbb{R}\}$$

یک ابرصفحه آفین است زیرا بازای هر $(x_0, 1) \in C$

$$C - (x_0, 1) = \{(x - x_0, 0); x \in \mathbb{R}\} = \{(x - x_0, 0); f(x - x_0, 0) = 0\}$$

یک ابرصفحه است.

تعریف ۵: برازش $\theta = (\theta_1, \dots, \theta_p)$ یک نابرازا نامیده می‌شود اگر فقط اگر یک ابرصفحه آفین V در فضای X وجود داشته باشد به طوری که هیچ یک از x_i ها متعلق به V نباشند و برای همه x_i ها در یکی از نیم‌فضاهای باز آن داشته باشیم $r_i(\theta) > 0$ و همچنین برای همه x_i ها در نیم‌فضای باز دیگر $r_i(\theta) < 0$ باشد.

شکل ۳ مثالی از یک صفحه نابرازا است. در این مثال $p = 3$ است، بنابراین θ به یک صفحه متعلق است. فضای X به صورت صفحه افقی با $y \equiv 0$ که شامل خط V است، دیده می‌شود. (با این تعریف، هر η به طوری که همه مانده‌ها مثبت یا همه مانده‌ها منفی باشند یک نابرازا است، زیرا کافی است V به صورتی انتخاب شود که همه x_i ها در یک طرف آن واقع شوند). رگرسیون ژرفای هر $\theta \in \mathbb{R}^p$ متناسب با $Z_n \subset \mathbb{R}^p$ همانند تعریف ۲ به صورت زیر تعریف می‌شود:

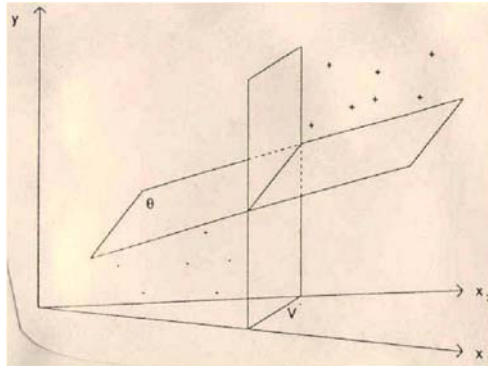
تعریف ۶: رگرسیون ژرفای یک برازش $\theta \in \mathbb{R}^p$ متناسب با مجموعه داده $Z_n \subset \mathbb{R}^p$ عبارت است از کوچکترین تعداد مشاهداتی که وقتی از مجموعه داده‌ها خارج شود θ ، نابرازا گردد. به عبارت دیگر $rdepth(\theta, Z_n)$ عبارت است از کمترین تعداد مانده‌هایی که باید تغییر علامت داده شوند تا θ نابرازا شود.

$rdepth(\theta, Z_n)$ همانند رابطه (۱) بیان می‌شود. بدین صورت که هر ابرصفحه آفین V در فضای X مشاهدات را به دو مجموعه تقسیم می‌کند که آن‌ها را با $L(V)$ و $R(V)$ نشان می‌دهیم. مجموعه‌های $L^+(V)$ و $L^-(V)$ به ترتیب تعداد مشاهدات در $L(V)$ با مانده‌های مثبت و تعداد مشاهدات در $L(V)$ با مانده‌های منفی است. $R(V)$ نیز بطور مشابه به دو مجموعه $R^+(V)$ و $R^-(V)$ تقسیم می‌گردد. بنابراین

$$rdepth(\theta, Z_n) = \min_V (\min\{L^+(V) + R^-(V), R^+(V) + L^-(V)\})$$

که V شامل همه ابرصفحه‌های آفین در فضای X است. برای محاسبه رگرسیون ژرفای یک برازش، جستجو را به مجموعه‌ای متناهی از ابرصفحه‌های V محدود می‌کنیم. الگوریتم‌های مناسب برای محاسبه رگرسیون ژرفا در فصل بعد ارائه می‌شود.

لازم به ذکر است که حالت تعمیم یافته قضیه ۱.۲ یعنی در حالت چندگانه (p بعدی) زمانی که $Z_n = \{(X_i, y_i); i = 1, \dots, n\}$ که $X_i = (x_{i1}, \dots, x_{i,p-1}) \in \mathbb{R}^{p-1}$ و



شکل ۳: مثالی از نابرازا $\theta \in \mathbb{R}^3$. ابرصفحه آفین V در فضای $X (\mathbb{R}^2)$ ، مشاهدات با مانده‌های مثبت را از مشاهدات با مانده‌های منفی جدا می‌کند.

$\theta = (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$ باشد نیز برقرار است.

قضیه ۱: خاصیت دقیق - برازش. اگر تعداد مشاهداتی که روی θ قرار گیرند k باشد (که $0 \leq k \leq n$) آنگاه

$$k \leq rdepth(\theta, Z_n) \leq \left\lfloor \frac{n+k}{2} \right\rfloor$$

برای $k = n$ داریم $rdepth(\theta, Z_n) = n$.

□
توضیح ۱: از آن جا که برآوردکننده L^1 همیشه از حداقل p مشاهده می‌گذرد، لذا رگرسیون ژرفای آن حداقل p است (بلموفیلد و استیگر^۴، ۱۹۸۳، را ببینید). همچنین این خاصیت برای برآوردکننده کمترین توان‌های دوم پیراسته (LTS) که به وسیله روسیف مطرح شده، برقرار است (روسیف ۱۹۸۴).
توضیح ۲: برآورد کننده کمترین توان‌های دوم (LS) هرگز نابرازا نیست. در حقیقت اگر Z_n به صورتی باشد که $X_n = \{x_1, x_2, \dots, x_n\}$ رتبه کامل داشته باشد و $n \geq 2p$ سپس $rdepth(\theta_{LS}, Z_n) \geq 1$.
 □

۴ رگرسیون ژرفای ماکسیمال

کران‌های بالا برای رگرسیون ژرفای ماکسیمال در قضیه زیر ارائه می‌شود. زیر مجموعه‌ای از \mathbb{R}^p در موقعیت عمومی قرار دارد هرگاه بیشتر از p مشاهده در هر زیرفضای آفین $(p-1)$ بعدی قرار نگیرد.

4) Bloomfield and Steiger

قضیه ۲: اگر (x_i, y_i) در موقعیت عمومی باشند، آنگاه

$$\max_{\theta} rdepth(\theta, Z_n) \leq \left\lfloor \frac{n+p}{2} \right\rfloor$$

□

حس ۱: برای هر مجموعه داده $Z_n \subset \mathbb{R}^p$ داریم

$$\max_{\theta} rdepth(\theta, Z_n) \geq \left\lfloor \frac{n}{p+1} \right\rfloor$$

۵ ژرفترین رگرسیون

ژرفترین رگرسیون T_r^* بعنوان θ ایی که $rdepth(\theta, Z_n)$ را بیشینه می‌کند، تعریف شده است. ژرفترین رگرسیون T_r^* "متعادل‌ترین" برازش است. T_r^* را می‌توان به وسیله محاسبه رگرسیون ژرفای همه برازش‌هایی که از میان p مشاهده می‌گذرد، به دست آورد.

تعریف ۷: برآوردکننده ژرفترین رگرسیون $T_r^*(Z_n)$ عبارت است از برازش θ با بزرگترین رگرسیون ژرفا، یعنی

$$T_r^*(Z_n) = \arg \max_{\theta} rdepth(\theta, Z_n)$$

توجه ۳: در رگرسیون چندگانه، ژرفترین رگرسیون برخلاف L^1 و LTS برای تبدیلات یکنوای y ، هم‌وردا است.

مثال ۳: می‌خواهیم رابطه بین میزان بخار مصرفی ماهانه توسط یک ماشین بخار (Y) با متوسط دمای ماهانه به فارنهایت (X_1) و تعداد روزهای کار در ماه (X_2) را مورد بررسی قرار می‌دهیم. بدین منظور ۲۵ مشاهده جمع‌آوری شده است که در جدول ۱ نشان داده شده‌اند (به دراپر و اسمیت^۵، ۱۹۹۸، مراجعه شود).

با استفاده از روش‌های رگرسیونی مختلف، ضرایب رگرسیونی را برآورد کرده و رگرسیون ژرفای هر یک از این برازش‌ها را به دست آورده‌ایم. نتایج در جدول ۲ آمده است. دیده می‌شود که برآوردکننده ژرفترین رگرسیون، رگرسیون ژرفای بالاتری نسبت به سایر برآوردکننده‌ها دارد و در نتیجه متعادل‌ترین برازش است.

۶ خواص بهینگی ژرفترین رگرسیون

در این بخش بعضی از خواص هم‌وردایی ژرفترین رگرسیون را بیان می‌کنیم. همچنین با معرفی کردن مقدار فروریزش به‌عنوان معیار پایداری برآوردکننده در برابر نقاط پرت، نشان می‌دهیم که روش ژرفترین رگرسیون با توجه به این معیار روشی استوار است.

5) Draper and Smith

جدول ۱: داده‌های ماشین بخار

| Y | X_2 | X_1 | Y | X_2 | X_1 |
|-------|-------|-------|-------|-------|-------|
| ۹,۵۷ | ۱۹ | ۳۹,۱ | ۱۰,۹۸ | ۲۰ | ۳۵,۳ |
| ۱۰,۹۴ | ۲۳ | ۴۶,۸ | ۱۱,۱۳ | ۲۰ | ۲۹,۷ |
| ۹,۵۸ | ۲۰ | ۴۸,۵ | ۱۲,۵۱ | ۲۳ | ۳۰,۸ |
| ۱۰,۰۹ | ۲۲ | ۵۹,۳ | ۸,۴۰ | ۲۰ | ۵۸,۸ |
| ۸,۱۱ | ۲۲ | ۷۰,۰ | ۹,۲۷ | ۲۱ | ۶۱,۴ |
| ۶,۸۳ | ۱۱ | ۷۰,۰ | ۸,۷۳ | ۲۲ | ۷۱,۳ |
| ۸,۸۸ | ۲۳ | ۷۴,۵ | ۶,۳۶ | ۱۱ | ۷۴,۴ |
| ۷,۶۸ | ۲۰ | ۷۲,۱ | ۸,۵۰ | ۲۳ | ۷۶,۷ |
| ۸,۴۷ | ۲۱ | ۵۸,۱ | ۷,۸۲ | ۲۱ | ۷۰,۷ |
| ۸,۸۶ | ۲۰ | ۴۴,۱ | ۹,۱۴ | ۲۰ | ۵۷,۵ |
| ۱۰,۳۶ | ۲۰ | ۳۳,۴ | ۸,۲۴ | ۲۰ | ۴۶,۴ |
| ۱۱,۸۸ | ۲۲ | ۲۸,۶ | ۱۲,۱۹ | ۲۱ | ۲۸,۹ |
| | | | ۱۱,۸۸ | ۲۱ | ۲۸,۱ |

جدول ۲: مقایسه رگرسیون ژرفای برآوردکننده‌های رگرسیونی با برآوردکننده ژرفترین رگرسیون

| روش | LS | LAD | برآورد هوپر | برآوردکننده ژرفترین رگرسیون |
|-----------------|---------|---------|-------------|-----------------------------|
| $\hat{\beta}_0$ | ۹,۲۱۳۰ | ۹,۷۳۷۰ | ۹,۳۸۲۰ | ۴,۰۹۰۰ |
| $\hat{\beta}_1$ | -۰,۰۷۴۸ | -۰,۰۷۳۸ | -۰,۰۷۶۱ | -۰,۰۸۲۴ |
| $\hat{\beta}_2$ | ۰,۲۰۶۰ | ۰,۱۹۲۰ | ۰,۲۰۴۰ | ۰,۴۷۶۰ |
| رگرسیون ژرفا | ۵ | ۵ | ۶ | ۱۱ |

۱.۶ هم‌وردایی ژرفترین رگرسیون

بررسی خواص فوق در قضیه زیر آمده است:
قضیه ۳.۴: مجموعه داده $Z_n = \{(X_i, Y_i); i = 1, \dots, n\}$ که در آن

$$X_i = (x_{i,1}, \dots, x_{i,p-1}) \in \mathbb{R}^{p-1}$$

و برآوردکننده ژرفترین رگرسیون $T_r^*(Z_n) = \arg \max_{\theta} rdepth(\theta, Z_n)$ را در نظر می‌گیریم:
(a) برآوردکننده ژرفترین رگرسیون، رگرسیون هم‌وردا است، یعنی برای هر بردار ستونی V

$$T_r^*(Z'_n) = T_r^*(Z_n) + V$$

که $T_r^*(Z'_n) = \arg \max_{\theta} rdepth(\theta, Z'_n)$ و $Z'_n = \{(X_i, y_i + X_i V); i = 1, \dots, n\}$
(b) برآوردکننده ژرفترین رگرسیون، مقیاس هم‌وردا است، یعنی برای هر c ثابت

$$T_r^*(Z'_n) = c T_r^*(Z_n)$$

که $T_r^*(Z'_n) = \arg \max_{\theta} rdepth(\theta, Z'_n)$ و $Z'_n = \{(X_i, cy_i); i = 1, \dots, n\}$ **(c)**

برآوردکننده ژرفترین رگرسیون، آفین هم‌وردا است، یعنی برای هر ماتریس مربع ناویژه A

$$T_r^*(Z'_n) = A^{-1} T_r^*(Z_n)$$

که $T_r^*(Z'_n) = \arg \max_{\theta} rdepth(\theta, Z'_n)$ و $Z'_n = \{(X_i A, y_i); i = 1, \dots, n\}$

۲.۶ استواری ژرفترین رگرسیون

یک معیار معروف برای پایداری برآورد کننده در برابر نقاط پرت، مقدار فروریزش است. بنا به تعریف داناوو و هوبر^۶ (۱۹۸۳)، مقدار فروریزش هر برآوردکننده T_n به صورت زیر است:

$$\varepsilon_n^*(T_n, Z_n) = \min \left\{ \frac{k}{n}; \sup_{Z'_n} \|T_n(Z'_n) - T_n(Z_n)\| = \infty \right\}$$

که Z'_n شامل همهٔ مجموعه داده‌های بدست آمده با جایگزینی هر k مشاهده دلخواه Z_n با مقادیر اختیاری است. بنابراین مقدار فروریزش، کوچکترین کسری از مشاهدات تبدیل شده (مغشوش شده) ای است که می‌تواند برآوردکننده را به دلخواه منحرف کند. لازم به ذکر است که

6) Donoho and Huber

این اغتشاش فقط به نقاط پرت در y_i محدود نمی‌شود، بلکه Z'_n می‌تواند شامل نقاط پرت در x_i ها نیز باشد.

حدس ۲: برای هر مجموعه داده $Z_n \subset \mathbb{R}^p$ با x_i ها در موقعیت عمومی باشند آنگاه

$$\varepsilon_n^*(T_r^*, Z_n) \geq \frac{1}{n} \left(\left\lceil \frac{n}{p+1} \right\rceil - p + 1 \right) \approx \frac{1}{p+1} \quad (2)$$

رابطه (۲) بیان می‌کند که مقدار فروریزش ژرفترین رگرسیون همواره مثبت است. این مقدار، زمانی برابر با کران پایین $\frac{1}{p+1}$ می‌شود که داده‌های اصلی خود، ویژه باشند مثلاً وقتی آن‌ها روی منحنی گشتاور قرار گیرند.

بنابراین ژرفترین رگرسیون در برابر نقاط نافذ همانند نقاط دورافتاده عمودی، استوار است. بعلاوه ژرفترین رگرسیون متفاوت از رگرسیون L^1 تعریف شده به صورت

$$L^1(Z_n) = \arg \min_{\theta} \sum_{i=1}^n |r_i(\theta)|$$

است که در نتیجه آسیب‌پذیری آن در برابر نقاط نافذ می‌باشد.

۷ مثال‌ها

در این بخش برآوردکننده ژرفترین رگرسیون را برای داده‌های موجود بدست آورده و نشان می‌دهیم که در مقایسه با برآوردکننده‌های دیگر دارای بیشترین رگرسیون ژرفاست. این مقایسه با محاسبه رگرسیون ژرفای برآوردکننده‌ها انجام می‌گیرد. بدین منظور از الگوریتم‌هایی که برنامه مربوط به آنها با استفاده از نرم‌افزار Splus و زبان برنامه‌نویسی فرترن نوشته شده است و از وب سایت <http://win-www.uia.ac.be/u/statis/> قابل دسترسی است، استفاده می‌گردد.

۱.۷ مثال ۱

در یک تحقیق در سال ۱۹۷۶ میزان آلودگی Y رودخانه در ایالت نیویورک بررسی شده است. هدف، بررسی سهم هر نوع از زمین‌های مورد استفاده در اطراف رودخانه بر روی آلودگی آب رودخانه‌هاست. این آلودگی براساس میانگین غلظت نیترژن (میلی‌گرم) اندازه‌گیری می‌شود. نوع زمین‌های مورد استفاده در اطراف رودخانه‌ها و میزان آلودگی آن‌ها در زیر بیان شده است. داده‌ها در جدول ۳ نشان داده شده‌اند (به چاترجی و همکاران^۷، ۲۰۰۰ مراجعه شود).

Y -- میانگین غلظت نیترژن (میلی‌گرم بر لیتر) براساس نمونه‌هایی که در فواصل منظم بهار، تابستان و ماه‌های پاییز گرفته شده است.

7) Chatterjee and et al.

جدول ۳: داده‌های رودخانه‌های نیویورک

| Y | X_4 | X_3 | X_2 | X_1 | Y | X_4 | X_3 | X_2 | X_1 |
|-------|-------|-------|-------|-------|------|-------|-------|-------|-------|
| ۱,۰۰۱ | ۰,۰۹ | ۰,۷ | ۵۷ | ۲۹ | ۱,۱۵ | ۰,۲۹ | ۱,۲ | ۶۳ | ۲۶ |
| ۱,۰۰ | ۱,۹۸ | ۱,۹ | ۸۴ | ۲ | ۱,۹۰ | ۰,۵۸ | ۱,۸ | ۲۶ | ۵۴ |
| ۱,۴۲ | ۰,۵۶ | ۳,۴ | ۶۱ | ۱۹ | ۱,۹۹ | ۳,۱۱ | ۲۹,۴ | ۲۷ | ۳ |
| ۱,۶۵ | ۰,۲۴ | ۱,۳ | ۴۳ | ۴۰ | ۲,۰۴ | ۱,۱۱ | ۵,۶ | ۶۰ | ۱۶ |
| ۱,۲۱ | ۰,۲۳ | ۰,۹ | ۶۰ | ۲۶ | ۱,۰۱ | ۰,۱۵ | ۱,۱ | ۶۲ | ۲۸ |
| ۰,۷۵ | ۰,۱۶ | ۰,۷ | ۷۵ | ۱۵ | ۱,۳۳ | ۰,۱۸ | ۰,۹ | ۵۳ | ۲۶ |
| ۰,۸۰ | ۰,۳۵ | ۰,۸ | ۸۱ | ۳ | ۰,۷۳ | ۰,۱۲ | ۰,۵ | ۸۴ | ۶ |
| ۰,۸۷ | ۰,۱۵ | ۰,۵ | ۸۲ | ۶ | ۰,۷۶ | ۰,۳۵ | ۰,۷ | ۸۹ | ۲ |
| ۰,۸۷ | ۰,۱۸ | ۰,۴ | ۷۵ | ۴ | ۰,۸۰ | ۰,۲۲ | ۰,۹ | ۷۰ | ۲۲ |
| ۱,۲۵ | ۰,۱۳ | ۱,۱ | ۴۹ | ۴۰ | ۰,۶۶ | ۰,۱۳ | ۰,۵ | ۵۶ | ۲۱ |

جدول ۴: برآورد پارامترهای مدل و رگرسیون ژرفای برآوردکننده‌های مختلف رگرسیونی

| روش | LS | LAD | برآورد هوپر | برآوردکننده ژرفترین رگرسیون |
|--------------|---------|---------|-------------|-----------------------------|
| β_0 | ۱,۷۲۲۲ | ۲,۷۴۰۰ | ۲,۶۵۶۰ | ۲,۹۹۰۰ |
| β_1 | ۰,۰۰۵۸ | -۰,۰۰۵۲ | -۰,۰۰۵۳ | -۰,۰۱۰۶ |
| β_2 | -۰,۰۱۳۰ | -۰,۰۲۴۳ | -۰,۰۲۳۲ | -۰,۰۲۸۹ |
| β_3 | -۰,۰۰۷۲ | -۰,۰۲۱۶ | -۰,۰۲۲۳ | ۰,۰۴۵۶ |
| β_4 | ۰,۳۰۵۰ | ۰,۱۷۸۹ | ۰,۲۱۶۰ | ۰,۴۴۰۰ |
| رگرسیون ژرفا | ۳ | ۴ | ۴ | ۹ |
| $\sum e_i^2$ | ۱,۰۵۲۷ | ۱,۱۸۰۰ | ۱,۱۲۷۰ | ۹,۱۷۶۰ |
| $\sum e_i $ | ۳,۲۱۴۶ | ۲,۸۱۸۰ | ۲,۹۲۴۰ | ۵,۳۵۷۷ |

X_1 -- کشاورزی: درصد مساحت زمین مورد استفاده برای کشاورزی

X_2 -- جنگلی: درصد زمین جنگلی

X_3 -- مسکونی: درصد مساحت زمین مورد استفاده برای سکونت

X_4 -- تجاری/صنعتی: درصد مساحت زمین مورد استفاده تجاری یا صنعتی.

با بررسی داده‌های فوق دیده می‌شود که مشاهدات ۴ و ۵ دور افتاده و نافذ و در عین حال مؤثر هستند. بنابراین انتظار می‌رود که روش‌های LS ، LAD و برآورد هوپر تحت تأثیر این نقاط قرار گیرند و بطور قابل ملاحظه رگرسیون ژرفای کوچکتری نسبت به برآوردکننده ژرفترین رگرسیون داشته باشند. با نگاه کردن به جدول ۴ مشخص می‌شود که برآوردکننده رگرسیونی LAD و برآورد هوپر نزدیک به هم هستند و دارای رگرسیون ژرفای ۴ هستند. چون برآوردکننده ژرفترین رگرسیون نسبت به نقاط دورافتاده و نافذ استوار است لذا تحت تأثیر این نقاط قرار نمی‌گیرد و به نظر می‌رسد که برازشی مناسب برای داده‌هاست. با توجه به اینکه دارای بیشترین رگرسیون ژرفا یعنی ۹ است.

جدول ۵: جدول ۵: برآورد پارامترهای مدل و رگرسیون ژرفای برآوردکننده های مختلف رگرسیونی

| روش | LS | LAD | برآورد هویر | برآوردکننده ژرفترین رگرسیون |
|-----------------|----------|----------|-------------|-----------------------------|
| $\hat{\beta}_0$ | -۹,۶۹۵۵ | -۹,۴۹۷۴ | -۹,۶۶۳۳ | -۹,۹۵۰ |
| $\hat{\beta}_1$ | ۱,۱۲۰۸ | ۱,۱۱۴۵ | ۱,۱۳۷۲ | ۱,۱۳۰ |
| $\hat{\beta}_2$ | ۱,۹۶۴۳ | ۱,۸۸۵۹ | ۱,۹۹۸۲ | ۲,۰۱۰ |
| $\hat{\beta}_3$ | -۴۹,۹۱۱۲ | -۴۹,۸۵۱۳ | -۴۹,۸۸۵۹ | -۵۰,۱۴۰ |
| $\hat{\beta}_4$ | ۰,۵۲۸۸ | ۰,۴۷۶۲ | ۰,۴۸۷۲ | ۰,۴۹۶ |
| رگرسیون ژرفا | ۱۳ | ۱۰ | ۱۴ | ۱۸ |
| $\sum e_i^2$ | ۳۶,۴۷۲۰ | ۳۸,۶۳۱۰ | ۳۶,۶۹۸۰ | ۳۹,۹۴۱ |
| $\sum e_i $ | ۳۱,۳۲۷۰ | ۲۹,۸۸۹۰ | ۳۰,۸۹۱۰ | ۳۱,۷۱۸ |

۲.۷ مثال ۲- شبیه سازی

برای بررسی مناسب بودن برآوردکننده ژرفترین رگرسیون، ۴۰ نقطه در ۵ بعد مطابق با مدل $y = -10 + x_1 + 2x_2 - 50x_3 + 0.5x_4 + e$ تولید می کنیم که x_1, x_2, x_3, x_4 و e از توزیع نرمال استاندارد گرفته شده اند. پس از محاسبه y با استفاده از مقادیر تولید شده x_1, x_2, x_3, x_4 و e برآوردکننده های رگرسیونی بدست آمده از روش های مختلف و برآوردکننده ژرفترین رگرسیون را برای این مقادیر بدست می آوریم. نتایج بدست آمده در جدول زیر آمده اند.

با مقایسه ضرایب برآورد شده بدست آمده از روش های مختلف دیده می شود که ضرایب بدست آمده از روش ژرفترین رگرسیون نسبت به سایر روش ها به ضرایب مدل مفروض نزدیک تر است و در عین حال دارای بیشترین رگرسیون ژرفا است. به عبارت دیگر برآوردکننده ژرفترین رگرسیون

$$y = -9.95 + 1.13x_1 + 2.01x_2 - 50.14x_3 + 0.496x_4$$

دارای رگرسیون ژرفای تقریبی ۱۸ است.

۸ نتیجه گیری

وقتی که استواری و تشخیص نقاط دورافتاده اهمیت دارند روش LMS روشی مناسب است. اما زمانی که نقاط دورافتاده زیادی وجود ندارد (می توان به وسیله روش های LMS یا LTS بررسی کرد) و یکنوایی اهمیت دارد روش ژرفترین رگرسیون T_n^* انتخابی مناسب است.

مراجع

- [1] Bloomfield, P., and Steiger, W. (1983). Least Absolute Deviations: Theory, Applications, and Algorithms, Boston: Birkhauser.
- [2] Chatterjee, S., and Hadi, A.S., and Price, B. (2000), Regression Analysis by Example, New York: John Wiley.
- [3] Donoho, D.L., and Huber, P.J. (1983), "The notion of breakdown point", in A Festschrift for Erich Lehmann, eds. Bickel, P., Doksum, K., and Hodges, J.L., Belmont: Wadsworth.
- [4] Draper, N., Smith, H. (1998) Applied Regression Analysis, third edition, John Wiley & Sons.
- [5] Rousseeuw, P.J., (1984) "Least Median of Squares Regression", Journal of the American Statistical Association, **79**, 871-880.
- [6] Rousseeuw, P.J., and Hubert, M. (1999), "Regression Depth", Journal of the American Statistical Association, **94**, 388-402.

استفاده از معیارهای کولبک - لیبلر و فاصله چرنوف برای آنالیز ممیزی سری‌های زمانی چند متغیره

سارا شفيعی بابایی^۱، رحيم چيني پرداز^۲

^۱ دانشگاه شهيد چمران، دانشکده علوم ریاضی و کامپیوتر، گروه آمار

^۲ دانشگاه شهيد چمران، دانشکده علوم ریاضی و کامپیوتر، گروه آمار

چکیده: هدف از آنالیز ممیزی خطی در سری‌های زمانی بدست آوردن تابعی خطی می‌باشد که مطابق با آن تابع یک مشاهده سری زمانی را بتوان به یکی از دو یا چند جامعه مستقل سری زمانی تخصیص داد. روش‌های کلاسیک برای بدست آوردن تابع ممیزی مبتنی بر ماکزیمم کردن فاصله ماهالونوبیس است. در سالهای اخیر فاصله‌های کولبک - لیبلر و چرنوف نیز بوسیله آماردانان مورد توجه قرار گرفته است. در این مقاله فرض شده است که جوامع سری زمانی به صورت مستقل و دارای توزیع نرمال هستند. این جوامع ممکن است دارای میانگین‌های برابر و ماتریس‌های کواریانس برابر و یا متفاوت از هم باشند. تابع ممیزی خطی با استفاده از ماکزیمم کردن فاصله کولبک - لیبلر و چرنوف بدست آمده است. با توجه به اینکه استفاده از حوزه زمانی این تابع منجر به روش‌های پیچیده و محاسبات مشکل می‌شود از تقریب طیفی سری‌های زمانی در حوزه فرکانس استفاده شده است. سپس تابع ممیزی خطی با همان روش ولی برای سری‌های زمانی چند متغیره گسترش داده شده است. در انتها با روش شبیه‌سازی برای مدل‌های $ARMA(p, q)$ نتایج حاصل مورد بررسی و مقایسه قرار داده شده‌اند.

واژه‌های کلیدی: آنالیز ممیزی، تابع ممیزی خطی، آنالیز طیفی سری‌های زمانی، معیار ممیزی کولبک - لیبلر، معیار ممیزی چرنوف

۱ مقدمه

مساله آنالیز ممیزی و رده‌بندی از جمله مسایل پر اهمیت آماری است که به شکل کاربردهای مختلف و در بسیاری از علوم اهمیت خاصی پیدا کرده است. یکی از مهمترین و اساسی‌ترین کاربردهای آنالیز ممیزی در روش‌های کلاسیک تشخیص الگو در داده‌های آزمایشی سری‌های زمانی است که در موارد مختلف قابل توجه می‌باشد. به عنوان مثال گوینز و دیگران (۱۹۷۵) کاربردهای آنالیز ممیزی را در داده‌های دریافت شده از دستگاه نوار مغز بررسی کردند. در این زمینه همچنین می‌توان به گرش و دیگران (۱۹۷۹) آلاگان (۱۹۸۶، ۱۹۸۹) راجی (۱۹۸۵) اشاره نمود. در زمینه کاربرد توابع ممیزی در داده‌های مربوط

به گویش‌های مختلف گزارش شده می‌توان به ولف (۱۹۷۶) مارکل و دیگران (۱۹۷۶) و دویت (۱۹۸۵) اشاره کرد.

در زمین‌شناسی نیز مساله آنالیز ممیزی بین داده‌های لرزه‌ای مربوط به دو منبع مختلف ایجاد موج؛ یعنی زمین لرزه و انفجار هسته‌ای بسیار مورد توجه واقع بوده است. چنانچه در این راستا می‌توان به تی جی استیم (۱۹۷۶) درگاهی - نوبری و لیکاک (۱۹۸۱) درگاهی - نوبری (۱۹۹۵) شاموی (۱۹۹۶؛ ۱۹۸۸) بلند فورد (۱۹۹۳) کاکیزاوا و دیگران (۱۹۹۸) اشاره نمود. پیکالو (۱۹۹۰) نیز از آنالیز ممیزی در داده‌های اقتصادی استفاده کرده است. همچنین چینی‌پرداز (۲۰۰۴) با استفاده از بردارهای اتوکوریانس تابع ممیزی را برای داده‌های مربوط به زمین لرزه و انفجار هسته‌ای پیشنهاد کرده است.

هدف این مقاله؛ تعیین و بدست آوردن توابع ممیزی خطی در دامنه فرکانسی سری‌های زمانی است. توابع ممیزی در حالت یک متغیره و چند متغیره با استفاده از معیار اطلاع ممیزی کولپک - لیبار و معیار اطلاع ممیزی چرنوف بدست آورده می‌شوند و سپس تقریب‌های طیفی مناسب برای این توابع پیشنهاد می‌گردند. در پایان با یک مثال عددی نتایج استفاده از معیارهای ممیزی مقایسه می‌شوند.

۱.۱ توابع ممیزی در دامنه فرکانس

اگر مشاهدات سری زمانی به صورت $x_t = (x_0, x_1, \dots, x_{T-1})$ تحت دو فرضیه H_j برای $j = 1, 2$ باشد به طوریکه:

$$H_j : x \sim N_T(\circ, R_j)$$

در آن صورت تابع ممیزی بر اساس نسبت درستنمایی چگالی دو جامعه به صورت زیر خواهد بود:

$$d_Q(x) = x'(R_2^{-1} - R_1^{-1})x$$

و تخصیص مشاهده x به یکی از دو جامعه به این صورت می‌باشد که هرگاه $d_Q(x) \geq 0$ مشاهده x به جامعه اول تخصیص داده می‌شود و هرگاه $d_Q(x) < 0$ مشاهده x به جامعه دوم تخصیص داده می‌شود. اما در این حالت توزیع $d_Q(x)$ تحت هر یک از فرضیه‌ها بسیار پیچیده می‌باشد؛ بنابراین با استفاده از تبدیلات فوریه گسسته بحث آنالیز ممیزی به شکل بهتر و ساده‌تری در دامنه فرکانس دنبال می‌شود. در این حالت متغیرهای تولید شده در دامنه فرکانس؛ متغیرهای تقریباً نرمال و ناهمبسته با ماتریس‌های کواریانس قطری می‌باشند؛ بطوریکه بین مشاهدات پشت سر هم همبستگی ناچیزی وجود دارد.

فرض کنید $X_t = (X(0), X(1), \dots, X(T-1))'$ بردار تبدیل فوریه متناهی از مشاهدات x_t در دامنه فرکانس باشد؛ که به صورت زیر تعریف می‌شود:

$$X(k) = \frac{1}{\sqrt{T}} \sum_{t=0}^{T-1} x_t \exp\{i\lambda_k t\} \quad \lambda_k = 2\pi k T^{-1}$$

با بردار میانگین:

$$M(k) = \frac{1}{\sqrt{T}} \sum_{k=0}^{T-1} \mu_j(t) \exp \{i \lambda_k t\} \quad \lambda_k = 2\pi k T^{-1}$$

و تابع چگالی طیفی، یا توان طیف برابر با:

$$f_x(\lambda_k) = \sum_{k=0}^{T-1} R_k \exp \{i \lambda_k t\} \quad \lambda_k = 2\pi k T^{-1}$$

باشد. در حقیقت دو رابطه اخیر تبدیل فوریه متناهی از میانگین $\mu_j(t)$ و ماتریس کواریانس R_j برای $j = 1, 2$ در دامنه فرکانس می‌باشد.

لگاریتم تقریبی تابع چگالی توأم تمام مقادیر تبدیل فوریه گسسته به صورت زیر خواهد بود:

$$\ln p_j(x) = \frac{1}{2} \sum_{k=0}^{T-1} \ln f_j(\lambda_k) - \frac{1}{2} \sum_{k=0}^{T-1} \frac{|X(k) - M_j(k)|^2}{f_j(\lambda_k)}$$

اگر فرض شود که میانگین‌های جامعه‌ها نابرابر هستند؛ آنگاه تحت این فرض که $f_j(\lambda_k) = f(\lambda_k)$ برای $j = 1, 2$ تابع ممیزی خطی زیر تابع ممیزی بین دو جامعه خواهد بود:

$$d_L(x) = \mu'_j R^{\circ-1} x - \frac{1}{2} \mu'_j R^{\circ-1} \mu_j$$

$$= \sum_{k=0}^{T-1} \frac{\overline{M_j(k)X(k)}}{f(\lambda_k)} - \frac{1}{2} \sum_{k=0}^{T-1} \frac{|M_j(k)|^2}{f(\lambda_k)}$$

که در آن $R_j^{\circ} = \{r_j^{\circ}(s-t), t, s = 0, 1, \dots, T-1\}$ و $R^{\circ-1}$ دارای عناصری به فرم

$$r^{\circ-1}(s-t) = \frac{1}{\sqrt{T}} \sum_{k=0}^{T-1} \exp \{i \lambda_k (s-t)\} f(\lambda_k)$$

می‌باشند. [لیگت (۱۹۷۱)؛ شاموی و آنگر (۱۹۷۴) و شاموی (۱۹۸۲)].
به علاوه زمانیکه توابع کواریانس (مقادیر طیف) دو جامعه نابرابر باشند نیز؛ تابع رده‌بندی درجه دوم به صورت زیر بدست می‌آید:

$$d_Q^{\circ}(x) = x'(R_{\Psi}^{\circ-1} - R_{\Psi}^{\circ-1})x$$

$$= \sum_{k=0}^{T-1} |X(k)|^2 \left((f_{\gamma}^{-1}(\lambda_k) - f_{\gamma}^{-1}(\lambda_k)) \right)$$

و بر اساس اینکه $d_Q^\circ(x) \geq 0$ یا $d_Q^\circ(x) < 0$ باشد؛ سری زمانی x_t به جامعه اول و یا دوم تخصیص داده می‌شود.

۲ آنالیز ممیزی بر اساس معیارهای جداسازی

در این بخش با بکارگیری تبدیلات فوریه گسسته و تقریب‌های طیفی؛ آنالیز ممیزی را در دامنه فرکانسی سری‌های زمانی بر اساس معیارهای ممیزی دنبال می‌کنیم.

۱.۲ آنالیز ممیزی بر اساس معیارهای ممیزی کولبک - لیبلر و چرنوف در حالت یک متغیره

فرض کنید مساله رده‌بندی؛ تخصیص یک فرایند ایستای نرمال $T \times 1$ بعدی $x_t = (x_1, x_2, \dots, x_t)'$ با بردار میانگین $\mu_j = (\mu_{j0}, \mu_{j1}, \dots, \mu_{jT-1})'$ و ماتریس کواریانس $T \times T$ بعدی $R_j = \{\sigma_{j(s-t)}, t, s = 0, 1, \dots, T-1\}$ تحت فرضیه‌های H_j برای $j = 1, 2$ به یکی از دو جامعه نرمال باشد. یکی از معیارهای کلاسیک جداسازی بین دو جامعه چند متغیره؛ معیار کولبک - لیبلر می‌باشد که به صورت زیر تعریف می‌شود: کولبک - لیبلر (۱۹۵۲)؛ کولبک (۱۹۷۸).

$$I(1, 2; x) = \frac{1}{T} E_1 \left\{ \log \frac{p_1(x)}{p_2(x)} \right\} \quad (1)$$

که در آن $p_1(x)$ و $p_2(x)$ به ترتیب توابع چگالی مربوط به جامعه‌های نرمال تحت فرضیه‌های H_j برای $j = 1, 2$ در نظر گرفته می‌شود و T طول بردار مشاهدات سری زمانی است. به سادگی می‌توان رابطه (۱) را به صورت رابطه زیر بدست آورد:

$$I(1, 2; x) = \frac{1}{T} \left(\left(\text{tr}(R_1 R_2^{-1}) \log \frac{|R_2|}{|R_1|} - T \right) + \delta' R_2^{-1} \delta \right) \quad (2)$$

که در آن $\delta_t = \mu_{2t} - \mu_{1t}$ برای $t = 0, 1, \dots, T-1$ می‌باشد. به دلیل نامتقارن بودن این معیار؛ با تعریف رابطه زیر مقدار متقارنی از معیار کولبک - لیبلر حاصل می‌شود:

$$J(1, 2; x) = I(1, 2; x) + I(2, 1; x) \quad (3)$$

همچنین معیار جداسازی دیگری که بین دو جامعه چند متغیره به عنوان یک معیار ممیزی بکار برده می‌شود؛ معیار چرنوف [چرنوف (۱۹۵۲)؛ رینی (۱۹۶۱)] می‌باشد که پارزن (۱۹۹۰) آن را در ممیزی بین سری‌های زمانی پیشنهاد نموده و به صورت زیر تعریف می‌گردد.

$$Q_r(1, 2; x) = -\frac{1}{2T} \log E_1 \left\{ \log \left(\frac{p_1(x)}{p_2(x)} \right)^r \right\} \quad 0 < r < 1$$

$$= \frac{1}{2T} \left(\log \frac{|rR_1 + (1-r)R_2|}{|R_2|} - r \log \frac{|R_1|}{|R_2|} \right) \quad (4)$$

با استفاده از این خاصیت که $Q_r(1, 2; x) = Q_{1-r}(2, 1; x)$ می‌توان معیار متقارنی از $Q_r(1, 2; x)$ به صورت زیر تعریف نمود:

$$J_{Q_r}(1, 2; x) = Q_r(1, 2; x) + Q_r(2, 1; x) \quad (5)$$

با توجه به اینکه در دامنه زمان استفاده از این معیارها مشکل می‌باشد؛ لذا برای ممیزی بین دو جامعه نرمال از روش‌های طیفی استفاده می‌گردد. بدین ترتیب فرض کنید سری زمانی x_t تحت فرضیه‌های H_j برای $j = 1, 2$ دارای میانگین μ_j و ماتریس کواریانس $T \times T$ بعدی $R_j = \{\sigma_{j(s-t)}, t, s = 0, 1, \dots, T-1\}$ باشد. آنگاه رابطه معکوس تبدیل فوریه زیر برای تابع کواریانس وجود دارد و به شکل زیر است:

$$\sigma_j(s-t) = \int_{-\pi}^{\pi} f_j(\lambda) e^{i\lambda(s-t)} \left(\frac{d\lambda}{2\pi} \right)$$

که در آن $f_j(\lambda)$ طیف دو جامعه می‌باشد و فرض می‌شود که بر روی فاصله $[-\pi, \pi]$ مثبت؛ پیوسته و به طور مطلق انتگرال پذیر است. فولر (۱۹۹۶).
در این صورت فرض می‌شود دنباله $\delta_t = \mu_2 t - \mu_1 t$ برای $t = 0, 1, \dots, T-1$ در شرایط زیر صدق کند:

$$(I) \quad \sup_t |\delta_t| = c < \infty$$

$$(II) \quad \rho_T(\tau) = \frac{1}{T} \sum_{t=0}^{T-1-|\tau|} \delta_{t+|\tau|} \delta_t$$

که $\rho_t(\tau)$ تابع خود همبستگی سری زمانی x_t می‌باشد و دارای حد داده شده در رابطه زیر است:

$$\rho(\tau) = \lim_{T \rightarrow \infty} \rho_T(\tau) = \int_{-\pi}^{\pi} e^{i\lambda \tau} \left(\frac{dM\lambda}{2\pi} \right)$$

که $M(\lambda)$ تابع یکنوای اکیداً غیر نزولی است. به گونه‌ای که $M(-\pi)$ و از راست پیوستگی دارد. همچنین $\rho(\circ) > \circ$. فولر (۱۹۹۶).

iii) اگر $x_t = (x_1, x_2, \dots, x_t)$ فرایند سری زمانی ایستای نرمال با میانگین μ_j و ماتریس کواریانس R_{jz} تحت فرضیه‌های H_j برای $j = 1, 2$ باشد؛ آنگاه اگر تابع کواریانس $\sigma_j(s-t)$ دارای تابع چگالی مطلقاً انتگرال‌پذیر و مثبت و پیوسته $f_j(\lambda)$ در فاصله $[-\pi, \pi]$ باشد. همچنین تابع δ_t در شرایط (i) و (ii) صدق کند؛ در این صورت حد اطلاع ممیزی رابطه (۲) و (۴) به صورت زیر داده می‌شود: [شاموی و آنگر (۱۹۷۴)؛ کاکیزاوا و دیگران (۱۹۹۸)].

$$I(1, 2; x) = \frac{1}{2T} \int_{-\pi}^{\pi} \left(\frac{f_1(\lambda)}{f_2(\lambda)} - \log \frac{f_1(\lambda)}{f_2(\lambda)} - 1 \right) \frac{d\lambda}{2\pi} + \frac{1}{2} \int_{-\pi}^{\pi} \frac{1}{f_2(\lambda)} \frac{dM\lambda}{2\pi}$$

و اگر میانگین‌های دو جامعه برابر با صفر در نظر گرفته شوند؛ قسمت دوم از سمت راست رابطه بالا برابر با صفر خواهد بود.

لذا:

$$I(1, 2; x) = \frac{1}{2T} \int_{-\pi}^{\pi} \left(\frac{f_1(\lambda)}{f_2(\lambda)} - \log \frac{f_1(\lambda)}{f_2(\lambda)} - 1 \right) \frac{d\lambda}{2\pi} \quad (۶)$$

همچنین:

$$Q_r(1, 2; x) = \frac{1}{2T} \int_{-\pi}^{\pi} \left(\log \frac{|rf_1(\lambda) + (1-r)f_2(\lambda)|}{|f_2(\lambda)|} - r \log \frac{|f_1(\lambda)|}{|f_2(\lambda)|} \right) \frac{d\lambda}{2\pi} \quad (۷)$$

در نتیجه:

$$I(1, 2; x) = \lim_{T \rightarrow \infty} I_T(1, 2, x)$$

و به همین صورت:

$$I(2, 1; x) = \lim_{T \rightarrow \infty} I_T(2, 1, x)$$

و بنابراین:

$$J(1, 2; x) = \lim_{T \rightarrow \infty} J_T(1, 2, x)$$

بدست می‌آید.

بنابراین حد اطلاع ممیزی که به عنوان تقریب طیفی برای معیارهای ممیزی کولبک - لیبار و چرنوف در نظر گرفته می‌شود؛ به معیارهای اصلی ممیزی همگرا می‌باشد. یعنی اینکه صورت‌های حدی برای تقریب معیارهای ممیزی و معیار واگرایی J با فرم‌های معادل برای فرایند اصلی x_t مشابه می‌باشند.

پس می‌توان نوشت: [شاموی و آنگر (۱۹۷۴)؛ کاکیزاوا و دیگران (۱۹۹۸)؛ کازاکاس - پاپ‌آنتونی کازاکس (۱۹۸۰)].

$$I(1, 2; x^\circ) = \frac{1}{2T} \sum_{k=0}^{T-1} \left[\text{tr} \left(f_1(\lambda_k) f_2^{-1}(\lambda_k) \right) - \log \frac{f_1(\lambda_k)}{f_2(\lambda_k)} - 1 \right] \quad (۸)$$

و همچنین:

$$Q_r(\lambda, \nu; x^\circ) = \frac{1}{\nu T} \sum_{k=0}^{T-1} \left(\log \frac{|rf_\lambda(\lambda_k) + (1-r)f_\nu(\lambda_k)|}{|f_\nu(\lambda_k)|} - r \log \frac{|f_\lambda(\lambda_k)|}{|f_\nu(\lambda_k)|} \right) \quad (9)$$

ملاحظه می‌شود که این روابط؛ حتی برای مقادیر بزرگ T نیز به سادگی قابل محاسبه می‌باشد. بنابراین:

$$\lim_{T \rightarrow \infty} I_T(\lambda, \nu, x^\circ) = I(\lambda, \nu; x)$$

که در آن $I(\lambda, \nu; x)$ در رابطه (۶) داده شده است. همچنین:

$$\lim_{T \rightarrow \infty} Q_{rT}(\lambda, \nu, x^\circ) = Q_r(\lambda, \nu; x)$$

به همین ترتیب می‌توان رابطه‌های بالا را برای $I(\nu, \lambda; x^\circ)$ و $J(\lambda, \nu; x^\circ)$ بدست آورد. همچنین آنگر (۱۹۷۳) با استفاده از نتایجی که در واها (۱۹۶۸) و لیگت (۱۹۷۱) اثبات شده است نشان داده است که:

$$|I(\lambda, \nu, x) - I(\lambda, \nu, x^\circ)| = O(T^{-1})$$

به شرطی که:

$$\sum_{-\infty}^{+\infty} |t|^{1+\beta} |\sigma_j(t)| \leq M_j^\beta < \infty \quad 0 < \beta < 1 \quad j = 1, 2$$

که در آن $M > 0$ یک عدد ثابت می‌باشد. اکنون اگر $f_T(\lambda)$ به عنوان برآوردیاب ناریب ماتریس طیفی هر سری برداری نمونه‌ای در نظر گرفته شود که به صورت ناپارامتری بدست آمده است و $f(\lambda)$ به عنوان طیف جامعه Z ام برای $j = 1, 2$ در نظر گرفته شود. در آن صورت تقریب زیر برای معیارهای ممیزی کولیک - لیبلر و چرنوف به ترتیب مطابق با رابطه‌های (۸) و (۹) به صورت زیر بکار برده می‌شوند:

$$I(f_T, f_j) = \frac{1}{\nu T} \sum_{k=0}^{T-1} \left[\left(tr(f_T(\lambda_k) f_j^{-1}(\lambda_k)) \right) - \log \frac{f_T(\lambda_k)}{f_j(\lambda_k)} - 1 \right] \quad (10)$$

و همچنین:

$$Q_r(f_T, f_j) = \frac{1}{\nu T} \sum_{k=0}^{T-1} \left(\log \frac{|rf_T(\lambda_k) + (1-r)f_j(\lambda_k)|}{|f_j(\lambda_k)|} - r \log \frac{|f_T(\lambda_k)|}{|f_j(\lambda_k)|} \right) \quad (11)$$

لذا با مینیمم کردن اطلاع ممیزی بین چگالی‌های طیفی هر یک از دو نمونه سری زمانی با جامعه؛
قاعده ممیزی به صورت زیر فراهم می‌شود:

سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب
 $I(f_T, f_2) \geq I(f_T, f_1)$ و $I(f_T, f_2) < I(f_T, f_1)$ حال اگر گروه ماتریس طیفی بوسیله
متوسط ماتریس طیفی نمونه‌ای برای جامعه J ام به صورت زیر برآورد شود:

$$\bar{f}_j(\lambda_k) = \frac{1}{n_j} \sum_{l=1}^{n_j} f_T^{(l)}(\lambda_k)$$

که در آن $f_T^{(l)}(\lambda_k)$ برآورد ناپارامتری ماتریس طیفی l امین نمونه از جامعه J ام می‌باشد؛ می‌توان
با جایگزین کردن \bar{f}_j با f_j قاعده ممیزی بالا را به صورت زیر بدست آورد:
سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب
 $I(f_T, \bar{f}_2) \geq I(f_T, \bar{f}_1)$ و $I(f_T, \bar{f}_2) < I(f_T, \bar{f}_1)$ برقرار باشد.

۲.۲ آنالیز ممیزی بر اساس معیارهای ممیزی کولبک - لیبلر و چرنوف در حالت چند متغیره

فرض کنید مساله رده‌بندی؛ تخصیص سری‌های زمانی بردای $pT \times 1$
بعدی $x = (x'_0, x'_1, \dots, x'_{T-1})'$ که شامل p بردار سری زمانی یک متغیره
 $p \times 1$ بعدی $x_t = (x_0, x_1, \dots, x_{T-1})'$ با ماتریس کواریانس $pT \times pT$ بعدی
 $R_j(s-t)$ برای $s, t = 0, 1, \dots, T-1$ که خود شامل ماتریس‌های $p \times p$ بعدی
 $R_j = \{\sigma_{j(s-t)}, t, s = 0, 1, \dots, T-1\}$ تحت فرضیه‌های جدای H_j برای $j = 1, 2$
باشد. در این حالت معیارهای ممیزی کولبک - لیبلر و چرنوف به صورت زیر می‌باشند.

$$I(1, 2; x) = \frac{1}{2T} \left(\text{tr}(R_1 R_2^{-1}) - \log \frac{|R_1|}{|R_2|} - pT \right) \quad (12)$$

$$Q_r(1, 2; x) = -\frac{1}{2T} \left(\log \frac{|rR_1 + (1-r)R_2|}{|R_2|} - r \log \frac{|R_1|}{|R_2|} \right) \quad (13)$$

مانند حالت یک متغیره می‌توان معیار واگرایی J را به دلیل نامتقارن بودن معیارهای ممیزی کولبک
- لیبلر و چرنوف به صورت زیر تعریف کرد:
کاکیزاوا و دیگران (۱۹۹۸).

$$J(1, 2; x) = I(1, 2; x) + I(2, 1; x)$$

$$J_{Q_r}(1, 2; x) = Q_r(1, 2; x) + Q_r(2, 1; x)$$

ملاحظه می‌گردد که برای محاسبه رابطه‌های (۱۲) و (۱۳) لازم به محاسبه وارون و دترمینان ماتریس‌هایی با ابعاد بزرگ می‌باشد؛ و مخصوصاً زمانی که مقدار T بزرگ باشد این کار بسیار دشوار است. در این حالت نیز استفاده از تقریب‌های طیفی بسیار مفید می‌باشد؛ چنانچه در تبدیل ماتریس‌های کواریانس به ماتریس‌های طیفی نه تنها هیچ اطلاعی از دست نخواهد رفت بلکه بعد مجموعه ماتریس‌های از مرتبه $pT \times pT$ به ماتریس‌های از مرتبه $p \times p$ کاهش می‌یابد. در این زمینه می‌توان به هانن (۱۹۷۰)؛ پینسکر (۱۹۶۴)؛ شاموی و آنگر (۱۹۷۴) کاکیزاوا و دیگران (۱۹۹۸) مراجعه نمود. کازاکس - پاپ‌آنتونی کازاکس (۱۹۸۰) نیز خواص حدی رابطه‌های (۱۵) و (۱۶) را به صورت زیر فراهم نمودند:

$$\lim_{T \rightarrow \infty} \frac{1}{T} I(1, 2; x_t) \quad (14)$$

$$= I(f_1, f_2; x_t) = \frac{1}{2T} \int_{-\pi}^{\pi} \left(tr \left(f_1(\lambda) f_2^{-1}(\lambda) \right) - \log \frac{f_1(\lambda)}{f_2(\lambda)} - p \right) \frac{d\lambda}{2\pi}$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} Q_R(1, 2; x_t) \quad (15)$$

$$= Q_r(f_1, f_2; x_t) = \frac{1}{2T} \int_{-\pi}^{\pi} \left(\log \frac{|r f_1(\lambda) + (1-r) f_2(\lambda)|}{|f_2(\lambda)|} - r \log \frac{|f_1(\lambda)|}{|f_2(\lambda)|} \right) \frac{d\lambda}{2\pi}$$

که در آن $f_j(\lambda_k)$ ماتریس‌های طیفی $p \times p$ بعدی جامعه‌ها برای $j = 1, 2$ می‌باشند. با استفاده از تبدیل انتگرال به مجموع ریمان رابطه‌های (۱۴) و (۱۵) را می‌توان به صورت زیر بدست آورد:

$$I(1, 2; x) = \frac{1}{2T} \sum_{k=0}^{T-1} \left[tr \left(f_1(\lambda_k) f_2^{-1}(\lambda_k) \right) - \log \frac{f_1(\lambda_k)}{f_2(\lambda_k)} - p \right] \quad (16)$$

و همچنین:

$$Q_r(1, 2; x) = \frac{1}{2T} \sum_{k=0}^{T-1} \left(\log \frac{|r f_1(\lambda_k) + (1-r) f_2(\lambda_k)|}{|f_2(\lambda_k)|} - r \log \frac{|f_1(\lambda_k)|}{|f_2(\lambda_k)|} \right) \quad (17)$$

ملاحظه می‌شود که رابطه‌های (۱۶) و (۱۷) هر دو تابعی از ماتریس $tr \left(f_1(\lambda_k) f_2^{-1}(\lambda_k) \right)$ می‌باشند و در موارد خاصی به شکل عمومی‌تر برای بعضی از مقادیر تابع $H(\cdot)$ ؛ که یک اندازه جداسازی عمومی‌تر به شکل زیر است می‌توان تعریف نمود: [کاکیزاوا و دیگران (۱۹۹۸)].

$$D_H(f_1, f_2) = \frac{1}{2} \int_{-\pi}^{\pi} H \left(f_1(\lambda) f_2^{-1}(\lambda) \right) \frac{d\lambda}{2\pi} \quad (18)$$

به طوریکه:

$$H(\lambda) = f_1(\lambda)f_2^{-1}(\lambda)$$

لازم به ذکر است که فاصله عمومی $D_H(f_1, f_2)$ در رابطه (۱۸) دارای خاصیت شبه فاصله‌ای است. همچنین تابع $H(\cdot)$ را می‌توان به شکل‌های مختلفی برگزید به طوریکه خاصیت شبه فاصله‌ای آن حفظ گردد. در ادامه تنها برای دو رابطه (۱۲) و (۱۳) این تابع در نظر گرفته می‌شود که می‌توان آن را به صورت زیر نوشت: [کاکیزاوا و دیگران (۱۹۹۸)].

$$H_I(Z) = \text{tr}(Z) - \log |Z| - p \quad (19)$$

$$H_Q(Z) = \log |rZ + (1-r)I_p| - r \log |Z|$$

نکته‌ای که باید توجه کرد این است که عموماً فاصله $D_H(f_1, f_2)$ متقارن نیست و با تعریف رابطه زیر مقدار متقارنی از $D_H(f_1, f_2)$ به صورت زیر بدست می‌آید:

$$\widetilde{H}(z) = H(Z) + H(Z^{-1}) \quad (20)$$

و با قرار دادن رابطه‌های (۱۹) در رابطه (۲۰) معیار متقارنی از معیارهای کولبک - لیبار و چرنوف فراهم می‌شود. کاکیزاوا و دیگران (۱۹۹۸). با جمع بستن بر روی فرکانس‌هایی به شکل $k = 0, 1, \dots, T-1; \lambda_k = 2\pi kT^{-1}$ رابطه زیر بدست می‌آید:

$$D_H(f_1, f_2) \approx \frac{1}{\sqrt{T}} \sum_k \left(f_1(\lambda_k) f_2^{-1}(\lambda_k) \right) \quad (21)$$

اکنون اگر $f_T(\lambda)$ به عنوان برآورد ماتریس طیفی هر سری برداری نمونه‌ای باشد که به صورت ناپارامتری بدست آمده است و $f_j(\lambda)$ به عنوان طیف جامعه زام برای $j = 1, 2$ در نظر گرفته شود. در آن صورت تقریب زیر برای معیارهای ممیزی کولبک - لیبار و چرنوف به ترتیب مطابق با رابطه‌های (۱۶) و (۱۷) به صورت زیر بکار برده می‌شوند:

$$I(f_T, f_j) = \frac{1}{\sqrt{T}} \sum_{k=0}^{T-1} \left[\left(\text{tr}(f_T(\lambda_k) f_j^{-1}(\lambda_k)) \right) - \log \frac{f_T(\lambda_k)}{f_j(\lambda_k)} - p \right] \quad (22)$$

و همچنین:

$$Q_r(f_T, f_j) = \frac{1}{\sqrt{T}} \sum_{k=0}^{T-1} \left(\log \frac{|r f_T(\lambda_k) + (1-r) f_j(\lambda_k)|}{|f_j(\lambda_k)|} - r \log \frac{|f_T(\lambda_k)|}{|f_j(\lambda_k)|} \right) \quad (23)$$

با این تعریف این رابطه‌ها؛ یک روش عمومی بر اساس رابطه‌های (۱۵) و (۱۶) می‌توان بدست آورد. با این روش اندازه جداسازی بین طیف نمونه‌ای x_t و جامعه z_j بوسیله $D_H(f_T, f_j)$ اندازه‌گیری می‌شود و این اندازه به عنوان تابع ممیزی در نظر گرفته می‌شود. بر اساس این معیار تابع ممیزی به صورت زیر می‌باشد:

$$D_H = D_H(f_T, f_2) - D_H(f_T, f_1) \geq 0 \quad (24)$$

سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب $D_H \geq 0$ و $D_H < 0$ برقرار باشد.

به ازای $H_I(Z)$ داده شده در رابطه (۱۹) مقدار $D_H(f_T, f_j)$ شامل مینیمم اطلاعات ممیزی می‌باشد که اولین بار توسط کولبک (۱۹۵۱) بیان شده است. می‌توان آنالیز ممیزی را با استفاده از معیار ممیزی کولبک - لیبار نیز انجام داد که با مینیمم کردن اطلاع ممیزی بین چگالی‌های طیفی هر یک از دو نمونه سری زمانی با جامعه؛ تابع ممیزی با استفاده از رابطه‌های (۲۲) و (۲۴) به صورت زیر بدست آورده می‌شود:

$$D_H(I) = I(f_T, f_2) - I(f_T, f_1)$$

و قاعده ممیزی بدین صورت می‌باشد که سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب $D_H(I) \geq 0$ و $D_H(I) < 0$ برقرار باشد. همچنین بر اساس معیار چرنوف نیز می‌توان تابع ممیزی را با استفاده از رابطه‌های (۲۳) و (۲۴) به صورت زیر بدست آورد:

$$D_H(I) = Q_r(f_T, f_2) - Q_r(f_T, f_1)$$

و قاعده ممیزی بدین صورت می‌باشد که سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب $D_H(Q_r) \geq 0$ و $D_H(Q_r) < 0$ برقرار باشد. حال اگر گروه ماتریس طیفی بوسیله متوسط ماتریس طیفی نمونه‌ای برای جامعه z_j به صورت زیر برآورد شود:

$$\bar{f}_j(\lambda_k) = \frac{1}{n_j} \sum_{l=1}^{n_j} f_T^{(l)}(\lambda_k)$$

که در آن $f_T^{(l)}(\lambda_k)$ برآورد ناپارامتری ماتریس طیفی l امین نمونه از جامعه z_j می‌باشد؛ می‌توان با جایگزین کردن \bar{f}_j با f_j قاعده ممیزی بالا را به صورت زیر بدست آورد:

سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب $I(f_T, \bar{f}_2) - I(f_T, \bar{f}_1) \geq 0$ و $I(f_T, \bar{f}_2) - I(f_T, \bar{f}_1) < 0$ برقرار باشد.

همچنین بر اساس معیار چرنوف سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب $Q_r(f_T, \bar{f}_2) - Q_r(f_T, \bar{f}_1) \geq 0$ و $Q_r(f_T, \bar{f}_2) - Q_r(f_T, \bar{f}_1) < 0$ برقرار باشد.

$Q_r(f_T, \bar{f}_2) - Q_r(f_T, \bar{f}_1) < 0$ برقرار باشد.

کاکیزاوا و دیگران (۱۹۹۸) استفاده از روش‌های درست‌نمایی را در یافتن تابع ممیزی بکار بردند. این روش یکی از روش‌های معمول در یافتن تابع ممیزی بین دو سری زمانی ایستای نرمال با میانگین صفر ماتریس‌های طیفی $f_j(\lambda)$ برای $j = 1, 2$ می‌باشد و به صورت زیر است:

$$L(f_j) = \log p_j(x) = \frac{1}{2T} \sum_{k=0}^{T-1} \left[j(\lambda_k) - \frac{|X(k)|^2}{f_j(\lambda_k)} \right]$$

که در آن:

$$|X(k)|^2 = I_T(\lambda) = \frac{1}{2\pi T} \left\{ \sum_{k=0}^{T-1} x_t \exp\{-i\lambda_k t\} \right\} \left\{ \sum_{k=0}^{T-1} x_t \exp\{-i\lambda_k t\} \right\}^*$$

اگر فرایند نرمال باشد: در نتیجه تقریب طیفی لگاریتم راست‌نمایی x_t به صورت زیر ساده می‌شود: [وایتل (۱۹۵۴)].

$$L(f_j) = \frac{1}{2T} \sum_{k=0}^{T-1} \left[\log f_j(\lambda_k) + tr \left\{ I_T(\lambda_k) f_j^{-1}(\lambda_k) \right\} \right]$$

و بوسیله تقریب مجموع به حالت انتگرال رابطه بالا به صورت رابطه زیر خلاصه می‌شود:

$$L(f_j) = \frac{1}{2T} \int_{-\pi}^{\pi} \left[\log f_j(\lambda_k) + tr \left\{ I_T(\lambda_k) f_j^{-1}(\lambda_k) \right\} \right] \frac{d\lambda}{2\pi}$$

برای $k = 1, 2, \dots, T-1$, $\lambda_k = 2\pi k T^{-1}$

قاعده ممیزی بر اساس نسبت راست‌نمایی به صورت زیر می‌باشد که سری زمانی x_t به ترتیب به جامعه اول و یا دوم تخصیص داده می‌شود؛ اگر به ترتیب $L(f_1) - L(f_2) \geq 0$ و $L(f_1) - L(f_2) < 0$ برقرار باشد.

۳ مقایسه روش‌های ممیزی برای دو فرایند $AR(1)$

در این بخش به منظور مقایسه روش‌های ممیزی در استفاده از تابع ممیزی کولیک - لیبار و تابع ممیزی چرنوف؛ در ممیزی بین دو فرایند $AR(1)$ با دو پارامتر مختلف ϕ_1 و ϕ_2 به ترتیب برای فرضیه اول و دوم؛ نرخ‌های خطا یا رده‌بندی نادرست؛ محاسبه می‌شوند. برای این کار تعداد ۲۰۰ سری زمانی ایستای $AR(1)$ هر کدام به طول ۱۰۰ مشاهده در نرم‌افزار $S-plus 2000$ شبیه‌سازی شده است. جدول‌های (۱) و (۲) به ترتیب نشان دهنده نرخ‌های خطا یا رده‌بندی

نادرست به ترتیب برای تابع ممیزی کولبک - لیپلر و تابع ممیزی چرنوف می‌باشند. چنانچه در جداول (۱) و (۲) دیده می‌شود: هر دو تابع ممیزی در ممیزی بین دو جامعه به طور مناسب و خوبی عمل کرده‌اند. چنانچه مقادیر خطای رده‌بندی نادرست برای آنها تقریباً کم می‌باشد. همچنین همان طور که دیده می‌شود تابع ممیزی کولبک - لیپلر در مقایسه با تابع ممیزی چرنوف در ممیزی بین دو جامعه بهتر عمل کرده است؛ چنانچه خطای رده‌بندی نادرست برای آن کمتر می‌باشد.

جدول ۱: مقادیر نرخ خطا (رده‌بندی نادرست) برای دو فرایند $AR(1)$ با پارامترهای مختلف ϕ_1 و ϕ_2 با استفاده از تابع ممیزی چرنوف

| | -۰/۸ | -۰/۷ | -۰/۶ | -۰/۵ | -۰/۴ | -۰/۳ | -۰/۲ | -۰/۱ | ۰/۱ | ۰/۲ | ۰/۳ | ۰/۴ | ۰/۵ | ۰/۶ | ۰/۷ | ۰/۸ |
|------|------|------|------|------|------|------|------|------|-----|-----|-----|-----|-----|-----|-----|-----|
| -۰/۸ | * | ۲۸ | ۱۵ | ۷ | ۹ | ۲ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ |
| -۰/۷ | ۳۲ | * | ۳۴ | ۱۵ | ۱۳ | ۸ | ۴ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ |
| -۰/۶ | ۱۹ | ۳۷ | * | ۴۱ | ۲۲ | ۱۹ | ۸ | ۳ | ۱ | ۱ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ |
| -۰/۵ | ۱۳ | ۳۰ | ۳۲ | * | ۳۵ | ۲۵ | ۱۸ | ۱۴ | ۲ | ۱ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ |
| -۰/۴ | ۱۱ | ۱۶ | ۲۵ | ۳۶ | * | ۳۶ | ۳۶ | ۱۶ | ۴ | ۶ | ۳ | ۲ | ۰ | ۰ | ۰ | ۰ |
| -۰/۳ | ۸ | ۹ | ۱۷ | ۳۶ | ۳۵ | * | ۳۳ | ۲۸ | ۸ | ۱۰ | ۵ | ۳ | ۳ | ۰ | ۰ | ۰ |
| -۰/۲ | ۶ | ۵ | ۱۶ | ۲۷ | ۲۸ | ۳۹ | * | ۴۱ | ۲۳ | ۱۶ | ۷ | ۴ | ۲ | ۰ | ۰ | ۰ |
| -۰/۱ | ۳ | ۹ | ۷ | ۲۰ | ۲۳ | ۲۹ | ۴۵ | * | ۲۷ | ۱۶ | ۱۵ | ۱۲ | ۱۰ | ۲ | ۱ | ۰ |
| ۰/۱ | ۰ | ۴ | ۲ | ۴ | ۱۳ | ۹ | ۲۲ | ۳۲ | * | ۳۸ | ۳۰ | ۱۸ | ۱۰ | ۱۲ | ۷ | ۷ |
| ۰/۲ | ۰ | ۱ | ۰ | ۳ | ۷ | ۱۶ | ۱۸ | ۱۳ | ۴۴ | * | ۴۱ | ۲۸ | ۲۰ | ۱۵ | ۴ | ۶ |
| ۰/۳ | ۰ | ۰ | ۲ | ۳ | ۳ | ۳ | ۵ | ۹ | ۲۵ | ۳۹ | * | ۳۹ | ۲۵ | ۲۷ | ۲۰ | ۴ |
| ۰/۴ | ۰ | ۰ | ۱ | ۰ | ۰ | ۱ | ۵ | ۳ | ۱۵ | ۱۸ | ۳۷ | * | ۴۴ | ۳۳ | ۱۹ | ۷ |
| ۰/۵ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۲ | ۴ | ۱۱ | ۱۲ | ۲۹ | ۴۱ | * | ۳۷ | ۲۰ | ۸ |
| ۰/۶ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۵ | ۶ | ۱۱ | ۲۹ | ۳۷ | * | ۴۱ | ۲۵ |
| ۰/۷ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۴ | ۵ | ۱۲ | ۱۱ | ۲۰ | * | ۳۳ |
| ۰/۸ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۰ | ۱ | ۲ | ۹ | ۱۹ | ۲۷ | * |

مراجع

- [1] Alogon, J.(1989), Spectral Discrimination for Tow Groups of Time Series, J. of Time Series Analysis, vol 10, 203-214.
- [2] Blandford, R. R. (1993), Discrimination of Earthquakes and Explosions at Regional Distances Using Complexity, AFTACTR-93-044, HQ. Air Force Technical Applications Center, Patrick Air Force Base, FL.
- [3] Chernoff, H. (1952), A Measure of Asymptotic Efficiency for Tests a Hypothesis Based on The Sum of The Observations, Ann. Math. Stat, vol 25, 573-578.
- [4] Chinipardaz, R. (2004), Discriminatin in Seismology Based On Autocovariances, Far East J, Theo. Stat. vol 12(1), 101-115.
- [5] Dargahi- Noubary, G. R. and Laycock, P. J. (1981), Spectral Ratio Discrimination Theory, J. of Time Series Analysis, vol 2(2), 71-86.
- [6] Dargahi- Noubary, G. R. (1992), Discrimination Between Gaussian Time Series Based on Their Spectral Differences, Comm. Stat. Theory and Method, vol 21(9), 2439-2458.
- [7] Dewitt, P. E. (1985), His Master's (Digital) Voise, Time 125, No 13, 84-85.
- [8] Geresh, W. Martinelli, F. Yonemoto, J. Lew, M. D. & McEwan, J. A. (1979), Automatic Classification of Electroencephalograms: Kullback- Leibler Nearest neighbor Rules, Science, vol 205, 193-195.
- [9] Gevins, A. S. Veager, C. L. Diamond, S. Spire, J. Zeitlin, G. & Gevins, A. (1975), Automated Analysis of The Electerical Activity of The Human Brain (EEG), A Progress Report, Proc. IEEE 63, 1382-1399.
- [10] Fuller, W. A. (1996), Introduction to Statistical Time Series, John Wiley, New York.
- [11] Hannan, E. J. (1968), Multiple Time Series, Wiley, New York. (1992), Discrimination Between Gaussian Time Series Based on Their Spectral Differences, Comm. Stat. Theory and Method, vol 21(9), 2439-2458. (1992)
- [12] Kakizava, Y., Shumway, R. H. and Taniguchi, M. (1998), Discrimination and Clustering for Multivariate Time Series, JASA, vol 93, NO, 441,theory and Methods.
- [13] Kullback, S. & Leibler, R. A. (1951), An Information and Sufficiency, Ann. Math. Stat, vol 22, 76-86.

- [14] Kullback, S. (1952), An Application of Information Theory to Multivariate Analysis, *Ann. Math. Stat*, vol 23, 88-102.
- [15] Liggett, W. S. (1971), On The Asymptotic Optimality of Spectral Analysis for Testing Hypotheses about Time Series , *Ann. Math. Stat*, vol 42, 1348-1358.
- [16] Kazakos, D. & Popantoni-Kazakos, P. (1980), Spectral Distance Measuring Between Gaussian Processes *IEEE, Transactions on Automatic Control*, AC-25, 950-959.
- [17] Markel, J. D. & Gray, A. H., Jr. (1976), *Linear Prediction of Speech*, Springer, Berlin.
- [18] Parzen, E. (1990), *Time Series Statistics and Information*, IMA.
- [19] Piccolo, D. (1990), A Distance Measure for Classifying ARIMA Models, *J. of Time Series Analysis*, vol 11(21), 153-164.
- [20] Pinsker, M. S. (1964), *Information and Information Stability of Random Variables and Processes*, Holden-Day, San Francisco.
- [21] Renyi, A. (1961), On Measure on Entropy and Information, in *Proceeding of 4th Berkeley Sumposium on Mathematical Statistics and Probability*. Berkeley: University of California Press, PP. 547-561.
- [22] Raji, E. (1986), Evidence for a Peripheral Origin of The Brainstem Auditory Evoked Potentials, *Electronyogr. Clin. Neurophysiol*, vol 26, 13-16.
- [23] Shumway, R. H. & Unger, A. N. (1974), LInear Discrimination Function for Stationary Time Series, *JASA*, vol 69, 948-956.
- [24] Shumway, R. H. (1982). *Discriminant Analysis for Time Series*, *Handbook of Statistics*, vol 2, 1-46, eds, Krishnaiah, P. R. & kanal, L. N. Amsterdam, North-Holland.
- [25] Shumway, R. H. (1988), *Applied Statistical Time Series Analysis*, Prentice Hall.
- [26] Shumway, R. H. (1996), *Statistical Approaches to Seismic Discrimination*, in *Monitoring a Comprehensive Test Ban Treaty* eds. Husebye, E. S. & Dainty, A. M. Doordrecht: Kluwer, PP, 791-803.
- [27] Tjostheim, D. (1974), Autoregressive Representation of Seismic P-wave Signals with an Application to The Problem of Short Period Discriminants, *Geophys. J. Roy. Astron. Soc.* vol 43, 269-291.
- [28] Wahba, G. (1968), On The Distribution of Some Statistic Useful in the Analysis of Jointly Stationary Time Series, *Ann. Math. Stat*, vol 38, 1849-1868.

- [29] Whittle, P. (1954), Estimation and Information in Stationary Time Series, Arkiv for Matematik, vol 2, 423-434.
- [30] Wolf, J. J. (1976), Speech Recognition and Understanding Digital Pattern Recognition, K. S. Fu, Ed., Springer -Verlag, Berlin , 167-203.

مدل تاثیر تصادفی با داده‌هایی با پاسخهای آمیخته

علی صادقی^۱، مجتبی گنجعلی^۲

^۱ کارشناسی ارشد آمار اقتصادی - اجتماعی دانشگاه شهید بهشتی

^۲ عضو هیئت علمی دانشگاه شهید بهشتی

چکیده: در کاربردهای مختلف آمار چند متغیره در شاخه‌های مختلف علوم با مدل‌هایی تاثیر تصادفی مواجهیم که متغیر پاسخ برداری از داده‌های آمیخته پیوسته و گسسته است. در چنین موارد شناسایی همبستگی میان متغیرهای پیوسته و گسسته اهمیت بسیاری دارد. هنگامی که تمامی متغیرهای پاسخ پیوسته باشند می‌توان از تحلیل عاملی استفاده نمود و هنگامی که کلیه متغیرها گسسته باشند استفاده از روشهای چند متغیره گسسته می‌تواند سودمند باشد. در این مقاله با استفاده از روش تجزیه تک عاملی روشی ارائه گردیده است که می‌توان آمیخته‌ای از متغیرهای پیوسته و گسسته را با به دست آوردن همبستگی‌ها در مدل اثرات تصادفی مورد بررسی قرار داد. در این مدل متغیرهای پیوسته از توزیع نرمال و متغیرهای گسسته از توزیع دودویی تبعیت می‌کنند، شده‌اند. با استفاده از روش تجزیه تک عاملی علاوه بر بدست آوردن همبستگی میان متغیرهای پاسخ، تاثیر عوامل محیطی و ژنتیکی را نیز در نظر می‌گیریم. در این مقاله روشی مبتنی بر استفاده از تاثیرهای تصادفی را برای مدل‌بندی پاسخهای آمیخته نرمال و دودویی معرفی و نحوه تحلیل تاثیر برخی متغیرهای تبیینی بر این پاسخها را نیز شرح می‌دهیم.

واژه‌های کلیدی: اثرات تصادفی، تجزیه عاملی، متغیر پنهان، پاسخهای آمیخته

۱ مقدمه

در بسیاری از تحقیقات آماری با حالتی مواجه هستیم که آمیخته‌ای از داده‌های پیوسته و گسسته به عنوان متغیر پاسخ مد نظر است و بعضی از متغیرها دارای اثرات تصادفی هستند این مدلها توسط لیتل^۱ و رابین^۲ (۲۰۰۲) بسیار مورد توجه قرار گرفته است. در چنین حالتی بر خلاف هنگامی که تمامی اثرات ثابت هستند نمی‌توان از روشهای معمولی مورد استفاده در علم آمار مدل مناسب را به داده‌ها برازش داد [آگرستی^۳ (۲۰۰۲)]. روشهایی که تاکنون برای برازش مدل‌های با اثرات تصادفی برای داده‌های آمیخته مورد بررسی قرار گرفته‌اند به دلیل استفاده از روشهای عددی در برازش مدل مناسب از تئوری پیچیده‌ای استفاده کرده‌اند و بسیاری از آنها تنها در شرایط بخصوصی به نتیجه مورد نظر می‌رسند. در این مقاله سعی گردیده است که با استفاده از تئوری ساده و با استفاده از روشهای عددی ساده‌تر این بحث مورد توجه قرار گیرد.

1) Little 2) Rubin 3) Agresti

۲ مدل اثرات تصادفی برای پاسخهای آمیخته

مانند مدل مورد استفاده شامل و ریان^۴ (۱۹۹۶) یک مدل اثرات تصادفی را با M برآمد مشاهده شده که تعداد M_1 تا از این برآمدها دارای توزیع دودویی و $M_2 = M - M_1$ تا از برآمدهای باقی مانده از توزیع نرمال پیروی می‌کنند در نظر می‌گیریم. ذکر این نکته ضروری است که متغیرهای تصادفی در مدل مورد بحث مانند اثرات متغیرهای پنهان رفتار می‌کنند [بکر^۵ (۱۹۹۲)]. برای این منظور مدل اثرات تصادفی را برای آزمودنی i -ام به صورت زیر در نظر می‌گیریم:

$$\begin{aligned} y_{im}^* &= X_{im}^t \beta + \varepsilon_{im} & m = 1, \dots, M_1 \\ z_{im} &= U_{im}^t \beta + \varepsilon_{im} & m = M_1 + 1, \dots, M \end{aligned} \quad (۱)$$

که در آن X_{im} و U_{im} ماتریس متغیرهای تبیینی برای آزمودنی i -ام در برآمد m -ام است. در این مدل y_{im} مقدار مشاهده شده یک متغیر پنهان y_{im}^* با توزیع نرمال استاندارد است که فقط کوچکتر یا بزرگتر از عدد \circ بودن آنرا در دست داریم و اگر $y_{im}^* > \circ$ آنگاه $y_{im} = 1$ و در غیر این صورت $y_{im} = 0$ خواهد بود.

ماتریس واریانس کواریانس مدل فوق دارای پارامترهای زیادی حتی با در نظر گرفتن مقادیر کوچک M خواهد بود. برای درک بهتر این مسئله حالتی را در نظر می‌گیریم که $M_1 = 4$ و $M_2 = 2$ است. در این صورت مدل به صورت زیر بیان می‌شود.

$$\begin{aligned} y_{im}^* &= X_{im}^t \beta + \varepsilon_{im} & m = 1, 2, 3, 4 \\ z_{im} &= U_{im}^t \beta + \varepsilon_{im} & m = 5, 6 \end{aligned} \quad (۲)$$

در این صورت ساختار ماتریس واریانس کواریانس غیرساختاری به صورت زیر در خواهد آمد

$$\Sigma = \begin{pmatrix} 1 & \rho_{12} & \rho_{13} & \rho_{14} & \sigma_1 \rho_{15} & \sigma_2 \rho_{16} \\ \rho_{12} & 1 & \rho_{23} & \rho_{24} & \sigma_1 \rho_{25} & \sigma_2 \rho_{26} \\ \rho_{13} & \rho_{23} & 1 & \rho_{34} & \sigma_1 \rho_{35} & \sigma_2 \rho_{36} \\ \rho_{14} & \rho_{24} & \rho_{34} & 1 & \sigma_1 \rho_{45} & \sigma_2 \rho_{46} \\ \sigma_1 \rho_{15} & \sigma_1 \rho_{25} & \sigma_1 \rho_{35} & \sigma_1 \rho_{45} & \sigma_1^2 & \sigma_1 \sigma_2 \rho_{56} \\ \sigma_2 \rho_{16} & \sigma_2 \rho_{26} & \sigma_2 \rho_{36} & \sigma_2 \rho_{46} & \sigma_1 \sigma_2 \rho_{56} & \sigma_2^2 \end{pmatrix}$$

با وجود اثرات تصادفی در مدل نمی‌توان از روشهای معمول در آمار به منظور برآزش مدل با وجود آمیخته‌ای از برآمدهای پیوسته و گسسته استفاده کرد. برای بدست آوردن پارامترها از روشهای عددی استفاده می‌گردد اما یکی از مشکلات اصلی استفاده از روشهای عددی آن است که هنگامی که تعداد کمی داده در اختیار داریم و مدل دارای تعداد زیادی برآمد است نمی‌توانند به خوبی عمل کنند. برای مثال هنگامی که ۴ پاسخ گسسته و ۲ پاسخ پیوسته مانند مدل بالا در

4) Sammel and Ryan 5) Baker

دست باشد احتیاج به برآورد ۱۷ پارامتر در ماتریس کواریانس خواهیم داشت و تعداد کم داده‌ها می‌تواند مشکل ساز گردد و دقت مورد نظر در برآورد پارامترها حاصل نخواهد شد. اولین بار شامل، لوتیز^۶ و ریان (۱۹۹۷) این مدل را مورد بررسی قرار دادند اما روش مورد استفاده آنها از تئوری عددی پیچیده‌ای برای بدست آوردن پارامترهای مدل استفاده می‌کند و دیگر اشکال آن حساسیت بسیار زیاد مدل در برابر داده‌های کوچک است. یکی از راههای برخورد با چنین مشکلی کاهش تعداد پارامترهای مدل است که برای اولین بار در این مقاله برای چنین داده‌هایی مورد بررسی قرار گرفته است. برای این منظور استفاده از روش تجزیه تک عاملی که اولین بار توسط هکمن (۱۹۸۱) ارائه شده است می‌تواند بسیار سودمند باشد.

۳ مدل اثرات تصادفی با استفاده از تجزیه تک عاملی

همانگونه که بیان شد، روش تک عاملی می‌تواند در کاهش تعداد پارامترهای مدل مناسب باشد. برای توضیح این روش ابتدا جملات خطا ε_{im} را برای $m = 1, \dots, M$ به دو جزء با اثرات تصادفی (b_{im}) و (v_{im}) به صورت زیر تقسیم می‌کنیم:

$$\varepsilon_{im} = b_{im} + v_{im} \quad m = 1, \dots, M$$

در این صورت مدل (۱) با بکار بردن روش تجزیه تک عاملی که توسط هکمن^۷ (۱۹۸۱) برای اثرات تصادفی مورد استفاده قرار گرفته است، به صورت زیر در خواهد آمد

$$\begin{aligned} y_{im}^* &= X_{im}'\beta_m + b_{im} + v_{im} & m &= 1, \dots, M \\ z_{im} &= U_{im}'\beta_m + b_{im} + v_{im} & m &= M+1, \dots, M \end{aligned} \quad (3)$$

آنگاه برای قابل تشخیص بودن مدل خواهیم داشت:

$$var(\varepsilon_{im}) = var(b_{im} + v_{im}) = \sigma_{b_m}^2 + \sigma_{v_{im}}^2 \quad m = 1, \dots, M$$

که در آن

$$var(\varepsilon_{im}) = var(b_{im} + v_{im}) = \sigma_{b_m}^2 + \sigma_{v_{im}}^2 = 1 \quad m = 1, \dots, M$$

این نوع تعریف برای اثرات تصادفی، یک ساختار برای ماتریس واریانس، کواریانس همانند Σ ایجاد می‌کند.

می‌توان تعداد اثرات تصادفی را با لحاظ محدودیت روی ساختار کواریانس توسط تجزیه تک عاملی

6) Louis 7) Heckman

که توسط هکمن (۱۹۸۱) بیان شده است کاهش داد به این معنی که فرض کنیم $b_{im} = \lambda_m b_i$ برای $m = 1, \dots, M$ که برای برآمدهای دودویی $1 < \lambda_m < -1$ و برای برآمدهای نرمال $-\infty < \lambda_m < +\infty$ خواهد بود که در آن دارای توزیع نرمال استاندارد است و λ_m پارامترهای مقیاس هستند. در این حالت ساختار همبستگی به صورت زیر است:

$$cov(y_m, y_{m'}) = \lambda_m \lambda_{m'} \quad m, m' = 1, \dots, M_1 \quad m \neq m'$$

و

$$cov(z_m, z_{m'}) = \lambda_m \lambda_{m'} \quad m, m' = M_1, \dots, M_2 \quad m \neq m'$$

و همچنین

$$cov(z_{m'}, y_m) = \lambda_m \lambda_{m'} \quad m = 1, \dots, M_1, \quad m' = M_1, \dots, M_2$$

در این صورت تعداد پارامترهای مدل و ماتریس واریانس کواریانس به میزان قابل توجهی کاهش می‌یابد. برای مثال با $M = 6$ و $M_1 = 4$ و $M_2 = 2$ ساختار ماتریس واریانس کواریانس به صورت زیر در خواهد آمد:

$$\Sigma = \begin{pmatrix} 1 & \lambda_1 \lambda_2 & \lambda_1 \lambda_3 & \lambda_1 \lambda_4 & \lambda_1 \lambda_5 & \lambda_1 \lambda_6 \\ \lambda_1 \lambda_2 & 1 & \lambda_3 \lambda_2 & \lambda_2 \lambda_4 & \lambda_2 \lambda_5 & \lambda_2 \lambda_6 \\ \lambda_1 \lambda_3 & \lambda_3 \lambda_2 & 1 & \lambda_3 \lambda_4 & \lambda_3 \lambda_5 & \lambda_3 \lambda_6 \\ \lambda_1 \lambda_4 & \lambda_2 \lambda_4 & \lambda_3 \lambda_4 & 1 & \lambda_5 \lambda_4 & \lambda_6 \lambda_4 \\ \lambda_1 \lambda_5 & \lambda_2 \lambda_5 & \lambda_3 \lambda_5 & \lambda_4 \lambda_5 & \sigma_5^2 + \lambda_5^2 & \lambda_5 \lambda_6 \\ \lambda_1 \lambda_6 & \lambda_2 \lambda_6 & \lambda_3 \lambda_6 & \lambda_4 \lambda_6 & \lambda_5 \lambda_6 & \sigma_6^2 + \lambda_6^2 \end{pmatrix}$$

با استفاده از روش تک عاملی برای $M = 6$ تعداد پارامترهای مدل از ۱۷ پارامتر به ۸ پارامتر کاهش می‌یابد.

۴ مثال کاربردی

لیتل و شلوچتر^۸ (۱۹۸۵) تحقیق روانشناسی بر روی نمونه‌ای از ۶۹ خانوار انجام دادند. هدف از این تحقیق شناسایی تاثیر ناهنجاری‌های والدین بر روی ناهنجاریهای فرزندان خانوار است. آنجا که نمی‌توان اثر شدت ناهنجاری والدین بر روی فرزندان را مشاهده کرد، از مدل اثرات تصادفی در برآزش داده‌های فوق استفاده شده است. در مدل فوق از ۳ برآمد که هر یک به نوعی نشان دهنده (انعکاس دهنده) ناهنجاریهای فرزندان خانوار هستند استفاده شده‌اند هدف از این تحقیق آن است که بدانیم شدت ناهنجاریهای والدین چه تاثیری بر روی هر یک از برآمدها که برآمدهای مشاهده شده از ناهنجاریهای فرزندان هستند دارد. در اینجا یک برآمد دودویی و دو برآمد باقی مانده دارای توزیع نرمال هستند استفاده شده است. که عبارتند از:

8) Little and Schluchter

جدول ۱: خصوصیات توصیفی برآمدها

| متغیر | Low Risk | | H/M Risk | | Total | |
|-------|----------|--------------|----------|--------------|---------|--------------|
| | میانگین | انحراف معیار | میانگین | انحراف معیار | میانگین | انحراف معیار |
| y | ۰/۳۷۵ | ۰/۵۱۷ | ۰/۷۷۷ | ۰/۴۲۷ | ۰/۶۵۳ | ۰/۴۸۵ |
| z_1 | ۱۱۰/۸۷۵ | ۱۵/۹۷۷ | ۹۴/۷۷۷ | ۱۱/۸۱۵ | ۱۰۳/۱۷۳ | ۱۵/۳۶۲ |
| z_2 | ۱۳۶/۷۵ | ۳۰/۳۴۹ | ۹۶/۵۵۵ | ۳۱/۳۲۴ | ۱۱۵/۲۱۷ | ۳۲/۷۶۴ |

متغیر y که دارای توزیع دودویی است و مقدار ۱ را می‌پذیرد اگر فرزند خانواده دارای علائم ناهنجاری در رفتار باشد و مقدار ۰ را می‌پذیرد اگر فرزند خانواده دارای علائم ناهنجاری در رفتار نباشد. توجه به این نکته ضروری به نظر می‌رسد که این برآمد خود برآمد مشاهده شده از یک متغیر پنهان y^* با توزیع نرمال است.

متغیر z_1 که دارای توزیع نرمال است و نشان دهنده نمرات کسب شده فرزند خانوار در آزمون درک مطلب است.

متغیر z_2 که دارای توزیع نرمال بوده و نشان دهنده نمرات کسب شده فرزند خانوار در یک آزمون قدرت محاوره است.

در این مدل از یک متغیر کمکی X که یک متغیر دودویی است استفاده می‌کنیم. این متغیر کمکی مقدار ۱ را می‌پذیرد هنگامی که والدین دارای ناهنجاریهای روانی اندکی باشند (LR) و ۰ را می‌پذیرد هنگامی که حداقل یکی از والدین دارای سابقه بیماری روانی و یا ناهنجاری روانی در حد متوسط یا زیاد (H/MR) باشند. تحقیقات افراد دیگر بر روی داده‌های فوق نشان داده است که تاثیر دو گروه با ناهنجاری متوسط و زیاد یکسان است به همین دلیل این دو گروه را در این مقاله با هم ترکیب کرده‌ایم. (کنجلی^۹)

۵ تحلیل توصیفی داده‌ها

جدول ۱ خصوصیات توصیفی برآمدها را بیان می‌کند. همانگونه که در جدول ۱ اطلاعات توصیفی مربوط به برآمدها مشخص است، نمره درک مطلب در بین کودکانی که والدین آنها در معرض کمتر ناهنجاری روانی (LR) هستند از کودکانی که والدین آنها در معرض ناهنجاری متوسط و شدید روانی (H/MR) هستند بیشتر است. همچنین برای نمره آزمون قدرت محاوره، چنانکه در جدول اطلاعات توصیفی برآمدها مشخص است. نمره کسب شده توسط کودکان در گروه ناهنجاری کمتر (LR) بیشتر از گروه (H/MR) است.

این بدان معنی است که در صورت معنی‌داری، نمره در بین کودکانی که والدین آنها در معرض ناهنجاری متوسط و شدید روانی (H/MR) هستند نسبت به کودکانی که در معرض کمتر ناهنجاری والدین هستند (LR)، کمتر است و میزان تاثیر ناهنجاری والدین در این دو گروه متفاوت خواهد بود.

9) Gangali

جدول ۲: نتایج بدون در نظر گرفتن همبستگی‌ها

| پارامتر | برآورد | سطح معنی داری |
|---------------|--------|---------------|
| $\beta_{0,1}$ | ۰,۷۶۵ | -- |
| $\beta_{1,1}$ | -۱,۰۸۳ | ۰,۳۵۳ |
| $\beta_{0,2}$ | ۹۴,۷۷۸ | -- |
| $\beta_{1,2}$ | ۱۶,۱ | ۰,۰۰۸ |
| $\beta_{0,3}$ | ۹۶,۵۵۶ | -- |
| $\beta_{1,3}$ | ۴۰,۱۹۴ | ۰,۰۰۶ |

۶ تجزیه و تحلیل مدل اثرات تصادفی

مدل اثرات تصادفی مطابق مدل (۱) به صورت زیر در خواهد آمد:

$$\begin{aligned} y_{i\lambda}^* &= \beta_{0,m} + \beta_{1,m}X + \varepsilon_{im} & m = 1 & \quad (4) \\ z_{im} &= \beta_{0,m} + \beta_{1,m}X + \varepsilon_{im} & m = 2, 3 & \end{aligned}$$

در این مدل هنوز همبستگی میان متغیرها (λ ها) را در نظر نگرفته‌ایم. برآورد پارامترهای مدل در جدول ۲ داده شده است.

چنانچه در جدول ۲ مشخص است میان ناهنجاری والدین و ناهنجاری فرزندان رابطه معنی داری وجود ندارد و شدت ناهنجاری والدین تاثیری در ناهنجاری در بین فرزندان در سطح $\alpha = 0/05$ ندارد اما نمره درک مطلب و نمره آزمون محاوره در بین کودکانی که والدین آنها دارای ناهنجاری کمتری (LR) هستند بیشتر از فرزندان است که والدین آنها دارای ناهنجاری متوسط یا زیاد (H/MR) هستند. حال همبستگی میان دو متغیر پیوسته را نیز با وارد کردن اثرات تصادفی و استفاده از روش تک عاملی وارد مدل می‌کنیم در این صورت مدل به صورت زیر در خواهد آمد:

$$\begin{aligned} y_{i\lambda}^* &= \beta_{0,m} + \beta_{1,m}X + v_{i\lambda} & m = 1 & \quad (5) \\ z_{im} &= \beta_{0,m} + \beta_{1,m}X + \lambda_m b_i + v_{im} & m = 2, 3 & \end{aligned}$$

برآورد پارامترهای مدل در جدول ۳ داده شده است. چنانچه در جدول ۳ نیز ملاحظه می‌گردد همان نتایج مدل قبل نیز در این مدل بدست می‌آید علاوه بر آن می‌توان مقدار همبستگی میان برآمدها را نیز بدست آورد. حال حالت کلی را یعنی با وارد کردن اثر تصادفی برای برآمد دودویی بیان می‌کنیم در حقیقت با در نظر گرفتن اثر تصادفی همبستگی میان متغیر گسسته دودویی و متغیرهای پیوسته را نیز در نظر می‌گیریم. در این صورت مدل نهایی به صورت زیر در خواهد آمد:

جدول ۳: نتایج تنها با در نظر گرفتن همبستگی میان متغیرهای پیوسته

| پارامتر | برآورد | سطح معنی داری |
|-------------------|--------|---------------|
| $\beta_{\circ 1}$ | ۰,۷۶۵ | -- |
| β_{11} | -۱,۰۸۳ | ۰,۳۷۱ |
| $\beta_{\circ 2}$ | ۹۳,۷۱۸ | -- |
| β_{12} | ۱۴,۱ | ۰,۰۰۵ |
| λ_2 | ۱۳,۴۷۹ | ۰,۰۰۰ |
| $\beta_{\circ 3}$ | ۹۴,۹۳۷ | -- |
| β_{13} | ۳۷,۱۳۷ | ۰,۰۰۶ |
| λ_3 | ۴,۹۱۲ | ۰,۰۰۰ |
| σ_{v_2} | ۳,۰۷ | -- |
| σ_{v_3} | ۲۳,۱۹۷ | -- |

جدول ۴: نتایج کلی مدل

| پارامتر | برآورد | سطح معنی داری |
|-------------------|--------|---------------|
| $\beta_{\circ 1}$ | ۰,۵۰۳ | -- |
| β_{11} | -۱,۹۴ | ۰,۰۰۵ |
| λ_1 | ۰,۹۹ | ۰,۰۰۰ |
| $\beta_{\circ 2}$ | ۹۳,۳۹۴ | -- |
| β_{12} | ۱۴,۳۹۴ | ۰,۰۰۰ |
| λ_2 | ۱۳,۸۰۳ | ۰,۰۰۰ |
| $\beta_{\circ 3}$ | ۹۴,۴۹۷ | -- |
| β_{13} | ۳۷,۶۶۶ | ۰,۰۰۰ |
| λ_3 | ۲۰,۴۹۹ | ۰,۰۰۰ |
| σ_{v_2} | ۲,۹۰۷ | -- |
| σ_{v_3} | ۲۳,۵۶۵ | -- |

$$y_{i1}^* = \beta_{\circ m} + \beta_{1m}X + \lambda_{1i}b_i + v_{i1} \quad m = 1 \quad (6)$$

$$z_{im} = \beta_{\circ m} + \beta_{1m}X + \lambda_{mi}b_i + v_{im} \quad m = 2, 3$$

برآورد پارامترهای این مدل در جدول ۴ داده شده است. این برآوردها با استفاده از برنامه NAG (۱۹۹۶) و زیر روال $E^{\circ}UCF$ بدست آمده‌اند. چنانچه از سطح معنی داری پارامترهای برآورد شده نیز مشخص است تمامی پارامترها در سطح $\alpha = 0.05$ معنی دار هستند و این بدان معنی است که شدت ناهنجاری والدین بر روی ناهنجاری فرزندان تاثیر دارد و هر چه والدین دارای ناهنجاری بیشتری باشند خطر وجود ناهنجاری در فرزندان نیز بیشتر است همچنین شدت ناهنجاری در والدین در درک مطلب و قدرت محاوره فرزندان تاثیر گذار است به این صورت که هر چه والدین دارای ناهنجاری روانی بیشتری باشند قدرت درک

مطلب و محاوره فرزندان کاهش می‌یابد. با توجه به برآورد پارامترها ماتریس واریانس کواریانس به صورت زیر در می‌آید:

$$\Sigma = \begin{pmatrix} 1 & 13/526 & 20/273 \\ 13/526 & 198/973 & 282/947 \\ 20/273 & 282/947 & 975/518 \end{pmatrix}$$

حال با توجه به ماتریس واریانس کواریانس بالا کواریانس میان برآمدهای گسسته دودویی و پیوسته نرمال با وجود اثرات تصادفی در اختیار است.

۷ نتیجه‌گیری و پیشنهادات

چنانکه در جدول ۱ مشاهده می‌گردد هنگامی که همبستگی میان متغیرها را در مدل در نظر نگرفته‌ایم، میان ناهنجاری والدین و ناهنجاری فرزندان رابطه معنی‌داری وجود ندارد. اما هنگامی که در مدل نهایی همبستگی میان متغیرها را وارد مدل می‌کنیم مدل دقیقتر عمل نموده و تاثیر ناهنجاری والدین بر روی ناهنجاری فرزندان مشخص می‌گردد. نتایج حاصل شده از این مدل دقیقتر از مدل قبل می‌باشد.

در مدل اثرات تصادفی برای داده‌های آمیخته ما تغییرات واریانس برآمدها را درون دو گروه با ناهنجاری کمتر والدین و ناهنجاری متوسط یا بیشتر والدین ثابت در نظر گرفتیم. پیشنهاد می‌گردد هنگامی که تغییرات واریانس معنی‌دار است اثرات تغییرات واریانس را نیز وارد مدل کنیم به طور حتم این امر باعث برآزش بهتر مدل خواهد شد.

مراجع

- [1] Agresti, A. (2002). Categorical data Analysis. New York: Wiley.
- [2] Baker, F. B. (1992). Item response Theory: Parameter Estimation Techniques. New York: Dekker.
- [3] Gangali, M. (2003). A model for Mixed Continuous and Discrete Response with possible of Missing Response . Journal of Sciences, Islamic Republic of Iran, 14(1) , 53-60.
- [4] Heckman, J. (1981). Statistical models for discrete panel data in Manski, C.& Mc Fadden D., structural analysis of discrete data with econometric applications. 114-195, cambridge, mass: MIT press
- [5] Little, R. J. and Schelucher M. (1985). Maximum likelihood estimation for mixed continuous and categorical data with missing value. Biometrika, 72(3) , 497-512.

- [6] Little, R. J. A. and Rubin, D. B. (2002). Statistical analysis with Missing Data. New York: Wiley.
- [7] NAG, (1996). Numerical Algorithms Group Manual, Mark 16, Oxford, UK.
- [8] Sammel, M. D. and Ryan, L. M. (1996). Latent variable models with fixed effects. *Biometrics*, 52, 220-243.
- [9] Sammel, M. D. and Louise. M. and Ryan, L. M. (1997) Latent variable models for Mixed Discrete and Continuous Outcomes. *J. R. Statist. Soc. B*, 59, No 3, 667-678.

آزمون تصادفی شده برای میانگین توزیع نرمال بر پایه‌ی داده‌های نادقیق

سید محمود طاهری، تکتم بزرگوار

دانشکده علوم ریاضی، دانشگاه صنعتی اصفهان

چکیده: در این مقاله یک روش آزمون تصادفی شده برای انجام آزمون‌های مربوط به میانگین توزیع نرمال بر پایه‌ی داده‌های فازی مطرح، تشریح و بررسی می‌شود. در این روش تابع عضویت آماره‌ی آزمون بر اساس مقادیر مختلف α برش‌ها ساخته می‌شود. بر پایه‌ی تابع عضویت این آماره، احتمال رد فرض صفر محاسبه می‌شود. نشان داده می‌شود که در حالت خاص، این روش به روش آزمون فرض کلاسیک تبدیل می‌شود. شیوه‌ی فوق با چند مثال عددی مبتنی بر داده‌های واقعی، که به صورت اعداد فازی مثلثی بیان شده‌اند، تشریح می‌گردد.

واژه‌های کلیدی: داده‌های نادقیق، آزمون تصادفی شده، عدد فازی ذوزنقه‌ای

۱ مقدمه و تاریخچه

آزمون یک فرض آماری عبارت از به کارگرفتن مجموعه‌ی قواعدی است برای آن که تصمیم بگیریم که آیا یک فرض را بپذیریم یا آن را به نفع یک فرض مقابل رد کنیم. مثلاً فرض کنید که آماردانی می‌خواهد فرض صفر $\theta = \theta_0$ را در برابر فرض مقابل $\theta = \theta_1$ آزمون کند. یک نمونه تصادفی از جامعه موردنظر انتخاب می‌کند و بر اساس یک تابع آزمون تصمیم به رد یا پذیرش فرض صفر می‌گیرد. در آزمون‌های غیر تصادفی هر تابع آزمون فضای نمونه‌ای را به دو مجموعه، یک ناحیه قبول و یک ناحیه رد (بحرانی)، افراز می‌کند.

یکی از مفروضات مهم در آزمون‌های آماری آن است که داده‌های مشاهده شده دقیق هستند. اما گاهی اوقات با مشاهدات نادقیق و مبهم روبرو هستیم. در این موارد نمی‌توان از روش‌های متداول آماری استفاده کرد. یکی از روش‌های جانشین در این موارد به کارگیری مجموعه‌های فازی در فرمول‌بندی داده‌های نادقیق است. بدین ترتیب با مسئله آزمون فرض‌های آماری بر پایه‌ی داده‌های فازی روبرو هستیم.

آزمون فرض بر پایه‌ی داده‌های فازی نخستین بار توسط کازالس^۱ و همکاران [3,2] در سال ۱۹۸۶ مورد مطالعه قرار گرفت. آنها لم نیمن-پیرسن را برای حالت بالا تعمیم دادند و رهیافت بی‌زی را به این مسئله بررسی کردند [4]. سون^۲ و همکاران [10] نیز بر پایه‌ی تعمیم لم نیمن-پیرسن،

1) Casals 2) Son

مفهوم تواناترین آزمون فازی را تعریف و یک کاربرد از آن را بررسی کردند. زگورزسکی^۳ [6] رهیافتی را برای موضوع بالا پیشنهاد داده است که در آن، تصمیم درباره‌ی رد یا قبول فرض‌ها، به صورت فازی بیان می‌شود. طاهری و بهبودیان [11] مسئله آزمون فرض‌های فازی را بر پایه‌ی داده‌های فازی با رویکرد بیزی مطالعه نموده‌اند. رهیافت‌های دیگری نیز توسط رامر^۴ و کاندل^۵ [9] و واتانابه^۶ [12] به مسئله‌ی آزمون فرض با داده‌های فازی ارائه شده است. همچنین مونتنگرو^۷ و همکاران [8] آزمون مقایسه‌ی میانگین‌های دو جامعه را بر اساس داده‌های فازی مطالعه نموده‌اند. در این مقاله آزمون میانگین یک توزیع نرمال بر اساس داده‌های نادقیق بر پایه‌ی رهیافت لیو و شن^۸ [7] مورد مطالعه قرار می‌گیرد. نخست به معرفی مفاهیم مقدماتی مورد نیاز در این مقاله پرداخته و سپس در بخش‌های سوم و چهارم به بیان مسئله و تشریح شیوه کار می‌پردازیم. آزمون میانگین جامعه را در حالتی که واریانس جامعه معلوم و مجهول است، به ترتیب در بخش‌های پنجم و ششم مورد بررسی قرار می‌دهیم.

۲ مفاهیم مقدماتی

در این بخش فرض می‌شود که خواننده آشنایی اولیه با مجموعه‌های فازی دارد. در ادامه به چند مفهوم و نتیجه اساسی که در این مقاله مستقیماً به کار گرفته می‌شوند، می‌پردازیم [14,1].
تعریف ۱.۲ یک مجموعه فازی نرمال، محدب و تک نمائی از \mathbb{R} را که تابع عضویت آن قطعه به قطعه پیوسته باشد، عدد فازی گوئیم. مجموعه تمام اعداد فازی از \mathbb{R} را با $F(\mathbb{R})$ نشان می‌دهیم.
تعریف ۲.۲ عدد فازی \tilde{N} را یک عدد فازی LR گوئیم هرگاه تابع عضویت آن به شکل زیر باشد

$$\tilde{N}(x) = \begin{cases} L\left(\frac{m-x}{\alpha}\right) & x \leq m \\ R\left(\frac{x-m}{\beta}\right) & x > m \end{cases}$$

که در آن $L, R : \mathbb{R}^+ \rightarrow [0, 1]$ را توابع مرجع، m را مقدار میانه (نما) و α, β (اعداد حقیقی مثبت) را پهناهای چپ و راست \tilde{N} نامیم. \tilde{N} را با نماد $\tilde{N} = (m, \alpha, \beta)$ نشان می‌دهیم. چنانچه برای عدد فازی $LR = (m, \alpha, \beta)$ داشته باشیم

$$L(x) = R(x) = \begin{cases} 1 - x & 0 \leq x \leq 1 \\ 0 & otherwise \end{cases}$$

آنگاه \tilde{N} را عدد فازی مثالی گوئیم و با نماد $(m, \alpha, \beta)_T$ نشان می‌دهیم.
تعریف ۳.۲ عدد فازی \tilde{M} را یک عدد فازی دوزنقه‌ای LR گوئیم هرگاه تابع عضویت آن به

3) Grzegorzewski 4) Romer 5) Kandel 6) Watanabe 7) Montenegro
 8) Leu and Chen

شکل زیر باشد

$$\tilde{M}(x) = \begin{cases} L\left(\frac{m_1-x}{\alpha}\right) & m_1 > x \\ 1 & m_1 \leq x \leq m_2 \\ R\left(\frac{x-m_2}{\beta}\right) & m_2 < x \end{cases}$$

که در آن L و R توابعی هستند که در تعریف عدد فازی LR صدق می‌کنند. در این حالت \tilde{M} با نماد $(m_1, m_2, \alpha, \beta)_{LR}$ نشان داده می‌شود. واضح است که عدد فازی مثالی حالت خاصی از عدد فازی دوزنقه‌ای است.

تعریف ۴.۲ مجموعه عناصری از X را که درجه عضویت آنها در مجموعه فازی \tilde{A} حداقل به بزرگی α ($\alpha > 0$) باشد، α -برش (مجموعه تراز) \tilde{A} گوئیم و با \tilde{A}_α نشان می‌دهیم، یعنی

$$\tilde{A}_\alpha = \{x \in X; \tilde{A}(x) \geq \alpha\}$$

تعریف ۵.۲ فرض کنید (\mathbb{R}^n, F, P) یک فضای احتمال باشد. یک پیشامد فازی در \mathbb{R}^n عبارتست از یک مجموعه فازی A از \mathbb{R}^n که تابع عضویت آن اندازه پذیر بورل باشد.

تعریف ۶.۲ اگر \tilde{A} یک پیشامد فازی باشد، احتمال آن به صورت زیر تعریف می‌شود [13]

$$P(\tilde{A}) = \int_{\mathbb{R}^n} \tilde{A}(x) dP = E(\tilde{A}(x))$$

که در حالت گسسته علامت انتگرال به علامت سیگما تبدیل می‌شود.

عملگرهای جبری بر اعداد فازی بر پایه اصل گسترش و به صورت زیر تعریف می‌شوند.

تعریف ۷.۲ اگر $\lambda \in \mathbb{R} - \{0\}$ ، آنگاه ضرب اسکالر λ در عدد فازی \tilde{M} به صورت یک عدد فازی با تابع عضویت زیر تعریف می‌شود

$$\lambda \tilde{M}(x) = \tilde{M}\left(\frac{x}{\lambda}\right)$$

تعریف ۸.۲ فرض کنید $\tilde{M}, \tilde{N} \in F(\mathbb{R})$ با توابع عضویت پیوسته باشند. حاصل $\tilde{M} \oplus \tilde{N}$ و $\tilde{M} \ominus \tilde{N}$ به صورت مجموعه‌های فازی با توابع عضویت زیر تعریف می‌شوند

$$(\tilde{M} \oplus \tilde{N})(z) = \sup_{z=x+y} \min[\tilde{M}(x), \tilde{N}(y)]$$

$$(\tilde{M} \ominus \tilde{N})(z) = \sup_{z=x-y} \min[\tilde{M}(x), \tilde{N}(y)]$$

قضیه ۱.۲ اگر $\tilde{M} = (m, \alpha, \beta)_{LR}$ و $\tilde{N} = (n, \delta, \gamma)_{LR}$ دو عدد فازی LR باشند و $\lambda \in \mathbb{R}$ یک اسکالر باشد، آنگاه

$$\begin{aligned} a) \lambda \otimes (m, \alpha, \beta)_{LR} &= (\lambda m, \lambda \alpha, \lambda \beta)_{LR} & \lambda > 0 \\ b) \lambda \otimes (m, \alpha, \beta)_{LR} &= (\lambda m, -\lambda \beta, -\lambda \alpha)_{RL} & \lambda < 0 \\ c) \tilde{M} \oplus \tilde{N} &= (m + n, \alpha + \delta, \beta + \gamma)_{LR} \\ d) \tilde{M} \ominus \tilde{N} &= (m - n, \alpha + \delta, \beta + \gamma)_{RL} \end{aligned}$$

نکته ۱.۲ اگر $\tilde{M} = (m, \alpha, \beta)$ و $\tilde{N} = (n, \delta, \gamma)$ دو عدد فازی LR باشند، آنگاه ممکن است که ضرب یا تقسیم آن‌ها، یک عدد فازی LR نباشد.

۳ بیان مساله

فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی از جامعه‌ای با توزیع نرمال با میانگین μ و واریانس σ^2 باشد. در حالتی که مقادیر X_i ها به صورت اعداد دقیق باشند با شیوه‌های متداول برای آزمون فرض درباره‌ی میانگین جامعه آشنا هستیم. به طور متداول آزمون فرض درباره‌ی میانگین جامعه در حالتی که واریانس جامعه معلوم است، بر پایه‌ی آماره $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$ و در حالتی که واریانس جامعه مجهول است بر پایه‌ی آماره $Z = \frac{\bar{X} - \mu_0}{\frac{\tilde{\sigma}}{\sqrt{n}}}$ انجام می‌شود [5].

اکنون فرض کنید که مشاهدات مربوط به X_i ها دقیق نباشند (یا دقیق گزارش نشده باشند) بلکه به صورت اعداد فازی باشند. می‌خواهیم شیوه‌های متداول برای آزمون فرض درباره‌ی μ را به حالتی که گفته شد، به صورت یک آزمون تصادفی شده، تعمیم دهیم. برای سادگی محاسبات فرض می‌کنیم مشاهدات به صورت اعداد فازی دوزنقه‌ای هستند. نخست حالتی را که σ معلوم است در نظر می‌گیریم، آنگاه به حالتی می‌پردازیم که σ مجهول است.

۴ تشریح شیوه‌ی کار

۱.۴ حالت σ معلوم

در این حالت آماره‌ی $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$ به شکل زیر تبدیل می‌شود

$$\tilde{Z} \approx \frac{\frac{1}{n} \sum_{i=1}^n \tilde{X}_i - \mu_0}{\frac{\sigma}{\sqrt{n}}} \quad (1)$$

چنانچه مشاهدات اعداد فازی دوزنقه‌ای باشند، با توجه به اصل گسترش، \tilde{Z} یک عدد فازی دوزنقه‌ای می‌باشد [7, 14]. نخست تابع عضویت آماره فوق را بر اساس α -برش‌های \tilde{X}_i به

دست می‌آوریم، و سپس از آن برای انجام آزمون استفاده می‌کنیم. α -برش‌های \tilde{X}_i را به صورت $(\tilde{X}_i)_\alpha = [(\tilde{X}_i)_\alpha^L, (\tilde{X}_i)_\alpha^U]$ نشان می‌دهیم که در آن

$$(\tilde{X}_i)_\alpha^L = [\min_{x_i} \{x_i \in X \text{ s.t. } \mu_{\tilde{X}_i}(x_i) \geq \alpha\}]$$

و

$$(\tilde{X}_i)_\alpha^U = [\max_{x_i} \{x_i \in X \text{ s.t. } \mu_{\tilde{X}_i}(x_i) \geq \alpha\}]$$

فرض کنید

$$(\tilde{Z})_\alpha^L = \frac{(\tilde{X})_\alpha^L - \mu_0}{\frac{\sigma}{\sqrt{n}}}, \quad (\tilde{Z})_\alpha^U = \frac{(\tilde{X})_\alpha^U - \mu_0}{\frac{\sigma}{\sqrt{n}}} \quad (۲)$$

که در آن $(\tilde{X})_\alpha^L = \frac{1}{n} \sum_{i=1}^n (\tilde{X}_i)_\alpha^L$ و $(\tilde{X})_\alpha^U = \frac{1}{n} \sum_{i=1}^n (\tilde{X}_i)_\alpha^U$. بنابراین α -برش‌های \tilde{Z} به صورت $(\tilde{Z})_\alpha = [(\tilde{Z})_\alpha^L, (\tilde{Z})_\alpha^U]$ است، که $(\tilde{Z})_\alpha^L$ و $(\tilde{Z})_\alpha^U$ توابع خطی بر حسب α -برش‌های \tilde{X} هستند. بنابراین تابع عضویت \tilde{Z} عبارت است از

$$\mu_{\tilde{Z}}(z) = \begin{cases} L(z) & \tilde{Z}_0^L \leq z \leq (\tilde{Z})_1^L \\ 1 & (\tilde{Z})_1^L \leq z \leq (\tilde{Z})_1^U \\ R(z) & (\tilde{Z})_1^U \leq z \leq (\tilde{Z})_0^U \end{cases} \quad (۳)$$

که در آن توابع مرجع $L(z)$ و $R(z)$ به صورت زیر هستند

$$L(z) = \frac{z - (\tilde{Z})_1^L}{(\tilde{Z})_1^L - (\tilde{Z})_0^L}, \quad R(z) = \frac{(\tilde{Z})_0^U - z}{(\tilde{Z})_0^U - (\tilde{Z})_1^U}$$

تذکر ۱.۴ اگر تمام داده‌های فازی به مقادیر حقیقی تبدیل شوند، معادلات (۲) ساده شده و به معادله حالت کلاسیک تبدیل می‌شوند.

۲.۴ حالت σ مجهول

در این حالت، آماره‌ی $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$ به صورت زیر تبدیل می‌شود

$$\tilde{T} \approx \sqrt{n} \frac{(\frac{1}{n} \sum_{i=1}^n \tilde{X}_i) - \mu_0}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (\tilde{X}_i - \frac{1}{n} \sum_{i=1}^n \tilde{X}_i)^2}} \quad (۴)$$

در این جا نیز فرض می‌کنیم مشاهدات، به صورت اعداد فازی دوزنقه‌ای باشند. ابتدا تابع عضویت آماره فوق را بر اساس α -برش‌ها به صورت زیر به دست می‌آوریم، سپس از آن برای انجام آزمون استفاده می‌کنیم. فرض کنید

$$\tilde{T}_\alpha^L = \min_{(\tilde{X}_i)_\alpha^L \leq X_i \leq (\tilde{X}_i)_\alpha^U} \left\{ \sqrt{n} \frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu_0}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \frac{1}{n} \sum_{i=1}^n X_i)^2}} \right\} \quad (1-5)$$

$$\tilde{T}_\alpha^U = \max_{(\tilde{X}_i)_\alpha^L \leq X_i \leq (\tilde{X}_i)_\alpha^U} \left\{ \sqrt{n} \frac{\frac{1}{n} \sum_{i=1}^n X_i - \mu_0}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \frac{1}{n} \sum_{i=1}^n X_i)^2}} \right\} \quad (2-5)$$

که $(\tilde{X}_i)_\alpha^L$ و $(\tilde{X}_i)_\alpha^U$ را بر پایه‌ی نمونه می‌توان محاسبه نمود. بنابراین α -برش‌های \tilde{T} به صورت $\tilde{T}_\alpha = [\tilde{T}_\alpha^L, \tilde{T}_\alpha^U]$ است، که در آن \tilde{T}_α^L و \tilde{T}_α^U توابع غیرخطی بر حسب مشاهدات می‌باشند. برای محاسبه \tilde{T}_α می‌توان از نرم‌افزارهای محاسباتی استفاده کرد. بنابراین، با استفاده از مقادیر مختلف α ، تابع عضویت \tilde{T} که به صورت عدد $(\tilde{T}_0^L, \tilde{T}_1^L, \tilde{T}_1^U, \tilde{T}_0^U)_{LR}$ می‌باشد، به صورت زیر به دست می‌آید

$$\mu_{\tilde{T}}(t) = \begin{cases} 0 & t < \tilde{T}_0^L \text{ or } t > \tilde{T}_0^U \\ L(t) & \tilde{T}_0^L \leq t \leq \tilde{T}_1^L \\ 1 & \tilde{T}_1^L \leq t \leq \tilde{T}_1^U \\ R(t) & \tilde{T}_1^U \leq t \leq \tilde{T}_0^U \end{cases} \quad (6)$$

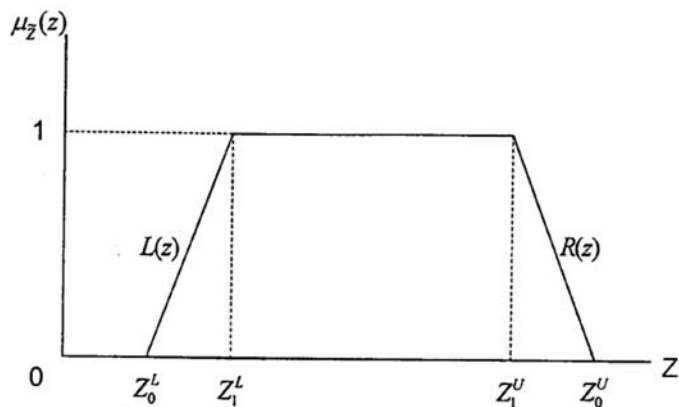
که $L(t), R(t)$ توابع مرجع می‌باشند. نکته ۱.۴ زمانی که همه‌ی داده‌ها مقادیر دقیق باشند روابط a, b یکی و به آماره حالت کلاسیک تبدیل می‌شوند.

۵ آزمون میانگین جامعه (σ معلوم)

همان طور که اشاره شد در این حالت آماره‌ی آزمون به صورت زیر خواهد بود

$$\tilde{Z} = \frac{\frac{1}{n} \sum_{i=1}^n \tilde{X}_i - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

تحت این فرض که مشاهدات \tilde{X}_i اعداد فازی دوزنقه‌ای هستند، مقدار \tilde{Z} یک عدد فازی دوزنقه‌ای LR و به صورت $\tilde{Z} = (\tilde{Z}_0^L, \tilde{Z}_1^L, \tilde{Z}_1^U, \tilde{Z}_0^U)_{LR}$ است [14]، که در آن $\tilde{Z}_0^L = \frac{\tilde{X}_0^L - \mu_0}{\frac{\sigma}{\sqrt{n}}}$



شکل ۱:

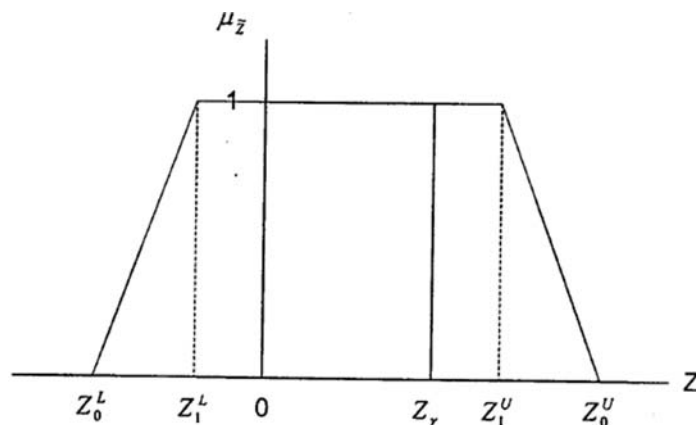
و $\tilde{Z}_1^L = \frac{\bar{X}_1^L - \mu_0}{\frac{\sigma}{\sqrt{n}}}$ و $\tilde{Z}_1^U = \frac{\bar{X}_1^U - \mu_0}{\frac{\sigma}{\sqrt{n}}}$ و $\tilde{Z}_1^U = \frac{\bar{X}_1^U - \mu_0}{\frac{\sigma}{\sqrt{n}}}$ را نشان می‌دهد.

تعریف ۱.۵ در شیوه‌ی آزمون تصادفی شده بر اساس داده‌های فازی، احتمال رد فرض صفر را به صورت احتمال قرار داشتن آماره آزمون در ناحیه بحرانی تقسیم بر احتمال مرتبط با آماره آزمون تعریف می‌کنیم و با P_0 نشان می‌دهیم. همچنین احتمال پذیرش فرض صفر را به صورت $1 - P_0$ در نظر می‌گیریم. در حالت خاص، $P_0 = 0$ نشان دهنده‌ی عدم رد فرض صفر و $P_0 = 1$ نشان دهنده‌ی رد قطعی فرض صفر می‌باشد.

تعریف ۲.۵ احتمال مرتبط با آماره آزمون فازی \tilde{Z} به صورت زیر تعریف می‌شود

$$\Delta = \int_{\tilde{Z}_1^L}^{\tilde{Z}_1^L} L(z)\phi(z)dz + \int_{\tilde{Z}_1^L}^{\tilde{Z}_1^U} \phi(z)dz + \int_{\tilde{Z}_1^U}^{\tilde{Z}_1^U} R(z)\phi(z)dz \quad (Y)$$

که در آن $\phi(z)$ تابع چگالی احتمال توزیع نرمال استاندارد Z است. فرض کنید که می‌خواهیم فرض صفر $H_0: \mu \leq \mu_0$ در برابر فرض مقابل $H_a: \mu > \mu_0$ آزمون کنیم. شکل ۲ رابطه‌ی بین تابع عضویت آماره آزمون فازی \tilde{Z} و ناحیه بحرانی $\{z: z > z_\gamma\}$ را، که γ سطح معنی داری آزمون و z_γ مقدار بحرانی است، نشان می‌دهد.



شکل ۲:

با توجه به شکل ۲ احتمال رد فرض صفر به صورت زیر به دست می‌آید

$$P_0 = \frac{1}{\Delta} \left[\int_{\tilde{z}_\gamma}^{\tilde{z}_1^U} \phi(z) dz + \int_{\tilde{z}_1^U}^{\tilde{z}_1^U} R(z) \phi(z) dz \right] \quad (\lambda)$$

در شکل ۳، پنج نمونه ممکن از توابع عضویت \tilde{Z} ، مربوط به آزمون یکطرفه راست فرض صفر $H_0: \mu \leq \mu_0$ در برابر فرض مقابل $H_a: \mu > \mu_0$ نشان داده شده است. با در نظر گرفتن ناحیه بحرانی $\{\tilde{Z} \text{ s.t. } \tilde{Z} > \tilde{z}_\gamma\}$ و بر اساس روابط (۷) و (۸) احتمال رد فرض صفر در حالت‌های a تا e عبارت است از

$$a) P_0 = 0$$

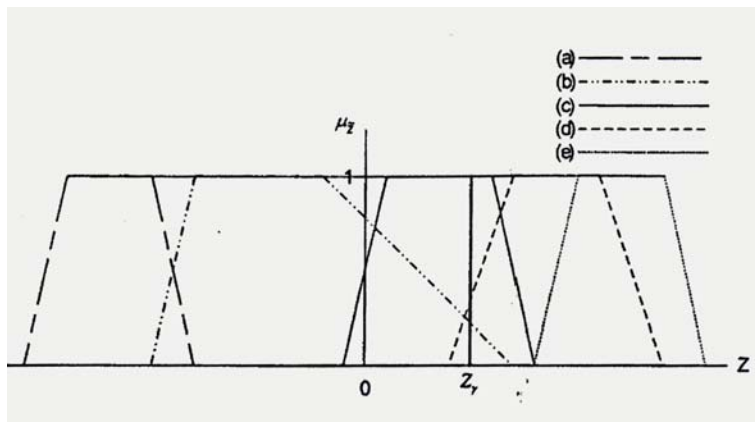
$$b) P_0 = \frac{1}{\Delta} \left[\int_{\tilde{z}_\gamma}^{\tilde{z}_1^U} R(z) \phi(z) dz \right]$$

$$c) P_0 = \frac{1}{\Delta} \left[\int_{\tilde{z}_\gamma}^{\tilde{z}_1^U} \phi(z) dz + \int_{\tilde{z}_1^U}^{\tilde{z}_1^U} R(z) \phi(z) dz \right]$$

$$d) P_0 = \frac{1}{\Delta} \left[\int_{\tilde{z}_\gamma}^{\tilde{z}_1^L} L(z) \phi(z) dz + \int_{\tilde{z}_1^L}^{\tilde{z}_1^U} \phi(z) dz + \int_{\tilde{z}_1^U}^{\tilde{z}_1^U} R(z) \phi(z) dz \right]$$

$$e) P_0 = 1$$

تذکر ۱.۵ تعریف مقدار احتمال P_0 در حالت یکطرفه چپ شبیه به حالت یکطرفه راست است. تذکر ۲.۵ از روش فوق می‌توان برای آزمون فرض‌های ساده در برابر مرکب نیز استفاده نمود [7].



شکل ۳:

مثال ۱.۵ فرض کنید $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_8$ ، هشت مشاهده فازی از توزیع نرمال $N(\mu, 1/618)$ به صورت اعداد فازی مثلثی متقارن زیر باشند

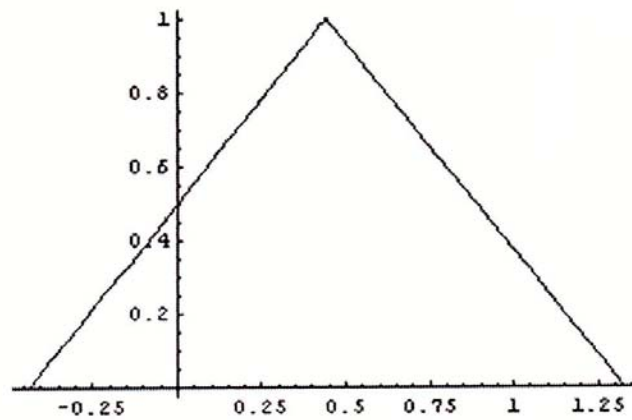
$$\begin{aligned} \tilde{X}_1 &= (3/08, 0/3, 0/3)_T, & \tilde{X}_2 &= (2/86, 0/3, 0/3)_T \\ \tilde{X}_3 &= (6/25, 0/6, 0/6)_T, & \tilde{X}_4 &= (4/11, 0/4, 0/4)_T \\ \tilde{X}_5 &= (2/71, 0/2, 0/2)_T, & \tilde{X}_6 &= (4/45, 0/4, 0/4)_T \\ \tilde{X}_7 &= (5/05, 0/5, 0/5)_T, & \tilde{X}_8 &= (5/23, 0/5, 0/5)_T \end{aligned}$$

می‌خواهیم آزمون یکطرفه فرض $H_0: \mu \leq 4$ را در برابر فرض مقابل $H_a: \mu > 4$ انجام دهیم. با توجه به این‌که هر عدد فازی مثلثی حالت خاصی از یک عدد فازی ذوزنقه‌ای است بنابراین تمام مشاهدات فوق را می‌توان نوعی اعداد فازی ذوزنقه‌ای در نظر گرفت. در این مثال $\sigma^2 = 1/618$ ، بنابراین آماره آزمون به شکل زیر می‌باشد

$$\tilde{Z} = \frac{\frac{1}{n} \sum_{i=1}^n \tilde{X}_i - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$

با کمک معادلات (۲) و (۳)، تابع عضویت \tilde{Z} عبارت است از

$$\mu_{\tilde{Z}}(z) = \begin{cases} \frac{1}{0/88} (z + 0/44) & -0/44 \leq z \leq 0/44 \\ \frac{1}{0/88} (1/33 - z) & 0/44 \leq z \leq 1/33 \end{cases} \quad (9)$$



شکل ۴:

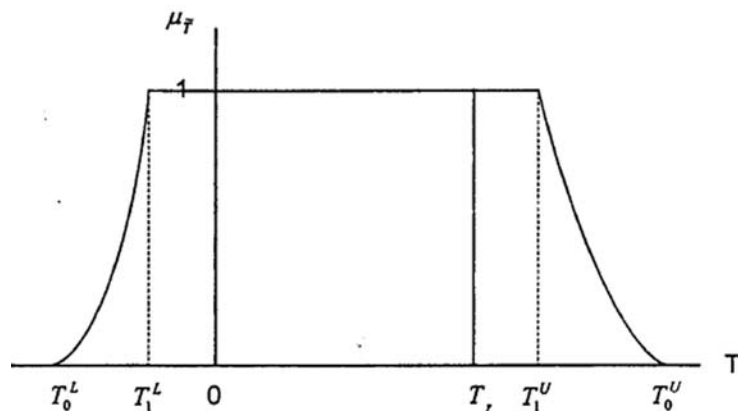
در شکل ۴ نمودار عدد فازی مثالی فوق رسم شده است. تکیه گاه آماره آزمون فازی، \tilde{Z} ، بین $1/33$ و $-0/44$ می باشد. این دامنه نشان دهنده این مطلب است که اگر چه آماره آزمون، یک عدد فازی است، اما مقدارش نمی تواند خارج از بازه $[1/33, -0/44]$ باشد. با توجه به اینکه ناحیه بحرانی این آزمون یکطرفه در سطح معنی داری $0/1$ به صورت $\{\tilde{Z} \text{ s.t. } \tilde{Z} > \tilde{Z}_{0/1}\}$ می باشد، با کمک معادله (۷)، احتمال مربوط به آماره آزمون فازی، \tilde{Z} ، به صورت زیر محاسبه می شود

$$\Delta = \int_{-0/44}^{0/44} \frac{1}{0/88} (z + 0/44) \phi(z) dz + \int_{0/44}^{1/33} \frac{1}{0/89} (1/33 - z) \phi(z) dz = 0/3042$$

سرانجام، با کمک معادله (۸)، احتمال رد فرض صفر عبارت است از

$$P_0 = \frac{1}{0/3042} \left[\int_{1/28}^{1/33} \frac{1}{0/89} (1/33 - z) \phi(z) dz \right] = 0/0008$$

بنابراین فرض صفر با احتمال $0/0008$ رد و با احتمال $0/9992$ پذیرفته می شود.



شکل ۵:

۶ آزمون میانگین جامعه (σ مجهول)

در این حالت، با توجه به مطالب بخش چهارم، داریم

$$\tilde{T} \approx \sqrt{n} \frac{(\frac{1}{n} \sum_{i=1}^n \tilde{X}_i) - \mu_0}{\sqrt{\frac{1}{n-1} \sum_{i=1}^n (\tilde{X}_i - \frac{1}{n} \sum_{i=1}^n \tilde{X}_i)^2}}$$

که در آن

$$S^2 = \frac{\sum_{i=1}^n (\tilde{X}_i - \frac{1}{n} \sum_{i=1}^n \tilde{X}_i)^2}{n-1}$$

توجه کنید که، عطف به نکته ۱.۲، رابطه‌ی موجود بین آماره‌ی آزمون فازی و داده‌های فازی در آماره \tilde{T} غیر خطی می‌باشد. بنابراین به دست آوردن تابع عضویت $\mu_{\tilde{T}}$ به طور دقیق تقریباً غیر ممکن است.

فرض کنید می‌خواهیم فرض صفر $H_0: \mu \leq \mu_0$ را در برابر فرض مقابل $H_a: \mu > \mu_0$ آزمون کنیم. شکل ۵ رابطه‌ی بین تابع عضویت $\mu_{\tilde{T}}$ و ناحیه بحرانی $\{\tilde{T} \text{ s.t. } \tilde{T} > \tilde{T}_\gamma\}$ را توصیف می‌کند، که در آن سطح معنی داری آزمون و مقدار بحرانی می‌باشد.

تعریف ۱.۶ احتمال مرتبط با آماره آزمون فازی، \tilde{T} ، به صورت زیر تعریف می‌شود

$$\Delta = \int_{\tilde{T}^L}^{\tilde{T}^U} L(t)p(t)dt + \int_{\tilde{T}^L}^{\tilde{T}^U} p(t)dt + \int_{\tilde{T}^U}^{\tilde{T}^U} R(t)p(t)dt \quad (۱۰)$$

که در آن $p(t)$ تابع چگالی احتمال استیودنت می‌باشد. در این صورت احتمال رد فرض صفر با توجه به شکل ۵ به صورت زیر تعریف می‌شود

$$P_0 = \frac{1}{\Delta} \left[\int_{\tilde{T}_\gamma}^{\tilde{T}^U} p(t)dt + \int_{\tilde{T}_\gamma}^{\tilde{T}^U} R(t)p(t)dt \right] \quad (۱۱)$$

تذکر ۱.۶ احتمال رد فرض صفر برای آزمون‌های یکطرفه و دو طرفه به طور مشابه تعریف می‌شوند.

در شکل ۶، تعداد ۵ نمونه ممکن از توابع عضویت \tilde{T} ، مربوط به آزمون یکطرفه راست فرض صفر $H_0: \mu \leq \mu_0$ در برابر فرض مقابل $H_a: \mu > \mu_0$ نشان داده است. با در نظر گرفتن ناحیه بحرانی $\{\tilde{T} \mid \tilde{T} > \tilde{T}_\gamma\}$ و بر اساس روابط (۱۰) و (۱۱) احتمال رد فرض صفر در حالت‌های a تا e به صورت زیر به دست می‌آیند

$$a) P_0 = 0$$

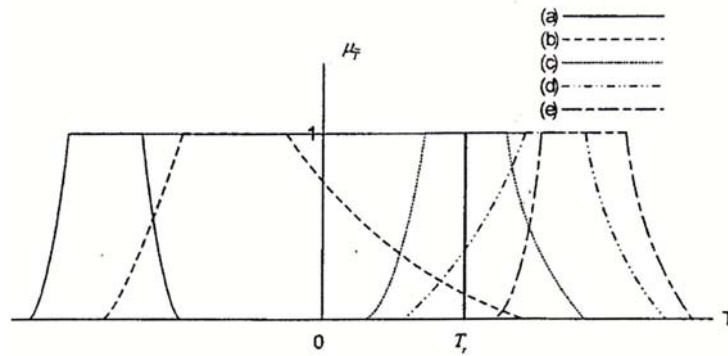
$$b) P_0 = \frac{1}{\Delta} \left[\int_{\tilde{T}_\gamma}^{\tilde{T}^U} R(t)p(t)dt \right]$$

$$c) P_0 = \frac{1}{\Delta} \left[\int_{\tilde{T}_\gamma}^{\tilde{T}^U} p(t)dt + \int_{\tilde{T}_\gamma}^{\tilde{T}^U} R(t)p(t)dt \right]$$

$$d) P_0 = \frac{1}{\Delta} \left[\int_{\tilde{T}_\gamma}^{\tilde{T}^L} L(t)p(t)dt + \int_{\tilde{T}_\gamma}^{\tilde{T}^U} p(t)dt + \int_{\tilde{T}_\gamma}^{\tilde{T}^U} R(t)p(t)dt \right]$$

$$e) P_0 = 1$$

چون بر اساس معادله (۴) رابطه‌ی بین آماره آزمون و داده‌های فازی غیر خطی است، بنابراین به دست آوردن $L(t)$ و $R(t)$ به طور دقیق بسیار دشوار است. اما می‌توان از روش‌های عددی استفاده نمود و مقادیر تقریبی برای $\int_{\tilde{T}_\gamma}^{\tilde{T}^U} R(t)p(t)dt$ ، $\int_{\tilde{T}_\gamma}^{\tilde{T}^L} L(t)p(t)dt$ ، $\int_{\tilde{T}_\gamma}^{\tilde{T}^U} p(t)dt$ به دست آورد. آنگاه مطابق شیوه‌ی بیان شده در حالت σ معلوم عمل نموده و احتمال رد فرض صفر را محاسبه می‌کنیم.



شکل ۶:

مراجع

- [1] Buckley, J.J., Eslami, E. (2002), An Introduction to Fuzzy Logic and Fuzzy Sets, Springer.
- [2] Casals, R., Gil, M.A., Gil, P. (1986), On the use of Zadeh' probabilistic definition for testing statistical hypotheses from fuzzy information, Fuzzy Sets and Systems, 20: 175-190.
- [3] Casals, M.R., Gil, M.A., Gil, P. (1986), The fuzzy decision problem: an approach to the problem of testing statistical hypothesis with fuzzy information, Euro. J. Oper. Res., 27: 371-382.
- [4] Casals, R., Gil, M.A. (1989), A note on the operativeness of Neyman-Pearson tests with fuzzy information, Fuzzy Sets and Systems, 30: 215-220.
- [5] Casella, G., Berger, R.L. (2002), Statistical Inference, Sec. Ed., Buxbury.
- [6] Grzegorzewski, P. (2001), Fuzzy tests- defuzzification and randomization, Fuzzy Sets and Systems, 118: 437- 446.
- [7] Leu, C.H., Chen, C.C. (2004), Randomized tests of mean for normal population with fuzzy data, Preprint.
- [8] Montenegro, M., Casals, M.R., Lubiano, M.A., Gil, M.A. (2001), Two-sample hypothesis tests of means of a fuzzy random variable, Information Sciences, 133: 89- 100.
- [9] Romer, C., Kandel, A. (1995), Statistical tests for fuzzy data, Fuzzy Sets and Systems, 72: 1- 26.

- [10] Son, J.C., Song, I., Kim, H.Y. (1992), A fuzzy decision problem based on the generalized Neyman- Pearson criterion, *Fuzzy Sets and Systems*, 47: 65- 75.
- [11] Taheri, S.M., Behboodian, J. (2004), On Bayesian approach to fuzzy hypothesis testing with fuzzy data, *Ital. J. Pure and Appl. Math.*, 22.
- [12] Watanabe, N., Imaizumi, T. (1993), A fuzzy statistical test of fuzzy hypotheses, *Fuzzy Sets and Systems*, 53: 167- 178.
- [13] Zadeh, L.A. (1968), Probability measures of fuzzy events, *J. Math. Anal. Appl.*, 23: 421- 427.
- [14] Zimmerman, H.J. (1996), *Fuzzy Set Theory and Its Applications*, 3rd Edition, Kluwer, Boston.

تعیین جایگاه اقتصادی استان‌های کشور - ۱۳۷۹ از نظر ارزش افزوده و درصد نیروی شاغل در بخش‌های اقتصادی به روش تحلیل عاملی

وحید طیفوری

سازمان مدیریت و برنامه‌ریزی استان گیلان معاونت آمار و اطلاعات

چکیده: دو شاخص مهم در ارزیابی و مقایسه نحوه توزیع سرمایه و سایر عوامل اقتصادی در بین استان‌های کشور، ارزش افزوده ایجاد شده و درصد نیروی شاغل در هر یک از بخش‌های اقتصادی می‌باشد. بررسی هر یک از استان‌ها از بعد این دو شاخص بطور مجزا و توأم ضمن فراهم نمودن امکان تحلیل‌های اقتصادی، ما را در شناخت امکانات، پتانسیل‌ها، نقاط قوت و ضعف استانها، یاری نموده و ابزاری مناسبی برای سیاست‌گذاری‌های هدفمند در نحوه توزیع سرمایه و امکانات، بمنظور کاهش نابرابری‌های اقتصادی می‌باشد.

در این مقاله، استان‌های کشور از نظر ارزش افزوده ایجاد شده و درصد نیروی شاغل در ۱۵ بخش اقتصادی^۲ در سال ۱۳۷۹ بعنوان متغیرهای مورد مطالعه بروش تحلیل عاملی^۳ در دو بخش مورد بررسی قرار گرفته است. طبق نتایج بدست آمده در تحلیل ارزش افزوده بر اساس آزمون‌های انجام شده سه عامل با تبیین بیش از ۹۱ درصد واریانس کل پیشنهاد گردید که پس از محاسبات لازم، رتبه‌بندی استانها برای هر یک از سه عامل، مشخص گردیده است. در بررسی درصد نیروی شاغل در هر یک از بخش‌های اقتصادی نیز با استفاده از همین روش پس از حذف برخی متغیرها با برازش مدل مناسب در ۱۱ بخش (متغیر) باقیمانده، مدل سه عاملی با تبیین حدود ۷۹ درصد واریانس پیشنهاد گردید و پس از محاسبه ضرایب عامل، استانها بر حسب امتیاز بدست آمده رتبه‌بندی شدند و در پایان هر بخش تحلیل کلی از این بررسی بعمل آمده است.

واژه‌های کلیدی: تحلیل عاملی، ارزش افزوده، درصد نیروی شاغل

۱ بخش اول

۱.۱ مقدمه

یکی از دستاوردهای بسیار مهم مرکز آمار ایران که با همکاری گسترده کارشناسان سازمان مدیریت و برنامه‌ریزی کلیه استانها طی سالهای اخیر به تحقق پیوسته، تهیه و تدوین حسابهای منطقه‌ای

(۱) آخرین نشریه منتشر شده حساب تولید استانها مربوط به سال ۱۳۷۹ می‌باشد.

2) بر اساس ISIC Rev.3 3) Factor Analysis

(حساب تولید استان‌ها) بوده است. این طرح با در نظر گرفتن تعاریف و مفاهیم، طبقه‌بندیها و جداول استاندارد کلیه فعالیت‌های اقتصادی را در هر استان مورد پوشش قرار داده و در این چارچوب اطلاعات مورد نیاز را جمع‌آوری و حساب اقتصادی استانها با استفاده از روش تولید محاسبه می‌گردد. مهمترین شاخصهای استنتاج شده از این طرح محاسبه محصول ناخالص داخلی^۴ و ارزش افزوده هر یک از فعالیت‌های اقتصادی در سطح استان می‌باشد. این حسابها بعنوان ابزار لازم در مقوله برنامه‌ریزی منطقه‌ای، یکی از زیر ساخت‌ها محسوب شده و مبنای موثقی در ارزیابی عملکرد نحوه توزیع سرمایه و امکانات در کشور و حرکتی برای امکان تدوین برنامه‌های لازم بمنظور رفع کاستیهای توسعه اقتصادی در کشور می‌باشد.

بمنظور سهولت در تحلیل سازوکارهای اقتصادی توسط برنامه‌ریزان تعیین جایگاه اقتصادی استانها از بعد دو شاخص ارزش افزوده و درصد نیروی شاغل در هر یک از بخش‌های اقتصادی بصورت مجزا و نوام، می‌تواند در شناخت امکانات و پتانسیلها، نقاط قوت و ضعف استانها، مورد ارزیابی قرار گیرد.

بنابر اهمیت موضوع در گزارش حاضر استان‌های کشور در سال ۱۳۷۹ از دو بعد ارزش افزوده و درصد نیروی کار در هـ---ریک از فعالیتهای اقتصادی^۵ بروش تحلیل عاملی، در دو بخش بصورت مجزا رتبه‌بندی و مورد بررسی قرار گرفته و در بخش آخر تحلیلی بر رابطه این دو شاخص اقتصادی با یکدیگر بعمل آمده است. با توجه به اهمیت شاخصهای یاد شده و کاربرد متنوع آنها در تحلیل‌های اقتصادی، گزارش حاضر از جمله اهداف زیر را می‌تواند پوشش دهد.

-- فراهم نمودن زمینه لازم در خصوص مطالعه منابع اقتصادی یا شناسایی و تعیین محورهای توسعه استانها.

-- شناخت مهمترین فعالیت (فعالیتها) در هر استان.

-- تعیین نقش و جایگاه هر استان در ساختار اقتصادی کشور.

-- مبنای تحقیقات تفصیلی در نحوه توزیع سرمایه و امکانات در جهت حذف فاصله‌های عمیق اقتصادی و همچنین تمرکز یا جابجایی سرمایه و نیروی کار در بخشهای مستعد و یا پر بازده.

۲.۱ روشهای رتبه‌بندی

الف -- روش وزن‌دهی

یکی از روشهایی که به رتبه‌بندی فعالیتها بر اساس برخی صفات که به صورت شاخصهای کمی یا کیفی بیان می‌شوند روش وزنی است. در این روش هر فعالیت به تناسب مقداری از صفت که به خود اختصاص می‌دهد امتیاز می‌گیرد. و در نهایت، جمع امتیازها، وضعیت فعالیت را نسبت به سایر فعالیتها مشخص می‌کند.

ویژگیهای این روش: سهولت و قابلیت کاربرد در تحقیقات مختلف، امکان استفاده از شاخصهای

4) Gross Domestic Products Region (GDPR) 5) International Standard Industrial Classification Of All Economic Active (ISIC Rev3)

کیفی، عدم لزوم همسویی شاخصها با یکدیگر و محدودیت‌های آن تورش نسبت به شاخصهایی که بایکدیگر همبستگی دارند، حساسیت زیاد این روش نسبت به چگونگی وزن دهی به شاخصها می‌باشد.

ب -- روش تاکسونومی عددی

تاکسونومی عددی روشی است برای طبقه‌بندی و رتبه‌بندی موضوعات، به گونه‌ای که بین عناصر هر طبقه بیشترین تشابه وجود داشته باشد و در عین حال بیشترین اختلاف را با عناصر طبقات دیگر داشته باشد. به این طریق می‌توان مجموعه عناصر همگن را در یک طبقه قرار داد و سپس بر اساس امتیاز کسب شده رتبه‌بندی را انجام داد. از مزایای این روش می‌توان به سهولت در انجام محاسبات و امکان محاسبه شاخصهای توسعه یافتگی اشاره نمود اما محدودیتها این روش لزوم تعیین شاخصها با توجه به انتخاب و نظر کارشناسان هر رشته مطالعاتی و همچنین لزوم همسویی شاخصها می‌باشد.

ج -- روش میزان انحراف از بهینه

این روش برای رتبه‌بندی جوامع از نظر شاخصهای کمی مورد استفاده قرار می‌گیرد. مبنای این روش بر استفاده از اعداد و ارقام خام شاخصها و استاندارد نمودن آنها می‌باشد. از ویژگیهای این روش، سهولت در محاسبات و سادگی در درک آن توسط کاربران، حذف مقیاس شاخصها بدلیل استاندارد نمودن آنها، تعیین میزان انحراف شاخص از مقدار بهینه بصورت یک نسبت و از محدودیتهای این روش می‌توان به لزوم همسویی شاخصها، حساسیت نسبت به شاخصهای انتخابی، تورش نسبت به همبستگی بین شاخصها، اهمیت یکسان کلیه شاخصها نسبت به یکدیگر اشاره نمود.

د -- روش تحلیل عاملی

کلی‌ترین شکل الگوی عاملی امکان اندازه‌گیری متغیرهایی که نمی‌توانند اندازه‌گیری شوند را فراهم می‌نماید. الگوی عاملی مجموعه‌ای از معادلات خطی هستند که به ویژه در علوم رفتاری، تصمیم در خصوص سوددهی شرکتها، ممیزی در استخدامها، اقتصاد و ... کاربرد دارد. این روش عموماً برای خلاصه نمودن اطلاعات مورد استفاده قرار می‌گیرد. بدین ترتیب با تلخیص و تلفیق متغیرهای مختلف اثر آنها را در قالب چند عامل که بیشترین اختلاف را با یکدیگر دارند ارایه می‌شود. از این روش به منظور طبقه‌بندی، رتبه‌بندی و تجزیه و تحلیل در علوم مختلف استفاده می‌شود. از مزایای این روش می‌توان به استفاده از روشها و تکنیک‌های منطقی و پیشرفته برای برآورد مقادیر عاملها (نظیر روش کمترین توانهای دوم موزون و یا روش رگرسیون)، امکان استاندارد کردن متغیرها، استفاده از تکنیک دوران به منظور سهولت در تفسیر عوامل اشاره نمود. محدودیت‌های این روش عبارت است از نیاز به تصمیم‌گیری‌های زیاد در مطالعه تحلیل عاملی (نظیر تعیین تعداد عاملها، تشخیص دوران و ...) حجم زیاد محاسبات، حساسیت نسبت به چگونگی انتخاب و کیفیت اطلاعات و همچنین عدم تضمین در وجود عوامل معنی‌دار در تمامی تحلیل‌های آماری داده‌ها.

۳.۱ تعریف متغیرها

متغیرهای در نظر گرفته شده در بخش دوم عبارتند از ارزش افزوده ایجاد شده در هر یک از فعالیت‌های اقتصادی سال ۱۳۷۹، به تفکیک استان و در قالب ماتریسی تعریف شده است، که عنصر a_{ij} این ماتریس بیانگر ارزش افزوده فعالیت اقتصادی j ام ($j = 1, \dots, 15$) در استان i ام ($i = 1, \dots, 28$) می‌باشد. متغیرهای بخش سوم عبارت است از برآورد درصد نیروی شاغل در هر یک از فعالیت‌های اقتصادی به تفکیک استان در سال ۱۳۷۹، که در قالب یک ماتریس تعریف شده است که عنصر b_{ij} این ماتریس بیانگر برآورد درصد نیروی شاغل در فعالیت اقتصادی j ام ($j = 1, \dots, 15$) در استان i ام ($i = 1, \dots, 28$) می‌باشد.

۴.۱ مراحل کلی انجام کار بروش تحلیل عاملی

الف -- تعیین ماتریس ضرایب همبستگی بین متغیرها و تفسیر آن: ماتریس همبستگی و آزمونهای منتج شده از آن ما را در تایید و یا رد استفاده از روش تحلیل عاملی یاری می‌کند. مهمترین شاخص مورد استفاده، شاخص KMO^6 می‌باشد که برای مقایسه مقادیر ضرایب همبستگی ساده و جزئی بروی کلیه متغیرها مورد استفاده قرار می‌گیرد و مقادیر کمتر از ۰/۵ آن عدم کارایی این روش را بیان می‌کند و رابطه آن عبارت است از:

$$KMO = \frac{\sum \sum_{i \neq j} r_{ij}^2}{\sum \sum_{i \neq j} r_{ij}^2 + \sum \sum_{i \neq j} d_{ij}^2}$$

که در آن

r_{ij} : ضریب همبستگی ساده

d_{ij} : ضریب همبستگی جزئی

ب -- تعیین تعداد عاملها، محاسبه عاملها و بررسی میزان برازش مدل: در صورت تایید استفاده از روش تحلیل عاملی، یکسری عوامل اولیه به یکی از روشهای موجود نظیر: مولفه‌های اصلی^۷، بیشینه درستی^۸، کمترین مربعات ناموزون^۹ و . . . تعیین می‌شود. اغلب انتخاب تعداد عاملها مبتنی بر چند ترکیب: ۱-- میزان تبیین واریانس هر عامل ۲-- دانش مورد بحث ۳-- عقلایی بودن نتایج، است. (البته لازم بذکر است که هر چه میزان تبیین واریانس در عوامل انتخابی بیشتر باشد، کارایی نتایج بهتر خواهد بود). در این گزارش برای تعیین تعداد عاملها از روش مولفه‌های اصلی استفاده شده و با تکیه بر موارد مطرح شده، تعداد عاملهای نهایی تعیین شده است.

6) Kaiser-Meyer-Olkin Measure of Sampling Adequacy 7) Principal Component
8) Maximum Likelihood 9) Unweighted Least Square

ج -- دوران:

هنگامیکه یک یا چند متغیر خاص توسط بیش از یک عامل تبیین شود آنگاه تفسیر عوامل مشکل خواهد بود لذا در شرایط مذکور و بمنظور سهولت در تفسیر، با انجام عمل دوران به یکی از روشهای متعامد یا مایل امکان تفسیر بهتر عوامل فراهم می‌گردد.

د -- تعیین مقدار متناظر هر عامل برای متغیرهای مورد بررسی و تهیه جداول نهایی: در این مرحله مشخص می‌شود که هر یک از عاملهای تعیین شده بر چه متغیر (متغیرهایی) تاکید دارند. پس از انطباق عوامل و با محاسبه ضرایب عاملها و بکارگیری متغیرهای استاندارد شده می‌توانیم بروش زیر امتیاز نهایی هر استان را بر اساس هر یک از عوامل تعیین نماییم.

$$A_{jk} = \sum_{i=1}^q w_{ji} z_{ik}$$

$$j = 1, 2, \dots, m$$

$$k = 1, 2, \dots, n$$

m : تعداد عاملها

n : تعداد استانها

q : تعداد متغیرها

A_{jk} : امتیاز نهایی عامل j ام برای استان k ام

w_{ji} : مقدار عامل j ام برای فعالیت (بخش) i ام

z_{ik} : مقدار استاندارد شده فعالیت (بخش) i ام برای استان k ام

۵.۱ تعاریف و مفاهیم

ارزش افزوده^{۱۰}

عبارت است از تفاوت بین ارزش کالاها و خدمات تولید شده و ارزش کالاها و خدمات بکار رفته در جریان تولید و یا به بیان دیگر ارزش افزوده حاصل ارزش ستاننده منهای ارزش داده است.

ستاننده^{۱۱}

مجموع ارزش کالاها و خدماتی است که در یک واحد تولیدی، تولید شده و برای استفاده در خارج آن واحد در دسترس قرار می‌گیرد.

داده^{۱۲} (مصارف واسطه)^{۱۳}

عبارت است از ارزش کالاها و خدماتی است که به صورت داده یک فرآیند تولید، در آن فرآیند، مصرف می‌شوند. به استثنای مصرف داراییهای ثابت (استهلاک).

شاغل

کلیه کسانی (ده ساله و بیشتر) که در هفت روز گذشته قبل از مراجعه مأمور آمارگیر حداقل دو روز به کاری اشتغال داشته‌اند و یا شاغل بوده‌اند و به عللی (از قبیل مرخصی و یا بیماری و

10) Value Added 11) Output 12) Input 13) Intermediate Comsumption

(... کار نکرده‌اند، شاغل محسوب می‌شوند.

درصد شاغلین

عبارت است از نسبت تعداد شاغلان ده ساله و بیشتر به جمعیت فعال ده ساله و بیشتر (شاغل و بیکار) ضربدر ۱۰۰.

طبقه‌بندی فعالیت‌های اقتصادی

در این طبقه‌بندی فعالیت‌های اقتصادی بر اساس تشابه ساختار تولید و بعضی روابط موجود آنها، از اهمیت استراتژیک برخوردار می‌باشند. که جنبه عمده و جوه تشابه مورد نظر عبارت است از:

-- مشخصه کالاها و خدمات تولید شده.

-- موارد استفاده از کالا و خدمات.

-- داده‌های بکار رفته در تولید، فرآیند تولید و تکنولوژی مورد استفاده در آن.

ساختار کالی طبقه‌بندی استاندارد فعالیت‌های اقتصادی بر اساس ISIC(rev3) در سطح رده‌های الفبایی بشرح زیر می‌باشد:

الف -- کشاورزی، شکار و جنگلداری

ب -- ماهیگیری

پ -- معدن

ت -- صنعت

ث -- تامین برق، گاز و آب

ج -- ساختمان

چ -- عمده فروشی، خرده فروشی و تعمیر وسایل نقلیه موتوری و کالاهای شخصی و خانگی

ح -- هتل و رستوران

خ -- حمل و نقل، انبارداری و ارتباطات

د -- واسطه‌گریهای مالی

ذ -- مستغلات، کرایه و خدمات کسب و کار

ر -- اداره امور عمومی، دفاعی و تامین اجتماعی

ز -- آموزش

ژ -- بهداشت و درمان و مددکاری اجتماعی

س -- سایر فعالیت‌های خدمات عمومی، اجتماعی و شخصی

۲ بخش دوم

تعیین و بررسی تحلیلی جایگاه استانها از نظر ارزش افزوده فعالیت‌های اقتصادی «بروش تحلیل عاملی»

۱.۲ تحلیل همبستگی

نتایج حاصل شده از ماتریس همبستگی ارزش افزوده بخش‌های مختلف اقتصادی به عنوان متغیرهای مورد بررسی نشان می‌دهد که اکثر متغیرهای دارای همبستگی مثبت و بالایی بوده و فقط سه متغیر کشاورزی و ...، ماهیگیری و معدن دارای کمترین میزان همبستگی با سایر متغیرها می‌باشد. با بررسی ماتریس سطح معنی‌داری (ماتریس مذکور به واقع بیانگر میزان احتمال ارتکاب اشتباه نوع اول می‌باشد) آزمون معنی‌دار بودن فرض $H_0: r_{ij} = 0$ برای کلیه متغیرها رد شده و لذا همبستگی بین متغیرها مورد تایید قرار گرفت. در گام بعدی با بررسی اعتبار روش تحلیل عاملی، با استفاده از آزمون‌های موجود به شرح زیر داریم.

با توجه به اینکه دترمینان ماتریس همبستگی برابر $10^{-23} \times 7,432$ بوده و چون به اندازه کافی کوچک است پس روش تحلیل عاملی در این مرحله تایید می‌شود. در بررسی ماتریس باقیمانده‌ها نیز چون بیش از ۹۲ درصد مقادیر این ماتریس کمتر از ۵ درصد می‌باشد و بدین ترتیب برازش مدل تحلیل عاملی پذیرفته می‌شود. همچنین با توجه به عناصر روی قطر ماتریس تصویر معکوس همبستگی که در واقع بیانگر مقدار MSA_i ها برای هر یک از متغیرها می‌باشد و چون بزرگ بودن این عناصر دلیل بر امکان حضور متغیر مورد بررسی در مدل می‌باشد، لذا در نتایج بدست آمده از این ماتریس مشاهده می‌شود که مقدار MSA برای متغیرهای معدن و ماهیگیری نسبتاً پایین بوده اما بدلیلی اهمیت این دو بخش در برخی از استانها و همچنین با توجه به پایین بودن میزان باقیمانده‌ها برای متغیرهای مذکور از حذف آنها خودداری گردید (جدول ۲). مقدار بدست آمده برای شاخص KMO برابر $0,814$ بوده که با توجه به بزرگ بودن مقدار آن، استفاده از روش تحلیل عاملی در این بررسی تایید می‌شود.

۲.۲ تعیین تعداد و روش محاسبه عاملها و میزان برازش مدل

برای تعیین تعداد عاملها از روش مولفه‌های اصلی استفاده شده است و لذا بر اساس مقادیر ویژه بدست آمده برای ۱۵ متغیر، چون مقادیر ویژه سه عامل نخست بزرگتر از یک بوده و بیشتر از ۹۱ واریانس کل را تبیین می‌کند (جدول ۳)، که در نتیجه روش تحلیل عاملی با سه عامل مورد پذیرش واقع می‌شود.

در ادامه بمنظور مشخص نمودن اینکه هر عامل بر چه متغیرهایی تاکید دارد اقدام نموده که در محاسبات اولیه بدلیل عدم تفکیک پذیری عوامل بر حسب متغیرهای مطرح شده انجام عمل دوران ضروری به نظر رسید. با توجه به نتیجه بدست آمده و برای سهولت در تفسیر عمل دوران متعامد بروش واریماکس^{۱۴} انجام شد که نتیجه آن در جدول ۴ مشاهده می شود. بر اساس نتایج بدست آمده، تعبیر هر یک از سه عامل به شرح زیر می باشد:

ویژگی عامل اول: این عامل عمدتاً بر فعالیتهای غیر از کشاورزی، شکار و جنگلداری و معدن و ماهیگیری تکیه دارد.

ویژگی عامل دوم: این عامل صرفاً بر اهمیت بخش کشاورزی، شکار و جنگلداری تاکید می کند.

ویژگی عامل سوم: این عامل با توجه به ضرایب مثبت بزرگ برای فعالیت های معدن و ماهیگیری، بر اهمیت فعالیتهای مورد نظر، دلالت دارد.

۳.۲ جداول نهایی

پس از تعیین ویژگی عاملها، به روش رگرسیونی جدول ضرایب نهایی عوامل برای هر یک از متغیرها در جدول ۵ محاسبه شده است. مقادیر این جدول از رابطه زیر که در آن از متغیرهای استاندارد شده استفاده شده، بدست آمده است.

$$A_{jk} = \sum_{i=1}^{15} w_{ji} z_{ik} \quad j = 1, 2, \dots, 28, \quad k = 1, 2, 3$$

بعنوان مثال در محاسبه امتیاز نهایی استان تهران تحت عامل اول داریم:

$$\begin{aligned} F_{11} &= \sum_{i=1}^{15} w_{1i} z_{i1} = w_{11} z_{11} + w_{12} z_{21} + \dots + w_{1,15} z_{1,15} \\ &= 0,117 \left(\frac{3692 - 2608,5}{2009,28} \right) + \dots + 0,107 \left(\frac{2665 - 231,08}{4937,4} \right) \\ &= 4,917 \end{aligned}$$

z_i ها مقادیر استاندارد شده هر یک از متغیرها بر اساس میانگین و انحراف معیار متناظر با همان متغیر می باشند.

۴.۲ نتیجه گیری

بمنظور سهولت در تحلیل پس از تعیین امتیاز هر استان، در جدول ۶ و ۷ و ۸ بترتیب نزولی مرتب شده است. بنابراین ملاحظه می شود که تحت عامل اول استان تهران به فاصله نسبتاً

14) Varimax with Kaiser Normalization

زیادی نسبت به سایر استان‌های کشور از نظر اهمیت و افزوده‌های فعالیتهای غیر از کشاورزی و ... معدن و ماهیگیری در مقام نخست قرار دارد و تمرکز بخشهایی با ارزش افزوده بالا نظیر: صنعت عمده‌فروشی و خرده‌فروشی و ... و امور عمومی در این استان وجود داشته و منجر به شکاف اقتصادی عمیق نسبت به سایر استانها شده است. از طرفی بمنظور طبقه‌بندی فرضی استان‌های تقریباً همگن تحت این عامل می‌توان استانها را به چهار گروه تقسیم نمود:

- گروه اول: تهران
- گروه دوم: اصفهان، خراسان، خوزستان
- گروه سوم: آذربایجان غربی، هرمزگان، مازندران، گیلان، فارس
- گروه چهارم: سایر استانها

در بررسی استانها تحت عامل دوم دیده می‌شود که در استان‌های کرمان، خوزستان، فارس و اصفهان اهمیت بخش کشاورزی، شکار و جنگلداری نسبت به سایر استانها بیشتر بوده و در مقابل استان‌های هرمزگان، بوشهر و سیستان و بلوچستان کمترین رتبه را داشته‌اند. تحت این عامل بر اساس یک دسته‌بندی فرضی می‌توان استان‌های تقریباً همگن را به چهار گروه بشکل زیر دسته‌بندی نمود:

- گروه اول: کرمان، خراسان، خوزستان، فارس، اصفهان
- گروه دوم: آذربایجان غربی، مازندران، قزوین، همدان، آذربایجان شرقی
- گروه سوم: سایر استانها
- گروه چهارم: سیستان و بلوچستان، بوشهر، هرمزگان

در بررسی عامل سوم که تاکید بر رونق فعالیت‌های معدن و ماهیگیری دارد ملاحظه می‌شود که استان‌های خوزستان و هرمزگان بترتیب دارای بالاترین رتبه بوده و در این خصوص استان‌های کرمان و اردبیل در این بخشها پایین‌ترین میزان اهمیت را داشته‌اند. تحت عامل سوم استان‌های همگن را به سه گروه همگن بشکل زیر دسته‌بندی می‌شوند:

- گروه اول: خوزستان، هرمزگان
- گروه دوم: بوشهر، کهگیلویه و بویراحمد، گیلان، مازندران، سیستان و بلوچستان، گلستان، اصفهان
- گروه سوم: سایر استانها

جدول ۱: ارزش افزوده فعالیتهای اقتصادی کشور به تفکیک استان -- ۱۳۷۹[†] (میلیارد ریال)

| استان | کشاورزی و ... | ماهیگیری | معادن | صنعت | تامین آب و برق و ... | ساختمان | عمده فروشی و ... |
|---------------------|---------------|----------|-------|-------|----------------------|---------|------------------|
| آذربایجان شرقی | ۳۱۷۸ | ۴۸ | ۴۹ | ۱۸۳۱ | ۱۴۱ | ۷۳۷ | ۱۸۶۹ |
| آذربایجان غربی | ۳۹۴۲ | ۲ | ۶۰ | ۵۶۹۱ | ۷۲۵ | ۷۱۵ | ۵۵۵۳ |
| اردبیل | ۱۹۱۶ | ۳ | ۹ | ۶۷۲ | ۶۱ | ۲۵۵ | ۱۲۷۰ |
| اصفهان | ۳۶۰۴ | ۷ | ۱۹۳ | ۱۲۸۳۷ | ۱۹۷۳ | ۲۸۶۸ | ۶۳۵۲ |
| ایلام | ۴۰۱ | ۲ | ۱۴۶ | ۱۶۲۶ | ۲۲ | ۱۳۳ | ۳۶۲ |
| بوشهر | ۵۶۵ | ۲۸۴ | ۳۲۴ | ۸۵۲ | ۳۲۲ | ۲۷۲ | ۹۱۶ |
| تهران | ۳۶۹۲ | ۵ | ۳۰۷ | ۲۹۰۷۶ | ۲۴۸۳ | ۸۱۷۴ | ۳۰۳۰۹ |
| چهارمحال و بختیاری | ۱۰۱۶ | ۸ | ۳۱ | ۴۲۲ | ۴۳ | ۲۷۴ | ۶۲۰ |
| خراسان | ۶۷۷۳ | ۳ | ۲۶۵ | ۵۷۹۳ | ۹۷۲ | ۲۶۹۲ | ۷۷۹۶ |
| خوزستان | ۴۲۰۸ | ۲۱۲ | ۵۴۴۸۶ | ۱۱۱۴۷ | ۱۶۵۴ | ۱۴۱۴ | ۴۶۵۱ |
| زنجان | ۱۴۲۸ | ۳ | ۹۷ | ۱۲۹۶ | ۱۱۶ | ۲۷۰ | ۷۳۲ |
| سمنان | ۱۱۳۵ | ۱ | ۱۰۰ | ۷۵۵ | ۸۱ | ۳۵۴ | ۳۸۹ |
| سیستان و بلوچستان | ۱۰۲۹ | ۱۸۰ | ۲۷ | ۳۹۲ | ۱۱۰ | ۳۲۷ | ۱۰۵۰ |
| قزوین | ۶۲۳۹ | ۱۳ | ۳۹۷ | ۳۲۷۵ | ۷۶۵ | ۱۵۵۴ | ۳۳۱۸ |
| قم | ۲۴۴۵ | ۲ | ۲۳ | ۳۳۹۳ | ۷۷۹ | ۵۰۴ | ۱۳۶۷ |
| کردستان | ۷۰۰ | ۱ | ۲۹ | ۱۴۰۶ | ۷۶ | ۴۶۶ | ۱۳۴۵ |
| کرمان | ۱۱۰۲ | ۶ | ۳۶ | ۴۶۶ | ۵۷ | ۳۰۱ | ۱۴۱۹ |
| کرمانشاه | ۸۱۰۲ | ۲ | ۶۴۰ | ۳۳۰۳ | ۲۸۵ | ۸۱۲ | ۱۹۸۱ |
| کهگیلویه و بویراحمد | ۱۵۳۰ | ۲ | ۴۰ | ۷۸۰ | ۲۲۳ | ۶۱۶ | ۱۲۵۷ |
| گاسفهان | ۵۷۳ | ۴ | ۳۰۰۲۳ | ۱۴۸ | ۳۷ | ۱۹۶ | ۴۱۱ |
| گیلان | ۲۷۴۲ | ۱۷۳ | ۳۰ | ۵۶۱ | ۵۷ | ۴۶۶ | ۲۰۳۴ |
| لرستان | ۲۶۴۰ | ۲۴۹ | ۷۵ | ۲۳۰۲ | ۵۶۸ | ۱۲۲۸ | ۳۳۸۶ |
| مازندران | ۲۱۸۱ | ۹ | ۱۷۵ | ۱۰۴۳ | ۵۷ | ۲۵۹ | ۱۵۳۱ |
| مرکزی | ۵۲۷۰ | ۲۱۷ | ۱۱۸ | ۲۲۱۲ | ۶۰۸ | ۱۱۹۵ | ۴۲۷۱ |
| هرمزگان | ۱۵۰۷ | ۱ | ۷۸ | ۵۸۱۲ | ۳۸۵ | ۷۵۹ | ۲۳۲۴ |
| همدان | ۱۳۶۱ | ۴۹۳ | ۸۸ | ۱۴۰۸ | ۵۷۳ | ۲۳۹ | ۱۳۰۰ |
| یزد | ۲۵۹۰ | ۳ | ۳۲ | ۱۰۷۱ | ۴۷۲ | ۳۷۱ | ۱۵۴۹ |
| | ۱۱۶۷ | ۱ | ۴۲۱ | ۱۳۸۷ | ۱۷۷ | ۵۲۲ | ۸۸۸ |

[†] نتایج حساب تولید استان های کشور -- ۱۳۷۹

| استان | هتل و رستوران | حمل و نقل و ... | واسطه‌گری مالی و ... | کسب و کار | امور عمومی و ... | آموزش | بهداشت و ... | سایر خدمات |
|---------------------|---------------|-----------------|----------------------|-----------|------------------|-------|--------------|------------|
| آذربایجان شرقی | ۲۰۷ | ۹۶۶ | ۱۵۳ | ۱۷۲۹ | ۱۵۸۷ | ۷۲۱ | ۵۰۶ | ۱۶۰ |
| آذربایجان غربی | ۳۶۸ | ۲۰۶۹ | ۲۷۷ | ۳۲۰۲ | ۱۴۱۴ | ۱۰۱۹ | ۷۰۶ | ۲۶۸ |
| اردبیل | ۸۶ | ۴۸۸ | ۶۱ | ۶۰۶ | ۳۹۴ | ۳۷۲ | ۲۳۶ | ۷۶ |
| اصفهان | ۵۰۸ | ۲۵۶۵ | ۶۵۶ | ۳۸۶۲ | ۳۰۷۶ | ۱۴۶۸ | ۱۱۲۳ | ۳۹۸ |
| ایلام | ۲۵ | ۱۸۹ | ۳۶ | ۱۹۶ | ۳۰۹ | ۲۲۴ | ۱۳۰ | ۲۶ |
| پوشهر | ۵۶ | ۷۰۱ | ۹۳ | ۸۳۲ | ۶۴۹ | ۳۰۹ | ۲۱۰ | ۴۳ |
| تهران | ۱۹۵۰ | ۱۲۲۱۷ | ۸۲۵۶ | ۳۴۵۹۶ | ۱۱۶۰۱ | ۵۰۰۹ | ۴۸۶۷ | ۲۶۶۵ |
| چهارمحال و بختیاری | ۲۳ | ۲۸۷ | ۶۴ | ۴۱۵ | ۲۴۳ | ۳۰۹ | ۱۷۲ | ۴۲ |
| خراسان | ۶۴۷ | ۳۲۹۸ | ۵۳۱ | ۳۸۲۱ | ۲۵۵۸ | ۲۱۱۴ | ۱۴۶۵ | ۵۲۷ |
| خوزستان | ۲۸۱ | ۲۸۹۳ | ۳۴۸ | ۲۹۸۳ | ۲۱۸۳ | ۱۱۷۸ | ۹۹۱ | ۱۹۴ |
| زنجان | ۶۱ | ۲۵۱ | ۵۶ | ۵۱۶ | ۳۷۹ | ۳۰۸ | ۱۶۰ | ۶۷ |
| سمنان | ۲۵ | ۲۳۵ | ۷۷ | ۴۲۳ | ۲۲۹ | ۲۵۶ | ۱۶۲ | ۵۰ |
| سیستان و بلوچستان | ۳۶ | ۸۲۳ | ۱۰۶ | ۹۵۲ | ۸۷۶ | ۵۷۰ | ۳۰۸ | ۸۳ |
| فارس | ۳۰۴ | ۲۵۲۵ | ۴۹۲ | ۳۱۸۲ | ۱۷۱۳ | ۱۲۱۳ | ۹۲۷ | ۳۰۶ |
| قزوین | ۷۸ | ۴۲۵ | ۸۷ | ۹۳۲ | ۲۰۵ | ۲۳۷ | ۲۰۸ | ۳۳ |
| قم | ۹۰ | ۳۱۰ | ۶۷ | ۱۰۳۲ | ۳۴۰ | ۲۲۵ | ۳۰۳ | ۲۵۸ |
| کردستان | ۶۱ | ۵۶۱ | ۶۹ | ۶۹۹ | ۷۰۲ | ۳۸۸ | ۲۹۸ | ۷۶ |
| کرمان | ۱۲۱ | ۱۲۲۹ | ۱۸۵ | ۱۶۷۷ | ۱۰۴۲ | ۸۸۵ | ۲۷۶ | ۹۵ |
| کرمانشاه | ۱۲۶ | ۵۳۶ | ۱۰۲ | ۱۲۵۹ | ۱۳۷۰ | ۵۶۸ | ۳۶۱ | ۹۲ |
| کهگیلویه و بویراحمد | ۱۹ | ۱۳۶ | ۳۱ | ۲۳۳ | ۲۳۶ | ۲۷۷ | ۱۵۱ | ۳۳ |
| گلستان | ۵۲ | ۴۶۸ | ۱۰۵ | ۸۲۹ | ۶۴۳ | ۴۷۳ | ۲۹۱ | ۶۹ |
| گیلان | ۲۴۶ | ۱۱۸۰ | ۲۰۰ | ۱۷۶۰ | ۸۶۲ | ۸۱۱ | ۵۹۳ | ۱۴۴ |
| لرستان | ۵۶ | ۷۷۷ | ۷۷ | ۹۱۴ | ۷۳۳ | ۵۱۲ | ۲۸۳ | ۶۵ |
| مازندران | ۳۵۸ | ۲۱۲۷ | ۳۱۴ | ۲۵۰۲ | ۹۳۱ | ۱۰۵۱ | ۸۷۳ | ۲۹۹ |
| مرکزی | ۵۳ | ۱۰۸۱ | ۹۸ | ۱۰۵۲ | ۲۸۷ | ۳۸۷ | ۳۵۵ | ۱۰۵ |
| هرمزگان | ۷۲ | ۲۹۵۵ | ۱۲۳ | ۹۱۵ | ۶۲۵ | ۳۲۴ | ۲۷۰ | ۵۱ |
| همدان | ۱۰۵ | ۶۳۹ | ۹۷ | ۶۳۵ | ۶۳۵ | ۵۳۹ | ۵۳۳ | ۱۱۱ |
| یزد | ۵۲ | ۵۵۳ | ۱۲۲ | ۷۰۱ | ۳۹۶ | ۳۷۲ | ۳۴۶ | ۱۰۶ |

ارزش افزوده فعالیتهای اقتصادی کشور به تفکیک استان -- ۱۳۷۹ (میلیارد ریال) (دنباله جدول ۱)

۳ بخش سوم

تعیین و بررسی تحلیلی جایگاه استانها از نظر درصد نیروی کار در فعالیتهای اقتصادی «بروش تحلیل عاملی»

۱.۳ تحلیل همبستگی

در این بخش نیز نظیر آنچه قبلاً گفته شده با استفاده از اطلاعات برآورد درصد شاغلان در هر یک از فعالیتهای اقتصادی استانها که در جدول ۹ آمده است، ماتریس ضرایب همبستگی بین متغیرها محاسبه گردید. در بررسی اولیه این ماتریس مشاهده شد که بین اکثر متغیرها همبستگی لازم وجود داشته و در تعداد معدودی از متغیرها به دلیل پایین بودن میزان همبستگی و با توجه به نتایج آزمونهای مربوطه نسبت به حذف آنها اقدام گردید^{۱۵} و بررسی نهایی روی یازده متغیر باقیمانده صورت گرفت.

بمنظور ارزیابی کارایی استفاده از روش تحلیل عاملی در این بررسی از ملاکهای مطرح شده

(۱۵) توضیح اینکه چهار متغیر حذف شده عبارتند از:

ساختمان -- عمده‌فروشی، خرده‌فروشی و ...، هتل و رستوران و سایر فعالیتهای خدمات عمومی و اجتماعی و ...

جدول ۲: محاسبه مقادیر MSA منتج شده از ماتریس همبستگی ارزش افزوده به تفکیک فعالیت

| مقادیر MSA | فعالیت |
|------------|----------------------|
| ۰,۵۷۴ | کشاورزی و ... |
| ۶,۵۶۸E-2 | ماهیگیری |
| ۹,۸۸۴E-2 | معدن |
| ۰,۷۴۲ | صنعت |
| ۰,۸۲۲ | تامین آب و برق و ... |
| ۰,۸۰۰ | ساختمان |
| ۰,۸۰۰ | عمده فروشی و ... |
| ۰,۸۷۲ | هتل و رستوران |
| ۰,۸۹۰ | حمل و نقل و ... |
| ۰,۸۲۵ | واسطه‌گری مالی و ... |
| ۰,۸۴۸ | کسب و کار و ... |
| ۰,۸۶۷ | امور عمومی و ... |
| ۰,۷۹۸ | آموزش ... |
| ۰,۹۰۷ | بهداشت و ... |
| ۰,۸۸۰ | سایر خدمات |

جدول ۳: مقادیر ویژه و درصد واریانس و درصد تجمعی واریانس تبیین شده توسط هر یک از عاملها

| عاملهای اولیه | مقادیر ویژه | درصد واریانس | درصد واریانس تجمعی |
|---------------|-------------|--------------|--------------------|
| ۱ | ۱۱,۴۶۰ | ۷۶,۳۹۷ | ۷۶,۳۹۶ |
| ۲ | ۱,۲۸۳ | ۸,۵۵۴ | ۸۴,۹۵۱ |
| ۳ | ۱,۰۱۰ | ۶,۷۳۶ | ۹۱,۶۸۷ |
| ۴ | ۰,۸۳۹ | ۵,۵۹۳ | ۹۷,۲۸۰ |
| ۵ | ۰,۲۷۸ | ۱,۸۵۵ | ۹۹,۱۳۵ |
| ۶ | ۵,۳۶۰E-2 | ۰,۳۵۷ | ۹۹,۴۹۲ |
| ۷ | ۲,۰۲۷E-2 | ۰,۱۳۵ | ۹۹,۶۲۷ |
| ۸ | ۱,۹۱۸E-2 | ۰,۱۲۸ | ۹۹,۷۵۵ |
| . | . | . | . |
| . | . | . | . |
| ۱۵ | ۳,۲۸۶E-۴ | ۲,۱۹۱E-۳ | ۱۰۰,۰۰۰ |

جدول ۴: ماتریس ضرایب عامل دوران یافته بروش واریماکس به تفکیک فعالیت (در تحلیل ارزش افزوده)

| عامل ۳ | عامل ۲ | عامل ۱ | فعالیت |
|-----------|-----------|-----------|----------------------|
| -۵,۰۶۰E-۳ | ۰,۸۶۸ | ۰,۲۱۸ | کشاورزی و ... |
| ۰,۷۴۵ | -۰,۳۵۱ | -۲,۱۵۳E-۲ | ماهگیری |
| ۰,۷۵۴ | ۰,۳۰۴ | -۲,۶۶۶E-۲ | معادن |
| ۰,۱۱۵ | ۰,۲۴۵ | ۰,۹۲۳ | صنعت |
| ۰,۳۰۵ | ۰,۳۸۰ | ۰,۷۸۰ | تامین آب و برق و ... |
| -۳,۲۸۳E-۲ | ۰,۲۰۷ | ۰,۹۶۸ | ساختمان |
| -۲,۷۹۱E-۲ | ۰,۱۳۰ | ۰,۹۸۷ | عمده فروشی و ... |
| -۴,۲۲۳E-۲ | ۰,۱۹۲ | ۰,۹۷۲ | هتل و رستوران |
| ۷,۹۰۹E-۲ | ۹,۵۴۵E-۲ | ۰,۹۸۱ | حمل و نقل و ... |
| -۷,۳۹۹E-۲ | -۳,۱۸۹E-۲ | ۰,۹۸۳ | واسطه‌گری مالی و ... |
| -۶,۳۹۳E-۲ | ۴,۹۴۱E-۲ | ۰,۹۹۱ | کسب و کار |
| -۱,۸۸۲E-۲ | ۰,۱۰۷ | ۰,۹۸۴ | امور عمومی و ... |
| -۲,۵۵۷E-۲ | ۰,۲۸۷ | ۰,۹۴۸ | آموزش |
| -۹,۰۰۵E-۲ | ۰,۱۸۰ | ۰,۹۸۱ | بهداشت و ... |
| -۹,۴۹۴E-۲ | ۵,۸۵۲E-۲ | ۰,۹۸۸ | سایر خدمات |

جدول ۵: ماتریس ضرایب نهایی عوامل به تفکیک فعالیت بروش رگرسیونی (در تحلیل ارزش افزوده)

| فعالیت | عامل ۱ | عامل ۲ | عامل ۳ |
|----------------------|--------|--------|--------|
| کشاورزی و ... | -۰٫۱۱۷ | ۰٫۷۸۰ | -۰٫۰۵۳ |
| ماهیگیری | ۰٫۰۶۳ | -۰٫۳۷۰ | ۰٫۶۱۵ |
| معادن | -۰٫۰۴۶ | ۰٫۲۴۵ | ۰٫۵۸۴ |
| صنعت | ۰٫۰۷۱ | ۰٫۰۷۱ | ۰٫۰۸۷ |
| تامین آب و برق و ... | ۰٫۰۳۳ | ۰٫۲۱۰ | ۰٫۲۲۹ |
| ساختمان | ۰٫۰۸۱ | ۰٫۰۳۶ | -۰٫۰۲۸ |
| عمده‌فروشی و ... | ۰٫۰۹۶ | -۰٫۰۳۹ | -۰٫۰۲۰ |
| هتل و رستوران | ۰٫۰۸۴ | ۰٫۰۲۲ | -۰٫۰۳۵ |
| حمل و نقل و ... | ۰٫۱۰۲ | -۰٫۰۷۷ | ۰٫۰۶۷ |
| واسطه‌گری مالی و ... | ۰٫۱۲۲ | -۰٫۱۸۸ | -۰٫۰۴۷ |
| کسب و کار | ۰٫۱۰۹ | -۰٫۱۱۴ | -۰٫۰۴۴ |
| امور عمومی و ... | ۰٫۱۰۰ | -۰٫۰۶۱ | -۰٫۰۱۱ |
| آموزش | ۰٫۰۶۵ | ۰٫۱۱۵ | -۰٫۰۲۸ |
| بهداشت و ... | ۰٫۰۸۷ | ۰٫۰۰۸ | -۰٫۰۰۸ |
| سایر خدمات | ۰٫۱۰۷ | -۰٫۱۰۳ | -۰٫۰۶۹ |

استفاده و نتایج آن بشرح زیر می‌باشد.

دترمینان ماتریس همبستگی با مقدار $۱۰^{-۵} \times ۴/۸۷۹$ بوده و بعلت کوچک بودن مقدار بدست آمده، استفاده از روش تحلیل عاملی تایید شد. بررسی مقدار MSAها بااستثنای بخش صنعت که خیلی قوی نیست بر حضور سایر متغیرها تاکید دارد، اما بدلیل اهمیت بخش صنعت، حضور آن در مدل الزامی به نظر رسید. مقدار شاخص KMO بدست آمده نیز تقریباً معادل $۰/۷$ می‌باشد که مقدار این شاخص نیز استفاده از روش تحلیل عاملی را تایید می‌کند.

۲.۳ تعیین تعداد و روش محاسبه عاملها و میزان برازش مدل

برای تعیین تعداد عاملها در این بررسی با توجه به مقادیر ویژه بدست آمده در جدول ۱۲ و میزان واریانس تبیین شده توسط هر یک از عاملهای اولیه مشاهده شد که عاملهای اول تا سوم بدلیل آنکه دارای بیشترین مقادیر ویژه و با تبیین حدود ۸۰ درصد واریانس کل، بعنوان عاملهای منتخب پذیرفته شدند. در ادامه برای سهولت در تفسیر عاملها، با استفاده از دوران متعامد ماتریس ضرایب عاملها برای هر یک از سه عامل محاسبه و نتایج آن در جدول ۱۳ آمده است. در تعبیر و تفسیر هر یک از عاملها با استفاده از مقادیر جدول ۱۳ بشرح زیر داریم:

جدول ۶: رتبه‌بندی ارزش افزوده استانها بر حسب عامل اول -- ۱۳۷۹ (تاکید بر اهمیت بخشهای غیر از کشاورزی، معدن و ماهیگیری)

| رتبه | استان | امتیاز کسب شده |
|------|----------------------|----------------|
| ۱ | تهران | ۴,۹۲ |
| ۲ | اصفهان | ۰,۵۹ |
| ۳ | خراسان | ۰,۴۶ |
| ۴ | خوزستان | ۰,۱۳ |
| ۵ | آذربایجان غربی | ۰,۰۵ |
| ۶ | هرمزگان | ۰,۰۳ |
| ۷ | مازندران | ۰,۰۳ |
| ۸ | گیلان | ۰,۰۲ |
| ۹ | فارس | ۰,۰۱ |
| ۱۰ | بوشهر | -۰,۱۵۲ |
| ۱۱ | آذربایجان شرقی | -۰,۱۹۲ |
| ۱۲ | سیستان و بلوچستان | -۰,۱۹۴ |
| ۱۳ | مرکزی | -۰,۱۹۶ |
| ۱۴ | کرمانشاه | -۰,۲۴۱ |
| ۱۵ | قم | -۰,۲۶۳ |
| ۱۶ | یزد | -۰,۳۰۷ |
| ۱۷ | گلستان | -۰,۳۰۸ |
| ۱۸ | کردستان | -۰,۳۱۹ |
| ۱۹ | همدان | -۰,۳۳۲ |
| ۲۰ | قزوین | -۰,۳۴۶ |
| ۲۱ | لرستان | -۰,۳۵۷ |
| ۲۲ | اردبیل | -۰,۳۸۹ |
| ۲۳ | سمنان | -۰,۳۹۱ |
| ۲۴ | زنجان | -۰,۳۹۲ |
| ۲۵ | چهارمحال و بختیاری | -۰,۴۰۰ |
| ۲۶ | ایلام | -۰,۴۱۰ |
| ۲۷ | کرمان | -۰,۵۱۴ |
| ۲۸ | کهگیلویه و بویر احمد | -۰,۵۲۶ |

جدول ۷: رتبه‌بندی ارزش افزوده استانها بر حسب عامل دوم -- ۱۳۷۹ (تاکید بر اهمیت بخش کشاورزی)

| رتبه | استان | امتیاز کسب شده |
|------|---------------------|----------------|
| ۱ | کرمان | ۲,۳۱ |
| ۲ | خراسان | ۱,۹۴ |
| ۳ | خوزستان | ۱,۷۷ |
| ۴ | فارس | ۱,۶۱ |
| ۵ | اصفهان | ۱,۰۹ |
| ۶ | آذربایجان غربی | ۰,۷۵ |
| ۷ | مازندران | ۰,۵۸ |
| ۸ | قزوین | ۰,۲۷ |
| ۹ | همدان | ۰,۱۹ |
| ۱۰ | آذربایجان شرقی | ۰,۱۵ |
| ۱۱ | کهگیلویه و بویراحمد | -۰,۰۹ |
| ۱۲ | لرستان | -۰,۱۳ |
| ۱۳ | اردبیل | -۰,۲۱ |
| ۱۴ | مرکزی | -۰,۲۲ |
| ۱۵ | کرمانشاه | -۰,۳۱ |
| ۱۶ | زنجان | -۰,۳۷ |
| ۱۷ | گلستان | -۰,۴۰ |
| ۱۸ | یزد | -۰,۴۵ |
| ۱۹ | گیلان | -۰,۴۸ |
| ۲۰ | سمنان | -۰,۵۰ |
| ۲۱ | تهران | -۰,۵۳ |
| ۲۲ | کردستان | -۰,۵۵ |
| ۲۳ | چهارمحال و بختیاری | -۰,۵۷ |
| ۲۴ | قم | -۰,۷۱ |
| ۲۵ | ایلام | -۰,۷۷ |
| ۲۶ | سیستان و بلوچستان | -۱,۰۸ |
| ۲۷ | بوشهر | -۱,۵۰ |
| ۲۸ | هرمزگان | -۱,۸۰ |

جدول ۸: رتبه‌بندی ارزش افزوده استانها بر حسب عامل سوم - ۱۳۷۹ (تاکید بر اهمیت بخشهای معدن و ماهیگیری)

| رتبه | استان | امتیاز کسب شده |
|------|----------------------|----------------|
| ۱ | خوزستان | ۳,۸۱ |
| ۲ | هرمزگان | ۲,۱۳ |
| ۳ | بوشهر | ۰,۹۶ |
| ۴ | کهگیلویه و بویر احمد | ۰,۹۵ |
| ۵ | گیلان | ۰,۷۶ |
| ۶ | مازندران | ۰,۵۲ |
| ۷ | سیستان و بلوچستان | ۰,۳۲ |
| ۸ | گلستان | ۰,۲۲ |
| ۹ | اصفهان | ۰,۰۷ |
| ۱۰ | قزوین | -۰,۳۲ |
| ۱۱ | مرکزی | -۰,۴۰ |
| ۱۲ | آذربایجان شرقی | -۰,۴۱ |
| ۱۳ | آذربایجان غربی | -۰,۴۳ |
| ۱۴ | فارس | -۰,۴۶ |
| ۱۵ | تهران | -۰,۴۷ |
| ۱۶ | همدان | -۰,۴۸ |
| ۱۷ | ایلام | -۰,۴۹ |
| ۱۸ | یزد | -۰,۵۲ |
| ۱۹ | خراسان | -۰,۵۴ |
| ۲۰ | چهارمحال و بختیاری | -۰,۵۵۱ |
| ۲۱ | زنجان | -۰,۵۵۴ |
| ۲۲ | کرمانشاه | -۰,۵۶۲ |
| ۲۳ | لرستان | -۰,۵۶۴ |
| ۲۴ | سمنان | -۰,۵۶۶ |
| ۲۵ | کردستان | -۰,۵۷۱ |
| ۲۶ | قم | -۰,۵۹۸ |
| ۲۷ | اردبیل | -۰,۵۹۹ |
| ۲۸ | کرمان | -۰,۶۴۹ |

ویژگی عامل اول: این عامل عمدتاً بر فعالیتهای ماهیگیری، حمل و نقل، تامین آب و برق و امور عمومی تاکید دارد و دارای ضرایب منفی نسبتاً بزرگ برای فعالیت صنعت می باشد. به بیان دیگر، استان های که تحت این عامل از رتبه بالایی برخوردارند سهم اشتغال بخش صنعت در آنها ضعیف است.

ویژگی عامل دوم: این عامل بر اهمیت فعالیتهای بهداشت و مددکاری اجتماعی و . . . ، آموزش و معدن تاکید می کند.

ویژگی عامل سوم: فعالیتهایی که توسط این عامل بر اهمیت آنها تاکید شده عبارتند از کسب و کار، واسطه گریهای مالی و صنعت می باشد و بعکس با ضریب منفی نسبت به فعالیت کشاورزی کمترین میزان اهمیت را در مورد فعالیت کشاورزی نشان می دهد.

۳.۳ جداول نهایی

پس از تعیین تعداد و ویژگیهای هر عامل امتیاز نهایی هر یک از استانها را بر اساس فرمول زیر محاسبه می نمایم که پس از مرتب نمودن امتیازات مذکور به صورت نزولی در جدول ۱۵ و ۱۶ و ۱۷ نشان داده شده است.

$$B_{jk} = \sum_{i=1}^{11} w_{ji} z_{ik} \quad j = 1, 2, \dots, 28 \quad k = 1, 2, 3$$

۴.۳ نتیجه گیری

در بررسی جدول ۱۵ مشاهده می شود که تحت عامل اول استان های هرمزگان، بوشهر، و خوزستان در ارتباط با اهمیت و تمرکز نیروی کار در رشته فعالیتهای ماهیگیری، حمل و نقل، تامین آب و برق و . . . ، امور عمومی بترتیب بالاترین رتبه را کسب نموده و در مقابل استان های یزد، زنجان در این خصوص در پایین ترین رتبه قرار داشته اند.

برای دسته بندی استان های همگن تحت عامل اول می توان آنها را در سه گروه فرضی بشکل زیر تقسیم نمود:

- گروه اول: هرمزگان، بوشهر، خوزستان
- گروه دوم: سیستان و بلوچستان، کرمانشاه، فارس، آذربایجان
- گروه سوم: سایر استانها

پس از مرتب نمودن استانها در ارتباط با عامل دوم ملاحظه می شود که استان های کهگیلویه و بویر احمد، سمنان، کرمان، ایلام و خوزستان در بخشهای بهداشت و مددکاری و . . . ، آموزش و معدن دارای بالاترین رتبه بوده اند و آذربایجان شرقی و غربی و قزوین در این خصوص کمترین رتبه را داشته اند.

تحت این عامل استان‌های تقریباً همگن عبارتند از:

- گروه اول: کهگیلویه و بویر احمد، سمنان، کرمان، ایلام
- گروه دوم: خوزستان، بوشهر، تهران، یزد، فارس، کرمانشاه، لرستان، مازندران
- گروه سوم: سایر استانها
- گروه چهارم: مرکزی، آذربایجان غربی و شرقی، قزوین

در تحلیل استانها بر اساس عامل سوم که بر اهمیت درصد نیروی کار در بخشهای صنعت، کسب و کار و واسطه‌گریهای مالی تاکید دارد نشان می‌دهد که استان‌های تهران، یزد، اصفهان و قم بترتیب بالاترین رتبه را در این خصوص دارا بوده و در مقابل استان‌های ایلام، کهگیلویه و بویر احمد و اردبیل پایین‌ترین رتبه را داشته‌اند.

طبقه‌بندی استان‌های تقریباً همگن تحت عامل سوم عبارت است از:

- گروه اول: تهران، یزد، اصفهان، قم
- گروه دوم: بوشهر، سمنان، مرکزی، آذربایجان شرقی، خوزستان، هرمزگان، چهارمحال و بختیاری
- گروه سوم: قزوین، فارس، کرمان، مازندران، لرستان، زنجان، گیلان، همدان
- گروه چهارم: سایر استانها

۵.۳ نتیجه‌گیری کلی

-- بر اساس رتبه‌های کسب شده استانها در خصوص دو شاخص ارزش افزوده و درصد نیروی کار در فعالیتهای اقتصادی، بارزترین نکته وضعیت استانها در دو بخش صنعت و خدمات بوده بطوریکه ملاحظه می‌شود استان تهران در ارزیابی این دو بخش فاصله زیادی نسبت به سایر استانها داشته است. به این ترتیب با تجمع و تمرکز سهم زیادی از رشته فعالیتهای عمده صنعتی و خدماتی عملاً منجر به جذب سرمایه و نیروی کار و در نهایت ایجاد شکاف اقتصادی با دیگر استانها شده است. بعنوان بارزترین مثال در بیان این موضوع نگاهی به سرانه تولید ناخالص داخلی نشان می‌دهد که در سال ۱۳۷۹ مقدار این شاخص در کشور معادل ۱۰/۱۴ میلیون ریال بوده و در استان تهران برابر ۱۳/۸۶ میلیون ریال می‌باشد و برای مقایسه بهتر چنانچه ارزش افزوده فعالیت استخراج نفت و گاز طبیعی که جنبه ملی دارد را از سر جمع تولید ناخالص داخلی کشور و استان تهران کسر نماییم، آنگاه سرانه مورد نظر در کشور با ۱/۷۴ واحد کاهش به ۸/۴۰ میلیون ریال نزول می‌کند در حالیکه این سرانه در استان تهران فقط با کسر ۰/۲ واحد به ۱۳/۸۴ میلیون ریال می‌رسد. بهر حال با توجه به اینکه به تجربه دیده شده بخش صنعت بعنوان یکی از ضرورت‌های توسعه اقتصادی پایدار که ناشی از وجود رابطه خطی مثبت سهم این بخش، از GDP (GDPR) با تولید ناخالص داخلی سرانه می‌باشد، همواره مطرح بوده است. بنابراین هر چه سهم این بخش افزایش یابد، دسترسی به توسعه پایدار در منطقه تسهیل شده و لذا به منظور کاهش بخشی از عدم تعادل‌ها که این امر لزوم توجه جدی به تقویت بخش صنعت بخصوص در

جدول ۹: برآورد درصد شاغلان[†] هر یک از فعالیتهای اقتصادی[‡] به تفکیک استان -- ۱۳۷۹

| استان | کشاورزی | ماهیگیری | معدن | صنعت | تامین آب و برق و ... | ساختمان | عمده فروشی و ... |
|---------------------|---------|----------|------|-------|----------------------|---------|------------------|
| آذربایجان شرقی | ۲۵٫۳۸ | -/۰۰ | ۰/۱۴ | ۲۸٫۵۸ | ۰/۶۵ | ۹٫۹۷ | ۱۲٫۱۲ |
| آذربایجان غربی | ۳۷٫۷۲ | -/۰۳ | ۰/۱۷ | ۱۲٫۶۷ | ۰/۶۸ | ۱۰/۷۰ | ۱۲٫۱۸ |
| اردبیل | ۳۷٫۶۸ | -/۰۱ | ۰/۱۸ | ۱۰/۸۰ | ۰/۶۷ | ۱۳٫۹۱ | ۱۲٫۰۴ |
| اصفهان | ۱۳٫۵۲ | -/۰۱ | ۰/۶۵ | ۳۲٫۲۴ | ۰/۱۳ | ۱۰/۳۶ | ۱۳٫۲۶ |
| ایلام | ۳۶٫۳۵ | -/۰۰ | ۰/۱۵ | ۵٫۷۱ | ۰/۸۷ | ۱۱٫۶۱ | ۹٫۳۶ |
| بوشهر | ۱۳٫۳۵ | ۵٫۷۴ | ۲٫۰۳ | ۷٫۸۰ | ۱٫۶۲ | ۱۰/۲۸ | ۱۱٫۲۷ |
| تهران | ۳٫۶۸ | -/۰۱ | ۰/۸۳ | ۲۳٫۹۴ | ۰/۹۸ | ۹٫۴۰ | ۱۹٫۳۷ |
| چهارمحال و بختیاری | ۲۳٫۴۷ | -/۰۷ | ۰/۲۸ | ۲۱٫۶۱ | ۰/۶۵ | ۱۸٫۷۵ | ۸٫۵۶ |
| خراسان | ۲۹٫۸۰ | -/۰۰ | ۰/۲۸ | ۱۹٫۷۸ | ۰/۷۴ | ۱۰/۵۳ | ۱۲٫۰۴ |
| خوزستان | ۲۰/۱۱ | ۱٫۳۸ | ۳٫۴۱ | ۱۴٫۲۲ | ۲٫۸۶ | ۱۰/۸۸ | ۱۲٫۷۸ |
| زنجان | ۳۶٫۸۲ | -/۰۰ | ۰/۷۶ | ۲۲٫۱۰ | ۰/۵۹ | ۹٫۰۲ | ۸٫۴۳ |
| سمنان | ۲۱٫۷۷ | -/۰۰ | ۲٫۸۰ | ۱۵٫۷۵ | ۱٫۱۴ | ۹٫۶۲ | ۹٫۷۹ |
| سیستان و بلوچستان | ۲۷٫۳۵ | ۱٫۹۸ | ۰/۳۴ | ۸٫۰۰ | ۰/۹۷ | ۱۵٫۲۹ | ۱۱٫۶۸ |
| فارس | ۲۴٫۹۷ | -/۰۱ | ۰/۵۷ | ۱۲٫۷۵ | ۱٫۰۶ | ۱۳٫۰۰ | ۱۳٫۹۶ |
| قزوین | ۳۰/۵۹ | -/۰۳ | ۰/۱۷ | ۲۵٫۲۵ | ۱٫۶۸ | ۶٫۲۶ | ۱۰/۷۷ |
| قم | ۸٫۲۲ | -/۰۰ | ۰/۵۴ | ۲۸٫۰۰ | ۰/۷۰ | ۱۲٫۸۶ | ۱۶٫۲۹ |
| کردستان | ۳۲٫۷۲ | -/۱۰ | ۰/۲۵ | ۱۳٫۲۵ | ۰/۶۱ | ۱۳٫۳۸ | ۱۳٫۴۰ |
| کرمان | ۳۰/۱۲ | -/۰۰ | ۲٫۲۵ | ۱۰/۵۹ | ۱٫۰۰ | ۱۱٫۲۵ | ۱۰/۷۹ |
| کرمانشاه | ۲۹٫۳۴ | -/۰۲ | ۰/۳۶ | ۷٫۷۲ | ۱٫۰۲ | ۱۲٫۴۶ | ۱۳٫۳۲ |
| کهگیلویه و بویراحمد | ۳۰/۰۴ | -/۰۱ | ۲٫۵۶ | ۴٫۷۳ | ۱٫۰۸ | ۱۸٫۶۰ | ۸٫۵۰ |
| گاستان | ۳۹٫۹۵ | -/۰۷ | ۰/۵۴ | ۱۶٫۹۱ | ۰/۷۳ | ۷٫۹۹ | ۱۰/۲۳ |
| گیلان | ۳۷٫۳۶ | -/۰۱ | ۰/۲۱ | ۱۳٫۶۶ | ۰/۸۲ | ۵٫۱۲ | ۱۳٫۷۳ |
| لرستان | ۲۹٫۲۲ | -/۰۱ | ۰/۲۸ | ۸٫۴۶ | ۰/۶۶ | ۱۸٫۶۳ | ۱۱٫۶۹ |
| مازندران | ۳۶٫۷۷ | -/۲۰ | ۰/۴۸ | ۱۱٫۲۸ | ۰/۸۲ | ۹٫۲۹ | ۱۲٫۷۱ |
| مرکزی | ۲۵٫۹۳ | -/۰۰ | ۰/۷۲ | ۲۷٫۰۵ | ۰/۹۰ | ۱۱٫۰۲ | ۱۰/۳۳ |
| هرمزگان | ۱۷٫۹۱ | ۷٫۰۰ | ۰/۴۹ | ۶٫۰۷ | ۲٫۶۸ | ۱۲٫۴۹ | ۱۲٫۱۶ |
| همدان | ۳۲٫۱۳ | -/۰۰ | ۰/۱۹ | ۱۳٫۹۲ | ۱٫۰۷ | ۱۳٫۷۴ | ۱۳٫۳۷ |
| یزد | ۱۳٫۲۹ | -/۰۲ | ۱٫۴۱ | ۳۰/۳۸ | ۰/۸۲ | ۱۱٫۱۹ | ۱۱٫۵۹ |

[†] برآورد بر اساس نتایج سرشماری نفوس و مسکن ۱۳۷۵ و نتایج طرح آمارگیری از ویژگیهای اشتغال و بیکاری کشور -- ۱۳۷۹
[‡] توضیح اینکه برخی ارقام مندرج در ستون فعالیت ماهیگیری بعلافت گرد شدن بصورت صفر نمایش داده شده است.

جدول ۱۰: برآورد درصد شاغلان هر یک از فعالیتهای اقتصادی به تفکیک استان - ۱۳۷۹ (دنباله)

| استان | هتل و رستوران | حمل و نقل و ... | واسطه‌گری مالی | کسب و کار | امور عمومی و ... | آموزش | بهداشت و ... | سایر خدمات |
|---------------------|---------------|-----------------|----------------|-----------|------------------|-------|--------------|------------|
| آذربایجان شرقی | ۰٫۷۵ | ۶٫۲۴ | ۰٫۶۴ | ۰٫۵۹ | ۶٫۳۱ | ۵٫۱۳ | ۱٫۴۱ | ۱٫۴۳ |
| آذربایجان غربی | ۰٫۵۵ | ۵٫۸۴ | ۰٫۵۹ | ۰٫۶۰ | ۹٫۰۲ | ۵٫۳۶ | ۱٫۵۱ | ۱٫۷۰ |
| اردبیل | ۱٫۰۷ | ۵٫۷۳ | ۰٫۶۶ | ۰٫۴۸ | ۵٫۹۳ | ۶٫۳۲ | ۱٫۵۰ | ۲٫۰۸ |
| اصفهان | ۰٫۶۱ | ۶٫۷۳ | ۰٫۹۴ | ۰٫۹۴ | ۸٫۴۸ | ۷٫۰۲ | ۱٫۹۷ | ۱٫۵۰ |
| ایلام | ۰٫۴۷ | ۵٫۶۳ | ۰٫۹۴ | ۰٫۳۹ | ۱۳٫۸۳ | ۹٫۶۱ | ۲٫۲۸ | ۱٫۶۸ |
| بوشهر | ۰٫۳۳ | ۱۳٫۵۰ | ۱٫۲۴ | ۱٫۷۴ | ۱۷٫۶۳ | ۸٫۴۸ | ۲٫۴۱ | ۱٫۷۱ |
| تهران | ۰٫۸۰ | ۸٫۸۷ | ۲٫۲۰ | ۲٫۵۴ | ۱۲٫۹۲ | ۶٫۸۸ | ۲٫۷۱ | ۳٫۵۱ |
| چهارمحال و بختیاری | ۰٫۳۵ | ۶٫۰۴ | ۰٫۶۸ | ۰٫۵۹ | ۶٫۸۳ | ۸٫۶۶ | ۱٫۹۲ | ۰٫۷۷ |
| خراسان | ۰٫۵۸ | ۵٫۷۳ | ۰٫۸۹ | ۰٫۸۳ | ۶٫۸۶ | ۶٫۸۶ | ۱٫۶۵ | ۲٫۵۶ |
| خوزستان | ۰٫۶۱ | ۷٫۶۷ | ۱٫۰۱ | ۱٫۰۸ | ۱۱٫۶۰ | ۷٫۱۷ | ۲٫۵۴ | ۱٫۸۶ |
| زنجان | ۰٫۵۱ | ۲٫۶۲ | ۰٫۶۰ | ۰٫۲۵ | ۶٫۱۴ | ۶٫۵۴ | ۱٫۳۱ | ۱٫۶۴ |
| سمنان | ۰٫۳۶ | ۷٫۱۷ | ۱٫۷۸ | ۰٫۴۶ | ۱۲٫۲۵ | ۹٫۳۳ | ۲٫۶۵ | ۲٫۲۱ |
| سیستان و بلوچستان | ۰٫۴۲ | ۸٫۲۱ | ۰٫۷۸ | ۰٫۴۲ | ۱۲٫۲۳ | ۷٫۳۶ | ۱٫۹۳ | ۲٫۰۷ |
| فارس | ۰٫۵۷ | ۹٫۱۸ | ۰٫۸۳ | ۰٫۶۶ | ۸٫۹۷ | ۸٫۵۶ | ۲٫۱۷ | ۱٫۹۴ |
| قزوین | ۰٫۵۹ | ۷٫۶۹ | ۰٫۷۰ | ۰٫۵۶ | ۵٫۰۹ | ۷٫۶۱ | ۱٫۳۴ | ۱٫۲۴ |
| قم | ۰٫۹۷ | ۷٫۲۴ | ۰٫۷۸ | ۰٫۶۰ | ۷٫۹۶ | ۸٫۰۰ | ۱٫۸۰ | ۳٫۷۷ |
| کردستان | ۰٫۶۵ | ۵٫۷۶ | ۰٫۵۷ | ۰٫۵۴ | ۸٫۴۴ | ۶٫۴۶ | ۱٫۷۰ | ۱٫۶۶ |
| کرمان | ۰٫۴۱ | ۶٫۵۶ | ۱٫۱۲ | ۰٫۶۷ | ۱۰٫۲۲ | ۹٫۸۱ | ۲٫۴۸ | ۱٫۷۹ |
| کرمانشاه | ۰٫۴۳ | ۷٫۹۹ | ۰٫۷۳ | ۰٫۷۴ | ۱۳٫۴۱ | ۷٫۲۱ | ۲٫۲۶ | ۱٫۶۸ |
| کهگیلویه و بویراحمد | ۰٫۳۵ | ۴٫۸۲ | ۰٫۹۳ | ۰٫۵۳ | ۱۰٫۷۹ | ۱۲٫۱۳ | ۳٫۰۹ | ۱٫۰۹ |
| گلستان | ۰٫۴۴ | ۵٫۰۴ | ۰٫۸۳ | ۰٫۳۹ | ۴٫۷۱ | ۵٫۵۱ | ۲٫۱۷ | ۲٫۳۲ |
| گیلان | ۱٫۱۳ | ۶٫۵۹ | ۰٫۸۸ | ۰٫۷۹ | ۶٫۵۸ | ۶٫۹۴ | ۲٫۲۶ | ۲٫۳۳ |
| لرستان | ۰٫۴۴ | ۶٫۵۲ | ۰٫۸۶ | ۰٫۸۱ | ۱۰٫۴۳ | ۷٫۷۳ | ۱٫۹۴ | ۱٫۳۷ |
| مازندران | ۰٫۶۳ | ۶٫۲۰ | ۰٫۸۸ | ۰٫۹۸ | ۶٫۵۸ | ۷٫۳۳ | ۲٫۳۲ | ۲٫۲۸ |
| مرکزی | ۰٫۴۱ | ۶٫۲۹ | ۰٫۹۷ | ۰٫۵۱ | ۵٫۷۳ | ۶٫۰۵ | ۱٫۵۵ | ۱٫۷۷ |
| هرمزگان | ۰٫۶۹ | ۱۲٫۱۳ | ۱٫۰۵ | ۱٫۲۹ | ۱۲٫۸۸ | ۶٫۲۲ | ۲٫۷۰ | ۳٫۰۳ |
| همدان | ۰٫۵۱ | ۶٫۶۲ | ۰٫۸۳ | ۰٫۲۰ | ۷٫۵۰ | ۶٫۲۴ | ۱٫۷۲ | ۱٫۳۸ |
| یزد | ۰٫۵۰ | ۷٫۵۹ | ۱٫۴۶ | ۰٫۸۱ | ۷٫۷۷ | ۸٫۲۳ | ۲٫۳۹ | ۱٫۹۶ |

جدول ۱۱: محاسبه مقادیر MSA منتج شده از ماتریس همبستگی درصد شاغلین به تفکیک فعالیت

| فعالیت | مقادیر MSA |
|----------------------|------------|
| کشاورزی و ... | ۰٫۵۵۷ |
| ماهنگیری | ۰٫۸۰۹ |
| معدن | ۰٫۷۱۴ |
| صنعت | ۰٫۳۴۵ |
| تامین آب و برق و ... | ۰٫۶۹۹ |
| حمل و نقل و ... | ۰٫۷۵۸ |
| واسطه‌گری مالی | ۰٫۶۶۴ |
| کسب و کار | ۰٫۷۵۴ |
| امور عمومی و ... | ۰٫۶۹۰ |
| بهداشت و ... | ۰٫۷۲۵ |
| آموزش | ۰٫۶۱۰ |

جدول ۱۲: مقادیر ویژه و درصد واریانس و درصد تجمعی واریانس تبیین شده توسط هر یک از عاملها

| عاملهای اولیه | مقادیر ویژه | درصد واریانس | درصد واریانس تجمعی |
|---------------|-------------|--------------|--------------------|
| ۱ | ۴,۹۰۵ | ۴۴,۵۹۲ | ۴۴,۵۹۲ |
| ۲ | ۲,۰۴۰ | ۱۸,۵۴۹ | ۶۳,۱۴۱ |
| ۳ | ۱,۷۳۶ | ۱۵,۷۸۴ | ۷۸,۹۲۵ |
| ۴ | ۰,۸۹۴ | ۷,۷۱۴ | ۸۶,۶۳۹ |
| ۵ | ۰,۴۵۳ | ۴,۱۱۹ | ۹۰,۷۵۸ |
| . | . | . | . |
| . | . | . | . |
| ۱۱ | ۴,۰۱۶E-۲ | ۰,۳۶۵ | ۱۰۰,۰۰۰ |

جدول ۱۳: ماتریس ضرایب عامل دوران یافته بر روش واریماکس به تفکیک فعالیت برای درصد شاغلان

| فعالیت | عامل ۱ | عامل ۲ | عامل ۳ |
|----------------------|--------|-----------|----------|
| کشاورزی و ... | -۰,۲۳۱ | -۰,۱۵۸ | -۰,۸۹۰ |
| ماهیکیری | ۰,۹۲۹ | -۱,۲۲۱E-۲ | ۶,۶۵۴E-۲ |
| معادن | ۰,۱۲۰ | ۰,۷۷۰ | ۰,۲۱۰ |
| صنعت | -۰,۵۱۹ | -۰,۴۲۱ | ۰,۶۴۷ |
| تامین آب و برق و ... | ۰,۷۳۰ | ۰,۱۸۸ | ۰,۱۴۹ |
| حمل و نقل و ... | ۰,۸۴۷ | ۵,۰۷۴E-۲ | ۰,۴۰۱ |
| واسطه‌گری مالی | ۰,۱۶۵ | ۰,۵۱۵ | ۰,۷۰۸ |
| کسب و کار | ۰,۴۳۳ | ۹,۲۲۰E-۲ | ۰,۷۵۴ |
| امور عمومی و ... | ۰,۶۱۳ | ۰,۵۸۶ | ۰,۱۶۳ |
| بهداشت و ... | ۰,۳۵۷ | ۰,۸۲۴ | ۰,۲۲۷ |
| آموزش | -۰,۱۳۱ | ۰,۸۷۶ | -۰,۱۰۷ |

جدول ۱۴: ماتریس ضرایب نهایی عوامل بروش رگرسیونی به تفکیک فعالیت برای درصد شاغلان

| فعالیت | عامل ۱ | عامل ۲ | عامل ۳ |
|----------------------|--------|--------|--------|
| کشاورزی و ... | ۰/۰۴۰ | ۰/۰۱۷ | -۰/۳۶۳ |
| ماهیکیری | ۰/۳۷۰ | -۰/۱۳۴ | ۰/۰۸۷ |
| معدن | -۰/۰۸۱ | ۰/۲۹۱ | ۰/۰۳۵ |
| صنعت | -۰/۲۳۶ | -۰/۱۴۲ | ۰/۳۸۱ |
| تامین آب و برق و ... | ۰/۲۵۲ | -۰/۰۲۹ | -۰/۰۳۶ |
| حمل و نقل و ... | ۰/۲۸۶ | -۰/۱۱۶ | ۰/۰۷۱ |
| واسطه‌گریهای مالی | -۰/۰۹۲ | ۰/۱۵۱ | ۰/۲۶۸ |
| کسب و کار | ۰/۰۶۸ | -۰/۰۶۴ | ۰/۲۸۰ |
| امور عمومی و ... | ۰/۱۴۹ | ۰/۱۵۱ | -۰/۰۳۸ |
| بهداشت و ... | ۰/۰۰۵ | ۰/۲۸۰ | ۰/۰۱۰ |
| آموزش | -۰/۱۵۷ | ۰/۳۸۷ | -۰/۰۸۲ |

استانهایی که پتانسیل بیشتری در این خصوص دارند را طلب می‌نماید.

-- مقایسه برخی رتبه‌های کسب شده در استانها بیانگر عدم تطابق الگوی ارزش افزوده فعالیتها در ارتباط با درصد نیروی کار در فعالیتهای متناظر می‌باشد. به بیان دیگر بر خلاف آنکه انتظار می‌رفت یک ارتباط نسبتاً قوی بین دو شاخص مذکور وجود داشته باشد (باستثنای بخشهایی که بر تکنولوژی‌های نوین که از نیروی انسانی کمتری استفاده می‌کنند، ملاحظه نمی‌شود).

-- برای بررسی این موضوع، طی بررسی تحلیل آماری همبستگی و آزمون‌های مربوطه بین این دو شاخص نتایجی بشرح زیر بدست آمد:

در بخش ماهیکیری مقدار همبستگی بدست آمده حدود ۰/۷ بوده و نتایج آزمون همبستگی، وجود ارتباط بین دو متغیر ارزش افزوده و درصد نیروی شاغلین را در این بخش تایید می‌کند.

یکی از نکات جالب بدست آمده آن بود که رابطه همبستگی بین دو شاخص مورد نظر در بخشهای بهداشت و درمان و ساختمان منفی و بترتیب حدود ۰/۱۵- و ۰/۳- بوده و در بخش آموزش حدود ۰/۱ می‌باشد که البته نتایج آزمون در سطح ۵ درصد همبستگی دو شاخص را رد می‌نماید.

این بررسی در مورد سایر بخشها نیز بیانگر یک ارتباط ضعیف ولی مثبت بوده، که مقدار همبستگی در آنها مجموعاً بین ۰/۱ تا ۰/۶ در نوسان بوده و در سطح ۵ درصد تایید می‌شوند.

موارد مذکور در واقع به نوعی مؤید پایین بودن میزان بهره‌وری درکشور و عدم بکارگیری مفید و بهینه از ظرفیتهای موجود در کلیه استانها بوده و لزوم بررسی همه جانبه آن از ضروری‌ترین

جدول ۱۵: رتبه بندی استانها بر حسب عامل اول برای درصد نیروی شاغل در فعالیتهای اقتصادی
 -- ۱۳۷۹ (تاکید بر اهمیت بخشهای ماهیگیری، حمل و نقل، تامین آب و ... ، امور عمومی و
 (...)

| رتبه | استان | امتیاز کسب شده |
|------|----------------------|----------------|
| ۱ | هرمزگان | ۳,۴۹ |
| ۲ | بوشهر | ۲,۶۷ |
| ۳ | خوزستان | ۱,۰۶ |
| ۴ | سیستان و بلوچستان | ۰,۸۴ |
| ۵ | کرمانشاه | ۰,۵۵ |
| ۶ | فارس | ۰,۱۷ |
| ۷ | آذربایجان غربی | ۰,۰۰ |
| ۸ | همدان | -۰,۰۲۸ |
| ۹ | گیلان | -۰,۰۲۸ |
| ۱۰ | تهران | -۰,۰۴۵ |
| ۱۱ | لرستان | -۰,۰۵۹ |
| ۱۲ | مازندران | -۰,۰۹۳ |
| ۱۳ | ایلام | -۰,۱۰۱ |
| ۱۴ | قزوین | -۰,۱۲۵ |
| ۱۵ | کردستان | -۰,۲۱۶ |
| ۱۶ | اردبیل | -۰,۲۵۰ |
| ۱۷ | کرمان | -۰,۴۴۶ |
| ۱۸ | سمنان | -۰,۴۷۸ |
| ۱۹ | خراسان | -۰,۵۴۱ |
| ۲۰ | آذربایجان شرقی | -۰,۵۷۸ |
| ۲۱ | گلستان | -۰,۵۸۷ |
| ۲۲ | اصفهان | -۰,۶۱۱ |
| ۲۳ | مرکزی | -۰,۶۷۹ |
| ۲۴ | قم | -۰,۶۹۳ |
| ۲۵ | کهگیلویه و بویر احمد | -۰,۷۰۵ |
| ۲۶ | چهارمحال و بختیاری | -۰,۷۶۶ |
| ۲۷ | زنجان | -۰,۸۳۰ |
| ۲۸ | یزد | -۰,۹۲۰ |

جدول ۱۶: رتبه‌بندی استانها بر حسب عامل دوم برای درصد نیروی شاغل در فعالیتهای اقتصادی -- ۱۳۷۹ (تاکید بر اهمیت بخشهای بهداشت و مددکاری و ... ، آموزش و معدن)

| رتبه | استان | امتیاز کسب شده |
|------|----------------------|----------------|
| ۱ | کهگیلویه و بویر احمد | ۲٫۸۵ |
| ۲ | سمنان | ۲٫۱۲ |
| ۳ | کرمان | ۱٫۶۲ |
| ۴ | ایلام | ۱٫۰۷ |
| ۵ | خوزستان | ۰٫۹۴ |
| ۶ | بوشهر | ۰٫۶۱ |
| ۷ | تهران | ۰٫۴۵ |
| ۸ | یزد | ۰٫۴۴ |
| ۹ | فارس | ۰٫۲۰ |
| ۱۰ | کرمانشاه | ۰٫۱۵ |
| ۱۱ | لرستان | ۰٫۰۷ |
| ۱۲ | مازندران | ۰٫۰۰ |
| ۱۳ | چهارمحال و بختیاری | -۰٫۱۲ |
| ۱۴ | سیستان | -۰٫۱۷ |
| ۱۵ | گیلان | -۰٫۲۷ |
| ۱۶ | گلستان | -۰٫۳۹ |
| ۱۷ | قم | -۰٫۴۴ |
| ۱۸ | اصفهان | -۰٫۵۱ |
| ۱۹ | خراسان | -۰٫۶۱ |
| ۲۰ | کردستان | -۰٫۶۱ |
| ۲۱ | هرمزگان | -۰٫۶۵ |
| ۲۲ | همدان | -۰٫۶۷ |
| ۲۳ | اردبیل | -۰٫۸۰ |
| ۲۴ | زنجان | -۰٫۸۱ |
| ۲۵ | مرکزی | -۰٫۹۰ |
| ۲۶ | آذربایجان غربی | -۰٫۹۹ |
| ۲۷ | قزوین | -۱٫۰۳ |
| ۲۸ | آذربایجان شرقی | -۱٫۵۳ |

جدول ۱۷: رتبه‌بندی استانها بر حسب عامل سوم برای درصد نیروی شاغل در فعالیتهای اقتصادی -- ۱۳۷۹ (تاکید بر اهمیت بخشهای کسب و کار، واسطه‌گریهای مالی و صنعت)

| رتبه | استان | امتیاز کسب شده |
|------|----------------------|----------------|
| ۱ | تهران | ۳,۳۲ |
| ۲ | یزد | ۱,۵۵ |
| ۳ | اصفهان | ۱,۳۹ |
| ۴ | قم | ۱,۳۲ |
| ۵ | بوشهر | ۰,۸۹ |
| ۶ | سمنان | ۰,۵۴ |
| ۷ | مرکزی | ۰,۵۰ |
| ۸ | آذربایجان شرقی | ۰,۴۳ |
| ۹ | خوزستان | ۰,۳۵ |
| ۱۰ | خراسان | ۰,۰۷ |
| ۱۱ | هرمزگان | ۰,۰۳ |
| ۱۲ | چهارمحال و بختیاری | ۰,۰۱ |
| ۱۳ | فزوین | -۰,۰۲ |
| ۱۴ | فارس | -۰,۲۱ |
| ۱۵ | کرمان | -۰,۳۹ |
| ۱۶ | مازندران | -۰,۴۸ |
| ۱۷ | لرستان | -۰,۵۱ |
| ۱۸ | زنجان | -۰,۵۲ |
| ۱۹ | گیلان | -۰,۵۲ |
| ۲۰ | همدان | -۰,۵۵ |
| ۲۱ | کرمانشاه | -۰,۶۵ |
| ۲۲ | کردستان | -۰,۷۴ |
| ۲۳ | سیستان و بلوچستان | -۰,۸۱ |
| ۲۴ | گلستان | -۰,۸۱ |
| ۲۵ | آذربایجان غربی | -۰,۸۵ |
| ۲۶ | اردبیل | -۰,۹۹ |
| ۲۷ | کهگیلویه و بویر احمد | -۱,۰۶ |
| ۲۸ | ایلام | -۱,۲۸ |

گامهای مطالعاتی مطالعات می‌باشد. در انتها بر اساس نتایج این تحقیق لزوم انجام مطالعات کلان اقتصادی در زمینه‌های: امکان برقراری توزیع عادلانه ثروت در کشور، استفاده از سازکارهای نوین در بکارگیری بهینه پتانسیل‌های اقتصادی موجود در کلیه استانها، اجرایی نمودن الگوهای افزایش بهره‌وری و کاهش هزینه‌ها و ... ضروری به نظر می‌رسد.

مراجع

- [۱] استنباط آماری چند متغیره، نارایان سی. جری، ترجمه ابوالقاسم بزرگ نیا - موسسه چاپ و انتشارات آستان قدس رضوی. ۱۳۶۶.
- [۲] آشنایی با روشهای چند متغیره، بی. اف. جی. مانلی، مترجم دکتر محمد مقدم، مهندس سید ابوالقاسم محمدی شوطی، مهندس مصطفی آقایی برزه، انتشارات پیشتاز علم، ۱۳۷۳.
- [۳] حسابهای منطقه‌ای، اصول نظری و کاربرد (ویرایش دوم)، دفتر حسابهای اقتصادی، مرکز آمار ایران، ۱۳۷۹.
- [۴] تحلیل آمار چند متغیره کاربردی، ریچارد. آ. جانسون، دین. دبلیو. دیچرن، ترجمه دکتر حسینعلی نیرومند، مشهد، دانشگاه فردوسی (مشهد)، ۱۳۷۸.
- [۵] حسابهای ملی ایران. حسابهای منطقه‌ای، حساب تولید استان‌های کشور ۱۳۷۹، دفتر حسابهای اقتصادی، مرکز آمار ایران، ۱۳۸۱.
- [۶] آمارگیری از ویژگیهای اشتغال و بیکاری خانوار ۱۳۷۹، مرکز آمار ایران، ۱۳۷۹.
- [۷] سالنامه آماری کشور - ۱۳۷۹، مرکز آمار ایران، اسوه ۱۳۸۰.

برآورد فاصله‌ای برای خانواده توزیع‌های نمایی طبیعی

محسن عارفی، غلامرضا محتشمی برزادران

گروه آمار دانشگاه بیرجند

چکیده: در این مقاله یک برآورد فاصله‌ای برای میانگین، در توزیع‌های نمایی طبیعی که دارای تابع واریانس درجه دوم بر حسب میانگین می‌باشد، مورد توجه قرار گرفته است. این خانواده توزیع‌های نمایی طبیعی شامل توزیع‌های دوجمله‌ای، دوجمله‌ای منفی، پواسن، نرمال، گاما و سکانت هذلولوی تعمیم یافته می‌باشد. در ضمن توزیع گوسین معکوس که جزء خانواده توزیع‌های نمایی طبیعی با تابع واریانس درجه سوم بر حسب میانگین می‌باشد، نیز مورد بررسی قرار گرفته است. برای این توریها، فاصله استاندارد والد و چهار فاصله اطمینان دیگر و مشخصه‌های مرتبط با آنها مورد ارزیابی قرار گرفته است. نتایج و محاسبات نشان می‌دهد که در میان این فواصل، فاصله نسبت درستنمایی کارایی بهتری را نشان می‌دهد.

۱ مقدمه

هدف از انجام این مقاله، بدست آوردن یک برآورد فاصله‌ای برای میانگین، در خانواده توزیع‌های نمایی طبیعی با تابع واریانس مورد نظر می‌باشد. بدین منظور در بخش دوم، خانواده توزیع‌های نمایی طبیعی و تابع واریانس درجه دوم و تابع واریانس درجه سوم را معرفی کرده‌ایم و ضرایب این تابع واریانسها را نیز مورد محاسبه قرار داده‌ایم.

در این مقاله ما بیشتر روی خانواده توزیع‌های نمایی طبیعی گسسته متمرکز می‌شویم و برای آنها انواع فواصل اطمینان را بدست می‌آوریم. همچنین در بین خانواده توزیع‌های نمایی طبیعی پیوسته، توزیع گوسین معکوس را که دارای تابع واریانس درجه سوم نسبت به میانگین می‌باشد، را بررسی می‌نمائیم.

در بخش سوم، انواع فواصل اطمینان را همراه با طول مورد انتظار و خاصیت اریبی را برای حالت گسسته با واریانس درجه دوم مورد ارزیابی قرار می‌دهیم. بخش چهارم شامل توزیع گوسین معکوس با تابع واریانس درجه سوم و فواصل مربوط به این توزیع، همراه با طول مورد انتظار می‌باشد.

۲ خانواده توزیع‌های نمایی طبیعی^۱

در این مقاله ما به دنبال بررسی یک برآورد فاصله‌ای برای میانگین، در خانواده توزیع‌های نمایی (NEF) با واریانس درجه دوم نسبت به میانگین و برای یک توزیع خاص با واریانس درجه سوم نسبت به میانگین می‌باشیم. خانواده توزیع‌های نمایی طبیعی با واریانس درجه دوم (NEF-QVF) شامل ۶ توزیع مهم «دوجمله‌ای، دوجمله‌ای منفی، پواسون، نرمال، گاما و توزیع سکانت هذلولوی تعمیم یافته (NEF-GHS)» می‌باشد (برای بررسی بیشتر در مورد توزیع سکانت هذلولوی تعمیم یافته به نوشته‌های موریس (۱۹۸۲) و برون (۱۹۸۶) مراجعه نمایید). همچنین خانواده توزیع‌های نمایی طبیعی با واریانس درجه سوم (NEF-CVF) علاوه بر ۶ توزیع فوق شامل ۶ توزیع دیگر می‌باشد، که ما در اینجا تنها توزیع گوسین معکوس را مورد بررسی قرار داده‌ایم (برای بررسی این توزیع‌ها به نوشته‌های لٹاک و مورا (۱۹۹۰) مراجعه نمائید). یک خانواده توزیع‌های نمایی طبیعی را به صورت

$$f(x/\xi) = \exp[\xi x - \psi(\xi)]h(x); \quad (۱)$$

نمایش می‌دهند، که به ξ پارامتر طبیعی می‌گویند. برای این خانواده توزیع‌ها

$$\mu = \psi'(\xi) \quad , \quad \sigma^2 = \psi''(\xi).$$

به ترتیب میانگین و واریانس می‌باشند. این واریانس تنها به میانگین μ بستگی دارد، که به صورت زیر نشان داده می‌شود:

$$\sigma^2 \equiv \text{var}(\mu) = a_0 + a_1\mu + a_2\mu^2 + a_3\mu^3; \quad (۲)$$

که برای تابع واریانس درجه دوم $a_3 = 0$ و a_0, a_1, a_2 می‌باشد. در نتیجه برای توزیع‌های مختلف این خانواده داریم:

• توزیع برنولی، $Bin(\lambda, p)$:

برای این توزیع نسبت‌های $\xi = \log \frac{p}{q}$ ، $\psi(\xi) = \log \frac{1}{q} = \log(\lambda + e^\xi)$ و $h(x) = 1$ با میانگین $\mu = p$ و تابع واریانس درجه دوم $v(\mu) = pq = \mu - \mu^2$ بدست می‌آید که ضرایب تابع واریانس به صورت $a_0 = 0$ ، $a_1 = 1$ و $a_2 = -1$ می‌باشد.

• توزیع دوجمله‌ای منفی، $NBin(\lambda, p)$ (تعداد موفقیت‌ها قبل از اولین شکست: x):
در این توزیع نسبت‌های زیر را داریم:

$$\xi = \log(p) \quad , \quad \psi(\xi) = -\log(\lambda - p) = -\log(\lambda - e^\xi) \quad , \quad h(x) = 1;$$

1) Natural Exponential Family (NEF)

همچنین میانگین به صورت $\mu = \frac{p}{q}$ و تابع واریانس بر حسب این میانگین به صورت $v(\mu) = \frac{p}{q^2} = \mu + \mu^2$ با ضرایب $a_0 = 0$ ، $a_1 = 1$ و $a_2 = 1$ می‌باشد.

• توزیع پواسن، $Poi(\lambda)$:

نسبت‌های فوق در این توزیع به صورت زیر می‌باشد:

$$\xi = \log \lambda, \quad \psi(\xi) = \lambda = e^\xi, \quad h(x) = \frac{1}{x!};$$

که تابع واریانس به صورت $v(\mu) = \lambda = \mu$ با ضرایب $a_0 = 0$ ، $a_1 = 1$ و $a_2 = 0$ می‌باشد.

• توزیع نرمال $N(\mu, \sigma^2)$ (σ^2 معلوم):

در این توزیع نسبت‌های مورد نظر به صورت زیر می‌باشد:

$$\xi = \frac{\mu}{\sigma^2}, \quad \psi(\xi) = \frac{\mu^2}{2\sigma^2} = \frac{\sigma^2}{2}\xi, \quad h(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{x^2}{2\sigma^2}\right\};$$

در این حالت تابع واریانس به صورت $v(\mu) = \sigma^2$ با ضرایب $a_0 = \sigma^2$ ، $a_1 = a_2 = 0$ می‌باشد.

• توزیع گاما، $Gamma(r, \lambda)$ (r معلوم):

روابط مورد نظر برای این توزیع به صورت

$$\xi = -\frac{1}{\lambda}, \quad \psi(\xi) = r \log \lambda = -r \log(-\xi), \quad h(x) = \frac{x^{r-1}}{\Gamma(r)};$$

با میانگین و تابع واریانس و با ضرایب $a_0 = a_1 = 0$ و $a_2 = \frac{1}{r}$ می‌باشد.

• توزیع سکانت هذلولوی تعمیم یافته، $GHS(r, \lambda)$ (r معلوم):

برای بررسی تابع چگالی این توزیع می‌توانید به نوشته‌های موریس (۱۹۸۲) و برون (۱۹۸۶) مراجعه نمایید. با توجه به تابع چگالی آن، نسبت‌های زیر بدست می‌آید:

$$\xi = \tan^{-1}(\lambda), \quad \psi(\xi) = r \log(1 + \lambda^2) = -r \log(\cos \xi), \quad h(x) = f_{r,0}(x);$$

برای این توزیع میانگین و تابع واریانس به صورت زیر بدست می‌آید:

$$\mu = \psi'(\xi) = r \tan(\xi) = r\lambda, \quad \nu(\mu) = \frac{r}{\cos^2(\xi)} = r + \frac{\mu^2}{r};$$

که ضرایب این تابع واریانس به صورت $a_0 = r$ ، $a_1 = 0$ و $a_2 = \frac{1}{r}$ می‌باشد.

• توزیع گوسین معکوس، $IG(m, \lambda)$ (λ معلوم):
 در این توزیع روابط زیر را می‌توان از فرم نمایی آن نتیجه گرفت:

$$\xi = \frac{-\lambda}{\sqrt{2}m}, \quad \psi(\xi) = \frac{-\lambda}{m} = -(-2\lambda\xi)^{1/2}, \quad h(x) = \sqrt{\frac{\lambda}{2\pi x^3}} \exp\left\{\frac{-\lambda}{2x}\right\};$$

که میانگین و تابع واریانس درجه سوم آن به صورت زیر محاسبه می‌گردد:

$$\mu = \psi'(\xi) = m, \quad \nu(\mu) = \psi''(\xi) = \frac{m^3}{\lambda} = \frac{\mu^3}{\lambda}$$

در نتیجه ضرایب تابع واریانس برابر $a_3 = \frac{1}{\lambda}$ و $a_0 = a_1 = a_2 = 0$ می‌باشد.

اکنون ما روی توزیع‌های نمایی گسسته با واریانس درجه دوم متمرکز می‌شویم و از میان توزیع‌های پیوسته، توزیع گوسین معکوس را که دارای تابع واریانس درجه سوم می‌باشد، نیز مورد بررسی قرار می‌گیرد.

۳ خانواده توزیع‌های نمایی طبیعی گسسته با واریانس درجه دوم

این خانواده توزیعها شامل سه توزیع دوجمله‌ای، دوجمله‌ای منفی، پواسن می‌باشد. برای این توزیعها ما انواع فواصل اطمینان و بعضی روابط خاص را بررسی می‌کنیم.

۱.۳ بسط اریبی

فاصله استاندارد والد به صورت $\hat{\mu} \pm z_{\alpha} n^{-1/2} \sqrt{a_0 + a_1 \hat{\mu} + a_2 \hat{\mu}^2}$ می‌باشد، که بر اساس کمیت محوری زیر پایگذاری می‌شود:

$$W_n = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sqrt{a_0 + a_1 \hat{\mu} + a_2 \hat{\mu}^2}} \xrightarrow{L} N(0, 1)$$

که در آن $\hat{\mu} = \bar{X}$ ، یک برآورد MLE از μ می‌باشد. در واقع به طور کلی W_n دارای توزیع مجانبی نرمال استاندارد می‌باشد. از آنجایی که $Z = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sqrt{a_0 + a_1 \hat{\mu} + a_2 \hat{\mu}^2}}$ دارای توزیع نرمال استاندارد می‌باشد، می‌توانیم میزان اریبی را محاسبه نمائیم. از یک عبارت جبری ساده رابطه بین W_n و Z در توزیعهایی با تابع واریانس درجه دوم به صورت زیر در می‌آید:

$$W_n = Z(1 + (a_1 + 2a_2\mu)\sigma^{-1}n^{-1/2}Z + a_2n^{-1}Z^2)^{-1/2}$$

این فرآیند را برای سه توزیع گسسته، می‌توان در نوشته‌های برون، کی وداس گیوپتا (۱۹۹۹b) به صورت زیر ملاحظه نمود، که این روابط براساس سری مکلاورن برای $f(Z) = W_n$ بدست می‌آید.

• برنولی، $Bin(1, p)$:

$$E(W_n) = \frac{p - \frac{1}{2}}{\sqrt{npq}} \left(1 + \frac{p}{2n} + \frac{9(p - \frac{1}{2})^2}{2npq} \right) + O(n^{-2}) \quad (3)$$

• دوجمله‌ای منفی، $NBin(1, p)$:

$$E(W_n) = -\frac{1+p}{2\sqrt{np}} \left(1 + \frac{1}{n} + \frac{9q^2}{\lambda np} \right) + O(n^{-2}) \quad (4)$$

• پواسن، $Poi(\lambda)$:

$$E(W_n) = -\frac{1}{2\sqrt{n\lambda}} \left(1 + \frac{9}{\lambda n \lambda} \right) + O(n^{-2}) \quad (5)$$

با توجه به رابطه (۳) معلوم می‌گردد که در حالت دوجمله‌ای، W_n دارای اریبی منفی برای $p < \frac{1}{2}$ و دارای اریبی مثبت برای $p > \frac{1}{2}$ می‌باشد و در نتیجه فاصله والد حول $\frac{1}{2}$ نوسان می‌کند. از معادله (۴) و (۵) این طور به نظر می‌رسد که W_n همواره دارای اریبی منفی می‌باشد.

۲.۳ انواع فاصله‌های اطمینان

فرض کنید $\hat{\mu} = \bar{X}$ برآورد MLE از μ باشد. با توجه به قضیه حد مرکزی^۲ و قضیه اسلوتسکی^۳ روابط زیر حاصل می‌شود:

$$W_n = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\hat{\sigma}} = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sqrt{a_0 + a_1\hat{\mu} + a_2\hat{\mu}^2}} \xrightarrow{L} N(0, 1); \quad (6)$$

$$Z = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sigma} = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sqrt{a_0 + a_1\mu + a_2\mu^2}} \xrightarrow{L} N(0, 1); \quad (7)$$

بر اساس دو تقریب (۶) و (۷) ما برخی از فاصله‌های اطمینان را به صورت زیر پایگذاری می‌کنیم: فاصله #۱: فاصله استاندارد والد بر اساس رابطه (۶) به صورت زیر بدست می‌آید:

$$CI_S = \hat{\mu} \pm k\hat{\sigma}n^{-\frac{1}{2}} = \hat{\mu} \pm k(a_0 + a_1\hat{\mu} + a_2\hat{\mu}^2)^{\frac{1}{2}}n^{-\frac{1}{2}}; \quad (8)$$

2) Central Limit Theorem 3) Slutsky's Theorem

فاصله # ۲: فاصله اسکور بر اساس تقریب (۷) پایگذاری می‌شود. این فاصله به وسیله وارونی از آزمون اسکور با دمه‌های برابر را‌تو تحت فرض $\mu = \mu_0$: H_0 بدست می‌آید. بنابراین فرض H_0 پذیرفته می‌شود، اگر μ در فاصله مورد نظر باشد. با حل یک معادله درجه دوم براساس تقریب (۷) فاصله اسکور بدست می‌آید:

$$CI_R = \tilde{\mu} \pm \frac{kn^{\frac{1}{\nu}}}{n - k^2 a_{\nu}} = (a_0 + a_1 \tilde{\mu} + a_2 \tilde{\mu}^2 + \frac{k^2}{4n} \Delta)^{\frac{1}{\nu}}; \quad (9)$$

که $\tilde{\mu} = \frac{n\bar{\mu} + k^2 a_1 / 2}{n - k^2 a_{\nu}}$ می‌باشد.

فاصله # ۳: فاصله اگرستی-کویل (AC) شکلی شبیه فاصله استاندارد والد دارد، با این تفاوت که به جای $\hat{\mu}$ مقدار $\tilde{\mu}$ و به جای n مقدار تصحیح شده $\tilde{n} = n - k^2 a_{\nu}$ را قرار می‌دهیم:

$$CI_{AC} = \tilde{\mu} \pm k(\nu(\tilde{\mu}))^{\frac{1}{\nu}} \tilde{n}^{-\frac{1}{\nu}} = \tilde{\mu} \pm k(a_0 + a_1 \tilde{\mu} + a_2 \tilde{\mu}^2)^{\frac{1}{\nu}} \tilde{n}^{-\frac{1}{\nu}}; \quad (10)$$

فاصله # ۴: فاصله نسبت درست‌نمایی به وسیله وارونی از آزمون نسبت درست‌نمایی با پذیرش فرض صفر $\mu = \mu_0$: H_0 ، اگر $-2 \log \Lambda_n \leq \chi_{\alpha, 1}^2$ باشد، بدست می‌آید. که در آن

$$\Lambda_n = \frac{L(\mu_0)}{\sup L_{\mu}(\mu)}$$

نسبت درست‌نمایی و L تابع درست‌نمایی بر اساس n مشاهده مستقل و هم‌توزیع می‌باشد. برای مطالعه بیشتر در این زمینه می‌توانید به نوشته‌های را‌تو (۱۹۷۳) و سرفلینگ (۱۹۸۰) مراجعه نمائید.

فاصله # ۵: فاصله جفریز:

فرض کنید $b(\cdot) = (\psi')^{-1}(\cdot)$ باشد، به وضوح آشکار است که $b(\mu) = \xi$ می‌باشد. فرض کنیم:

$$I(\mu) = n(a_0 + a_1 \mu + a_2 \mu^2)^{-1}$$

تابع پیشین جفریز را متناسب با $I^{1/2}(\mu)$ می‌گیریم، بنابراین تابع پسین به صورت زیر بدست می‌آید:

$$f(\mu|x) \propto \exp\{xb(\mu) - n\psi(b(\mu)) - \frac{1}{\nu} \log(a_0 + a_1 \mu + a_2 \mu^2)\}$$

فاصله جفریز با دمه‌های برابر برای μ از تابع پسین جفریز، به صورت زیر بدست می‌آید:

$$CI_J = \left[J_{\frac{\alpha}{\nu}}, J_{1-\frac{\alpha}{\nu}} \right]$$

که $J_{\frac{\alpha}{q}}$ و $J_{1-\frac{\alpha}{q}}$ چندکهای $\frac{\alpha}{q}$ و $1 - \frac{\alpha}{q}$ تحت توزیع پسین جفریز و بر مبنای n مشاهده می‌باشد (برای مطالعه بیشتر به نوشته‌های کس و واسرمن (۱۹۹۶) مراجعه نمائید). اکنون فاصله جفریز را برای سه توزیع گسسته محاسبه می‌نمائیم:

دوجمله‌ای: در این حالت $\psi(\xi) = \log(1 + e^\xi)$ و $b(\mu) = \log \frac{\mu}{1-\mu}$ می‌باشد. بنابراین تابع پیشین جفریز دارای توزیع $\beta(\frac{1}{q}, \frac{1}{q})$ و تابع پسین دارای توزیع $\beta(X + \frac{1}{q}, n - X + \frac{1}{q})$ با $X = \sum_{i=1}^n X_i$ می‌باشد. بنابراین فاصله جفریز با دمهای برابر برای در سطح $100(1 - \alpha)$ درصد به صورت زیر بدست می‌آید:

$$CI_J = [p_L, p_U] = [B_{\{\frac{\alpha}{q}, X + \frac{1}{q}, n - X + \frac{1}{q}\}}, B_{\{1 - \frac{\alpha}{q}, X + \frac{1}{q}, n - X + \frac{1}{q}\}}]$$

دوجمله‌ای منفی: در این حالت $\psi(\xi) = -\log(1 - e^\xi)$ و $b(\mu) = \log \frac{\mu}{1+\mu}$ می‌باشد. تابع پیشین برای μ متناسب با $\mu^{-\frac{1}{q}}(1 + \mu)^{-\frac{1}{q}}$ می‌باشد، که یک تابع بتا پریم نامیده می‌شود (برای اطلاعات بیشتر به جانسون (۱۹۹۵) مراجعه نمائید). برای حل این مشکل ابتدا یک فاصله اطمینان برای p بدست می‌آوریم. فاصله پیشین برای p متناسب با $p^{-\frac{1}{q}}(1 - p)^{-1}$ و تابع پسین دارای توزیع $Beta(X + 1/2, n)$ می‌باشد. پس فاصله اطمینان برای p در سطح $100(1 - \alpha)$ درصد برابر است با:

$$CI_J^p = [p_L, p_U] = [B_{\{\alpha/2, X + 1/2, n\}}, B_{\{1 - \alpha/2, X + 1/2, n\}}]$$

از آنجایی که $\mu = p/(1 - p)$ می‌باشد، پس فاصله جفریز برای μ عبارتست از:

$$CI_J = [p_L/(1 - p_L), p_U/(1 - p_U)]$$

یواسن: در این حالت $\psi(\xi) = e^\xi$ و $b(\mu) = \log \mu$ می‌باشد. تابع پیشین متناسب با $\lambda^{-\frac{1}{q}}$ و تابع پسین دارای توزیع $Gamma(X + 1/2, 1/n)$ می‌باشد. بنابراین فاصله اطمینان جفریز با دمهای برابر برای λ در سطح $100(1 - \alpha)$ درصد برابر با

$$CI_J = [\lambda_L, \lambda_U] = [G_{\{\frac{\alpha}{q}, X + \frac{1}{q}, \frac{1}{n}\}}, G_{\{1 - \frac{\alpha}{q}, X + \frac{1}{q}, \frac{1}{n}\}}]$$

خواهد بود.

۳.۳ طول فواصل اطمینان مورد انتظار

در میان فاصله اطمینانهای بدست آمده، کمترین طول برای یک فاصله اطمینان، همواره موضوع مورد اهمیتی بوده است. بنابراین در میان کلیه فواصل اطمینان، امید طولهای آنها را محاسبه و با یکدیگر مورد مقایسه قرار می‌دهیم.

قضیه: فرض کنید CI یک نمایش کلی برای هر یک از پنج فاصله اطمینان $CI_R, CI_S, CI_{LR}, CI_{AC}$ و CI_J باشد، امید طولها برابر است با:

$$L(n, \mu) = E(L_{CI}) = \sqrt{2}k(\mu + a\sqrt{\mu^2})^{1/2}n^{-1/2}\left(1 - \frac{\delta(k, \mu)}{\sqrt{2}n(\mu + a\sqrt{\mu^2})}\right) + O(n^{-2})$$

که $\delta(k, \mu)$ به صورت زیر محاسبه می شود:

$$\begin{aligned} \delta(k, \mu) &= 9\Delta \text{ for } CI_S; \\ &= 9(1 - k^2)\Delta - \sqrt{2}k^2a\sqrt{\mu + a\sqrt{\mu^2}} \text{ for } CI_R; \\ &= 9(1 - 2k^2)\Delta - 10k^2a\sqrt{\mu + a\sqrt{\mu^2}} \text{ for } CI_{AC}; \\ &= (9 - 2k^2)\Delta - 26k^2a\sqrt{\mu + a\sqrt{\mu^2}} \text{ for } CI_{LR}; \\ &= (5 - 2k^2)\Delta - 2(13k^2 + 17)a\sqrt{\mu + a\sqrt{\mu^2}} \text{ for } CI_J; \end{aligned}$$

نتیجه (۱): حالت پواسن را در نظر بگیرید. طول مورد انتظار برای انواع فاصله اطمینانها عبارتست از:

$$\begin{aligned} E(L_S) &= \sqrt{2}k\lambda^{1/2}n^{-1/2}\left(1 - \frac{9}{\sqrt{2}n\lambda}\right) + O(n^{-2}); \\ E(L_{LR}) &= \sqrt{2}k\lambda^{1/2}n^{-1/2}\left(1 + \frac{9(k^2 - 1) - \sqrt{2}k^2}{\sqrt{2}n\lambda}\right) + O(n^{-2}); \\ E(L_J) &= \sqrt{2}k\lambda^{1/2}n^{-1/2}\left(1 + \frac{9(k^2 - 1) + 4 - \sqrt{2}k^2}{\sqrt{2}n\lambda}\right) + O(n^{-2}); \\ E(L_R) &= \sqrt{2}k\lambda^{1/2}n^{-1/2}\left(1 + \frac{9(k^2 - 1)}{\sqrt{2}n\lambda}\right) + O(n^{-2}); \\ E(L_{AC}) &= \sqrt{2}k\lambda^{1/2}n^{-1/2}\left(1 + \frac{9(2k^2 - 1)}{\sqrt{2}n\lambda}\right) + O(n^{-2}); \end{aligned}$$

با توجه به فرآیندهای بالا معلوم می شود که برای هر

$$k > \frac{2}{\sqrt{7}} = 0,76$$

فواصل CI_{AC} و $CI_R, CI_J, CI_{LR}, CI_S$ به ترتیب دارای کوچکترین تا بزرگترین طول می باشند. بنابراین فاصله نسبت درستی دارای کوچکترین طول می باشد.

نتیجه (۲): برآورد فاصله‌ای برای $\mu = p/q$ را در حالت دوجمله‌ای در نظر بگیرید. طول مورد انتظار برای فواصل مختلف به صورت زیر می‌باشد:

$$\begin{aligned}
 E(L_s) &= 2kp^{1/2}q^{-1}n^{-1/2}\left(1 - \frac{9q^2}{\sqrt{2}np}\right) + O(n^{-2}); \\
 E(L_{LR}) &= 2kp^{1/2}q^{-1}n^{-1/2}\left(1 - \frac{9q^2 - 2k^2(1 + 11p + p^2)}{\sqrt{2}np}\right) + O(n^{-2}); \\
 E(L_J) &= 2kp^{1/2}q^{-1}n^{-1/2} \\
 &\quad \left(1 - \frac{9q^2 - 2k^2(1 + 11p + p^2) - 2(2 + 13p + 2p^2)}{\sqrt{2}np}\right) + O(n^{-2}); \\
 E(L_R) &= 2kp^{1/2}q^{-1}n^{-1/2}\left(1 - \frac{9q^2 - 9k^2(1 + 6p + p^2)}{\sqrt{2}np}\right) + O(n^{-2}); \\
 E(L_{AC}) &= 2kp^{1/2}q^{-1}n^{-1/2}\left(1 - \frac{9q^2 - 18k^2(1 + 4p + p^2)}{\sqrt{2}np}\right) + O(n^{-2});
 \end{aligned}$$

با توجه به بسط‌های بالا نتیجه می‌شود که برای هر $p > 0$ و برای هر $k > \sqrt{17/23} = 0,86$ فواصل $CI_{AC}, CI_R, CI_J, CI_{LR}, CI_S$ به ترتیب دارای کوچکترین تا بزرگترین طول می‌باشند. بنابراین فاصله نسبت درستی‌مایی دارای کوچکترین طول مورد انتظار می‌باشد. بر خلاف دو حالت پواسن و دوجمله‌ای منفی، یک ترتیب یکنواختی برای طول مورد انتظار در حالت دوجمله‌ای وجود ندارد. ولی اگر نسبت به p از طول‌های مورد انتظار انتگرال بگیریم، یک ترتیب واضحی به صورت زیر وجود دارد. نتیجه (۳): حالت دوجمله‌ای را در نظر بگیرید. انتگرالی از طول‌های مورد انتظار برای فواصل مختلف به صورت زیر می‌باشد:

$$\begin{aligned}
 \int_0^1 E(L_J)dp &= \frac{k\pi}{4}n^{-1/2} - \left(\frac{37}{36} + \frac{5k^2}{36}\right)\frac{k\pi}{4}n^{-3/2} + O(n^{-2}); \\
 \int_0^1 E(L_{LR})dp &= \frac{k\pi}{4}n^{-1/2} - \left(1 + \frac{5k^2}{36}\right)\frac{k\pi}{4}n^{-3/2} + O(n^{-2}); \\
 \int_0^1 E(L_R)dp &= \frac{k\pi}{4}n^{-1/2} - \frac{k\pi}{4}n^{-3/2} + O(n^{-2}); \\
 \int_0^1 E(L_s)dp &= \frac{k\pi}{4}n^{-1/2} - \frac{k\pi}{4}n^{-3/2} + O(n^{-2}); \\
 \int_0^1 E(L_{AC})dp &= \frac{k\pi}{4}n^{-1/2} - \left(\frac{k^2}{2} - 1\right)\frac{k\pi}{4}n^{-3/2} + O(n^{-2});
 \end{aligned}$$

بنابراین فواصل CI_J ، CI_{LR} ، CI_R ، CI_S و CI_{AC} به ترتیب از کوچکترین به بزرگترین طول می‌باشد. به طور تقریبی فواصل CI_J و CI_{LR} با هم برابر می‌باشند، پس فاصله نسبت درستی طول بهتری دارد.

۴ فاصله اطمینان برای تابع توزیعی با واریانس درجه سوم

خانواده توزیع‌های نمایی طبیعی با واریانس درجه سوم (NEF-CVF) شامل ۶ توزیع علاوه بر توزیع‌های دیگر می‌باشد که ما در میان این ۶ توزیع تنها بعضی خواص توزیع گوسین معکوس را مورد بررسی قرار می‌دهیم (برای بررسی این توزیعها به نوشته‌های لتاک و مورا (۱۹۹۰) مراجعه نمائید).

برای توزیع گوسین معکوس $IG(m, \lambda)$ که λ معلوم می‌باشد، میانگین برابر $\mu = m$ و واریانس تحت میانگین به صورت زیر محاسبه می‌شود:

$$var(\mu) = \frac{m^3}{\lambda} = \frac{\mu^3}{\lambda};$$

در نتیجه ضرایب تابع واریانس درجه سوم برابر $a_0 = a_1 = a_2 = 0$ و $a_3 = \frac{1}{\lambda}$ می‌باشد. با توجه به اینکه $W_n = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\hat{\sigma}^2}$ می‌باشد، براساس یک عبارت جبری ساده و با توجه به تابع واریانس درجه سوم W_n به صورت زیر تبدیل می‌شود:

$$W_n = Z(1 + (a_1 + 2a_2\mu + 3a_3\mu^2)\sigma^{-1}n^{-1/2}Z + (a_2 + 3a_3\mu)n^{-1}Z^2 + a_3n^{-3/2}Z^3\sigma)^{-1/2}$$

با توجه به ضرایب تابع واریانس درجه سوم در توزیع گوسین معکوس داریم:

$$W_n = Z(1 + 3a_3\mu^2\sigma^{-1}n^{-1/2}Z + 3a_3n^{-1}Z^2\mu + a_3n^{-3/2}Z^3\sigma)^{-1/2}$$

از سری مکلاورن برای $f(Z) = W_n$ می‌توانیم میزان آریبی را محاسبه نمود:

$$E(W_n) = -\frac{3}{4}n^{-1/2}\mu^{1/2}\lambda^{-1/2}\left(1 + \frac{35\mu}{\lambda n\lambda}\right) + O(n^{-2});$$

پس آریبی در توزیع گوسین معکوس همواره منفی می‌باشد.

۵ انواع فاصله اطمینانهای برای توزیع گوسین معکوس

برای این توزیع ما سه فاصله اطمینان را به صورت زیر بدست می‌آوریم:
 □ فاصله استاندارد والد: با توجه به تقریب اسلوتسکی این فاصله به صورت زیر بدست می‌آید:

$$CI_S = \hat{\mu} \pm kn^{-\frac{1}{2}} \hat{\sigma} = \hat{\mu} \pm kn^{-\frac{1}{2}} \hat{\mu}^{\frac{1}{2}} \lambda^{-\frac{1}{2}}$$

که $\hat{\mu} = \bar{X}$ برآورد MLE از μ می‌باشد.
 □ فاصله نسبت درستنمایی: این فاصله با توجه به نسبت درستنمایی Λ_n و با پذیرش فرض $H_0: \mu = \mu_0$ و بر اساس رابطه $-\log \Lambda_n \leq k^2$ به صورت زیر بدست می‌آید:

$$\begin{aligned} \Lambda_n &= \frac{L(\mu_0)}{\sup_{\mu} L(\mu)} = \frac{\exp\{\frac{-\lambda}{\mu_0} \sum_{i=1}^n X_i + \frac{n\lambda}{\mu_0}\} \prod_{i=1}^n h(X_i)}{\exp\{\frac{-\lambda}{\mu} \sum_{i=1}^n X_i + \frac{n\lambda}{\mu}\} \prod_{i=1}^n h(X_i)} \\ &= \exp\{\frac{-n\lambda}{\mu_0} \bar{X} + \frac{n\lambda}{\mu_0} - \frac{n\lambda}{\mu}\}; \end{aligned}$$

در نتیجه می‌توان نوشت:

$$-\log \Lambda_n = -\log \left\{ \frac{-n\lambda}{\mu_0} \bar{X} + \frac{n\lambda}{\mu_0} - \frac{n\lambda}{\mu} \right\} \leq k^2$$

با حل یک معادله درجه دوم فاصله مورد نظر برای $\frac{\mu_0}{\mu} \leq k^2 \leq \frac{n\lambda}{\mu_0}$ به صورت زیر بدست می‌آید:

$$CI_{LR} = \frac{n\hat{\mu}}{n - k^2 a_{\mu} \hat{\mu}} \pm \frac{k\hat{\mu} \sqrt{na_{\mu} \hat{\mu}}}{n - k^2 a_{\mu} \hat{\mu}}$$

□ فاصله اسکور: این فاصله اطمینان بر اساس قضیه حد مرکزی بر اساس آزمون اسکور با دماهای برابر و با پذیرش فرض صفر $H_0: \mu = \mu_0$ محاسبه می‌شود. با توجه به قضیه حد مرکزی می‌توان نوشت:

$$Z = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sigma} = \frac{\sqrt{n\lambda}(\hat{\mu} - \mu)}{\sqrt{\mu^3}} \leq k$$

برای بدست آوردن این فاصله باید یک معادله درجه سوم را محاسبه نمود. ابتدا تغییر متغیر $y = \sqrt{\mu}$ را در رابطه بالا اعمال می‌کنیم:

$$\frac{\sqrt{n\lambda}(\hat{\mu} - y^2)}{y^3} \leq k \Rightarrow \sqrt{n\lambda}(\hat{\mu} - y^2) \leq ky^3 \Rightarrow ky^3 + \sqrt{n\lambda}y^2 - \sqrt{n\lambda}\hat{\mu} \geq 0;$$

با تغییر متغیر $y = \frac{1}{x}$ معادله درجه سوم زیر حاصل می‌شود:

$$x^3 - \frac{1}{\hat{\mu}} - \frac{k}{\sqrt{n\lambda\hat{\mu}}} \leq 0$$

با توجه به مقدار $\frac{n\lambda}{\mu} \leq k^2 \leq \frac{n\lambda}{\hat{\mu}}$ داریم:

$$\Phi = \frac{-\mu}{\hat{\mu}^3} + \frac{2\sqrt{k}^2}{n\lambda\hat{\mu}^2} > 0$$

از دستور کاردان برای معادلات درجه سوم رابطه زیر حاصل می‌گردد:

$$x < \left(\frac{k}{\sqrt{n\lambda\hat{\mu}}} + \sqrt{\frac{\Phi}{10\lambda}} \right)^{1/3} + \left(\frac{k}{\sqrt{n\lambda\hat{\mu}}} - \sqrt{\frac{\Phi}{10\lambda}} \right)^{1/3}$$

اکنون اگر تغییر متغیرهای فوق را اعمال نمائیم، یک فاصله اطمینان یکطرفه به صورت زیر بدست می‌آید:

$$\left[\left(\frac{k}{\sqrt{n\lambda\hat{\mu}}} + \sqrt{\frac{\Phi}{10\lambda}} \right)^{1/3} + \left(\frac{k}{\sqrt{n\lambda\hat{\mu}}} - \sqrt{\frac{\Phi}{10\lambda}} \right)^{1/3} \right]^{-2} \leq \mu$$

۱.۵ طول فاصله اطمینان مورد انتظار برای توزیع گوسین معکوس

• فاصله استاندارد والد: طول این فاصله اطمینان عبارتست از:

$$L_s = 2kn^{-1/2}\lambda^{-1/2}\hat{\mu}^{3/2} = 2kn^{-1/2}\lambda^{-1/2}\left(Z\frac{\sigma}{\sqrt{n}} + \mu\right)^{3/2}$$

با استفاده از سری مکاورن برای رابطه بالا، طول مورد انتظار به صورت زیر محاسبه می‌شود:

$$E(L_s) = 2kn^{-1/2}\lambda^{-1/2}\mu^{3/2}\left(1 + \frac{3\mu}{\lambda n\lambda}\right) + O(n^{-2})$$

• فاصله اطمینان نسبت درست‌نمایی: طول این فاصله اطمینان عبارتست از:

$$L_{LR} = 2kn^{-1/2}\lambda^{-1/2}\hat{\mu}^{3/2}(n - k^2\lambda^{-1}\hat{\mu})^{-1}$$

در نتیجه بر اساس سری مکاورن، طول مورد انتظار برابر است با:

$$E(L_{LR}) = \frac{2kn^{-1/2}a_3^{3/2}\mu^{3/2}}{n - k^2a_3} \left(1 + \frac{3\mu}{4} + \frac{3k^2a_3\mu^2}{n - k^2a_3\mu} + \frac{2k^4a_3^2\mu^3}{(n - k^2a_3\mu)^2}\right) + O(n^{-2})$$

طول فاصله یکطرفه اسکور بدست آمده، بینهایت می‌باشد. با مقایسه طولهای مورد انتظار بالا به نظر می‌رسد که فاصله نسبت درست‌نمایی دارای کوچکترین طول می‌باشد.

۶ نتیجه گیری

نتایج و محاسبات بالا نشان می دهد که فاصله استاندارد والد خیلی ضعیف عمل می کند. بدست آوردن دو فاصله اطمینان نسبت درستی و جفریز کار ساده ای نیست، ولی این دو فاصله اطمینان کارایی بهتری را نشان می دهند و همچنین دارای کوتاهترین طول نیز می باشند. فاصله های اطمینان اسکور و AC در مقایسه با دو فاصله قبلی برای توزیعی با تابع واریانس درجه دوم، راحتتر بدست می آید، ولی برای توزیعی با تابع واریانس درجه سوم، نیاز به حل معادلات درجه سوم می باشد، که به سادگی قابل انجام نمی باشد.

مراجع

- [1] Agresti, A. & Coull, B.A. (1998). Approximate is better than exact for interval estimation of binomial proportions. *The American Statistician* 52, 119-126.
- [2] Brown, L.D. (1986). *Fundamental of statistical exponential families with applications in statistical decision theory*. Lecture Notes-Monograph Series, Institute of Mathematical Statistics, Hayward, California.
- [3] Brown, L.D., Cai, T. & DasGupta, A. (1999b). Confidence interval for a binomial proportion and Edgeworth expansions. Technical Report, Department of Statistical, Purdue University. Submitted.
- [4] Brown, L.D., Cai, T. & DasGupta, A. (2003). Interval estimation in exponential families. Technical Report, University of Pennsylvania and Purdue University. Submitted.
- [5] Johnson, N.L., Kotz, S. & Balakrishnan, N. (1995). *Continuous univariate distributions*. second edition, Vol. 2. Wiley, New York.
- [6] Kass, R.E. & Wasserman, L. (1996). The selection of prior distribution by formal rules. *J. Amer. Stat. Assoc.*, 91, 1343-70.
- [7] Letac, G. & Mora, M. (1990). Natural real exponential families with cubic variance functions.
- [8] Morris, C.N. (1982). Natural exponential families with quadratic variance functions. *Ann. Statist.* 10, 65-80.
- [9] Morris, C.N. (1983). Natural exponential families with quadratic variance functions: statistical theory. *Ann. Statist.* 11, 515-529.
- [10] Rao, C.R. (1973). *Linear statistical inference and its applications*. Wiley, New York.

- [11] Serfling, R.J. (1980). Approximation theorems of mathematical statistics. Wiley, New York.

روش درست‌نمایی وزنی در برآوردیابی برای ناحیه کوچک

ملیحه عباس نژاد مشهدی

گروه آمار، دانشکده علوم ریاضی دانشگاه فردوسی مشهد

چکیده: برآوردیابی برای ناحیه کوچک^۱ در سالهای اخیر مورد توجه زیادی قرار گرفته است، به دلیل این که تقاضا برای برآوردگرهای معتبر نواحی کوچک هم از جانب بخش‌های دولتی و هم از جانب بخش‌های خصوصی افزایش یافته است. از بررسیهای نمونه‌ای می‌توان برآوردهای مستقیم معتبری برای پارامترهای کل جامعه به دست آورد اما برآوردگرهای مستقیم برای هر ناحیه کوچک یعنی برآوردگرهایی که تنها بر مبنای داده‌های موجود در ناحیه به دست می‌آیند به دلیل کوچکی حجم نمونه در هر ناحیه به خطاهای استاندارد بزرگی منجر می‌شوند. بنابراین ضروری به نظر می‌رسد روشهایی برای این منظور جایگزین شوند که علاوه بر اطلاعات موجود در خود ناحیه، اطلاعات سایر نواحی (یا اطلاعات سرشماری) را نیز مورد استفاده قرار دهند. در این مقاله ابتدا به غیر مجاز بودن برآوردگر درست‌نمایی ماکزیمم میانگین نرمال در برآوردیابی همزمان^۲ یعنی \bar{X} اشاره‌ای کوتاه خواهیم داشت، روش درست‌نمایی وزنی^۳ را در خانواده‌ی نمایی به طور کلی بیان کرده و سپس آن را در مسأله برآوردیابی برای ناحیه کوچک مورد بحث قرار می‌دهیم.

واژه‌های کلیدی: ناحیه کوچک، درست‌نمایی وزنی، برآوردگر جیمز استاین^۴

۱ مقدمه

ناحیه کوچک (یا ناحیه محلی) به یک ناحیه جغرافیایی کوچک مانند یک استان، شهر یا یک بخش سرشماری اطلاق می‌گردد. همچنین می‌توان یک زیرجامعه کوچک همانند یک گروه از مردم درون یک ناحیه جغرافیایی بزرگ را که از سن و جنس و نژاد معینی هستند به عنوان یک ناحیه کوچک در نظر گرفت.

مسأله برآوردیابی برای ناحیه کوچک مسأله‌ای است که در سالهای اخیر توجه زیادی به آن شده است. این امر به این علت است که تقاضا برای آماره‌های نواحی کوچک چه از جانب بخش‌های دولتی و چه از جانب بخش‌های خصوصی رشد زیادی داشته است. به عنوان مثال بخش دولتی برای برنامه‌ریزی کلان در سطح کشور نیاز به آماره‌هایی از قبیل نرخ بیکاری، متوسط درآمد خانوارها و . . . در استانها و شهرهای مختلف دارد. به عنوان مثال فرض کنید که بعضی

1) Small Area 2) Simultaneous Estimation 3) Weighted Likelihood 4) James-stein Estimator

از نواحی جغرافیایی یا زیرگروههایی از جامعه از قبیل بعضی استانها و یا قشر خاصی از افراد از نظر سطح درآمد پائینتر از یک میانگین قابل قبول و در حد انتظار باشند. واضح است که قبل از انجام هر کاری و هر نوع برنامه‌ریزی لازم است که چنین نواحی یا زیرگروههایی مشخص شوند. داده‌های به دست آمده از بررسیهای نمونه‌ای را می‌توان برای استخراج برآوردهای مستقیم معتبر برای نواحی بزرگ یا قلمروهای بزرگ (نواحی با نمونه‌های بزرگ) استفاده نمود. اما برآوردهای مستقیم یعنی برآوردهایی که تنها بر مبنای داده‌هایی که از نمونه اخذ شده از آن ناحیه استخراج می‌شوند قرار دارند دارای خطای بسیار بزرگی خواهند بود که ناشی از کوچک بودن بیش از اندازه حجم نمونه در آن نواحی است. دلیل این امر این است که در ابتدا دقت مورد نظر برای برآورد پارامتر کل جامعه در نظر گرفته شده و حجم نمونه اصلی بر مبنای آن تعیین گردیده است. یعنی در هنگام برنامه‌ریزی و تعیین حجم نمونه فقط حجم نمونه کل مورد نیاز محاسبه شده است و طبعاً سهم هر ناحیه کوچک از کل حجم نمونه ناچیز است، در نتیجه برای نواحی کوچک این برآوردها دقت کافی را نخواهند داشت.

واضح است که از هر ناحیه نیز نمی‌توان نمونه‌ای بزرگ اختیار کرد. لذا ناچار باید به دنبال راهی بود که دقت برآوردها را تا آنجا که ممکن است بالا برد. برای این منظور می‌توان از داده‌های نواحی دیگر و اطلاعاتی که در آن نواحی براساس نمونه به دست آمده است در جهت بهبود بخشیدن به برآوردهای پارامترهای نواحی کوچک استفاده نمود. به این طریق حجم نمونه مؤثر افزایش یافته و دقت برآوردها افزایش می‌یابد. در مقاله حاضر روش درست‌نمایی وزنی را به این منظور مورد بحث قرار می‌دهیم.

۲ روش درست‌نمایی وزنی

روش درست‌نمایی بدون شک یکی از قدرتمندترین ابزارها در مسائل برآوردیابی و آزمون فرض می‌باشد؛ اما زمانی که در برآوردیابی همزمان با مسأله ترکیب اطلاعات روبرو هستیم روش درست‌نمایی کلاسیک کارآمد نخواهد بود. اهمیت این مسأله در مقاله استاین (۱۹۵۶) دیده شد. در واقع مقاله او نشان داد که در برآوردیابی همزمان میانگین‌های جوامع نرمال برآوردگر درست‌نمایی ماکزیمم یعنی میانگین نمونه غیر مجاز می‌باشد. او نشان داد که برای برآورد میانگین‌های جوامع نرمالی که از یک جهت شبیه یکدیگرند می‌توان برای به دست آوردن برآوردگری که دقت بیشتری از میانگین‌های نمونه دارد از سایر نمونه‌های مرتبط کمک گرفت.

مسأله ترکیب اطلاعات از منابع گوناگون خارج از روش درست‌نمایی قرار گرفته بود تا زمانی که تیبشیرانی و هاستی (۱۹۸۷) روش درست‌نمایی محلی^۵ را در محدوده رگرسیون ناپارامتری ارائه نمودند. هو روش درست‌نمایی محلی را خارج از حوزه رگرسیون ناپارامتری مورد بررسی قرار داد که این روش درست‌نمایی وزنی مرتبط^۶ (REWL) نامیده شد. علاوه بر این هو و زیدک (۲۰۰۰)

5) Local Likelihood 6) Relevance Weighted Likelihood

نشان دادند که چطور می‌توان برآوردگر معروف جیمز-استاین را به عنوان برآوردگر REWL استخراج نمود مشروط بر این که وزنها از نمونه برآورد شوند. بنابراین می‌توان روش درست‌نمایی را برای ترکیب اطلاعات نمونه‌های جوامع مختلف توسعه داد.

روش درست‌نمایی وزنی با روش درست‌نمایی کلاسیک متفاوت است. در روش کلاسیک فرض بر این است که مشاهدات از یک جامعه همسان خارج شده‌اند. در عوض درست‌نمایی وزنی در استنباط پارامتری زمانی پیش می‌آید که علاوه بر نمونه مورد نظر از جامعه مورد مطالعه نمونه‌های مرتبط اما مستقل از جوامع دیگر نیز در دسترس هستند که این همان مسأله مطلوب در برآوردیابی برای ناحیه کوچک می‌باشد. روش درست‌نمایی وزنی با اختصاص وزنهایی به این نمونه‌ها بر اساس میزان ارتباطشان اطلاعات آنها را ترکیب می‌کند.

در این مقاله به تولید برآوردگرهای همزمان میانگین‌های چندین توزیع می‌پردازیم که هر یک از این توزیعها متعلق به یک خانواده نمایی طبیعی با توابع واریانس درجه دوم (NEF-QVF)^۷ هستند.

۳ روش درست‌نمایی وزنی در خانواده نمایی

فرض کنید m جامعه داریم. بردارهای نمونه‌های مستقل با $(y_{i1}, \dots, y_{in_i})'$ $i = 1, \dots, m$ نشان داده می‌شوند که $n_i \geq 0$. در عمل همیشه این امکان وجود دارد که از یک ناحیه محلی هیچ مشاهده‌ای گرفته نشود. می‌دانیم Y دارای توزیع نمایی با پارامتر θ است اگر

$$f(y, \theta) = \exp[\theta y - \phi(\theta) + h(y)]$$

بنابراین اگر $\theta = (\theta_1, \dots, \theta_m)'$ و ω_{ij}^* وزن اختصاص یافته به نمونه j ام در برآورد پارامتر i ام باشد تابع درست‌نمایی به صورت زیر تعریف می‌شود:

$$\begin{aligned} WL(\theta) &= WL(\theta_1, \dots, \theta_m) \\ &= \prod_{i=1}^m \prod_{j=1}^m \prod_{l=1}^{n_j} \exp[(\theta_i y_{jl}) - \phi(\theta_i) + h(y_{jl})]^{\omega_{ij}^* n_j^{-1}} \end{aligned}$$

با فرض

$$\bar{y}_j = \frac{1}{n_j} \sum_{l=1}^{n_j} y_{jl} \quad j = 1, \dots, m$$

7) Natural Exponential Family with Quadratic Variance Function

داریم:

$$WL(\theta) = \exp\left[\sum_{i=1}^m [\theta_i \sum_{j=1}^m \omega_{ij}^* \bar{y}_j - \phi(\theta_i) \sum_{j=1}^m \omega_{ij}^* + \sum_{j=1}^m \omega_{ij}^* n_j^{-1} \sum_{l=1}^{n_j} h(y_{jl})]\right] \quad (1)$$

بنابراین اگر ω_{ij}^* ها معلوم باشند

$$\frac{\partial \ln WL(\theta)}{\partial \theta_i} = \sum_{j=1}^m \omega_{ij}^* \bar{y}_j - \phi'(\theta_i) \sum_{j=1}^m \omega_{ij}^* = 0$$

$$\frac{\partial^2 \ln WL(\theta)}{\partial \theta_i^2} = -\phi''(\theta_i) \sum_{j=1}^m \omega_{ij}^* \leq 0$$

برآوردگرهای درست‌نمایی ماکزیمم وزنی (MWLE's) برای میانگینهای جامعه یعنی $\phi_i = \phi'(\theta_i)$ عبارتند از:

$$\hat{\phi}_{i_{WL}} = \frac{\sum_{j=1}^m \omega_{ij}^* \bar{y}_j}{\sum_{j=1}^m \omega_{ij}^*} = \sum_{j=1}^m \omega_{ij} \bar{y}_j \quad (2)$$

که

$$\omega_{ij} = \frac{\omega_{ij}^*}{\sum_{j=1}^m \omega_{ij}^*}$$

اما در عمل چون وزنها معلوم نیستند باید از روی داده‌ها برآورد شوند.

۱.۳ روش درست‌نمایی وزنی در مدل تعویض پذیر

در این حالت فرض کنید

$$\omega_{ii} = \omega \quad i = 1, \dots, m \quad \omega_{ij} = \omega^* \quad j \neq i = 1, \dots, m$$

چنین فرضی این ایده را در بردارد که نمونه‌های سایر جوامع نسبت به اطلاعاتی که در مورد i امین پارامتر دارند تعویض پذیر^۸ هستند. این فرض صریحاً به برآوردگر جیمز استاین در حالت نرمال زمانی که وزنها با استفاده از داده‌ها برآورد می‌شوند منجر می‌شود. تحت این فرض

$$\hat{\phi}_{i_{WL}} = \omega \bar{y}_i + \omega^* \sum_{j \neq i} \bar{y}_j$$

8) Exchangable

و چون مجموع وزنها برابر یک می باشد پس

$$\omega + (m - 1)\omega = 1 \Rightarrow \omega^* = \frac{1 - \omega}{m - 1}$$

و

$$\hat{\phi}_{iWL} = \omega \bar{y}_i + \frac{1 - \omega}{m - 1} \sum_{j \neq i} \bar{y}_j \quad (3)$$

برای پیدا کردن ω بهینه باید $E(\hat{\phi}_i - \phi_i)^2$ نسبت به ω می نیمم شود. با فرض

$$u_i = V(\bar{y}_i) = \frac{1}{n_i} \phi''(\theta_i) \quad \bar{u} = \frac{1}{m} \sum_{i=1}^m u_i$$

$$s_\phi^2 = \frac{1}{m - 1} \sum_{i=1}^m (\phi_i - \bar{\phi})^2$$

داریم:

$$g(\omega) = m[\omega^2 \bar{u} + (1 - \omega)^2 (m - 1)^{-1} \bar{u} + m(1 - \omega)^2 (m - 1)^{-1} s_\phi^2] \quad (4)$$

$$\Rightarrow \frac{\partial g(\omega)}{\partial \omega} = m[2\omega \bar{u} - 2(1 - \omega)(m - 1)^{-1} \bar{u} - 2m(1 - \omega)(m - 1)^{-1} s_\phi^2] = 0$$

و لذا

$$\omega_{opt} = (\bar{u} + m s_\phi^2) [m(\bar{u} + s_\phi^2)]^{-1} \quad (5)$$

با جایگذاری (۵) در رابطه (۳) داریم:

$$\hat{\phi}_{iWL}^{opt} = (1 - \beta_{opt}) \bar{y}_i + \beta_{opt} \bar{y} \quad (6)$$

که

$$\bar{y} = \frac{1}{m} \sum_{i=1}^m \bar{y}_i \quad \beta_{opt} = \frac{\bar{u}}{\bar{u} + s_\phi^2}$$

باز هم \bar{u} و s_ϕ^2 نامعلومند و باید برآورد شوند. فرض کنید:

$$s_y^2 = \frac{1}{m - 1} \sum_{i=1}^m (y_i - \bar{y})^2$$

آنگاه $E(s_y^2) = \bar{u} + s_\phi^2$. اما چون جوامع دارای توزیع نمایی با توابع واریانس درجه دوم هستند لذا

$$V(y_{jl}) = \nu_0 + \nu_1 \phi_j + \nu_2 \phi_j^2 = Q(\phi_j) \quad l = 1, \dots, n_j, \quad j = 1, \dots, m$$

$$E[Q(\bar{y}_j)] = \nu_0 + \nu_1 \phi_j + \nu_2 [V(\bar{y}_j) + \phi_j^2]$$

$$\Rightarrow E[Q(\bar{y}_j)] = Q(\phi_j) + \nu_2 n_j^{-1} Q(\phi_j) = (n_j + \nu_2) u_j \quad (7)$$

پس

$$E\left[\frac{1}{m} \sum_{j=1}^m (n_j + \nu_2)^{-1} Q(\bar{y}_j)\right] = u$$

فرض کنید

$$\hat{u} = \frac{1}{m} \sum_{j=1}^m (n_j + \nu_2)^{-1} Q(\bar{y}_j)$$

بنابراین β_{opt} را به وسیله $\hat{\beta}$ برآورد می‌کنیم که عبارتست از:

$$\hat{\beta} = \min\left\{\frac{\hat{u} + 1/m}{s_y^2 + 1/m}, 1\right\} \quad (8)$$

اکنون ϕ_i به صورت زیر برآورد می‌شود:

$$\hat{\phi}_{iWL} = (1 - \hat{\beta}) \bar{y}_i + \hat{\beta} \bar{y} \quad i = 1, \dots, m \quad (9)$$

تذکره: حالت خاصی را در نظر بگیرید که نمونه‌هایی با حجم مساوی n به طور مستقل از توزیعهای $N(\theta_i, \sigma^2)$ $i = 1, \dots, m$ استخراج شده‌اند. و $\sigma^2 > 0$ معلوم است. در این صورت $\nu_0 = \sigma^2, \nu_1 = \nu_2 = 0$ پس $\hat{u} = \frac{\sigma^2}{n}$. بنابراین با جایگذاری در (۸) داریم:

$$\hat{\phi}_{iWL} = \left(1 - \frac{\sigma^2/n}{s_y^2}\right) \bar{y}_i + \frac{\sigma^2/n}{s_y^2} \bar{y} \quad i = 1, \dots, m$$

۲.۳ روش درست‌نمایی وزنی در مدل تعویض پذیر محلی

فرض کنید $\omega_{ii} = \omega_i, \omega_{ij} = \omega_{i*} \quad j \neq i = 1, \dots, m$ یعنی فرض تعویض پذیری سراسری بر داشته شده اما فرض تعویض پذیری درون هر ناحیه محلی باقی می‌ماند. آنگاه

$$\hat{\phi}_{iWL} = \omega_i \bar{y}_i + \omega_{i*} \sum_{j \neq i} \bar{y}_j$$

و چون $1 = \omega_i + (1 - m)\omega_{i^*}$ داریم:

$$\hat{\phi}_{iWL} = \omega_i \bar{y}_i + (1 - \omega_i)(1 - m)^{-1} \sum_{j \neq i} \bar{y}_j$$

برای به دست آوردن ω_i های بهینه اکنون $g(\omega_i) = E(\hat{\phi}_{iWL} - \phi_i)^2$ را نسبت به ω_i به ازای هر $i = 1, \dots, m$ می‌نیمیم.

$$g(\omega_i) = \omega_i^2 u_i + (1 - \omega_i)^2 (m - 1)^{-1} (m\bar{u} - u_i) + m^2 (m - 1)^{-2} (1 - \omega_i)^2 (\phi_i - \bar{\phi})^2$$

$$g'(\omega_i) = 2\omega_i u_i - 2(1 - \omega_i)(m - 1)^{-1} (m\bar{u} - u_i) - 2m^2 (m - 1)^{-2} (1 - \omega_i)(\phi_i - \bar{\phi})^2$$

$$g''(\omega_i) = 2u_i + \frac{2}{m - 1} (m\bar{u} - u_i) + \frac{2m^2}{(m - 1)^2} (\phi_i - \bar{\phi})^2 > 0$$

پس

$$\omega_{i_{opt}} = \frac{m\bar{u} - u_i + m^2 (\phi_i - \bar{\phi})^2}{(m - 1)^2 u_i + m\bar{u} - u_i + m^2 (\phi_i - \bar{\phi})^2}$$

بنابراین

$$\hat{\phi}_{iWL}^{opt} = (1 - \beta_{i_{opt}}) \bar{y}_i + \beta_{i_{opt}} \bar{y} \quad (10)$$

که

$$\beta_{i_{opt}} = \frac{(m - 1)u_i}{(m - 2)u_i + \bar{u} + m(\phi_i - \bar{\phi})^2}$$

اما در عمل β_i ها معلوم نیستند و باید برآورد شوند. توجه کنید که

$$\begin{aligned} E(\bar{y}_i - \bar{y})^2 &= V(\bar{y}_i - \bar{y}) + (\phi_i - \bar{\phi})^2 \\ &= \frac{1}{m} [(m - 2)u_i + \bar{u} + m(\phi_i - \bar{\phi})^2] \end{aligned}$$

همچنین

$$E[Q(\bar{y}_i)] = \nu_0 + \nu_1 \phi_i + \nu_2 (n_i + \nu_2) = (n_i + \nu_2) u_i$$

بنابراین $E[Q(\bar{y}_i)] / (n_i + \nu_2) = u_i$. از این رو $\beta_{i_{opt}}$ به وسیله $\hat{\beta}_i$ برآورد می‌شود که

$$\hat{\beta}_i = \max \left[\frac{m - 1}{m} \frac{Q(\bar{y}_i)(n_i + \nu_2)^{-1} + m^{-1}}{(\bar{y}_i - \bar{y})^2 + m^{-1}}, 1 \right]$$

و لذا

$$\hat{\phi}_{iWL} = (1 - \hat{\beta}_i) \bar{y}_i + \hat{\beta}_i \bar{y}$$

۳.۳ روش درست‌نمایی وزنی بدون فرض تعویض پذیری

حالتی را در نظر بگیرید که هیچگونه فرض تعویض پذیری وجود ندارد. اگر $\hat{\phi}_i = \sum_{j=1}^m \omega_{ij} \bar{y}_j$ پس ω_{ij} ‌های بهینه با می‌نیمم کردن $E(\sum_{j=1}^m \omega_{ij} \bar{y}_j - \phi_i)^2$ نسبت به ω_{ij} به دست می‌آیند با این شرط که $\sum_{j=1}^m \omega_{ij} = 1$. فرض کنید $\Omega_i = (\omega_{i1}, \dots, \omega_{im})'$ و فرض کنید

$$g(\Omega_i) = E\left(\sum_{j=1}^m \omega_{ij} \bar{y}_j - \phi_i\right)^2 - 2\lambda_i \left(\sum_{j=1}^m \omega_{ij} - 1\right)$$

که λ_i ‌ها ضرایب لاگرانژ می‌باشند.

$$g(\Omega_i) = \sum_{j=1}^m \omega_{ij}^2 u_j + \left(\sum_{j=1}^m \omega_{ij} \phi_j - \phi_i\right)^2 - 2\lambda_i \left(\sum_{j=1}^m \omega_{ij} - 1\right)$$

بنابراین

$$\frac{\partial g}{\partial \omega_{ij}} = 2\omega_{ij} u_j + 2\left(\sum_{j=1}^m \omega_{ij} \phi_j - \phi_i\right) \phi_j - 2\lambda_i$$

$$\frac{\partial^2 g}{\partial \omega_{ij}^2} = 2u_j + 2\phi_j^2 > 0 \quad \frac{\partial^2 g}{\partial \omega_{ij} \partial \omega_{kl}} = 0 \quad j \neq l$$

از این رو ماتریس همسایان برای g معین مثبت است و ω_{ij} ‌های بهینه با حل معادلات $\frac{\partial g}{\partial \omega_{ij}} = 0 \quad i, j = 1, \dots, m$ به دست می‌آیند. با استفاده از نماد ماتریسی این معادلات به صورت زیر نوشته می‌شوند:

$$M\Omega_i = \lambda_i \mathbf{1}_m + \phi_i \Phi$$

که

$$M = D + \Phi\Phi', \quad D = \text{Diag}(u_1, \dots, u_m), \quad \Phi = (\phi_1, \dots, \phi_m)'$$

بنابراین Ω_i بهینه به صورت زیر به دست می‌آید.

$$\Omega_i^{opt} = M^{-1}(\lambda_i \mathbf{1}_m + \phi_i \Phi) \quad (11)$$

چون $\sum_{j=1}^m \omega_{ij} = 1$ بنابراین

$$1 = \mathbf{1}_m' \Omega_i^{opt} = \lambda_i \mathbf{1}_m' M^{-1} \mathbf{1}_m + \phi_i \mathbf{1}_m' M^{-1} \Phi$$

$$\lambda_i = \frac{1 - \phi_i \lambda'_m M^{-1} \Phi}{\lambda'_m M^{-1} \lambda_m} \quad (12)$$

اکنون با استفاده از (۱۰) و (۱۱) داریم:

$$\Omega_i^{opt} = \frac{1 - \phi_i \lambda'_m M^{-1} \Phi}{\lambda'_m M^{-1} \lambda_m} (M^{-1} \lambda_m) + \phi_i M^{-1} \Phi$$

و

$$M^{-1} = D^{-1} - \frac{D^{-1} \Phi \Phi' D^{-1}}{1 + \Phi' D^{-1} \Phi}$$

حال چون $E(\bar{y}_i) = \phi_j$, $E[(n_j + \nu_r)^{-1} Q(\bar{y}_j)] = u_j$

$$Q(\bar{y}_j) = \nu_0 + \nu_1 \bar{y}_j + \nu_2 \bar{y}_j^2$$

داریم:

$$\hat{u}_j = (n_j + \nu_r)^{-1} Q(\bar{y}_j)$$

$$\hat{\Omega}_i^{opt} = \frac{1 - \bar{y}_i \lambda'_m \hat{M}^{-1} Z}{\lambda'_m \hat{M}^{-1} \lambda_m} (\hat{M}^{-1} \lambda_m) + \bar{y}_i \hat{M}^{-1} Z$$

که

$$Z = (\bar{y}_1, \dots, \bar{y}_m)' , \quad \hat{M} = (\hat{D} + ZZ')^{-1} , \quad \hat{D} = \text{Diag}(\hat{u}_1, \dots, \hat{u}_m)$$

مراجع

- [1] Ghosh, M., and Rao, J.N.K. (1994), Small area estimation: an appraisal (with discussion), *Statistical Science*, 9, 65-93.
- [2] Hu, F. and Zidek, J.V. (2000), The relevance weighted likelihood with applications. In *Empirical Bayes and Likelihood Inference. Lecture Notes in Statistics*. Eds. S.E. Ahmed and N. Reid. Springer-Verlag, New York, pp 211-235.
- [3] Stein, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *proceedings of the Third Berkeley Symposium in Mathematical Statistics and Probability*, V1. University of California Press, pp 197-206.
- [4] Tibshirani. R. and Hastie, T. (1987). Local likelihood estimation. *J. Amer. Statist. Assoc.*, 82, 559-567.

رگرسیون چندک

فرهاد فتاحی^۱، عباس گرامی^۲

^۱ دانشگاه تربیت مدرس

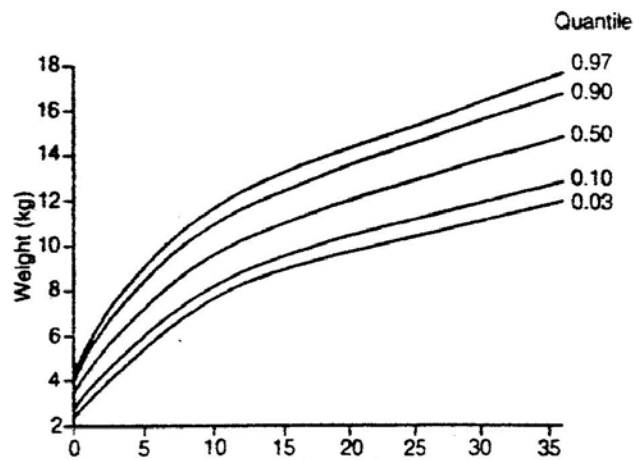
^۲ پژوهشکده آمار

چکیده: رگرسیون چندک را می‌توان بعنوان تعمیم طبیعی برآورد حداقل مربعات کلاسیک مدل‌های میانگین شرطی در برآورد مجموعه‌ای از مدل‌ها، برای چندک شرطی مد نظر قرار داد. مورد خاص و عمده برآوردگر رگرسیون میانه می‌باشد که مجموع قدرمطلق خطاها را مینیمم می‌نماید. سایر توابع چندک شرطی با مینیمم کردن یک مجموع موزون نامتقارن از قدر مطلق خطاها برآورد می‌گردد. با کنار هم قرار گرفتن کلیه توابع برآورد شده چندک شرطی می‌توان دید بسیار کامل‌تری در مورد اثر متغیرهای کمکی بر روی پارامترهای موقعیت، مقیاس و شکل توزیع متغیر پاسخ بدست آورد.

واژه‌های کلیدی: رگرسیون چندک، تابع چندک تجربی، چندک رگرسیونی، چندک نمونه‌ای، چندک شرطی

۱ مقدمه

در شکل (۱) چندکهای مختلف وزن نوزادان پسر در مقابل سن‌شان رسم شده است. مسئله مورد علاقه ممکن است فقط بر ارتباط متوسط این دو متغیر متمرکز نشده باشد، بلکه ممکن است بر ارتباط بین چندکهای آخری و سن مبتنی باشد. ممکن است چندکها بعنوان کرانه‌های مناسب مشکلات تغذیه‌ای بالقوه نوزادان را منعکس سازند بنابراین وقتی که فرایند افزایش وزن را در طول زمان مورد آزمایش قرار می‌دهیم بسیار سودمند واقع می‌گردد. واضح است که در این حالت و بسیاری حالات مشابه، ممکن است یک مدل رگرسیون استاندارد برای پیش‌بینی ارتباط بین کلیه چندکهای یک متغیر وابسته Y از یک سو و بردار متغیرهای مستقل X از سوی دیگر مناسب نباشد. به همین دلیل کانکر، باست و هاگ روش رگرسیون (MAD) یعنی میانگین قدر مطلق انحرافات را برای رگرسیون چندک تعمیم دادند که در آن تحت فرض خطی بودن برآوردگرها، روشی را برای مدل‌بندی چندکهای تابع توزیع شرطی $F(y|X)$ بیان کردند. در تشریح این روش ضرورت دارد چگونگی ارتباط بین چندکهای معمولی به عنوان جوابی برای معادلات برآورد مشخص شود که بواسطه آن یک تابع زیان نمونه‌ای خاص مینیمم می‌گردد. این امر ما را در بدست آوردن برآورد مدل‌های رگرسیونی برای چندکهای شرطی قادر می‌سازد. با در دست داشتن نمونه‌ای n تایی از



شکل ۱: نمودار سن و وزن

مشاهدات یک متغیره $\{y_1, y_2, \dots, y_n\}$ ، p امین چندک نمونه‌ای $\hat{\mu}_p$ ، $0 < p < 1$ باید در رابطه زیر صدق کند.

$$\sum_{i=1}^n [pI(y_i \geq \hat{\mu}_p) - (1-p)I(y_i < \hat{\mu}_p)] = 0 \quad (1)$$

در رابطه بالا I تابع نشانگر بوده که در صورتی که شرط درون پرانتز برقرار باشد مقدار یک می‌گیرد در غیر اینصورت مقدار صفر را اختیار می‌کند. توجه کنید که $\hat{\mu}_p$ برآوردی از چندک جامعه متناظر خود یعنی μ_p می‌باشد که در رابطه زیر صدق می‌کند:

$$p[1 - F(\mu_p)] - (1-p)F(\mu_p) = 0$$

رابطه (۱) را می‌توان بصورت یک معادله برآوردکننده نیز نشان داد که به باقیمانده‌های مثبت $r_i = y_i - \hat{\mu}_p$ وزن p و به باقیمانده‌های منفی وزن $1-p$ داده می‌شود و مقدار مورد نظر توسط علامت باقیمانده اندازه‌گیری می‌شود بنابراین رابطه (۱) را می‌توان به فرم $\sum_{i=1}^n \psi_p(r_i) = 0$ نشان داد. که در آن

$$\psi_p(r) = \begin{cases} p\psi(r) & r > 0 \\ (1-p)\psi(r) & o.w \end{cases} \quad (2)$$

گاهی اوقات $\psi(r)$ را بعنوان تابع تأثیر در برآورد چندکهای رگرسیونی با تابع زیان نامتقارن متناظر با آن که به فرم زیر است بکار می‌برند.

$$\rho_p(r) = \begin{cases} p\rho(r) & r > 0 \\ (1-p)\rho(r) & o.w \end{cases} \quad (3)$$

که در آن $\rho(r) = |r|$. رابطه برآوردکننده فوق و به تبع آن $\hat{\mu}_p$ ، را می‌توان با مینیمم کردن متوسط تابع زیان نمونه‌ای زیر بدست آورد.

$$n^{-1} \sum_{i=1}^n \rho_p(y_i - \mu_p) \quad (4)$$

برای اینکه روش بالا را به رگرسیون چندک تعمیم دهیم، روشی همانند رگرسیون کمترین مربعات خطی را در پیش می‌گیریم و μ_p موجود در رابطه (۴) را به $\mu_p(X) = \beta x' \beta_p$ تبدیل می‌کنیم، که در آن β_p یک بردار ضرائب رگرسیونی برای p امین چندک شرطی است. می‌توان برآورد β_p را با حل رابطه زیر بدست آورد.

$$\sum_{i=1}^n \psi_p(y_i - \underline{x}'_i \hat{\beta}_p) \underline{x}_i = 0 \quad (5)$$

در حالت خاص $p = 0.5$ ، سطح برازش داده شده معادل با میانه رگرسیون y نسبت به X است که رابطه زیر را مینیمم می‌کند (که عبارتست از متوسط قدرمطلق انحرافات باقیمانده‌ها).

$$n^{-1} \sum_{i=1}^n |y_i - X'_i \hat{\beta}_{0.5}| \quad (6)$$

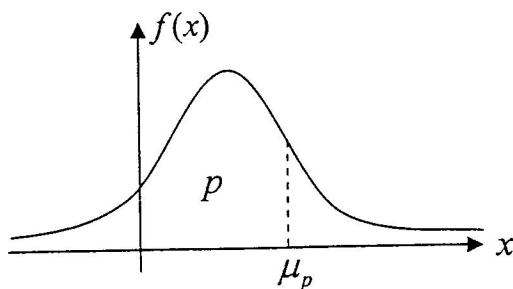
۲ معرفی چندک

فرض کنید متغیر تصادفی X ، دارای تابع توزیع تجمعی $F(x)$ باشد پارامتر μ_p را چندک مرتبه p برای $F(x)$ یا برای X ، می‌نامیم هرگاه نامساوی دو طرفه زیر را داشته باشیم.

$$p(X < \mu_p) \leq p \leq p(X \leq \mu_p) \quad 0 < p < 1$$

معنی این نامساوی دو طرفه این است که مقدار احتمال در فاصله باز $(-\infty, \mu_p)$ حداکثر p و در فاصله نیم باز $(-\infty, \mu_p]$ حداقل p می‌باشد (شکل ۲).

μ_p را که در آن $0 < p < 1$ است، چندک مرتبه p می‌نامند هرگاه تقریباً $100p$ درصد داده‌ها کوچکتر از آن باشند.



شکل ۲: نمودار چنک p

چندکهاکلی تر از میانه می‌باشند و در حقیقت $p = 0.5$ همان میانه است به زبان هندسی اگر از نقطه μ_p خطی به موازات محور y ها رسم کنیم مساحت زیر منحنی فراوانی که در سمت چپ این خط قرار دارد برابر p می‌باشد. چندکهای معروف عبارتند از:

الف) چارکها که به ازای $p = 0.25, 0.5, 0.75$ بدست می‌آیند و آنها را به ترتیب با $\mu_{0.25}, \mu_{0.5}, \mu_{0.75}$ نشان می‌دهند.

ب) دهکها که به ازای $p = 0.1, \dots, 0.9$ بدست می‌آیند و آنها را به ترتیب با D_1, \dots, D_9 نشان می‌دهند.

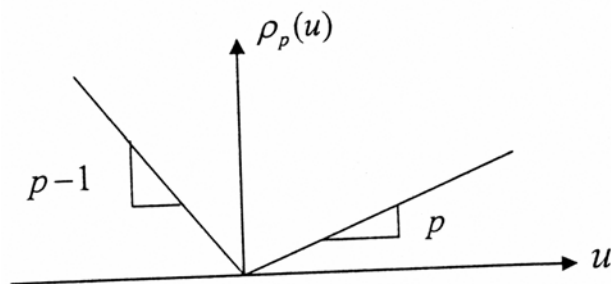
ج) صدکها که به ازای $p = 0.01, \dots, 0.99$ بدست می‌آیند و آنها را به ترتیب با p_1, \dots, p_{99} نشان می‌دهند.

هر متغیر تصادفی را می‌توان توسط تابع توزیعش $F(x) = P(X \leq x)$ مشخص کرد و برای هر $0 < p < 1$ چندک p ام متغیر X برابر $F^{-1}(p) = \inf\{x : F(x) \geq p\}$ که در حالت $p = \frac{1}{4}$ میانه $F^{-1}(\frac{1}{4})$ را خواهیم داشت.

چندکها از یک مسئله بهینه‌سازی ساده ناشی می‌شوند به این صورت که اگر به ازای هر $p \in (0, 1)$ تابع زیان را بصورت $\rho_p(u) = u(p - I(u < 0))$ تعریف کنیم شکل این تابع زیان بصورت (شکل ۳) خواهد بود.

در اینجا هدف یافتن برآورد نقطه‌ای x یعنی \hat{x} است، بطوریکه با انتخاب \hat{x} ، امید زیان مینیمم شود بنابراین اگر امید زیان بصورت زیر تعریف شود:

$$M = E\rho_p(x - \hat{x}) = (p - 1) \int_{-\infty}^{\hat{x}} (x - \hat{x})dF(x) + p \int_{\hat{x}}^{\infty} (x - \hat{x})dF(x)$$



شکل ۳: شکل تابع زیان

با مشتق‌گیری نسبت به \hat{x} داریم:

$$\frac{dM}{d\hat{x}} = (\lambda - p) \int_{-\infty}^{\hat{x}} dF(x) - p \int_{\hat{x}}^{\infty} dF(x) = F(\hat{x}) - p$$

چون تابع F یکنواست بنابراین هر عضو مجموعه $\{x : F(x) = p\}$ امید زیان را مینیمم می‌کند. در صورتیکه $\hat{x} = F^{-1}(p)$ یکتا باشد جواب مورد نظر بدست می‌آید. در غیر اینصورت فاصله‌ای از چندکهای p ام را داریم و با توافق به این قرار داد که تابع چندک تجربی از طرف چپ پیوسته است. پس کوچکترین عضو این مجموعه را به عنوان جواب ممکن انتخاب می‌کنیم. پیدا کردن بهترین برآوردگر نقطه‌ای برای زیان خطی نامتقارن، منجر به مفهوم چندکها می‌شود. در صورت متقارن بودن زیان، قدر مطلق آن، میانه را به خوبی مشخص خواهد کرد. برای مثال اگر مقادیر کوچکتر از برآورد، سه برابر مقادیر بزرگتر از برآورد باشد، \hat{x} را به گونه‌ای انتخاب می‌کنیم که مقدار احتمال $P(X \leq x)$ سه برابر مقدار احتمال $P(X > x)$ باشد یعنی \hat{x} را صدک ۷۵ام توزیع F انتخاب می‌کنیم. اکنون اگر تابع توزیع تجربی،

$$F_n(x) = n^{-1} \sum_{i=1}^n I(X_i \leq x)$$

جایگزین F شود می‌توان \hat{x} را با مینیمم کردن امید زیان،

$$\int (x - \hat{x}) dF_n(x) = n^{-1} \sum_{i=1}^n \rho_p(x_i - \hat{x})$$

بدست آورد در اینجا وقتی pn یک عدد صحیح باشد، بازه‌ای از جوابهایی را که به صورت $\{x : F_n(x) = p\}$ هستند، خواهیم داشت بنابراین جواب، دقیق و روشن نخواهد بود و خواهیم دید که این نتیجه زیاد سودمندی هم نیست. اکنون حالت ساده‌ای از چندکهای نمونه‌ای مرسوم را با یک مقدار جزئیات بیشتر شرح می‌دهیم. می‌توان مسئله پیدا کردن چندک نمونه‌ای p ام را بصورت:

$$\min_{\mu \in R} \sum_{i=1}^n \rho_p(y_i - \mu) \quad (۷)$$

که در آن $\rho_p(u) = u(p - I(u < \circ))$ است بیان کرد که ممکن است بصورت زیر بعنوان یک مسئله برنامه‌ریزی خطی با معرفی $2n$ متغیر غیر واقعی u_i و v_i ، $\{u_i, v_i : i = 1, \dots, n\}$ ، که u_i و v_i قسمتهای مثبت و منفی بردار باقیمانده‌ها را نشان می‌دهند فرموله شود.

$$\min_{(\mu, u, v) \in R \times R^{+2n}} \{p \sum_{i=1}^n u_i + (1-p) \sum_{i=1}^n v_i \mid I_n p + u - v = y\} \quad (۸)$$

$\sum_{i=1}^n$ برداری n بعدی از یک‌ها می‌باشد. بوضوح می‌بینیم که در (۸) یک تابع خطی روی یک مجموعه محدودیت چند وجهی مینیمم می‌شود. این مجموعه محدودیت چند وجهی شامل اشتراک ابر صفحه‌های $(2n+1)$ بعدی تولید شده، با شرایط تساوی و مجموعه $R \times R^{2n}$ است. در تابع هدف (۷) در صورتی جواب بهینه خواهیم داشت که مشتقات راست و چپ $R(\mu) = \sum_{i=1}^n \rho_p(y_i - \mu)$ که در زیر بدست آمده‌اند هر دو نامنفی باشند.

$$\dot{R}(\mu, +1) = \lim_{h \rightarrow \circ} \frac{(R(\mu + h) - R(\mu))}{h} = \sum_{i=1}^n I((y_i < \mu + h) - \mu) \quad (۹)$$

$$\dot{R}(\mu, -1) = \lim_{h \rightarrow \circ} \frac{(R(\mu - h) - R(\mu))}{-h} = \sum_{i=1}^n (\mu - I((y_i < \mu - h))) \quad (۱۰)$$

یعنی اگر عدد np در فاصله بسته $[N^-, N^+]$ قرار داشته باشد $\{y_i < \mu \pm \circ\} = N^\pm$ در اینصورت $\hat{\mu}_p$ بین دو آماره ترتیبی مجاور هم قرار می‌گیرد و وقتی مقدار این آماره‌های ترتیبی مجاور برابر شود در اینصورت $\hat{\mu}_p$ یکتا است.

۳ چندکهای نمونه‌ای

اگر بخواهیم اطلاعات بیشتری را در مورد یک مجموعه از مشاهدات ارائه دهیم، ساده‌ترین کار، محاسبه چندکهای نمونه‌ای است. در اینجا هدف معرفی ابتدائی چندکهای نمونه‌ای است که همانا در مرتب شدن مشاهدات نقش بسزایی دارند. اگر $\{Y_t : t = 1, \dots, n\}$ یک نمونه

تصادفی از متغیر تصادفی Y با تابع توزیع $F(\circ)$ باشد آنگاه برای $0 < p < 1$ ، p امین چندک نمونه‌ای، عبارتست از هر جوابی که از مینیمم کردن عبارت زیر بدست می‌آید.

$$\min_{\mu_p \in R} \left[\sum_{t \in \{t: y_t \geq \mu\}} p(y_t - \mu_p) + \sum_{t \in \{t: y_t \leq \mu\}} (p - 1)(y_t - \mu_p) \right] \quad (11)$$

در حالت خاص $p = \frac{1}{2}$ ، میانه نمونه، 50° امین چندک نمونه‌ای است. حال این تعریف را به مدل‌های خطی تعمیم داده و چندکهای رگرسیونی را تعریف می‌کنیم.

۴ چندکهای رگرسیونی

رابطه ریاضی مدل رگرسیونی خطی:

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + u_i \quad x_{i1} = 1, i = 1, 2, \dots, n$$

را در نظر می‌گیریم که u_i ها دارای تابع توزیع $F(\circ)$ باشند. فرم ماتریسی رابطه عبارتست از:

$$y = X\beta + u$$

که $u = (u_1, \dots, u_n)'$ بردار خطا، $\beta = (\beta_1, \dots, \beta_p)'$ بردار پارامترهای مجهول، X یک ماتریس $n \times p$ معلوم رتبه کامل ستونی و $y = (y_1, \dots, y_n)'$ بردار مشاهدات متغیر وابسته است. اگر $\{X_t, t = 1, \dots, n\}$ نشان دهنده دنباله‌ای از بردارها از سطرهاى ماتریس X و $\{y_t, t = 1, \dots, n\}$ نشان دهنده اعداد تصادفی از فرایند رگرسیونی $u_t = y_t - x_t'\beta$ با تابع توزیع $F(\circ)$ باشد آنگاه برای $0 < p < 1$ ، p امین چندک رگرسیونی بصورت،

$$\min_{\beta \in R^p} \left[\sum_{t \in \{t: y_t \geq x_t'\beta\}} p(y_t - x_t'\beta) + \sum_{t \in \{t: y_t < x_t'\beta\}} (p - 1)(y_t - x_t'\beta) \right] \quad (12)$$

بدست می‌آید. که آن را با نماد $\beta(p)$ نشان می‌دهیم. در اینجا ساختار عبارت (۱۲) را بر اساس تعریف چندکهای نمونه‌ای (۱۱) بیان کردیم. ولی در عمل برای ساده‌تر نشان دادن عبارت (۱۲) ابتدا تابع زیان $\rho_p(\circ)$ را بصورت $\rho_p(u) = u(p - I(u < \circ))$ تعریف می‌کنیم. سپس $\beta(p)$ را از مینیمم کردن تابع زیر روی β در فضای R^p برای مقادیر ثابت p ، $0 < p < 1$ بدست می‌آوریم.

$$\min_{\beta \in R^p} \sum_{i=1}^n \rho_p(y_i - x_i'\beta) \quad (13)$$

در عمل عبارت (۱۳) با عبارت (۱۴) یکسان است. با در نظر گرفتن $\hat{\beta}(p)$ ، برآورد $\beta(p)$ ، برآورد تابع چندک شرطی p ، $\hat{\Phi}_Y(p|x) = x' \hat{\beta}(p)$ بدست می‌آید. می‌توان مسئله رگرسیون چندک (۱۳) را بصورت یک مسئله برنامه‌ریزی خطی مشابه برنامه‌ریزی خطی بخش قبل به صورت زیر فرموله کرد.

$$\min_{(\beta, u, v) \in R^p \times R^{+2n}} \{pI_n' u + (\lambda - p)I_n' v | X\beta + u - v = y\} \quad (14)$$

که X ماتریس طرح رگرسیون $n \times p$ است. در اینجا هم دوباره تابع خطی روی یک مجموعه محدودیت چند وجهی مینیمم می‌شود و بیشتر خصوصیات شناخته شده $\hat{\beta}(p)$ از خصوصیات شناخته شده جوابهای برنامه‌ریزی خطی پیروی می‌کند.

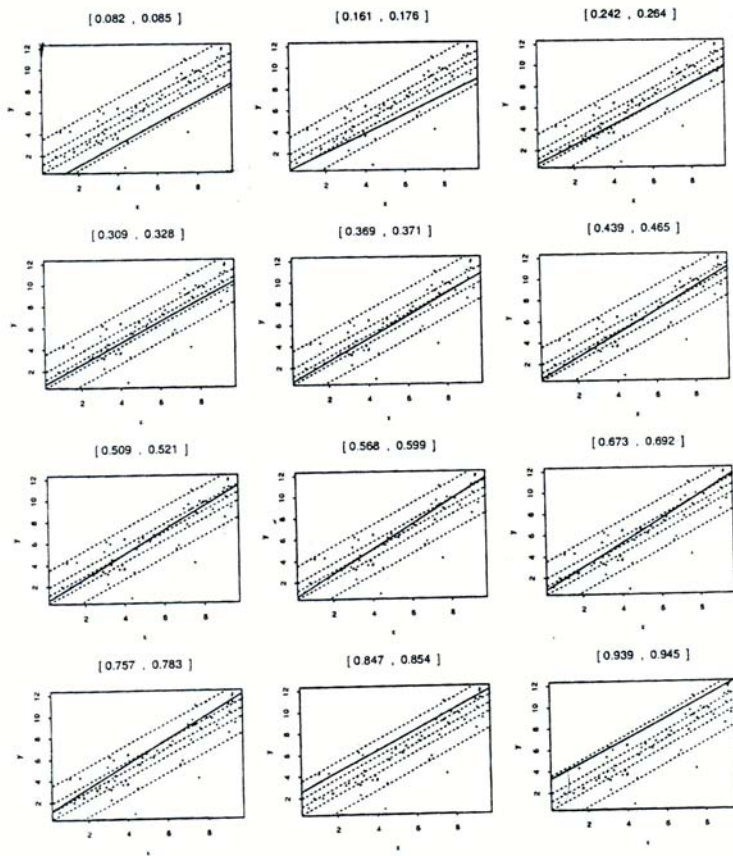
۵ مثال ۱

یک مدل رگرسیون بصورت، $y_i = \beta_0 + x_i' \beta_1 + u_i$ داریم که در آن خطاها هم توزیع مستقل از هم هستند. بنابراین تابع چندک y_i بصورت $\Phi(p|x) = \beta_0 + x' \beta_1 + F_u^{-1}(p)$ خواهد بود که F_u تابع توزیع مشترک خطاها است. $\hat{\beta}(p)$ برآورد پارامترهای جامعه می‌باشد. در شکل (۴)، 60° مشاهده تولید شده از توزیع نرمال با خطاهای مستقل و هم توزیع، به همراه خطوط چندک رگرسیون برازش شده متناظر با $\{0/0.5, 0/0.25, 0/0.5, 0/0.75, 0/0.95\}$ نشان داده شده است. خط ممتد در هر شکل، چندک شرطی برآورد شده را برای فاصله p (که در بالای هر شکل نشان داده شده است) شرح می‌دهد. می‌بینیم که با افزایش p ، این خط ممتد با شیبی تقریباً برابر با شیب توابع چندک شرطی، در بین داده‌ها تغییر می‌کند و به سمت خطوط چندک شرطی بالاتر می‌رود. در بالای هر شکل بازه‌ای برای هر p نشان داده شده است که پاسخ متناظر آن بهینه است.

۶ اساس کار رگرسیون و رگرسیون چندک

تعبیر شهودی نحوه کار رگرسیون معمولی بیشتر بر اساس تعبیر هندسی طرح حداقل مربعات است. مینیمم کردن فاصله اقلیدسی $\|Y - \hat{Y}\|$ روی تمام y ها، در فضای خطی تولید شده ستونهای X را می‌توان به این تشبیه کرد که تویی با مرکز y آنقدر باد می‌شود تا به زیر فضای تولید شده توسط X مماس شود. جایگزینی توپهای اقلیدسی با لوزیهای چند وجهی فاصله ρ_p ، بطوریکه، $d_p(y, \hat{y}) = \sum_{i=1}^n \rho_p(y_i - \hat{y}_i)$ بعضی مسائل جدید را مطرح می‌کند اما از بسیاری جهات بصورت قبل عمل می‌کند. در روش حداقل مربعات برای بدست آوردن معادلات نرمال عبارت زیر را مینیمم می‌کردیم:

$$\|y - \hat{y}(\beta)\|^2 = (y - X\beta)'(y - X\beta)$$



شکل ۴: برازش چندکهای رگرسیونی

با مشتق گیری نسبت به β :

$$\nabla \|y, \hat{y}(\beta)\|^2 = X'(y - X\beta) = 0 \quad \rightarrow \quad \hat{\beta} = (X'X)^{-1}X'y$$

$\hat{\beta}$ بدست می‌آید. در رگرسیون چندک هم با شیوه‌ای مشابه عمل می‌کنیم یعنی از تابع:

$$R(\beta) = d_p(y, \hat{y}(\beta)) = \sum_{i=1}^n \rho_p(y_i - x_i\beta) \quad (15)$$

بر حسب β مشتق می‌گیریم و با توجه به اینکه ممکن است مشتقات، به جهت بستگی داشته باشند، مشتقات سوئی زیر را در نظر می‌گیریم:

$$\begin{aligned} \nabla R(\beta, w) &= \frac{d}{dt} R(\beta + tw)|_{t=0} \\ &= \frac{d}{dt} \sum_{i=1}^n u_i(\beta + tw) [p - I_{(u_i(\beta+tw) < 0)}] \\ &= - \sum_{i=1}^n \psi^*(y_i - x_i'\beta, -x_i'w) x_i'w \end{aligned} \quad (16)$$

که در آن:

$$\psi^*(u, v) = \begin{cases} p - I(u < 0) & \text{if } u \neq 0 \\ p - I(v < 0) & \text{if } u = 0 \end{cases} \quad (17)$$

اگر به ازای تمام $w \in R^p$ و با فرض $\|w\| = 1$ ، داشته باشیم $\nabla R(\hat{\beta}, w) \geq 0$ ، در این صورت $\hat{\beta}$ تابع $R(\beta)$ را مینیمم می‌کند.

اساس تعریف چندکهای رگرسیونی:

همانگونه که جستجو برای یافتن میانه، ما را به یک مشاهده میانی یکتا یا یک جفت از مشاهدات وسطی مجاور هدایت می‌کند که هر دو جواب، مجموع قدر مطلق باقی مانده‌ها را مینیمم می‌کنند. در رگرسیون چندک هم به جستجوی زیر مجموعه‌های p عضوی که جواب را مشخص می‌کنند، رهنمون می‌شویم. در اصطلاحات برنامه‌ریزی خطی، به این زیر مجموعه‌های p عضوی، جوابهای پایه گفته می‌شود. در بررسی این زیر مجموعه‌های p عضوی از مشاهدات چند نماد را معرفی می‌کنیم $B(p)$ را مجموعه تمام نقاطی در نظر می‌گیریم که عبارت $\sum_{i=1}^n \rho_p(y_i - x_i\beta)$ را مینیمم می‌کند و $\beta(p)$ را یک عضو این مجموعه تعریف می‌کنیم. همچنین فرض می‌کنیم

که اعضای آن را با h و متمم مجموعه h را با $h^- = N - h$ نشان می‌دهیم. ماتریس X و بردار y را بر این اساس تقسیم بندی می‌کنیم. $X(h)$ نشان دهنده یک ماتریس $p \times p$ با سطرها $\{x_i : i \in h\}$ است و حال آنکه $y(h)$ یک بردار p بعدی با مختصات $\{y_i : i \in h\}$ است. ip یک بردار p بعدی با عناصر یک است یعنی $i'p = (1, 1, \dots, 1)$ و در نهایت H^* ، مجموعه زیر مجموعه‌های p عضوی از N می‌باشد که به ازای هر $X(h), h \in H$ ، رتبه کامل باشد. یعنی

$$H^* = \{h \in H \mid \text{rank} X(h) = p ; |X(h)| \neq 0\}$$

اکنون با بکارگیری این نمادها، می‌توان هر جواب پایه‌ای را که از نقاط $\{(x_i, y_i) : i \in h\}$ ، می‌گذرد با فرض ناویژه بودن، ماتریس $X(h)$ بصورت $b(h) = X^{-1}(h)y(h)$ بیان کرد. تعریف: در مشاهدات رگرسیونی (y, X) ، اگر به ازای هر p مشاهده، یک برازش دقیق و منحصر بفرود داشته باشیم در موقعیت عمومی هستند یعنی برای هر $h \in H^* = \{h \in H ; |X(h)| \neq 0\}$ داشته باشیم:

$$y_i - x_i b(h) = 0 \quad \forall i \in h$$

خواص چندکهای رگرسیونی:

قضیه ۱: اگر X دارای رتبه p باشد و برای بعضی از عناصر H^* ، مجموعه $B(p)$ ، حداقل یک عضو به صورت $b(h) = X^{-1}(h)y(h)$ داشته باشد، آنگاه $B(p)$ ، یک پوسته محدب می‌باشد و تمام اعضای آن دارای همین شکل هستند. این قضیه شکل چندکهای رگرسیونی را تحت یک قانون کلی بیان می‌کند.

قضیه ۲: اگر (y, X) در موقعیت عمومی باشند، آنگاه یک جواب منحصر بفرود برای مسئله چندک رگرسیونی بصورت $b(h) = X^{-1}(h)y(h)$ وجود دارد اگر و فقط اگر به ازای یک مقدار $h \in H^* = \{h \in H ; |X(h)| \neq 0\}$ داشته باشیم:

$$(p-1)I_p \leq \mu_h \leq pI_p$$

که در آن:
$$\mu_h = - \sum_{i \in h} \psi_p(y_i - x_i' b(h)) \quad , \quad - x' i X^{-1}(h) v x_i' X^{-1}(h) \quad , \quad \psi_p = p - I(u < 0) \quad , \quad I' p = (1, 1, \dots, 1)_{1 \times p}$$

علاوه بر این، $b(h)$ جواب یکتاست اگر و فقط اگر نامعادله‌ها اکید باشند در غیر اینصورت مجموعه جواب $B(p)$ ، پوسته محدبی از جوابهای متعددی به شکل $b(h)$ است. توجه: تباهیدگی اولیه در مسئله رگرسیون چندک به این معنی است که مشاهدات رگرسیونی (y, X) در موقعیت عمومی نباشند یعنی بیشتر از p باقیمانده صفر داشته باشیم.

قضیه ۳: فرض کنید N^+, N^-, N^0 بترتیب برابر تعداد مولفه‌های مثبت، منفی و صفر بردار

باقیمانده $y - X'\hat{\beta}(p)$ ، باشند اگر $a \in R^p$ وجود داشته باشد بطوریکه $Xa = I_n$ باشد، پس برای هر $\hat{\beta}(p)$ که جواب معادله:

$$\min_{\beta \in R^p} \sum_{i=1}^n \rho_p(y_i - x'_i \beta)$$

می‌باشد داریم: $N^+ \leq n(1-p) \leq N^+ + N^0$ و $N^- \leq np \leq N^- + N^0$ نتیجه: اگر $N^0 = p$ باشد، در این صورت هیچ تباهی وجود ندارد بنابراین نسبت باقیمانده‌های منفی تقریباً برابر p و نسبت باقیمانده‌های مثبت تقریباً برابر $(1-p)$ است یعنی:

$$\frac{N^-}{n} \leq p \leq \frac{N^- + p}{n}, \quad \frac{N^+}{n} \leq (1-p) \leq \frac{N^+ + p}{n}$$

قضیه ۴: فرض کنید $A_p \times p$ ، یک ماتریس ناویژه در R^p و $a > 0$ باشد

اگر $\beta(p) \in B(p)$ باشد آنگاه جوابهای $R(\beta) = \sum_{i=1}^n \rho_p(y_i - x'_i \beta)$ به ازای هر $p \in (0, 1)$ دارای خواص زیر هستند:

- (i) $\hat{\beta}(p; ay, X) = a\hat{\beta}(p; y, X)$
- (ii) $\hat{\beta}(p; -ay, X) = a\hat{\beta}(1-p; y, X)$
- (iii) $\hat{\beta}(p; y + X\gamma, X) = \hat{\beta}(p; y, X) + \gamma$
- (iv) $\hat{\beta}(p; y, XA) = A^{-1}\hat{\beta}(p; y, X)$

خواص (i) و (ii) پایایی مقیاس و (iii) پایایی انتقال یا پایایی رگرسیونی و (iv) پایایی دوباره پارامتری کردن طرح نامیده می‌شود.

کانکر و باست برای مدل رگرسیون خطی $y_i = x'_i \beta + u_i \quad i = 1, 2, \dots, n$ جواب:

$$\hat{\beta}_n(p) = \min_{\beta \in R^p} \sum_{i=1}^n \rho_p(y_i - x'_i \beta)$$

را بعنوان برآورد چندک رگرسیونی p ام معرفی کردند. این چندکهای رگرسیونی p بعدی یک دنباله از ابر صفحه‌ها را تعیین می‌کنند که توابع چندک شرطی متغیر پاسخ y را برآورد می‌کنند همانند مدل یک نمونه‌ای، زمانی که p بین $(0, 1)$ تغییر می‌کند، یک برنامه‌ریزی خطی پارامتری داریم که می‌تواند به سادگی حل شود.

۷ مثال ۲

انگل در سال ۱۸۵۷، رابطه بین هزینه غذا و درآمد خانوار را مورد تجزیه و تحلیل قرار داد. او یک رگرسیون ساده از داده‌ها را که مربوط به ۲۳۵ خانوار اروپایی از رده کارگران می‌باشد، در

نظر گرفت. در اینجا ساده بدین مفهوم است که بردار متغیر تصادفی، تابعی یک بعدی است و داده‌ها را می‌توان به وسیله نمودار پراکنندگی زوجهای (x_i, y_i) نشان داد که y_i میزان هزینه غذا و x_i میزان درآمد خانوار را نشان می‌دهد. در شکل (۵)، نمودار ۸ خط رگرسیون برازش داده شده نشان داده شده است که ۷ خط مربوط به خطوط چندکهای رگرسیونی متناظر با $p \in \{0/05, 0/1, 0/25, 0/5, 0/75, 0/9, 0/95\}$ و یک خط هم مربوط به برآورد خط حداقل مربعات تابع میانگین شرطی می‌باشد که به صورت خط چین رسم شده است. برازش میانه در $p = 0/5$ ، توسط خط ممتد تیره‌تر نشان داده شده است. نمودار به روشنی معلوم می‌کند که هزینه غذا با افزایش درآمد خانوار افزایش پیدا می‌کند فضای خطوط رگرسیون چندک، همچنین معلوم می‌کند که توزیع شرطی هزینه غذا به طرف چپ متمایل شده است. فضای کم چندکهای بالاتر (مرتبه بالاتر) نشان می‌دهد که فشردگی داده‌ها زیاد است، یعنی در این فضا داده‌ها بیشترند و فضای زیاد بین چندکهای پایین‌تر (مرتبه پایین‌تر) نشان می‌دهد که در این فضا فشردگی داده‌ها کمتر است. یعنی در این قسمت داده‌ها کمتر هستند.

برازشهای میانگین و میانه شرطی در این مثال کاملاً متفاوتند و برازش خط حداقل مربعات یک برآورد نسبتاً ضعیفی از میانگین شرطی برای فقیرترین خانواده‌ها در نمونه می‌باشد چون از بالای اکثر مشاهدات با درآمد خیلی پایین می‌گذرد.

۸ بحث

وقتی که خطاها دارای توزیع نرمال هستند برآوردهای حداقل مربعات کاراترین برآوردهای ناریب برای ضرایب رگرسیون هستند اما وقتی توزیع خطاها دمه‌ای طولانی داشته باشد این برآوردها نمی‌توانند خیلی کارا باشند تحت این شرایط ما انتظار داریم که در داده‌ها نقاط دور افتاده وجود دارد. در حالتی که داده‌ها دارای نقاط دور افتاده باشند برازش حداقل مربعات روش خوبی نیست. به همین دلیل است که روش‌های برازش دیگری ارائه گردیده است که نسبت به نقاط دور افتاده حساس نیستند.

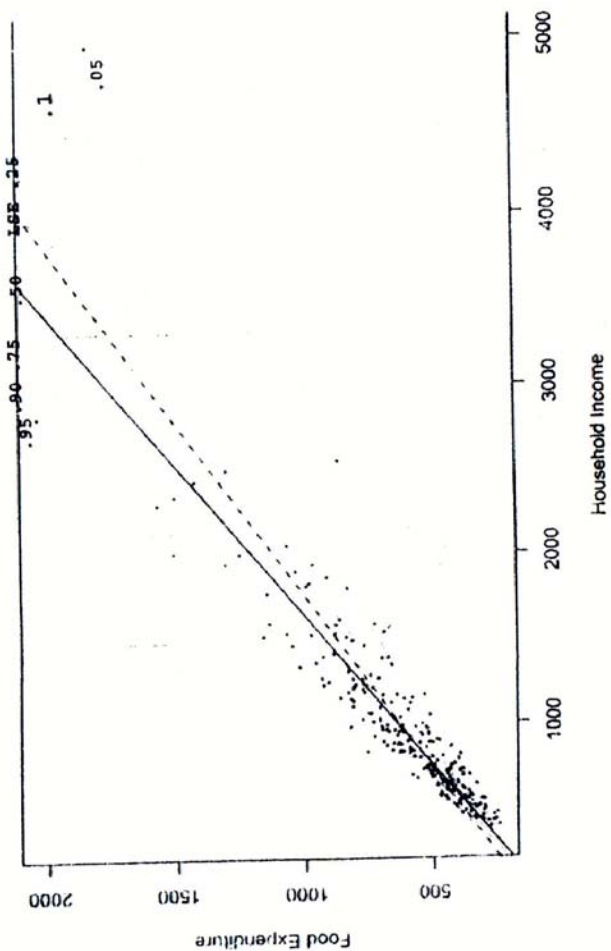
حساسیت کمترین مربعات به نقاط دور افتاده وابسته به دو عامل است.

(۱) اگر ما از معیار مربعات باقیمانده برای اندازه‌گیری استفاده کنیم هر باقیمانده بزرگ در مقایسه با سایر باقیمانده‌ها اندازه خیلی بزرگ خواهد داشت.

(۲) با استفاده از یک معیار مکانی مانند میانگین که نیرومند نیست هر مربع خطای بزرگ تأثیر خیلی قوی روی معیار، خواهد داشت در نتیجه در داده‌های انتهایی تأثیر نامناسبی بر برازش خواهد داشت.

دوره ساده برای حل این مسئله وجود دارد.

(۱) ما می‌توانیم به جای معیار اندازه e^2 از تابع دیگری از e یعنی $\rho(e)$ استفاده کنیم که اندازه باقیمانده‌ها را با کرانگین کمتری انعکاس دهد برای اینکه این تابع اندازه محسوس باشد تابع ρ باید: (الف) متقارن باشد. $\rho(e) = \rho(-e)$ (ب) مثبت باشد. $\rho(e) \geq 0$ (ج) یکنوا باشد. یعنی



شکل ۵: برازش حداقل مربعات میانگین و چندکهای رگرسیونی داده‌های انگل

$|e_1| \geq |e_2|$ if $\rho(|e_1|) \geq \rho(|e_2|)$ این نظریه ما را به سمت برآورد M (برآوردی که تابعی از $\rho(e)$ به کار می‌برد که ترکیبی از e^2 و $|e|$ می‌باشد) می‌برد.

(۲) می‌توانیم به جای مجموع (یا بطور معادل میانگین) از معیار مکانی نیرومندتری از قبیل میانه (حالت خاص رگرسیون چندک) یا میانگین پیراسته (*trimmed mean*) استفاده می‌کنیم. اساس روشهای رگرسیونی برگرفته از این ایده شامل حداقل میانه مربعات و حداقل میانگین پیراسته مربعات می‌شود.

مهمترین انگیزه برای پیدایش رگرسیون چندک نیرومندی ذاتی آن در مقابل مشاهدات دورافتاده در متغیر پاسخ بوده است. در حالیکه رگرسیون معمولی به یک نقطه دور افتاده حساسیت دارد تأثیر مشاهدات دور افتاده بر $\hat{\beta}(p)$ محدود است در واقع دور شدن مشاهدات از چندک رگرسیونی برازش داده شده (مرتبه y ام) هیچ تأثیری روی برازش مدل ندارد (که در مثال ذکر شده نمایان است).

مراجع

- [1] Georgea, F. Seber and Alanj, Lee (2003). Linear Regression Analysis Sonc. Inc. Allrights reserved Auckland , New Zealand
- [2] Gilbert Basset, JR. and Roger Koenker, June (1982) An Emprical Quantile Function for Linear Models With iid Error Journal of the American Statistical Association ,Volume 77 ,Nvmber 378 Theory and Methods Section
- [3] Hogg, R.V (1975) Estimates of percentile regression lines using salary data journal of the American statistical Association 70. 56-59
- [4] Koenker, R. and Bassett, G.S., (1978). Regression Quantiles. *Econometrica* **46**, 33-50.
- [5] Roger Koenker June 23, (2002). Quantile Regression Another Introduction Research was Partially supported by NSF grant SBR-9911184.
- [6] Roger Koenker October 25, (2000). Quantile regression International Encyclopedia of the social sciences
- [7] Roger Koenker November 10, (2000).Quantile Regression Department of Economhcs University of Illinois Urbana-champaign champaign61820 USA

میانگین‌گیری بیزی مدل‌های رگرسیونی

افشین فلاح، محسن محمدزاده

گروه آمار دانشگاه تربیت مدرس

چکیده: در روش‌های مدل‌سازی که بسته به معیار انتخاب مدل و نظر تحلیلگر مدل‌های متفاوت انتخاب می‌شوند، عدم حتمیت در فرآیند مدل‌سازی نادیده گرفته می‌شود. میانگین‌گیری بیزی مدل‌ها شیوه‌ای توانمند در مدل‌سازی است، که در آن به هر مدل به تناسب میزان حمایتی که از جانب داده‌ها صورت می‌پذیرد، وزنی اختصاص داده می‌شود و میانگین وزنی همه مدل‌ها بعنوان مدل نهایی بکار گرفته می‌شود. در این مقاله نشان داده می‌شود که روش میانگین‌گیری بیزی، خطای مدل‌سازی را کاهش و کارایی را افزایش می‌دهد.

واژه‌های کلیدی: عدم حتمیت، توزیع پیش‌بین، میانگین‌گیری بیزی

۱ مقدمه

معمولاً در مدل‌سازی، کلاسی از مدل‌ها در نظر گرفته می‌شود، سپس یکی از آنها بر اساس یک یا چند معیار ارزیابی بعنوان بهترین مدل انتخاب می‌شود. در حالی که ممکن است رقبای بسیار خوبی برای مدل انتخابی در فضای مدل وجود داشته باشند. این شیوه مدل‌سازی دارای نارسایی‌هایی است که مهمترین آنها در نظر نگرفتن عدم حتمیت مدل‌های انتخابی در فرآیند مدل‌سازی است. میانگین‌گیری بیزی مدل‌ها^۱ (*BMA*)، شیوه‌ای است که در آن از تمام مدل‌های موجود در فضای مدل برای دستیابی به مدلی مناسب استفاده می‌شود. در این روش بر اساس میزان حمایت داده‌ها از هر مدل، وزنی به آن اختصاص داده می‌شود. سپس مدل حاصل از میانگین وزنی همه مدل‌ها برای انجام استنباط و پیش‌بینی به کار گرفته می‌شود. ریشه تاریخی این روش به مقاله لیمر (۱۹۷۸) برمی‌گردد، که به دلیل محاسبات دشوار و پیچیده در آن زمان چندان مورد توجه قرار نگرفت. با پیشرفت رایانه‌ها و فنون محاسبات تقریبی در اوایل دهه نود این روش دوباره مطرح و در کانون توجهات قرار گرفت. مادیکان و رفتری (۱۹۹۴) و مادیکان و یورک (۱۹۹۵) دو روش پایه‌ای برای اجرای این شیوه مدل‌سازی ارائه کردند. هوئینگ و همکاران (۱۹۹۴، ۱۹۹۶، ۱۹۹۹) به نحوه انجام محاسبات و اجرای این روش پرداختند. روبرت (۲۰۰۰) و لیکوویچ (۲۰۰۲) میانگین‌گیری بیزی را در حالت چند متغیره مورد توجه قرار دادند. در این مقاله روش میانگین‌گیری بیزی مدل‌ها، نحوه محاسبه مؤلفه‌های و روش‌های اجرای آن در بخش ۲ شرح داده

1) Bayesian Model Averaging

می‌شود و سپس در بخش ۳ به ارزیابی روش *BMA* پرداخته و نهایتاً بحث و نتیجه‌گیری در بخش ۴ ارائه خواهد شد.

۲ میانگین‌گیری بیزی مدلها

فرض کنید برای متغیر وابسته Y و مجموعه‌ای از متغیرهای پیش‌بین X_1, \dots, X_k ، هدف یافتن بهترین مدل از بین همه مدل‌های خطی

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X_j + \varepsilon \quad (1)$$

باشد، که در آن X_{i1}, \dots, X_{ip} زیرمجموعه‌ای از متغیرهای X_1, \dots, X_k هستند. بسته به حضور یا عدم حضور هر متغیر پیش‌بین در مدل، 2^k مدل رگرسیونی وجود دارد، که فرض می‌شود همگی در فضای مدل \mathcal{M} قرار دارند. اگر $\mathcal{M} = \{M_1, \dots, M_T\}$ نشان دهنده فضای مدل و Δ کمیتی مورد علاقه، مانند مشاهده‌ای در آینده باشد. در اینصورت استنباط بیزی بر اساس توزیع پسین Δ صورت می‌پذیرد، که با فرض مشاهده مجموعه داده D ، بنا بر قاعده احتمال کل، آمیخته‌ای از احتمال‌های پسین همه مدلها بصورت

$$Pr(\Delta|D) = \sum_{k=1}^T Pr(\Delta|M_k, D) \cdot Pr(M_k|D) \quad (2)$$

است. میانگین و واریانس پسین Δ به ترتیب بصورت

$$\begin{aligned} E[\Delta|D] &= \sum_{k=1}^T E[\Delta|M_k, D] \cdot Pr(M_k|D) \\ &= \sum_{k=1}^T \hat{\Delta}_k \cdot Pr(M_k|D) \end{aligned} \quad (3)$$

و

$$\begin{aligned} Var[\Delta|D] &= E_M[Var(\Delta|D, M)] + Var_M[E(\Delta|D, M)] \\ &= \sum_{k=1}^T (Var[\Delta|D, M_k] + \hat{\Delta}_k^2) \times Pr(M_k|D) - E[\Delta|D]^2 \end{aligned} \quad (4)$$

هستند (رفتری، ۱۹۹۳)، که در آن مؤلفه $Var_M[E(\Delta|D, M)]$ عدم حتمیت بین مدلها را نشان می‌دهد. عبارت (۲) یک میانگین وزنی احتمال پسین است، که در آن هر مدل با احتمال

پسین متناظر خود، یعنی $P(M_k|D)$ ، وزن دار شده است و توزیع پیش‌بین Δ به شرط مدل M_k بصورت

$$Pr(\Delta|M_k, D) = \int Pr(\Delta|\theta_k, M_k, D).Pr(\theta_k) d\theta_k \quad (5)$$

است، که در آن θ_k بردار پارامترهای مدل M_k را نشان می‌دهد. احتمال پسین مدل M_k با استفاده از قاعدهٔ بیز بصورت

$$Pr(M_k|D) = \frac{Pr(D|M_k).Pr(M_k)}{\sum_{j=1}^T Pr(D|M_j).Pr(M_j)} \quad (6)$$

بدست می‌آید، که در آن $Pr(M_k)$ احتمال پیشین درست بودن مدل M_k و

$$Pr(D|M_k) = \int Pr(D|\theta_k, M_k).Pr(\theta_k|M_k) d\theta_k \quad (7)$$

درست‌نمایی جمع بستهٔ^۲ مدل M_k و $Pr(D|\theta_k, M_k)$ درست‌نمایی مدل M_k است. برای محاسبهٔ $Pr(\Delta|D)$ باید مؤلفه‌های تشکیل دهندهٔ آن را مشخص و جایگزین نمود. بنابراین لازم است توزیع پیشین پارامترها، احتمالهای پیشین و پسین هر مدل و توزیع پیش‌بین کمیت مورد علاقه را برای همهٔ مدل‌های موجود در فضای مدل، مشخص کرد.

الف -- تعیین احتمال پیشین پارامترهای مدل: یکی از مسائل دشوار در *BMA* تخصیص پیشین به پارامترهای مدل است. استفاده از پیشین‌های نامناسب^۳، منجر به توزیع‌های پسین نامناسب می‌شود، که در اینصورت نمی‌توان این احتمالها را به عنوان احتمال مدل و نسبت آنها را به عنوان بیز فاکتور تعبیر نمود. به همین دلیل بسیاری از محققان گونه‌های مختلفی از پیشین‌های آگاهی بخش^۴ را پیشنهاد کرده‌اند. هوئتینگ (۱۹۹۴) استفاده از پیشین‌های مناسبی^۵ را پیشنهاد کرده است که در قسمتهایی از فضای پارامتر با درست‌نمایی بزرگ هموار باشند. مدل (۱) را می‌توان بصورت

$$Y = X\beta + \varepsilon \quad (8)$$

نوشت، که در آن $X_{n \times (p+1)}$ ماتریس مشاهدات، Y بردار n بعدی متغیرهای وابسته و ε بردار خطا است. فرض می‌شود ε ها برای مشاهدات مختلف، مستقل و دارای توزیع نرمال با میانگین صفر و واریانس σ^2 هستند. بردار $\beta = (\beta_0, \dots, \beta_p)$ و پارامتر σ^2 نامعلوم هستند. توزیعهای پیشین این پارامترها بایستی بگونه‌ای تعیین شوند که عدم حتمیت آنها را به خوبی منعکس سازند. هوئتینگ (۱۹۹۴) ردهٔ پیشین‌های مزدوج نرمال-گاما را بصورت

$$\beta \sim N_{p+1}(\mu, \sigma^2 V) \quad , \quad \frac{\nu \cdot \lambda}{\sigma^2} \sim \chi_\nu^2 \quad (9)$$

2) Integrated Likelihood 3) Improper Prior 4) Informative 5) Proper Prior

در نظر گرفت، که در آنها ν ، λ ، ماتریس $V_{(p+1) \times (p+1)}$ و بردار $p+1$ بعدی μ ، ابر پارامترهایی هستند که باید برآورد شوند.

ب -- تعیین احتمال پیشین هر مدل: وقتی هیچ اطلاع پیشینی در مورد مدلها وجود ندارد یا میزان اطلاعات پیشین در مورد مدلها اندک است، معمولاً فرض می شود توزیع مدلها یکنواخت است. یعنی $Pr(M_j) = \frac{1}{T}$ ، در اینصورت احتمال پسین مدل M_k بصورت

$$Pr(M_k|D) = \frac{Pr(D|M_k)}{\sum_{j=1}^T Pr(D|M_j)} \quad (10)$$

خواهد بود. اگر بعضی از مدلها در مقایسه با سایرین دارای احتمال بیشتری باشند یا اطلاعات خوبی در مورد آنها در دسترس باشد، لازم است این اطلاعات برای تعدیل احتمالهای پیشین مدلها بکار گرفته شوند تا از پیشینهای آگاهی بخش تر استفاده شود. در مسائل مربوط به انتخاب متغیرهای پیشین، اطلاعات پیشین بشکل شواهد قبلی برای در نظر گرفتن یک متغیر مورد استفاده قرار می گیرند. فرض کنید تنها این نوع از اطلاعات پیشین در اختیار باشد و مدل M_k توسط بردار $(\delta_{k1}, \dots, \delta_{kp})$ مشخص شود، که در آن

$$\delta_{ki} = \begin{cases} 1 & X_i \in M_k \\ 0 & X_i \notin M_k \end{cases} \quad i = 1, \dots, p$$

توابع نشانگر هستند. حال اگر احتمال مؤثر بودن متغیر X_i بر Y را نشان دهد و بپذیریم که اطلاعات پیشین در مورد متغیرهای متفاوت تقریباً مستقل هستند، می توان

$$Pr(M_k) = \prod_{i=1}^p [\pi_i^{\delta_{ki}} \times (1 - \pi_i)^{1-\delta_{ki}}] \quad (11)$$

را به عنوان احتمال پیشین صحیح بودن مدل M_k در نظر گرفت. چون این توزیع به متغیری که مهمتر است احتمال بزرگتری تخصیص می دهد، توسط مادیکان و رفتری (۱۹۹۱) توزیع پیشین متغیر نامیده شده است.

ج -- تعیین احتمال پسین هر مدل: برای محاسبه احتمال پسین هر مدل لازم است $Pr(D|M_k)$ از رابطه (۷) که یک انتگرال با بعد برابر با تعداد پارامترهای مدل M_k است، محاسبه شود. بنابراین محاسبه دقیق این احتمال به دلیل پیچیده بودن انتگرال مربوطه تنها در حالات بسیار خاص و ساده امکان پذیر است و در سایر موارد از روشهای تقریبی و محاسباتی استفاده می شود. در اینجا از معیار اطلاع بیز^۷ (BIC) برای این منظور بهره می بریم. معیار BIC برای رگرسیون خطی بصورت

$$BIC_j = n \log(1 - R_j^2) + k_j \log n \quad (12)$$

6) Variable Prior 7) Bayesian Information Criteria

تعریف می‌شود، که در آن n تعداد مشاهدات، R_j^* ضریب تعیین تعدیل شده مدل و k_j تعداد متغیرهای پیش‌بین موجود در مدل j ام و r_j^* را نشان می‌دهد. در اینصورت درست‌نمایی جمع بسته مدل j ام را می‌توان بصورت تقریبی

$$Pr(D|M_j) \propto e^{-\circ/\delta BIC_j} = e^{-\circ/\delta(n \log(1-r_j^*)+k_j \log n)} \quad (۱۳)$$

خواهد بود. بر این اساس احتمال پسین مدل k ام با استفاده از قاعدهٔ بیز بصورت

$$Pr(M_k|D) = \frac{\exp\{-\circ/\delta BIC_k\} \cdot Pr(M_k)}{\sum_{j=1}^T \exp\{-\circ/\delta BIC_j\} \cdot Pr(M_j)} \quad (۱۴)$$

نوشت. چون معیار BIC برای بسیاری از مدلها دارای شکل بسته و ساده‌ای است، استفاده از این تقریب موجب سهولت محاسبه و افزایش سرعت می‌شود.

د -- تعیین توزیع پیش‌بین: انتگرال تشکیل دهندهٔ توزیع پیش‌بین (۵) از دو جزء تشکیل شده است. معمولاً برای جزء دوم آن از تقریب

$$Pr(\Delta|M_k, D) \approx Pr(\Delta|M_k, \hat{\theta}_k, D) \quad (۱۵)$$

استفاده می‌شود (رفتاری و همکاران، ۱۹۹۴)، که در آن $\hat{\theta}_k$ برآورد حداکثر درست‌نمایی بردار پارامترهای مدل M_k است.

اکنون با فرض مشخص بودن همهٔ مولفه‌های لازم برای اجرای روش BMA محاسبهٔ مجموع (۲) بواسطهٔ تعداد زیاد جملات آن عملاً امکان‌پذیر نیست و لازم است زیر مجموعه‌ای از محتمل‌ترین مدلها انتخاب شود. مادیگان و رفتاری (۱۹۹۴) روش پنجره اوکام^۸ را برای این منظور پیشنهاد کرده‌اند، که از دو اصل کلی پیروی می‌کنند. بنابر اصل اول مدلهایی که در مقایسه با محتمل‌ترین مدل خیلی کم شانس هستند، کنار گذاشته می‌شوند. بنابر اصل دوم که تیغ اوکام^۹ نامیده می‌شود، مدلهایی که نسبت به زیرمدلهای ساده‌تر خود کمتر از جانب داده‌ها حمایت می‌شوند، کنار گذاشته می‌شوند. با اجرای اصل اول مدلهایی که در مجموعه

$$\mathcal{A}' = \left\{ M_k : \frac{\text{Max}_{M_l \in \mathcal{M}} \{Pr(M_l|D)\}}{Pr(M_k|D)} > C_1 \right\}$$

قرار دارند و همچنین بنابر اصل تیغ اوکام، مدلهایی که در مجموعه

$$\mathcal{B} = \left\{ M_k : \exists M_l \in \mathcal{A} \text{ s.t. } M_l \subset M_k, \frac{Pr(M_l|D)}{Pr(M_k|D)} > C_2 \right\}$$

8) Occam's Window 9) (OW)Occam's Razor

باشند، از مجموع (۲) خارج می‌شوند، که در آن C_1 و C_2 توسط تحلیل‌گر تعیین می‌شوند. در اینصورت مجموع (۲) را می‌توان بصورت

$$Pr(\Delta|D) = \frac{\sum_{M_k \in A} Pr(\Delta|M_k, D) \cdot Pr(D|M_k) \cdot Pr(M_k)}{\sum_{M_k \in A} Pr(D|M_k) \cdot Pr(M_k)} \quad (۱۶)$$

بازنویسی کرد، که در آن $A = A' - B \in M$ مجموعه پذیرش است.

۳ مثال کاربردی

در این بخش کاربردی روش BMA در تشخیص متغیرهای پیش‌بین مؤثر در مدل و برآورد پارامترهای مدل با روش مدل‌سازی گام به گام بعنوان یکی از روشهای مدل‌سازی مرسوم مورد مقایسه قرار می‌گیرد. با توجه به اینکه بطور طبیعی کاربردی روش BMA با افزایش تعداد متغیرهای پیش‌بین، افزایش می‌یابد، برای نمایش بهتر کاربردی روش BMA از مثالی با تعداد متغیرهای پیش‌بین کم استفاده شده تا حتی‌المکان از بالا بودن کاربردی این روش بدلیل تعداد زیاد متغیر پیش‌بین اجتناب شود.

جدول ۱: داده‌های مربوط به کاربردی شغلی پرستاران

| | | | | | | | | | | | | | | | |
|-------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Y | ۴۵ | ۶۵ | ۷۳ | ۶۳ | ۸۳ | ۴۵ | ۶۰ | ۷۳ | ۷۴ | ۶۹ | ۶۶ | ۶۹ | ۷۱ | ۷۰ | ۷۹ |
| X_1 | ۷۴ | ۶۵ | ۷۱ | ۶۲ | ۷۹ | ۵۶ | ۶۸ | ۷۶ | ۸۳ | ۶۲ | ۵۴ | ۶۱ | ۶۳ | ۸۴ | ۷۸ |
| X_2 | ۲۹ | ۵۰ | ۶۷ | ۴۴ | ۵۵ | ۴۸ | ۴۱ | ۴۹ | ۷۱ | ۴۴ | ۵۲ | ۴۶ | ۵۶ | ۸۲ | ۵۳ |
| X_3 | ۴۰ | ۶۴ | ۷۹ | ۵۷ | ۷۶ | ۵۴ | ۶۶ | ۶۵ | ۷۷ | ۵۷ | ۶۷ | ۶۶ | ۶۷ | ۶۸ | ۸۲ |
| X_4 | ۶۶ | ۶۸ | ۸۱ | ۵۹ | ۷۶ | ۵۹ | ۷۱ | ۷۵ | ۷۶ | ۶۷ | ۶۳ | ۶۴ | ۶۰ | ۶۴ | ۶۴ |
| X_5 | ۹۳ | ۷۴ | ۸۷ | ۸۵ | ۸۴ | ۵۰ | ۶۹ | ۶۷ | ۸۴ | ۸۱ | ۶۸ | ۷۵ | ۶۴ | ۷۸ | ۷۸ |
| X_6 | ۴۷ | ۴۹ | ۳۳ | ۳۷ | ۳۳ | ۴۲ | ۳۷ | ۴۳ | ۳۳ | ۴۳ | ۳۶ | ۴۳ | ۳۵ | ۳۷ | ۳۹ |
| Y | ۸۳ | ۷۵ | ۶۷ | ۶۷ | ۵۲ | ۵۲ | ۶۶ | ۵۵ | ۴۲ | ۶۵ | ۶۸ | ۸۰ | ۵۰ | ۸۷ | ۸۴ |
| X_1 | ۶۵ | ۸۶ | ۶۱ | ۷۱ | ۵۹ | ۷۱ | ۶۲ | ۶۷ | ۶۵ | ۵۵ | ۷۸ | ۷۶ | ۵۸ | ۸۶ | ۸۳ |
| X_2 | ۴۹ | ۶۳ | ۶۴ | ۴۵ | ۶۷ | ۳۲ | ۵۱ | ۵۱ | ۴۱ | ۶۵ | ۵۷ | ۴۳ | ۴۳ | ۷۰ | ۳۸ |
| X_3 | ۸۲ | ۷۹ | ۷۵ | ۶۷ | ۶۴ | ۴۴ | ۷۲ | ۶۰ | ۴۵ | ۵۸ | ۷۳ | ۸۴ | ۵۵ | ۸۱ | ۸۳ |
| X_4 | ۸۴ | ۶۵ | ۸۴ | ۶۰ | ۸۰ | ۶۹ | ۴۸ | ۷۱ | ۶۸ | ۵۵ | ۷۱ | ۹۳ | ۸۵ | ۵۶ | ۸۲ |
| X_5 | ۵۵ | ۸۰ | ۸۱ | ۸۶ | ۷۹ | ۶۵ | ۸۱ | ۸۱ | ۵۸ | ۷۶ | ۷۷ | ۷۹ | ۸۴ | ۷۵ | ۷۹ |
| X_6 | ۳۸ | ۴۱ | ۴۵ | ۴۸ | ۵۴ | ۴۳ | ۴۳ | ۳۹ | ۵۱ | ۳۵ | ۴۲ | ۳۵ | ۴۰ | ۳۰ | ۴۱ |

جدول ۱ مشاهدات یک مطالعه رگرسیونی را نشان می‌دهد، که هدف آن بررسی تأثیر ۶ متغیر مستقل شامل قاطعیت (X_1)، علاقه‌مندی (X_2)، بلند همتی (X_3)، روابط اجتماعی (X_4)، هنرمندی در حل مشکلات (X_5) و ابتکار (X_6) بر کاربردی شغلی پرستاران است. بدلیل وجود ۶ متغیر پیش‌بین، فضای مدل شامل $2^6 = 64$ مدل است. روش پنجره اوکام ۱۸ مدل را بر اساس مقایسه احتمالهای پسین مدلها به عنوان برترین مدلها موجود مورد استفاده قرار می‌دهد. مدلها حاصل به همراه احتمال پسین، مقدار R_{adj}^2 به درصد و معیار اطلاع بیز آنها در جدول ۲ ارائه شده‌اند. تعداد زیاد مدلهایی که در پنجره اوکام قرار گرفته‌اند (بیش از $\frac{1}{p}$ مدلها)، نشان

جدول ۲: مدل‌های انتخابی توسط روش پنجره اوکام

| شماره مدل | مدل | احتمال پسین | R_{adj}^2 | BIC |
|-----------|-------------------------------|-------------|-------------|----------|
| ۱ | X_3 X_6 | ۰٫۱۹۷۶ | ۸۲٫۲۰۴ | -۴۱٫۷۵۹۷ |
| ۲ | X_2 X_3 X_6 | ۰٫۱۴۷۳ | ۸۲٫۹۲۲ | -۴۱٫۱۵۲۱ |
| ۳ | X_1 X_2 X_3 X_6 | ۰٫۱۱۷۵ | ۸۴٫۵۵۵ | -۴۰٫۶۰۰۵ |
| ۴ | X_1 X_3 X_6 | ۰٫۰۹۸۹ | ۸۲٫۴۴۷ | -۴۰٫۳۵۶۵ |
| ۵ | X_3 | ۰٫۰۶۸۸ | ۷۷٫۲۹۶ | -۳۹٫۶۲۸۹ |
| ۶ | X_3 X_5 X_6 | ۰٫۰۴۸۱ | ۸۱٫۵۵۲ | -۳۸٫۹۱۴۳ |
| ۷ | X_2 X_3 X_5 X_6 | ۰٫۰۴۲۱ | ۸۳٫۴۲۳ | -۳۸٫۶۴۸۳ |
| ۸ | X_3 X_4 X_6 | ۰٫۰۳۸۳ | ۸۱٫۲۶۱ | -۳۸٫۴۶۰۴ |
| ۹ | X_2 X_3 | ۰٫۰۳۵۰ | ۷۸٫۸۲۰ | -۳۸٫۲۷۶۶ |
| ۱۰ | X_1 X_3 | ۰٫۰۳۴۹ | ۷۸٫۸۱۶ | -۳۸٫۲۷۱۲ |
| ۱۱ | X_2 X_3 X_4 X_6 | ۰٫۰۳۰۳ | ۸۳٫۰۴۱ | -۳۷٫۹۸۷۶ |
| ۱۲ | X_1 X_2 X_3 | ۰٫۰۲۶۱ | ۸۰٫۷۵۴ | -۳۷٫۶۸۶۲ |
| ۱۳ | X_1 X_2 X_3 X_5 X_6 | ۰٫۰۲۵۰ | ۸۷٫۷۰۰ | -۳۷٫۶۰۵۷ |
| ۱۴ | X_1 X_2 X_3 X_4 X_6 | ۰٫۰۲۲۳ | ۸۴٫۵۸۲ | -۳۷٫۳۸۲۹ |
| ۱۵ | X_1 X_3 X_5 X_6 | ۰٫۰۱۹۷ | ۸۲٫۵۳۴ | -۳۷٫۱۳۳۳ |
| ۱۶ | X_1 X_3 X_4 X_6 | ۰٫۰۱۹۱ | ۸۲٫۴۹۵ | -۳۷٫۰۶۸۶ |
| ۱۷ | X_3 X_5 | ۰٫۰۱۶۲ | ۷۷٫۶۶۳ | -۳۶٫۷۳۴۳ |
| ۱۸ | X_3 X_4 | ۰٫۰۱۲۸ | ۷۷٫۳۰۴ | -۳۶٫۲۷۱۸ |

دهنده زیاد بودن عدم حتمیت است، یعنی بیش از یک مدل صحیح وجود دارد، بنابراین انتخاب یک مدل به عنوان مدل نهایی، منطقی بنظر نمی‌رسد. ملاحظه می‌شود که مقادیر R_{adj}^2 و معیار اطلاع بیز BIC ، برای مدل‌هایی که در پنجره اوکام قرار گرفته‌اند، نزدیک به یکدیگر هستند. بیشترین مقدار R_{adj}^2 مربوط به مدل شماره ۱۳ است، در حالی که بر اساس مقدار کمیته معیار اطلاع بیز، مدل ۱ بهترین است. انتخاب یک مدل و مینا قرار دادن آن به معنی آن است که احتمال صحیح بودن آن یک می‌باشد، در حالی که بیشترین احتمال پسین مربوط به مدل شماره ۱ و حدود ۰٫۲ است و احتمال پسین سایر مدل‌ها همگی از ۰٫۲ کمتر است. جدول ۲ نشان می‌دهد که متغیر X_3 در همه مدل‌ها و متغیر X_6 در ۱۲ مدل منتخب روش پنجره اوکام حضور دارند، که نشان دهنده تأثیرگذار بودن این دو متغیر است. برآورد ضرایب متغیرهای مستقل که در جدول ۳ ارائه شده نیز اهمیت دو متغیر X_3 و X_6 را تأیید می‌کند. همانطور که در جدول ۳ ملاحظه می‌شود، برآوردهای بیز ضرایب این دو متغیر که با استفاده از روش BMA بدست آمده‌اند، به خوبی اهمیت این دو متغیر را منعکس می‌سازند. جدول ۴ برآورد پارامترها و خطای معیار آنها را که به روش گام به گام بدست آمده‌اند را نشان می‌دهد. همانطور که ملاحظه می‌شود، خطای معیار متناظر با برآوردهای این دو پارامتر در روش BMA بصورتی قابل ملاحظه از مقادیر مشابه در روش گام به گام بزرگتر هستند. این یکی از برتریهای روش BMA است که میزان واقعی عدم حتمیت را نشان می‌دهد. مطلب دیگری که می‌توان به آن اشاره کرد این است

جدول ۳: برآورد ضرایب در روش BMA

| پارامتر رگرسیون | $Pr[\beta_i \neq 0 D]$ | میانگین پسین | انحراف معیار پسین |
|--------------------|--------------------------|--------------|----------------------|
| β_1 | ۰/۳۶۴ | ۰/۰۶۲۸ | ۰/۱۱۰۶ |
| β_2 | ۰/۴۴۶ | -۰/۰۷۴۸ | ۰/۱۰۷۸ |
| β_3 | ۱ | ۰/۸۰۷۶ | ۰/۱۲۵۸ |
| β_4 | ۰/۱۲۳ | ۰/۰۰۱۹ | ۰/۰۴۹۷ |
| β_5 | ۰/۱۵۱ | ۰/۰۱۰۳ | ۰/۰۴۸۲ |
| β_6 | ۰/۸۰۶ | ۰/۳۷۹۱ | ۰/۲۵۸۱ |

جدول ۴: برآورد ضرایب در مدل کامل ($R_{adj}^2 = ۰/۸۴۸۴$) و روش گام به گام

| مدل کامل | برآورد پارامترها | خطای معیار | روش گام به گام | برآورد پارامترها | خطای معیار |
|-------------|---------------------|---------------|-------------------|---------------------|---------------|
| β_0 | ۲۲/۹۱۳۵ | ۱۴/۱۶۶۶۵ | * | ۳۱/۹۶ | ۱۲/۶۴۸۵ |
| β_1 | ۰/۲۳۳۳ | ۰/۱۲۸۳ | --- | --- | --- |
| β_2 | -۰/۱۸۹۵ | ۰/۱۰۱۶ | --- | --- | --- |
| β_3 | ۰/۸۵۸۹ | ۰/۱۳۳۲ | * | ۰/۷۸۷۲ | ۰/۰۹۶۸ |
| β_4 | ۰/۰۷۲۱ | ۰/۱۴۷۰ | --- | --- | --- |
| β_5 | ۰/۰۲۴۰ | ۰/۱۱۰۱ | --- | --- | --- |
| β_6 | ۰/۴۱۷۷ | ۰/۱۹۷۰ | * | -۰/۴۴۶ | ۰/۲۰۲۷ |

که بخت حضور متغیر X_2 در مدل بیش از ۰/۸ است، که بر اساس قواعد سرانگشتی جفریز (۱۹۶۱) شواهد مثبتی برای حضور متغیر X_2 در مدل ارائه می‌کند، در حالی که اطلاعات موجود در این متغیر، با حذف آن در مدل حاصل از روش گام به گام نادیده گرفته می‌شوند. برآورد ضرایب متغیرهای در روش BMA در مقایسه با سایر روشها و خصوصاً در مورد متغیرهای X_4 و X_1 کاملاً کوچک و نزدیک صفر هستند. با این وجود کاملاً نادیده گرفته نمی‌شوند و از اطلاعات آنها استفاده می‌شود.

۴ بحث و نتیجه‌گیری

در هنگام مدلسازی غالباً نمی‌توان مدلی یافت که بصورت کامل به داده‌ها برازش داشته باشد. در این حالات انتخاب یک مدل به معنی از بین رفتن اطلاعات سایر مدلها است. میانگین‌گیری بیزی مدلها از اطلاعات همه مدلها، یا دسته‌ای از بهترین مدلها استفاده می‌کند و برخلاف روشهای مرسوم عدم حتمیت را به خوبی منعکس می‌کند. از نقطه نظر کاربری پیش‌بین نیز نتیجه حاصل از روش BMA از نتایج حاصل از هر یک از مدلهای موجود در فضای مدل مطلوب‌تر است.

عملکرد این روش با افزایش عدم حتمیت بهبود می‌یابد. مثلاً هر چه تعداد متغیرهای پیش‌بین در مدل رگرسیونی بیشتر باشد، کارایی این روش بیشتر می‌شود. در بدترین حالت که عدم حتمیت اندک باشد، بهتر بودن روش BMA نسبت به سایر روشها چشمگیر نمی‌باشد و بکارگیری روشهای مرسوم بدلیل سادگی مقرون به صرفه‌تر است.

مراجع

- [1] Good, I. J., (1950). "Probability and The Weighing of Evidence", Griffin London.
- [2] Hoeting, J., (1994). "Accounting for Uncertainty in Linear Regression Models", Ph.D Disertation, Department of Statistics, University of Washington.
- [3] Hoeting, J., Raftery, A. and Madigan, D., (1996). "A Method for Simultaneous Variable Selection and Outlier Identification in Linear Regression Models". *J. Comput, Statist.* **22**, 251-271.
- [4] Hoeting, J., Raftery, A. and Madigan, D., (1999). "Bayesian Simultaneous Variable Selection and Transformation Selection in Linear Regression Models". Technical Report 9905, Dept. Statistics, Colorado State Univ. Available at www.Colostate.edu.
- [5] Lipkovich, I, A., (2002). "Bayesian Model Averaging and Variable Selection in Multivariate Ecological Models", Ph.D Disertation, Faculty of The Virginia Polytechnic Institue and State University, Blacksburg, Virginia.
- [6] Madigan, D., Gavrin, J. and Raftery, A. E., (1995). "Eliciting Prior Information To Enhace The Preformance of Bayesian Graphical Models", *Comm. Statist. Thory Methods*, **24**, 2271-2292.
- [7] Madigan, D and Raftery, A., (1991). "Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam's Window", Technical Reports 213, Univ. Washington, Seattle.
- [8] Madigan, D. and Raftery, A., (1994). "Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam's Window", *Journal of American Statistical Association*, **89**, 1535-1546.
- [9] Miller A., (1990). "Subset Selection in Regression Variables", New York, Chpman-hal.
- [10] Raftery, A. E., (1995). "Bayesian Model Selection in Social Research(With Discussion)", in *Sociological Methodology 1995*(P. V. Marsden, ed.) 111-195. Blakwell. Cambridge, MA.

- [11] Raftery, A. E., (1996). "Approximate Bayes Factor and Accounting for Model Uncertainty in Generalized Linear Models", Technical Report, Biometrika, **83**, 351-266.
- [12] Robert, B. Nobel, J. (2000). "Multivariate applications of Bayesian Model averaging", Ph.D disertation, Faculty of the virginia polytechnic institue and state university, Blacksburg, virginia.

نمونه‌های گردان^۱، جایگزینی برای سرشماری‌های جمعیتی

زهره فلاح محسن‌خانی، فاطمه هرندی، فرشید جمشیدی

گروه پژوهشی طرح‌های فنی و روش‌های آماری پژوهشکده آمار

چکیده: در آمارگیری‌های نمونه‌ای متداول، معمولاً نمونه به یک جامعه مشخص و به یک مقطع زمانی خاص منتسب می‌شود. در سال‌های اخیر طرح نمونه‌گیری احتمالی با عنوان طرح نمونه‌گیری گردان معرفی شده است که برآوردهای حاصل از نمونه‌های این طرح نتایجی فراتر از نتایج آمارگیری‌های نمونه‌ای متداول ارائه می‌کنند تا آنجا که انجام این نمونه‌گیری در طی زمان بعنوان جایگزینی برای سرشماری‌ها معرفی شده است.

در یک طرح نمونه‌گیری گردان، k نمونه احتمالی غیر متداخل همزمان انتخاب می‌شوند. اندازه نمونه‌ها، یکسان و برابر با کسر ثابتی از کل جامعه می‌باشد. از هر نمونه در یک دوره زمانی محدود مصاحبه به عمل می‌آید و ضمن آن که اطلاعات آن دوره برآورد می‌شود، با متوسط‌گیری از اطلاعات دوره‌های زمانی مختلف، برآوردهایی برای کل دوره‌ها نیز ارائه می‌شود. سئوالی که در این روش مطرح است این است که با وجود اختلاف بین دوره‌ها چگونه می‌توان از نمونه‌های دوره‌های مختلف زمانی، متوسط‌گیری کرد و به نتایج حاصل اطمینان نمود؟

در این مقاله طرح نمونه‌گیری گردان بعنوان دیدگاهی جدید در تولید آمارهای رسمی بویژه بعنوان جایگزینی برای سرشماری‌های دهساله ارائه می‌شود.

واژه‌های کلیدی: نمونه‌گیری احتمالی، نمونه‌گیری ادواری، نمونه‌های گردان، نمونه احتمالی متداخل، نمونه احتمالی غیر متداخل

۱ مقدمه

نمونه‌های گردان به وسیله کیش^۲ در یک سری از مقالات شرح داده شده است. (کیش ۱۹۷۹، ۱۹۸۱، ۱۹۸۳، ۱۹۸۷، ۱۹۹۰، ۱۹۹۷، ۱۹۹۸، و کیش و ورما^۳ ۱۹۸۳ و ۱۹۸۶). کیش در مقاله‌های مذکور اصول تجمع^۴ اطلاعات در فضا و زمان را با استفاده از داده‌های حاصل از نمونه‌گیری گردان تشریح کرده است وی معتقد به استفاده از این روش برای اهداف گوناگون (کیش ۱۹۹۹)، به‌خصوص در کشورهای در حال توسعه (کیش ۱۹۷۹) و نیز در سرشماری ایالات متحده (کیش ۱۹۸۱) است. استفاده او از مفهوم نمونه‌های گردان، حداقل به سال ۱۹۵۸ تحت عنوان «نمونه‌گیری مستمر» (کیش، لاجوی^۵ و راکوف^۶، ۱۹۶۱) باز می‌گردد. وی همچنین پیشنهاد

1) Rolling Sample 2) Kish 3) Verma 4) Cumulating 5) Lovejoy
6) Rackow

نموده است نمونه‌های مذکور برای سرشماری‌های نمونه‌ای ملی توسط اداره سرشماری آمریکا استفاده شوند. (الکساندر^۷ ۲۰۰۰).

آمارگیری جامعه آمریکا^۸ (ACS)، که به عنوان جایگزینی برای سرشماری در نظر گرفته شده است، از یک نمونه‌گیری گردان استفاده می‌کند. این مقاله چگونگی به‌کارگیری مفهوم نمونه‌های گردان را با استفاده از ملاحظات عملیاتی و اهداف خاص در ACS شرح می‌دهد.

۲ نمونه‌های گردان

در یک طرح «نمونه گردان» همزمان k نمونه احتمالی غیر متداخل (پانل) انتخاب می‌شود که هر یک شامل $1/F$ کل جامعه می‌باشد. هر پانل در یک دوره زمانی، مورد مصاحبه قرار می‌گیرد تا اینکه همه نمونه‌ها بعد از k دوره مصاحبه شوند. بر حسب دقت مورد نیاز، اطلاعات یک پانل $1/F$ می‌تواند برای به دست آوردن برآوردهای خوبی از جامعه و چند محدوده وسیع کافی باشد. برای محدوده‌های کوچکتر یا برای حصول دقت بیشتر در محدوده‌های بزرگ، از تجمع چند پانل متوالی تا میزان k/F جمعیت، استفاده می‌شود. یک طرح نمونه گردان که در آن $k = F$ باشد، «سرشمای گردان^۹» نامیده می‌شود. برای یک نمونه گردان با دوره‌های ماهانه، طبیعی است که F مضربی از دوازده و تجمع‌ها، فصلی، نیمه‌ساله، سالانه و چند ساله می‌باشد.

۳ آمارگیری جامعه آمریکا (ACS)

سرشماری ده سال یکبار همراه با نمونه‌گیری مسکن با نرخهای $1/6$ تا $1/10$ ، منبع اصلی جمع‌آوری داده‌های ایالتی در زمینه مشخصه‌های جمعیتی و مسکن ایالات متحده است. همچنین آمارگیری جاری جمعیت (CPS) که آمارگیری نیروی کار ماهانه ایالات متحده است، با نرخ نمونه‌گیری یک از هزار، ماهانه انجام می‌شود، اما طرح آمارگیری آن بگونه‌ای است که نمونه نمی‌تواند به شکل قابل استفاده، آن چنان که در نمونه گردان مشاهده می‌شود در زمان تجمع پیدا کند (به دلیل حجم قابل توجه تداخل نمونه بین دوره‌ها). بنابراین دولت فدرال، در فاصله بین دو سرشماری، اطلاعات نسبتاً اندکی راجع به مشخصه‌های جمعیتی در سطح نواحی کشور ارائه می‌نماید. بعد از بحث‌ها و نشست‌های مفصل در فاصله زمانی سالهای ۱۹۹۱ تا ۲۰۰۲ نباشد در صورت تأمین منابع مالی توسط کنگره، آمارگیری جامعه آمریکا (ACS) توسط دفتر سرشماری در سال ۲۰۰۳ با نمونه ماهانه‌ای حدود ۲۵۰۰۰۰ آدرس، آغاز شود و در صورت موفقیت، جایگزین سرشماری ده سال یکبار گردد. ACS همان موضوعات سرشماری را با فراهم کردن جزئیات بیشتر در زمینه اجتماعی، اقتصادی و ویژگی‌های مسکن از کل ایالات متحده پوشش می‌دهد.

7) Alexander 8) American Community Survey 9) Rolling Census

ACS یک آمارگیری مبتنی بر نمونه گردان است (کیش ۱۹۹۸) و مطابق با یک نمونه گردان ماهانه با نرخ متوسط تقریباً $F = 48^\circ$ یا یک نمونه سالانه با $F = 40^\circ$ است. آمارگیری از $k = 60$ استفاده می‌کند که کوتاهترین تجمع قابل انتشار آن، برآوردهای سالانه است. نمونه‌گردان با مجموعه متفاوتی از آدرسها در هر ماه تا زمانی که در تمام جامعه گردش پیدا کند، تماس برقرار می‌کند. آدرسها یک نمونه سیستماتیک از آدرس شهروندان هستند که از فایلهای بهنگام شده آدرسهای اصلی استخراج می‌شوند که در تمام مناطق در هرماه تهیه می‌شود و در واقع بدون وجود آن اجرای طرح ACS امکان‌پذیر نمی‌باشد. نمونه سالانه تقریباً شامل ۳ درصد جامعه می‌شود که در طی یک دوره ۵ ساله به حدود ۱۵ درصد بالغ می‌گردد. ACS از طریق پست اجرا و بی‌پاسخی‌ها به وسیله تلفن پیگیری می‌شود. یک نمونه تصادفی شامل یک سوم بی‌پاسخی‌های باقیمانده، به صورت حضوری پیگیری می‌شود.

در محدوده‌هایی با نرخ پاسخگویی متوسط و F ماهانه برابر 48° ، خطای معیار برای یک برآورد متوسط ۵ ساله حاصل از ACS، قدری بیشتر از برآورد متناظر حاصل از نمونه‌گیری همراه سرشماری خواهد بود (حدوداً $1/33$ برابر بزرگتر). با توجه به امتیاز به موقع بودن و اینکه انتظار می‌رود نرخ داده‌های گم شده بدلیل استفاده از آمارگیران دائمی پائین باشد، پیش‌بینی می‌شود که در اکثر موارد برآوردها به قدر کافی نزدیک به واقعیت باشند. در نواحی با نرخهای پاسخ پائین‌تر از متوسط نرخ پاسخ زیر نمونه‌گیری برای پیگیری بی‌پاسخیها، اندازه نمونه مؤثر را کاهش خواهد داد. این امر نه تنها به دلیل کاهش تعداد مصاحبه‌ها بلکه به این دلیل است که وزنه‌های نامساوی نوعاً باعث بالا بردن اثر طرح 1° می‌شود (کیش ۱۹۶۵). برای جبران این مسئله، ACS از یک نرخ نمونه‌گیری بالاتر برای بی‌پاسخی در محدوده‌هایی با نرخ پاسخ پایین‌تر از حد متوسط استفاده می‌کند. [چارلز اچ الکساندر ۱۱ ۲۰۰۱]

۴ استفاده از متوسط‌ها در ACS

مسئله مهمی که در ارتباط با ACS مطرح می‌شود درک مسائلی است که استفاده از متوسط‌های چند ساله در موارد زیر ایجاد می‌کند:

- ۱) تفسیر و استفاده از متوسط‌های چند ساله برای توصیف و مقایسه نواحی
 - ۲) استفاده از متوسط‌های چند ساله برای تعیین مقدار تخصیص بودجه
 - ۳) استفاده از متوسط‌های چند ساله برای اندازه‌گیری تغییرات و روندها
- توصیه شده که برای تولید آمارهای توصیفی پایه برای نواحی کوچک، کاربران داده‌ها را برای چند سال (بسته به جمعیت نواحی)، انباشته کنند. این توصیه برای بسیاری از کاربران ACS خصوصاً آن دسته از کاربران که عادت به استفاده از داده‌های سرشماری دارند، مفید است. برای پیش‌بینی مقدار تخصیص بودجه نیز متوسط‌ها مفید هستند، گرچه بررسیهای اولیه، تأکید بیشتری

بر استفاده از متوسط‌های سه سالانه، حتی برای نواحی کوچک، دارند. برای اندازه‌گیری تغییرات در طول زمان، کاربر باید در کل سری‌های زمانی برآوردهای سالانه را تحلیل کند.

۱.۴ مدل‌هایی برای استفاده‌های توصیفی اساسی از داده‌های ACS

سه مدل اساسی وجود دارند که تحت آنها متوسط مقادیر سال‌های گذشته قابل استفاده است و کاربر داده‌ها بایستی استفاده و تفسیر متوسط‌ها را با توجه به یکی از این مدلها انجام دهد. این مدلها عبارتند از:

۱-- میانگین برای دوره خاص زمانی

۲-- مدل «کاربرد نوعی سرشماری^{۱۲}»

۳-- مدل «اغتشاش تصادفی^{۱۳}».

برای تشریح این سه مدل، هم از داده‌های سرشماری ۲۰۰۰ و هم از داده‌های ۲۰۰۲-۱۹۹۸ در سال ۲۰۰۳ استفاده می‌شود. فرض کنید X_t مقدار واقعی ناحیه‌ای خاص در سال t باشد، برای سال جاری، $t = ۲۰۰۳$ است.

به جای استفاده از مقدار واقعی، آنچه مشاهده می‌کنیم $\hat{X}_t = X_t + \epsilon_t$ است، که در آن ϵ_t خطای نمونه‌گیری است. در سرشماری، \hat{X}_t تنها هر ۱۰ سال یک بار مشاهده می‌شود اما دارای واریانس نمونه‌گیری کوچکتری در مقایسه با برآورد سالانه ACS می‌باشد.

۱.۱.۴ مدل ۱: میانگین برای دوره‌ی خاص زمانی

در این مدل هدف برآورد میانگین زیر است:

$$(X_{۱۹۹۸} + \dots + X_{۲۰۰۲})/۵$$

در این حالت، متوسط ACS برآوردگر آشکاری است و می‌تواند تحت یک رده از مدل‌های سری زمانی $\{X_t\}$ توجیه گردد.

گرچه این کاربرد مهمی است، اما غیر معمول می‌باشد و کاربرد معمول، بررسی «وضعیت کنونی» است که در مدل ۲ منعکس شده است.

۲.۱.۴ مدل ۲: مدل «کاربرد متداول سرشماری»

به‌عنوان یک مثال درباره‌ی استفاده‌ی «متداول» از داده‌های سرشماری، نوع نمایش‌های معمول در نقشه محدود‌های سرشماری را در نظر بگیرید که درصد مردم یا واحدهای مسکونی دارای مشخصه‌های خاص را در محدوده مورد نظر، با رنگ‌های مختلف نشان می‌دهد. چنانچه هدف نمایش محل تجمع نژادهای مختلف در شهرستان باشد، سؤال این است که در سال ۲۰۰۳

12) Typical Census Users 13) Random Noise

از نقشه‌های سرشماری ۲۰۰۰ استفاده شود، چه فرضیاتی راجع به تغییر طی زمان باید به‌کار گرفته شود و اگر به‌جای اطلاعات سرشمادی ۲۰۰۰ از متوسط اطلاعات سال‌های ۱۹۹۸-۲۰۰۲ استفاده شود تفسیرها چه تفاوتی می‌کنند؟ با توجه به اینکه موضوع مورد نظر «وضعیت کنونی» است، از برآوردهای سرشماری در مدل استفاده می‌شود با این فرض که این برآوردها، وضعیت کنونی را توصیف می‌کنند. هرچند ممکن است برای بسیاری از نواحی، مقادیر از سال ۲۰۰۰ کاهش یا افزایش یافته باشد اما هیچ تعدیل مشخصی برای تصحیح این تغییرات انجام نمی‌شود.

به بیان دیگر مدل ۲ عبارتست از:
 • فرض مجازی (مدل پیش‌گزینه)

$$X_{1998} = X_{1999} = \dots = X_{2003} = \mu$$

• گزینه ۱ (روند)

$$X_t = X_{1998} + C(t - 1998) \quad ; \quad C \neq 0$$

• گزینه ۲ (جهش ناگهانی)

$$\begin{aligned} X_t &= X_{1998} & ; & \quad t < I \\ X_t &= X_{1998} + C & ; & \quad t \geq I \end{aligned}$$

ACS دارای این مزیت است که با استفاده از سری‌های برآورد تک تک سالها، اطلاعاتی راجع به نواحی‌ای که دارای روند و جهش هستند، ارائه می‌دهد.

۳.۱.۴ مدل ۳: «اغتشاش تصادفی»

گاهی تغییرات بی‌قاعده‌تر از «روندها» و «جهش‌های ناگهانی» هستند. فرض کنید که مقدار واقعی برای ناحیه مورد بررسی به صورت زیر باشد

$$\begin{aligned} X_t &= \mu + \eta_t \\ E(\eta_t) &= 0 \end{aligned}$$

و η_t ها غیر همبسته باشند.

می‌توان یک شهر کوچک شامل ۲۰ واحد مسکونی را در نظر گرفت که در سال اول ۳ خانوار، سال دوم ۶ خانوار، سال سوم ۴ خانوار و... در آن فقیر باشند. در این حالت، متوسط ۵ ساله به عنوان برآورد μ و η به عنوان «اغتشاش» غیر دلخواه تعبیر می‌شود که مربوط به تغییرات مقدار واقعی در جامعه است.

۲.۴ استفاده از متوسط‌های چند سالانه برای اندازه‌گیری مقدار تخصیص بودجه

بسیاری از فرمولها که برای تخصیص بودجه برنامه‌های دولت فدرال بکار برده می‌شود از داده‌های آخرین سرشماری استفاده می‌کنند. به این منظور از رویکردهای گوناگونی استفاده می‌شود که در این جا مدل ساده زیر ارائه شده است. نیاز واقعی به بودجه در سال جاری بوسیله متغیر X_t برای ناحیه جغرافیایی خاص اندازه‌گیری می‌شود. از داده‌های در دسترس برای ارزیابی نیاز جاری به صورت زیر استفاده می‌شود:

$$\hat{A}_t = f(\dots, \hat{X}_{t-i}, \dots, \hat{X}_{t-1})$$

که در آن f یک تابع و \hat{X}_{t-i} برآورد نمونه در سال $t - i$ می‌باشد. فرض شده است که برآورد سال جاری \hat{X}_t در زمانی که به آن نیاز است در دسترس نیست. توزیع سری زمانی $\{X_t\}$ برای هر ناحیه نامعلوم است و ممکن است دارای الگوهای متفاوتی در مناطق مختلف و برای متغیرهای مختلف باشد.

اگر توزیع سری زمانی $\{X_t\}$ را بتوان برای ناحیه و متغیر مورد نظر تعیین کرد، تابع پیش‌بینی بهینه f تعیین می‌گردد. گرچه مقدار بهینه f ممکن است بسته به ناحیه تغییر کند. مقدار متوسط می‌تواند به عنوان یک پیش‌بینی کننده «توافقی»^{۱۴} برای همه ناحیه‌ها استفاده شود. گزینه‌های زیر را می‌توان برای برآورد نیاز جاری به کمک متوسط‌های ACS در نظر گرفت.

$$\begin{aligned} \hat{A}_t &= \hat{X}_{t-1} && \text{۱ ساله:} \\ \hat{A}_t &= 1/3(\hat{X}_{t-3} + \hat{X}_{t-2} + \hat{X}_{t-1}) && \text{۳ ساله:} \\ \hat{A}_t &= 1/5(\hat{X}_{t-5} + \dots + \hat{X}_{t-1}) && \text{۵ ساله:} \end{aligned}$$

اما در صورت عدم استفاده از متوسط‌های ACS و استفاده از سرشماری قبلی، $\hat{A}_t = \hat{X}_0$ ، $t = 1, \dots, 11$ که در آن $t = 0$ مشخص کننده سال آخرین سرشماری است. بسیاری پیشنهاد کرده‌اند که بهتر است بجای میانگین ساده از میانگین وزنی با وزنهای بزرگ‌تر برای سالهای اخیر استفاده شود. (کارلرز^{۱۵} ۲۰۰۱)

۳.۴ اندازه‌گیری تغییرات در طی زمان

ACS دارای کاربردهای بالقوه‌ای برای اندازه‌گیری تغییرات در طول زمان است. در صورتی که بدون داشتن نمونه‌های گردان امکان اندازه‌گیری این تغییرات وجود ندارد. در نواحی کوچک با نمونه‌های غیرمتداخل، احتیاط در تعبیر نتایج به دلیل انحراف استاندارد زیاد، لازم است. همچنین بدیهی است که در اندازه‌گیری تغییرات در طی زمان، تغییرات ناشی از مسایلی همچون تورم (در متغیر درآمد)، و گذر زمان (در متغیر سن) و ... باید از تغییرات واقعی جدا شوند.

۵ مزایا و معایب سرشمای گردان در مقایسه با شیوه سنتی سرشماری جمعیتی

همان گونه که قبلاً توضیح داده شد یک طرح نمونه‌گیری گردان که در آن $k = F$ است، سرشماری گردان نامیده می‌شود. در این قسمت مزایا و معایب این نوع سرشماری در مقایسه با سرشماری سنتی جمعیت تشریح می‌شود. مزایا و معایب سرشماری گردان نسبت به شیوه سنتی سرشماری جمعیتی را از جنبه‌های: خروجی‌های سرشماری، پذیرش دیدگاه، کار، استمرار عملیات، پوشش، قوانین و هزینه‌ها می‌توان بررسی نمود.

۱.۵ خروجی‌ها

در سرشماری گردان، بدلیل وجود اطلاعات دوره‌های زمانی متوالی، کیفیت آمارهای جمعیتی بهبود می‌یابد. عیب اصلی خروجی‌ها این است که نمی‌تواند تصویری لحظه‌ای از کل جامعه ارائه کند. دو مزیت مهم خروجی این نوع سرشماری را به صورت ذیل می‌توان بیان کرد:

- ۱- کیفیت برتر برآوردها در مقایسه با برآوردهای حاصل از طرحهای آمارگیری معمول در فواصل بین سرشماری‌ها
 - ۲- سطح مناسب برای مدل‌بندی داده‌ها به نحوی که با برآوردهای کل جامعه سازگار شود.
- مزیت دیگر خروجی‌های این سرشماری، امکان پذیر بودن تغییر محتوی سئوالات در طی دوره است که موجب انطباق بهتر نتایج آن با نیازهای در حال تغییر کاربران می‌باشد.

۲.۵ پذیرش دیدگاه کاربر

بدلیل امکان‌پذیر بودن تغییر محتوای سئوالات در طی دوره، دولت‌های محلی از جمله مخالفان اصلی این نوع سرشماری هستند زیرا معتقدند که با این روش، مقایسه پذیری بین نواحی از بین می‌رود. برای جلب موافقت آنها لازم است آگاهی آنها از این روش را بالا برد.

۳.۵ استمرار عملیات

با توجه به تجربه آمارگیری جاری جمعیت آمریکا^{۱۶} (CPS) استمرار عملیات، امکان استفاده از نیروی میدانی ثابت یا نیمه ثابت را فراهم می‌آورد که موجب انباشتگی تجربه و توان کارشناسی و در نتیجه بهبود کیفیت اجرا می‌شود.

16) Current Population Survey

۴.۵ پوشش

پوشش در سرشماری‌های گردان با مشکلات متفاوتی به صورت ذیل مواجه است:

- ۱-- شمول بیش از یکبار خانوارها در دوره ۱۰ ساله بدلیل عدم ایستایی جمعیت
- ۲-- عدم شمول بعضی از خانوارها در دوره ۱۰ ساله بدلیل عدم ایستایی جمعیت
- ۳-- پوشش مضاعف یا عدم پوشش خانوارهای دارای بیش از یک اقامتگاه
- ۴-- مشکلات پوشش ناشی از کم توجهی عمومی به این نوع سرشماری در مقایسه با سرشماری سنتی که به دلیل منحصر به فرد بودن آن یک واقعه‌ی مهم ملی تلقی می‌شود.

۵.۵ قوانین

این نوع سرشماری ممکن است نیازمند وضع قوانین جدید یا اصلاح قوانین موجود در زمینه فواصل زمانی اجرا، پوشش موضوعی و جغرافیایی و ... باشد.

۶.۵ هزینه

با توجه به بررسی‌های اولیه انجام شده، سرشماری گردان پرهزینه‌تر از سرشماری سنتی است (حدود ۱۰ درصد بیشتر) البته در طولانی مدت و با استمرار عملیات، صرفه اقتصادی بخصوص با تغییر شیوه‌های اجرایی (از جمله آمارگیری به روش پستی) قابل مشاهده خواهد بود.

با توجه به مطالب فوق مزایا و معایب سرشماری گردان را می‌توان به شرح زیر برشمرد.

مزایا

- وجود اطلاعات مستمر
- انعطاف پذیری سرفصل‌های اطلاعات
- کارایی ناشی از استمرار عملیات
- استفاده از کارکنان میدانی ثابت و مکان ثابت برای پردازش داده‌ها

معایب

- از دست رفتن تصویر لحظه‌ای جمعیت
- تضعیف اهمیت سرشماری از دیدگاه جامعه
- افزایش خطر تأثیر سوّ جریان‌ات خارجی بر کار میدانی
- نیاز به متقاعد نمودن کاربر و وضع قوانین جدید
- داشتن هزینه بیش از سرشماری‌های سنتی

۶ نتیجه‌گیری

با توجه به مزیت‌هایی که این شیوه نسبت به سرشماری سنتی جمعیت داراست و با توجه به این که تجربیات ناشی از اجرای چندسال این شیوه با تغییر شیوه‌های اجرایی می‌تواند بسیاری از معایب

این روش را کاهش داده یا برطرف نماید، انتظار می‌رود این روش بتواند جایگزینی مناسب برای سرشماری سنتی جمعیت باشد.

مراجع

- [1] Alexander, C. H. (1998), "Recent Developments In the American Community Survey", Proceedings of the Survey Research Methods Section, American Statistical Association, pp. 92-100.
- [2] Alexander, C. H. and Wetrogan, S (2000), "Integrating the American Community Survey and the Intercensal Demographic Estimates Program.", Integrating the American Community Survey and the Intercensal Demographic, American Statistical Association, pp. 295-300.
- [3] Alexander, C. H. (2001), "Still Rolling: Leslie Kish 'Rolling Samples' And The American Community Survey", Estimates Program.", Proceedings of Statistics Canada Symposium, pp. 1-10.
- [4] Kish, L., Lovejoy, W., and Rackow, P. (1961), "A Multi-Stage Probability Sample for Sample for Traffic Surveys", Proceedings of the Social Statistics Section, American Statistical Association, pp. 227-230
- [5] Kish, L. (1965), Survey Sampling .John Wiley and Sons, New York.
- [6] Kish, L.(1979), "Rolling Samples Instead of Census", Asia and Pacific Census Forum, G(1), August 1979, pp.12-13.
- [7] Kish, L.(1981), "Using Cumulated Rolling Samples to Integrate Census and Survey Operations of the Census Bureau", Washington, D.C., U.S. Government Printing Office.
- [8] Kish, L.(1983), "Data Collection for Details over Space and Time", In T. Wright, (ed.), Statistical Methods and the Improvement of Data Quality, New York: Academic, pp. 72-84.
- [9] Kish, L., and Verma, V (1983), Census Plus Samples: Combined Uses and Designs, Bulletin of the International Statistical Institute, pp. 66-82.
- [10] Kish, L., and Verma, V., (1986), "Complete Census and Samples", Journal of Official Statistics, 2, pp. 381-93. the Census Bureau", Washington, D.C., U.S. Government Printing Office.
- [11] Kish, L.(1987), "Statistical Research Design, John Wiley and Sons, New York.
- [12] Kish, L.(1990), "Rolling Samples and Censuses", Survey Methodology, 16, 63-79.

- [13] Kish, L.(1997),”Periodic and Rolling Samples and Census”, Chapter 7 In Statistics and Public Policy, Bruce D. Spencer, ed., Clarendon Press, Oxford.
- [14] Kish, L.(1998),”Space/Time Variations and Rolling Samples”,Journal of Official Statistics, 14, 1, 1998, pp. 31-46.

فواصل پیش‌بینی برای مدل‌های اتورگرسیو خطی با استفاده از تکنیکهای بوت استرپ

ملیحه قابل^۱، صادق رضایی^۲

^۱ دانشجوی کارشناسی ارشد دانشگاه صنعتی امیرکبیر
^۲ گروه آمار دانشگاه صنعتی امیرکبیر

چکیده: در این مقاله ابتدا به محاسبه فواصل پیش‌بینی برای مدل‌های اتورگرسیو خطی می‌پردازیم و روشهای بهبود این فواصل را بررسی می‌کنیم. سپس تکنیکهای بوت استرپ را برای برآورد پارامتر نامعین و کاهش اریبی این برآورد در نمونه‌های کوچک بکار می‌بریم. همچنین تصحیح اریبی پارامتر برآورد شده را برای محاسبه فواصل پیش‌بینی بکار گرفته و در پایان بوسیله شبیه‌سازی مونت کارلو برای مدل $AR(1)$ فواصل اطمینان معرفی شده را با سطوح پوشش متفاوت مقایسه می‌کنیم.

واژه‌های کلیدی: فواصل پیش‌بینی، اتورگرسیو، بوت استرپ، تصحیح اریبی

۱ مقدمه

فاصله پیش‌بینی یک حد پایین و بالایی است با این احتمال که فرآیند تحقق یافته بعدی باید بین این دو حد قرار بگیرد. یک فاصله باید روی نمونه مشاهده شده تحقق یافته شرطی شود. بعضی از موضوعات وابسته به هم، در ارتباط با محاسبه فواصل پیش‌گویی با سطوح پوشش متناسب می‌باشند. که شامل برآورد پارامتر نامعین (یعنی پارامترهای مدل برآورد شده متغیر تصادفی‌اند) اریبی نمونه کوچک از پارامتر برآورد شده، پیش‌بینی توابع غیر خطی از پارامتر برآورد شده و . . . نشان خواهیم داد که این بررسی‌ها با تکنیکهای بوت استرپ متناسب می‌توانند در ارتباط باشند. روش فاصله اطمینان BOX-JENKINS که توسط Jenkins و Box در سال ۱۹۷۶ اعمال شده شاید ساده‌ترین روش باشد. که در بخش دوم این مقاله به معرفی این روش پرداخته‌ایم. سپس فاصله اطمینان چندک Efron و فاصله اطمینان چندک Hall و فواصل t -چندک را معرفی کرده و این فواصل را در مدل‌های اتورگرسیو بکار می‌بریم. در سال ۲۰۰۱، Romo, Ruiz و Pascual شیوه‌ای را اعمال کردند که با وجود نمایش رو به عقب فرآیند، نیازی به دوباره نمونه‌گیری نداشت. بنابراین این شیوه برای مدل‌هایی با مولفه‌های میانگین متحرک به سرعت بکار گرفته شد. در پایان بخش دوم فواصل اطمینان PRR را معرفی می‌کنیم. بخش سوم این مقاله شامل تجویز برآورد نامعین و تصحیح اریبی پارامتر برآورد شده می‌باشد. و بخش چهارم شامل شبیه‌سازی

مونته کارلو برای فواصل پیش‌بینی PRR با تصحیح اربیبی پارامتر برآورد شده می‌باشد و در پایان تخمین مونته کارلو را برای مدل AR(1) شبیه‌سازی کرده و با سطوح پوشش متفاوت آنها را مقایسه می‌کنیم.

۲ محاسبه فواصل پیش‌بینی

مدل AR(p) را در نظر بگیرید.

$$Y_t = \delta + \phi_1 Y_{t-1} + \dots + \phi_p Y_{t-p} + a_t \quad (1)$$

که a_t یک دنباله iid با تابع توزیع F_a ، $E[a_t] = 0$ ، $E[a_t^2] = \sigma_a^2 < \infty$ ، $E[a_t^3] = 0$ ، $E[a_t^4] = \sigma_a^4$ می‌باشد. میانگین مربع خطا (MMSE) پیش‌گویی k مرحله بعدی Y_{t+k} وقتی پارامترهای مدل معلوم هستند بصورت:

$$\tilde{y}_{T+k} = \delta + \phi_1 \tilde{y}_{T+k-1} + \dots + \phi_p \tilde{y}_{T+k-p} \quad (2)$$

می‌باشد که $\tilde{y}_{T+i} = y_{T+i}$ برای $i \leq 0$ که y_t یک مقدار تحقق یافته از Y_T است. فاصله پیش‌بینی استاندارد Box-jenkins (BJ) در سطح $(1 - \alpha)$ بصورت زیر تعریف می‌شود:

$$\{\tilde{y}_{T+k} + z_{\frac{\alpha}{2}} \hat{\sigma}_k, \tilde{y}_{T+k} - z_{1-\frac{\alpha}{2}} \hat{\sigma}_k\} \quad (3)$$

که $\hat{\sigma}_k$ برآورد واریانس خطای پیش‌بینی k -مرحله بعدی است و $\gamma = \phi(Z_\gamma)$ که ϕ تابع توزیع نرمال استاندارد است. بنابراین روش BJ یک فاصله را روی نمونه مشاهده شده تحقق یافته شرطی می‌کند. حال در مدل بالا اگر $\{\delta, \phi_1, \dots, \phi_p, \sigma_a^2\}$ معلوم باشند و F_a تابع توزیع گوسین باشد آنگاه چگالی پیش‌گویی یک مرحله بعدی در $T+1$ به شرط اطلاعات قابل دسترس تا T بصورت $(\tilde{y}_{T+1}, \sigma_a^2) | I_T = N(\tilde{y}_{T+1}, \sigma_a^2)$ می‌باشد که

$$I_T = [Y_T = y_t, \dots, Y_{T-p+1} = y_{T-p+1}], \tilde{y}_{T+1} = \delta + \phi_1 y_T + \dots + \phi_p y_{T-p+1}$$

در این حالت فاصله پیش‌بینی $(1 - \alpha)$ ، BJ بصورت: $\{\tilde{y}_{T+1} \pm Z_{\frac{\alpha}{2}} \cdot \sigma_a\}$ می‌باشد.

$$P_r[Y_{T+1} \in \tilde{y}_{T+1} \pm Z_{\frac{\alpha}{2}} \cdot \sigma_a] = 1 - \alpha$$

$$P_r[Z_{\frac{\alpha}{2}} < \frac{Y_{T+1} - \tilde{y}_{T+1}}{\sigma_a} < Z_{\frac{\alpha}{2}}] = 1 - 2\phi(Z_{\frac{\alpha}{2}}) = 1 - \alpha$$

فواصل چند مرحله‌ای به روش مشابه بدست می‌آیند.

۱.۲ پیش‌گویی‌های بوت استرپ

چگونگی یک روش بوت استرپ استاندارد را بررسی می‌کنیم. مدل (۱) را برای $t = p + 1, \dots, T$ با داشتن مقادیر $\hat{\sigma}_a^2, \hat{\phi}_1, \dots, \hat{\phi}_p, \hat{\delta}$ و باقیمانده‌های \hat{a}_t برآورد می‌کنیم. اگر باقیمانده‌ها دارای میانگین صفر نباشند آنها را مرکزی کرده و در صورت امکان با عامل $\sqrt{(T-p)/(T-2p)}$ دوباره مقیاس‌بندی می‌کنیم. تابع توزیع تجربی باقیمانده‌های مرکزی شده و مقیاس‌بندی شده را با $F_{\hat{a}}$ نمایش می‌دهیم. تابع توزیع تجربی برای هر مقدار $\{\hat{a}_t\}$ و مقادیر $t = p + 1, \dots, T$ برابر با $\frac{1}{T-p}$ است. بطوریکه مستقلاً از $F_{\hat{a}}$ بصورت با جایگذاری متناظر با نمونه‌گیری خطاها استخراج شده‌اند. یک نمونه بوت استرپ توسط نمونه‌گیری از $F_{\hat{a}}$ تولید کرده که باقیمانده‌ها بصورت $\{a_{b,t}^*\}$ برای $t = T + 1, \dots, T + k$ مشخص می‌شوند. معادله اتورگرسیو در این حالت بصورت زیر می‌باشد:

$$y_{b,T+k}^* = \hat{\delta} + \hat{\phi}_1 y_{b,T+k-1}^* + \dots + \hat{\phi}_p y_{b,T+k-p}^* + a_{b,T+k}^* \quad (4)$$

که $y_{b,T+s}^* = y_{T+s}$ برای $s \leq 0$. با تکرار این عملیات B مرتبه یک نمونه بوت استرپ y_b^* برای $b = 1, \dots, B$ خواهیم داشت که $y_b^* = [y_{b,T+1}^*, \dots, y_{b,T+k}^*]'$.

۲.۲ فاصله اطمینان چندک

افرون (Efron 1979) روش فاصله اطمینان چندک را بصورت زیر محاسبه کرد. فرض کنید یک مجموعه داده X^* طبق رابطه $\hat{P} \rightarrow X^*$ را تولید کرده‌ایم. تکرارهای بوت استرپ $\hat{\theta}^* = S(X^*)$ محاسبه شده‌اند. اگر \hat{G} تابع توزیع تجمعی $\hat{\theta}^*$ باشد، فاصله چندک $1 - 2\alpha$ توسط چندکهای α و $1 - \alpha$ از \hat{G} بصورت زیر تعریف می‌شود:

$$[\hat{\theta}_{Lo}, \hat{\theta}_{Up}] = [\hat{G}^{-1}(\alpha), \hat{G}^{-1}(1 - \alpha)] \quad (5)$$

چون طبق تعریف $\hat{G}^{-1}(\alpha) = \hat{\theta}_{\beta}^{*(\alpha)}$ و $\hat{G}^{-1}(1 - \alpha) = \hat{\theta}_{\beta}^{*(1-\alpha)}$ امین چندک تجربی از مقادیر $\hat{\theta}^*(b)$ است. یعنی مقدار $\beta \cdot \alpha$ ، آنگاه فاصله چندک $(1 - 2\alpha)$ بصورت تقریبی برابر است با:

$$[\hat{\theta}_{Lo}, \hat{\theta}_{Up}] \approx [\hat{\theta}_{\beta}^{*(\alpha)}, \hat{\theta}_{\beta}^{*(1-\alpha)}] \quad (6)$$

بخاطر محبوبیت روش چندک افرون، مختصراً روش چندک Hall را معرفی می‌کنیم. Hall(1988)، هفت روش را برای محاسبه فواصل پیش‌بینی بررسی کرد. Hall روش چندک دیگری را معرفی کرد و توزیع تجربی چندکها را برای تعریف نقاط بحرانی بالا و پایین فواصل پیش‌بینی بکار گرفت. فرض کنید F_k^* تابع توزیع تجربی از $\{y_{b,T+k}^*, b=1, \dots, B\}$ باشد آنگاه فاصله بصورت

$[F_k^{*-1}(\frac{\alpha}{\gamma}), F_k^{*-1}(\frac{1-\alpha}{\gamma})]$ می‌باشد. به عنوان مثال یک فاصله با پوشش ۹۵ درصد و $B = 1000$ تکرار بصورت زیر می‌باشد.

$$\alpha = 0.05, F_k^{*-1}(0.025) = L_k = y_{T+k}^{*(25)} F_k^{*-1}(0.975) = U_k = y_{T+k}^{*(975)}$$

که مقادیر L_k, U_k متناسب با آماره‌های ترتیبی از توزیع تجربی می‌باشند. بنابراین فاصله چندک Hall بصورت زیر محاسبه می‌شود:

$$[\hat{y}_{T+k} - t_{1-\frac{\alpha}{\gamma}}, \hat{y}_{T+k} - t_{\frac{\alpha}{\gamma}}] \quad (7)$$

که \hat{y}_{T+k} پیش‌بینی k مرحله بعدی بر اساس برآوردهای پارامتر نمونه y_T, \dots, y_{T-p+1} است یعنی:

$$\hat{y}_{T+k} = \hat{\delta} + \hat{\phi}_1 \hat{y}_{T+k-1} + \dots + \hat{\phi}_p \hat{y}_{T+k-p} \quad (8)$$

که $\hat{y}_{T+s} = y_{T+s}$ برای $s \leq 0$ و $t_{1-\frac{\alpha}{\gamma}}$ و $t_{\frac{\alpha}{\gamma}}$ چندکهای توزیع بوت استرپ از $Y_{T+k}^* - \hat{y}_{T+k}$ است. مقادیر تحقق یافته Y_{T+k}^* بصورت $\{y_{b,T+k}^*, b = 1, \dots, B\}$ است. محاسبه فواصل پیش‌بینی و تصحیح آریبی برآوردهای پارامتر با حل معادله زیر می‌تواند بررسی شود.

$$E[f_t(F_0, F_1) | F_0] = 0 \quad (9)$$

که F_0 تابع توزیع جامعه و F_1 تابع توزیع تجربی از نمونه $\{y_1, \dots, y_T\}$ است. χ تابعی از کلاس $\{f_t, t \in \tau\}$ می‌باشد که $t = \{t_1, t_2\}$ جوابهای معادله‌اند. امید نسبت به جامعه F_0 شرطی شده است. برای محاسبه فواصل اطمینان $(1 - \alpha)$ درصد برای آماره $(F_0, \theta) = \theta(F_0)$ می‌تواند بصورت زیر باشد.

$$f_t(F_0, F_1) = I\{\theta(F_1) - t_1 \leq \theta(F_0) \leq \theta(F_1) + t_2\} - (1 - \alpha)$$

که $I\{A\}$ تابع نشانگر است و مقدار ۱ را می‌گیرد اگر پیشامد A رخ دهد و در غیر اینصورت برابر صفر است. $\theta_0 = Y_{T+k}$ و $\theta_1 = \theta(F_1)$ برآوردی از پیش‌بینی k -مرحله بعدی بر اساس نمونه داده‌ها بصورت \hat{y}_{T+k} است. بنابراین فاصله اطمینان بصورت $\{\theta_1 - t_1, \theta_1 + t_2\}$ یا $\{\hat{y}_{T+k} - t_1, \hat{y}_{T+k} + t_2\}$ می‌باشد. به جای حل معادله جامعه که شامل مجهولهای $F_0, \theta(F_0), t$ می‌باشد، معادله نمونه را حل می‌کنیم. یعنی:

$$E[f_t(F_1, F_2) | F_1] = 0 \quad (10)$$

که با جایگذاری F_0 بوسیله F_1 و F_1 بوسیله F_2 بدست آمده که F_2 تابع توزیع تجربی از داده‌های χ که دوباره نمونه‌گیری شده‌اند، می‌باشد که جواب معادله بالا $\hat{t} = \{\hat{t}_1, \hat{t}_2\}$ است، که در این حالت تابع f_t بصورت:

$$f_{\hat{t}}(F_1, F_2) = I\{\theta(F_2) - \hat{t}_1 \leq \theta(F_1) \leq \theta(F_2) + \hat{t}_2\} - (1 - \alpha)$$

و فاصله اطمینان بصورت $\{\theta_1 - \hat{t}_1, \theta_1 + \hat{t}_2\}$ می‌باشد. که جواب معادله نمونه باید تقریب خوبی برای معادله جامعه باشد. که Hall آن را به عنوان قاعده بوت استرپ معرفی می‌کند. که در حالت پیش‌بینی معادله زیر را حل می‌کنیم.

$$P(\{Y_{T+k}^* - \hat{t}_1 \leq \hat{y}_{T+k} \leq Y_{T+k}^* + \hat{t}_2\}) = 1 - \alpha$$

برای یک فاصله با دمه‌ای برابر داریم:

$$P_r(Y_{T+k}^* - \hat{y}_{T+k} \leq \hat{t}_1) = 1 - \frac{\alpha}{2}$$

$$P_r(Y_{T+k}^* - \hat{y}_{T+k} \leq -\hat{t}_2) = \frac{\alpha}{2}$$

بنابراین \hat{t}_1 چندک $(1 - \frac{\alpha}{2})$ از توزیع بوت استرپ $Y_{T+k}^* - \hat{y}_{T+k}$ و $-\hat{t}_2$ چندک $\frac{\alpha}{2}$ ام می‌باشد.

۳.۲ فواصل t - چندک

علاوه بر فواصل چندک Hall می‌توان فواصل t - چندک را نیز محاسبه کرد. فاصله اطمینان t - چندک (Efron, Tibshirani 1993) بصورت $(\hat{\theta} - \hat{t}^{1-\alpha} \cdot \hat{S}e, \hat{\theta} + \hat{t}^\alpha \cdot \hat{S}e)$ می‌باشد. که فاصله t - چندک در مدل اتورگرسیو بصورت زیر محاسبه می‌شود.

$$[\hat{y}_{T+k} - t_{1-\frac{\alpha}{2}} \cdot \hat{\sigma}(\hat{y}_{T+k}), \hat{y}_{T+k} + t_{\frac{\alpha}{2}} \hat{\sigma}(\hat{y}_{T+k})] \quad (11)$$

که $t_{\frac{\alpha}{2}}, t_{1-\frac{\alpha}{2}}$ در اینجا چندکهای بوت استرپ از $\frac{Y_{T+k}^* - \hat{y}_{T+k}}{\hat{\sigma}(Y_{T+k}^*)}$ می‌باشند. یک ایده اینست که بوت استرپ استودنت شده $\frac{\hat{y}_{T+k} - Y_{T+k}}{\hat{\sigma}(\hat{y}_{T+k})}$ بهتر از $(\hat{y}_{T+k} - Y_{T+k})$ است. هر چند فواصل اطمینان بر اساس آماره‌های محوری، ممکن است بطور میجانبی دقیق‌تر از آماره‌های غیر محوری نباشند. حال بررسی می‌کنیم که آیا استودنت کردن دقت فواصل پیش‌بینی را بهبود می‌بخشد؟ از رابطه (۸) پیش‌بینی k -مرحله بعدی برای $p = 1$ می‌تواند بصورت زیر نوشته شود.

$$\hat{y}_{T+k} = \sum_{i=0}^{k-1} \hat{\phi}^i \hat{\delta} + \hat{\phi}^k y_T$$

که $\hat{y}_{T+s} = y_{T+s}$ برای $S \leq 0$. و خطای استاندارد پیش‌بینی بصورت

$$\sigma(\hat{y}_{T+k}) = \sqrt{E[(\hat{y}_{T+k} - E(\hat{y}_{T+k}))^2]}$$

می‌باشد. با استفاده از تقریب $E[\hat{y}_{T+k}] \approx \sum_{i=0}^{k-1} \phi^i \delta + \phi^k y_T$ خواهیم داشت
 در (Clement, Hendry) $\hat{y}_{T+k} - E(\hat{y}_{T+k}) = \sum_{i=0}^{k-1} (\hat{\phi}^i \delta - \phi^i \delta) + (\hat{\phi}^k - \phi^k) y_T$
 سال ۱۹۹۸ رابطه زیر را اثبات کردند.

$$\sigma^2(\hat{y}_{T+k}) = \sigma_a^2 \left(\frac{1 - \phi^k}{1 - \phi} \right)^2 T^{-1} + k^2 \phi^{2(k-1)} (1 - \phi^2) y_T^2 T^{-1} \quad (12)$$

برای فرآیندهای ایستا و $\delta = 0$ یک عرض از مبدا برآورد شده است. برای ساختن این عملگر
 σ_a^2 را با یک برآورد نااریب و $\hat{\phi}$ را با جایگزین می‌کنیم.

$$\hat{\sigma}^2(\hat{y}_{T+k}) = \hat{\sigma}_a^2 \left(\frac{1 - \hat{\phi}^k}{1 - \hat{\phi}} \right)^2 T^{-1} + k^2 \hat{\phi}^{2(k-1)} (1 - \hat{\phi}^2) y_T^2 T^{-1} \quad (13)$$

$\hat{\sigma}(Y_{T+k}^*)$ بطور مشابه برآورد می‌شود بطوری که برای نمونه b بوت استرپ داریم:

$$\hat{\sigma}^2(y_{b,T+k}^*) = \hat{\sigma}_{b,a}^{2*} \left(\frac{1 - (\phi_b^*)^k}{1 - \phi_b^*} \right)^2 T^{-1} + k^2 (\phi_b^*)^{2(k-1)} (1 - \phi_b^{*2}) y_T^2 T^{-1} \quad (14)$$

که ϕ_b^* برآورد پارامتر بر اساس داده‌های بوت استرپ (دوباره نمونه‌گیری) و $\sigma_{b,a}^{2*}$ واریانس خطای
 برآورد شده بر اساس داده‌های بوت استرپ می‌باشد. اثر برآورد پارامتر را روی شکل چگالی‌های
 پیش‌بینی با استفاده از روش بوت استرپ که توسط Pascual, Romo, Ruiz در سال ۱۹۸۸
 برای مدل‌های ARIMA(p,d,q) اعمال شده را به نام فواصل اطمینان PRR معرفی می‌کنیم.
 ابتدا پارامترها توسط روش شبه درست‌نمایی ماکزیمم شرطی (QML) برآورد می‌شوند. روش
 بوت استرپ که توسط Pascual در سال ۱۹۸۸ برای ساخت فواصل پیش‌گویی برای مقادیر
 بعدی از سری تولید شده در فرآیند ARIMA(p,d,q) در نظر می‌گیریم.

$$\nabla^d y_t = \phi_0 + \phi_1 \nabla^d y_{t-1} + \dots + \phi_p \nabla^d y_{t-p} + a_t + \theta_1 a_{t-1} + \dots + \theta_q a_{t-q} \quad (15)$$

که a_t یک فرآیند اغتشاش خالص است و ∇ عملگر تفاضلی بصورت $\nabla y_t = y_t - y_{t-1}$
 که $(\phi_0, \phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q)$ پارامترهای مجهول هستند. از سری مشاهده شده
 $\{y_1, \dots, y_T\}$ پارامترها توسط یک برآورد سازگار مانند QML شرطی می‌توانند برآورد شوند.
 با داشتن مقادیر $(\hat{\phi}_0, \hat{\phi}_1, \dots, \hat{\phi}_p, \hat{\theta}_1, \dots, \hat{\theta}_q)$ باقیمانده‌ها بصورت زیر محاسبه می‌شوند.

$$\hat{a}_t = \nabla^d y_t - \hat{\phi}_0 - \hat{\phi}_1 \nabla^d y_{t-1} - \dots - \hat{\phi}_p \nabla^d y_{t-p} - \hat{\theta}_1 \hat{a}_{t-1} - \dots - \hat{\theta}_q \hat{a}_{t-q} \quad (16)$$

که باقیمانده‌ها متناظر با پر یودهایی از زمان $t = 0, -1, -2, \dots$ مجموعه‌ای برابر صفر هستند.
 \hat{F}_a را تابع توزیع تجربی از باقیمانده‌های مرکزی شده در نظر می‌گیریم. یک مجموعه از $p + q$

مقدار اولیه از متغیر y_t یعنی $\{y_1, \dots, y_{p+d}\}$ داده شده، یک تکرار بوت استرپ از سریهای $\{y_1^*, \dots, y_T^*\}$ توسط معادله زیر ساخته می‌شود.

$$\nabla_{y_t^*}^d = \hat{\phi}_0 + \sum_{j=1}^p \hat{\phi}_j \nabla_{y_{t-j}^*}^d + \sum_{j=1}^q \hat{\theta}_j \hat{a}_{t-j}^* + \hat{a}_t^*, t = p+d+1, \dots, T \quad (17)$$

که $y_t^* = y_t$ برای $t = 1, \dots, p+q$ و $\hat{a}_{1+p+d-p}^*, \dots, \hat{a}_T^*$ تصادفاً از \hat{F}_a استخراج شده‌اند. یکبار دیگر پارامترها از این سریهای بوت استرپ برآورد می‌شوند. یعنی $(\hat{\phi}_0^*, \hat{\phi}_1^*, \dots, \hat{\phi}_p^*, \hat{\theta}_1^*, \dots, \hat{\theta}_q^*)$ پیش‌گویی k -مرحله بعدی بصورت زیر می‌باشد.

$$\nabla_{y_{T+k}^*}^d = \hat{\phi}_0^* + \sum_{j=1}^p \hat{\phi}_j^* \nabla_{y_{T+k-j}^*}^d + \sum_{j=1}^q \hat{\theta}_j^* \hat{a}_{T+k-j}^* + \hat{a}_{T+k}^*, k = 1, 2, \dots \quad (18)$$

که

$$\begin{aligned} k \leq j, y_{T+k-j}^* &= y_{T+k-j} \\ k \leq j, a_{T+k-j}^* &= \hat{a}_{T+k-j} \end{aligned}$$

یعنی $p+d$ مشاهده آخر از سریها و q باقیمانده آخر در مرتبه‌ای که چگالی پیش‌گویی شرطی روی داده‌های مشاهده شده بدست می‌آید، ثابت است. حدود پیش‌گویی به صورت چندکهای تابع توزیع بوت استرپ از y_{T+k}^* تعریف می‌شوند. یعنی اگر $G^*(h) = P(y_{T+k}^* \leq h)$ تابع توزیع y_{T+k}^* باشد و برآورد مونت کارلو بصورت $G_B^*(h) = \# \frac{(y_{T+k}^* \leq h)}{B}$ می‌باشد که یک فاصله پیش‌بینی $1-\alpha\%$ برای Y_{T+k}^* بصورت زیر تعریف می‌شود.

$$[L_B^*, U_B^*] = [Q_B^*(\frac{1-\alpha}{2}), Q_B^*(\frac{1+\alpha}{2})] \quad (19)$$

که $Q_B^* = G_B^{*-1}$. این روش را که تغییرات مربوط به برآورد پارامتر را در نظر می‌گیرد PRR می‌نامیم. بنابراین روش بوت استرپ PRR به آسانی برای مدلهایی با میانگین متحرک استفاده می‌شود.

۳ تجویز برآورد نامعین^۱

در بخش قبل خطاها نرمال فرض نشده بودند و تغییرات در برآوردهای پارامتر نادیده گرفته شده بود. یک نمونه بوت استرپ از T مشاهده با نمونه‌گیری از \hat{F}_a بصورت زیر تولید می‌شود.

$$y_{b,t}^* = \hat{\delta} + \hat{\phi}_1 \hat{y}_{b,t-1}^* + \dots + \hat{\phi}_p y_{b,t-p}^* + a_{b,t}^* \quad (20)$$

1) Allowing Estimation Uncertainty

(خطاها با $\{a_{b,t}^*\}$ برای $b = p + 1, \dots, T$ مشخص می‌شوند). برای $t = p + 1, \dots, T$ به شرط مقادیر مشاهده شده اولیه $\{y_1, \dots, y_p\}$ بطوریکه $y_{b,t}^* = y_t$ برای $t = 1, \dots, p$ که نمونه $\{y_{b,p+1}^*, \dots, y_{b,T}^*\}$ را می‌دهد. با تولید B نمونه به این روش برآورد پارامتر در مدل (۱) بصورت $\theta_b^* = \{\hat{\delta}_b, \hat{\phi}_{1,b}, \dots, \hat{\phi}_{p,b}\}$ خواهند بود و معادله بصورت زیر می‌باشد.

$$y_{b,T+k}^* = \hat{\delta}_b + \hat{\phi}_{1,b} y_{b,T+k-1}^* + \dots + \hat{\phi}_{p,b} y_{b,T+k-p}^* + a_{b,T+k}^* \quad (21)$$

$$k = 1, 2, \dots, b = 1, \dots, B$$

که $\{a_{b,T+1}^*, \dots, a_{b,T+k}^*\}$ همانند قبل از $F_{\hat{\alpha}}$ استخراج شده‌اند. توزیع تجربی از $\{y_b^*\}$ برآورد بوت استرپ، از توزیعهای پیش‌بینی مجهول است. اثر برآورد نامعین اینک متشکل از هر یک از تکرارهای بوت استرپ داده قبلی که دوباره استخراج شده و مدل برآورد شده روی آن است.

۱.۳ تصحیح اریبی پارامتر برآورد شده

دیدیم که تکرارهای بوت استرپ از برآورد پارامتر اصلی تولید شدند. Kilian (1998) نشان داد که تصحیح اریبی برآوردهای $\{\hat{\delta}, \hat{\phi}_1, \dots, \hat{\phi}_p\}$ در ساخت فواصل اطمینان برای پاسخهای ناگهانی می‌تواند مفید باشد. شباهت بین پیش‌بینی‌های چند مرحله‌ای و پاسخهای ناگهانی یک استراتژی مشابهی را پیشنهاد می‌کند، به ویژه برای k بزرگ، پارامترهای با توان بالا ممکن است اریبی را بیشتر کنند. وقتی حجم نمونه کوچک است و ریشه‌ها به ناحیه نایستا نزدیک هستند، نادیده گرفتن اریبی نمونه کوچک ممکن است در فواصل با پوشش ضعیف نتیجه شوند. فرض کنید $\theta' = \{\delta, \phi_1, \dots, \phi_p\}$ و ψ مقدار اریبی در برآورد θ باشد، $\psi = E[\hat{\theta} - \theta]$ که یک برآورد بر اساس بوت استرپ بصورت $\hat{\psi} = \bar{\theta}^* - \hat{\theta}$ است که $\bar{\theta}^* = B^{-1} \sum_{b=1}^B \theta_b^*$ و θ_b^* برآوردهایی از پارامترهای مدل (۱) برای هر B تکرار بوت استرپ هستند که از $F_{\hat{\alpha}}$ تولید شده‌اند. شکل تابع $f_t(F_0, F_1) = \theta(F_1) - \theta(F_0) - \psi$ در این حالت بصورت $\psi = \theta(F_1) - \theta(F_0) - \psi$ می‌باشد که $t = \psi$ اریبی را در برآورد پارامترهای جامعه مشخص می‌کند. پارامترهای جامعه $\theta(F_0)$ و برآوردهای OLS بر اساس نمونه $\chi = \{y_1, \dots, y_T\}$ بصورت $\hat{\theta} = \theta(F_1)$ هستند. بنابراین معادله جامعه بصورت زیر می‌باشد:

$$E[\theta(F_1) - \theta(F_0) - \psi | F_0] = 0 \quad (22)$$

که برای به دست آوردن برآوردی از $\hat{\psi}$ که $\hat{\psi} = B^{-1} \sum_{b=1}^B \theta_b^*$ $E[\theta(F_1) | F_1] = \bar{\theta}^*$ معادله نمونه را حل می‌کنیم.

$$E[\theta(F_1) - \theta(F_1) - \psi | F_1] = 0 \quad (23)$$

که $\theta_b^* = \theta(F_{1b})$ تابع توزیع تجربی از χ_b^* می‌باشد که از F_1 استخراج شده است. امید ریاضی θ^* بوسیله میانگین تعداد نمونه‌های زیاد بوت استرپ تقریب زده می‌شود. این تصحیح

اگر لازم باشد، در شکل ساده‌ای می‌تواند تکرار شود. به عنوان مثال اگر در (۲۳)، $\theta(F_1)$ با $\hat{\psi} = \hat{\theta} - \tilde{\theta}$ جایگزین شود آنگاه F_{2b} تابع توزیع تجربی از یک نمونه تولید شده از تصحیح اریبی برآوردهاست. برای مدل‌هایی با ریشه‌های نزدیک به ریشه‌های نزدیک به دایره واحد، باید مطمئن شویم که تصحیح اریبی تعداد ریشه‌های واحد را تغییر نمی‌دهد. در $|\hat{\phi}(z)| = 0$ (که $\hat{\phi} = 1 - \hat{\phi}_1 L - \dots - \hat{\phi}_p L^p$) فرضاً $\hat{\theta}$ درون ناحیه ایستا قرار بگیرد با تصحیح اریبی $\tilde{\theta} = \hat{\theta} - \hat{\psi}$ آنگاه $\tilde{\theta}$ به سادگی در 2° جایگزین $\hat{\theta}$ برای تولید B نمونه بوت استرپ می‌شود. برآوردهای پارامتر روی این نمونه‌ها به دست می‌آیند و با θ_b^* و تصحیح اریبی $\hat{\psi} = \theta_b^* - \hat{\psi}$ مشخص می‌شوند. در عمل عوامل تصحیح اریبی را برای هر θ_b^* می‌توان برآورد کرد. در انجام تصحیح‌های اریبی باید مطمئن شویم که تصحیح باعث پیش‌بینی‌های ایستا به ناپایستا نمی‌شود.

۲.۳ الگوریتم مونت کارلو برای فاصله پیش‌بینی PRR

بطور خلاصه یک تخمین مونت کارلو از فواصل پیش‌بینی PRR با تصحیح اریبی برآوردهای پارامتر بصورت زیر اعمال می‌شود.

مرحله (۱) شبیه‌سازی سری $\{y_1, \dots, y_T\}$ به طول T از مدل $AR(p)$ با توزیع خطای گوسین $F(\hat{a})$ شبیه‌سازی R دنباله از سریها به طول K با بیشترین حد پیش‌بینی. بکارگیری مقادیر پارامتر که از توزیع خطای $F(a)$ گرفته شده و روی y_{T-p+1}, \dots, y_T برای یک فرآیند مرتبه p ام شرطی شده است. این دنباله‌ها مقادیر تحقق یافته آینده هستند و برای برآوردهای فواصل پوشش بکار می‌روند.

مرحله (۲) برآورد $\hat{\theta}$ و F_a روی $\{y_1, \dots, y_T\}$ و انجام بوت استرپ برای تصحیح اریبی $\hat{\theta}$ (مرحله ۳) شبیه‌سازی B تکرار بوت استرپ به طول T با بکارگیری $\hat{\theta}$ و F_a و برآورد (۱) برای بدست آوردن θ_b^* . برآورد اریبی را از مرحله ۲ برای تصحیح اریبی این برآوردها برای داشتن θ_b^* بکار می‌بریم.

مرحله (۴) با شرطی کردن روی $\{y_{T-p+1}, \dots, y_T\}$ برای هر $\hat{\theta}_b^*$ گرفته شده از $F_{\hat{a}}$ و دنباله‌ای از سریها تولید می‌کنیم. که $\{y_{1,T+k}^*, \dots, y_{B,T+k}^*\}$ را برای $k = 1, \dots, K$ می‌دهد. فرض کنید L_k^* و U_k^* فواصل انتهایی در روش چندک Hall را مشخص کنند آنگاه $\lambda_k^* = U_k^* - L_k^*$ طول فاصله و پوشش $\beta_k^* = \frac{\#\{L_k^* \leq y_{T+k}^* \leq U_k^*\}}{R}$ که اندیس r, r امین عضو از R دنباله شبیه‌سازی مرحله ۱ است.

مرحله (۵) تکرار مراحل ۴-۱، مرتبه، مقادیر بدست آمده روی هر تکرار با m اندیس گذاری می‌کنیم $m = 1, \dots, M$.

مرحله (۶) محاسبه طول و پوشش برآوردهای مونت کارلو که تغییرات این برآوردها بصورت زیر می‌باشد.

$$\bar{\beta}_k^* = \frac{\sum_{m=1}^M \beta_{k,m}^*}{M}$$

$$S.E(\bar{\beta}_k^*) = \left\{ \frac{[\sum_{m=1}^M (\beta_{k,m}^* - \bar{\beta}_k^*)^2]}{M-1} \right\}^{\frac{1}{2}}$$

$$\bar{\lambda}_k^* = \frac{\sum_{m=1}^M \lambda_{k,m}^*}{M}$$

$$S.E(\bar{\lambda}_k^*) = \left\{ \frac{[\sum_{m=1}^M (\lambda_{k,m}^* - \bar{\lambda}_k^*)^2]}{M-1} \right\}^{\frac{1}{2}}$$

۳.۳ شبیه‌سازی مونت کارلو برای یک مدل AR(1)

تخمین مونت کارلو را برای مدل $AR(1)$ ، $Y_t = 0.95Y_{t-1} + a_t$ ، محاسبه می‌کنیم. برآوردهایی به حجم $T = \{25, 50\}$ با ماکزیمم حد پیش‌بینی $k = 10$ را در نظر می‌گیریم. فواصل اطمینان (BJ) Box-Jenkins و فواصل PRR و فواصل PRR تصحیح‌اریبی شده k را با سطوح پوشش $c = 80\%$ و $c = 90\%$ شبیه‌سازی و مقایسه می‌کنیم.

مراجع

- [1] Michael P. Clements, Nick Taylor.(2001), Bootstrapping prediction intervals for autoregressive models, International Journal of Forecasting 17, 247-267.
- [2] Pascual, Romo, Ruiz.(2001) Effects of parameter estimation on prediction densities :a bootstrap approach International Journal of Forecasting 17, 83-103.
- [3] Thombs, Schucany.(1990) Bootstrap prediction intervals for autoregression Journal of the American Statistical Association 85, NO 410.
- [4] Efron, Tibshirani.(1993) An introduction to the Bootstrap.
- [5] Box, Jenkins.(1976) Time Series Analysis: Forecasting and control.

تعیین تعداد خوشه‌ها در تحلیل‌های آمیخته با استفاده از آنتروپی نرمال شده

محمد قربانی^۱، محسن محمدزاده^۲

^۱ گروه آمار دانشگاه تبریز

^۲ گروه آمار دانشگاه تربیت مدرس

چکیده: یکی از مسائل مهم در تحلیل خوشه‌ای، تعیین تعداد خوشه‌هاست. علی‌رغم اینکه این موضوع توسط محققان زیادی مورد بحث واقع شده ولی همچنان به عنوان یک مشکل خوشه‌بندی مطرح است. در این مقاله هدف، تعیین تعداد خوشه‌ها با استفاده از آنتروپی نرمال شده در تحلیل آمیخته است که این معیار از ارتباط بین درست‌نمایی کلی و درست‌نمایی رده‌بندی به دست می‌آید. همچنین با استفاده از شبیه‌سازی مونت کارلو نشان داده خواهد شد که این معیار بهتر از معیارهایی مانند BIC و AIC عمل می‌کند.

واژه‌های کلیدی: تحلیل آمیخته، تعداد خوشه‌ها، آنتروپی نرمال شده، درست‌نمایی رده‌بندی

۱ مقدمه

یکی از مسائل مهم در بسیاری از مطالعات ژنتیکی و پزشکی، خوشه‌بندی داده‌ها است، که در آن N مشاهده با M ویژگی به g گروه همگن افراز می‌شوند، به طوری که تشابهات درون گروه‌ها ماکسیمم گردند. روشهای مختلفی از جمله، روشهای سلسله مراتبی برای خوشه‌بندی کردن داده‌ها به کار می‌روند که در تعریف «فاصله بین دو خوشه» با هم تفاوت دارند (هارتیگان، ۱۹۷۵). همچنین به دلیل اینکه این روشها براساس مدل خاصی بنا نشده‌اند، استنباط آماری بر اساس آنها امکان‌پذیر نیست و تعداد خوشه‌ها نیز به صورت ابتکاری با تعریف آستانه‌ای دلخواه تعیین می‌شود. برای رفع این مشکلات لازم است روش خوشه‌بندی حتی الامکان مبتنی بر سلیقه محقق نبوده و بر اساس یک مدل یا توزیع احتمالی باشد تا بتوان در مورد آن استنباط آماری انجام داد. معمولاً مجموعه مشاهدات تحت بررسی، همگی از یک جامعه خاص نیستند و برای تشخیص این که هر مشاهده از کدام جامعه آمده است، منطقی است فرض شود که هر مشاهده بر اساس ویژگی‌ها و خصوصیاتش دارای توزیع احتمال خاصی است. بنابراین جامعه‌ای مرکب از چند زیرجامعه، دارای توزیع احتمالی آمیخته بصورت

$$f(x|\psi) = \lambda_1 f_1(x|\theta_1) + \dots + \lambda_g f_g(x|\theta_g)$$

است، که در آن برای $j = 1, \dots, g$ تابع چگالی مولفه‌ها، $0 < \lambda_j \leq 1$ و معمولاً $\psi = (\lambda, \theta)$ و $\theta = (\theta_1, \dots, \theta_g)$ ، $\lambda = (\lambda_1, \dots, \lambda_g)$ ، $\sum_{j=1}^g \lambda_j = 1$ فرض می‌شود مولفه‌ها دارای توزیع تک مدی $n(\mu_j, \sigma_j^2)$ است و هدف از خوشه‌بندی، تجزیه مولفه‌های چند بعدی مبهم و آمیخته به مولفه‌های ساده تک مدی است. یکی از مسائل مهم خوشه‌بندی تعیین تعداد خوشه‌ها (مولفه‌ها) است. علی‌رغم اینکه این موضوع توسط محققان زیادی مورد بحث واقع شده ولی همچنان به عنوان یک مشکل خوشه‌بندی مطرح است. در این مقاله از معیاری بنام معیار آنتروپی نرمال شده برای تعیین تعداد خوشه‌ها استفاده می‌شود و بکمک تکنیک شمشیه‌سازی مونت کارلو نشان داده خواهد شد این معیار بهتر از معیارهایی مانند AIC و BIC عمل می‌کند.

۲ تعیین معیارهای خوشه‌بندی بر اساس مدل

فرض کنید X_1, \dots, X_n یک نمونه تصادفی از جامعه Π با زیر جامعه‌های Π_1, \dots, Π_g باشند و $\phi(x|\mu_j, \Sigma_j)$ تابع چگالی متغیر تصادفی X_i در زیر جامعه Π_j باشد. در این صورت اگر λ_j احتمال تعلق X_i به جامعه Π_j باشد، توزیع X_i عبارت است از

$$f(x; \psi) = \sum_{j=1}^g \lambda_j \phi(x|\mu_j, \Sigma_j) \quad (1)$$

که در آن $\theta_j = (\mu_j, \Sigma_j)$ ، $\theta = (\theta_1, \dots, \theta_g)$ ، $\psi = (\lambda_1, \dots, \lambda_g, \theta_1, \dots, \theta_g)$ ، $\sum_{j=1}^g \lambda_j = 1$ و $\lambda_j > 0$ تابع درست‌نمایی نمونه تصادفی به صورت

$$L(\psi|x) = \prod_{i=1}^n \sum_{j=1}^g \lambda_j f_j(x|\theta_j) \quad (2)$$

خواهد بود. برای $C_j = \{i, X_i \in \Pi_j\}$ معیار خوشه‌بندی حاصل از ماکسیمم کردن تابع درست‌نمایی

$$L(x|C) = \prod_{j=1}^g \lambda_j^{n_j} \prod_{X_i \in C_j} f(x_i|\theta_j) \quad (3)$$

معادل معیار حاصل از ماکسیمم کردن (۲) خواهد بود (مک‌لن، ۱۹۸۲ و فرالی و رافتری، ۱۹۹۸ و وحیدی و همکاران، ۱۳۸۱). معمولاً مولفه اصلی X_i ها معلوم نیستند و برای مشخص کردن مولفه اصلی X_i متغیرهای گروه‌بندی Z_{ij} به صورت

$$Z_{ij} = \begin{cases} 1 & X_i \in \Pi_j \\ 0 & X_i \notin \Pi_j \end{cases}$$

تعریف می‌شوند. بر اساس مشخصه‌های گروه‌بندی لگاریتم تابع درستنمایی را می‌توان به صورت

$$\log L(\psi|x, Z) = \sum_{j=1}^g \sum_{i=1}^n z_{ij} \log \{\lambda_j f(x_i/\theta_j)\} \quad (4)$$

نوشت. اگر z_{ij} ها مشخص باشند، عناصر خوشه j ام به صورت $C_j = \{j; z_{ij} > z_{ij'} \quad j \neq j'\}$ خواهد بود. ولی اگر z_{ij} ها معلوم نباشند برای انجام تحلیل خوشه‌ای از برآورد مقدار مورد انتظار آنها استفاده می‌شود. پس تحلیل خوشه‌ای را می‌توان به عنوان برآورد مقدار مورد انتظار z_{ij} ، یعنی

$$E(Z_{ij}) = P(X_i \in \Pi_j) = \frac{\lambda_j f(x_i|\theta_j)}{\sum_{j=1}^g \lambda_j f(x_i|\theta_j)}$$

تلقی نمود، که لازم است برای هر j پارامترهای نامعلوم λ_j و θ_j برآورد شوند. الگوریتم امید ریاضی و ماکسیم‌سازی^۱ (EM) روشی کاربردی در محاسبات تکراری، برای به دست آوردن برآورد ماکسیم درستنمایی پارامترها در توزیع‌های آمیخته است. فرض کنید $\psi^{(0)}$ مقدار اولیه ψ باشد، در این صورت مراحل الگوریتم EM به صورت زیر خواهد بود. (مک‌لن و کریشنان، ۱۹۹۷):

مرحله E : محاسبه امید ریاضی لگاریتم تابع درستنمایی در نقطه $\psi^{(0)}$ به شرط مشاهده داده‌های کامل،

$$Q(\psi, \psi^{(0)}) = E_{\psi^{(0)}} \{ \log L_c(\psi|x) \}$$

مرحله M : بدست آوردن مقداری مانند ψ^* برای ψ به طوری که

$$Q(\psi^*, \psi^{(0)}) = \max_{\psi \in \Omega} (Q(\psi, \psi^{(0)}))$$

مراحل E و M تا زمانی تکرار می‌شوند که شرط همگرایی $|L(\psi^{(k+1)}) - L(\psi^{(k)})| < \epsilon$ برقرار شود. (جی اف وو، ۱۹۸۳).

فرض کنید X_1, \dots, X_n داده‌های ناکامل و $Y_i = (X_i, Z_i)$ داده‌های کامل در الگوریتم EM باشند. در این صورت برآورد پارامترهای (۴) با استفاده از این الگوریتم به صورت زیر خواهند بود.

$$\lambda_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n z_{ji}^{(k)}$$

1) Expectation-Maximization Algorithm

$$\mu_i^{(k+1)} = \frac{\sum_{j=1}^n z_{ji}^{(k)} X_j}{\sum_{j=1}^n z_{ij}^{(k)}}$$

$$(\Sigma)^{(k+1)} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^n z_{ij}^{(k)} (X_j - \mu_i^{(k+1)})(X_j - \mu_i^{(k+1)})'$$

$$\hat{z}_{ij}^{(k)} = \frac{\pi_i^{(k)} f(x_j | \theta_i^{(k)})}{\sum_{i=1}^k \pi_i^{(k)} f(x_j | \theta_i^{(k)})}$$

با برآورد z_{ij} به وسیله الگوریتم EM می توان تحلیل خوشه‌ای را متناسب با بیشترین مقدار $\hat{z}_{ij}^{(k)}$ انجام داد (وحیدی اصل، محمدزاده و قربانی، ۱۳۸۱ و مک‌لن و کریشنان، ۱۹۹۷). ولی یکی از مسائل اساسی خوشه‌بندی، تعیین تعداد خوشه‌هاست این موضوع با استفاده از معیارهای مختلف توسط اغلب محققان (وحیدی و محمدزاده و قربانی، ۱۳۸۱، فرالی و رافتری، ۱۹۹۸) بررسی شده ولی همچنان به عنوان یک مشکل خوشه‌بندی مطرح است.

۳ تعیین تعداد خوشه‌ها

برای بیان معیار آنتروپی نرمال شده در تعیین تعداد خوشه‌ها نخست به تعریف آنتروپی می‌پردازیم. متغیر تصادفی X با تابع جرم احتمالی $P(x) = P(X = x)$ و مقادیر $\{x_i, i = 1, \dots, n\}$ در نظر بگیرد، عدم قطعیت مرتبط با برآمد یک مشاهده از متغیر X را آنتروپی متغیر تصادفی X نامیده و با نماد $H(X)$ نشان می‌دهند که عبارت است از

$$H(X) = - \sum P(x) \log P(x)$$

حال لگاریتم تابع درستنمایی (۲) را در نظر بگیرید

$$L(g) = \sum_{i=1}^n \ln \sum_{j=1}^g \hat{\lambda}_j f_j(x | \hat{\theta}_j)$$

که در آن $\hat{\lambda}_j$ و $\hat{\theta}_j$ برآوردهای ماکزیمم درستنمایی λ_j و θ_j می‌باشند. با استفاده از محاسبات مستقیم می‌توان نشان داد

$$L(g) = C(g) + E(g)$$

که در آن

$$C(g) = \sum_{j=1}^g \sum_{i=1}^n \hat{z}_{ij} \log \{ \hat{\lambda}_j f(x_i | \hat{\theta}_j) \}$$

$$E(g) = - \sum_{j=1}^g \sum_{i=1}^n \hat{z}_{ij} \ln \hat{z}_{ij} \geq 0$$

بر اساس روابط فوق، لگاریتم تابع درست‌نمایی $L(g)$ ، به لگاریتم تابع درست‌نمایی رده‌بندی $C(g)$ و معیار آنتروپی $E(g)$ ، تفکیک شده است. اگر چه آنتروپی $E(g)$ یک معیار اطلاع است ولی نمی‌توان مستقیماً از آن برای تعیین تعداد خوشه‌ها استفاده کرد زیرا $L(g)$ یک تابع افزایشی از g است و باید نرمالیزه شود. داریم

$$1 = \frac{C(g) - C(1)}{L(g) - L(1)} + \frac{E(g) - E(1)}{L(g) - L(1)} \quad g > 1 \quad (5)$$

و معیار آنتروپی نرمال شده که باید برای تعیین تعداد خوشه‌ها می‌نیم شود عبارت است از

$$NEC(g) = \frac{E(g)}{L(g) - L(1)}, \quad E(1) = 0$$

اما $NEC(1)$ خوش تعریف نیست زیرا $NEC(1) = \infty$ (بایرناسکی، سیلوکس و گوورت ۱۹۹۹). لذا به طور مستقیم قادر به مقایسه حالت $g = 1$ در برابر $g > 1$ در استفاده از $NEC(g)$ نخواهیم بود. لذا می‌بایست $NEC(g)$ را برای رفع این مشکل خاص بسط دهیم. برای تصمیم‌گیری بین $g = 1$ و $g > 1$ می‌توان از این نقطه نظر که $C(g)$ معیاری برای اندازه‌گیری دقت افراز داده‌ها در g خوشه و $C(1) = L(1)$ معیاری برای اندازه‌گیری دقت برازش تک خوشه به داده‌هاست، استفاده کرد. هنگام مقایسه دو تعداد از خوشه‌ها، g و g' ممکن است یک مدل آمیخته با تعداد پارامترهای زیاد مناسب به نظر رسد ولی از دیدگاه ساده، منطقی است اگر برای $g > g'$ ، $C(g) < C(g')$ باشد g بر g' ترجیح داده شود. بنابراین برای انتخاب $g > 1$ در برابر $g = 1$ لازم است $C(g) > C(1) = L(1)$. در صورتی که $C(g) > L(1)$ ، تمام جزئیات معادله (۵) نامنفی خواهند بود که در این صورت $1 \leq NEC(g) \leq \infty$ ، که تنها حالت $g > 1$ این شرط را فراهم می‌کند پس اگر $NEC(g) \leq 1$ نباشد دلیلی بر انتخاب بیش از یک خوشه وجود ندارد.

۴ مقایسه معیارها

در این بخش معیارهای NEC ، BIC و AIC بر اساس تکنیک شبیه‌سازی مونت کارلو مورد مقایسه قرار می‌گیرند. برای این منظور نمونه‌هایی به حجم n از توزیع نرمال تک متغیره با واریانس‌های برابر یک و همچنین نرمال آمیخته دو متغیره با ماتریس‌های واریانس برابر I تولید

می شود.

برای توزیع نرمال تک متغیره چهار نوع توزیع با پارامترهای مختلف زیر در نظر گرفته شده است:
الف) توزیع نرمال استاندارد.

ب) توزیع نرمال آمیخته دو مولفه‌ای بامیانگین‌های $\mu_1 = 0$ و $\mu_2 = 2$ و نسبت‌های آمیختگی برابر.

ج) توزیع آمیخته نرمال دو مولفه‌ای با میانگین‌های $\mu_1 = 0$ و $\mu_2 = 2$ و نسبت‌های آمیختگی $\lambda_1 = 0.3$ و $\lambda_2 = 0.7$.

د) توزیع آمیخته نرمال سه مولفه‌ای با میانگین‌های $\mu_1 = 0$ ، $\mu_2 = 2$ و $\mu_3 = 4$ و نسبت‌های آمیختگی مساوی.

برای حالت دو متغیره سه نوع توزیع زیر در نظر گرفته شده است.

الف) توزیع نرمال دو متغیره استاندارد.

ب) توزیع آمیخته نرمال دو مولفه‌ای با بردار میانگین $(0, 0)$ و $(2, 2)$ و نسبت‌های مساوی.

ج) توزیع آمیخته نرمال سه مولفه‌ای با بردار میانگین $(0, 0)$ ، $(2, 2)$ و $(2, -2)$ و نسبت‌های مساوی.

از هر یک از توزیع‌ها دو نمونه با حجم‌های $n = 50$ و $n = 200$ هر کدام به تعداد 100 بار شبیه‌سازی شده است. بر اساس نمونه‌های تولید شده، پارامترهای مدل‌های آمیخته بر اساس الگوریتم EM برآورد شده‌اند. در جدول ۱ نتایج مربوط به میانگین و انحراف معیار آنتروپی NEC بیان شده است که در آن d نشان دهنده بعد فضای نمونه‌ای و g نشان دهنده تعداد خوشه‌هاست. بر اساس نتایج به دست آمده NEC برای $n = 200$ بهتر از $n = 50$ عمل می‌کند. به عنوان مثال در جدول ۱ برای حالت الف) توزیع آمیخته نرمال تک متغیره میانگین آنتروپی برای $g = 1$ کمتر از بقیه حالت‌هاست لذا فرض تک خوشه بودن پذیرفته می‌شود. همچنین برای حالت ب) توزیع دو متغیره وقتی که $n = 200$ است، میانگین معیار آنتروپی برای $g = 2$ برابر 6.24 است که کمتر از مقدار آن برای تعداد خوشه‌های $g = 1$ ، $g = 3$ و $g = 4$ است و فرض دو مولفه‌ای بودن را تایید می‌کند. جدول ۲ درصد فراوانی انتخاب توزیع آمیخته g مولفه‌ای بر اساس استفاده از معیارهای AIC ، BIC و NEC برای حالت‌های مختلف مذکور و برای تعداد خوشه‌های مختلف بیان می‌کند.

به عنوان مثال برای حالت اول توزیع نرمال تک متغیره در 85 درصد اوقات NEC جواب صحیح داده است در حالیکه در 70 درصد اوقات AIC پاسخ صحیح می‌دهد. در حالت کلی بر اساس جدول ۲ نتیجه می‌شود که کیفیت تصمیم‌گیری NEC بین AIC و BIC قرار دارد. معیار AIC تعداد مولفه‌های توزیع آمیخته را اندکی بیش تخمین می‌کند و بر عکس BIC تعداد مولفه‌ها را اندکی کمتر تخمین می‌زند.

جدول ۱: میانگین و انحراف معیار برای معیار NEC بر اساس توزیع‌های آمیخته نرمال تک متغیره و دو متغیره

| d | n | پارامترهای توزیع | تعداد خوشه‌ها | | | |
|---|-----|--|------------------|-----------------|------------------|------------------|
| | | | g = ۱ | g = ۲ | g = ۳ | g = ۴ |
| | | $p_1 = 1$ $\mu_1 = 0$ | ۱۶/۴۱ (۶/۱۵) | ۱۷/۹۳ (۷/۹۷) | ۲۲/۸۳ (۸/۶۱) | ۲۹/۳۲ (۱۰/۳۵) |
| ۱ | ۲۰۰ | $p_1 = 0.5$ $\mu_1 = 0, \mu_2 = 2$ | ۱۴/۸۳ (۵/۸۴) | ۱۰/۹۸ (۴/۳۶) | ۱۵/۷۳۵ (۵/۴۹) | ۲۲/۱۰ (۶/۸۱) |
| | | $p_1 = 0.7$ $\mu_1 = 0, \mu_2 = 2$ | ۱۵/۹۸ (۶/۱۲) | ۱۵/۲۸ (۵/۹۵) | ۲۲/۶۷ (۷/۳۱) | ۳۰/۵۲ (۱۰/۱۱) |
| | | $p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۲۰/۴۵ (۱۲/۸۸) | ۷/۴۲ (۲/۳۵) | ۵/۸۹ (۰/۹۷) | ۹/۱۶ (۳/۴۱) |
| | | $p_1 = 1$ $\mu_1 = 0$ | ۱۷/۹۲ (۷/۲۱) | ۱۸/۲۲ (۸/۹۱) | ۲۳/۲۵ (۱۰/۱۲) | ۳۰/۴۹ (۱۲/۷۴) |
| ۱ | ۵۰ | $p_1 = 0.5$ $\mu_1 = 0, \mu_2 = 2$ | ۱۴/۲۷ (۷/۴۲) | ۱۲/۲۳ (۵/۸۲) | ۱۷/۱۸ (۷/۷۹) | ۲۴/۹۴ (۶/۹۱) |
| | | $p_1 = 0.7$ $\mu_1 = 0, \mu_2 = 2$ | ۱۸/۵۴ (۱۰/۱۷) | ۱۶/۲۶ (۷/۶۴) | ۲۴/۵۲ (۹/۸۶) | ۳۲/۸۶ (۱۳/۲۲) |
| | | $p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۲۸/۳۸ (۱۱/۸۲) | ۵/۸۲ (۲/۱۰) | ۵/۱۵ (۰/۶۲) | ۷/۵۵ (۲/۸۴) |
| | | $p_1 = 1$ $\mu_1 = [0, 0]$ | ۳/۴۱ (۳/۰۲) | ۱۶/۸۴ (۶/۱۲) | ۱۸/۶۹ (۶/۶۴) | ۲۴/۴۷ (۷/۷۶) |
| ۲ | ۲۰۰ | $p_1 = 0.5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۸/۸۶ (۷/۴۲) | ۶/۲۴ (۵/۸۲) | ۸/۱۲ (۷/۷۹) | ۱۰/۴۲ (۶/۹۱) |
| | | $p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۱۱/۱۹ (۱۱/۸۲) | ۳/۵۹ (۲/۱۰) | ۲/۰۲ (۰/۶۲) | ۲/۹۰ (۲/۸۴) |
| | | $p_1 = 1$ $\mu_1 = [0, 0]$ | ۹/۶۷ (۴/۱۵) | ۱۱/۵۸ (۶/۹۹) | ۲۳/۹۶ (۷/۴۴) | ۳۰/۹۵ (۸/۹۴) |
| ۲ | ۵۰ | $p_1 = 0.5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۱۰/۸۱ (۶/۶۸) | ۸/۹۷ (۴/۱۲) | ۱۰/۲۴ (۵/۸۳) | ۱۲/۷۰ (۷/۵۶) |
| | | $p_1 = p_2 = p_3 = \frac{1}{3}$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۲/۹۲ (۰/۸۸) | ۲/۳۵ (۲/۱۰) | ۴/۴۳ (۰/۶۲) | ۱۱/۹۷ (۲/۸۴) |

جدول ۲: درصد فراوانی انتخاب g خوشه بر اساس توزیع‌های آمیخته نرمال تک متغیره و دو متغیره

| d | n | پارامترهای توزیع | تعداد خوشه‌ها | معیارها | | |
|-----|-----|--|---------------|---------|-------|-------|
| | | | | AIC | BIC | NEC |
| ۱ | ۲۰۰ | $p_1 = 1$ $\mu_1 = 0$ | ۱ | ۷۰ | ۹۵ | ۸۵ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۲ | ۳۰ | ۵ | ۱۵ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۱ | ۱۰ | ۴۰ | ۱۵ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۲ | ۸۰ | ۶۰ | ۸۰ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۳ | ۱۰ | ۰ | ۵ |
| | | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۱ | ۲۵ | ۵۰ | ۲۵ |
| ۱ | ۵۰ | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۲ | ۶۵ | ۵۰ | ۷۰ |
| | | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۳ | ۱۰ | ۰ | ۵ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۱ | ۰ | ۱۰ | ۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۲ | ۵ | ۱۵ | ۱۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۳ | ۸۰ | ۷۵ | ۸۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۴ | ۱۵ | ۰ | ۱۰ |
| ۱ | ۲۰۰ | $p_1 = 1$ $\mu_1 = 0$ | ۱ | ۶۰ | ۸۵ | ۸۰ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۲ | ۳۰ | ۱۵ | ۲۰ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۳ | ۱۰ | ۰ | ۰ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۱ | ۲۵ | ۵۵ | ۳۵ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۲ | ۶۰ | ۴۵ | ۵۵ |
| | | $p_1 = 0/5$ $\mu_1 = 0, \mu_2 = 2$ | ۳ | ۱۵ | ۰ | ۱۰ |
| ۱ | ۵۰ | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۱ | ۰ | ۲۰ | ۰ |
| | | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۲ | ۱۵ | ۲۰ | ۲۰ |
| | | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۳ | ۶۵ | ۵۰ | ۷۰ |
| | | $p_1 = 0/7$ $\mu_1 = 0, \mu_2 = 2$ | ۴ | ۲۰ | ۰ | ۱۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۱ | ۰ | ۲۰ | ۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = 0, \mu_2 = 2, \mu_3 = 4$ | ۲ | ۱۵ | ۲۰ | ۲۰ |
| ۲ | ۲۰۰ | $p_1 = 1$ $\mu_1 = [0, 0]$ | ۱ | ۸۰ | ۱۰۰ | ۹۵ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۲ | ۲۰ | ۰ | ۵ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۱ | ۵ | ۳۵ | ۱۰ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۲ | ۸۰ | ۶۰ | ۸۰ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۳ | ۱۵ | ۵ | ۱۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۱ | ۰ | ۵ | ۰ |
| ۲ | ۵۰ | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۲ | ۵ | ۲۵ | ۵ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۳ | ۷۵ | ۷۰ | ۷۵ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۴ | ۲۰ | ۰ | ۲۰ |
| | | $p_1 = 1$ $\mu_1 = [0, 0]$ | ۱ | ۶۰ | ۸۰ | ۸۰ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۲ | ۴۰ | ۲۰ | ۲۰ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۱ | ۱۰ | ۶۰ | ۱۰ |
| ۲ | ۲۰۰ | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۲ | ۶۵ | ۴۰ | ۷۰ |
| | | $p_1 = 0/5$ $\mu_1 = [0, 0], \mu_2 = [2, 2]$ | ۳ | ۲۵ | ۰ | ۲۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۱ | ۰ | ۲۵ | ۰ |
| ۲ | ۵۰ | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۲ | ۲۰ | ۳۰ | ۱۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۳ | ۵۵ | ۴۵ | ۶۰ |
| | | $p_1 = p_2 = p_3 = 1/3$ $\mu_1 = [0, 0], \mu_2 = [2, 2], \mu_3 = [2, -2]$ | ۳ | ۲۵ | ۰ | ۳۰ |

مراجع

- [1] Biernacki, C., Celeux, G. and Govert, G. (1999), An Improvement of the NEC Criterion for Assessing the Number of Clusters in a Mixture model, *Pattern Recognition Letters*, 20, 267-272.
 - [2] Celeux, G. and Soromenho, G. (1996), An Entropy Criterion for Assessing the Number of Clusters in a Mixture Model. *Journal of Classification* , 13, 195-212.
 - [3] Fray, C. and Raftery, A. E. (1998), How Many Clusters? Which Clustering Method? Answer via Model-Based Cluster Analysis. Technical Report, No. 329. Seattle: Department of Statistics, University of Washington.
 - [4] Fraley, C. and Raftery, A. E. (1999), MCLUST: Software for Model-Based Cluster Analysis, *J. Classification*, 16, 297-306.
 - [5] Hartigan, J. A. (1975), *Clustering Algorithms*. Wiley, New York.
 - [6] Jeff C. F., (1983), On the Convergence of the EM Algorithm, *Annals of Statistics*, 11, 95-103.
 - [7] McLachlan, G. J. (1982), The Classification and Mixture Maximum Likelihood Approaches to Cluster Analysis. In Krishnan, P. R. and Kanal, L. N. (eds), *Handbook of Statistics*, 2, 199-208. North-Holland, Amsterdam.
 - [8] McLachlan, G. J. and Krishnan, T. (1997), *The EM Algorithm and Extensions*. Wiley, New York.
- [۹] محمدقاسم وحیدی اصل، محسن محمدزاده و محمد قربانی (۱۳۸۱)، تعیین مدل خوشه‌بندی احتمالاتی بر اساس معیار اطلاع بیزی، مجموعه مقالات ششمین کنفرانس بین‌المللی آمار ایران، ص ۴۶۵-۴۵۵.

شناخت اشیاء (شکل و ساختار) با استفاده از آمار

هادی گنجی^۱، محمود صفارزاده^۲

^۱ فوق لیسانس آمار اقتصادی - اجتماعی، کارشناس ارشد پژوهشکده حمل و نقل
^۲ دکتری برنامه‌ریزی حمل و نقل، معاون پژوهشی پژوهشکده حمل و نقل

چکیده:

۱ مقدمه

در فضای وسیعی از اشیاء منتظم، اندازه‌گیری، تحلیل و مقایسه شکل اشیاء با یکدیگر از اهمیت بالایی برخوردار است [1]. لیدز^۱ از جمله اولین افرادی است که از آنالیز شکل برای مشخص کردن قضیه محل مرکزی در جغرافیا آن را به کار برد [1]. ماردیا^۲ در سال (۱۹۷۷) توزیع مربوط به تولید مثلث‌ها که به وسیله تعداد مشخصی نقطه تولید می‌شوند را استخراج کرد. در یک نگاه تخصصی‌تر پاسخ برای این سؤال که آیا روستاها در (شهرها) یک منطقه به صورت یک قاعده مشخص و با فاصله مساوی از روستاهای (شهرهای) اطراف ایجاد شده‌اند [2]، از موضوعات مطرح شده در آنالیز شکل می‌باشد.

تلفیق آنالیز شکل با آمار در سال ۱۹۸۶ به وسیله اولینگس^۳ و جانسون^۴ بوده است که در مقاله مربوطه اسکلت موش‌ها از نظر شکل ساختاری، اندازه مورد تجزیه و تحلیل قرار گرفته است [3]. در این مقاله در آن هستیم که ضمن معرفی آنالیز و کاربردهای آن، کاربرد آمار را در آنالیز شیء مورد بررسی قرار داده و تکنیک‌های مختلف آماری که در این موضوع کاربرد دارد، معرفی شود.

۲ آنالیز شکل

اشیاء در همه جا هستند چه انسان آن را به وجود آورده باشد و یا اینکه محصول طبیعت باشد. امروزه دانشمندان به دنبال این موضوع هستند که با جمع‌آوری اطلاعات هندسی این اشیاء به مطالعه خصوصیت شکل آنها بپردازند. در حقیقت این موضوع با جمع‌آوری اطلاعات مربوط به نقاط موجود در شیء انجام می‌گیرد. با این توضیح اولیه مشخص شد که جامعه‌ای از نقاط که شیء را تعریف کند، وجود دارد.

با روش‌های مختلف اطلاعات مربوط به برخی از نقاط را که دارای اهمیت می‌باشند جمع‌آوری کرده و بعد از آن به تجزیه و تحلیل پرداخته می‌شود. در حقیقت یک کار آماری است که شامل

1) Leeds 2) Mardia 3) O'ltiggins 4) Jonson

روش نمونه‌گیری، استنباط آماری راجع به پارامتر آن و نتیجه‌گیری است. جهت بررسی این موضوع به تعاریف اولیه مورد نیاز، پرداخته می‌شود.

شکل: به کلیه اطلاعات هندسی که از یک شیء بعد از حذف اثرات اندازه، دوران و محل قرارگیری باقی می‌ماند اطلاق می‌گردد [4]. از این تعریف اینگونه استنتاج می‌گردد که شکل اشیاء نسبت به تغییرات هندسی از قبیل تغییر اندازه یا دوران ثابت باقی می‌ماند. در حقیقت در آنالیز آماری شکل با داشتن اطلاعات مربوط به تعدادی از زوایا و یال‌ها، کل شکل برآورد می‌گردد. البته امروزه با پیشرفت صنایع عکس‌برداری و ثبت اطلاعات لازمه، موضوع از این نیز فراتر رفته و دیگر خصوصیت اشیاء، از قبیل جنس تشکیل دهنده نیز مورد بررسی قرار می‌گیرد.

از این تعریف اینگونه استنباط می‌شود که دو شکل یکسان هستند، اگر با تغییر اندازه و یا دوران یکی را به دیگری بتوان تبدیل کرد. در عمل نیز می‌توان با این ویژگی اشکال گوناگون اشیاء را با یکدیگر مقایسه کرد.

اشکال به وسیله تعدادی محدودی از نقاط مشخص می‌شود که به این نقاط مشخص نقطه علامت^۵ گفته می‌شود و ویژگی آن این است که مرز بین نقاط بیرونی و درونی جامعه می‌باشد. به طور مثال در آناتومی جمجمه گوشه چشم را می‌توان به عنوان یک نقطه علامت در نظر گرفت. این نقاط معمولاً به وسیله متخصصین رشته مربوطه مشخص می‌گردد. با استفاده از اطلاعات این نقاط است که استنباط بر روی شکل صورت می‌گیرد. دو روش جهت تجزیه و تحلیل این اطلاعات وجود دارد: ۱- تجربی، ۲- هندسی.

۱.۲ روش تجربی

در روش تجربی می‌توان نسبت بین فاصله نقاط علامت یا زوایا را مشخص کرده و این داده‌ها را با استفاده از آنالیز چند متغیره تحلیل کرد که به طور خاص در بیولوژی به آن چند متغیره morphometries گفته می‌شود که ریمنت^۶ (۱۹۸۴) به آن پرداخته است. در کل در این روش به تفاوت در فاصله بین طول‌ها یا پهناها در میان نقاط علامت پرداخته می‌شود. به طور مثال در بررسی شباهت بین جمجمه‌ها پاره خط‌های ایجاد شده میان نقاط علامت به وسیله آنالیز چند متغیره طبقه‌بندی می‌شود. مثلاً طبقه‌بندی قسمت‌های مختلف جمجمه با توجه به جنسیت که از اطلاعات نقاط علامت بدست آمده شامل می‌شود.

یک ابزار بسیار مفید در این مورد تحلیل عاملی است که جنبه‌های مختلف اندازه یا خصوصیات شکل را بیان می‌کند. معمولاً عامل اول در برگزیده اطلاعات مربوط به کلیه متغیرها (اطلاعات مربوط به نقاط علامت) بوده و اندازه کلی شکل از آن استخراج می‌گردد [4].

دسته دیگری از مطالعات تحلیل تفاوت بین اشکال با اندازه‌های یکسان است. این کار معمولاً با برآزش یک معادله غیر خطی بر روی طول نقاط علامت همراه است. البته از اطلاعات مربوط به زوایا، نسبت بین طول‌ها و یا استفاده از مقادیر واقعی بر اساس محورهای مختصات نیز استفاده می‌گردد. هنوز هم از روش‌های چند متغیره در بیولوژی و شکل‌شناسی استفاده بسیاری

5) Lamdmark 6) Reyment

می‌گردد که می‌توان به کار برد در طبقه‌بندی شکل اشاره کرد، همچنین باید ذکر کرد که برقراری رابطه خطی بین نسبت اندازه طول‌ها و زوایا کار مشکلی است.

۲.۲ روش‌های هندسی

در دو دهه اخیر تعریف دقیق‌تری از نقاط علامت ارائه شده است که این خود کمک‌کننده در توسعه آنالیز شکل می‌باشد. البته با گسترش تکنولوژی این کار به دیجیتالی شدن اشیاء منجر شده است. اگر هیچ پیش شرطی بر روی این نقاط علامت وجود نداشته باشد روش کار همان استفاده از آنالیز چند متغیره می‌باشد که با توجه به اینکه مختصات این نقاط اقلیدسی نبوده، این کار مشکل می‌باشد [5].

ایده اصلی در روش هندسی بدین صورت می‌باشد که بجای استفاده از مقادیر عددی (جبر) به اطلاعات هندسی اشیاء پرداخته می‌شود که این کار معمولاً با تبدیل قسمتی از شکل همراه است. در این روش نیز هدف اصلی آنالیز شکل همان استنباط در رابطه با تفاوت اشکال اشیاء است. به طور مثال بیولوژیک با استفاده از روش‌های هندسی در ترمیم و یا تشریح تغییرات شکل، یک راهبرد اساسی است. همچنین این روش در باستان‌شناسی، ستاره‌شناسی و جغرافی کاربرد زیادی دارد.

۳ کاربردهای اختصاصی آنالیز شکل [15,11]

جهت مشخص شدن کاربردهای وسیع این موضوع به تعدادی از موارد انجام شده اشاره می‌شود. به طور مثال در پاسخ به این سؤال که چه تغییراتی در شکل به هنگام رشد رخ می‌دهد و یا اینکه شکل اشیاء چگونه به اندازه آنها مرتبط می‌شود یا چه تأثیری بیماری بر روی شکل شیء مورد نظر دارد. ارتباط شکل با دیگر متغیرها از قبیل سن، جنسیت و یا شرایط محیطی چگونه می‌باشد. چگونه می‌توان شکل مورد نظر را طبقه‌بندی و شاخص‌گذاری کرد و در آخر این که چگونه می‌توان تغییرپذیری شکل را توضیح داد. در تمام موارد اشاره شده، روش‌های آنالیز چند متغیره پاسخگو می‌باشند اما یکی از مشکلات اساسی در مورد مطالعه اینگونه موارد حجم کوچک نمونه با تعداد متغیر (نقاط علامت) زیاد می‌باشد. در ذیل به پاره‌ای از این کاربردها به صورت طرح مسأله اشاره می‌شود.

۱.۳ بیولوژی

در یک آزمایش برای مشخص کردن تأثیر وزن بر روی اسکلت و پیدا کردن توضیح برای این تفاوت، موش‌ها را به سه گروه مجزا تقسیم می‌کنند. گروه کنترل، کوچک و بزرگ. گروه کنترل هیچ مشخصه خاصی ندارند ولی در گروه بزرگ موش‌ها با اندام‌های بزرگ و وزن بالا و در گروه کوچک

بالعکس انتخاب شدند. نتیجه آزمایش مشخص کننده این موضوع می باشد که این تأثیر چگونه است.

۲.۳ داروسازی

شناسایی تفاوت شکلی در مغز انسانها در دو گروه که اول شامل افراد سالم و گروه دوم شامل افراد با بیماری اسکیزوفرنی تفاوتی در شکل وجود دارد.

۳.۳ کشاورزی [6]

تعیین و برآورد نحوه توزیع ماهیان در زیر آب و دریافت اطلاعات فیزیکی مربوط به آنها.

۴.۳ جغرافیا [6]

قضیه جغرافیای مرکزی به این موضوع اشاره دارد که شهرها در یک شش ضلعی منتظم در یک منطقه هموار توزیع می شوند. ایده اولیه این موضوع از ۴۴ نقشه مربوط به ۶ کشور به دست آمده است. اگر درستی این قضیه فرض گردد، می توان این آزمون فرض را انجام داد که مثلث های ساخته شده تحت شهرهای همسایه متساوی الاضلاع هستند یا خیر؟
تمام موارد اشاره شده با استفاده از عکس برداری، ثبت اطلاعات نقاط علامت و استفاده از آنالیز چند متغیره مورد بررسی قرار گرفته اند. با توجه به اشاره به برخی از کاربردهای آنالیز شکل در علوم مختلف به یک مورد از کاربردهای آن با آنالیز تصویر پرداخته می شود.

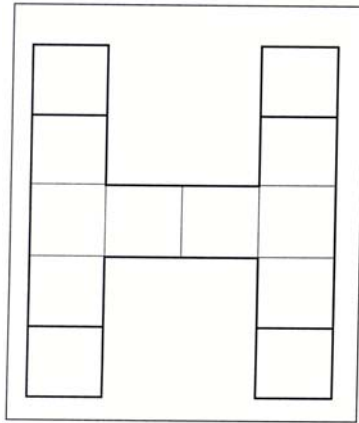
۴ آنالیز تصویر

آنالیز تصویر معمولاً با سنجش از راه دور همراه است. سنجش از دور به عنوان تکنولوژی و علمی تعریف می شود که به وسیله آن می توان با گرفتن تصویر از اشیاء یا پدیده ها آن ها را تعیین، اندازه گیری یا تجزیه و تحلیل نمود. منبع داده های سنجش از دور معمولاً تابش الکترومغناطیسی است که از یک پدیده یا شیء، بازتاب یا گسیل می شود. کاربردهای سنجش از دور شامل مواردی همچون [7]:

- ۱- طبقه بندی پوشش زمین
- ۲- آشکارسازی تغییرات پوشش زمین
- ۳- پوشش و بررسی تصویری کیفیت آبها
- ۴- اندازه گیری دمای سطح دریاها
- ۵- بررسی و مطالعه برف
- ۶- پوشش و بررسی تصویری اجزای جوی
- ۷- آشکارسازی خط واره ها
- ۸- تعبیر زمین شناسی
- ۹- اندازه گیری ارتفاعی (تولید DEM)

جهت مطالعه یک تصویر آن را به گروهی از سلولها (المانهای) مساوی تقسیم کرده و به هر کدام از این سلولها (المانها) یک پیکسل گفته می شود. در حقیقت با انجام این کار یک

ماتریس $r \times c$ تولید می‌شود. مقادیر مربوط به عناصر این ماتریس، روشنایی رنگ خاکستری بازتاب شده از هر سلول می‌باشد که با اشعه X به آن تابش شده است، این مقدار به موقعیت این پیکسل‌ها بستگی دارد. دامنه تغییرات این رنگ خاکستری از سطح صفر (مشکی) تا ۲۵۵ (سفید) است. بنابراین یک پیکسلی که با عدد 5° مشخص می‌شود خاکستری تیره و ۱۲۸ خاکستری متوسط و 20° خاکستری روشن را در عمل از خود نشان می‌دهد. ساین معمول تصاویر $(c = 256) \times (r = 256)$ یا $(c = 512) \times (r = 512)$ است. دلیل آن هم $2^9 = 512$ و $2^8 = 256$ است (یعنی در نظر گرفتن ۸ سطح از رنگ خاکستری یا ۹ سطح از آن) که با بیشتر شدن این سطوح کار با آنها بسیار مشکل است. برای تصاویر رنگی نیز به همین منوال می‌باشد، بدین ترتیب که تصاویر رنگی در سه گروه از خاکستری نمایش داده می‌شود. برای هر باند آبی، قرمز و سبز جداگانه که یک تصویر 256×256 با $256 \times 3 = 768$ پیکسل نمایش داده می‌شود. چند متغیره بسیار بزرگ که بعضی از آنها نیز غیرقابل‌دسی می‌باشد. به طور مثال تصویر حرف H و نمایش ماتریس مقادیری آن به صورت زیر است:



⇓

$$\begin{pmatrix} \cdot & 211 & 2 & 1 & 243 & \cdot \\ \cdot & 200 & 5 & 8 & 251 & \cdot \\ \cdot & 210 & 241 & 251 & 254 & \cdot \\ \cdot & 236 & 2 & 12 & 204 & \cdot \\ \cdot & 251 & 7 & 12 & 218 & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}$$

بعد از معرفی اولیه مربوط به آنالیز تصویر نحوه نمونه‌گیری مشخص کردن جامعه، بحث بر روی برآورد مقادیر مربوطه می‌باشد که از آن طریق بتوان طبقه مورد نظر شکل و یا ظاهر آن را بتوان

تخمین زد. در این برآورد دو روش کلاسیک و بیز در استنباط آماری وجود دارد. در حقیقت با این روش می‌توان پیش‌بینی کرد که خصوصیت شکل و یا شکل ظاهری آن به چه ترتیبی می‌باشد.

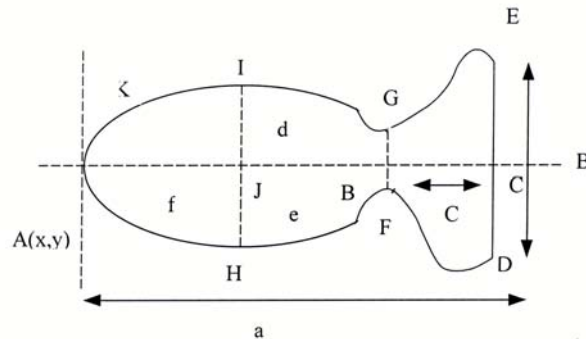
تجزیه و تحلیل تصاویر از روش آماری بیز تقریباً از سال ۱۹۸۰ آغاز گردیده است. برای استفاده از این روش احتیاج به یک مدل پیشین می‌باشد که نشان دهنده اطلاعات اولیه از اشیاء در موقعیت خاص است. همچنین احتیاج به مدل توأم از سطوح مختلف رنگ خاکستری در تصویر است که به محل اشیاء نیز بستگی دارد. به طور مثال در یک تصویر پزشکی می‌توانیم این اطلاعات اولیه را داشته باشیم که تصویر گرفته شده از مغز یا قلب می‌باشد. اطلاعات اولیه می‌تواند در قالب مدل اولیه S_0 ارائه شود. توجه شود که S_0 می‌تواند شامل اطلاعات مربوط به یک شی یا کلبه اشیاء موجود در تصویر باشد. تابع توزیعی اولیه‌ای برای پارامترهای جامعه (شکل) در نظر گرفته می‌شود و با توزیع $\pi(s)$ که کلمه تغییرات S_0 را نشان می‌دهد، مشخص می‌گردد.

بنابراین S تابعی از پارامترهای $\theta_1 \dots \theta_p$ است. مدل کل تصویر نیز مورد نیاز می‌باشد. تصویر I را با ماتریس $x_i \in \{1 \dots c\} \times \{1 \dots r\}$ فرض کنید [1] تابع توزیعی توأم از سطوح مختلف خاکستری که از شیء S مشخص می‌شود و به صورت $L(I/S)$ نمایش داده می‌شود. طبق تعریف بیز توزیع پسین شیء موجود در تصویر با اطلاعات حاصل از تصویر به شکل زیر خواهد بود: $\pi(S/I) \propto L(I/S) \pi(S)$ که در اینجا از روش‌های خاص برآورد از قبیل ماکزیمم کردن و یا شبیه‌سازی استفاده می‌گردد. دو روش نیز جهت برآورد مدل پیشین ۱- هندسی ۲- آماری وجود دارد.

۱.۴ روش هندسی

در روش هندسی مشخصات هندسی شیء از قبیل خطوط، دایره‌ها، بیضی‌ها و قوس‌ها و . . . مورد توجه قرار می‌گیرد. مدل‌های به کار رفته در این روش معمولاً دارای دو پارامتر اساسی است. ۱- محل ۲- کل شیء مورد نظر. پارامتر کل، بزرگی شیء را نشان می‌دهد و پارامتر محل اطلاعات ریزتری را بیان می‌دارد. برای مشخص شدن موضوع به ارائه یک مثال می‌پردازیم.

به طور مثال شکل یک ماهی را در نظر بگیرید. روش هندسی جهت تدوین توزیع پیشین را می‌توان به صورت زیر در نظر گرفت. شکل را به دو سهمی و یک مثلث به گونه‌ای که طول ماهی a و طول دم آن b است، تقسیم می‌کنیم.



مقادیر a, b, c, d, e, f مشخص کننده اندازه و شکل ماهی است. فرض کنید که مختصات $A(0, 0)$ باشد. بنابراین مختصات نقطه $c(a, 0)$ مشخص می‌گردد. شرط دیگر اینکه $GF < ED$ و نقاط B و C وسط CF و ED قرار دارند و $GFDE$ دوزنقه است. بنابراین D و E دارای مختصات $(a, \frac{b}{4})$ ، $(a, -\frac{b}{4})$ می‌باشد. معادله سهمی AIG به وسیله $I(f, d)$ مشخص می‌گردد به طوری که معادله AHF از نقطه $H(f, -e)$ ، $(0, 0)$ می‌گذرد. در نتیجه مختصات $B(a - c, 0)$ می‌باشد که مختصات عرض نقاط F و G از این سهمی‌ها مشخص می‌گردد. به علاوه مختصات J نقطه تلاقی IH و AC است. در نتیجه:

$$A = (0, 0), B = (a - c, 0), C = (a, 0), D = (a, -\frac{b}{4}), E = (a, \frac{b}{4})$$

$$F = (a - c, -e + k_2(a - c - f)^2), G = (a - c, d - k_1(a - c - f)^2)$$

$$H = (f, -e), I = (f, d)$$

و مقادیر $k_1 = \frac{d}{f}$ و $k_2 = \frac{e}{f^2}$ است. همچنین هر نقطه‌ای مثل $x = k$ بر روی سهمی AIG دارای مختصات $(a - k, d - k_1(a - k - f)^2)$ است. برای مثال برای یک تصویر 120×120 مقادیر قابل قبول برای میانگین و میانگین شرطی به صورت زیر است:

$$E(a) = 110, \quad E(d/a) = \frac{a}{9}, \quad E(b/d) = 2d$$

$$E(e/d) = \frac{2d}{3}, \quad E(c/a) = \frac{a}{6}, \quad E(f/a) = \frac{2a}{16}$$

تغییر شکل ماهی می‌تواند با تغییر مختصات نقطه A و طول a تغییر کند که دیگر نقاط نیز با توجه به این دو مقدار تغییراتی خواهند داشت. برای مثال می‌توان توابع توزیع زیر را در نظر گرفت:

$$a \sim N(110, 20), \quad d/a \sim N\left(\frac{a}{9}, 3\right), \quad b/d \sim N(2d, 1)$$

$$e/d \sim N\left(\frac{2d}{3}, 3\right), \quad \frac{c}{a} \sim N\left(\frac{a}{6}, 1\right)$$

به این ترتیب $\frac{f}{a} \sim N\left(\frac{\sqrt{a}}{\sqrt{6}}, 0\right)$ است که f دارای مقدار ثابت $\frac{\sqrt{a}}{\sqrt{6}}$ است اگر a مشخص باشد.

۲.۴ روش غیرهندسی

فرض کنید که مختصات نقطه $A(x, y)$ یک متغیر پیوسته یکنواخت بر روی تصویر $N \times N$ باشد. تابع توزیع توأم (x, y) و (a, b, c, d, e) از معادله ۱ به صورت زیر در نظر گرفته می‌شود.

$$\pi(x, y, a, b, c, d, e) = \phi_1(a) \phi_2(d/a) \phi_3(b/d) \phi_4(e/d) \phi_5(c/a) \cdot \frac{1}{N^5}$$

اگر $0 < x, y < N$ و $J = 1, \dots, 5$ تابع توزیع نرمال با میانگین و واریانس معادلات ۱ باشد. روش ساده دیگر این است که برای $\theta_1 = x$, $\theta_2 = y$, $\theta_3 = a$ توزیع پیشین زیر در نظر گرفته شود.

$$n(\theta_1, \theta_2, \theta_3) = \begin{cases} \exp\left\{-\frac{1}{\sigma^2}(\theta_3 - 110)^2\right\} & 0 < \theta_1, \theta_2 < N \\ 0 & o.w \end{cases}$$

هدف از دو روش فوق نحوه استفاده از پارامترهای هندسی در برآورد مدل پیشین می‌باشد که می‌توان برای تمام اشیاء تعمیم داد.

مثال: در انتها بعد از ارائه تعاریف اولیه و مفهوم کلی آنالیز شکل و تعیین نقش اساسی آمار در این موضوع مثال واقعی که مورد بررسی قرار گرفته است ارائه می‌شود. البته لازم به ذکر است که به دست آوردن یک چنین داده‌هایی از داخل ایران بسیار مشکل می‌باشد. هدف این مقاله نیز ارائه یکی دیگر از کاربردهای وسیع آمار می‌باشد.

در طبقه‌بندی سطح زمین در ناحیه نیوجرسی که یک منطقه مشخص خاص آن مد نظر می‌باشد، عکسی از آنجا تهیه شده و مورد بررسی قرار گرفته است. این ناحیه شامل بخش مزرعه، جنگل، آب و منطقه شهرنشینی می‌باشد [9,11]. در ابتدا باید به این موضوع اشاره داشت که در نقاط مرزی بین این چهار طبقه است که معمولاً تشخیص و برآورد آنها مفهوم پیدا می‌کند و از اهمیت بیشتری برخوردار است. می‌توان مدل‌های پیش‌بینی تجربی را برای این نقاط مشخص کرد که معمولاً مدل پیشین مناسب در این موارد توزیع Gibbs می‌باشد که شامل دو پارامتر است.

$$\bar{\pi}(S) = \frac{1}{Z(\beta)} \exp\{-\beta V(x)\} \quad \beta > 0$$

به طوری که $Z(\beta)$ ثابت پایدار^۷ است.

$$Z(\beta) = \int \exp\{-\beta V(x)\} dx$$

$V(x)$ همان سطوح مختلف رنگ خاکستری انعکاس شده از پیکسل‌ها است. با توجه به بحث همسایگی در پیکسل‌های در یک شکل و در انواع مختلف تابع توزیع‌های پیشین دیگری نیز مطرح می‌باشد به طور مثال تابع log-Cauchy است. به صورت $\pi(s) = \left(\frac{u^s}{1+u^s}\right)$ و یا پیشنهاد دیگری مبنی بر استفاده از توزیع log-Cosh که دارای فرمت کلی $\bar{\pi}(S) = c \setminus \log \cosh(c\setminus, \mu)$ می‌باشد.

با مشخص کردن تابع توزیع‌های پیشین و استفاده از الگوریتم EM در برآورد پارامترهای مدل، نقاط نامشخص تصویر و یا به عبارت بهتر با استفاده از قانون ماکزیمم احتمال پیکسل‌ها را به یکی از چهار طبقه اختصاص داده می‌شود.

البته در تشخیص این طبقه‌بندی می‌توان به یک سری از اطلاعات اولیه موجود در خصوص منطقه نیز اشاره داشت و آنها را کلاسه‌بندی کرد:

- ۱- پیکسل رودخانه از نظر روشنایی رنگ خاکستری، کمترین شباهت را به پیشکل مزرعه دارد. همین مورد برای پیکسل مربوط به جاده نیز صادق است.
- ۲- در صورت وجود شیب در یک جاده، پیکسل مربوطه شبیه به مزرعه نمی‌باشد.
- ۳- پیکسل‌های جنوب بیشتر شبیه به جنگل می‌باشند.
- ۴- پیکسل‌های جاده و آب بسیار شبیه به هم هستند.
- ۵- پیکسل‌های جاده و مناطق شهرنشینی کمترین شباهت را به هم دارا می‌باشند.
- ۶- پیکسل‌های جاده و جنگل بسیار شبیه به هم می‌باشند.

(منظور از شبیه به هم این است که سطوح خاکستری منعکس شده از این مناطق تقریباً نزدیک به هم می‌باشند). مورد بالای اشاره شده در ۶ بند در حقیقت ترکیب تجزیه و انتخاب مدل‌های پیشین است.

توزیع‌ها و موارد اشاره شده در حقیقت به صورت شرطی است که هر کدام برآورد پارامترها را در آن سطح طبقه در بردارد. علاوه برآن فرض می‌شود که اشعه‌های بازتاب از شکل از فرآیند پواسن پیروی می‌کند.

$$Y_t = P \left(\sum_S a_{ts} x_s \right)$$

که a_{ts} اثر پیکسل‌های همسایه بر مقدار بازتاب‌ها در S طبقه مشخص می‌باشد. با توجه به اطلاعات داده شده و در نظر گرفتن اطلاعات توزیع پیشین، جداول زیر حاصل می‌گردد که:

همانطور که قبلاً نیز بدان اشاره شد، هدف این مقاله انشعاب و یا ارائه روش آماری جهت برآورد نبوده، بلکه مشخص کردن یکی از روش‌های استنباط آماری در یک کار عملی می‌باشد.

7) Normalizing Constant

جدول ۱: خلاصه اندازه نیکویی برازش

| | | | RMSE | | | MAD | | |
|----------|----------------|------------|------|------|-----|-----|------|-----|
| | | | Min | Mean | Max | Min | Mean | Max |
| مزرعه | C _۰ | Log-Cauchy | ۷,۳ | ۷,۶ | ۷,۸ | ۲,۷ | ۲,۸ | ۲,۹ |
| | | Log Cosh | ۷,۸ | ۸ | ۸,۲ | ۲,۹ | ۲,۹ | ۳ |
| جنگل | C _۱ | Log-Cauchy | ۶,۵ | ۷,۱ | ۸,۵ | ۲,۴ | ۲,۶ | ۲,۷ |
| | | Log Cosh | ۶,۷ | ۷ | ۸ | ۲,۵ | ۲,۷ | ۲,۹ |
| آب | C _۳ | Log-Cauchy | ۵,۸ | ۶ | ۶,۳ | ۲ | ۲,۲ | ۲,۳ |
| | | Log Cosh | ۶,۳ | ۶,۷ | ۷,۲ | ۲,۲ | ۲,۴ | ۲,۴ |
| شهرنشینی | C _۲ | Log-Cauchy | ۶,۱ | ۶,۳ | ۶,۵ | ۲,۱ | ۲,۲ | ۲,۳ |
| | | Log Cosh | ۶,۲ | ۶,۴ | ۶,۶ | ۲,۳ | ۲,۴ | ۲,۵ |

Rsm_e = میانگین توان دوم خطاها

MAD = قدر مطلق توان دوم خط

همانطور که از جدول می‌توان استنباط کرد، با توجه به استفاده از مدل‌های مختلف توزیع پیشین و برآورد منطقه مورد نظر به طور مثال مشخص کردن منطقه مزرعه با توجه به دو تابع توزیع پیشین به ترتیب مقادیر خطاها ۶، ۷ و یا ۸ می‌باشد و به همین ترتیب برای سایر مناطق، که این می‌تواند روشی جهت پیش‌بینی نقاط غیرمشخص در تصویر باشد.

مراجع

- [1] LANL. DRYDEN and K.V (1998), "STATISTICAL SHAPE ANALYSIS", WILEY
- [2] Goodall, C.R and Lange, N, (1982), "Growth Cure models for Cor-related triangular shapes". In proceeding of the zist INTERFACE Symposium.
- [3] Kendell D.G. (1983), "The shape of poisson-Delaunay triangles" studies in probability and Related topics.
- [4] Mardia, K.V. (1989b) , "Markor models and Bayesian methods in image Amalysis", Journal of Applied statistics, 10: 125-130.
- [5] H.-L. (1991a) "On geodesics in Eoclideam shap space". Journal of the London Mathematical Society, 44:360-372.
- [6] Mardia, K.V. and Dryden, I.L. (1989b) "The statistical analysis of shap data", Biometrika, 76: 71-282.
- [7] Mardia, K.V. and Kanji, G.K. editors (1993), "Statisties and Images", Vol 1, oxford. Carfax.

- [8] Mardia, K.V. , Kent, J.T., and Walder, A.N. (1991). "Statistical shape models in Image analysis" Computer science and statistics.
- [9] Green, P. (1990), "Bayesian Reconstructions from Emission Tomography Data using a Modified EM Algorithm". IEEE Transactions on Medical Imaging.
- [10] H/tchcock, D., Glasbey, A. (1997). "Binary Image Restriction at subpixel Resolution" Biometrics, 1040-1053.
- [11] Frigessi, A. and stamder, J. (1994), "Informative Priors for the Buysian classification of satellite images". JASA, Vol. 89. No. 426

مقایسه روش‌های مختلف نمونه‌گیری با استفاده از الگوریتم‌های مونت کارلو و بوت استرپ

عباس محمدخانی^۱، نصراله ایران‌پناه^۲

^۱ گروه آمار موسسه آموزش عالی جهاد دانشگاهی

^۲ گروه آمار دانشگاه تربیت مدرس

چکیده: در نظریه نمونه‌گیری از جامعه محدود، هدف برآورد پارامترهای مورد نظر مانند میانگین، درصد و نسبت (میانگین دو جامعه) با استفاده از یک روش نمونه‌گیری مناسب و همچنین مقایسه بین روش‌های مختلف نمونه‌گیری از نظر اندازه‌های دقت برآوردگرها مانند اریبی و خطای معیار است.

در این مقاله مقایسه بین روش‌های نمونه‌گیری مانند تصادفی ساده، طبقه‌بندی و خوشه‌ای به همراه برآوردگرهای نسبتی، رگرسیونی، هارتلی - راس و نسبت را با استفاده از دو الگوریتم مونت کارلو و بوت استرپ بر روی یک مثال واقعی (سرشماری دهستانهای کشور در سال ۱۳۷۵) ارائه می‌شود.

در روش مونت کارلو با نمونه‌گیری‌های مکرر از جامعه موجود و در روش بوت استرپ با نمونه‌گیری‌های مکرر از تنها یک نمونه موجود و محاسبه برآورد مورد نظر اندازه‌های دقت برآوردگرها را در روش‌های مختلف نمونه‌گیری برآورد و مقایسه می‌شود. الگوریتم‌های دو روش مونت کارلو و بوت استرپ با استفاده از نرم‌افزار SPLUS برنامه‌نویسی شده است.

همچنین نرم‌افزاری با Visual Basic طراحی گردیده است که روش‌های مختلف نمونه‌گیری را برای برآورد پارامترهای مختلف جامعه به همراه فاصله اطمینان انجام می‌دهد و این روش‌ها را از نظر اندازه دقت با هم مقایسه می‌کند.

واژه‌های کلیدی: مونت کارلو، بوت استرپ، فاصله اطمینان صدکی

۱ مقدمه

نظریه آمار همواره تلاش می‌کند به سه پرسش اساسی زیر پاسخ دهد:

(۱) داده‌ها را چگونه جمع‌آوری شود؟

(۲) داده‌های جمع‌آوری شده را چگونه خلاصه و تحلیل شود؟

(۳) دقت خلاصه داده‌ها چقدر است؟

پرسش ۱ هدف روش‌های نمونه‌گیری و پرسش ۳ از جمله مباحث مهم استنباط آماری است. در نمونه‌گیری از جامعه محدود به حجم N بر اساس اطلاعات موجود از ساختار جامعه، هزینه نمونه‌گیری، سادگی اجرا، هدف مورد بررسی و به خصوص دقت نمونه‌گیری ممکن است روش‌های مختلف نمونه‌گیری قابل اجرا باشد. از جمله روش‌های نمونه‌گیری می‌توان به تصادفی ساده، طبقه‌بندی، خوشه‌ای (یک مرحله‌ای، دو مرحله‌ای و ...)، سیستماتیک و با احتمال متغیر اشاره کرد. میانگین، درصد و نسبت (دو میانگین) از جمله پارامترهای مهم جامعه است که بر اساس روش‌های مختلف نمونه‌گیری به دنبال برآورد مناسب آنها هستیم. اریبی و خطای معیار از مهمترین اندازه‌های دقت برآوردگرها در نظریه نمونه‌گیری هستند.

برآوردگرهای میانگین و درصد نمونه در اکثر روش‌های نمونه‌گیری ناریب و خطای معیار آنها دارای شکل بسته دقیق است. اما برآوردگرهایی مانند نسبت (دو میانگین) نمونه در ساده‌ترین روش نمونه‌گیری یعنی نمونه‌گیری تصادفی ساده، اریب و خطای معیار آن دارای شکل بسته دقیق نمی‌باشد. از شبیه‌سازی مونت کارلو و الگوریتم بوت استرپ می‌توان برای برآورد اریبی و خطای معیار برآوردگرها استفاده کرد.

افزون [۵] و [۶] روش بوت‌استرپ را برای برآورد توزیع و اندازه‌های دقت مانند اریبی و خطای معیار آماره/برآورد ارائه کرد. این روش بر اساس نمونه‌گیری تصادفی ساده با جایگذاری از جامعه است و در روشهای دیگر نمونه‌گیری کاربرد ندارد. بیگل و فریدمن [۴] کاربرد روش بوت‌استرپ را در نمونه‌گیری طبقه‌بندی ارائه کردند.

در بخش ۲ روش‌های مختلف نمونه‌گیری به همراه برآوردگرهای مختلف ارائه می‌شود. در بخش‌های ۳ و ۴ الگوریتم‌های مونت کارلو و بوت استرپ ارائه می‌گردد. در بخش ۵ دو الگوریتم مونت کارلو و بوت استرپ در روش‌های مختلف نمونه‌گیری و در بخش ۶ یک مثال کاربردی برای اجرای دو الگوریتم ارائه می‌گردد و سرانجام در بخش ۷ نرم‌افزار طراحی شده برای انجام و مقایسه روش‌های نمونه‌گیری ارائه می‌شود.

۲ روش‌های نمونه‌گیری

جامعه‌ای محدود به حجم N را با مقادیر صفت (x_1, \dots, x_N) در نظر می‌گیریم. از جمله اهداف مهم نظریه نمونه‌گیری برآورد پارامترهای میانگین جامعه \bar{x}_N به روش‌های مختلف نمونه‌گیری و برآورد اندازه دقت آن است. در زیر مروری بر برآوردگرها و دقت آنها در روش‌های مختلف نمونه‌گیری خواهیم داشت.

(۱) نمونه‌گیری تصادفی ساده: فرض می‌کنیم (X_1, \dots, X_n) یک نمونه تصادفی ساده با

جایگذاری از جامعه محدود (x_1, \dots, x_N) باشد در این صورت:

$$\bar{x}_n = \hat{X}_N = \frac{1}{n} \sum_{i=1}^n X_i, \quad var(\hat{x}_N) = \frac{\sigma^2}{n}, \quad \hat{var}(\hat{x}_N) = \frac{s_n^2}{n}.$$

برآوردگر میانگین جامعه و برآورد واریانس آن ناریب هستند.

۲) نمونه‌گیری طبقه‌بندی: در نمونه‌گیری طبقه‌بندی فرض می‌شود، واحدهای جامعه در k طبقه مجزا قرار دارند (N_i حجم طبقه i ام است، به طوری که $\sum_{i=1}^k N_i = N$). یکی از روش‌های تعیین اندازه نمونه در طبقات روش تخصیص متناسب است. در این روش از هر طبقه به طور مجزا نمونه‌گیری تصادفی ساده به حجم n_i ($n_i = \frac{N_i}{N}n$) استخراج و با وزن مناسب ($W_i = \frac{N_i}{N}$) نتایج طبقات بصورت

$$\begin{aligned} \hat{x}_N &= \bar{X}_{st} = \sum_{i=1}^k W_i \bar{X}_{n_i}, \\ var(\hat{x}_N) &= \sum_{i=1}^k W_i^2 \frac{\sigma_i^2}{n_i}, \\ \hat{var}(\hat{x}_N) &= \sum_{i=1}^k W_i^2 \frac{s_{n_i}^2}{n_i}. \end{aligned}$$

ادغام می‌شوند. برآوردگرهای میانگین جامعه و برآورد واریانس آن ناریب هستند. دقت نمونه‌گیری طبقه‌بندی همواره بیشتر از تصادفی ساده است.

۳) نمونه‌گیری خوشه‌ای دو مرحله‌ای: در نمونه‌گیری خوشه‌ای فرض می‌شود، واحدهای جامعه در N خوشه مجزا قرار دارند (M_i حجم خوشه i ام جامعه است به طوری که $M_0 = \sum_{i=1}^N M_i$ حجم جامعه و $\bar{M} = \frac{M_0}{N}$ متوسط حجم هر خوشه در جامعه است). در نمونه‌گیری خوشه‌ای دو مرحله‌ای، در مرحله اول از بین N خوشه جامعه n خوشه نمونه‌گیری و سپس از داخل n خوشه نمونه، نمونه‌گیری مجددی انجام می‌شود (نمونه‌گیری

در هر دو مرحله تصادفی ساده است). نتایج هر خوشه به طور متناسب به صورت

$$\hat{x}_N = \bar{X}_{cl} = \frac{1}{n} \sum_{i=1}^n \frac{M_i}{\bar{M}} \bar{X}_i,$$

$$Var(\hat{x}_N) = \frac{\sigma_b^2}{n} + \frac{1}{nN} \sum_{i=1}^N \left(\frac{M_i}{\bar{M}}\right)^2 \frac{\sigma_i^2}{m_i},$$

$$\widehat{Var}(\hat{x}_N) = \frac{s_b^2}{n} + \frac{1}{nN} \sum_{i=1}^n \left(\frac{M_i}{\bar{M}}\right)^2 \frac{s_i^2}{m_i}.$$

ادغام می‌شود. که در آن σ_b^2 و s_b^2 به ترتیب واریانس بین میانگین‌های خوشه‌های جامعه و نمونه هستند. برآوردگرهای میانگین جامعه و برآورد واریانس آن ناریب هستند. نمونه‌گیری خوشه‌ای در کاربرد ساده ولی معمولاً دقت بالایی ندارد.

در زیر با استفاده از یک صفت کمکی (Y) که همبستگی خطی بالایی با صفت اصلی (X) دارد و در مورد جامعه آن اطلاع کافی وجود دارد می‌توان با روش نمونه‌گیری تصادفی ساده برآوردهای دقیق‌تری برای میانگین جامعه ارائه کرد. فرض کنید زوجی $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ یک نمونه تصادفی ساده با جایگذاری از جامعه محدود زوجی $\{(x_1, y_1), \dots, (x_n, y_n)\}$ است و در مورد صفت کمکی Y اطلاع کامل وجود دارد.

(۴) برآورد نسبتی (نوع ۱):

$$\hat{x}_N = \bar{X}_{r1} = R_n \cdot \bar{y}_N, \quad (R_n = \frac{\bar{X}_n}{\bar{Y}_n}).$$

(۵) برآورد رگرسیونی:

$$\hat{x}_N = \bar{X}_{re} = \bar{X}_n + \hat{\beta}(\bar{y}_N - \bar{Y}_n), \quad (\hat{\beta} = \frac{S(X, Y)}{S_y^2}).$$

(۶) برآورد نسبتی (نوع ۲):

$$\hat{x}_N = \bar{X}_{r2} = \bar{R}_n \cdot \bar{y}_N, \quad (\bar{R}_n = \frac{1}{n} \sum_{i=1}^n \frac{X_i}{Y_i}).$$

(۷) برآورد هارتلی - راس:

$$\hat{x}_N = \bar{X}_{HR} = \bar{R}_n \cdot \bar{y}_N + \frac{n(N-1)}{N(n-1)} (\bar{X}_n - \bar{R}_n \cdot \bar{Y}_n).$$

(۸) برآورد نسبت (میانگین دو جامعه):

$$r_{\hat{N}} = r_n = \frac{\bar{X}_n}{\bar{Y}_n}.$$

از بین برآوردگرهای ۴ تا ۸ تنها برآوردگر هارتلی - راس نارایب است و سایر برآوردگرها مجاناً نارایب هستند. همچنین واریانس برآوردگرهای ۴ تا ۸ جواب دقیق ندارد و تنها جواب جانبی دارند. برای آشنایی بیشتر با روش‌های نمونه‌گیری به مرجع [۲] رجوع کنید.

۳ شبیه‌سازی مونت کارلو

فرض کنید X_1, \dots, X_n متغیرهای تصادفی مستقل و هم توزیع (*i.i.d.*) با تابع توزیع تجمعی معلوم F و T آماره (برآوردگر) مورد نظر باشد. با استفاده از شبیه‌سازی مونت کارلو می‌توان توزیع، اریبی، خطای معیار و ... T را بصورت مراحل زیر تقریب کرد.

(۱) نمونه تصادفی مستقل و هم‌توزیع (X_1, \dots, X_n) از توزیع F بدست می‌آید.

(۲) آماره T بر روی نمونه مرحله ۱ بصورت $T \equiv T(X_1, \dots, X_n)$ محاسبه می‌گردد.

(۳) مراحل ۱ و ۲ بار تکرار می‌شود و T_1, \dots, T_B بدست می‌آید.

(۴) برای B ‌های بزرگ تابع توزیع، اریبی و خطای معیار T بصورت

$$F_T(t) \simeq \frac{1}{B} \sum_{i=1}^B I(T_i \leq t)$$

$$bias(T) \simeq \frac{1}{B} \sum_{i=1}^B T_i - \theta$$

$$SE(T) \simeq \left\{ \frac{1}{B} \sum_{i=1}^B [T_i - \frac{1}{B} \sum_{j=1}^B T_j]^2 \right\}^{\frac{1}{2}}$$

با هر میزان دقتی تقریب می‌شود. در نمونه‌گیری از جامعه محدود به حجم N نیز فرض می‌شود (X_1, \dots, X_n) یک نمونه تصادفی ساده (باجایگذاری یا بدون جایگذاری) از جامعه (x_1, \dots, x_N) است، می‌توان توزیع و مشخصات T را مشابه الگوریتم بالا بدست آورد. فقط به جای نمونه‌گیری به حجم n از F در مرحله ۱، نمونه‌گیری از (x_1, \dots, x_N) انجام می‌شود.

۴ الگوریتم بوت استرپ

فرض کنید X_1, \dots, X_n متغیرهای تصادفی مستقل و هم‌توزیع (*i.i.d.*) با تابع توزیع تجمعی نامعلوم F و یا (X_1, \dots, X_n) یک نمونه تصادفی ساده با جایگذاری از جامعه محدود و نامعلوم (x_1, \dots, x_n) باشد. با استفاده از روش بوت استرپ می‌توان توزیع، اربیی، خطای معیار و ... آماره $T \equiv T(X_1, \dots, X_n)$ را با داشتن تنها یک نمونه مشاهده شده (x_1, \dots, x_n) بصورت مراحل زیر تقریب کرد.

۱) نمونه بوت استرپ به حجم n ، (X_1^*, \dots, X_n^*) بصورت تصادفی ساده با جایگذاری از (x_1, \dots, x_n) بدست می‌آید.

۲) آماره T بر روی نمونه بوت استرپ مرحله ۱ بصورت $T^* \equiv T(X_1^*, \dots, X_n^*)$ محاسبه می‌گردد.

۳) مراحل ۱ و ۲ B بار تکرار می‌شود و T_1^*, \dots, T_B^* بدست می‌آید.

۴) برای B ‌های بزرگ تابع توزیع، اربیی و خطای معیار T بصورت

$$\hat{F}_T(t) = \frac{1}{B} \sum_{i=1}^B I(T_i^* \leq t),$$

$$\hat{bias}(T) = \frac{1}{B} \sum_{i=1}^B T_i^* - T,$$

$$\hat{SE}(T) = \left\{ \frac{1}{B} \sum_{i=1}^B [T_i^* - \frac{1}{B} \sum_{j=1}^B T_j^*]^2 \right\}^{\frac{1}{2}},$$

برآورد می‌گردد. برای آشنایی بیشتر با الگوریتم بوت استرپ کلاسیک به افرون و تییشیرانی [۷] و در یک مرجع فارسی به [۱] رجوع کنید.

۵ الگوریتم مونت کارلو و بوت استرپ در روشهای نمونه‌گیری

این دو الگوریتم برای ۳ روش نمونه‌گیری تصادفی ساده، طبقه‌بندی و خوشه‌ای دو مرحله‌ای بصورت زیر خلاصه می‌شود. قابل ذکر است که در الگوریتم مونت کارلو باید کل جامعه (به حجم N) موجود ولی در الگوریتم بوت استرپ تنها یک نمونه (به حجم n) کافی است.

۱) نمونه‌گیری تصادفی ساده: دو الگوریتم برای برآورد نسبت (میانگین دو جامعه) حالت ۸) ارائه می‌شود. برای سایر برآوردهای ۱ و ۴ تا ۷ بطور مشابه قابل تکرار است.

الف) الگوریتم مونت کارلو:

از جامعه $\{(x_1, y_1), \dots, (x_N, y_N)\}$ به روش نمونه‌گیری تصادفی ساده با جایگذاری نمونه $\{(x_1, y_1), \dots, (x_n, y_n)\}$ بدست می‌آید و سپس برآورد نسبت $r_n = \frac{\bar{x}_n}{\bar{y}_n}$ محاسبه می‌گردد. اگر این عمل B بار تکرار شود، آنگاه $(r_{n,1}, \dots, r_{n,B})$ بدست می‌آید. توزیع، اریبی و خطای معیار برآوردگر R_n به ازای B های بزرگ با جایگذاری برآوردگر R_n به جای T در بخش ۳ قابل تقریب است.

ب) الگوریتم بوت استرپ:

از جامعه زوجی به حجم N تنها یک نمونه تصادفی ساده با جایگذاری به حجم n ، $\{(x_1, y_1), \dots, (x_n, y_n)\}$ موجود است. به روش نمونه‌گیری تصادفی ساده با جایگذاری نمونه بوت استرپ به حجم n ، $\{(x_1^*, y_1^*), \dots, (x_n^*, y_n^*)\}$ بدست می‌آید. برآورد نسبت $r_n^* = \frac{\bar{x}_n^*}{\bar{y}_n^*}$ در نمونه بوت استرپ محاسبه می‌شود. اگر این عمل B بار تکرار شود، آنگاه $(r_{n,1}^*, \dots, r_{n,B}^*)$ بدست می‌آید. توزیع، اریبی و خطای معیار برآوردگر R_n به ازای B های بزرگ با جایگذاری برآوردگر R_n^* به جای T^* در بخش ۴ قابل تقریب است.

(۲) نمونه‌گیری طبقه‌بندی:

الف) الگوریتم مونت کارلو: با داشتن واحدهای جامعه به حجم N که در k طبقه مجزا قرار دارند، اگر B بار نمونه‌گیری به حجم n مشابه نمونه‌گیری طبقه‌بندی بخش ۲ تکرار شود، آنگاه $(\bar{x}_{st,1}, \dots, \bar{x}_{st,B})$ بدست می‌آید. توزیع، اریبی و خطای معیار برآوردگر \bar{X}_{st} به ازای B های بزرگ با جایگذاری برآوردگر \bar{X}_{st} به جای T در بخش ۳ قابل تقریب است. ب) الگوریتم بوت استرپ: فرض کنید نمونه‌ای به حجم n (از هر طبقه بطور مجزا نمونه‌ای به حجم n_i به روش تصادفی ساده با جایگذاری ساده با جایگذاری بدست آمده است و $n = \sum_{i=1}^k n_i$) بدست آمده است. از هر طبقه بطور مجزا نمونه‌ای به حجم n_i به روش تصادفی ساده با جایگذاری استخراج و میانگین کل که حاصل یک نمونه بوت استرپ است بصورت $\bar{x}_{st}^* = \sum_{i=1}^k W_i \bar{x}_{n_i}^*$ محاسبه می‌شود. اگر این عمل B بار تکرار شود، آنگاه $(\bar{x}_{st,1}^*, \dots, \bar{x}_{st,B}^*)$ بدست می‌آید. توزیع، اریبی و خطای معیار برآوردگر \bar{X}_{st} به ازای B های بزرگ با جایگذاری برآوردگر \bar{X}_{st}^* به جای T^* در بخش ۴ قابل تقریب است.

(۳) نمونه‌گیری خوشه‌ای دومرحله‌ای:

الف) الگوریتم مونت کارلو: با داشتن واحدهای جامعه به حجم M که در N خوشه مجزا قرار دارند، اگر B بار نمونه‌گیری به حجم m ، مشابه نمونه‌گیری خوشه‌ای دومرحله‌ای بخش ۲ تکرار شود، آنگاه $(\bar{x}_{cl,1}, \dots, \bar{x}_{cl,B})$ بدست می‌آید. توزیع، اریبی و خطای معیار برآوردگر \bar{X}_{cl} به ازای B های بزرگ با جایگذاری برآوردگر \bar{X}_{cl} به جای T در بخش ۳ قابل تقریب است.

ب) الگوریتم بوت استرپ: فرض کنید از بین N خوشه جامعه n خوشه و در داخل هر خوشه نمونه، نمونه‌گیری مجددی به حجم m_i (هر دو مرحله نمونه‌گیری تصادفی ساده با جایگذاری است) بدست آمده است. از بین n خوشه نمونه به روش تصادفی ساده با جایگذاری n خوشه انتخاب و در داخل هر خوشه نمونه، نمونه‌گیری مجدد با حجم m_i به روش تصادفی ساده با جایگذاری انتخاب می‌شود. میانگین کل که حاصل یک نمونه بوت استرپ است بصورت $\bar{x}_{cl}^* = \frac{1}{n} \sum_{i=1}^n \frac{M_i}{M} \bar{x}_i^*$ محاسبه می‌شود. اگر این عمل B بار تکرار شود، آنگاه $(\bar{x}_{cl,1}^*, \dots, \bar{x}_{cl,B}^*)$ بدست می‌آید. توزیع، اریبی و خطای معیار برآوردگر \bar{X}_{cl} به ازای B های بزرگ با جایگذاری برآوردگر \bar{X}_{cl}^* به جای T^* در بخش ۴ قابل تقریب است.

۶ یک مثال کاربردی

دو الگوریتم مونت کارلو و بوت استرپ در روشهای مختلف نمونه‌گیری و برآوردهای مختلف در مثال کاربردی زیر اجرا می‌شود.

در سر شماری عمومی نفوس و مسکن سال ۱۳۷۵ از مناطق روستایی کشور [۳] جامعه را $N = ۲۲۱۳$ دهستان در نظر گرفته و هر دهستان بصورت یک واحد نمونه‌گیری در نظر گرفته می‌شود. دهستانهای کشور در $k = ۲۶$ استان تقسیم گردیده‌اند. دو صفت اصلی X تعداد جمعیت با سواد و صفت کمکی Y تعداد خانوار هر دهستان تعریف می‌شود. واحد اندازه‌گیری هزار نفر برای دو صفت تعریف می‌شود. با توجه به جامعه تعریف شده

$$\bar{x}_N = ۶/۱۷۳, \quad \bar{y}_N = ۱/۹۹۳, \quad \sigma_x = ۵/۶۵۱, \quad \sigma_y = ۱/۶۵۷, \quad \rho_{xy} = ۰/۹۸,$$

است. در جدول ۱، مقدار برآورد اریبی و خطای معیار برآوردگر میانگین جامعه و همچنین یک فاصله اطمینان صدکی $۰/۹۵$ برای میانگین جامعه به دو روش مونت کارلو و بوت استرپ بر اساس نمونه‌ای تصادفی به حجم $n = ۵۰۰$ نشان داده شده است. نمونه‌گیری به سه روش تصادفی ساده با جایگذاری، طبقه بندی تخصیص متناسب و خوشه‌ای دومرحله‌ای و همچنین برآورد نسبتی و رگرسیونی در نمونه‌گیری تصادفی ساده با جایگذاری است.

در نمونه‌گیری طبقه بندی تعداد طبقات $k = ۲۶$ استان و در نمونه‌گیری خوشه‌ای تعداد خوشه‌ها $N = ۲۶$ استان در نظر گرفته می‌شود. در نمونه‌گیری خوشه‌ای از بین $N = ۲۶$ خوشه تعداد $n = ۱۰$ خوشه انتخاب و حجم هر خوشه در نمونه (m_i) نیز متناسب با حجم هر خوشه انتخاب شده است. قابل ذکر است، برای انجام دو روش نمونه‌گیری طبقه بندی و خوشه‌ای و برای مقایسه بین روشها تعداد نمونه بزرگ $n = ۵۰۰$ از جامعه انتخاب شده است. نتایج حاصل از شبیه‌سازی مونت کارلو با روابط دقیق نظری (و یا مجانبی) برای مقادیر اریبی و خطای معیار برابر است. نتایج حاصل از الگوریتم بوت استرپ نیز حاکی از دقت بالای روش در مقایسه با الگوریتم مونت کارلو است.

جدول ۱: مقایسه اربیبی، خطای معیار و فاصله اطمینان برآورد میانگین جامعه در روش‌های مختلف نمونه‌گیری ($n = 500$) با دو الگوریتم مونت کارلو و بوت استرپ.

| فاصله اطمینان | خطای معیار | | اربیبی | | نمونه‌گیری |
|----------------------------|------------|------------|-----------|------------|-----------------|
| | بوت استرپ | مونت کارلو | بوت استرپ | مونت کارلو | |
| (بوت استرپ) (۶,۱۹,۶,۷۳) | ۰,۲۵۰ | ۰,۲۵۲ | ۰,۶۰۰ | ۰,۶۰۰ | تصادفی ساده |
| (۵,۶۲,۶,۵۸) | ۰,۲۳۹ | ۰,۲۳۶ | ۰,۶۰۰۲ | ۰,۶۰۰۰ | طبقه بندی |
| (۲,۷۳,۸,۱۸) | ۱,۳۸۸ | ۱,۵۰۴ | ۰,۶۰۰۲ | ۰,۶۰۰۰ | خوشه‌ای |
| (۶,۱۲,۶,۲۳) | ۰,۰۵۰ | ۰,۰۵۷ | ۰,۶۰۰۰ | ۰,۶۰۰۰ | برآورد نسبی |
| (۶,۰۹,۶,۲۹) | ۰,۰۴۷ | ۰,۰۵۰ | ۰,۶۰۰۰ | ۰,۶۰۰۰ | برآورد رگرسیونی |

جدول ۲: مقایسه اریبی، خطای معیار و فاصله اطمینان برآوردهای مختلف میانگین و نسبت جامعه در یک نمونه تصادفی ساده ($n = 10$) با دو الگوریتم مونت کارلو و بوت استرپ.

| فاصله اطمینان | | خطای معیار | | اریبی | | برآورد معمولی |
|---------------|--------------|------------|------------|-----------|------------|------------------|
| بوت استرپ | مونت کارلو | بوت استرپ | مونت کارلو | بوت استرپ | مونت کارلو | |
| (۴,۷۰,۸,۵۳) | (۳,۶۰,۱۰,۶۹) | ۰,۹۸۲ | ۱,۷۹۶ | ۰,۰۰۷ | ۰,۰۰۱ | نسبتی ۱ |
| (۵,۶۴,۶,۴۲) | (۵,۴۲,۶,۸۷) | ۰,۳۷۶ | ۰,۳۷۴ | ۰,۰۳۱ | ۰,۰۲۷ | نسبتی ۲ |
| (۵,۶۳,۶,۵۸) | (۵,۲۲,۷,۰۴) | ۲,۴۲۹ | ۴,۶۵۲ | ۰,۱۹۷ | ۰,۲۵۳ | هارتلی راس |
| (۵,۶۴,۶,۴۲) | (۵,۳۸,۷,۰۹) | ۰,۹۲۱ | ۱,۳۶۶ | ۰,۰۱۰ | ۰,۰۰۶ | رگرسیون |
| (۵,۵۵,۶,۴۴) | (۵,۳۹,۶,۹۰) | ۰,۲۲۷ | ۰,۳۸۹ | ۰,۰۲۳ | ۰,۰۳۱ | نسبت |
| (۲,۸۳,۳,۲۲) | (۲,۷۲,۳,۴۵) | ۰,۱۹۸ | ۰,۱۸۸ | ۰,۰۰۵ | ۰,۰۱۴ | |

در جدول ۲، با استفاده از تنها یک نمونه کوچک به حجم $n = 10$ به روش تصادفی ساده با جایگذاری برآورد مقدار اریبی و خطای معیار برآوردگر میانگین جامعه و یک فاصله اطمینان صدکی 95% به روش مونت کارلو و بوت استرپ نشان داده شده است. برآورد میانگین جامعه به روشهای معمولی، نسبتی نوع ۱ و ۲، هارتلی - راس، رگرسیونی و برآورد نسبت (میانگین دو جامعه) است.

قابل ذکر است که به جز برآوردگر معمولی میانگین جامعه برای برآوردگرهای دیگر مقدار اریبی و خطای معیار جواب دقیقی نظری ندارد و به صورت شبیه سازی مونت کارلو این مقادیر دقیق محاسبه گردیده است. نتایج حاصل از الگوریتم بوت استرپ حاکی از دقت بالای روش در مقایسه با الگوریتم مونت کارلو است. قابل توجه است که در الگوریتم مونت کارلو بازنمونه‌گیری به حجم $n = 10$ از جامعه $N = 2213$ و در الگوریتم بوت استرپ بازنمونه‌گیری به حجم $n = 10$ از تنها یک نمونه به حجم $n = 10$ بدست آمده است.

۷ طراحی نرم افزار روش های نمونه گیری

با استفاده از Visual Basic نرم افزاری تهیه گردیده است که برآورد میانگین، درصد، مقدار کل، تعداد کل و نسبت (میانگین دو جامعه) را به روشهای مختلف نمونه‌گیری تصادفی ساده (با و بدون جایگذاری)، طبقه‌بندی، خوشه‌ای، سیستماتیک و با احتمال متغیر محاسبه می‌کند.

در بانک اطلاعاتی نرم افزار مقادیر دو صفت اصلی و کمکی X و Y در جامعه‌ای به حجم N قرار می‌گیرد. در نرم افزار به روشهای مختلف نمونه‌گیری، نمونه‌های به حجم n بصورت تصادفی استخراج می‌شود. در نرم افزار علاوه بر برآورد میانگین، درصد، مقدار کل، تعداد کل و نسبت (میانگین دو جامعه)، برآورد مقدار واریانس برآوردگر و همچنین فاصله اطمینان با ضریب اطمینان دلخواه نیز محاسبه می‌شود.

هدف از طراحی این نرم افزار مقایسه روشهای مختلف نمونه‌گیری با استفاده مثالهای واقعی از جامعه است. این نرم افزار در نسخه‌های مختلف Windows بصورت خودکار (Auto Run) قابل نصب و اجرا است. این نرم افزار از نویسنده اول قابل دریافت است.

مراجع

[۱] ایران پناه، نصراله و پاشا، عین‌اله (۱۳۷۶). آشنایی با الگوریتم بوت استرپ. اندیشه آماری، سال دوم، شماره ۱، ۳۳-۴۶.

[۲] عمیدی، علی (۱۳۷۸). نظریه نمونه‌گیری و کاربردهای آن. جلد اول و دوم، مرکز نشر دانشگاهی.

[۳] شناسنامه دهستانهای کشور، سرشماری عمومی نفوس و مسکن، ۱۳۷۵. جلدهای ۱-۲۶،
جداول ۴ و ۵.

- [4] Bickel, P. J. and Freedman, D. A. (1984), Asymptotic Normality and the Bootstrap in Stratified Sampling, *Annals of Statistics*, **12**, 470-482.
- [5] Efron, B. and Tibshirani, R. (1993), *An Introduction to the Bootstrap*, Chapman and Hall, New York.
- [6] Efron, B. (1979), *The Jackknife, the Bootstrap and Other Resampling Plans*, Society for Industrial and Applied Mathematics, Philadelphia, PA.
- [7] Efron, B. and Tibshirani, R. (1993), *An Introduction to the Bootstrap*, Chapman and Hall, New York.

پیشگویی فضایی بیزی با استفاده از روشهای مونت کارلو

محسن محمدزاده، مجید جعفری خالدی

گروه آمار دانشگاه تربیت مدرس

چکیده: اغلب، تجزیه و تحلیل داده‌های فضایی با فرض گوسی بودن میدان تصادفی صورت می‌پذیرد. اما در بعضی مسائل با داده‌های مثبتی مواجه‌ایم که لگاریتم آنها گوسی یا تقریباً گوسی است. در این مقاله با فرض آنکه میدان تصادفی تبدیل یافته، گوسی با ساختار میانگین و کوواریانس پارامتری است و پارامترها متغیرهایی تصادفی با توزیع بیشین مناسب هستند پیشگویی فضایی با استفاده از رهیافت بیزی تعیین می‌شود. چون در مسائل کاربردی معمولاً محاسبه این پیشگو دشوار است بکمک روشهای مونت کارلو، یک تقریب برای توزیع پیشگو ارائه و براساس آن پیشگویی فضایی بیزی تقریبی معرفی خواهد شد.

واژه‌های کلیدی: داده‌های فضایی، میدان تصادفی لگ - گوسی، پیشگویی فضایی بیزی، روشهای مونت کارلو

۱ مقدمه

برای مدل‌بندی داده‌های فضایی معمولاً فرض می‌شود داده‌ها تحقیقی از یک میدان تصادفی $\{Z(t); t \in D \subseteq R^d\}$ ، $d \geq 1$ ، هستند، که در آن D یک مجموعه اندیس‌گذار است. پیشگویی فضایی مقدار نامعلوم میدان تصادفی در یک یا چند موقعیت مشخص و بر اساس مشاهدات بدست آمده در سایر موقعیتها، در علوم مختلف از جمله هواشناسی، محیط زیست، همه‌گیرشناسی، زمین‌شناسی، جغرافیا، زمین‌شناسی، جغرافی، کشاورزی و . . . که بنوعی با داده‌های فضایی سر و کار دارند، یکی از مسائل مهم بشمار می‌رود. اغلب با فرض گوسی بودن میدان تصادفی، پیشگویی فضایی تعیین می‌شود. اما در بعضی مسائل با داده‌های مثبتی مواجه‌ایم که از مدل گوسی تبعیت نمی‌کنند، بلکه تبدیل لگاریتمی آنها گوسی یا تقریباً گوسی است. در این صورت میدان تصادفی مورد نظر لگ - گوسی^۱ نامیده می‌شود. بعنوان مثال در هواشناسی مشاهدات مربوط به میزان بارندگی عموماً از این نوع بشمار می‌روند. در این مقاله با فرض آنکه میدان تصادفی تبدیل یافته، گوسی با ساختار میانگین و کوواریانس پارامتری است و پارامترها متغیرهایی تصادفی با توزیع بیشین مناسب هستند، پیشگویی فضایی با استفاده از رهیافت بیزی تعیین می‌شود. اما در عمل محاسبه این پیشگو دشوار است، لذا در این مقاله با

1) Log-Gaussian

استفاده از روشهای مونت کارلو، توزیع تقریبی پیشگو و پیشگوی فضایی بیزی ارائه خواهند شد. در سالهای اخیر بدنبال توسعه روشهای مونت کارلو در آمار بیزی، استفاده از روشهای بیزی برای تجزیه و تحلیل داده‌های فضایی رشد و توسعه چشمگیری یافته است. برگر و همکاران (۲۰۰۱) توزیع‌های پیشینی که پسین‌های نامناسب^۲ نتیجه می‌دهند، را مطالعه کردند. برای حالتی که با داده‌های گمشده مواجه‌ایم، روش بیزی ارائه شده توسط اوج و همکاران (۲۰۰۲) می‌تواند رهگشا باشد. آلدورث و کرسی (۲۰۰۳) پیشگویی بیز تجربی توابع غیر خطی از میدان تصادفی با مشاهدات نویز دار را بررسی کردند. هنگامی که مشاهدات فاقد نویز هستند، جعفری و محمدزاده (۲۰۰۳) پیشگوی فضایی بیزی را برای یک میدان تصادفی گوسی ارائه کردند. جعفری و محمدزاده (۲۰۰۴a) نیز پیشگویی برای میدانهای تصادفی گوسی با مشاهدات نویز دار را ارائه نمودند. سپس جعفری و محمدزاده (۲۰۰۴b) پیشگوی فضایی بیزی برای یک میدان تصادفی که تبدیلی نامعلوم از آنها گوسی باشد را با استفاده از توزیع‌های پیشین نامناسب تعیین کردند که بدنبال آن در این مقاله با استفاده از توزیع‌های پیشین مناسب و بکمک تکنیک مونت کارلو پیشگوی فضایی بیزی برای میدان تصادفی لگ - گوسی ارائه می‌شود.

۲ پیشگویی فضایی بیزی

فرض کنید $\{Z(t) ; t \in D\}$ یک میدان تصادفی از متغیرهای تصادفی مثبت و بردار $Z = (Z(t_1), \dots, Z(t_n))$ میدان تصادفی را در موقعیتهای t_1, \dots, t_n نمایش دهد و پیشگویی مقدار نامعلوم میدان تصادفی در نقطه t_0 ، یعنی $Z(t_0)$ ، براساس بردار مشاهدات $z = (z(t_1), \dots, z(t_n))$ مورد نظر باشد. با فرض آنکه میدان تصادفی

$$\{Y(t) = \ln Z(t) ; t \in D\}$$

گوسی حقیقی مقدار بترتیب با توابع میانگین و کوواریانس

$$\begin{aligned} E(Y(t)) &= f'(s)\beta \\ \text{Cov}(Y(s), Y(t)) &= \sigma^2 \rho(s, t; \theta) \end{aligned}$$

باشد، که در آن $f(t) = (f_1(t), \dots, f_q(t))'$ یک بردار $p \times 1$ از مولفه‌های غیرتصادفی معلوم، $\beta = (\beta_1, \dots, \beta_p) \in R^p$ پارامترهای رگرسیون، $\sigma^2 = \text{Var}(Y(t))$ واریانس و $\theta = (\theta_1, \dots, \theta_q) \in \Theta \subseteq R^q$ با بردار پارامترهای $Y(\cdot)$ تابع همبستگی میدان $\rho(\cdot, \cdot; \theta)$ است. بدین ترتیب

$$Y = \ln Z = (\ln Z(t_1), \dots, \ln Z(t_n)) \sim N_n(X\beta, \sigma^2 \Sigma_\theta)$$

است، که در آن ماتریس طرح $X_{n \times q}$ و ماتریس همبستگی فضایی نمونه Σ_θ بصورت

$$X = \begin{pmatrix} f'(t_1) \\ \vdots \\ f'(t_n) \end{pmatrix}, \quad \Sigma_\theta = (\rho(t_i, t_j; \theta))$$

تعریف می‌شوند. با فرض اینکه X ماتریسی رتبه کامل و Σ_θ ماتریسی معین مثبت باشد، تابع درستنمایی بعد از حذف ضرایب ثابت بصورت

$$L(\phi; z) = \left(\frac{1}{\sqrt{2\sigma^2\pi}}\right)^n |\Sigma_\theta|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2} (\ln z - X\beta)' \Sigma_\theta^{-1} (\ln z - X\beta)\right\} J$$

خواهد بود، که در آن $J = (\prod_{i=1}^n |z_i|)^{-1}$ ژاکوبین تبدیل و $\phi = (\beta, \sigma^2, \theta)$ بردار $p + q + 1$ بعدی پارامترهای مدل است. با توجه به گوسی بودن میدان تصادفی $Y(\cdot)$ ، توزیع توأم $(Y(t_0), Y)$ نرمال $n + 1$ متغیره است و می‌توان نشان داد توزیع $Y(t_0)$ بشرط z و ϕ ، نرمال بترتیب با میانگین و واریانس

$$\begin{aligned} E(Y(t_0)|z, \phi) &= (f'(t_0)\beta + r'_\theta \Sigma_\theta^{-1} (\ln z - X\beta)) \\ \text{Var}(Y(t_0)|z, \phi) &= \sigma^2 (1 - r'_\theta \Sigma_\theta^{-1} r_\theta) \end{aligned}$$

خواهد بود، که در آن $r_\theta = (\rho(t_0, t_i; \theta))_{n \times 1}$ است. بنابراین بشرط z و ϕ ، توزیع پیشگو برای $Z(t_0) = \exp(Y(t_0))$ بصورت

$$\begin{aligned} f(z_0|z, \phi) &= \left(\frac{1}{\sqrt{2\pi\sigma^2 z_0^2 (1 - r'_\theta \Sigma_\theta^{-1} r_\theta)}}\right)^{\frac{1}{2}} \\ &\times \exp\left\{-\frac{(\ln z_0 - E(Y(t_0)|z, \phi))^2}{2\sigma^2 (1 - r'_\theta \Sigma_\theta^{-1} r_\theta)}\right\} \quad z_0 > 0 \quad (1) \end{aligned}$$

خواهد بود. با فرض معلوم بودن پارامتر ϕ ، پیشگوی بهینه برای $Z(t_0)$ با تابع زیان قدر مطلق خطا، میانه توزیع (۱) یعنی

$$\text{Median of } [Z(t_0)|z, \phi] = \exp\{f'(t_0)\beta + r'_\theta \Sigma_\theta^{-1} (\ln z - X\beta)\}$$

می‌باشد. چون پیشگوی بهینه به پارامتر $p + q + 1$ بعدی ϕ بستگی دارد، هنگامی که نامعلوم است، پیشگوی فضایی به روش بیزی تعیین می‌گردد. بدین منظور توزیعهای پیشین مناسب برای پارامترها اختیار می‌شود و به جهت ساده‌تر شدن محاسبات، پارامترهای (β, σ^2) از پارامتر θ مستقل در نظر گرفته می‌شوند. بنابراین توزیع پیشین توأم پارامترهای مدل را می‌توان بفرم

$$\pi(\phi) = \pi(\beta, \sigma^2, \theta) = \pi(\beta|\sigma^2)\pi(\sigma^2)\pi(\theta)$$

نوشت. اکنون توزیعهای پیشین حاشیه‌ای

$$(\beta|\sigma^2) \sim N(\beta_0, \sigma^2 V_0) \quad \sigma^2 \sim \chi_{SCT}^2(a, b)$$

را در نظر می‌گیریم، که مقادیر β_0, V_0, a و b معلوم هستند و χ_{SCT}^2 نشاندهنده توزیع کای دوی معکوس می‌باشد. توزیع پیشین $\pi(\theta)$ نیز بطور دلخواه اختیار می‌شود. چون تعیین توزیع پسین پارامترهای مدل از رابطه

$$\pi(\beta, \sigma^2, \theta|z) \propto f(z|\phi)\pi(\beta|\sigma^2)\pi(\sigma^2)\pi(\theta) \quad (2)$$

دشوار می‌باشد، ابتدا توزیع پسین حاشیه‌ای β را با شرطی کردن بر پارامترهای σ^2 و θ تعیین می‌کنیم. بسادگی می‌توان نشان داد این توزیع بفرم

$$(\beta|z, \sigma^2, \theta) \sim N(\hat{\beta}_\gamma, \sigma^2 \hat{V}) \quad (3)$$

است، که در آن

$$\begin{aligned} \hat{\beta}_\gamma &= (V_0^{-1} + X'\Sigma_\theta^{-1}X)^{-1}(V_0^{-1}\beta_0 + X'\Sigma_\theta^{-1}y) \\ \hat{V} &= (V_0^{-1} + X'\Sigma_\theta^{-1}X)^{-1} \end{aligned}$$

می‌باشند. سپس توزیع پسین حاشیه‌ای σ^2 با شرطی کردن بر پارامتر θ بصورت

$$(\sigma^2|z, \theta) \sim \chi_{SCT}^2(a + n, b') \quad (4)$$

تعیین می‌شود، که در آن

$$b' = \frac{ab + n\hat{\sigma}^2 + \hat{\beta}'V_\beta^{-1}\hat{\beta} + \beta_0'V_0^{-1}\beta_0 - \mathbf{m}'\hat{V}\mathbf{m}}{a + n}$$

و $\hat{\sigma}^2, \hat{\beta}, V_\beta$ و \mathbf{m} بترتیب برابر

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n}(y - X\beta)' \Sigma_\theta^{-1} (y - X\beta) \\ \hat{\beta} &= (X'\Sigma_\theta^{-1}X)^{-1} X'\Sigma_\theta^{-1}y \\ V_\beta &= (X'\Sigma_\theta^{-1}X)^{-1} \\ \mathbf{m} &= (V_\beta^{-1}\hat{\beta} + V_0^{-1}\beta_0) \end{aligned}$$

می‌باشند. در مرحله بعد توزیع پسین حاشیه‌ای θ بصورت

$$\pi(\theta|z) \propto \frac{f(z|\phi)\pi(\beta|\sigma^2)\pi(\sigma^2)\pi(\theta)}{\pi(\beta|z, \sigma^2, \theta)\pi(\sigma^2|z, \theta)} \quad (5)$$

تعیین می‌شود. اکنون بر اساس قانون بیز و با جایگذاری توزیعهای پسین حاشیه‌ای (۳)، (۴) و (۵) در رابطه

$$\pi(\phi|z) = \pi(\beta, \sigma^2, \theta|z) = \pi(\beta|z, \sigma^2, \theta)\pi(\sigma^2|z, \theta)\pi(\theta|z) \quad (6)$$

توزیع پسین پارامترهای مدل قابل تعیین می‌باشد. از آنجا که عموماً تعیین $\pi(\theta|z)$ و به تبع آن توزیع پسین (۶) بسیار دشوار است، با استفاده از روشهای مونت کارلوی زنجیر مارکف (MCMC)^۳، یک نمونه تصادفی از توزیعهای (۵) و (۶) بدست آورده می‌شود و استنباط بیزی پارامترهای مدل براساس این نمونه صورت می‌پذیرد (جعفری و محمدزاده، ۲۰۰۳).
 اکنون می‌توان توزیع پیشگوی بیزی $Z(t_0)$ را از رابطه

$$f(z_0|z) = \int f(z_0|z, \phi)\pi(\phi|z)d\phi \quad (7)$$

بدست آورد. با توجه به آنکه محاسبه تحلیلی توزیع پیشگو از رابطه (۷) مشکل می‌باشد، ابتدا با انتگرال‌گیری از $f(z|\phi)\pi(\beta|z, \sigma^2, \theta)$ نسبت به β ، می‌توان نشان داد توزیع $Z(t_0)$ بشرط z, σ^2, θ - نرمال بصورت

$$(Z(t_0)|z, \sigma^2, \theta) \sim LN(\mu_\lambda, \Sigma_\lambda)$$

است، که در آن

$$\begin{aligned} \mu_\lambda &= (f'(t_0) - r'_\theta \Sigma_\theta^{-1} X)(V_0^{-1} + X' \Sigma_\theta^{-1} X)^{-1} V_0^{-1} \beta \\ &+ [r'_\theta \Sigma_\theta^{-1} + (f'(t_0) - r'_\theta \Sigma_\theta^{-1} X)(V_0^{-1} + X' \Sigma_\theta^{-1} X)^{-1} X' \Sigma_\theta^{-1}] y \\ \Sigma_\lambda &= \sigma^2 (\lambda - r'_\theta \Sigma_\theta^{-1} r) \\ &+ [(f'(t_0) - r'_\theta \Sigma_\theta^{-1} X)(V_0^{-1} + X' \Sigma_\theta^{-1} X)^{-1} (f'(t_0) - r'_\theta \Sigma_\theta^{-1} X)'] \end{aligned}$$

سپس با انتگرال‌گیری از $\pi(\sigma^2|z, \theta)f(z(t_0)|z, \sigma^2, \theta)$ نسبت به σ^2 ، توزیع $Z(t_0)$ بشرط z, θ - تی استیودنت بصورت

$$(Z(t_0)|z, \theta) \sim LT_{n+a}(\mu_\lambda, b' \Sigma_\lambda)$$

3) Monte Carlo Markov Chain

بدست می‌آید. اکنون می‌توان توزیع پیشگوی بیزی را بصورت

$$f(z_0|z) = \int \pi(\theta|z)f(z_0|z, \theta)d\theta \quad (۸)$$

نوشت، که یک انتگرال q گانه و تعیین آن ساده‌تر از رابطه (۷) می‌باشد. با این وجود محاسبه تحلیلی آن امکان‌پذیر نیست، لذا روش مونت کارلو جهت تقریب آن مورد استفاده قرار می‌گیرد. برای این منظور، از توزیع پیشین $\pi(\theta)$ یک نمونه تصادفی $i.i.d$ به حجم M بصورت $\{\theta_i\}_{i=1}^M$ انتخاب می‌شود. بدین ترتیب، برای مقادیر به اندازه کافی بزرگ M می‌توان

$$f(z_0|z) \approx \frac{\sum_{i=1}^M f(z(t_0)|z, \theta_i)f(z|\theta_i)}{\sum_{i=1}^M f(z|\theta_i)}$$

را بعنوان توزیع تقریبی پیشگو مورد استفاده قرار داد و بر اساس تابع زیان قدر مطلق خطا پیشگوی فضایی بیزی تقریبی برای $Z(t_0)$ میانه توزیع تقریبی پیشگو، یعنی

$$\hat{Z}(t_0) = \text{Median of } [Z(t_0)|z]$$

تعیین می‌گردد.

۳ بحث و نتیجه‌گیری

برای یک میدان تصادفی لگ - گوسی بر اساس رهیافت بیزی و بکمک روشهای مونت کارلو پیشگوی فضایی بیزی تعیین گردید. لازم بذکر است که روابط بیان شده برای میدانهای تصادفی که تبدیلی معلوم از آنها گوسی است، قابل توسعه می‌باشد. در این مقاله پیشگویی فضایی بیزی با فرض اینکه داده‌ها فاقد نویز هستند، مورد مطالعه قرار گرفت. اما در اغلب مسائل آمار فضایی، بدلائل مختلف از جمله خطای اندازه‌گیری، مشاهدات همراه با نویز می‌باشند. انجام این نوع پیشگویی بر اساس مشاهدات نویدار بعنوان یک مسئله جدید قابل بررسی است.

مراجع

- [1] Aldworth, J. and Cressie, N. (2003). Prediction of Nonlinear Spatial Functionals. Statistical Planning and Inference (Article in Press).
- [2] Berger, J. O., Deoliveira, V. and Sanso, B. (2001). Objective Bayesian Analysis of Spatially Correlated Data. JASA 96: 1351-1370.

- [3] Jafari, K. M. and Mohammadzadeh, M. (2003). Bayesian Spatial Prediction for a Gaussian Random Field. Proceedings of the Fourth Seminar of the Probability and Stochastic Processes.
- [4] Jafari, K. M. and Mohammadzadeh, M. (2004a). Bayesian Spatial Prediction for a Gaussian Random Field with Noisy Observations. To appear in Journal of Science, Isfahan University, Iran.
- [5] Jafari, K. M. and Mohammadzadeh, M. (2004b). Bayesian Spatial Prediction for a Transformed Random Field. To appear in Journal of Science, Tehran University, Iran.
- [6] Oh, M.S., Shin, W. D. and Kim, J. H. (2002). Bayesian Analysis of Regression Models with Spatially Corelated Errors and Missing Observation. Computational Statistics and Data Analysis, 39: 387-400.

پیش‌بینی بیزی طول عمر برای مدل پارتو با حجم نمونه تصادفی

محسن محمدزاده^۱، منصور زرگر^۲^۱ دانشگاه تربیت مدرس، گروه آمار^۲ مؤسسه آموزش عالی نجف آباد، گروه آمار

چکیده: یکی از مسائل مهم کاربردی، پیش‌بینی طول عمر یا زمانهای شکست $n - r$ مؤلفه در یک نمونه n تایی بر اساس زمانهای شکست r مؤلفه ترتیبی اول نمونه می‌باشد. چون در اینگونه مسائل اغلب با مواردی مواجه می‌شویم که بدلائل مختلف، از جمله محدودیت‌های شرایط آزمایش، حجم نمونه ثابت باقی نمی‌ماند، لازم است پیش‌بینی طول عمر بر اساس حجم نمونه متغیر صورت پذیرد. در این مقاله برای یک مدل طول عمر پارتو پیش‌بینی بیزی در حالتی که حجم نمونه یک متغیر تصادفی با توزیع پواسن یا دوجمله‌ای باشد، مورد بررسی قرار می‌گیرد. سپس فواصل پیش‌بینی بر اساس دو حالت حجم نمونه ثابت و تصادفی با تکنیک شبیه‌سازی مورد مقایسه عددی قرار گرفته، نشان داده می‌شود فواصل پیش‌بینی که بر اساس حجم نمونه تصادفی بدست می‌آیند بطور متوسط دارای طول کوتاهتری هستند و می‌توان با این روش پیش‌بینی‌های دقیقتری را برای زمان شکست هر یک از مؤلفه‌های $r + 1$ به بعد نمونه بدست آورد.

واژه‌های کلیدی: مدل طول عمر پارتو، پیش‌بینی بیزی، اندازه نمونه تصادفی

۱ مقدمه

توزیع پارتو بوسیله پارتو برای توزیع درآمد معرفی شد. اهمیت این توزیع به کاربردهای آن در بسیاری از مطالعات اقتصادی اجتماعی به قرن نوزدهم به بعد مربوط می‌شود. توزیع پارتو نقش مهمی در بررسی اندازه جمعیت شهرها، حوادث با منابع طبیعی، تغییرات قیمت سهام، توزیعهای درآمد، ریسک بیمه، شکستهای شغلی، خطای خوشه‌بندی کردن در مدارهای ارتباطی و غیره بازی کرده است. آرنولد (۱۹۸۲) یک بررسی تاریخی گسترده‌ای از کاربرد آن در زمینه توزیع درآمد را می‌دهد.

در این مقاله به پیش‌بینی در مدل طول عمر پارتو، که یکی از مدل‌های مهم طول عمر است، می‌پردازیم و پیش‌بینی را برای مشاهدات آینده در یک نمونه انجام می‌دهیم (پیش‌بینی تک نمونه‌ای). فرض کنید n شیء را مورد آزمایش عمر قرار داده‌ایم و پس از بدست آمدن r زمان شکست اول آزمون را خاتمه داده مقادیر طول عمر آنها بصورت مقادیر ترتیبی در دسترس هستند (داده‌های سانسور شده نوع II). اکنون قصد داریم بر اساس این r مشاهده ترتیبی برای طول عمر مشاهدات آینده در همین نمونه یعنی $n - r$ حالت باقیمانده پیش‌بینی داشته باشیم. برای

مثال فرض کنید در یک شرکت بیمه n شغل بیمه شده است و پس از گذشت ۶ ماه تعداد r تا از این شغلها ورشکسته شده‌اند و برای گرفتن خسارت به شرکت بیمه مراجعه کرده‌اند. پس مدت زمان فعالیت این r شغل تا زمان ورشکستگی در دست می‌باشد و همچنین فرض کنید مدت زمان فعالیت شغلها تا ورشکستگی از یک مدل طول عمر مانند مدل طول عمر پارتو پیروی می‌کند. حال می‌خواهیم بر اساس این r زمان شکست که بصورت آماره‌های ترتیبی هستند، زمان شکست خوردن $n - r$ حالت باقیمانده را پیش‌بینی کنیم و فاصله اطمینانی برای پیش‌بینی مشاهدات بعدی بدست آوریم و همچنین مهمتر از همه اینکه، زمان شکست خوردن آخرین شغل یعنی $X_{(n)}$ را پیش‌بینی کنیم.

مسئله مهمی که باید مورد توجه قرار گیرد آنست که در اکثر کاربردها، اندازه نمونه ثابت فرض می‌شود. اما مسأله استنباط آماری در مواقعی که اندازه نمونه یک متغیر تصادفی است، خیلی مهم است، چرا که در عمل نمونه‌های با اندازه تصادفی، در یک حالت طبیعی، بطور فراوان اتفاق می‌افتد. حتی در بعضی کاربردها داشتن اندازه نمونه ثابت غیر ممکن است. زیرا همیشه بعضی از مشاهدات، به دلایل مختلف از بین می‌روند. برای مثال در بعضی از آزمایشات کشاورزی و زیست شناسی، تعدادی از عضوهای نمونه تحت بررسی به دلایل نامعلومی می‌میرند. پس مجبوریم اندازه نمونه را تصادفی فرض کنیم که در اینجا اندازه نمونه یک متغیر تصادفی و دارای یکی از توزیعهای دو جمله‌ای یا پواسن می‌باشد. که در همین بخش بطور توأم، بعد از بدست آوردن فواصل پیش‌بینی با اندازه نمونه ثابت به پیش‌بینی با اندازه نمونه تصادفی می‌پردازیم. فرض کنید n شیء را مورد آزمون عمر قرار داده‌ایم، که طول عمر آنها از یک توزیع پارتو با تابع توزیع زیر پیروی می‌کند.

$$F(x|\alpha, \sigma) = 1 - \left(\frac{\sigma}{x}\right)^\alpha \quad x \geq \sigma > 0, \alpha > 0 \quad (1)$$

این تابع توزیع دارای دو پارامتر شکل α و مقیاس σ می‌باشد. و تابع بقا بصورت:

$$\bar{F}(x) = \left(\frac{\sigma}{x}\right)^\alpha, \quad x \geq \sigma$$

تعریف می‌شود. و تابع چگالی احتمالی آن برابر است با:

$$f(x|\alpha, \sigma) = \frac{\alpha \sigma^\alpha}{x^{\alpha+1}} \quad (2)$$

اکنون اگر r مشاهده ترتیبی اول در دسترس باشند، بر اساس این r آماره ترتیبی اول از یک نمونه تصادفی از توزیع پارتو (۲)، یعنی $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(r)}$ ، تابع درستنمایی برای $X = (X_{(1)}, X_{(2)}, \dots, X_{(r)})$ برابر است با:

$$L(\alpha, \sigma; \mathbf{x}) \propto \prod_{i=1}^r f(x_{(i)}) \{1 - F(x_{(r)})\}^{n-r}$$

که برای توزیع پارتو داریم:

$$L(\alpha, \sigma; x) \propto \alpha^r \exp\{-\alpha\{\ln(x_{(r)}^{n-r} \prod_{i=1}^r x_{(i)}) - n \ln \sigma\}\} \quad (3)$$

اگر $Y = X_{r+s}$ مولفه مرتبه $(r+s)$ ام باشد و بعنوان s امین مشاهده بعدی تلقی گردد، آنگاه چگالی شرطی Y بشرط x عبارتست از:

$$f(y|x) = \frac{f(y, x)}{f(x)} = \frac{s \binom{n-r}{s} \{F(y) - F(x_{(r)})\}^{s-1} \{1 - F(y)\}^{n-r-s} f(y)}{\{1 - F(x_{(r)})\}^{n-r}} \quad (4)$$

$y > x_{(r)}$.
و این تابع چگالی شرطی برای توزیع پارتو عبارتست از:

$$f(y|x) = \frac{\alpha s \binom{n-r}{s}}{y} \sum_{j=s}^{s-1} (-1)^j \binom{s-1}{j} \left(\frac{x_{(r)}}{y}\right)^{\alpha(n-r-s+j+1)} \quad y > x_{(r)} \quad (5)$$

۲ پیش‌بینی بیزی در مدل طول عمر پارتو

از آنجا که توزیع پارتو دارای دو پارامتر می‌باشد، پیش‌بینی در ۳ حالت زیر انجام می‌گیرد.
الف) حالت پارامتر شکل نامعلوم
ب) حالت پارامتر شکل و مقیاس نامعلوم
ج) حالت بی‌اطلاعی

۱.۲ حالت پارامتر شکل نامعلوم

در حالتی که σ معلوم و α نامعلوم است، اگر در تابع درست‌نمایی (۳)، $r = n$ قرار دهیم آنگاه آماره $Z = \ln \prod_{i=1}^n \left(\frac{X_i}{\sigma^n}\right)$ برای α بسنده و دارای توزیع $\Gamma(n, \alpha)$ خواهد بود. بنابراین یک خانواده پیشین مزدوج طبیعی برای α را خانواده گاما در نظر می‌گیریم. این مطلب را آرنولد و پرس (۱۹۸۳) مورد بررسی قرار دادند. بنابراین چگالی پیشین بصورت

$$\pi(\alpha|a, b) = \frac{\alpha^{a-1}}{\Gamma(a)} b^a \exp(-b\alpha) \quad \alpha > 0, a > 0, b > 0 \quad (6)$$

و چگالی پسین α بصورت

$$\begin{aligned} \pi(\alpha|x, b) &\propto \pi(\alpha|a, b) \cdot f(x|\alpha, \sigma) \\ &= \frac{\alpha^{r+a-1}}{\Gamma(r+a)} H^{r+a} \exp(-\alpha H) \quad \alpha > 0 \end{aligned} \quad (7)$$

خواهد بود، که در آن $H = b + \ln(\sigma^{-n} x_{(r)}^{n-r} \prod_{i=1}^r x_i)$ است. با ضرب معادلات (5) و (7) و انتگرالگیری روی α ، چگالی پیش‌بینی Y بشرط x موقعی که اندازه نمونه ثابت است بصورت

$$\begin{aligned} f(y|x, n) &= \int_{\alpha} f(y|x, \alpha) \pi(\alpha|x, b) d\alpha \\ &= \int_{\alpha} f(y, \alpha|x) d\alpha \\ &= \frac{s(r+a)H^{r+a}}{y} \sum_{j=0}^{s-1} \frac{(-1)^j \binom{n-r}{s} \binom{s-1}{j}}{\{H + a_j \ln(\frac{y}{x_{(r)}})\}^{r+a+1}} \end{aligned} \quad (8)$$

بدست می‌آید، که در آن $a_j = n - r - s + j + 1$ می‌باشد. تابع بقاء پیش‌بینی Y به شرط x, n عبارتست از:

$$P_r(Y \geq t|x, n) = sH^{r+a} \cdot \sum_{j=0}^{s-1} \frac{(-1)^j \binom{n-r}{s} \binom{s-1}{j}}{a_j \{H + a_j \ln(\frac{t}{x_{(r)}})\}^{r+a}} = \tau \quad (9)$$

حال اگر $\tau = (1+\gamma)/2$ و $\tau = (1-\gamma)/2$ ، آنگاه مرزهای بالایی و پایینی برای $y = x_{(r+s)}$ یعنی s امین مشاهده بعدی موقعی که اندازه نمونه ثابت است، بدست خواهد آمد. حال فرض کنید که اندازه نمونه N ، یک متغیر تصادفی بواسن با تابع جرم احتمال زیر باشد:

$$p(n) = \frac{e^{-\lambda} \cdot \lambda^n}{n!}, \quad n = 0, 1, 2, \dots \quad (10)$$

با استفاده از کنسول (۱۹۸۴) و گوپتا و گوپتا (۱۹۸۴) چگالی پیش‌بینی Y بشرط x برای موقعی که اندازه نمونه تصادفی است، بشکل

$$\begin{aligned} f(y|x) &= \frac{1}{P_r(N \geq r+s)} \cdot \sum_{n=r+s}^{\infty} f(y|x, n) P(n) \\ &= \frac{cs(r+a)}{y} \cdot \sum_{n=r+s}^{\infty} \Omega(s, j) \frac{\lambda^n H^{r+a}}{n! \{H + a_j \ln(\frac{y}{x_{(r)}})\}^{r+a+1}} \end{aligned} \quad (11)$$

است، که در آن $c = \frac{e^{-\lambda}}{1 - \sum_{m=0}^{r+s+1} \frac{e^{-\lambda} \lambda^m}{m!}}$ و $a_j = n - r - s + j + 1$ و $\Omega(s, j) = \binom{n-r}{s} \sum_{j=0}^{s-1} (-1)^j \binom{s-1}{j}$ و تابع بقاء پیش‌بینی برای Y بشرط x عبارتست از:

$$P_r(Y \geq t|x) = cs \cdot \sum_{n=r+s}^{\infty} \frac{\lambda^n H^{r+a} \Omega(s, j)}{a_j n! \{H + a_j \ln(\frac{t}{x_{(r)}})\}^{r+a}} \quad (۱۲)$$

برای بدست آوردن فاصله پیش‌بینی 100γ درصد برای y ، روشهای عددی تکراری لازم هستند، که با پیدا کردن t از معادله (۱۲) برای یک مقدار $\tau = P_r(Y \geq t)$ داده شده با $\tau = (1 + \gamma)/2$ و $\tau = (1 - \gamma)/2$ ، قابل محاسبه هستند. برای پیش‌بینی زمان شکست اولین مشاهده بعدی یعنی $y = x_{(r+1)}$ داریم:

$$P_r(Y \geq t|x) = \frac{e^{-\lambda}}{1 - \sum_{m=0}^r \frac{e^{-\lambda} \lambda^m}{m!}} \cdot \sum_{n=r+1}^{\infty} \frac{\lambda^n}{n!} \left\{ 1 + \frac{(n-r) \ln(\frac{t}{x_{(r)}})}{H} \right\}^{-(r+a)} \quad (۱۳)$$

و اگر اندازه نمونه N یک متغیر تصادفی از توزیع دوجمله‌ای با پارامترهای M و p باشد، آنگاه

$$f(y|\sigma, x) = \frac{ks(r+a)}{y} \cdot \sum_{n=r+s}^M \Omega(s, j) \frac{\binom{n}{M} p^n q^{M-n} H^{r+a}}{\{H + a_j \ln(\frac{y}{x_{(r)}})\}^{r+a+1}}$$

که در آن $k^{-1} = 1 - \sum_{W=0}^{r+s-1} \binom{M}{W} p^W q^{M-W}$ و تابع بقاء پیش‌بینی

$$P_r(Y \geq t|\sigma, x) = ks \sum_{n=r+s}^M \Omega(s, j) \frac{\binom{M}{n} p^n q^{M-n} H^{r+a}}{a_j \{H + a_j \ln(\frac{t}{x_{(r)}})\}^{r+a}} \quad (۱۴)$$

می‌باشد. با قرار دادن $s = 1$ ، در معادله (۱۴)، تابع بقاء پیش‌بینی $y = x_{(r+1)}$ عبارتست از:

$$P_r(x_{(r+1)} \geq t|x) = \frac{1}{1 - \sum_{W=0}^r \binom{M}{W} p^W q^{M-W}} \cdot \sum_{n=r+1}^M \frac{\binom{M}{n} p^n q^{M-n} H^{r+a}}{\{H + (n-r) \ln(\frac{t}{x_{(r)}})\}^{r+a}} \quad (۱۵)$$

همچنین برای بدست آوردن فاصله پیش‌بینی y ، روشهای عددی تکراری لازم هستند.

۲.۲ حالت پارامترهای مقیاس و شکل نامعلوم

در این بخش، در حالتی که هر دو پارامتر مقیاس و شکل مجهولند، فاصله پیش‌بینی محاسبه می‌شود. در این حالت $X_{(1)} = \text{Min}(X_1, X_2, \dots, X_n)$ و $\sum_{i=1}^r \ln X_i = \ln \prod_{i=1}^r X_i$ برای α و σ ، آماره‌های بسنده توأم مینیمال هستند. در نتیجه توزیع پیشین مزدوج توأم طبیعی تعمیم یافته توسط آرنولد و پرس (۱۹۸۳) را بعنوان پیشین برای α و σ استفاده می‌کنیم، که در آن فرض می‌شود α دارای توزیع گاما است و توزیع شرطی σ بشرط α از شکل تابع توانی می‌باشد. این پیشین بفرم

$$\pi(\sigma, \alpha) = g_1(\sigma|\alpha) \cdot g_2(\alpha)$$

است، که در آن

$$g_1(\sigma|\alpha) = \zeta \alpha \sigma^{\zeta \alpha - 1} \sigma^{-\zeta \alpha}, \quad \sigma_0 < \sigma < \sigma_1$$

و

$$g_2(\alpha) = \frac{\alpha^{\delta-1}}{\Gamma(\delta)} \{\ln \mu - \zeta \ln \sigma_0\}^\delta \cdot \exp\{-\alpha(\ln \mu - \zeta \ln \sigma_0)\}$$

یعنی

$$g_2(\alpha) \sim \Gamma(\delta, \ln \mu - \zeta \ln \sigma_0)$$

بنابراین

$$\pi(\sigma, \alpha) = A^{-1} \cdot \alpha^\delta \cdot \sigma^{\zeta \alpha - 1} \cdot \mu^{-\alpha}, \quad \sigma_0 < \sigma < \sigma_1, \quad \alpha > 0 \quad (16)$$

و پارامترهای $\delta, \zeta, \mu, \sigma_0$ مثبتند و A تعدیل‌کننده ثابتی است که بوسیله فرمول زیر تعیین می‌شود

$$A = \zeta^{-1} \cdot \Gamma(\delta) \cdot \{\ln \mu - \zeta \ln \sigma_0\}^{-\delta}$$

پیشین (۱۶) توزیع گاما-توانی نامیده می‌شود و با نماد $P\Gamma(\delta, \zeta, \mu, \sigma_0)$ نشان داده می‌شود و نتایج بدست آمده در بخش قبل را می‌توان با در نظر گرفتن $\delta = \zeta = 0, \mu = 1, \sigma_0 = \infty$ بدست آورد.

با بکارگیری قضیه بییزی و با استفاده از تابع درست‌نمایی (۳)، ملاحظه می‌شود که چگالی پسین σ, α بشرط x در همان خانواده گاما-توانی با پارامترهای $(r + \delta, n + \zeta, \mu(x_{(r)}^{n-r} \prod_{i=1}^r x_i), \text{Min}(\sigma_0, x_1))$ قرار دارد. بنابراین:

$$\pi(\sigma, \alpha|x) = A'^{-1} \alpha^{r+\delta} \cdot \left(\frac{\sigma^{n+\zeta-\frac{1}{\alpha}}}{\mu \cdot x_{(r)}^{n-r} \cdot \prod_{i=1}^r x_i} \right)^\alpha = \pi(\sigma, \alpha) \cdot f(x|\sigma, \alpha)$$

$$\alpha > 0, \quad \sigma_0 < \sigma < L^*$$

بعد از مقداری محاسبات ریاضی، چگالی پسین می‌تواند بشکل

$$\pi(\sigma, \alpha|x) = \frac{n + \zeta}{\Gamma(r + \delta)} \cdot (\alpha u)^{r+\delta} \cdot \exp(-\alpha v) \quad \alpha > 0, \quad 0 < \sigma < L^* \quad (17)$$

نوشته می‌شود که در آن

$$u = \ln\left\{(\mu \cdot x_{(r)}^{n-r} \cdot \prod_{i=1}^r x_i) / L^{*n+\zeta}\right\}$$

و

$$v = \ln\left\{(\mu \cdot x_{(r)}^{n-r} \cdot \prod_{i=1}^r x_i) / \sigma^{n+\zeta-\frac{1}{\alpha}}\right\}$$

با ضرب معادلات (۵) و (۱۷) و انتگرالگیری روی σ, α ، تابع چگالی پیش‌بینی بیزی Y بشرط x با اندازه نمونه ثابت بشکل زیر می‌باشد:

$$f(y|x) = \frac{s(r + \delta) \binom{n-r}{s} u^{r+\delta}}{y} \sum_{j=0}^{s-1} \frac{(-1)^j \binom{s-1}{j}}{\{u + a_j \ln(\frac{y}{x_{(r)}})\}^{r+\delta+1}} \quad y > x_{(r)} \quad (18)$$

چگالی پیش‌بینی Y موقعی که اندازه نمونه N یک متغیر تصادفی با تابع جرم احتمال (۱۰) می‌باشد، برابر است با:

$$f(y|x) = \frac{cs(r + \delta)}{y} \cdot \sum_{n=r+s}^{\infty} \Omega(s, j) \cdot \frac{\lambda^n \cdot u^{r+\delta}}{n! \{u + a_j \ln(\frac{y}{x_{(r)}})\}^{r+\delta+1}} \quad (19)$$

و تابع بقا پیش‌بینی عبارتست از:

$$P_r(Y \geq t|x) = cs \sum_{n=r+s}^{\infty} \Omega(s, j) \frac{1}{n! a_j} \left\{ \frac{u}{u + a_j \ln(\frac{t}{x_{(r)}})} \right\}^{r+\delta} \quad (20)$$

برای $s = 1$ ، تابع بقا پیش‌بینی $Y = X_{(r+1)}$ عبارتست از:

$$P_r(x_{(r+1)} \geq t|x) = \frac{e^{-\lambda}}{1 - \sum_{m=0}^r \frac{e^{-\lambda} \lambda^m}{m!}} \cdot \sum_{n=r+1}^{\infty} \frac{\lambda^n}{n!} \cdot \left\{ \frac{u}{u + a_j \ln(\frac{t}{x_{(r)}})} \right\}^{r+\delta} \quad (21)$$

دوباره با این فرض که اندازه نمونه N ، یک متغیر تصادفی با توزیع دوجمله‌ای است، تابع چگالی پیش‌بینی Y بصورت

$$f(y|x) = \frac{ks.(r + \delta)}{y} \cdot \sum_{n=r+s}^M \Omega(s, j) \frac{\binom{M}{n} p^n q^{M-n} u^{r+\delta}}{\{u + a_j \ln(\frac{y}{x_{(r)}})\}^{r+\delta+1}} \quad (22)$$

است و تابع بقاً پیش‌بینی عبارتست از:

$$P_r(Y \geq t|x) = ks \sum_{n=r+s}^M \Omega(s, j) \frac{\binom{M}{n} p^n q^{M-n}}{a_j} \cdot \left\{ \frac{u}{u + a_j \ln(\frac{t}{x_{(r)}})} \right\}^{r+\delta} \quad (23)$$

با قراردادن $s = 1$ در معادله (۲۳)، فاصله پیش‌بینی 100γ درصد برای $x_{(r+1)}$ ، جوابهای t از معادلات

$$\frac{1}{1 - \sum_{W=0}^r \binom{M}{W} p^W q^{M-W}} \sum_{n=r+1}^M \binom{M}{n} p^n q^{M-n} \left\{ \frac{u}{u + (n-r) \ln(\frac{t}{x_{(r)}})} \right\}^{r+\delta} = \tau$$

بدست می‌آید که در آن $\tau = \frac{1-\gamma}{\gamma}$ و $\tau = \frac{1+\gamma}{\gamma}$ می‌باشند.

۳.۲ بدون اطلاعات قبلی از α, σ

یک پیشین بدون اطلاع (NIP) پیشینی است که اثراتش بر همه مقادیر پارامترها بی‌اهمیت است. و هیچ اطلاعی به آنچه بوسیله داده‌های تجربی بدست آمده است، نمی‌افزاید. بنابراین، عموماً، استنباط بیزی بر اساس یک NIP فقط جنبه نظری دارد و اینچنین استنباط بیزی، در حل مسائل پیش‌بینی، موقعی که پیدا کردن جوابهای کلاسیک فوق العاده مشکل می‌باشد، مهم است. زیرا تقریب بیزی بر اساس یک NIP می‌تواند شبیه یک روش ریاضی برای بدست آوردن فاصله پیش‌بینی کلاسیک عمل کند. با استفاده از استدلال باکس و تیو (۱۹۷۳) و نیو و همدی (۱۹۸۷) می‌توان نشان داد که NIP برای پارامترهای α, σ بشکل زیر است:

$$\pi(\sigma, \alpha) \propto \frac{1}{\sigma \cdot \alpha}, \quad \sigma > 0 \quad (24)$$

چگالی پسین توأم α, σ بشرط x عبارتست از:

$$\begin{aligned} \pi(\sigma, \alpha|x) &= f(x|\alpha, \sigma) \cdot \pi(\sigma, \alpha) \\ &= \frac{n(\alpha z)^{r-1} \cdot \sigma^{n\alpha-1}}{\Gamma(r-1) \cdot x_{(r)}^{(n-r)\alpha} \cdot \prod_{i=1}^r x_i^\alpha}, \quad \sigma < x_{(1)} \end{aligned} \quad (25)$$

که در آن $z = \sum_{i=1}^r \ln x_i + \ln(x_{(r)}^{-r}/x_{(1)}^n)$ همانند بخش قبلی می توان با ضرب معادلات (۵) و (۲۵) و انتگرالگیری روی σ و α ، چگالی پیش بینی $y = x_{(r+s)}$ را بصورت

$$f(y|x) = \frac{cs(r-1)}{y} \cdot \sum_{n=r+s}^{\infty} \Omega(s, j) \cdot \frac{z^{r-1}}{n! \{z + a_j \ln(y/x_{(r)})\}^r} \quad (26)$$

بدست آورد، که متناظراً تابع بقاء پیش بینی بفرم

$$P_r(Y \geq t|x) = cs \sum_{n=r+s}^{\infty} \Omega(s, j) \cdot \frac{\lambda^n z^{r-1}}{a_j \cdot n! \{z + a_j \ln(t/x_{(r)})\}^{r-1}} \quad (27)$$

خواهد بود. برای $s = 1$ ، تابع بقا پیش بینی $x_{(r+1)}$ عبارتست از:

$$P_r(x_{(r+1)} \geq t|x) = \frac{e^{-\lambda}}{1 - \sum_{m=0}^r \frac{e^{-\lambda} \lambda^m}{m!}} \cdot \sum_{n=r+1}^{\infty} \frac{\lambda^n}{n!} \left\{ 1 + \frac{(n-r) \ln(\frac{t}{x_{(r)}})}{z} \right\}^{r-1} \quad (28)$$

دوباره اگر اندازه نمونه N ، دارای توزیع دوجمله ای باشد، تابع چگالی پیش بینی Y بشرط x بصورت

$$f(y|x) = \frac{ks(r-1)}{y} \sum_{n=r+s}^M \Omega(s, j) \frac{\binom{M}{n} p^n q^{M-n} z^{r-1}}{\{z + a_j \ln(\frac{y}{x_{(r)}})\}^r} \quad (29)$$

خواهد شد و تابع بقاء پیش بینی عبارتست از:

$$P_r(Y \geq t|x) = ks \sum_{n=r+s}^M \Omega(s, j) \frac{\binom{M}{n} p^n q^{M-n} z^{r-1}}{a_j \{z + a_j \ln(\frac{t}{x_{(r)}})\}^{r-1}} \quad (30)$$

برای $s = 1$ ، تابع بقاء پیش بینی $x_{(r+1)}$ برابر است با:

$$P_r(x_{(r+1)} \geq t|x) = \frac{1}{1 - \sum_{W=0}^r \binom{M}{W} p^W q^{M-W}} \cdot \sum_{n=r+1}^{\infty} \frac{\binom{M}{n} p^n q^{M-n} z^{r-1}}{\{z + (n-r) \ln(\frac{t}{x_{(r)}})\}^{r-1}} \quad (31)$$

مثال ۱: حالت پارامتر شکل نامعلوم

فرض کنید مدت زمانی که یک شغل فعالیت می‌کند تا شکست بخورد بر حسب سال توزیع پارتو با $\sigma = 1$ دارد و $1/01, 1/05, 1/08, 1/14, 1/28, 1/30, 1/33, 1/43, 1/59$ و $1/62$ طول عمر مفید 10° نمونه تصادفی از اینچنین شغلی باشند. اگر فقط شش زمان شکست مرتب شده اول در اختیار باشند، علاقمند به محاسبه فاصله پیش‌بینی برای طول عمر اولین شغل از $n - r$ تای باقیمانده هستیم (یعنی $x(r)$). جدول ۱، فواصل پیش‌بینی ۹۵ درصد برای $x(r)$ را نشان می‌دهد که با حل کردن معادلات (۹)، (۱۳) و (۱۵) بطور تکراری بدست آمده‌اند. با این فرض که توزیع پیشین α ، گاما با مقادیر مختلف $a = 1, 5$ و $b = 1, 2, 5$ است. برای یک اندازه نمونه تصادفی، با توزیعهای پواسن و درجمله‌ای با پارامترهای مختلف $\lambda = 1, 2^\circ, 5^\circ$ و $(1/1, 0/5, 5^\circ, 0/5), (2^\circ, 0/8), (M, P)$ ، محاسبات را انجام می‌دهیم. مقادیر مختلف λ, M, P و پارامترهای پیشین برای بررسی حساسیت فواصل به این پارامترها، انتخاب شده‌اند.

مثال ۲: حالت پارامتر مقیاس و شکل نامعلوم

برای تشریح نتایج بخش ۲.۲، از همان داده‌های مثال ۱ استفاده می‌کنیم (که α, σ نامعلومند). و با حل معادلات (۱۸)، (۲۱)، (۲۳)، (۲۸) و (۳۱) بطور عددی فاصله پیش‌بینی ۹۵ درصد برای $x(r)$ را با استفاده از همان مقادیر λ, M, P ، بدست می‌آوریم. پارامتر توزیعهای پیشین (شامل NIP) با مقادیر مختلف $L^* = x(1)$ و $\delta = 1, 5, 10$ و $\mu = 3, 7, 30$ و $\zeta = 2, 6, 20$ بکار گرفته شده‌اند و نتایج در جدول ۲ آمده است.

۳ مقایسه پیش‌بینی بیزی با اندازه‌های نمونه ثابت و تصادفی

۱.۳ پیش‌بینی در توزیع پارتو با پارامترهای α نامعلوم و σ معلوم

فرض کنید طول عمر اشیاء مورد آزمایش از توزیع پارتو با تابع توزیع

$$F(x|\alpha, \sigma) = 1 - \left(\frac{\sigma}{x}\right)^\alpha \quad \circ < \sigma \leq x, \alpha > \circ$$

و تابع چگالی احتمال

$$f(x|\alpha, \sigma) = \frac{\alpha\sigma^\alpha}{x^{\alpha+1}}$$

پیروی می‌کنند، که در آن پارامتر α نامعلوم و σ معلوم می‌باشد. اندازه نمونه اولیه $n = 10$ در نظر گرفته می‌شود و پس از بدست آمدن طول عمر $r = 6$ شیء اول آزمایش خاتمه می‌یابد (داده‌های سانسور شده نوع II). اکنون قصد داریم با فرض ثابت بودن اندازه نمونه، برای اولین مشاهده بعدی، یعنی $x = x(r)$ بر اساس r مشاهده اول، پیش‌بینی فاصله‌ای داشته باشیم (پیش‌بینی تک نمونه‌ای). در روش شبیه‌سازی ابتدا باید ۶ عدد از توزیع پارتو، بعنوان مقادیر مشاهده شده

جدول ۱: فاصله پیش‌بینی ۹۵ درصد برای $x(y)$ موقعی که α نامعلوم و σ معلوم

| پیشین | مرزها | اندازه نمونه ثابت ($n = 10$) | نتایج برای یک اندازه نمونه تصادفی (پواسن) با پارامتر λ | | نتایج برای یک اندازه نمونه تصادفی (دوجمله‌ای) با پارامترهای زیر M, P مقادیر | | | | |
|-------|-------|-----------------------------------|--|----------------|---|----------|-----------|-----------|-----------|
| | | | $\lambda = 1$ | $\lambda = 20$ | $M = 20$ | $M = 50$ | $P = 0/8$ | $P = 0/5$ | $P = 0/1$ |
| ۱ ۲ | LB | ۱,۳۱ | ۱,۳۳ | ۱,۳۱ | ۱,۳۰ | ۱,۳۱ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ |
| | UB | ۲,۱۷ | ۶,۱۳ | ۱,۷۰ | ۱,۵۵ | ۱,۷۵ | ۱,۶۲ | ۱,۵۵ | ۱,۵۵ |
| ۱ ۵ | LB | ۱,۳۲ | ۱,۳۵ | ۱,۳۱ | ۱,۳۰ | ۱,۳۱ | ۱,۳۱ | ۱,۳۰ | ۱,۳۰ |
| | UB | ۳,۲۳ | ۲۸,۲۰ | ۱,۹۴ | ۱,۶۱ | ۲,۰۷ | ۱,۷۷ | ۱,۶۱ | ۱,۶۱ |
| ۵ ۱ | LB | ۱,۳۰ | ۱,۳۱ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ |
| | UB | ۱,۶۲ | ۲,۴۰ | ۱,۴۸ | ۱,۴۳ | ۱,۵۰ | ۱,۴۶ | ۱,۴۳ | ۱,۴۳ |

جدول ۲: فاصله پیش‌بینی ۹۵ درصد برای $x_{(v)}$ و موقعی که α, σ نامعلومند

| پیش‌بینها با $L^* = x_{(1)}$ | مرزها | اندازه نمونه ثابت ($n = 10$) | نتایج برای یک اندازه نمونه تصادفی (پوسن) با پارامتر λ | | نتایج برای یک اندازه نمونه تصادفی (دو جمله‌ای) با پارامترهای زیر مقادیر M, P | |
|---------------------------------|---------|---|---|---------------|--|----------------|
| | | | $\lambda = 1$ | $\lambda = 2$ | $\lambda = 5$ | $\lambda = 10$ |
| δ | ζ | μ | $M = 2$ | $M = 5$ | $M = 50$ | $M = 500$ |
| | | | $P = 0/8$ | $P = 0/5$ | $P = 0/5$ | $P = 0/1$ |
| ۱ | ۲ | ۳ | ۱,۳۲ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ |
| | | | ۳,۷۰ | ۱,۶۱ | ۱,۶۵ | ۱,۵۲ |
| ۵ | ۶ | ۷ | ۱,۳۲ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ |
| | | | ۳,۰۶ | ۱,۵۱ | ۱,۵۳ | ۱,۴۳ |
| ۱۰ | ۲۰ | ۳۰ | ۱,۳۲ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ |
| | | | ۲,۹۶ | ۱,۴۷ | ۱,۴۹ | ۱,۴۰ |
| NIP | | | ۱,۳۱ | ۱,۳۰ | ۱,۳۰ | ۱,۳۰ |
| | | | ۲,۷۹ | ۱,۶۸ | ۱,۶۵ | ۱,۶۲ |

تولید کنیم. برای این امر ابتدا از توزیع یکنواخت $(U(0, 1))$ ، U مقدار u_1, u_2, \dots, u_6 را تولید می‌کنیم، سپس بر اساس قضیه احتمال انتگرال باگرفتن $F_X(X) = U$ و بدست آوردن معکوس این تابع بصورت $F_X^{-1}(U) = X$ ، مقادیر توزیع پارتو از

$$x = \sigma \cdot (1 - u)^{-1/\alpha}$$

حاصل می‌شود. اکنون برای $\alpha = 7$ ، $\sigma = 1$ ، مقدار برای x بدست می‌آید، که از توزیع پارتو پیروی می‌کنند. اکنون با استفاده از معادله (۹) با قرار دادن دو مقدار $\tau = 0/05$ و $\tau = 0/95$ و با قرار دادن $s = 1$ و حل آن با نرم‌افزار MATLAB، دو مقدار برای $x = x(\tau)$ بدست می‌آید، که همان کران بالا و پائین فاصله پیش‌بینی برای $x(\tau)$ ، موقعی که اندازه نمونه ثابت می‌باشد، هستند. بهمین طریق با استفاده از داده‌های تولید شده و با فرض اینکه اندازه نمونه متغیر تصادفی با تابع جرم احتمال پواسن یا دوجمله‌ای می‌باشد، می‌توان با استفاده از معادلات (۱۳) و (۱۵) فاصله پیش‌بینی برای اولین مشاهده بعدی $x(\tau)$ را وقتی اندازه نمونه تصادفی است، بدست آورد. نتایج بدست آمده در جدول ۳ خلاصه شده‌اند، که نتایج زیر حاصل می‌شوند.

- ۱- با زیاد شدن مقدار λ ، فواصل پیش‌بینی کوتاهتر می‌شوند.
- ۲- با زیاد شدن مقدار M ، فواصل پیش‌بینی کوتاهتر می‌شوند.
- ۳- برای $\lambda = 1$ ، طول فواصل پیش‌بینی با اندازه نمونه ثابت، کوتاهتر و بهتر از از اندازه نمونه تصادفی می‌باشد. اما برای $\lambda = 20$ و $\lambda = 50$ طول فواصل با اندازه نمونه تصادفی کوتاهتر می‌باشد.
- ۴- با زیاد شدن مقدار a ، فواصل کوتاهتر می‌شوند.
- ۵- با زیاد شدن مقدار b ، فواصل بزرگتر می‌شوند.

۲.۳ پیش‌بینی در توزیع پارتو با پارامترهای α و σ نامعلوم و پیشین بی اطلاع

در این بخش همانند بخش قبل ابتدا از توزیع پارتو ۶ عدد تولید می‌کنیم و بعد با فرض آنکه α و σ هر دو نامعلومند، برای اولین مشاهده بعدی، فواصل پیش‌بینی را بدست می‌آوریم. همچنین فرض می‌کنیم که هیچ اطلاع قبلی از پارامترهای توزیع نداریم و از پیشین بدون اطلاع استفاده می‌کنیم. میانگین فواصل پیش‌بینی شبیه‌سازی شده برای مقادیر مختلف λ را با استفاده از معادلات (۱۸)، (۲۱) و (۲۸) بدست می‌آوریم. نتایج این بررسی در جدول ۴ آمده و نتایج زیر حاصل می‌شوند.

- ۱- با زیاد شدن مقدار λ ، فواصل کوتاهتر می‌شوند.
- ۲- برای λ های بزرگ، فواصل پیش‌بینی با اندازه نمونه تصادفی بهتر از اندازه نمونه ثابت است.

جدول ۳: میانگین طول فواصل پیش‌بینی برای $x = x(\gamma)$ در توزیع پارتو با پارامتر نامعلوم α و با اندازه نمونه ثابت و تصادفی

| پارامترهای توزیع پیشین | اندازه نمونه ثابت | نتایج برای اندازه نمونه تصادفی با توزیع دوجمله‌ای و M, p مقادیر | | | نتایج برای اندازه نمونه تصادفی با توزیع پواسن و λ مقادیر | | |
|------------------------------|-------------------------|---|---------------|---------------|--|----------------------|-----------------------|
| | | $\lambda = 1$ | $\lambda = 2$ | $\lambda = 5$ | $M = 2$ $p = 0/8$ | $M = 5$ $p = 0/5$ | $M = 10$ $p = 0/5$ |
| a | b | $n = 10$ | | | | | |
| ۱ | ۱ | ۰,۸۵۰ | ۶,۳۱۴ | ۰,۴۳۹ | ۰,۱۸۹ | ۰,۷۳۷ | ۰,۲۷۱ |
| ۱ | ۲ | ۱,۱۶۵ | ۱۴,۷۳۶ | ۰,۴۹۲ | ۰,۲۸۵ | ۰,۷۷۰ | ۰,۵۳۳ |
| ۱ | ۵ | ۲,۴۸۸ | ۷۶,۴۵۳ | ۰,۸۲۷ | ۰,۳۵۵ | ۰,۷۵۸ | ۰,۴۸۹ |
| ۵ | ۱ | ۰,۴۳۷ | ۳,۳۲۱ | ۰,۳۰۷ | ۰,۳۰۳ | ۰,۶۹۴ | ۰,۲۶۷ |
| ۵ | ۲ | ۰,۶۰۷ | ۵,۳۱۷ | ۰,۳۴۰ | ۰,۳۱۵ | ۰,۶۳۶ | ۰,۳۵۴ |
| ۵ | ۵ | ۰,۹۹۶ | ۱۳,۲۱۰ | ۰,۵۰۹ | ۰,۳۳۰ | ۰,۵۱۷ | ۰,۴۲۶ |

جدول ۴: میانگین طول فواصل پیش‌بینی برای $x = x(\nu)$ در توزیع پارتو با پارامتر نامعلوم α و σ و با اندازه نمونه ثابت و تصادفی موقعی که پیش‌بین بی‌اطلاع داشته باشیم.

| پیش‌بین بی‌اطلاع | اندازه نمونه ثابت | اندازه نمونه تصادفی پواسن با مقادیر λ | | | |
|------------------|-------------------|---|----------------|----------------|----------------|
| | | $\lambda = 1$ | $\lambda = 10$ | $\lambda = 20$ | $\lambda = 50$ |
| | $n = 10$ | | | | |
| NIP | ۰٫۷۷۶ | ۱۰٫۵۵۹ | ۱٫۱۸۳ | ۰٫۲۳۷ | ۰٫۰۷۱ |

۴ بحث و نتیجه‌گیری

در مدل طول عمر پارتو، میانگین طول فواصل پیش‌بینی برای مشاهدات بعدی در یک نمونه که با استفاده از روش شبیه‌سازی بدست می‌آیند، به پارامترهای توزیع پیش‌بین و توزیع اندازه نمونه بستگی دارند. به همین دلیل در حالت کلی نمی‌توان هر یک از دو روش پیش‌بینی بر اساس اندازه نمونه ثابت و تصادفی را بر دیگری ترجیح داد. وقتی اندازه نمونه متغیری با توزیع پواسن $p(\lambda)$ در نظر گرفته می‌شود، برای مقادیر بزرگ λ ، طول فواصل پیش‌بینی حاصل از روش اندازه نمونه تصادفی، کوتاهتر از طول فواصل پیش‌بینی حاصل از اندازه نمونه تصادفی می‌باشد. دلیل این تفاوت می‌تواند ناشی از بزرگتر شدن حجم نمونه ثانوی برای مقادیر بزرگ λ باشد. همچنین بر اساس نتایج بدست آمده، فواصل پیش‌بینی به پارامترهای توزیع پیش‌بین حساس هستند. برای توزیع پارتو با پیش‌بین $Gamma(a, b)$ ، وقتی a زیاد می‌شود، فواصل پیش‌بینی برای اولین مشاهده ثانوی کوتاهتر و وقتی b زیاد می‌شود، فواصل طولانی‌تر می‌شوند.

مراجع

- [1] Arnold, B. C. (1982), Pareto Distribution, International Co-Operative Publishing House.
- [2] Arnold, B. C. and Press, S. J. (1983), Bayesian Inference for Pareto Populations, J. of Econometrics, 21, 287-306.
- [3] Arnold, B. C. and Press, S. J. (1989), Bayesian Estimation and Prediction for Pareto Data, J. Amer. Statist. Assoc., 84, 1079-1084.
- [4] Consul, P. C. (1984), On the Distribution of Order Statistics for Random Sample Size, Statist. Netherlands, 38, 249-256.

- [5] Dunsmore, I. R. (1974), The Bayesian Predictive Distribution in Life Testing Models, *Journal of Technometrics*, 16, 455-460.
- [6] Epstein, B. and Sobel, M. (1953), Life Testing, *J. Amer. Statist. Assoc.*, 48, 486-502.
- [7] Faulkenberry, D. G. (1973), A Method of Obtaining Prediction Intervals, *J. Amer. Statist. Assoc.*, 68, 433-435.
- [8] Gupta, D. and Gupta, R. C. (1984), On the Distribution of Order Statistics for a Random Sample Size, *Statist. Netherlands*, 38, 13-19.
- [9] Hewett, J. E. (1968), A Note on Prediction Intervals Based on Partial Observations in Certain Life Test Experiments, *Technometrics*, 10, 850-853.
- [10] Kulldorff, G. and Vannman, K. (1973), Estimation of the Location and Scale Parameters of a Pareto Distribution by Linear Functions of Order Statistics, *J. Amer. Statist. Assoc.*, 68, 218-227.
- [11] Nigm, A. M. and Hamdy, H. I. (1987), Bayesian Prediction Bounds for the Pareto Lifetime Model, *Commun. Statist., Theory and Methods*, 16, 1761-1772.
- [12] Soliman, A. A. (2000), Bayes Prediction in a Pareto Lifetime Model with Random Sample Size, *The Statistician*, 49, Part 1, 51-62.

شبیه‌سازی احتمالات ورشکستگی در فرآیندهای مخاطره بیمه

نادر مظاهری

کارشناس ارشد آمار شرکت توزیع برق شهرستان اصفهان

چکیده: یک مساله مهم در امور بیمه تعیین احتمالات ورشکستگی است. محاسبه احتمالات ورشکستگی معمولاً از روشهای تحلیلی میسر نبوده و به همین علت شبیه‌سازی این احتمالات دارای اهمیت می‌باشد. ورشکستگی پیشامدی با احتمال کم بوده، شبیه‌سازی معمولی آن به دلیل اینکه نیازمند تولید مسیره‌های نمونه‌ای بسیار زیاد است دارای کارایی لازم نمی‌باشد و شبیه‌سازی باید با روشهای خاصی انجام شود. در این مقاله به یک نوع از این روشها با نام نمونه‌گیری ویژه خواهیم پرداخت. در فرآیندهای بیمه با اندازه‌های خسارت زیر نمائی یک نمونه‌گیری ویژه تغییر اندازه احتمال بصورت چرخش زیر نمائی می‌باشد. در این مقاله شرطهائی برای بهینه بودن مجانبی برآوردکننده‌های ناشی از این نوع شبیه‌سازی برای احتمالات ورشکستگی با افق زمانی متناهی و نامتناهی را شرح می‌دهیم و با معرفی دو مدل پواسون مرکب و مارکف - مدول پواسون مرکب به شبیه‌سازی احتمال ورشکستگی در هر یک از این دو مدل می‌پردازیم.

واژه‌های کلیدی: احتمال ورشکستگی، پیشامد کم محتمل، چرخش نرخ خطر، چرخش نرخ خطر تأخیری وزنی، فرایند مخاطره بیمه، کارائی شبیه‌سازی، نمونه‌گیری ویژه

۱ مقدمه

فرض کنیم X_i اندازه خسارت n ام در زمان ورود η_n ($n \geq 1$) با $\eta_0 < \eta_1 < \dots$ باشد. زمان بین ورود n امین و $n-1$ امین ادعا را با $\xi_n = \eta_n - \eta_{n-1}$ نشان می‌دهیم. اگر $N(t)$ تعداد تصادفی از ادعاها در فاصله زمانی $[0, T]$ باشد فرایند مخاطره بیمه را به صورت زیر در نظر می‌گیریم:

$$U(t) = u + \Pi(t) - \sum_{i=1}^{N(t)} X_i$$

که در آن u , $\Pi(t)$ به ترتیب نشان دهنده ذخیره اولیه و تابع حق بیمه می‌باشند. در این مقاله تابع حق بیمه را به صورت ct با نرخ c در نظر می‌گیریم. برای اطلاعات بیشتر از خواص تابع حق بیمه می‌توان به اسموزن^۱ (۲۰۰۰) مراجعه کرد. دیده می‌شود که $U(t)$ ذخیره تا زمان t است.

1) asmussen

اگر فرایند مازاد خسارت را به صورت $u - U(t)$ تعریف کنیم و آن را با $S(t)$ نشان دهیم در این صورت احتمال ورشکستگی با افق زمانی متناهی T به صورت زیر داده می‌شود:

$$\psi(u, T) = P(U(t) < u, t < T) = \left(\sup_{t < T} S(t) > u \right) \quad (۱)$$

همچنین احتمال ورشکستگی با افق نامتناهی را به صورت $\psi(u) = \psi(u, \infty)$ تعریف می‌کنیم. فرض کنیم با احتمال یک

$$\lim_{t \rightarrow \infty} \frac{\sum_{i=1}^{N(t)} X_i}{t} = w\beta \quad , \quad \lim_{t \rightarrow \infty} \frac{N(t)}{t} = w$$

وقتی فرایند مازاد خسارت $S(t)$ با رانش منفی است (یعنی $S(t) \rightarrow -\infty$ وقتی $t \rightarrow \infty$) احتمال ورشکستگی با افق نامتناهی کمتر از ۱ خواهد بود که در این صورت $c > w\beta$ یا $\rho = \frac{w\beta}{c} < ۱$. اگر $S(t)$ با رانش غیر منفی باشد با توجه به تعریف (۱) ورشکستگی با افق متناهی با احتمال ۱ رخ می‌دهد.

در ساده‌ترین حالات نیز احتمالات ورشکستگی به صورت صریح امکان‌پذیر نمی‌باشد و اغلب در تقریب زدن آن به روشهای گوناگون سعی شده است. بخشی از این تقریب‌ها شامل برآوردهای ناشی از شبیه‌سازی می‌باشد. باید توجه داشت که وقتی $\psi(u) \rightarrow ۰$, $u \rightarrow \infty$ و $\psi(u, T) \rightarrow ۰$ یعنی با زیاد شدن u احتمال ورشکستگی کم می‌شود و بنابراین ورشکستگی یک پیشامد کم محتمل خواهد بود. از این رو شبیه‌سازی معمولی^۲ که برآورد $\psi(u, T)$ و $\psi(u)$ با میانگین نمونه‌ای از تکرارهای مستقل $I(\sup_{0 \leq t < T} S(t) > u)$ و $I(\sup_{t \geq ۰} S(t) > u)$ ($I(\cdot)$ متغیر تصادفی نشانگر می‌باشد) است دارای کارایی لازم نبوده و روشهای شبیه‌سازی خاصی مورد نیاز است. در شبیه‌سازی معمولی باید مسیرهای نمونه‌ای بسیاری تولید کرد تا به مسیر نمونه‌ای که روی آن ورشکستگی رخ می‌دهد دسترسی پیدا کنیم. تولید این مسیرهای نمونه‌ای ضعف این نوع شبیه‌سازی را نمایان می‌سازد.

نمونه‌گیری ویژه^۳ یک روش خاص شبیه‌سازی می‌باشد که تحت آن سرعت وقوع پیشامدهای کم محتمل بالاتر می‌رود [هیدل برگر^۴ (۱۹۹۵) و اسموزن (۲۰۰۰)]. در این روش پس از تغییر اندازه احتمال از نسبت درستنمایی برای تعیین برآورد ناریب استفاده می‌شود. تغییر اندازه احتمال باید به گونه‌ای باشد که واریانس برآورد کننده در مقایسه با شبیه‌سازی معمولی بزرگ نشود [لهتن^۵ و نیرهین^۶ (۱۹۹۲)]. در حالتی که توزیع اندازه خسارت با دم سبک باشد تغییر اندازه احتمال به صورت چرخش نمایی^۷ مفید می‌باشد [اسموزن (۱۹۸۵)]. بدین معنی که این رهیافت به طور

2) naive simulation 3) importance sampling 4) Heidelberger 5) Lehtonen
6) Nyrhinen 7) exponential twisting

مجانبی بهینه می‌باشد. به طور مجانبی بهینه بودن ملاکی استاندارد برای میزان کارایی شبیه‌سازی پیشامدهای کم محتمل می‌باشد.

هنگامی که اندازه خسارت با دم سنگین باشد به کار بردن چرخش نمایی برای حل مسأله مفید واقع نمی‌شود. توزیع با دم سنگین توزیعی است که تابع مولد گشتاور متناهی نداشته و به همین علت از چرخش نمایی که مبتنی بر تابع مولد گشتاور است نمی‌توان استفاده کرد. اسموزن و بینسوانگر^۸ (۱۹۹۷) اولین رهیافت برای شبیه‌سازی در حضور توزیع‌های زیرنمایی^۹ (توزیع‌های زیرنمایی کلاسی از توزیع‌های با دم سنگین هستند) را ارائه دادند. آنها در صدد شبیه‌سازی احتمال تجاوز یک مجموع هندسی از متغیرهای تصادفی $i.i.d.$ با توزیع زیرنمایی از یک مقدار مشخص بودند. جونجا^{۱۰} و شهاب‌الدین^{۱۱} (۱۹۹۹) با معرفی چرخش نرخ خطر^{۱۲} یا چرخش زیرنمایی به حل مسأله فوق پرداختند.

آنها همچنین با بهبود بخشیدن این روش به وسیله چرخش نرخ خطر تأخیری^{۱۳} نشان دادند که شبیه‌سازی مبتنی بر آن برای چند توزیع زیرنمایی به طور مجانبی بهینه می‌باشد. این تکنیکها روشهایی برای برآورد آن است که یک مجموع هندسی از متغیرهای تصادفی $i.i.d.$ از یک مقدار داده شده تجاوز کند. حل چنین مسأله‌ای منجر به برآورد احتمال ورشکستگی نیز می‌شود. زیرا از تبدیل Pollaczek-Khinchin احتمال ورشکستگی در فرایندهای مخاطره بیمه با ادعاهای پواسون و اندازه‌های خسارت $i.i.d.$ به مسأله مجموع هندسی ذکر شده تبدیل می‌شود. به هر حال در اکثر موارد حتی به کمک چنین تبدیلی نیز برآورد احتمالات ورشکستگی بسیار مشکل است.

بوتز^{۱۴} و شهاب‌الدین (۲۰۰۰) برآورد $\psi(u)$ را با زمانهای بین ورود $i.i.d.$ و بدون استفاده از تبدیل Pollaczek-Khinchin در نظر گرفتند. آنها با شرط اینکه اندازه‌های خسارت $i.i.d.$ بوده و از زمانهای بین ورود مستقل باشند مستقیماً به شبیه‌سازی فرایند مازاد خسارت پرداختند. بدین صورت که قدم تصادفی $S_i = S_{i-1} + (X_i - c\xi_i)$ به دست آمده از فرایند مازاد خسارت را شبیه‌سازی کردند.

آنها تغییر اندازه چرخش زیرنمایی برای اندازه‌های خسارت X_i را به کار بردند. سپس نشان دادند که برآوردکننده حاصل برای توزیع‌های اندازه خسارت وایبول به طور مجانبی در یک مجموعه بزرگ بهینه^{۱۵} می‌باشد. به طور مجانبی در یک مجموعه بزرگ بهینه بودن ملاک دیگری برای کارایی شبیه‌سازی پیشامدهای کم محتمل می‌باشد.

در این مقاله به تعمیم کارهای انجام شده توسط بوتز و شهاب‌الدین (۲۰۰۰) پرداخته می‌شود و خواهیم دید که با ارائه شرطهای کافی روی فرایند خسارت و توزیع‌های اندازه خسارت تغییر

8) Binswanger 9) subexponential 10) Juneja 11) Shahabuddin 12) hazard rate twisting hazard 13) delayed hazard rate twisting 14) Boots 15) larg set asymptotically optimal

اندازه‌های استفاده شده به طور مجانبی در یک مجموعه بزرگ بهینه می‌باشند. همچنین حالتی را در نظر می‌گیریم که اندازه‌های خسارت وابسته بوده و به شبیه‌سازی احتمالات ورشکستگی خواهیم پرداخت.

۲ نمادگذاری و مدلها

اگر دم توزیع F را با $\bar{F} = 1 - F$ چگالی آن را با f و تابع نرخ خطرش را با $\lambda(x) = \frac{f(x)}{\bar{F}(x)}$ نشان دهیم داریم: $\bar{F}(x) = e^{-\Lambda(x)}$ به طوریکه $\Lambda(x) = \int_0^x \lambda(s)ds$ تابع خطر است. توزیع دم انتگرالگیری شده F را به صورت $F_I(x) = \frac{\int_0^x \bar{f}(y)dy}{\int_0^\infty \bar{f}(y)dy}$ به شرط $(\int_0^\infty \bar{F}(y)dy < \infty)$ تعریف می‌کنیم و $\Lambda(x), \lambda_I(x)$ را به ترتیب تابع نرخ خطر و تابع خطر متناظر با F_I در نظر می‌گیریم. برای توابع دلخواه $z_1(x), z_2(x)$ و $z_1(x) \sim z_2(x)$ بدین معنی است که وقتی $x \rightarrow \infty$ به $\frac{z_1(x)}{z_2(x)}$ همگراست.

زمان ورشکستگی را به صورت $\sigma(u) = \inf\{t : S(t) > u\}$ تعریف می‌کنیم و قرار می‌دهیم $\tau(u) = N(\sigma(u))$. اگر ورشکستگی رخ دهد $\tau(u)$ تعداد ادعاها تا قبل از وقوع ورشکستگی می‌باشد. برای $T < \infty$ تعریف می‌کنیم $\tau(u, T) = N(\text{Min}(\delta(u), T))$. در این صورت داریم:

$$\psi(u, T) = P(\sigma(u) < T) \quad , \quad \psi(u) = P(\sigma(u) < \infty) = P(\tau(u) < \infty)$$

۱.۲ مدل مخاطره بیمه

مدل مخاطره بیمه را مانند بخش ۱ در نظر می‌گیریم. اندازه‌های خسارت X_i می‌توانند به طور تصادفی یکی از m نوع مختلف از توزیعهای F_1, F_2, \dots, F_m با میانگینهای متناهی را داشته باشند. فرض کنیم X_i ها مستقل بوده و حداقل یکی از F_i ها ($1 \leq i \leq m$) دارای توزیع زیرنمایی باشد. توزیعهای زیرنمایی رده‌ای خاص از توزیعهای با دم سنگین هستند. در مقابل توزیعهای با دم سنگین توزیعهای با دم سبک توزیعهایی هستند که دم آنها با نرخ نمایی (سریعتر از توزیعهای با دم سنگین) کاهش می‌یابد.

تعریف: اگر X_n دنباله‌ای از متغیرهای تصادفی *i.i.d.* و غیر منفی با توزیع F و n امین پیچش F^{*n} باشد F را زیرنمایی گوییم هرگاه برای هر $n \geq 2$ داشته باشیم:

$$\frac{\bar{F}^{*n}(u)}{n\bar{F}(u)} \sim \frac{P(X_1 + \dots + X_n > u)}{nP(X_1 > u)} \rightarrow 1 \quad , \quad (u \rightarrow \infty)$$

توزیع‌هایی مانند لگ نرمال، وایبول و پارتو متعلق به این خانواده هستند. تابع نرخ خطر بسیاری از توزیع‌های زیرنمایی نهایتاً نزولی می‌باشد. درامبرج^{۱۷} و کلاپلبرگ^{۱۸} (۱۹۹۷) جزئیات در مورد توزیع‌های زیرنمایی آمده است. حال در مدل‌های مخاطره بیمه دو مدل زیر را معرفی می‌کنیم.

۱.۱.۲ مدل پواسون مرکب

در این مدل زمانهای ورود ادعاها یک فرایند تجدید با زمانهای بین ورود نمایی هستند و بنابراین تعداد ادعاها تشکیل یک فرایند پواسون را می‌دهد. همچنین در این مدل X_i ها $i.i.d.$ بوده که از ξ_i ها نیز مستقل می‌باشند.

۲.۱.۲ مدل مارکوف - مدل پواسون مرکب

در این مدل زمانهای بین ورود ادعاها $i.i.d.$ مستقل از X_i ها با فرایند پواسون بوده شدت ورود و توزیع اندازه خسارت بر اساس فرایند مارکوف $J(t), t \geq 0$ با فضای حالت متناهی $\{1, 2, \dots, m\}$ همراه توزیع مانای یکتای $\pi_1, \pi_2, \dots, \pi_m$ ، $\pi_i > 0$ برای هر $1 \leq i \leq m$ مشخص می‌شود. وقتی $J(t)$ در حالت i است شدت ورود را w_i نرخ حق بیمه c را برابر ۱ و توزیع اندازه خسارت را F_i قرار می‌دهیم. در این صورت داریم $F_i = P(X_i \leq x | J(\eta_j) = i)$. در این مدل داریم:

$$\beta = \frac{\sum_{i=1}^m w_i \pi_i \beta_i}{\sum_{i=1}^m w_i \pi_i}, \quad w = \sum_{i=1}^m \pi_i w_i$$

که در آن β_i میانگین F_i می‌باشد. بنابراین داریم: $\rho = \sum_{i=1}^m \pi_i w_i \beta_i$. اگر تعریف کنیم

$$P_i(\cdot) = P(\cdot | J(\cdot) = i)$$

احتمال ورشکستگی قبل از زمان T با حالت اولیه i را به صورت $\psi_i(u, T)$ نشان می‌دهیم. در ادامه برای سادگی اندیس i را حذف می‌کنیم.

۳ شبیه‌سازی پیشامدهای کم محتمل

همانطور که در بخش ۱ گفته شد برای u بزرگ $\psi(u, T)$ کوچک بوده و شبیه‌سازی معمولی برای برآورد آن کاربرد چندانی ندارد. در اینجا نوع دیگری از شبیه‌سازی با نام نمونه‌گیری ویژه را معرفی می‌کنیم. فرض کنیم فرایند تصادفی که درصدد شبیه‌سازی آن هستیم در یک فضای احتمال با

17) Embrech 18) Kluppelbrg

اندازه احتمال P تعریف شده باشد. $A(u)$ را پیشامدی با احتمال $a(u) = P(A(u))$ طوری در نظر می‌گیریم که وقتی $u \rightarrow \infty$, $a(u) \rightarrow 0$. فرض کنیم Q اندازه احتمال دیگری روی همان فضای احتمال بوده به طوری که P نسبت به آن مطلقاً پیوسته باشد. در این صورت داریم:

$$a(u) = E_Q(I(A(u)dP/dQ)) \quad (۲)$$

در (۲) dP/dQ نسبت درستی است و اندیس Q نشان دهنده آن است که امید تحت اندازه Q گرفته می‌شود. نسبت درستی dP/dQ به صورت زیر تعریف می‌شود:

$$\frac{dP}{dQ} = L(X_1, X_2, \dots, X_n) = \prod_{i=1}^n \frac{f(X_i)}{g(X_i)}$$

که در آن f و g توابع چگالی X_i به ترتیب تحت اندازه احتمال P, Q می‌باشند. در نمونه‌گیری ویژه مسیرهای نمونه‌ای تحت اندازه احتمال Q را تولید کرده در هر حالت نسبت درستی را محاسبه می‌کنیم. سپس از میانگین نمونه‌ای $I\left(A(u)\left(\frac{dP}{dQ}\right)\right)$ ها برآورد نارایب برای $a(u)$ تشکیل می‌دهیم. مسأله مهم در نمونه‌گیری ویژه تعیین اندازه احتمال جدید Q می‌باشد به طوری که واریانس این برآورد نارایب در مقایسه با شبیه‌سازی معمولی مقدار بزرگی نباشد. یک ملاک استاندارد برای مشخص کردن کارایی شبیه‌سازی پیشامدهای کم محتمل به طور مجانبی بهینه بودن می‌باشد. **تعریف:** $\hat{a}(u)$ برآوردکننده به طور مجانبی بهینه برای $a(u)$ می‌باشد هرگاه:

$$\lim_{u \rightarrow \infty} \text{SUP} \frac{\log(\text{Var}[\hat{a}(u)])}{\log(a^2(u))} \geq 1$$

در بسیاری از موارد رسیدن به چنین ملاکی مشکل می‌باشد از این رو با صرف نظر از آریبی ملاک ضعیف‌تر زیر را به کار می‌بریم:

تعریف: بهینگی مجانبی در مجموعه بزرگ نرمال شده \mathcal{A} (w.a.o) فرض کنیم $\delta \in (0, 1)$ ثابت و مشخص باشد. تجزیه پارامتری $a(u)$ با پارامتر β را به صورت زیر در نظر می‌گیریم:

$$a(u) = \gamma_\beta(u) + \epsilon_\beta(u)$$

اگر داشته باشیم:

$$(۱) \quad \lim_{u \rightarrow \infty} \sup \frac{\epsilon_\beta(u)}{a(u)} \leq a \quad \text{شده } \beta \text{ داده}$$

$$(۲) \quad \lim_{\beta \rightarrow \infty} \sup \frac{\epsilon_\beta(u)}{a(u)} = 0 \quad \text{برای هر } u \text{ مشخص}$$

(۳) برای هر β داده شده یک برآورد ناریب برای $\gamma_\beta(u)$ وجود داشته به طوریکه:

$$\lim_{u \rightarrow \infty} \inf \frac{\log(\text{work}(u) \times \text{Var}[\hat{\gamma}_\beta(u)])}{\log(\gamma_\beta^2(u))} \geq 1$$

در این صورت $\hat{\gamma}_\beta(u)$ را برآورد به طور مجانبی بهینه در یک مجموعه بزرگ نرمال شده برای $a(u)$ می‌نامیم. $\text{work}(u)$ تابعی از u بوده که محاسبات در هر بار تکرار شبیه‌سازی را نشان می‌دهد. بوتز و شهاب‌الدین (۲۰۰۰). در تعریف فوق پارامتر β برای کنترل آریبی نسبتی برآوردکننده به کار می‌رود به طوریکه برای هر مقدار مشخص u آریبی نسبتی کمتر از δ باشد. برای سادگی از $\gamma(u)$ به جای $\gamma_\beta(u)$ و از $\epsilon(u)$ به جای $\epsilon_\beta(u)$ استفاده می‌کنیم.

اندازه احتمال جدید Q که از آن استفاده می‌کنیم چرخش نرخ خطر (HRT) روی هر یک از توزیعهای اندازه خسارت می‌باشد. چرخش نرخ خطر به صورت جایگزین کردن توزیع F_i با توزیع جدید $F_{j,\theta_u}(x) = 1 - e^{-\Lambda_j(x)(1-\theta_u)}$ ، $(0 \leq \theta_u < 1)$ می‌باشد و چگالی مطابق با F_{j,θ_u} به صورت $f_{j,\theta_u}(x) = (1 - \theta_u)\lambda_j(x)e^{-(1-\theta_u)\Lambda_j(x)}$ است. θ_u تابعی مناسب از u است که با ساختن کران بالا برای نسبت درستنمایی و می‌نیم کردن آن نسبت به θ_u به دست می‌آید. چرخش نرخ خطر تأخیری وزنی ${}^{\circ}\text{WDHRT}$ (تعمیمی برای HRT) می‌باشد به طوریکه برآوردکننده ناشی از آن بهینگی بیشتری نسبت به HRT دارد. بوتز و شهاب‌الدین (۱۹۹۹). چگالی چرخش نرخ خطر تأخیری وزنی را با معرفی پارامتر وزنی w_u و پارامتر تأخیری x_u^* (هر دو تابعی از u هستند) به صورت زیر تعریف می‌کنیم:

$$f_{j,\theta_u,x_u^*}(x) = \begin{cases} \frac{f_j(x)}{1+w_u} & x \leq x_u^* \\ \left(1 - \frac{F_j(x_u^*)}{1+w_u}\right) \frac{f_{j,\theta_u}(x)}{F_{j,\theta_u}(x_u^*)} & x > x_u^* \end{cases}$$

تابع توزیع آنرا نیز با F_{j,θ_u,x_u^*} نشان می‌دهیم.

انتخاب هر یک از این پارامترها باید به گونه‌ای باشد که در نمونه‌گیری ویژه خواص اندازه احتمال جدید بسیار نزدیک به خواص اندازه احتمال اصلی باشد.

۴ احتمالات ورشکستگی با افق نامتناهی

همانطور که گفته شد با استفاده از چرخش نرخ خطر برای برآورد احتمال ورشکستگی از:

$$E_Q(I(\tau(u) < \infty)L_Q(u))$$

استفاده می‌کنیم. که در آن L تابع نسبت درست‌نمایی می‌باشد. در اینجا با یک مسأله روبرو هستیم. تغییرات $L_Q(u)$ روی مجموعه مسیرهای نمونه‌ای زیاد است. اگر بتوانیم تابعی از u مانند $k_*(u)$ طوری پیدا کنیم که:

$$E(I(\tau(u) < \infty)) \sim E(I(\tau(u) \leq K_*(u))) = E_Q(I(\tau(u) \leq K_*(u)L_Q(u)))$$

در این حالت اگر L دارای واریانس کمی روی $(\tau(u) \leq k_*(u))$ باشد می‌توانیم بر این مسأله فائق آییم.

حال این سؤال پیش می‌آید که چه میزان اریبی را می‌توان تحمل کرد؟ اگر قرار دهیم:

$$\gamma(u) = E(I(\tau(u) \leq k_*(u)))$$

در این صورت اگر $\frac{E(I(\tau(u) \leq k_*(u)))}{E(I(\tau(u) < \infty))} < 1 - \delta$ می‌توان از اریبی به وجود آمده صرف نظر کرد [بوتز و شهاب‌الدین (۲۰۰۱)]. مسأله دیگر تعیین $k_*(u)$ می‌باشد. برای تعیین این تابع شرطهای زیر را در نظر می‌گیریم. خواهیم دید که این شرطها برای برقراری ملاک w.a.o مورد نیاز می‌باشد.

۱.۴ شرطهای انحرافات بزرگ

$$(۱) \quad \lim_{u \rightarrow \infty} \frac{-\log(\psi(u))}{\Lambda_I(u)} = 1 \quad \text{به گونه‌ای باشد که:}$$

$$(۲) \quad \text{برای هر } \delta > 0 \text{ داده شده تابع } k_*(u) \text{ به گونه‌ای باشد که:}$$

$$\lim_{u \rightarrow \infty} \frac{P(\tau(u) \leq k_*(u))}{P(\tau(u) < \infty)} \geq 1 - \delta$$

$$(۳) \quad \log(k_*(u)) = O(\Lambda_I(u))$$

$$(۴) \quad \text{ثابت } b > 1 \text{ به گونه‌ای باشد که: } \lim_{u \rightarrow \infty} \frac{k_*(u)}{(\Lambda(u))^{b-1}} = 0 \quad \text{و}$$

$$\lim_{u \rightarrow \infty} \frac{k_*(u)F^{+-}(\Lambda(u)^{-b})}{u} = 0$$

$$F^{+-} = \inf(x : F(x) = y)$$

بدین با شرطهای فوق محدودیتهایی را برای $k_*(u)$ قائل می‌شویم.

رده زیادی از فرایندهای مخاطره بیمه در شرطهای انحرافات بزرگ ذکر شده صدق می‌کنند. به عنوان مثال اگر F زیرنمایی باشد در مدل پواسون مرکب و مارکوف - مدل پواسون مرکب شرط (۱) برقرار است. همچنین توزیع وایبول با shape پارامتر کمتر از (۱) در شرطهای (۲) الی (۴) صدق می‌کنند.

۲.۴ الگوریتم شبیه‌سازی و کارایی آن

برای شبیه‌سازی از چرخش نرخ خطر تأخیری وزنی روی هر یک از توزیعهای اندازه خسارت استفاده می‌کنیم. پارامترها را به صورت زیر در نظر می‌گیریم:

$$\theta_u = 1 - \frac{1}{\Lambda(u)}, \quad w_u = \frac{c_1 \log \delta}{k_*(u)}, \quad x_u^* = F^{+-}(1 - \Lambda(u)^{-b}) \quad (*)$$

به طوری‌که $c_1 > 1$ و b ثابت ذکر شده در شرط (۴) می‌باشد. باید توجه داشت که از $\Lambda(x_u^*) = b \log(\Lambda(u))$ نتیجه می‌شود که وقتی $u \rightarrow \infty$ ، $x_u^* \rightarrow \infty$ زیرا $\Lambda(u) \rightarrow \infty$.
 الگوریتم ۱ (WDHRT) برای برآورد $\psi(u)$: به جای $F_i = 1 - e^{-\Lambda_i(x)}$ از WDHRT روی F_i با پارامترهای θ_u ، w_u ، x_u^* انتخاب شده به صورت (*) استفاده می‌کنیم. L را نسبت درستی‌مایی قرار می‌دهیم. شبیه‌سازی را با تکرار از

$$S = \sum_{i=1}^{\min(\tau(u), k_*(u))} (X_i - c\xi_i)$$

انجام می‌دهیم. اگر $S > u$ باشد قرار می‌دهیم $V = L$ و در غیر این صورت قرار می‌دهیم $V = 0$. میانگین K تکرار $i.i.d$ از V را به عنوان یک برآورد ناریب از $P(\tau(u) \geq k_*(u))$ در نظر می‌گیریم. از این برآورد کننده برای برآورد $\psi(u)$ استفاده می‌کنیم.
 قضیه: تحت شرطهای انحرافات بزرگ اگر توزیع F دارای تابع نرخ خطر نهایتاً نزولی با خاصیت نرمال شده برای $\psi(u)$ فراهم می‌کند (برای اثبات می‌توان به بوتز و شهاب‌الدین (۲۰۰۱) مراجعه کرد).

۵ احتمالات ورشکستگی با افق متناهی

a را کمینه روی هر a می‌گیریم به طوری‌که $P(N(T) \geq k) \leq a^k$ برای k به قدر کافی بزرگ برقرار باشد. شرط زیر را در نظر می‌گیریم:

$$\lim_{u \rightarrow \infty} -\frac{\log(\psi(u, T))}{\Lambda(u)} = 1$$

قضیه: با فرض نهایتاً نزولی بودن $\lambda(x)$ در مدل‌های تجدید و مارکوف - مدل یواسون مرکب $\psi(u, T)$ در شرط فوق صدق می‌کند.

۱.۵ الگوریتم شبیه‌سازی و کارایی آن

θ_u, x_u^* را مانند قبل انتخاب می‌کنیم. اما w_u را طوری انتخاب می‌کنیم که $w_u < \frac{1}{a} - 1$. اگر a نامعلوم باشد قرار می‌دهیم: $w_u = 0$.
 الگوریتم ۲ (WDHRT) برای برآورد $\psi(u, T)$: به جای $F_i = 1 - e^{-\Lambda_i(x)}$ از چرخش نرخ خطر تأخیری وزنی استفاده می‌کنیم. تکرارها را از

$$S = \sum_{i=1}^{\tau(u, T)} (X_i - \text{cmin}(\delta(u), T))$$

تولید می‌کنیم. اگر $S > u$ باشد قرار می‌دهیم $V = L$ و اگر $S \leq u$ باشد قرار می‌دهیم $V = 0$. میانگین K تکرار $i.i.d$ از V را به عنوان برآورد ناریب $\psi(u, T)$ در نظر می‌گیریم. قضیه: با فرض نهایتاً نزولی بودن تابع نرخ خطر خسارت الگوریتم ۲ با پارامترهای θ_u, x_u^*, w_u به طور مجانبی بهینه می‌باشد.

۶ ورشکستگی با خسارتهای وابسته

در بسیاری از مواقع فرض استقلال اندازه‌های خسارتهای غیر واقعی جلوه کرده به طوری که صرفنظر از وابستگی خسارتهای می‌تواند برای برآورد احتمال ورشکستگی گمراه کننده باشد. دیده شده است که وقتی خسارتهای وابسته هستند نسبت به حالتی که مستقل اند زمان ورشکستگی زودتر رخ می‌دهد [امیلیانو^{۲۱} و کلون^{۲۲} (۲۰۰۲)].
 تعریف: تابع زیر را از $[0, 1]^n$ به $[0, 1]$ برای $n \geq 1$ به صورت زیر در نظر می‌گیریم:

$$F_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = C(F_{X_1}, F_{X_2}, \dots, F_{X_n}) \quad (3)$$

به طوری که F_{X_1, X_2, \dots, X_n} تابع توزیع تجمعی بردار تصادفی از متغیرهای (X_1, \dots, X_n) و F_{X_i} ها $(1 \leq i \leq n)$ نشان دهنده توزیعهای حاشیه‌ای هستند. تابع C یک کاپولا^{۲۳} نامیده می‌شود. C یک تابع توزیع تجمعی توأم با توزیعهای حاشیه‌ای یک متغیره $U(0, 1)$ می‌باشد. کاپولاها در واقع توابعی هستند که ارتباط بین توزیع توأم و توزیعهای حاشیه‌ای را با پارامترهای توصیف کننده وابستگی متغیرها نشان می‌دهد. در رابطه (۳) تابع کاپولا منحصر بفرد نبوده ولی با فرض مطلقاً پیوسته بودن توزیعهای حاشیه‌ای یک متغیره تابع C تنها تابع با تعریف (۳) می‌باشد.
 مثال: بردار (u_1, u_2, \dots, u_n) را با توزیعهای حاشیه‌ای یک متغیره $U(0, 1)$ در نظر می‌گیریم.

تابع

$$C(u_1, u_2, \dots, u_n) = -\frac{1}{\log \eta} \log \left(1 + \frac{\prod_{k=1}^n (\eta^{u_k} - 1)}{(\eta - 1)^{n-1}} \right) \quad (4)$$

یک کاپولا بوده که موسوم به فرانک کاپولا^{۲۴} می‌باشد. تابع کاپولای ارائه شده در (۴) را می‌توان به صورت:

$$C(u_1, \dots, u_n) = \phi^{-1}(\phi(u_1) + \phi(u_2), \dots, \phi(u_n)) \quad (5)$$

نوشت. در (۵)، ϕ به صورت:

$$\phi(t) = \log \left(\frac{\eta^t - 1}{\eta - 1} \right) \quad (6)$$

می‌باشد. در (۴) و (۶)، η به گونه‌ای است که $1 - \eta$ پارامتر یک توزیع لگاریتمی گسسته می‌باشد [امیلیانو و کلومین (۲۰۰۲)].

۱.۶ شبیه‌سازی احتمالات ورشکستگی با خسارت‌های وابسته

در بخش‌های (۴) و (۵) الگوریتم‌های شبیه‌سازی مطرح شد. در آنجا دیدیم که مسیرهای نمونه‌ای از فرایند مورد نظر را تولید کرده و از بین آنها خود را روی مسیرهایی که ورشکستگی روی آنها رخ می‌دهد متمرکز کردیم. حال به تولید چنین مسیرهایی با فرض وابسته بودن خسارت‌ها می‌پردازیم. با استفاده از تابع کاپولای (۴) و مولد (۶) مسیرهای نمونه‌ای را به صورت زیر تولید می‌کنیم:

(۱) متغیر تصادفی Z را با توزیع لگاریتمی گسسته و با پارامتر $1 - \eta$ تولید می‌کنیم.

(۲) n متغیر تصادفی $i.i.d$ با توزیع $U(0, 1)$ تولید می‌کنیم. این متغیرهای تولید شده را به صورت $U^* = (u_1^*, u_2^*, \dots, u_n^*)$ نشان می‌دهیم.

(۳) قرار می‌دهیم $U = M_Z(z^{-1} \log U^*)$ به طوری که $M_Z(\cdot)$ تابع مولدگشتاور Z به صورت $M_Z(t) = \frac{\log(1 - (1 - \eta)e^t)}{\log(\eta)}$ می‌باشد و $\log U^* = (\log u_1^*, \log u_2^*, \dots, \log u_n^*)$

۷ نتایج عددی

در این بخش با استفاده از الگوریتم‌های ارائه شده به برآورد احتمالات ورشکستگی با افق زمانی متناهی و نامتناهی در دو مدل پواسون مرکب و مارکوف - مدل پواسون مرکب می‌پردازیم. نتایج عددی ارائه شده از برنامه نوشته شده به زبان پاسکال برای شبیه‌سازی به دست آمده است. وقتی F دارای توزیع وایبول است برای انتخاب w_u ، c_1 را به طور تجربی برای احتمال ورشکستگی

24) Frank's copula

با افق نامتناهی انتخاب می‌کنیم. این انتخاب را با \tilde{w}_u نشان می‌دهیم. وقتی F لگ نرمال است قرار می‌دهیم $w = 0$. از قضیه حد مرکزی نصف طول فاصله اطمینان $(1 - \eta)$ درصد عبارت است از:

$$Z_{1-\eta/2} \frac{(Var(\hat{a}(u)))^{\frac{1}{2}}}{a(u)}$$

که در آن Z_n ، n امین چندک توزیع نرمال استاندارد را نشان می‌دهد. خطای نسبی را به صورت زیر تعریف می‌کنیم:

$$RE[\hat{a}(u)] = Z_{1-\eta/2} \frac{(Var(\hat{a}(u)))^{\frac{1}{2}}}{a(u)}$$

۱.۷ احتمالات ورشکستگی با افق متناهی در مدل پواسون مرکب

در جدول ۱ برآوردهای $\psi(u, T)$ در مدل پواسون مرکب با اندازه‌های خسارت وایبول و پارامترهای $\beta = 2$, $\rho = 0.5$, $b = 2.1$ آمده است. تعداد تکرارها در شبیه‌سازی $n = 300000$ می‌باشد. با مقایسه خطاهای نسبی می‌بینیم که انتخاب \tilde{w}_u بر $w = 0$ ترجیح دارد.

جدول ۱: برآورد $\psi(u, T)$ در مدل پواسون مرکب با اندازه‌های خسارت وایبول $(1, 0.5)$

| u | پارامتر | برآورد $\psi(u, 10^6)$ | خطای نسبی به درصد |
|-----|-------------------|------------------------|-------------------|
| 200 | $w = 0$ | 1.3E-5 | 10.9 |
| 200 | $w = \tilde{w}_u$ | 1.49E-5 | 3.9 |
| 800 | $w = 0$ | 1.04E-11 | 17.5 |
| 800 | $w = \tilde{w}_u$ | 1.01E-11 | 4.7 |

جدول ۲: برآورد $\psi(u, T)$ در مدل مارکوف - مدل پواسون مرکب با وایبول $(1, 0.75)$ و وایبول $(1, 0.75)$

| u | برآورد $\psi(u)$ | خطای نسبی به درصد |
|-----|------------------|-------------------|
| 100 | 2.55E-4 | 3.4 |
| 400 | 1.69E-8 | 3.8 |

۲.۷ مدل مارکوف مدل - پواسون مرکب

جدولهای ۲ و ۳ برآوردهای $\psi(u, T)$ را در دو مدل مارکوف - مدل پواسون مرکب نشان می‌دهد. در جدول ۲ $n = 1000000$, $b = 2/1$, $w = \tilde{w}_n$ و در جدول ۳ $\tilde{w} = 0$ می‌باشد. در هر دو جدول F_2 و F_1 وایبول $(1, 0.75)$ با $\beta_2 = 1.1906$ بوده و $\beta_1 = 0.75$ و $\pi_1 = 0.5$, $\pi_2 = 0.5$, $w_1 = 0.2$, $w_2 = 0.3$ و $\rho = 0.38$ می‌باشند. در جدول ۲ F_1 وایبول $(1, 0.75)$ بوده و $\beta_1 = 2$ می‌باشد. در جدول ۳ F_1 پارتو $(1, 0.75)$ با $\beta_1 = 2$ است.

جدول ۳: برآورد $\psi(u, 50)$ در مدل مارکوف - مدل پواسون مرکب با پارتو $(2, 1)$ و وایبول $(1, 0.75)$

| u | برآورد $\psi(u)$ | خطای نسبی به درصد |
|-----|------------------|-------------------|
| 100 | 3.32E-5 | 6.8 |
| 400 | 8.53E-6 | 6.5 |

مراجع

- [1] Asmussen, S. (1985), Conjugate processes and the simulation of ruin problem, Stochastic processes and their applications, 20:213-229.
- [2] Asmussen, S. (2000), Ruin probabilities, World scientific, Singapore, New Jersey, London, Hong Kong.
- [3] Asmussen, S. and Binswanger, K. (1997), Simulation of ruin probabilities for subexponential claim, ASTIN Bulletin, 27(2):297-318.
- [4] Boots, N.K. and Shahabuddin, P. (2000), Simulation GI/GI/1 queues and insurance risk processes with subexponential distributions, Proceeding of the 2000 winter simulation conference 656-665, IEEE press, piscataway, New Jersey.
- [5] Boots, N.K. and Shahabuddin, P. (2001), A framework for simulating small ruin probabilities in insurance risk processes with subexponential distributions, Research Report, Dept. of Industrial Engineering and Operation Research, Columbia University, NY 10027.
- [6] Embrechts, P. Kluppelberg, C. and Mikosch, T. (1997), Modeling Extremal events, Springer-Verlag, Berlin, Heidelberg.
- [7] Emiliano, A. Valdez and Kelvin, Mo. (2002), Ruin probabilities with dependent claims, Faculty of Commerce and Economics, The University of New South Wales, Sydney, Australian 2052.

- [8] Heidelberger, P. (1995) Fast simulation of rare events in queueing and reliability models, *ACM Transaction on Modeling and Computer Simulation*, 6:43-85.
- [9] Juneja, S. and Shahabuddin, P. (1999), Simulation heavy-tailed processes using delayed hazard rate twisting, Research Report, Dept. of Industrial Engineering and Operation Research, Columbia University, NY 10027.
- [10] Lehtonen, T. and Nyrhinen, H. (1992), Simulating level- crossing probabilities by importance sampling, *Advances in applied probability* , 24:858-874.

تعیین دبی‌ها اوج^۱ با استفاده از داده‌های بالاتر از یک آستانه معین^۲

سید سعید موسوی ندوشنی

دانشکده صنعت آب و برق (شهید عباسپور)

چکیده: یکی از روش‌های برآورد سیلابها، تحلیل فراوانی^۳ است. در این روش با استفاده از توزیع آماری مناسب می‌توان چندک‌های^۴ سیل را محاسبه نمود. اما برای این کار ما به یک سری از داده نیاز داریم، که از کل آمار موجود اخذ می‌گردد. برای این گزینش دو روش پیشنهاد می‌گردد:

(۱) سری حداکثرهای سالانه (یک حداکثر در هر سال)

(۲) سری بالاتر از یک آستانه معین

در این مقاله بر روی نحوه نمونه‌گیری با استفاده از روش (۲) بحث خواهد شد، که معرف بهتری از جامعه آماری سیلابها خواهد بود. نحوه انتخاب آستانه و آزمون‌هایی نظیر همگنی، ایستایی و استقلال (بر روی داده‌های بالاتر از آستانه معین) مورد علاقه قرار می‌گیرد. انجام آزمون‌های فوق به ما مقدار تقریبی آستانه را دیکته می‌کند و در انتها ملاحظه خواهد شد که یافتن مقدار مذکور یک رویکرد چندگانه است.

واژه‌های کلیدی: سری، آستانه، تحلیل فراوانی، ایستایی، همگنی، نمونه‌گیری

۱ مقدمه

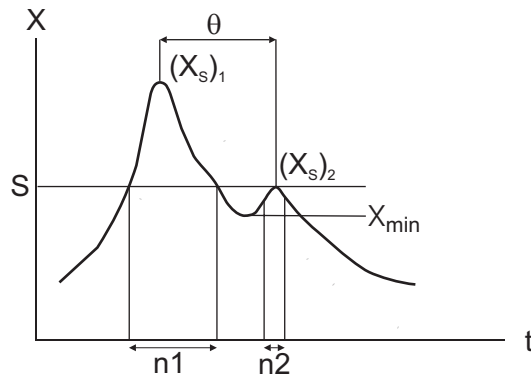
برنامه‌ریزی موثر و بهره‌برداری از طرح‌های آبی نیاز به آگاهی از رفتار احتمالی مقادیر حدی دارد. تحلیل فراوانی می‌تواند عهده‌دار تبیین چنین رفتاری باشد. تحلیل فراوانی با سریهای زیر قابل انجام است:

(۱) سری حداکثرهای سالانه (یک حداکثر در هر سال)

(۲) سری بالاتر از یک آستانه معین

با توجه به اینکه در روش آستانه معین، برای یکسال به انتخاب یک داده محدود نمی‌شویم، می‌توان مدعی شد که انتخاب عقلایی‌تری از جامعه سیلابها انجام شده است. با توجه به مقدار آستانه ممکن است که پاره‌ای از داده‌های روش حداکثر سالانه انتخاب نشود.

1) peak discharge 2) Peak Over Threshold 3) frequency analysis 4) quantiles



شکل ۱: نمودار تعیین ضابطه بازه زمانی بین سیلابها

روش دوم در مقایسه با روش اول، از اطلاعات بیشتری برخوردار است. استفاده از روش مذکور بالاخص در مواقعی که با داده‌ها کم روبرو هستیم، توصیه می‌شود.

۱.۱ نکات مورد بحث مقاله

در این مقاله ابتدا به نحوه نمونه‌گیری بالاتر از یک آستانه معین اشاره خواهد شد، سپس به مطالعه فرآیند پرداخته می‌شود و در خاتمه یک مطالعه موردی انجام خواهد شد.

۲ روش نمونه‌گیری با استفاده از آستانه معین

۱.۲ ضابطه استقلال دبی‌های اوج

در روش تحلیل فراوانی، فرض استقلال داده‌ها بررسی می‌شود. در ادبیات موضوع مورد مطالعه چندین ضابطه برای بررسی این فرض وجود دارد. مطابق (USWRC, 1976) بازه زمانی دو سیل متوالی پنج روز باضافه لگاریتم طبیعی سطح حوضه آبریز است (بر حسب مایل مربع)، علاوه بر این مطابق شکل ۱ جریانهای میانی بین دو نقطه اوج متوالی باید ۷۵ درصد کوچکترین دو دبی متوالی باشد. بنابراین نقطه اوج $(X_s)_2$ رد خواهد شد، چنانچه شرایط زیر برقرار شود.

$$X_{min} > \frac{3}{4} \min[Q_1, Q_2] \text{ یا } \theta < 5 \text{ days} + \ln(A)$$

کانن^۵ (۱۹۷۹) این ضابطه را ارائه می دهد که نقطه اوج $(X_s)_2$ رد می شود، اگر داشته باشیم:

$$\theta < 3T_p \text{ یا } X_{min} > \frac{2}{3}(X_s)_1$$

که در آن T_p متوسط اولین پنج هیدروگراف تیپ است. انتخاب ضابطه استقلال مساله پیچیده ای است، میکل^۶ (۱۹۸۴) حداقلی را برای θ بدست نمی دهد، اما برای آزمون استقلال ضرایب خودهمبستگی مرتبه اول و دوم برای سری زمانی مقادیر دبی های اوج را پیشنهاد می کند. اگر فرض صفر با سطح معنی دار بودن معینی رد شد، آنگاه مقدار بیشتری را برای بازه زمانی θ در نظر می گیریم.

۲.۲ انتخاب مقدار آستانه

برای این کار دو رویکرد متفاوت می توان اتخاذ نمود:

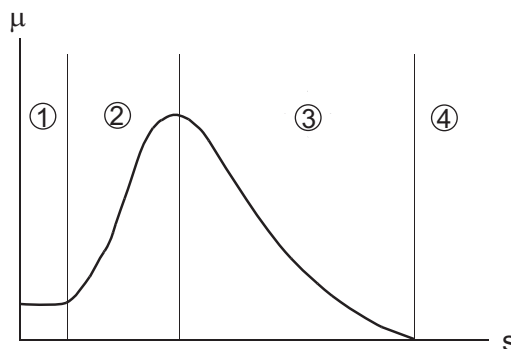
- ۱) استفاده از یک ضابطه فیزیکی مانند تعیین میزان سیلاب برای یک رودخانه مشخص
- ۲) استفاده از روش آماری، مقدار آستانه باید به گونه ای تعیین شود که دبی های اوج بالاتر از آن تشکیل یک سری داده ی مستقل از هم را بدهند و ضمناً فرآیند پواسون نیز آنرا تبیین کند. در اینجا آزمون های که برای انتخاب میزان آستانه، پیشنهاد شده است، مورد بررسی قرار می گیرد.

۱.۲.۲ آزمون ۱: متوسط تعداد داده های بالاتر از آستانه

شرط استقلال مشاهدات، یکی از فرضیات فرآیند پواسون است. این نکته قابل توجه است که انتخاب مقدار کوچک آستانه همیشه نمی تواند شرط استقلال را محقق سازد. انتخاب آستانه کوچک، نقاط اوج را بدست می دهد، که همیشه از هم مستقل نیستند. افزایش مقدار آستانه به تغییر شرایط کمک می کند. مطابق شکل ۱ اگر $(X_s)_2 < S < X_{min}$ باشد، آنگاه مقادیر وابسته را افزایش می دهد. بطور کلی با افزایش مقدار آستانه می توان چهار ناحیه را از هم تمیز داد. این نواحی در شکل ۲ نشان داده شده است که در آن μ متوسط تعداد نقاط اوج بالاتر از آستانه است. اکنون به شرح نواحی پرداخته می شود.

۱) ناحیه ۱: برای مقادیر اندک آستانه، کل سری زمانی (با طول NY ، یعنی N سال متوالی) بالاتر از آستانه قرار می گیرد.

۲) ناحیه ۲: در این ناحیه با افزایش مقدار آستانه، متوسط تعداد نقاط اوج بالاتر از آستانه افزایش می یابد.



شکل ۲: نمودار تغییرات متوسط مقادیر بالاتر از آستانه بر حسب مقدار آستانه

۳) ناحیه ۳: افزایش مقدار آستانه از یک اندازه به بعد، باعث کاهش متوسط تعداد نقاط اوج بالاتر از آستانه می‌گردد.

۴) ناحیه ۴: در این ناحیه هیچ نقطه‌ی اوجی بالاتر از آستانه قرار نمی‌گیرد، چون مقدار آستانه از کل سری بزرگتر است.

۲.۲.۲ آزمون ۲: متوسط مقادیر داده‌های بالاتر از آستانه

دیویسن و اسمیت^۷ (۱۹۹۰) و ندن و بلیس^۸ (۱۹۹۳) پیشنهاد کردند که انتخاب مقدار آستانه برای اینکه در دامنه مقادیر متوسط بالاتر از آستانه $(\bar{X}_s - S)$ قرار گیرد، باید تابع خطی از مقدار آستانه یعنی S باشد. این آزمون می‌تواند در انتخاب مقدار آستانه به گونه‌ای عمل کند که پارامترهای تابع چگالی احتمالی داده‌های بالاتر از آستانه را به حداکثر پایداری خود برسد. میکل^۹ (۱۹۸۴) توصیه می‌کند که تحلیل حساسیت چندک^{۱۰} برآورده شده را که تابعی از مقدار آستانه است، انجام شود.

۳.۲.۲ آزمون ۳: اندیس پراکنش

مقادیر آستانه بر اساس فرض پواسون بودن تعداد نقاط اوج، می‌تواند انتخاب شود. آشکار و رسل^{۱۱} (۱۹۸۷) این آزمون را بر اساس پیشنهاد کانن^{۱۲} (۱۹۷۹) انجام دادند. آستانه S طوری انتخاب می‌شود که اندیس پراکنش $\hat{i} = \frac{\sigma_{m_t}^2}{\mu_{m_t}}$ (که در آن m_t تعداد نقاط اوج در فاصله $[0, t]$)

7) Davison and Smith 8) Naden and Bayliss 9) Miquel 10) quantile
11) Ashkar and Rousselle 12) Cunnane

در فاصله اطمینان $[I_t(0,0.5), I_t(0,0.95)]$ قرار گیرد. شرح بیشتر این آزمون در ادامه مقاله خواهد آمد.

۴.۲.۲ بحث

رزبرگ و مدسن^{۱۳} (۱۹۹۲) توصیه نمودند وقتی که از آزمون ۱ و ۲ استفاده می‌شود که مقدار $\mu > 2$ باشد. در اینصورت توزیع GP یا توزیع نمایی منفی، دارای برآزش مناسبی است. متوسط مقادیر بالاتر از یک آستانه دارای تغییرات خطی نسبت به مقدار آستانه، در توزیع GP است و برای توزیع نمایی منفی تغییرات ثابت است.

(۱) توزیع GP

$$G_s(x) = \Pr[X_s < x] = 1 - [1 - (k/a)(x - S)]^{1/k}$$

$$\text{with } \hat{X}_S - S = \frac{a}{k+1} \quad (1)$$

$$G_{s^*}(x) = \Pr[X_{s^*} < x] = 1 - [1 - (k^*/a^*) \times (x - S^*)]^{1/k^*}$$

$$\text{with } k^* = k; a^* = a - k(S - S^*);$$

$$\bar{X}_{s^*} - S^* = \bar{X}_s - S - (S^* - S)k/(k+1)$$

(۲) توزیع نمایی منفی (توزیع GP با $k = 0$)

$$G_s(x) = 1 - \exp[(x - S)/a] \text{ با } \bar{X}_s - S = a \quad (2)$$

اگر آزمون شماره ۳ رد شود، می‌بایست از توزیعهایی نظیر دوجمله‌ای یا دوجمله‌ای منفی استفاده نمود. در نتیجه می‌توان ادعا نمود که روش کلی و واحد برای انتخاب مقدار آستانه وجود ندارد. بهر حال، انتخاب مقدار آستانه کاملاً به توزیع انتخاب شده و فرض استقلال بستگی دارد. بنابراین رویکردهای متفاوت می‌تواند به انتخاب آن کمک نماید. لذا خطوط کلی زیر که می‌تواند مقدار آستانه را معین سازد، بشرح زیر می‌باشد:

(۱) مشخص نمودن محدوده‌ایی از آستانه که آزمون‌های ۲ و ۳ را محقق می‌سازد.

(۲) انتخاب بزرگترین مقدار آستانه که در آن μ بزرگتر از ۲ یا ۳ باشد.

۳ مطالعه فرآیند وقوع

فرآیند وقوع حوادث E یا توسط مدت زمان بین دو حادثه متوالی θ تعیین می‌شود، یا تعداد حوادثی (m_t) است که در بازه زمانی $[0, t]$ رخ می‌دهد. در هر یک از حالات، می‌توان توزیع احتمال و مقدار میانگین را داشت.

(۱) فاصله بین دو حادثه

$$F(x) = \Pr[\theta < x],$$

$$f(x)dx = \Pr[x < \theta < x + dx]$$

دوره بازگشت حادثه بصورت زیر تعریف می‌شود:

$$T = E(\theta) = \int_0^{\infty} \theta f(\theta) d\theta \quad (۳)$$

(۲) تعداد m_t حادثه در فاصله زمانی $[0, t]$ رخ داده است که $w_k(t) = \Pr[m_t = k]$ می‌باشد. همچنین می‌توان $E(m_t)$ را روی فاصله زمانی $[0, t]$ تعریف نمود، و از اندیس پراکنش $I_t = \frac{\text{var}(m_t)}{E(m_t)}$ نیز می‌توان استفاده کرد.

لانگ^{۱۴} (۱۹۹۷) خواص توزیع‌های فرآیند را نشان داد. در این مرحله آزمون‌ها و رویه استفاده از آنها برای انتخاب توزیع و مقدار آستانه تشریح می‌شود.

۱.۳ آزمون اندیس پراکنش

از آزمون χ^2 می‌توان برای نیکویی برازش استفاده نمود، توصیه می‌شود که m_t را به حداقل پنج دسته با دست کم هفت و یا هشت عضو در هر دسته تقسیم نمود. معمولاً آمارهای دیگری برای رفع محدودیت آزمون χ^2 پیشنهاد می‌شود. کانن^{۱۵} (۱۹۷۹) برای توزیع پواسون اندیس پراکنش (I_t) را پیشنهاد نمود. فرض می‌کنیم که تعداد سیلابهای سالانه دارای توزیع نرمال است، مشروط بر آنکه پارامتر توزیع پواسون بیشتر از ۵ باشد و با توجه به اینکه میانگین و واریانس توزیع پواسون با هم برابر است، چنین بدست می‌آید:

(۱) عبارت $h = \sum_{i=1}^{NY} \{[m_{\lambda}(i) - E(m_{\lambda})]/[\text{var}(m_{\lambda})]^{1/2}\}^2$ از توزیع χ^2 پیروی می‌کند که دارای درجه آزادی NY است.

(۲) عبارت $\hat{h} = \sum_{i=1}^{NY} [m_{\lambda}(i) - \bar{m}_{\lambda}]^2 / \bar{m}_{\lambda} = (NY - 1)\hat{i}_{\lambda}$ از توزیع χ^2 پیروی می‌کند که دارای درجه آزادی $NY - 1$ است (یک درجه آزادی برای برآورد پارامتر توزیع

توزیع پواسون کم می‌شود) و $m_1(NY), m_1(2), \dots, m_1(1)$ تعداد حوادث در NY سال متوالی است. در حالت نظری مقدار اندیس پراکنش برابر واحد است.

لذا داریم:

(۱) اگر $I_1(\omega/2) < \hat{i}_1$ باشد، فرض پواسون بودن با سطح معنی‌دار $\omega\%$ حذف می‌شود و از توزیع دوجمله‌ای $\hat{i}_1 < 1$ استفاده می‌گردد.

(۲) اگر $\hat{i}_1 \in [I_1(\omega/2), I_1(1 - \omega/2)]$ باشد، فرض پواسون بودن پذیرفته می‌شود.

(۳) اگر $\hat{i}_1 > I_1(1 - \omega/2)$ از توزیع دوجمله‌ای منفی می‌توان استفاده نمود ($\hat{i}_1 > 1$)

با توجه به $I_1(p) = \chi^2(p)/(NY - 1)$ ($p = \omega/2$ or $p = 1 - \omega/2$) است.

۲.۳ آزمون ایستایی

لانگ^{۱۶} (۱۹۹۵) فاصله تولرانس m_t را در فاصله $[0, t]$ پیشنهاد کرد که در شکل شماره ۳ نشان داده شده است. اگر NF تعداد نقاط اوج مشاهده شده در فاصله $[0, t_{end}]$ باشد و $0 < t < t_{end}$ و $0 \leq m_t \leq NF$ ، آنگاه داریم:

$$-\mu t \leq \epsilon_t \leq NF - \mu t, \quad \epsilon_t = m_t - \mu t \quad (4)$$

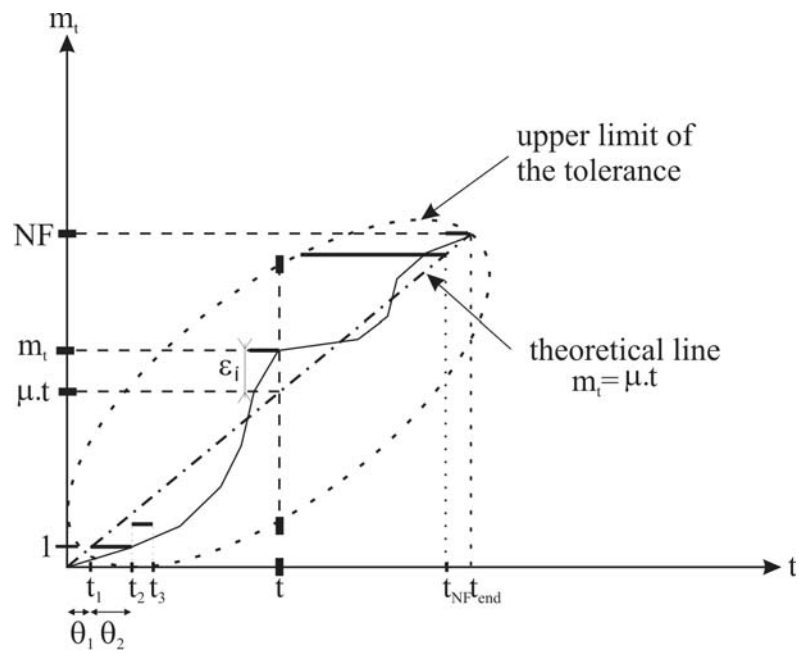
ما پیشنهاد می‌کنیم که فرآیند سیلاب توسط یک فرآیند پواسون همگن تبیین شود:

$$w_k(t) = \Pr(m_t = k) = e^{-\mu t} \frac{(\mu t)^k}{k!}, \quad E(m_t) = \mu t$$

اگر ما احتمال شرطی داشتن k نقطه اوج در فاصله $[0, t]$ را با $w_k^*(t)$ نشان دهیم، آنگاه NF نقطه اوج در فاصله $[0, t_{end}]$ خواهیم داشت و می‌توان نوشت:

$$\begin{aligned} w_k^*(t) &= \Pr[m_t = k | m_{t_{end}} = NF] \\ &= [w_k(t) w_{NF-k}(t_{end} - t)] / [w_{NF}(t_{end})] \\ &= \binom{NF}{k} (t/t_{end})^k (1 - t/t_{end})^{NF-k} \end{aligned} \quad (5)$$

با توجه به عبارت $w^*(t) = w_{NF-k}^*(t_{end} - t)$ ، فاصله تولرانس متقارن و به نقطه‌ای به مختصات $(t_{end}/2, NF/2)$ است.



شکل ۳: نمودار فاصله تولرانس تعداد مقادیر بالاتر از آستانه

۳.۳ آزمون فصلی

اردا و دیگران^{۱۷} (۱۹۹۳) نشان دادند که روند فصلی فرآیند جریان رودخانه می‌تواند تاثیر بسزایی روی توزیع مقادیر حدی بگذارد. در رویکرد فصلی دو عامل اهمیت دارد، مطالعه یک عامل که توام با ریسک است و دوم عدم ایستایی است که رخ می‌دهد.

۴ مدل نمودن نقاط اوج بالاتر از یک آستانه

اگر X یک متغیر تصادفی باشد، X_s مقدار حداکثر در هر حالت سیلابی را نشان می‌دهد، که بالاتر از مقدار آستانه S است. اکنون می‌توان تابع چگالی احتمال این متغیر را بدست آورد:

$$G_s(x) = \Pr[X_s < x] \quad (۶)$$

دوره بازگشت $T(x)$ را می‌توان بصورت زمان متوسط بین مقدار متوالی X_s که بزرگتر از x است، تعریف نمود. رزبرگ^{۱۸} (۱۹۸۵) رابطه بین دوره بازگشت و توزیع G_s بصورت رابطه (۷) تعریف نمود:

$$G_s(x) = 1 - \frac{1}{\mu T(x)} \quad (۷)$$

۱.۴ انتخاب توزیع نقاط اوج

اولین گام در تعیین توزیع این است که سری داده‌ها در شرایط استقلال، همگنی و ایستایی صدق کند دومین گام استفاده از یک توزیع مناسب و برآورد پارامترهای آن (به روش گشتاورها، حداکثر درست‌نمایی، و گشتاور وزنی احتمال).

توزیع نمایی و در حالت کلی GP دارای این خاصیت هستند که تابع چگالی شرطی آنها نیز به ترتیب برای هر مقدار آستانه S نمایی و GP می‌باشند. این خاصیت توسط آشکار و روسل^{۱۹} (۱۹۸۳) و ونگ^{۲۰} (۱۹۹۱) رویکردی برای استفاده از این توزیعها برای مقادیر بالاتر از یک آستانه شد. رویکرد شرطی بشرح زیر است:

$$\begin{aligned} G_{s^*}(x) &= \Pr[X_{s^*} < x] = \Pr[X_s < x | X_s > S^*] \\ &= \frac{\Pr[S^* < X_s < x]}{\Pr[X_s > S^*]} \end{aligned}$$

و

$$G_{s^*}(x) = \frac{G_s(x) - G_s(S^*)}{1 - G_s(S^*)} = 1 - \frac{1 - G_s(x)}{1 - G_s(S^*)} \quad (۸)$$

۵ توزیع حداکثر سیلاب سالانه

۱.۵ ارتباط بین توزیع‌های مقادیر بالاتراز آستانه و مقادیر حداکثر سالانه

متغیر X^* مقادیر حداکثر سالانه X با تابع توزیع F_x . روابط زیر توسط شان و لین^{۲۱} (۱۹۶۴) ارائه شده است:

$$F_{X^*}(x) = \Pr[X^* < x] = \sum_{k=0}^{\infty} w_k(\lambda) [G_s(x)]^k \quad (۹)$$

در حالت عمومی

$$F_{X^*}(x) = \exp\{-\mu[\lambda - G_s(x)]\} \quad (Poisson\ process) \quad (۱۰)$$

معادله شماره (۱۰) به نتایج زیر منجر می‌شود:

۱) مدل مقادیر بالاتر از آستانه با توزیع نمایی برای سیلابهای سالانه به سمت توزیع گامبل میل می‌کند:

$$\begin{aligned} G_s(x) &= \lambda - \exp[-(x - S)/a] \\ F_x(x) &= \exp\{-\mu \exp[-(x - S)/a]\} \end{aligned}$$

۲) مدل مقادیر بالاتر از آستانه باتوزیع پارتو تعمیم یافته^{۲۲} (GP) برای مقادیر حداکثر سالانه به سمت توزیع حدی تعمیم یافته^{۲۳} (GEV)

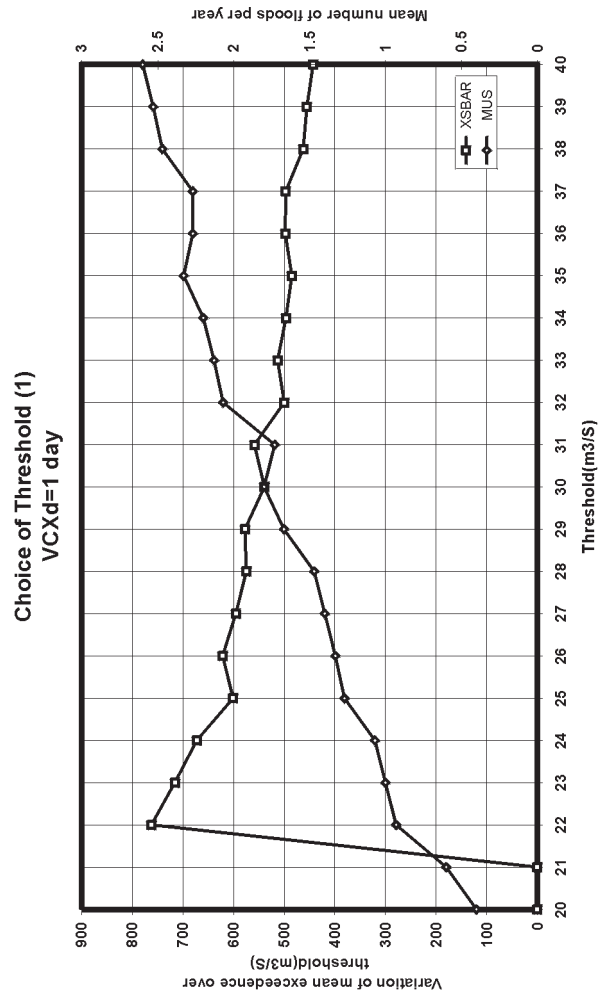
$$\begin{aligned} G_s(x) &= \lambda - [\lambda - (k/a)(x - S)]^{1/k} \quad (if\ k \neq 0) \\ F_x(x) &= \exp\{-\mu[\lambda - (k/a)(x - S)]^{1/k}\} \end{aligned}$$

۶ مطالعه حالتی

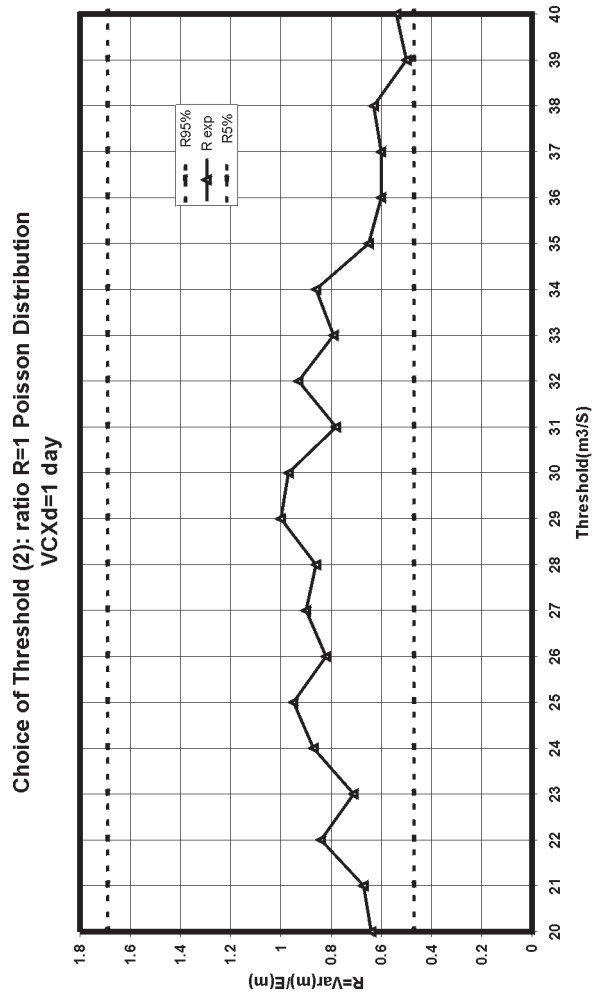
در این قسمت آمار دبی^{۲۴} متوسط روزانه یک ایستگاه هیدرومتری (در یک حوضه آبریز) اخذ شده است. در اولین بررسی سری زمانی داده‌ها رسم می‌شود، که از روی آن می‌توان محدوده‌ای را برای مقدار آستانه حدس زد. اما برای اینکه بتوان تغییرات متوسط تعداد داده‌ها و متوسط مقدار آنها را ملاحظه نمود، به کمک نرم‌افزارهای مربوطه، در محدوده‌ای که مدنظر است مقادیر متوسط‌ها محاسبه شده و در شکل شماره ۴، مشاهده می‌گردد.

شرط پواسون بودن تعداد داده‌های بالاتر از آستانه بر حسب مقادیر مختلف آن، توسط اندیس پراکنش کنترل شده است، که در شکل شماره ۵ ملاحظه می‌شود.

21) Shan and Lynn 22) Generalized Pareto 23) Generalized Extreme Value
24) discharge



شکل ۴: انتخاب مقدار آستانه (آزمون ۱: $\mu(S) = f_1(S)$ ، آزمون ۲: $\bar{X}_s - S = f_2(S)$)



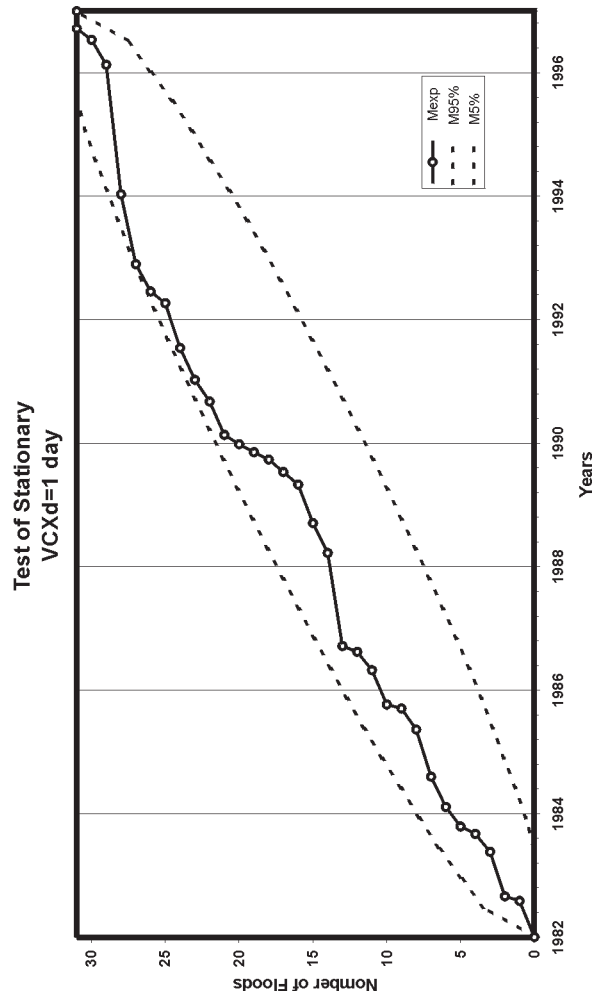
شکل ۵: آزمون اندیس پراکندگی (آزمون ۳)

از روی دو شکل ۴ و ۵، مقدار m^3/s را برای آستانه انتخاب می‌کنیم و سپس با استفاده از این مقدار می‌توان شرط ایستایی را کنترل نمود. مورد اخیر در شکل شماره ۶ قابل رویت است.

۷ نتیجه‌گیری

به علت رویکرد چندگانه روش داده‌های بالاتراز یک آستانه معین، چالش‌هایی در این روش مطرح شده است که موجب شهرت آن در بین کاربران شده است. در این مقاله سعی بر این شد که خطوط کلی که دارای همگرایی بیشتری بین صاحب نظران است برای استفاده از روش بالاتراز آستانه ارائه گردد. مسأله عمده این روش، انتخاب مقدار آستانه است، همانطور که ملاحظه شد نمی‌توان مقدار واحدی را برای آستانه بدست داد. اما می‌توان یک محدوده‌ای برای آستانه‌ها برگزید که در تجزیه و تحلیل‌های بعدی به نتایج یکسانی رسید. برای دستیابی به چنین محدوده‌ای باید آزمون‌های مطرح شده در این مقاله را مورد استفاده قرار داد.

سرانجام علیرغم مزایای این روش، توصیه می‌شود که تجزیه و تحلیل فراوانی برای روش مذکور و داده‌های حداکثر سالانه انجام شود، زیرا مقایسه نتایج آنها با یکدیگر مقایسه خوبی برای این دو روش می‌باشد.



شکل ۶: آزمون ایستایی فرآیند

مراجع

- [1] Ashkar, F., Rousselle, J. (1983) .The effect of certain restrictions imposed on the interarrival times of flood events on the poisson distribution used for modeling flood counts., Water Resources Res. 19(2), 481-485.
- [2] Ashkar, F., Rousselle, J. (1987). Partial duration series modeling under the assumption of a poissonian flood, count. J. Hydrol 90, 135-144.
- [3] Cunnane, C. (1979). A note on the poisson assumption in partial duration series models, Water Resources Res. 15(2), 489-494.
- [4] Davison, A.C., Smith, R.L. (1990) Models for exceedances over high threshold, J. R. Stat. Soc. B52(3), 393-442
- [5] Lang, M., (1995) Les chroniques en hydrologie. PhD thesis, Université Joseph Fourier Grenoble, Cemagref Lyon, France, 296 p.
- [6] Lang, M., Rassmussen, P.F., Oberlin, G., Bobée, B. (1997) Over-threshold sampling: modeling of occurrences by renewal processes, Revue des Sciences de l'Eau 3, 279-320
- [7] Miquel, J. (1984) Guide pratique d'estimation des probabilités de crue, Eyrolles, Paris 160 p.
- [8] Naden, P.S., Bayliss, A.C., (1993). Flood estimation: peak-over-threshold techniques., MaFF Conf. of River and Coastal Engineers, University of Loughborough, UK, 5-7 July, pp. 9.1.1-9.1.18
- [9] Ouarda, T.B, Ashkar, F., El-Jabi, N. (1993). Peaks over threshold model for seasonal flood variation, In: Kuo, C.Y. (Ed.), Engineering Hydrology (Pro. of symposium of San Francisco, California, 25-30 July), ASCE Publication pp. 341-346
- [10] Rosbjerg, D., (1985). Estimation in partial duration series with independent and dependent peak values, J. Hydrol. 76, 183-195.
- [11] Rosbjerg, D., Madsen, H., (1992). On the choice of threshold level in partial duration series, XVII Nordic Hydrological Conference, Alta, Norway, NHP Rep no. 30, pp. 604-615 23 (3), 245-256.
- [12] Shane, R.M., Lynn, W.R. (1964) Mathematical model for flood risk evaluation, J. Hydraul. Div. ASCE 90 (HY6), 1-20
- [13] Wang, Q.J. (1991) The POT model described by the generalized Pareto distribution with Poisson arrival rate., J. Hydrol. 129, 263-280

روش‌های نوین اقتصادسنجی در تحلیل داده‌های مکانی

رزیتا مویدفر، آذر ابراهیمی

سازمان مدیریت و برنامه‌ریزی استان اصفهان

چکیده: همواره دانشمندان علوم اقتصادی در جهت دستیابی به تحلیل‌های منطقی و نزدیکتر به واقعیت ویژگی‌های اقتصادی جوامع، به شیوه‌های کمی در پردازش داده‌های آماری توجه داشته‌اند. این مسأله شاخه‌ای را در علم اقتصاد بوجود آورد، به نام اقتصادسنجی که توانسته در چند دهه اخیر علاقمندان بسیاری را به خود مشغول و از ابعاد گوناگون ماهیت و روش، پیشرفتهای قابل توجهی نماید. در این رابطه مدل‌های اقتصادی به سه دسته تقسیم می‌شوند: ۱- ad hoc، ۲- time series و ۳- Econometrics.

استفاده از مدل‌های سری زمانی با فرض همگن بودن مشاهدات (individuals) با تغییر بخش‌ها یا واحدهای اقتصادی توأم می‌باشد و این مسأله بعضاً نتایج را دچار اربیب می‌سازد. در حالیکه مدل‌های panel data به عنوان یکی از پیشرفتهای اخیر در تخمین معادلات و تحلیل داده‌ها به روش اقتصادسنجی فرض ناهمگن بودن متغیر میان واحدهای مختلف اقتصادی را در یک روند زمانی امکان‌پذیر ساخته و با افزایش حجم مشاهدات و درجه آزادی در تخمین موجب کارایی بیشتر و واریانس کمتر مدل می‌گردد.

اما از طرف دیگر تخمین مدل‌های panel به روش اقتصادسنجی مرسوم به هنگام مواجهه با داده‌های مکانی، به دلیل وجود اثرات وابستگی فضایی و ناهمسانی فضایی منجر به نقض فروض گاس - مارکف در تخمین حداقل مربعات معمولی می‌گردد.

یکی از تحولات و پیشرفتهای ایجاد شده در بکارگیری روشهای کمی، تکامل شاخه اقتصادسنجی به اقتصادسنجی فضایی است. این تکنیک با وارد ساختن ماتریس مجاورت امکان اندازه‌گیری اثرات ناشی از مجاورتهای مکانی را فراهم کرده و مشکلات ناشی از تحلیل داده‌های مکانی به روش اقتصادسنجی مرسوم را برطرف می‌سازد.

این مقاله تلاش می‌کند تا علاوه بر ارائه ادبیات نظری مدل‌های panel، به مفهوم و موضوع اقتصادسنجی فضایی پرداخته و سپس با برقراری ارتباط بین این دو تکنیک، مدل مناسبی را در تحلیل داده‌های مکانی (منطقه‌ای) معرفی می‌نماید.

واژه‌های کلیدی: داده‌های مکانی، داده‌های پانل، اقتصادسنجی فضایی

۱ مقدمه

اقتصادسنجی به عنوان ترکیبی از تئوریهای اقتصادی، اقتصاد ریاضی، آمار اقتصادی و آمار ریاضی امکان کمی سازی مدل‌های اقتصادی را به منظور بررسی احکام و فرضیه‌های کیفی اقتصاد فراهم ساخته است.

هر چند به اعتقاد بسیاری از اقتصاددانان بیان روابط اقتصادی در قالب معادلات ریاضی و تخمین عددی آنها موجب ایجاد انحراف در پیش‌بینی مسیر اقتصاد گردیده و سیاستگذاران اقتصادی را در انتخاب سیاست‌های مناسب دچار انحراف می‌سازد، اما واقعیت آن است که روش اقتصادسنجی در تایید و توسعه تئوریهای اقتصادی نقش بسزایی داشته و از این رو به عنوان یکی از مباحث پژوهشی در علم اقتصاد مطرح می‌باشد. این امر گروهی از اقتصاددانان را بر آن داشته تا اساس تحقیقات خود را بر ابداع، اصلاح و پیشرفت روشهای اقتصادسنجی و کاربرد آن در تایید یا رد تئوریهای اقتصادی استوار سازند و بدین وسیله موجب گسترش حوزه نفوذ این بخش از اقتصاد در سایر بخش‌های آن گردند. از جمله این موارد می‌توان بکارگیری روشهای اقتصادسنجی در بررسی مدل‌های اقتصاد منطقه‌ای را نام برد. تئوریهای مطرح در این بخش از اقتصاد در بسیاری از موارد لزوم بررسی مشاهدات در دو بعد مکان و زمان را موجب می‌گردد و این مساله نتایج روشهای اقتصادسنجی معمول را دچار تورش می‌سازد. لذا اصلاح و پیشرفت این روشها به سمت روشهای نوین اقتصادسنجی در تحلیل داده‌های مکانی اجتناب پذیر بوده است. در این روند مدل‌های پانل و تکنیک اقتصادسنجی فضایی با هدف کاهش خطا در تخمین مدل‌های اقتصاد منطقه‌ای و تحلیل داده‌های مکانی شکل گرفت.

این مقاله می‌کوشد تا ضمن معرفی مباحث نظری این مدلها، مزیت‌های عملی بکارگیری آنها را طی تخمین یک مدل تجربی بر شمرد.

بخش اول این مقاله با ارائه مدل پایه اقتصادسنجی فروض اساسی تشکیل مدل‌های اقتصادسنجی را معرفی می‌نماید. در بخش دوم و سوم مباحث نظری مدل پانل و تکنیک اقتصادسنجی فضایی بر مبنای خدشه‌دار شدن فروض اساسی مدل‌های اقتصادسنجی معمولی در بکارگیری داده‌های مکانی شرح داده می‌شود و سپس بخش چهارم این مقاله نتایج بدست آمده در تخمین یک مدل تجربی اقتصاد منطقه‌ای را به روشهای اقتصادسنجی معمولی، پانل و اقتصادسنجی فضایی مورد مقایسه قرار می‌دهد و در نهایت بخش پنجم به نتیجه‌گیری از کل مباحث مطرح شده در این مقاله می‌پردازد.

۲ مدل پایه اقتصادسنجی

علم اقتصاد شکل گرفته بر مبنای تئوریها و فرضیه‌های اقتصادی، همواره به دنبال کمی سازی متغیرهای کیفی و تصریح فرضیه‌ها در قالب مدل‌های ریاضی به منظور فهم بهتر روابط حاکم بر اقتصاد جوامع بوده است.

در این رابطه مدل‌های اقتصادی در سه دسته قابل طبقه‌بندی گردید:

۱- Ad hoc: به معنی پذیرفتن فرضی معین بدون نیاز به سنجش و اندازه‌گیری آن
 ۲- time series: به معنی وابسته بودن یک متغیر به سابقه خودش که اساس مدل‌های خود رگرسیون (AR)، میانگین متحرک (MA) و یا ترکیبی از این دو (Arma و Arima) را تشکیل می‌دهد.

۳- Econometrics: به معنی تخمین کمی و تجربی مدل‌های اقتصادی ساخته شده بر مبنای روابط ریاضی که این مدل‌ها می‌تواند داده‌های سری زمانی یا مقطعی و یا ترکیبی از این دو را بکار گیرد و با استفاده از نتایج بدست آمده تئوریهای اقتصادی را رد یا اثبات نمایند. به عبارت دیگر اقتصادسنجی را می‌توان به عنوان تحلیل کمی پدیده‌های اقتصادی در دنیای واقع بر مبنای بسط و توسعه همپای تئوری و مشاهده که توسط روشهای متناسب استنتاجی به یکدیگر مرتبط شده‌اند، تعریف نمود.

ابزار اصلی تحلیل اقتصادسنجی رگرسیون است که اولین بار توسط فرانسیس گالتون^۱ (۱۸۸۶) و پس از آن کارل پیرسن^۲ (۱۹۰۳) مطرح شد که در بررسی ارتباط بین طول قد فرزندان و والدین مورد استفاده قرار گرفت. به طور کلی می‌توان گفت تحلیل‌های رگرسیون به مطالعه وابستگی یک متغیر (متغیر وابسته) به یک یا چند متغیر دیگر (متغیر توضیحی) می‌پردازد که با تخمین یا پیش‌بینی مقدار متوسط یا میانگین مقادیر متغیر نوع اول در حالتی که مقادیر متغیر نوع دوم معلوم یا معین شده باشد (در نمونه گیری تکراری)، صورت می‌پذیرد. با توجه به این تعریف می‌توان نوشت:

$$E(Y|x_i) = f(x_i)$$

که در آن $f(x_i)$ تابعی از متغیرهای X_i می‌باشد. این تابع صرفاً بیان می‌کند که میانگین توزیع Y بر حسب x_i معلوم، به طور تبعی به x_i مربوط است. به عبارت دیگر بیان می‌دارد که چگونه مقدار میانگین Y بر حسب x تغییر می‌کند. یک نوع ساده از این تابع رگرسیون خطی است که به شکل زیر نوشته می‌شود:

$$E(Y|x_i) = B_1 + B_2 x_i$$

از طرف دیگر می‌دانیم که y_i خاص در جامعه لزوماً منطبق بر میانگین شرطی خودش نبوده و دارای انحراف می‌باشد. لذا می‌توانیم انحراف y_i را در اطراف امید خودش همانند ذیل بدست آوریم:

$$U_i = Y_i - E(Y|x_i)$$

بنابراین با توجه به روابط قبلی خواهیم داشت:

$$Y_i = B_1 + B_2 x_i + u_i$$

اما واقعیت این است که در بسیاری از موارد دسترسی به مشاهدات جامعه کاری است دشوار و تنها مشاهدات نمونه‌ای از متغیرهای مورد نظر قابل اندازه‌گیری می‌باشند. پس تخمین تابع

1) Francis Galton 2) Karl Pearson

رگرسیون بر مبنای مشاهدات نمونه به صورت:

$$Y_i = \beta_1 + \beta_2 x_i + e_i$$

به عنوان یک تقریب، بیانگر رابطه تابعی رگرسیون جامعه خواهد بود. یکی از روشهای تخمین پارامترهای مدل روش حداقل مربعات مستقیم (OLS) می باشد که اولین بار توسط کارل فردریک گوس^۳ ریاضیدان آلمانی ارائه شد. همانطور که قبلاً گفته شد، مقادیر x_i داده شده در تابع رگرسیون مقادیری از Y_i را به صورت $Y_i = \beta_1 + \beta_2 x_i + e_i$ بدست می دهد که برابر با میانگین شرطی Y_i بوده و نسبت به مقادیر واقعی Y_i دارای مقدار انحرافی به اندازه e_i می باشد. روش حداقل مربعات معمولی بر مبنای حداقل نمودن مجموع مربعات عوامل خطا $\sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2$ دارای خصوصیات ویژه ای است که این روش را به صورت یکی از مشهورترین و فویترین روشهای تحلیل رگرسیون در آورده است. فرض اساسی حاکم بر تخمین زنده ها حداقل مربعات به صورت:

- ۱ - $E(u_i | x_i) = 0$
- ۲ - $cov(u_i, u_j) = 0$
- ۳ - $var(u_i | x_i) = \sigma^2$
- ۴ - $cov(u_i, x_i) = 0$

مدل رگرسیون دقیقاً تصریح شده است (عدم وجود خطای تصریح یا تورش) به گونه ای است که مطابق قضیه گوس - مارکف، تخمین زنده های حداقل مربعات را در بین تخمین زنده های خطی، بدون تورش و دارای حداقل واریانس می سازد یا به عبارت دیگر بهترین تخمین زن بدون تورش خطی (BLUE) می باشند. اما تغییر متغیرها در دو بعد مکان و زمان و فرض ناهمگنی بین واحدهای مختلف در یک مقطع از زمان و بین زمانهای مختلف وجود یک واحد، رابطه خطی رگرسیون را به صورت:

$$Y_{it} = \alpha_i + B_i x_{it} + u_{it}$$

تغییر داده و شرایطی را فراهم می سازد که موجب نقض قضیه گوس - مارکف گردیده و تخمین به روش حداقل مربعات معمولی را دچار تورش می سازد. در این حالت نتایج حاصل از تخمین به روش حداقل مربعات معمولی، دیگر BLUE نبوده و در اینجا روشهای جایگزین برای حداقل نمودن واریانس تخمین زنده ها مطرح می گردد. در این رابطه نه تنها مدل پانل به عنوان روشی برای تخمین ضرائب رگرسیون امکان لحاظ نمودن فروض ناهمگنی بین واحدهای مختلف را در طول زمان فراهم می سازد، بلکه تکنیک اقتصادسنجی فضایی نیز با وارد ساختن متغیر وابستگی فضایی، موجب کاهش خطای تخمین ناشی از وابستگی و ناهمسانی فضایی می گردد. بر این اساس در بخش های بعدی اساس مدل پانل و تکنیک اقتصادسنجی فضایی به تفصیل توضیح داده خواهد شد.

3) Kard Fredric Gauss

۳ مدل‌های شبکه‌ای (پانل)

در قسمت قبل نشان دادیم که یک معادله رگرسیون به دنبال تخمین رابطه‌ای علی بین یک متغیر وابسته و یک یا چند متغیر مستقل می‌باشد. آنچه تاکنون بدان پرداخته شد، تخمین آماری مشاهداتی بود که تغییرات آنها در یک بعد صورت می‌گرفت. به عبارت دیگر تغییرات یا در طول زمان مشاهده می‌شد و یا در بین واحدهای مختلف. مدل‌های پانل این امکان را فراهم می‌سازد که بتوانیم تفاوت بین داده‌ها در طول زمان را در میان واحدهای مختلف در نظر بگیریم. به این معنی که در یک سری زمانی مثلاً ۱۰ ساله برای هر سال می‌توان تغییرات مشاهدات را برای ۱۰ واحد ناهمگن مثلاً ۱۰ خانوار یا ۱۰ بنگاه اقتصادی و یا ۱۰ منطقه تجاری نیز مورد بررسی قرار داد. لذا در تخمین مدل‌های پانل به جای استفاده از متغیرها در حالت برداری از متغیرها به صورت ماتریس استفاده می‌شود. که البته این ماتریس‌ها بر اساس مشاهدات در دسترس می‌توانند متقارن یا نامتقارن باشند و در نتیجه پانل‌های متوازن (balance) یا نامتوازن (unbalance) را خواهیم داشت.

با توجه به این توضیحات می‌توان مزایای مدل‌های پانل را به صورت زیر طبقه‌بندی نمود:

۱- کنترل همگنی یا ناهمگنی بین واحدهای مختلف است که از تصریح نادرست (misspecification) در نتایج جلوگیری می‌کند.

۲- افزایش حجم مشاهدات که موجب افزایش درجه آزادی و کارایی بیشتر مدل و واریانس کمتر می‌گردد.

۳- مدل پانل پویایی‌های تعدیل (dynamics of adjustment) را مورد مطالعه قرار می‌دهد. مثلاً در سنجش بیکاری داده‌های مقطعی می‌توانند نسبتی از جمعیت را تخمین بزنند که در یک زمان خاص بیکار بوده اند ولی همان نسبت ممکن است در زمان دیگر بیکار نباشد. پس پانل قادر است رفتار افراد را از یک دوره به دوره دیگر مورد بررسی قرار دهد.

۴- در بسیاری موارد اثرات متغیرها روی متغیرهای دیگر به آسانی در مدل‌های مقطعی و سری‌های زمانی قابل شناسایی نیست. ولی در مدل پانل با در نظر گرفتن تغییرات مقطعی و زمانی این امکان بوجود می‌آید.

۵- در پانل به دلیل افزایش کارایی تخمین‌ها می‌توان مدل‌های رفتاری پیچیده‌تری را طراحی کرده و تخمین زد و در واقع قابلیت مدل از لحاظ آماری نیز افزایش می‌یابد.

۶- بررسی مسائل در سطح خرد موجب می‌شود در بسیاری از موارد داده‌ها به صورت تجمعی (aggregate) در نظر گرفته شوند که این موجب تورش یا اریب در داده‌ها در قلمرو خرد می‌شود. در مدل پانل امکان بررسی داده‌ها به صورت غیر تجمعی فراهم شده و تورش ناشی از تجمعی بودن داده‌ها بر طرف می‌شود.

حالت کلی در مدل‌های پانل به صورت:

$$y_{it} = \alpha_{it} + \beta x_{it} + u_{it}, \quad \begin{matrix} i = 1, \dots, N \\ t = 1, \dots, T \end{matrix}$$

می‌باشد.

در اینجا α_{it} عرض از مبدأ است x_{it} مجموعه متغیرهای مستقل به صورت it مشاهده روی k متغیر توضیحی می‌باشد. B یک بردار به صورت $1 \times K$ و نشان دهنده و در طول زمان ثابت است. مانند عامل مدیریت که برای واحدهای مختلف متفاوت ولی در زمان ثابت است. جزء بعدی شامل شیب پانل است. u_{it} جزء خطا و شامل سه قسمت است. جزئی از آن مربوط به متغیرهایی است که برای داده‌های مقطعی متفاوت متغیرهایی است در زمان مشابه، یکسان، اما در طول زمان متفاوت است. مانند قیمت، نرخ بهره و انتظارات نسبت به آینده. قسمت سوم جزئی است که نه تنها در طول زمان بلکه در مقطع هم تغییر می‌کند مثل سود که برای واحدهای مختلف و در زمان‌های مختلف، متفاوت است. در آزمون معنی‌داری ضرائب مدل پانل ۴ حالت قابل بررسی است:

حالت اول مبتنی بر همگن بودن شیب و نا همگن بودن عرض از مبدأ می‌باشد به صورت:

$$H_0 : \begin{matrix} \alpha_i^* \neq \alpha_j^* \\ \beta_i = \beta_j \end{matrix}$$

که در صورت تایید فرض H_0 مدل پانل به صورت $Y_{it} = \alpha_i^* + B'x_{it} + u_{it}$ خواهد بود. حالت دوم فرض می‌کند عرض از مبدأ همگن و شیب پانل ناهمگن می‌باشد به صورت:

$$H_0 : \begin{matrix} \alpha_i^* = \alpha_j^* \\ \beta_i \neq \beta_j \end{matrix}$$

که در صورت تایید فرض H_0 خواهیم داشت $Y_{it} = \alpha_i^* + b'_i x_{it} + u_{it}$ حالت سوم هر دو عرض از مبدأ و شیب پانل را ناهمگن فرض می‌کند و به صورت زیر بررسی می‌شود.

$$H_0 : \begin{matrix} \alpha_i^* \neq \alpha_j^* \\ \beta_i \neq \beta_j \end{matrix}$$

که در صورت تایید فرضیه H_0 مدل پانل به صورت $Y_{it} = \alpha_i^* + B_i^* x_{it} + u_{it}$ خواهد بود. و حالت چهارم فرض می‌کند که هر دو عرض از مبدأ و شیب پانل همگن بوده و فرضیه H_0 را به صورت زیر تعریف می‌کند:

$$H_0 : \begin{matrix} \alpha_i^* = \alpha_j^* \\ \beta_i = \beta_j \end{matrix}$$

که در صورت تایید فرضیه H_0 مدل پانل به صورت $Y_{it} = \alpha^* + B'x_{it} + u_{it}$ خواهد بود. در این حالت اساس همان مدل OLS تعمیم یافته است و اصطلاحاً به آن pooling data می‌گویند. نتایج مدل پانل از سه روش تخمین قابل محاسبه است:

- ۱- pooled data
- ۲- fixed effect
- ۳- Random effect

روش pooled data با فرض وجود همگنی بین ضرائب واحدهای مختلف به صورت OLS تعمیم یافته یا (general least square) GLS می‌باشد. در این روش فرضیه H_0 به صورت:

$$\alpha_1 = \alpha_2 = \Lambda = \alpha_n$$

$$\beta_1 = \beta_2 = \Lambda = \beta_n$$

تعریف می‌شود.

برای مقایسه و گزینش این روش نسبت به دو روش دیگر آماره F به صورت:

$$F(n-1, nT-n-k) = \frac{(R_u^2 - R_p^2) / (n-1)}{(\Lambda - R_u^2) / (nT-n-k)}$$

مورد بررسی قرار می‌گیرد.

در اینجا R_u^2 ضریب همبستگی حاصل از تخمین مدل به صورت دو روش دیگر و R_p^2 ضریب همبستگی حاصل از تخمین مدل به صورت روش pooled می‌باشد. در صورتی که آماره F از F جدول با تعداد مشاهدات (n) و تعداد متغیرها (k) و احتمال خطای معین بزرگتر باشد، فرضیه H_0 رد شده و روش‌های منتخب برای تخمین با نتایج کاراتر fixed effect یا Random effect خواهند بود.

در این دو روش فرض تخمین مدل مبتنی بر ناهمگنی بین ضرائب عرض از مبدأ و همگنی ضرائب شیب در میان مشاهدات می‌باشد، $\alpha_i \neq \alpha_j$ و $B_i = B_j$. فرض بر این است که α_i ها دارای تابع توزیع احتمال هستند، به میانگین α و δ_α^2 واریانس و مستقل از رگرسورها توزیع شده‌اند.

آزمون هاسمن دارای توزیع χ^2 (کای دو) بوده و امکان انتخاب روش کاراتر را از میان دو روش fixed و Random فراهم می‌سازد. در این آزمون فرضیه H_0 بر مبنای صحت روش Random به شکل زیر صورت می‌گیرد:

$E(u_{it}|x_{it}) = 0$ که در آن $\alpha_i = \alpha + \eta_i$ و $u_{it} = \varepsilon_{it} + \eta_i$ در اینجا x_{it} متغیرهای مستقل بوده و u_{it} عامل خطا می‌باشد که جزئی از آن به صورت η_i مربوط به متغیرهایی است که برای داده‌های مقطعی متفاوت و در طول زمان ثابت است بنابراین این جزء به صورت قسمتی از ضریب عرض از مبدأ (α_i) تعیین می‌شود.

\hat{B}_{RE} ضرائب شیبی است که به روش Random effect به دست آمده و ضرائب شیب به دست آمده از روش fixed effect می‌باشد.

با فرض $cov(\hat{B}_{RE}, \hat{q}) = 0$ خواهیم داشت:

$$var(\hat{q}) = var(\hat{B}_{FE}) - var(\hat{B}_{RE})$$

و بر این اساس آماره هاسمن به صورت:

$$H = \hat{q}' [var(\hat{q})]^{-1} \hat{q} \approx \chi_k^2$$

تعریف می‌شود که در اینجا k درجه آزادی و تعداد رگرورها می‌باشد. اگر آماره (H) کوچک باشد روش random کارا بوده و انتخاب می‌گردد و در غیر این صورت روش انتخابی اثر ثابت (fixed effect) برای تخمین ضرائب خواهد بود.

آنچه در معرفی مدل پانل به دست آمد، امکان بررسی تغییرات مشاهدات در طول زمان و در میان واحدهای مختلف به طور همزمان می‌باشد. یکی از حالت‌ها در این‌گونه تحلیل‌ها بررسی متغیرهای سری زمانی در میان مکان‌های مختلف است به عنوان مثال بررسی روابط تجاری در یک دوره ۲۰ ساله در میان کشورهای عضو یک منطقه تجاری مانند اکو (ECO) می‌باشد. در یک چنین تحقیقاتی علاوه بر آنکه استفاده از مدل پانل اجتناب ناپذیر می‌باشد، متغیر دیگری که نمایانگر اثرات حاصل از مکان وقوع پدیده‌های تجاری است، ضرورت حضور می‌یابد و این نقطه اتصال مدل‌های پانل به روش اقتصادسنجی فضایی است. از این رو در بخش بعدی به طور دقیق تکنیک اقتصادسنجی فضایی تشریح خواهد شد.

۴ اقتصادسنجی فضایی

گسترش کاربرد مدل‌های اقتصادسنجی در علوم منطقه‌ای به منظور تحلیل ارتباط موجود میان داده‌های کمی، تکنیک اقتصادسنجی فضایی را به عنوان زیر شاخه اقتصادسنجی به وجود آورد. به گونه‌ای که این تکنیک اثر متقابل فضایی (خود همبستگی فضایی) و ساختار فضایی (ناهمگنی فضایی) را بررسی نموده و مدل‌های رگرسیونی با داده‌های مقطعی یا پانل (panel) را مورد تخمین قرار می‌دهد.

همانطور که در بخش یک نیز توضیح داده شد مدل‌های حداقل مربعات معمولی بر پایه استقرار فروض گاوس - مارکوف به صورت استقلال همسانی واریانس مشاهدات دارای نتایجی بدون تورش خواهد بود. در حالیکه بکارگیری داده‌های مکانی در تخمین مدل‌های اقتصادسنجی مرسوم نتایج را با دو مسأله مواجه خواهد ساخت:

- ۱- وابستگی فضایی موجود میان مشاهدات
- ۲- ناهمسانی فضایی و این دو مسأله باعث نقض فروض گاوس - مارکوف شده و نتایج بدست آمده را توأم با تورش خواهد ساخت.

۱.۴ وابستگی فضایی

وابستگی فضایی در مجموعه‌ای از مشاهدات نمونه اشاره به این حقیقت دارد که یک مشاهده مربوط به یک مکان (i) در ارتباط با مشاهدات مکانهای j ($j \neq i$) قرار می‌گیرد. یعنی

$$Y_i = f(Y_j), \quad i = 1, \dots, n \quad i \neq j$$

i می‌تواند هر مقداری از ۱ تا n را بگیرد، یعنی در واقع وابستگی فضایی می‌تواند بین چند مشاهده باشد.

داده‌های مربوط به واحدهای فضایی مانند شهرها، استانها، ایالت‌ها، حوزه‌های سرشماری و . . . می‌توانند دارای خطای اندازه‌گیری باشند. این خطا به این دلیل رخ می‌دهد که مرزهای اداری برای جمع‌آوری اطلاعات دقیقاً منعکس کننده طبیعت فرآیند مربوط به ایجاد آن داده‌ها نیستند. مثلاً اندازه‌گیری نرخ بیکاری و اندازه نیروی کار، بر مبنای محل زندگی افراد می‌تواند نشانگر وابستگی فضایی باشد چون افراد برای یافتن شغل به شهرها یا استانهای مجاور می‌روند. از طرف دیگر بعد فضایی فعالیت‌های اجتماعی-دموگرافیک، اقتصادی یا منطقه‌ای جنبه مهمی از مسأله مدلسازی است. علوم منطقه‌ای بر مبنای این پیش فرض که مکان و فاصله نیروهای مهم دخیل در جغرافیای انسانی و فعالیت بازاریابی به عبارت دیگر فعالیت‌ها صرفاً تابعی از متغیرهای توضیحی شناخته شده نیست و تابعی از مجاورت و . . . است. همه این موارد در تئوری علوم منطقه‌ای موید وجود وابستگی فضایی و اثرات پراکندگی، سلسله مراتبی مکان‌ها و سرریزهای فضایی است.

۲.۴ ناهمسانی فضایی

اصطلاح ناهمسانی فضایی اشاره به تغییر روابط در طول فضا دارد. یعنی در هر نقطه از فضا انتظار وجود یک رابطه متفاوت را داریم:

$$y_i = x_i \beta_i + \epsilon_i, \quad i = 1, \dots, n$$

در اینجا i نشان‌دهنده نقاط موجود در فضا، x_i نشانگر بردار $(k \times 1)$ از متغیرهای توضیحی همراه با مجموعه پارامترهای β_i مربوط به آن، y_i متغیر وابسته در مکان i و ϵ_i بیانگر خطای تصادفی در رابطه مذکور است. به عبارت دیگر خواهیم داشت:

$$y_i = f_i(x_i, \epsilon_i)$$

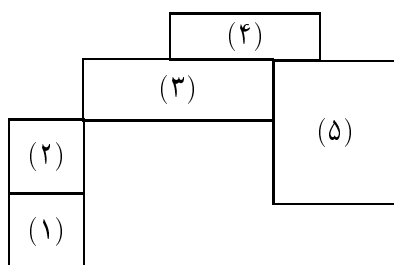
با توجه به رابطه مذکور هنگام حرکت در بین مشاهدات، توزیع داده‌های نمونه‌ای نشانگر میانگین و واریانس ثابتی نخواهد بود. بر این اساس کمی‌سازی بعد مکان و ورود آن به مدل‌های اقتصادسنجی به منظور رفع خطای ناشی از وابستگی فضایی و ناهمسانی واریانس مشاهدات مکانی، ما را به سمت طراحی مدل‌های اقتصادسنجی فضایی رهنمون می‌سازد.

۳.۴ کمی‌سازی مکان

- تعیین کمیت و مقدار عددی جنبه‌های مکانی به دو روش امکان‌پذیر خواهد بود.
- ۱- کمی‌سازی مجاورت فضایی
 - ۲- کمی‌سازی موقعیت فضایی

۱.۳.۴ کمی سازی مجاورت فضایی

مجاورت منعکس کننده موقعیت نسبی واحد منطقه ای یک مشاهده در فضا نسبت به سایر واحدهاست. اطلاعات مجاورت را می توان از روی نقشه بدست آورد. وابستگی فضایی بین واحدهای مجاور بیشتر است. برای مثال یک نمونه ۵ ناحیه ای به شکل مقابل را در نظر می گیریم.



ماتریس W با ابعاد 5×5 شامل ۲۵ عنصر صفر و یک است که نشان دهنده مجاورت میان این ۵ منطقه خواهد بود. برای تعریف مجاورت راههای مختلفی در نظر گرفته می شود.

-- مجاورت خطی: $W_{ij} = 1$ برای عناصری که یک کناره مشترک بلافاصله از سمت راست یا چپ ناحیه مورد نظر با آن داشته باشند. در این نمونه فرضی $W_{۵۳} = 1$ خواهد بود.

-- مجاورت رخ مانند: $W_{ij} = 1$ برای نواحی که با ناحیه مورد نظر یک ضلع مشترک داشته باشد. در ۵ ناحیه بالا $w_{۳۴} = 1$ و $w_{۳۵} = 1$ خواهد بود.

-- مجاورت فیل مانند: $W_{ij} = 1$ برای نواحی که یک گوشه مشترک با ناحیه مورد نظر دارند. در اینجا $w_{۳۳} = 1$ خواهد بود.

-- مجاورت خطی دو سویه: برای نواحی که بلافاصله در سمت چپ یا راست ناحیه مورد نظر قرار دارند. برای این مثال، همان نتایج مجاورت رخ مانند بدست خواهد آمد.

-- مجاورت رخ دو سویه: برای عناصر سمت راست، چپ، بالا و پایین $W_{ij} = 1$ خواهد بود.

-- مجاورت ملکه مانند: $W_{ij} = 1$ برای نواحی که یک ضلع یا زاویه مشترک با ناحیه مورد نظر داشته باشند. انتخاب هر کدام از این تعاریف بستگی به ماهیت مسأله و اطلاعات غیر نمونه ای دارد. در مثال فوق ماتریس مجاورت بر حسب تعریف مجاورت رخ به صورت زیر بدست خواهد آمد.

$$W = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{bmatrix}$$

این ماتریس یک ماتریس متقارن بوده و عناصر روی قطر اصلی صفر می باشد. عناصر ماتریس W غیرتصادفی و برونزا هستند. استانداردسازی عناصر این ماتریس بر حسب سطر به گونه ای

که مجموع عناصر هر سطر معادل یک شود موجب می‌شود تا از حاصلضرب ماتریس استاندارد مجاورت در بردار مربوط به متغیر وابسته برداری حاصل شود که عناصر آن میانگین مشاهدات نواحی مجاور باشد. به این ترتیب:

$$C = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & 0 & \frac{1}{3} & \frac{1}{3} & 0 \end{bmatrix}$$

$$\begin{bmatrix} y_1^* \\ y_2^* \\ y_3^* \\ y_4^* \\ y_5^* \end{bmatrix} = cy = \begin{bmatrix} y_2 \\ y_1 \\ 0,5y_4 + 0,5y_5 \\ 0,5y_3 + 0,5y_5 \\ 0,5y_3 + 0,5y_4 \end{bmatrix}$$

$$\begin{aligned} \rightarrow y_i &= f(y_j) \\ \rightarrow y &= Pcy + \epsilon \quad i \neq j \end{aligned}$$

در اینجا c ماتریس مجاورت استاندارد شده و y^* یک متغیر توضیحی در تعریف y می‌باشد و p پارامتر نشان دهنده وابستگی فضایی است. به این معنی که بخشی از کل تغییرات y در طول نمونه از طریق وابستگی آن به مناطق مجاورش توضیح داده می‌شود.

۲.۳.۴ کمی‌سازی موقعیت فضایی

مکان مربوط به مشاهدات در فضا در مدل‌سازی روابط ناهمسانی فضایی دارای اهمیت می‌شود. یکی از کسانی که در این زمینه به معرفی معرفی مدل پرداخت، کاستی^۴ می‌باشد. مدل مطرح شده به نام بسط فضایی است. به صورت زیر:

$$\beta = ZJ\beta, \quad Y = x\beta + \epsilon$$

در این مدل y بیانگر بردار $n * 1$ متغیر وابسته مربوط به مشاهدات فضایی و x یک ماتریس $n * nk$ شامل اقلام x_i نشان دهنده بردارهای $k * 1$ متغیر توضیحی است.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad x = \begin{bmatrix} x'_1 & 0 & \dots & 0 \\ 0 & x'_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & x'_n \end{bmatrix}$$

4) cosetti

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

اطلاعات مکانی در ماتریس z بیان شده که دارای عناصر zxi و zyi ، $i = 1, \dots, n$ است که مختصات طول و عرض هر مشاهده را نشان می‌دهد.

$$Z = \begin{bmatrix} Z_{x_1} \otimes I_k & Z_{y_1} \otimes I_k & \circ & \cdots \\ \circ & \vdots & \vdots & \circ \\ \vdots & \vdots & Z_{x_1} \otimes I_k & Z_{y_1} \otimes I_k \end{bmatrix}$$

$$J = \begin{bmatrix} I_k & \circ \\ \circ & I_k \\ \vdots & \vdots \\ \circ & I_k \end{bmatrix} \quad \beta_0 = \begin{bmatrix} \beta_x \\ \beta_y \end{bmatrix}$$

مدل نشان می‌دهد که پارامترها به صورت تابعی از مختصات طول و عرض تغییر می‌کنند. در اینجا واضح است که x و z و J اطلاعات موجود مشاهدات داده را نشان می‌دهند و تنها B_0 نشان دهنده پارامترهایی در مدل است که باید برآورد گردند.

مدل ناهمسانی فضایی را از طریق ایجاد امکان انحراف در رابطه محاسبه می‌کند، به گونه‌ای که گروههای مشاهدات مجاور یا همسایه مشخص شده با مختصات طول و عرض، مقادیر پارامتر مشابهی می‌گیرند.

هنگامی که مکان تغییر می‌یابد، رابطه رگرسیون تغییر می‌کند تا بر ارزش خطی محلی در میان گروه‌های مشاهداتی که تقریب نزدیکی برای یکدیگرند تطبیق یابد.

از دیگر روشهای مطرح شده برای برآورد انحراف در طول فضا که در زمینه اقتصادسنجی دارای کاربرد است، روش رگرسیون‌های وزنی جغرافیایی (GWR)^۵ است که توسط افرادی به نام کارلتون^۶، براندسون^۷ و فودرینگهام^۸، طراحی و معرفی شده در این مدل y نشان دهنده بردار $n * 1$ مشاهدات متغیر وابسته که از n نقطه در فضا به دست آمده، می‌باشد و x ماتریس $n * k$ متغیرهای توضیحی و بردار $n * 1$ خطاهای نرمال که دارای واریانس ثابت است. با فرض اینکه W_i نشانگر ماتریس قطری $n * n$ شامل وزنهایی بر مبنای فاصله برای مشاهده i باشد که منعکس کننده فاصله میان مشاهده i و سایر مشاهدات دیگر است می‌توان مدل GWR را به صورت زیر نوشت:

$$w_{iy} = w_i x \beta_i + w_i \epsilon_i$$

اندیس i در B_i نشان می‌دهد که بردار $k * 1$ پارامتر مربوط به مشاهده i است. مدل GWR، n مورد از چنین بردارهای مربوط به برآوردهای پارامترها را ایجاد می‌کند که هر یک برای یک

5) Geographically Weighted Regressions 6) Charlton 7) Brundson
8) Fotheringham

مشاهده است. این برآوردها با استفاده از رابطه زیر ایجاد می‌گردند:

$$B_i = (x'w_i'x)^{-1}(x'w_i'y)$$

۴.۴ یک مثال تجربی

با توجه به آنچه در دو قسمت قبل توضیح داده شد، هدف ما در اینجا این است که با ارائه نتایج یک تخمین تجربی بر این نکته متمرکز شویم، هرگاه در تئوریهای اقتصاد، مدل نظری علاوه بر زمان، تحت تأثیر مکان وقوع داده‌ها نیز قرار داشته باشد، و یا به عبارت دیگر امکان فرض همگنی در بین واحدهای مختلف برای یک یا چند متغیر وجود نداشته باشد، نه تنها پردازش مدل به روش پانل نتایج واقعی‌تری را به دست می‌دهد بلکه تکنیک اقتصادسنجی فضایی با در نظر گرفتن وابستگی‌های فضایی، تکنیک مناسب‌تری جهت تخمین مدل خواهد بود. همگرایی در درآمد سرانه بین مناطق مختلف مبنی بر تئوری رشد نئوکلاسیک به صورت:

$$\frac{\delta g_k}{\delta k} = \frac{sf(k)k - sf(k)}{k^2} = -s[f(k) - kf(k)]/k^2 < 0$$

$$\therefore \ln(y_{it}/y_{it-1}) = a - (1 - e^{-B})\ln y_{it-1} + u_{it}$$

یک نمونه از الگوهای است که تحت تأثیر زمان و مکان وقوع داده‌ها می‌باشد. در اینجا k سرمایه سرانه، s نرخ میل به پس‌انداز، g_k نرخ رشد سرمایه سرانه، B ضریب همگرایی و y_{it} تولید سرانه در مکان i در زمان t می‌باشد.

با توجه به الگوی ارائه شده، فرضیه همگرایی مطرح می‌کند که اقتصادها با سطوح پایین‌تر در درآمد سرانه به سمت رشد سریع‌تری در متغیرهای سرانه (درآمد سرانه و سرمایه سرانه) پیش می‌روند.

آزمون این فرضیه برای درآمد سرانه استانهای ایران طی دوره ۷۰-۱۳۸۰ به منظور مقایسه طی سه مرحله ۱- حداقل مربعات معمولی ۲- مدل پانل ۳- مدل پانل با استفاده از تکنیک اقتصادسنجی فضایی، صورت گرفت.

در مرحله اول میانگین نرخ رشد درآمد سرانه هر استان طی دوره ۷۰-۱۳۸۰ محاسبه گردید و سپس به روش حداقل مربعات معمولی با استفاده از مدل بین مقطعی (cross section) ارتباط بین میانگین نرخ رشد درآمد سرانه و درآمد سرانه سال پایه (۱۳۷۰) برای استانهای ایران محاسبه گردید. نتایج به دست آمده به صورت زیر مشاهده می‌شود:^۹

$$\ln(y_{i80}/y_{i70}) = \frac{0.206}{(0.51)} - \frac{(1 - e^{-0.17})}{(-0.47)} \ln(y_{i70}) + \epsilon$$

(۹) مراجعه شود به پیوست (۱)

$$R^2 = 0.9 \quad D.W = 1.86$$

هر چند ضریب مثبت B تأیید کننده وجود همگرایی در بین استانهای ایران می باشد، اما این ضریب معنی دار نبوده و تابع از قدرت توضیح دهندگی کافی برخوردار نمی باشد.

در مرحله دوم با استفاده از مدل پانل برای هر متغیر یک ماتریس که تعداد سطرهای آن برابر تعداد استانها و تعداد ستونها برابر تعداد سالهای مورد نظر بود در نظر گرفته شده و به این ترتیب حجم مشاهدات از ۲۴ مشاهده در مرحله اول به ۲۴۰ مشاهده در مرحله دوم افزایش یافت. تخمین مدل به روش غیر خطی در محیط Tsp 4.3 نتایج زیر را به دست می داد:

$$\ln(y_{it}/y_{it-1}) = \frac{3.32}{7.2} - (\lambda - e^{-0.36}) \ln(y_{it-1}) + \epsilon$$

$$R^2 = 0.19$$

همانطور که مشاهده می شود R^2 مدل نسبت به مرحله قبل افزایش یافته و معنی داری ضرایب به میزان قابل توجهی بالا رفته است و این نشان می دهد، روش تخمین روش مناسبی است و تنها متغیر مستقل از قدرت کافی برای توضیح تغییرات متغیر وابسته برخوردار نمی باشد.

در مرحله بعد با استفاده از تکنیک اقتصادسنجی فضایی، ماتریس مجاورت (w) به عنوان متغیر توضیحی وابستگی فضایی به صورت:

$$\ln(y_{it}/y_{it-1}) = a - (\lambda - e^{-B}) \ln y_{it-1} + Pw \ln(y_{it}/y_{it-1}) + \epsilon$$

وارد شد. نتایج به دست آمده به صورت زیر مشاهده می شود:

$$\ln(y_{it}/y_{it-1}) = \frac{3.02}{7.8} - (\lambda - e^{-0.31}) \ln y_{it-1} + 0.70 \ln(y_{it}/y_{it-1}) + \epsilon$$

$$R^2 = 0.47$$

در این مدل نه تنها معنی داری ضرائب افزایش یافته بلکه $R^2 = 0.47$ نیز نشان دهنده افزایش قدرت توضیحی متغیرهای مستقل در تغییرات متغیر وابسته می باشد از طرف دیگر ضریب مثبت متغیر وابستگی فضایی نمایانگر وجود اثرات مثبت ناشی از مجاورت در نرخ رشد در آمد سرانه هر استان بوده است.

به طور کلی میزان ضریب همگرایی این نتیجه را به دست می دهد که استانهای ایران با سرعت ۳۱ درصد در سال به سمت نقطه تعادل پایدار در رشد اقتصادی همگرا بوده و نیمی از شکاف رشد اقتصادی موجود بین آنها طی ۲/۵ سال از میان می رود.^{۱۰}

^{۱۰} (۱) مراجعه شود به پیوست (۱)

۵ نتیجه گیری

آنچه در این مقاله بدان پرداختیم معرفی مبنای نظری مدل پانل و کاربرد آن در تحلیل روابط بین داده‌هایی بود که در دو بعد بخش و زمان دارای ناهمگنی بوده و تخمین روابط خطی مابین آن‌ها جز در قالب ماتریس‌های پانل، نتایج را دچار تورش خواهد ساخت. هدف ما از طرح این مسأله پرداختن به داده‌هایی بود که به طور خاص نه تنها در طول زمان متغیر بوده بلکه تحت تأثیر مکان وقوع خود نیز می‌باشند و در اینجا بود که تکنیک اقتصادسنجی فضایی در اندازه‌گیری اثرات فضایی موضوعیت می‌یافت. در روش اقتصادسنجی معمول زمانی که داده‌ها در دو بعد زمان و مکان به صورت یک ماتریس در دسترس می‌باشند صرفنظر کردن از اثرات فضایی، نتایج را به دلیل وجود وابستگی و ناهمسانی واریانس فضایی بین داده‌ها، دچار خطا می‌سازد که این مسأله با وارد کردن متغیر فضایی با توجه به مبنای نظری روش اقتصادسنجی فضایی برطرف می‌گردد. این موضوع در آزمون تجربی همگرایی بین درآمد سرانه استان‌های ایران طی دوره ۷۰-۱۳۸۰ به خوبی ثابت گردید و نتایج نشان داد، کاربرد مدل پانل با استفاده از تکنیک اقتصادسنجی فضایی، بهترین گزینه در تخمین ضرائب توابعی است که ارتباط بین داده‌ها دو بعدی از نظر زمان و مکان می‌باشد.

مراجع

- [۱] راتو و میلر، اقتصادسنجی کاربردی، ترجمه حمید ابریشمی، مؤسسه تحقیقات پولی و بانکی، تهران، ۱۳۷۰.
- [۲] عسگری، علی و نعمت‌اله اکبری، روش‌شناسی اقتصادسنجی فضایی، تئوری و کاربرد، مجله پژوهشی دانشگاه اصفهان، جلد دوازدهم، شماره ۱ و ۲، ۱۳۸۰، صص ۹۳-۱۲۲.
- [۳] گجراتی، دامودار، مبانی اقتصادسنجی، ترجمه حمید ابریشمی، مؤسسه انتشارات و چاپ دانشگاه تهران، تهران، ۱۳۷۷.
- [4] Anselin, Luc, "Spatial Econometrics", Bruton Center School of Social Sciences, University of Texas at Dallas, Richardson, 1999
- [5] Baltaji, Badi. H, Econometric Analysis of Panel Data, John Wiley & sons Ltd, 2001
- [6] Bond, Stephen.R, "Dynamic Panel Data Models: A Guide to Micro Data Methods and Practice", Nuffield College, Oxford and Institute for Fiscal Studies, 2002
- [7] Greene, William. H, "Econometric Analysis", Prentice Hall, Fifth Edition, 2003

- [8] Levine , Ned , “Spatial Statistics and GIS , Software Tools to Quantify Spatial Patterns” , Computer Report , 1996
- [9] Rey , Sergio. J and Brett D. Montouri , “Us Regional Income Convergence : A Spatial Econometric Perspective , Department of Geography San Diego State University , San Diego , CA , USA , Regional Studies , vol 33 -2 , pp : 143-156 , 1997

پیوست (۱)

Dependent Variable: Y
 Method: Least Squares
 Date: 05/10/04 Time: 19:14
 Sample: 1 24
 Included observations: 24

| Variable | Coefficient | Std. Error | t-Statistic | Prob. |
|----------|-------------|------------|-------------|--------|
| C | 0.206257 | 0.404214 | 0.510267 | 0.6149 |
| X | -0.017040 | 0.036280 | -0.469672 | 0.6432 |

=====

| | | | |
|--------------------|-----------|-----------------------|-----------|
| R-squared | 0.009927 | Mean dependent var | 0.016429 |
| Adjusted R-squared | -0.035076 | S.D. dependent var | 0.028592 |
| S.E. of regression | 0.029089 | Akaike info criterion | -4.157281 |
| Sum squared resid | 0.018615 | Schwarz criterion | -4.059110 |
| Log likelihood | 51.88737 | F-statistic | 0.220592 |
| Durbin-Watson stat | 1.864353 | Prob(F-statistic) | 0.643211 |

=====

Current sample: 1 to 240

PANEL DATA ESTIMATION
 =====

Balanced data: NI= 24, T= 10, NOB= 240

TOTAL (plain OLS) Estimates:

Dependent variable: Y

Mean of dependent variable = .016429
 Std. error of regression = .131270
 Std. dev. of dependent var. = .143271 R-squared = .164017
 Sum of squared residuals=4.10118 Adjusted R-squared = .160504
 Variance of residuals = .017232

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.262174 | .037805 | -6.93481 |
| C | 2.91224 | .418877 | 6.95249 |

Standard Errors are heteroskedastic-consistent (HCTYPE=1).

BETWEEN (OLS on means) Estimates:

Dependent variable: Y

Mean of dependent variable = .016429
 Std. error of regression = .022971
 Std. dev. of dependent var. = .028592 R-squared = .382588
 Sum of squared residuals=.011609 Adjusted R-squared = .354524
 Variance of residuals = .527662E-03

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.101579 | .027512 | -3.69224 |
| C | 1.13841 | .303912 | 3.74586 |

WITHIN (fixed effects) Estimates:

Dependent variable: Y

Sum of squared residuals = 3.54057 R-squared = .278291
 Variance of residuals = .016468 Adjusted R-squared = .197729
 Std. error of regression = .128327

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.498626 | .062129 | -8.02563 |

Standard Errors are heteroskedastic-consistent (HCTYPE=1).

F test of A,B=Ai,B: F(23,215) = 1.4801, P-value = [.0791]
 Critical F value for diffuse prior (Leamer, p.114) = 6.4579

Variance Components (random effects) Estimates:

VWITH (variance of U_{it}) = 0.16468E-01
 VBET (variance of A_i) = 0.76408E-03
 (computed from small sample formula)
 THETA (0=WITHIN, 1=TOTAL) = 0.68307

Dependent variable: Y

Sum of squared residuals = 3.99415 R-squared = .185835
 Variance of residuals = .018577 Adjusted R-squared = .094951
 Std. error of regression = .136299

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.299537 | .041642 | -7.19321 |
| C | 3.32492 | .460056 | 7.22722 |

Hausman test of H0:RE vs. FE:CHISQ(1)=18.644, P-value=[.0000]

END OF OUTPUT.

| MEMORY USAGE: | ITEM: | DATA ARRAY | TOTAL MEMORY |
|-----------------------------------|-------|------------|--------------|
| UNITS: (4-BYTE WORDS) (MEGABYTES) | | | |
| MEMORY ALLOCATED | : | 500000 | 4.0 |
| MEMORY ACTUALLY REQUIRED | : | 7247 | 2.1 |
| CURRENT VARIABLE STORAGE | : | 51 | |

Current sample: 1 to 240

PANEL DATA ESTIMATION
=====

Balanced data: NI= 24, T= 10, NOB= 240

TOTAL (plain OLS) Estimates:

Dependent variable: Y

Mean of dependent variable = .016429
 Std. error of regression = .107360
 Std. dev. of dependent var. = .143271 R-squared = .443169
 Sum of squared residuals=2.73171 Adjusted R-squared = .438470
 Variance of residuals = .011526

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.235172 | .031862 | -7.38107 |
| G | .702933 | .075243 | 9.34215 |
| C | 2.60531 | .353358 | 7.37301 |

Standard Errors are heteroskedastic-consistent (HCTYPE=1).

BETWEEN (OLS on means) Estimates:

Dependent variable: Y

Mean of dependent variable = .016429
 Std. error of regression = .020700
 Std. dev. of dependent var. = .028592 R-squared = .521400
 Sum of squared residuals=.899863E-02 Adjusted R-squared=.475819
 Variance of residuals = .428506E-03

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.081598 | .026081 | -3.12866 |
| G | -.859328 | .348195 | -2.46795 |
| C | .928331 | .286796 | 3.23690 |

WITHIN (fixed effects) Estimates:

Dependent variable: Y

Sum of squared residuals = 2.31672 R-squared = .527762
 Variance of residuals = .010826 Adjusted R-squared = .472594
 Std. error of regression = .104047

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.411161 | .049500 | -8.30634 |
| G | .676947 | .074284 | 9.11299 |

Standard Errors are heteroskedastic-consistent (HCTYPE=1).

F test of A,B=Ai,B: F(23,214) = 1.6667, P-value = [.0329]
 Critical F value for diffuse prior (Leamer, p.114) = 6.4279

Variance Components (random effects) Estimates:

VWITH (variance of Uit) = 0.10826E-01
 VBET (variance of Ai) = 0.70043E-03
 (computed from small sample formula)
 THETA (0=WITHIN, 1=TOTAL) = 0.60716

Dependent variable: Y

Sum of squared residuals = 2.61332 R-squared = .467302
 Variance of residuals = .012212 Adjusted R-squared = .405070
 Std. error of regression = .110507

| Variable | Estimated Coefficient | Standard Error | t-statistic |
|----------|-----------------------|----------------|-------------|
| X | -.271065 | .034930 | -7.76019 |
| G | .701419 | .062810 | 11.1672 |
| C | 3.00178 | .385995 | 7.77674 |

Hausman test of H0:RE vs. FE: CHISQ(2)=16.861,P-value=[.0002]

END OF OUTPUT.

| MEMORY USAGE: | ITEM: | DATA ARRAY | TOTAL MEMORY |
|--------------------------|----------------|-------------|--------------|
| UNITS: | (4-BYTE WORDS) | (MEGABYTES) | |
| MEMORY ALLOCATED | : | 500000 | 4.0 |
| MEMORY ACTUALLY REQUIRED | : | 7948 | 2.1 |
| CURRENT VARIABLE STORAGE | : | 5266 | |

هم انباشتگی کسری در سری‌های زمانی و بررسی نرخ تورّم در ایران

سید محمود میرجلیلی^۱، قاسم تارمست^۲

^۱ دانشجوی کارشناسی ارشد دانشگاه شهید چمران اهواز*

^۲ هیئت علمی دانشگاه شهید چمران اهواز

چکیده: اغلب روشهای تجزیه و تحلیل سری‌های زمانی بر روی سری‌های زمانی ایستا انجام می‌شود در حالی که در بسیاری از زمینه‌ها مانند اقتصاد سربهائی به صورت نایستا وجود دارد. روشهای کلاسیک تجزیه و تحلیل سری‌های زمانی ابتدا سری‌های نایستا را به ایستا تبدیل می‌کنند که این کار با تفاضل‌گیری انجام می‌شود. روش تفاضل‌گیری در سری‌های بلند مدت باعث حذف روند بلند مدت سری می‌شود و در سری‌های کوتاه مدت باعث افزایش خطا می‌شود. برای حل این مشکل روش هم‌انباشتگی (cointegration) برای اولین بار توسط انگل - گرنجر (۱۹۸۷) معرفی شد بر اساس این روش هرگاه بتوان ترکیب خطی از چند سری نایستا با درجه نایستائی مشابه (انباشتگی Integration)، $I(d)$ را طوری پیدا کرد که آن ترکیب خطی ایستا باشد $I(0)$ به آن متغیرها هم‌انباشته می‌گویند. با این روش علی‌رغم نایستائی سری‌های تحت مطالعه، نتایج قابل تفسیر هستند و دیگر مسأله رگرسیون کاذب پیش نمی‌آید.

در روشهای سنتی سری‌های زمانی اندازه درجه انباشتگی یعنی d را یک عدد صحیح (اغلب عدد ۱ یا ۲) در نظر می‌گرفتند. در این تحقیق عدد d را یک عدد حقیقی در نظر می‌گیریم که، اگر d یک عدد حقیقی بین صفر تا دو باشد به سری انباشته کسری می‌گوئیم و برای اینکه تبدیل به سری ایستا شود باید تفاضل‌گیری کسری به اندازه d بر روی سری انجام شود. در رابطه هم‌انباشتگی نیز وقتی سری‌های نایستا با هم ترکیب می‌شوند باید ترکیب خطی آنها یک سری ایستا یعنی $I(0)$ باشد در حالی که در روش هم‌انباشتگی کسری ترکیب خطی آنها می‌تواند انباشته از درجه $\frac{1}{p}$ تا $\frac{1}{q}$ باشد.

در این تحقیق بعد از معرفی روش هم‌انباشتگی کسری، رابطه بین نرخ تورّم، نرخ ارز، حجم پول و تولید ناخالص داخلی از سال ۱۳۵۰ تا سال ۱۳۸۱ در ایران را بررسی می‌کنیم.

واژه‌های کلیدی: انباشتگی کسری، تفاضل‌گیری کسری، هم‌انباشتگی کسری

۱ مقدمه

سریهای زمانی به داده‌هایی گفته می‌شود که در طول زمان و با فواصل مساوی اندازه‌گیری می‌شوند. بطور کلی سری‌های زمانی را می‌توانیم به دو کلاس تقسیم کنیم کلاس اول شامل سری‌های است که طول زمان دارای میانگین ثابتی هستند و یا به عبارتی دیگر حرکت میانگین آنها به زمان بستگی ندارد، تابع خودهمبستگی آنها به سرعت کاهش می‌یابد، واریانس آنها محدود و ثابت است و در طول زمان تغییر نمی‌کند و چگالی طیفی آنها نیز برای تمام فرکانس‌ها وجود دارد، این‌گونه سریها را ایستای کواریانس می‌گویند.

کلاس دوم شامل سری‌های زمانی نایستا است، سری‌های نایستا در طول زمان میانگین ثابتی ندارند تابع خودهمبستگی آنها به کندی نزول می‌کند، دارای واریانس نامحدود هستند و تابع چگالی طیفی آنها برای تمام فرکانس‌ها وجود ندارد. اغلب روشهایی که برای تجزیه و تحلیل سری‌های زمانی وجود دارد بر روی سری‌های زمانی ایستا انجام می‌شود و در صورت مواجه شدن با سری‌های نایستا، با استفاده از روش تفاضل‌گیری آنها را تبدیل به سری‌های ایستا می‌کنند. اگر سری نایستا از درجه یک باشد با یکبار تفاضل‌گیری تبدیل به سری ایستا می‌شود و اگر نایستا از درجه دو باشد با دوبار تفاضل‌گیری و همین‌طور الی آخر. بنابراین اگر $\{x_t; t = 1, \dots, n\}$ یک سری زمانی باشد و داشته باشیم

$$(1 - B)^d x_t = \epsilon_t \quad \epsilon_t \sim WN(0, \sigma^2) \quad (1)$$

یعنی سری x_t با d بار تفاضل‌گیری به یک سری اغتشاش خالص ϵ_t که دارای توزیع نرمال با میانگین صفر و واریانس ثابت σ^2 است تبدیل می‌شود. اگر x_t بصورت معادله (۱) باشد آنرا انباشته از درجه d می‌گوئیم و آنرا با $x_t \sim I(d)$ نشان می‌دهیم. در اغلب کارهای کاربردی و روشهای سنتی سری‌های زمانی d را صفر، یک یا حداکثر دو در نظر می‌گیرند. اما در دو دهه اخیر یک کلاس دیگر از سری‌های زمانی توسط هاسکینگ^۱ (۱۹۸۱)، انگل^۲ (۱۹۸۳) و فاکس تاکو^۳ (۱۹۸۶) و غیره معرفی شده است. این کلاس از سری‌های زمانی بین سری‌های ایستا و نایستا قرار گرفته است که d یعنی پارامتر تفاضل‌گیری می‌تواند هر عدد طبیعی در بازه $(-\frac{1}{2}, 2)$ را بپذیرد. این دسته از سری‌های زمانی را سری‌های انباشته کسری می‌نامند. برای اینگونه سری‌ها تفاضل‌گیری در معادله (۱) را می‌توانیم بصورت زیر نشان دهیم:

$$(1 - B)^d = \sum_{j=0}^{\infty} \frac{\Gamma(j-d)}{\Gamma(-d)\gamma(j+1)} L^j \quad (2)$$

بطوریکه

$$\Gamma(a) = \begin{cases} \int_0^{\infty} x^{(a-1)} e^{-x} dx & a > 0 \\ \frac{(-1)^n}{\Gamma(a+1)} & a < 0 \\ \Gamma(a+1) = a\Gamma(a) & a = 1, \dots, n \end{cases}$$

1) Husking 2) Engle 3) Fax and Taquu

سریهایی که پارامتر تفاضل‌گیری آنها در بازه $\frac{1}{p} < d < \frac{1}{q}$ قرارگیرد را ایستا و سریهایی که $|d| \geq \frac{1}{p}$ باشد را نایستا می‌نامند.

روشهای چند متغیره سری‌های زمانی که ارتباط بین متغیرها را بررسی می‌کنند اغلب بر اساس فرض ایستایی سریها انجام می‌شوند و برای سری‌های نایستا نیز ابتدا آنها را با روش تفاضل‌گیری تبدیل به سری‌های ایستا می‌کنند، روش تفاضل‌گیری در سری‌های بلند مدت باعث حذف روند سری و در سری‌های کوتاه مدت باعث افزایش خطا می‌شود. برای رفع این مشکل و بدست آوردن رابطه بلند مدت بین متغیرها، روش هم‌انباشتگی برای اولین بار توسط انگل - گرنجر^۴ (۱۹۸۷) معرفی شد. بر اساس این روش اگر بتوان ترکیب خطی از متغیرهای نایستا را طوری پیدا کرد که آن ترکیب خطی ایستا باشد به آن متغیرها هم‌انباشته می‌گویند. فرض کنید X_t یک بردار $1 \times p$ متغیر از سری‌های زمانی باشد و داشته باشیم

$$\beta^l X_t = \epsilon_t \quad (3)$$

که β یک بردار $1 \times p$ از پارامترها باشد. در تعریف ابتدایی انگل و گرنجر (۱۹۸۷) سری‌های X_t همگی $I(1)$ بودند و ϵ_t می‌بایست $I(0)$ باشد تا رابطه هم‌انباشتگی وجود داشته باشد. آنها در مقاله خود اشاره‌ای به هم‌انباشتگی کسری نیز کردند. در رابطه هم‌انباشتگی کسری سری‌های X_t می‌توانند انباشته کسری، $I(d)$ باشند و باقیمانده‌های ϵ_t باید انباشته از درجه $I(d-b)$ ، که $b > 0$ است باشند.

حال مساله اصلی که پیش می‌آید تشخیص مقادیر d و b است. انگل و گرنجر (۱۹۸۷) برای اینکار از آزمونهای ریشه واحد دیککی و فولر^۵ و دیککی فولر تعمیم یافته^۶ (۱۹۸۱) استفاده می‌کردند. هدف اصلی ما نیز در این مقاله بررسی برآورد پارامتر تفاضل‌گیری یعنی d است. روشهای اولیه برای برآورد d ، ابتدا توسط جواک و پوتر هارک^۷ (۱۹۸۶) و کانچ^۸ (۱۹۸۷) معرفی شدند و سپس توسط ساول^۹ (۱۹۹۲) و رابینسون^{۱۰} (۱۹۹۵ a,b) و غیره گسترش یافتند. در زمینه هم‌انباشتگی کسری هم بعد از اشاره انگل و گرنجر (۱۹۶۷) مطالعه و تحقیق در این زمینه توسط چان وترین^{۱۱} (۱۹۹۵)، رابینسون (۱۹۹۴) رابینسون و ماریناسی^{۱۲} (۱۹۹۸) و ماریناسی و رابینسون (۲۰۰۱) و رابینسون و یاجیما^{۱۳} (۲۰۰۲) صورت گرفته است.

مدل‌بندی فرآیندهای کسری امکان بیشتری را برای وجود رابطه هم‌انباشتگی بین فرآیندها بوجود می‌آورند و تحقیقات جدید و خاصی را در این زمینه طلب می‌کنند. روش‌های مختلفی برای هم‌انباشتگی که در آنها درجه انباشتگی عدد صحیح فرض شده گسترش یافته‌اند، ولی این روشها، برای درجه‌های مختلف انباشتگی معتبر نیستند و باید برای فرآیندهای کسری تعمیم داده

4) Engle and Granger 5) Dicky-Fuller 6) Augmented Dicky-Fuller 7) Geweke and Poter Hudack 8) Kunsch 9) Sowell 10) Robinson ,P.M 11) Chan and Terrin 12) Robinson and Murrinaci 13) Yajimma

شوند. اما این که ما می‌توانیم بی نهایت مقدار برای d در نظر بگیریم کار را سخت می‌کند و برای حل این مشکل باید توسط آزمون‌هایی مقادیر d را از قبل تعیین کنیم. یک مساله پیچیده‌تر و سخت‌تر این است که هر یک از متغیرها دارای درجه انباشتگی متفاوت باشند و بخواهیم رابطه هم‌انباشتگی آنها را بررسی کنیم. که در این زمینه تحقیقات کمی صورت گرفته است. مقاله حاضر روشهایی را برای استنباط در مورد فرآیندهای کسری و نیز امکان رابطه هم‌انباشتگی کسری را بررسی می‌کند و در مورد برآورد پارامترهای هم‌انباشتگی بحث می‌کند. در ادامه این مقاله در بخش دوم تعاریف انباشتگی و هم‌انباشتگی کسری را مطرح می‌کنیم. در قسمت سوم برآورد پارامتر تفاضل‌گیری و آزمون فرض در مورد آن را با روشهایی که در حوزه فرکانس و برای فرکانس‌های نزدیک صفر، توسط رابینسون (۱۹۹۵)، لوبیتو^{۱۴} (۱۹۹۶) و لوبیتو و رابینسون (۱۹۹۸) معرفی شده‌اند را بررسی می‌کنیم. در قسمت چهارم برآورد بردار رگرسیون پارامترهای هم‌انباشتگی که در حوزه فرکانس و در یک باند کوچک حول فرکانس صفر توسط رابینسون و ماریناسی (۱۹۹۹) معرفی شده‌اند را بررسی می‌کنیم و بالاخره در قسمت پنجم، رابطه هم‌انباشتگی کسری بین نرخ تورم، نرخ ارز، حجم پول و تولید ناخالص داخلی در سالهای ۱۳۸۱ - ۱۳۵۰ در ایران را بررسی می‌کنیم.

۲ تعاریف انباشتگی و هم‌انباشتگی کسری

تعریفهای متفاوتی برای فرآیندهای انباشته کسری $I(d)$ وجود دارد یکی از آنها این طور بیان می‌کند که فرآیند $\{X_t; t \in Z\}$ را $I(d)$ می‌گوئیم اگر یک فرآیند $\{\epsilon_t; t \in Z\}$ که $I(0)$ است وجود داشته باشد بطوری‌که:

$$X_t = \mu + \Delta^{-d} \epsilon_t \mathbb{1}(t > 0) \quad (4)$$

که $\mathbb{1}(0)$ تابع نشانگر، $\Delta = (1 - B)$ ، عملگر وقفه و μ میانگین X_t است. فرض کنید X_t یک بردار $1 \times p$ متغیره باشد که هر کدام از آنها انباشته از درجه $d_i = 1, 2, \dots, p$ باشند، برای وجود رابطه هم‌انباشتگی حداقل باید دو سری زمانی داشته باشیم که درجه انباشتگی آنها مشابه باشد. حال اگر بیش از دو سری زمانی داشته باشیم درجه انباشتگی آنها می‌تواند متفاوت باشد. اما درجه انباشتگی حداقل دو سری که دارای بزرگترین درجه انباشتگی هستند باید با هم مشابه باشد. رابینسون و یاجیما (۲۰۰۲) بیان می‌کنند برای بررسی بهتر روابط هم‌انباشتگی بین سریها می‌توانیم آنها را با توجه به درجه انباشتگی‌شان در s مجموعه بصورت زیر افراز کنیم.

$$d_1 = \dots = d_{i_1} > d_{i_1+1} = \dots = d_{i_2} > \dots > d_{i_{s-1}+1} = \dots = d_{i_s} \quad (5)$$

که در آن $i_0 = 0$ و $1 \leq i_1 < i_2 < \dots < i_s = p$ است. فرض کنید

$$\beta = (\beta_1, \beta_2, \dots, \beta_p)' \quad (۶)$$

یک بردار p بعدی از ضرایب باشد. بردار β را نیز مانند (۵) افراز می‌کنیم.

$$\beta = (\beta^{(1)'}, \dots, \beta^{(s)'})' \quad (۷)$$

که در آن $p_l = i_l - i_{l-1}$ و $\beta^{(l)} = (\beta_{i_{l-1}+1}, \dots, \beta_{i_l})$, $l = 1, 2, \dots, s$ تعداد سری‌های زمانی در مجموعه i است و به همین صورت بردار X_t را افراز می‌کنیم.

$$X_t = (X_t^{(1)'}, \dots, X_t^{(s)'})', \quad X_t^{(l)} = (X_{i_{l-1}+1,t}, \dots, X_{i_l,t}) \quad (۸)$$

حال با توجه به افراز بالا تعاریف هم‌انباشتگی کسری را که توسط نویسندگان مختلف بیان شده است در اینجا می‌آوریم.

تعریف انگل - گرنجر (۱۹۸۷): اگر $d_1 = d_2 = \dots = d_p = d = 1$ باشد و یک بردار β وجود داشته باشد بطوریکه βX_t $I(d_u)$ باشد و $d_u = d - b = 0$ باشد آنگاه این سری هم‌انباشته از درجه $d - b$ یا $CI(d, b)$ هستند.

تعریف جوهانسون (۱۹۹۶): مؤلفه‌های بردار X_t را هم‌انباشته گوئیم اگر یک بردار β وجود داشته باشد بطوریکه βX_t ، $I(d_u)$ باشد و $d_u < d_1$ باشد. در این تعریف d_u که در درجه انباشتگی βX_t است باید از d_1 که بزرگترین درجه انباشتگی بین سریها است کوچکتر باشد.

تعریف فلورز و زافارز (۱۹۹۶): مؤلفه‌های بردار X_t را هم‌انباشته گوئیم اگر یک بردار β وجود داشته باشد بطوریکه $\beta(X_t)$ ، $I(d_u)$ باشد و $d_u < d_1$ باشد. در این تعریف ضریب متغیری که دارای بزرگترین درجه انباشتگی است نباید صفر باشد.

تعریف رابینسون و ماریناسی (۱۹۹۸): مؤلفه‌های بردار X_t را هم‌انباشته گوئیم اگر یک بردار $\beta \neq 0$ وجود داشته باشد بطوریکه $\beta(X_t)$ ، $I(d_u)$ باشد و $d_u < d_p$ باشد. در این تعریف d_u درجه انباشتگی βX_t باید از d_p که کوچکترین درجه انباشتگی بین سریها است کوچکتر باشد.

تعریف رابینسون و یاجیما (۲۰۰۲): مؤلفه‌های بردار X_t را هم‌انباشته گوئیم اگر یک بردار غیرصفر $\beta(l)$ وجود داشته باشد بطوریکه $\beta^{(l)}(l)X_t^{(l)}$ ، $I(d_u)$ باشد و $d_u < d_{i_l}$ باشد. در این تعریف بردار هم‌انباشتگی بصورت $\beta = (o', o', \dots, o', \beta^{(l)'}, o', \dots, o')'$ است. در ادامه سعی می‌کنیم با استفاده از یک مثال تعاریف بالا را نشان دهیم.

مثال: بردار $p = 5$ متغیره x_t را بصورت زیر در نظر بگیرید

$$x_t = (u_t + \varepsilon_{1t}, au_t + \varepsilon_{2t}, v_t + \varepsilon_{3t}, bv_t + \varepsilon_{4t}, \varepsilon_{5t}) \quad (۹)$$

که در آن $0 < e < d, v_t \sim I(e), u_t \sim I(d), a, b \neq 0$ (۱۹۹۶) $r = 1$ تعداد رابطه‌های هم‌انباشتگی است و بردار هم‌انباشتگی بصورت زیر است

$$\beta = (0, 0, -1, 0, 0)'$$

چون که درجه انباشتگی سربها متفاوت است تعریف انگل - گرنجر (۱۹۸۷) کاربردی ندارد. طبق تعریف جوهانسن (۱۹۹۶) $r = 4$ تعداد روابط هم‌انباشتگی است که بصورت زیر نشان می‌دهیم

$$\beta = (a, -1, 0, 0, 0)', \beta = (0, 0, b, -1, 0)', \beta = (0, 0, 0, 1, 0)', \beta = (0, 0, 0, 0, 1)'$$

طبق تعریف فلورز و زافارز (۱۹۹۶) تعداد روابط هم‌انباشتگی $r = 1$ است که بصورت زیر نمایش می‌دهیم

$$\beta = (0, -1, 0, 0, 0)'$$

طبق تعریف رابینسون و ماریناسی (۱۹۹۸) تعداد روابط هم‌انباشتگی $r = 0$ است. چون ترکیب خطی وجود ندارد درجه انباشتگی آن کمتر از کوچکترین درجه انباشتگی یعنی صفر باشد. طبق تعریف رابینسون و یاجیما (۲۰۰۴) تعداد روابط هم‌انباشتگی $r = 2$ است که بصورت زیر نشان می‌دهیم

$$\beta = (a, -1, 0, 0, 0)' \quad \beta = (0, 0, b, -1, 0)'$$

در تعریف رابطه هم‌انباشتگی باید به نکات زیر توجه کنیم:

الف) اگر یک بردار هم‌انباشتگی را بصورت نرمال تبدیل 16 کنیم آنگاه آن بردار منحصر به فرد می‌شود

ب) خطاهای حاصل از روابط هم‌انباشتگی می‌توانند انباشته از درجه‌های متفاوت باشند.

ج) تعاریف بالا فقط شامل ترکیبهای خطی می‌شوند و ترکیبهای غیر خطی را پوشش نمی‌دهند

۳ استنباط در مورد درجه انباشتگی

استنباط و برآورد مرتبه‌های انباشتگی کسری یک مساله کلیدی در هم‌انباشتگی کسری است. روشهای برآورد d همگی در حوزه فرکانس انجام می‌شوند. روشهایی که در حوزه زمان وجود دارند بیشتر برای آزمون فرض در مورد d بکار می‌روند در حوزه فرکانس نیز اغلب روشهایی بکار گرفته شده که بصورت نیمه پارامتری هستند و فرکانسهای نزدیک صفر در دوره نگار (برآورد تابع

چگالی طیف) را برای برآورد d بکار می‌برند روشهای برآورد نیمه پارامتری روشهای نیرومندی هستند که بدون اینکه نیازی به تعیین مدل داشته باشند می‌توانند سازگاری برآوردها را تامین کنند. اما در روشهای پارامتری مشخص نبودن پارامترهای مدل می‌تواند باعث ناسازگاری در برآوردها شود. تعریف انباشتگی که در معادله (۴) بیان کردیم فقط برای یک متغیر بود اما اگر یک بردار از متغیرها داشته باشیم باید به مدلی فکر کنیم که بطور توأم همه فرایندهای انباشته کسری را شامل شود. برآورد درجه انباشتگی برای سریهایی که دارای فرض ایستائی در کواریانس باشند و چگالی طیفی آنها موجود باشد بطور اساسی گسترش یافته‌اند. ما نیز در اینجا همین موضوع را ادامه می‌دهیم چون که خاصیت مجانبی سری‌های انباشته کسری با $\frac{1}{p} \leq |d|$ شبیه به سری‌های ایستا است. برای سری‌های نایستا $|d| > \frac{1}{p}$ نیز می‌توانیم با یک مرتبه تفاضل‌گیری خاصیت مجانبی سری‌های ایستای کواریانس را بوجود آوریم و پس از آنکه درجه انباشتگی را برآورد کردیم به آن مقدار یک اضافه کنیم [رابینسون (b) و (۱۹۹۵a) لوبیتو (۱۹۹۶)]. توجه کنید که در فرایندهای AR استنباط روی داده‌های تفاضل‌گیری شده باعث عدم کارائی می‌شود ولی در فرایندهای کسری وقتی فرض ما روی داده‌های تفاضل‌گیری شده باشد استنباط در مورد فرایندها هنوز دارای کارائی لازم است (رابینسون b) (۱۹۹۴).

برای برآورد درجه انباشتگی خطاهای e_t, d_e می‌توانیم روشهایمان را روی یک فرایند ξ_t که با $\hat{e}_t = y_t - \hat{\beta}x_t$ برآورد می‌شود انجام دهیم $\hat{\beta}$ هم با یکی از روشهایی که در بخش (۴) توضیح می‌دهیم برآورد می‌شود

یک بردار $1 \times p$ متغیره x_t که ایستا در کواریانس با تابع چگالی طیفی $f_x(\lambda)$ را در نظر بگیرد. (k,l) امین مولفه $f_x(\lambda)$ بصورت زیر است

$$f_{kl}(\lambda) \sim g_{kl} e^{i(\frac{\pi}{p})(\delta_k - \delta_l)} \lambda^{-\delta_k - \delta_l} \quad \lambda \rightarrow 0^+ \quad (10)$$

که در آن $k, l = 1, 2, \dots, p$. نشان می‌دهد که نسبت بخش حقیقی و موهومی سمت راست و چپ معادله به سمت یک میل می‌کند، برای $k = 1, 2, \dots, p$ $\frac{1}{p} < d < \frac{1}{p}$ ماتریس $G = [g_{kl}]$ برای حالتی که هم انباشتگی بین سریها وجود نداشته باشد معین مثبت و در غیر این صورت نیمه معین مثبت است و برای $k = 1, 2, \dots, p$ $g_{kk} > 0$ است.

دو راهکار کلی برای برآورد $d = (d_1, d_2, \dots, d_p)$ وجود دارد اولی براساس لگاریتم دوره نگار که توسط جواک و پوترهاک (۱۹۸۳) معرفی شده و خاصیت‌های مجانبی آن توسط رابینسون (a) (۱۹۹۵) و هرویچ و دیگران^{۱۷} بدست آمده است.

فرض کنید $\{x_t, t = 1, 2, \dots, n\}$ یک سری زمانی باشد تبدیل فوریه گسسته بصورت زیر انجام می‌شود

$$w_x(\lambda) = \frac{1}{(2\pi n)^{\frac{1}{2}}} \sum_{t=1}^n x_t e^{it\lambda} \quad (11)$$

همچنین اگر یک سری $\{y_t, t = 1, 2, \dots, n\}$ داشته باشیم دوره نگار متقابل بصورت زیر می شود

$$I_{xy} = w_x(\lambda)w_y(\lambda) \quad (12)$$

حال فرض کنید $I_{kk}(\lambda)$ ، k امین مولفه قطری $I_x(\lambda)$ باشد. برای اعداد صحیح $s, k = 1, 2, \dots, q$ که $y_{kj} = \log(I_{kk}(\lambda_j))$ ، $s < [\frac{n}{q}]$ ، $j = 1, 2, \dots, s$ علامت جزء صحیح و s پارامتر پنجره برای هموارسازی دوره نگار است را در نظر بگیرید، در این صورت بصورت زیر برآورد می شود

$$\tilde{d}_k = \frac{\sum_{j=1}^s \nu_j Y_{kj}}{\sum_{j=1}^s \nu_j Y_{kj}} \quad \nu_j = \log j - \frac{1}{s} \sum_{j=1}^s \log j \quad k = 1, \dots, q \quad (13)$$

توزیع مجانبی d_k توسط رابینسون (۱۹۹۵ a) بررسی شده که بصورت زیر است

$$\tilde{d}_k \sim N(d_k, \frac{\pi^2}{24s})$$

برای آزمون فرض $\Pi d = \rho$: H_0 می توانیم از آماره والد بصورت زیر استفاده کنیم

$$W = s(\Pi d - \rho)'(\Pi \Omega \Pi')^{-1}(\Pi d - \rho) \quad (14)$$

فرض صفر در صورتی رد می شود که این آماره بزرگ تر از χ_u^2 (چندک u ام توزیع χ^2 باشد). در این آماره $d = (d_1, d_2, \dots, d_p)$ و Ω یک برآورد سازگار از واریانس حدی ماتریس $2s(\tilde{d}_k - d)$ است (رابینسون ۱۹۹۵ a). توجه کنید که روشهای نیمه پارامتری 18 برای برآورد درجه انباشتگی دارای نرخ سازگاری $n^{\frac{1}{2}}$ هستند. بنابراین یک استنباط قابل اعتبار در این زمینه نیاز به حجم نمونه بزرگ دارد.

برای برآورد d می توانیم قیدهایی را روی d_k اعمال کنیم که در این صورت می توان برآوردهای کاراتری از d را که در ادامه کار رابینسون (۱۹۹۵ a) که توسط ماریتاسی و رابینسون (۲۰۰۱) معرفی شده اند بدست آورد. فرض کنید این قیود بصورت $d_1 = d_2 = \dots = d_p$ باشند مقدار مشترک d_* با برآوردهای GLS بصورت زیر بدست می آید

$$\tilde{d}_* = - \frac{\sum_{j=1}^s \nu_j \Omega^{-1} Y_j \nu_j}{\sum_{j=1}^s \nu_j^2 \Omega^{-1} \nu_j} \quad (15)$$

که در آن $y_j = (y_{1j}, \dots, y_{pj})$ است. توزیع مجانبی برآورد کننده \tilde{d}_* نیز توسط ماریناسی و رابینسون (۲۰۰۱) بصورت زیر بدست آمده است

$$\tilde{d}_* \sim N(d_*, \frac{1_p \tilde{\Omega}^{-1} 1_p}{\Psi_s})$$

بنابراین می‌توانیم آزمون والد را برای این مقدار مشترک d_* انجام دهیم.

کارایی برآوردهای \tilde{d} و d_* خیلی کمتر از برآوردهای کلاس دیگری از این برآوردها یعنی برآوردهای نیمه پارامتری گوسی، یا برآورد وایتل^{۱۹} که توسط کانچ^{۲۰} (۱۹۸۷) معرفی شد و سپس توسط رابینسون (b) (۱۹۹۵) و لوبیتو (۱۹۹۸) گسترش یافت. این برآوردها با استفاده از بهینه کردن تابع بزرگنمایی گوسی در فرکانسهای فوریه نزدیک به صفر بدست می‌آیند. برای برآوردهای تک متغیره d_* از تابع هدف زیر^{۲۱} که بصورت تک متغیره است استفاده می‌شود

$$R_k(d_k) = \log \left(\frac{1}{s} \sum_{j=1}^s I_{kk}(\lambda_j) j^{2d_k} \right) - \frac{2d_k}{s} \sum_{j=1}^s \log j \quad (۱۶)$$

و d_k بصورت $\bar{d}_k = \arg \min(R_k(d_k))$ برای $k = 1, 2, \dots, p$ با مینیمم کردن $R_k(d_k)$ برای $d \in (\frac{-1}{p}, \frac{1}{p})$ که ایستائی و وارون پذیری را بیان می‌کند برآورد می‌شود و هر کدام از \bar{d}_k ها دارای توزیع تقریبی $N(d_k, \frac{1}{\Psi_s})$ هستند، بنابراین این برآوردها کننده خیلی کاراتر از \tilde{d}_k است [رابینسون (b) (۱۹۹۵)]. یک برآورد کاراتر دیگر نیز توسط لوبیتو (۱۹۹۶) برای برآورد چند متغیره بصورت زیر معرفی شده است

$$R(\delta) = \log \left| \frac{1}{s} \sum_{j=1}^s A_j(d) \right| - \frac{2}{s} \sum_{k=1}^p d_k \sum_{j=1}^s \log j \quad (۱۷)$$

که در آن $\Lambda_j(d) = \text{diag}(e^{\frac{i\pi d_1}{p}} j^{d_1}, \dots, e^{\frac{i\pi d_p}{p}} j^{d_p})$ و $A_j(d) = \text{Re}(\Lambda_j(d) I_{xx}(\lambda_j) \Lambda_j(d))$ و برآورد \hat{d} بصورت $\hat{d} = \arg \min(R(d))$ در بازه فشرده $d \in (\frac{-1}{p}, \frac{1}{p})^p$ بدست می‌آید. \hat{d}_k و \tilde{d}_* برخلاف \bar{d}_k و \tilde{d} در یک فرم بسته تعریف نمی‌شوند. بهر حال یک روش دیگر برای برآورد d_k روش مرحله‌ای نیوتن است. در گام اول آن باید یک برآورد سازگار d با نرخ همگرایی $s^{\frac{1}{2}}$ بدست آوریم. در این روش هر چند در مراحل بعدی کارائی مرتبه اول بهبود پیدا نمی‌کند اما کارائی مراتب بالاتر بهبود پیدا می‌کند [رابینسون (۱۹۹۸)]. در این برآورد \hat{d} در مرحله $v + 1$ بصورت زیر بدست می‌آید

$$\hat{d}^{[v+1]} = \hat{d}^{[v]} - \left\{ 2 \left(I_q + \hat{G}^{[v]} \circ \hat{G}^{[v]-1} \right) \right\}^{-1} \frac{\partial R(\hat{d}^{[v]})}{\partial R} \quad (۱۸)$$

که در آن $v \geq 0$ و $\hat{d}^{[v]}$ برآورد مرحله قبلی آن، \circ نشان دهنده ضرب هادامارد^{۲۲}، $\hat{G}(d) = s^{-1} \sum_{j=1}^s A_j(d)$ و $\hat{G}^{[v]} = \hat{G}^{[v]}(\hat{d}^{[v]})$ در معادله (۱۷) ماتریس داخل براکت یک برآورد سازگار از احتمال حدی $\frac{\partial R(\hat{d}^{[v]})}{\partial d \partial d'}$ است. توجه کنید که $(I_p + \hat{G}^{[v]} \circ \hat{G}^{[v]^{-1}})$ وقتی که $p = 1$ باشد به Ψ کاهش پیدا می‌کند [ماریناسی و رابینسون (۲۰۰۱)].

یک انتخاب برای $\hat{d}, \hat{d}^{[1]}$ است، چون که مانند \hat{d} دارای نرخ سازگاری $s^{\frac{1}{2}}$ است با این حال باز دارای کارایی کمتری است. خصوصیتی که در بالا ذکر کردیم برای فرض عدم هم‌انباشستگی بود. اما در برآوردهای انفرادی که اول توضیح دادیم فرض عدم هم‌انباشستگی لازم نیست. در نهایت یک برآورد کارایی دیگر که با توجه به قید $d_1 = d_2 = \dots = d_p = d_*$ توسط ماریناسی و رابینسون (۲۰۰۱) معرفی شده و d_* در آن نامعلوم است بصورت زیر است

$$R_*(d) = \log \left| \frac{1}{s} \sum_{j=1}^s A_j(d_* \setminus p) \right| - \frac{p}{s} \sum_{j=1}^s \log j \quad (19)$$

و $\hat{d} = \arg \min(R_*(d))$ ، برآورد مرحله‌ای نیوتن با همین کارایی بصورت زیر بدست می‌آید

$$\hat{d}_*^{[v+1]} = \hat{d}_*^{[v]} - \frac{1}{p} \frac{\partial R(\hat{d}_*^{[v]} \setminus p)}{\partial d} \quad v \geq 0 \quad (20)$$

و برای $\hat{d}_*^{[1]}$ یک برآورد با $s^{\frac{1}{2}}$ سازگاری در نظر می‌گیریم. برای انجام آزمون فرض برای \hat{d}, \hat{d}_* و برآوردهای مرحله‌ای نیوتن مطابق باینها می‌توانیم از آماره والد استفاده کنیم

$$\hat{W}_{\Pi d = \rho} = (\Pi \hat{d} - \rho) \left[\Psi \Pi \left(I_p + \hat{G}(\hat{d}_*) \circ \hat{G}(\hat{d}_*)^{-1} \right) \Pi' \right] (\Pi \hat{d} - \rho) \quad (21)$$

که دارای توزیع حدی χ_u^2 است و به همین صورت برای آزمون d_* می‌توان با استفاده از تقریب $\hat{d}_* \sim N(d_*, (ps)^{-1})$ عمل نمود.

تابعهای هدف (۱۵)، (۱۶) و (۱۸) پیشنهاد می‌کنند از آزمونهایی بر اساس اصل چند گانه لاگرانژ^{۲۳} (LM) یا اصل نسبت درست‌نمایی^{۲۴} (LR) استفاده کنیم. یک آزمون LM برای $d_1 = d_2 = \dots = d_p$ که برای فرض مقابلی متفاوت از (۱۰) بکار می‌رود، توسط لوبیتو و رابینسون (۱۹۹۸) معرفی شده است و یک آزمون LR برای حالت $q = 1$ توسط رابینسون (۱۹۹۸) معرفی شده است. برای آزمون فرض $\Pi d = \rho$ از آماره LM زیر استفاده می‌کنیم

$$LM_{\Pi d = \rho} = s \frac{\partial R(d_*)}{\partial d} \left[\Psi \Pi \left(I_p + \hat{G}(\hat{d}_*) \circ \hat{G}(\hat{d}_*)^{-1} \right) \Pi \right] \frac{\partial R(d_*)}{\partial d} \quad (22)$$

که در آن \hat{d}_0 می‌تواند برآوردی باشد که $R(d)$ را تحت فرض $\Pi d = \rho$ مینیمم می‌کند یا می‌تواند یک تقریب نیوتن از d باشد. طبق تئوری مجانبی رابینسون (b ۱۹۹۵) و لوبیتو (۱۹۹۶) آماره (۲۲) دارای توزیع مجانبی χ_u^2 است. برای آزمون $H_0 : d_* = d_0$ یک آماره LM بصورت زیر است

$$LM_{d_* = d_0} = s \left(\frac{\partial R(d_0)}{\partial d} \right)^2 / 4p \quad (23)$$

که تحت فرض صفر دارای توزیع حدی χ_u^2 است. آماره LR برای فرضهای $H_0 : \Pi d = \rho$ و $H_0 : d_* = d_0$ بترتیب بصورت زیر است

$$LR_{\Pi d = \rho} = 2s \left(R(\hat{d}_0) - R(\hat{d}) \right) \quad LR_{d_* = d_0} = 2s \left(R(\hat{d}_{*0}) - R(\hat{d}_*) \right) \quad (24)$$

که تحت فرض صفر دارای توزیع مجانبی χ_u^2 و χ_u^2 هستند.

۴ برآورد بردارهای هم‌انباشتگی

بعد از برآورد درجه انباشتگی بین متغیرها برای بررسی رابطه بین آنها برآورد β را در رابطه زیر توضیح می‌دهیم

$$y_t = \beta' x_t + e_t \quad (25)$$

برای بردار قابل مشاهده $z_t = (x_t', y_t')$ و طبق رابطه‌های (۱۱) و (۱۲) دوره نگار میانگین‌گیری شده متقابل^{۲۵} را بصورت زیر تعریف می‌کنیم

$$F_{xy}(m) = 2Re \left[\frac{2\pi}{n} \sum_{j=1}^m I_{xy}(\lambda_j) \right] - \frac{2\pi}{n} I_{xy}(\pi) l(m = \frac{n}{2}) \quad (26)$$

که در آن $l(\cdot)$ تابع نشانگر و عدد صحیح m در شرط $1 \leq m \leq \frac{n}{2}$ صدق می‌کند. آخرین جمله معادله (۲۳) وقتی وجود دارد که n زوج باشد و m بیشترین مقدار خود یعنی $\frac{n}{2}$ را بگیرد وقتی که $m = [\frac{n}{2}]$ [.] تابع جزء صحیح است) عبارت زیر را می‌توانیم استنتاج کنیم

$$\hat{F}_{xy}([\frac{n}{2}]) = \frac{1}{n} \sum_{t=1}^n n(x_t - \bar{x})(y_t - \bar{y})' \quad (27)$$

که آن را کواریانس تصحیح شده میانگین^{۲۶} می‌نامیم و \bar{x} و \bar{y} میانگین نمونه هستند همانطور که مشاهده می‌کنیم \hat{F}_{xy} فرکانسهای $[\lambda_1, \lambda_m]$ را برای تولید کواریانس نمونه با هم سهیم می‌کند؛

معادلات (۷) و (۸) را برای حالتی که تصحیح میانگین صورت نگرفته باشد را نیز می‌توانیم بدست آوریم که بدلیل حجیم شدن مطلب از آوردن آن در اینجا خودداری می‌کنیم.

حال برآورد β به روش کمترین مربعات در حوزه فرکانس^{۲۷} (FDLS) را بصورت زیر معرفی می‌کنیم

$$\hat{\beta}_m = \hat{F}_{xx}(m)^{-1} \hat{F}_{xy}(m) \quad (28)$$

و فرض می‌کنیم که معکوس $\hat{F}_{xx}(m)$ وجود داشته باشد شاموی^{۲۸} (۱۹۸۶). یک حالت خاص برای معادله (۲۴) برآورد کمترین مربعات (OLS) بصورت زیر است

$$\hat{\beta}_{[\frac{n}{4}]} = \left(\sum_{t=1}^n (x_t - \bar{x})(x_t - \bar{x})' \right)^{-1} \sum_{t=1}^n (x_t - \bar{x})(y_t - \bar{y})' \quad (29)$$

وقتی $m \leq \frac{n}{4}$ است دو حالت برای رفتار مجانبی توسط ماریناسی و رابینسون (۲۰۰۱) برای این برآورد بصورت زیر در نظر گرفته شده است

$$m \sim Cn \quad 0 \leq C \leq \frac{1}{4} \quad (30)$$

و

$$\frac{1}{m} + \frac{m}{n} \rightarrow 0 \quad (31)$$

در حالت (۲۷) یک دسته از فرکانسهای غیر قابل تولید استفاده می‌شود. این حالت $\hat{\beta}_m$ بوسلیه هانن^{۲۹} (۱۹۶۳) معرفی شد و سپس بوسلیه رابینسون (۱۹۷۳) و انگل (۱۹۷۴) بررسی شد آنها این روش را «رگرسیون طیفی پنجره‌ای^{۳۰}» نامیدند. در حالت (۲۸) یکسری از فرکانسهای فوریه افزایشی دوباره استفاده می‌شوند. اما برآورد β فرکانسهای قابل تولید در اطراف فرکانس صفر انجام می‌شود.

برآورد OLS (۲۶) برای تحلیل سری‌های اتو رگرسیو (AR)، $I(1)$ و $I(2)$ که با هم هم‌انباشته باشند کاربرد فراوانی دارد. در ابتدا این برآورد، برآوردی مطلوب برای مبتکران هم‌انباشتگی [انگل و گرنجر (۱۹۸۷)] بود و سپس بعنوان یک برآورد اولیه برای محاسبه باقیمانده‌ها که برای برآوردهای کاراتری از β بکار می‌رفتند استفاده شد [مثلاً فیلیپس (۱۹۹۱)، فیلیپس و هانسن (۱۹۹۰)]. یک خصوصیت مطلوب برآورد FDLS این است که فرضهای متعامد بودن x_t ها با e_t که در برآوردهای رگرسیونی ضروری است در اینجا نیازی نیست با این وجود باز هم $\hat{\beta}_m$ سازگار است چونکه x_t بطور مجانبی بر e_t مشرف است. به هر حال عدم وجود تعامد بین x_t و e_t در برآورد کمترین مربعات و وقتی که x_t ایستا ($d \leq \frac{1}{4}$) است باعث بوجود آمدن آریبی در معادلات

27) Frequency Domain Least Square 28) Shumway 29) Hanan 30) Band spectrum regression

همزمان می‌شود. این باعث ایجاد انگیزه در رابینسون (a ۱۹۹۴) شد تا حالت (۲۸) را مورد بررسی قرار دهد. او نشان داد علاوه وجود هم بستگی بین x_t و e_t باز هم $\hat{\beta}_m$ برای β سازگار است. البته وقتی بررسی اطراف فرکانس صفر باشد و چگالی طیفی x_t بر e_t مشرف باشد. (این خصوصیت برای حالت (۲۷) وجود ندارد) بنابراین در سری‌های ایستا برآورد (۲۸) خیلی بهتر از برآورد OLS (۲۶) است. در شرایط کاربردی ممکن است از ایستایی یا نایستایی مطمئن نباشیم، مخصوصاً در حالتی که سریها انباشته کسری باشند تشخیص ایستایی و نایستایی خیلی است. به همین دلیل رابینسون و ماریناسی (۱۹۹۷، ۱۹۹۹، ۲۰۰۰) $\hat{\beta}_m$ در حالت نایستایی محض بررسی کردند. دانشمندان دیگری نیز بررسی این برآورد را آغاز کردند. چونکه رابطه هم‌انباشتگی که در بالا توضیح دادیم یک پدیده در فرکانسهای پائین است سهیم کردن فرکانسهای بالا در این برآورد ضرورتی ندارد. در واقع تکنیکهای اصلی که در این زمینه وجود دارند از وارد کردن فرکانسهای خیلی بزرگتر از صفر خودداری می‌کنند. بنابراین فرایندهای ورودی معمولاً باید $I(0)$ باشند. خصوصیات $\hat{\beta}_m$ وابسته به مرتبه‌های انباشتگی d_1, d_2, \dots, d_p و d_e هم است و این خصوصیات در هر کدام از نواحی d_1, d_2, \dots, d_p و d_e درای حالت متفاوتی هستند. برای اینکه توضیحات ما ساده‌تر شود دیگر به جزئیات و شرایط خاص نمی‌پردازیم و فرض می‌کنیم که مقادیر d برای تمام سریها مشترک باشند. در واقع توضیحات رابینسون و ماریناسی (a ۱۹۹۹)، که کار اصلی آنها نیز روی دوره نگار میانگین‌گیری شده بود، فقط برای حالت $p = 2$ بررسی شد. ما مؤلفه $\hat{\beta}_m$ را با $\hat{\beta}_{im}$ نشان می‌دهیم و فرض می‌کنیم که (۲۸) همچنان برقرار باشد. ما در اینجا فقط نرخ همگرایی را گزارش می‌کنیم که در آن $X_n \sim h(n)$ نشان دهنده این است که $x_n h(n)^{-1}$ در توزیع همگرا به توزیع یک متغیر تصادفی غیر قابل تولید است، که بیان می‌کند $X_n = O_p(h(n))$ است. حال خصوصیات مجانبی $\hat{\beta}_{im}$ که توسط ماریناسی و رابینسون (۲۰۰۱) بیان شده‌اند را در زیر می‌آوریم:

الف) نایستایی کمتر از ریشه واحد باشد $d_i > \frac{1}{p}, d_e \geq 0, d_i + d_e < 1$

$$\hat{\beta}_{i[\frac{n}{p}]} - \hat{\beta}_i \sim_d n^{1-d_i-d_{min}}, \quad \hat{\beta}_{im} - \hat{\beta}_i \sim_d n^{d_e-d_i} m^{1-d_{min}-d_e} \quad (32)$$

ب) حالت باندی^{۳۱} $\frac{1}{p} \leq d_i = 1 - d_e \leq 1$

$$\hat{\beta}_{i[\frac{n}{p}]} - \hat{\beta}_i \sim_d n^{2d_e-1} \log n \quad \hat{\beta}_{im} - \hat{\beta}_i \sim_d n^{2d_e-1} \log m \quad (34)$$

ج) حالت $d_1 = d_2 = \dots = d_p, I(0)/I(1)$

$$\hat{\beta}_{i[\frac{n}{p}]} - \hat{\beta}_i \sim_d n^{-1} \quad \hat{\beta}_{im} - \hat{\beta}_i \sim_d n^{-1} \quad (35)$$

31) Boundary case

$$d_1 = d_2 = \dots = d_p > \frac{1}{p} \text{ باشد واحد ریشه بزرگتر از ریشه واحد باشد}$$

$$d_e > 0 \quad d_i + d_e > 1$$

$$\hat{\beta}_{i[\frac{p}{2}]} - \hat{\beta}_i \sim_d n^{d_e - d_i} \quad \hat{\beta}_{im} - \hat{\beta}_i \sim_d n^{d_e - d_i} \quad (36)$$

در اینجا توزیع حدی متغیرهای تصادفی غیر استاندارد است، جزئیات بیشتر در این مورد را در رابینسون و ماریناسی (۱۹۹۹) ببینید، ماریناسی و رابینسون (۲۰۰۰) نیز تابع حدی آنها استفاده کرده‌اند. در اینجا حالت سوم، حالت عمومی $I(1)/I(0)$ است که در ممت‌های مربوط به هم‌انباشتگی بر اساس مدل AR وجود دارد. مطالب مربوط به FDLs نیز نشان می‌دهند که یکسری نقص‌هایی مانند اریبی درجه دوم مخصوصاً نسبت به حالت OLS در آنها وجود دارد. جزئیات در این مورد نیز توسط رابینسون و ماریناسی (۱۹۹۹) آمده است. معادله (۳۰) در اینجا نشان می‌دهد که نرخ همگرایی FDLs یا عبارتی توزیع حدی آن بهتر از روش OLS است. در مورد (۲۹) برتری روش FDLs واضح است. حالت (۳۲) نشان می‌دهد وقتی که از یک مقدار بزرگ m برای وارد کردن فرکانسها استفاده کرده‌ایم حذف فرکانسهای بالا تغییری در سرعت همگرایی برآوردها بوجود نمی‌آورد. نقص فرکانسهای بالا نه تنها نقصی ایجاد نمی‌کند حتی باعث بهبود برآوردها نیز می‌شود به همین دلیل است که حذف فرکانسهای بالا باعث ایجاد اریبی در معادلات همزمان می‌شود، به عبارتی دیگر در واریانس فقط فرکانسهای پائین شرکت می‌کنند. ما در اینجا حالتی که مؤلفه‌های قطعی مانند جمله ثابت، روند خطی و ... در مدل وجود داشته باشند را بررسی نمی‌کنیم. ولی رابینسون و ماریناسی (۲۰۰۰) این حالت را بررسی کرده‌اند. اگر مؤلفه‌های قطعی تحت تأثیر مؤلفه‌های تصادفی باشند نتایج حالات (۳۲-۲۸) برای این حالت نیز برقرار است.

۵ بررسی نرخ تورم در ایران

در این قسمت برای نشان دادن کاربرد رابطه هم‌انباشتگی رابطه بین شاخص قیمت کالاها CPI، نرخ ارز غیر رسمی ER، تولید ناخالص داخلی GDP و حجم پول M در ایران را بررسی کرده‌ایم. داده‌های مورد استفاده در این مقاله مربوط به لگاریتم سری‌های نام برده از سال ۱۳۸۱ - ۱۳۵۰ در ایران است.

برای بررسی ایستائی سریها ابتدا از آزمونهای استاندارد دیکي فولر و دیکي فولر تعمیم یافته که در کارهای عملی بسیار از آنها استفاده می‌شود را بکار برده‌ایم. مقادیر آماره آزمون برای آزمون دیکي فولر تعمیم یافته در جدول (۱) آورده‌ایم. هیچکدام از این مقادیر بزرگتر مقادیر آزمون ADF نیستند تا بتوانیم فرض صفر نایستائی را رد کنیم. بنابراین هر چهار سری نایستا از درجه یک هستند. پس از اینکه متوجه شدیم هر چهار سری نایستا از درجه مشابه یک هستند می‌توانیم آزمون هم‌انباشتگی به روش انگل - گرنجر و جوهانسون^{۳۲} (۱۹۹۱) را روی این سری‌ها انجام دهیم.

32) Juohonsen

جدول ۱: مقادیر آماره آزمون ADF برای لگاریتم سری‌های شاخص قیمت کالاها CPI، نرخ ارز غیر رسمی ER، تولید ناخالص داخلی GDP و حجم پول M

| سری زمانی | آماره آزمون ADF | مقادیر بحرانی |
|-----------|-----------------|---------------|
| Lcpi | -۲,۱۹ | -۳,۵۷ |
| LM | -۳,۱۴ | -۳,۵۷ |
| LER | -۲,۷ | -۲,۹ |
| LGDP | -۱,۷ | -۳,۵۷ |

در روش انگل - گرنجر آزمون ADF را روی باقیمانده‌های حاصل از معادله رگرسیون زیر انجام می‌دهیم

$$LCPI = ۲,۶ + ۰,۱۶۴LER - ۰,۰۸LGDP + ۰,۶۸LM + \varepsilon_t \quad (۳۷)$$

آماره ADF برای باقیمانده‌ها $-۲,۴$ و مقدار بحرانی آن $-۲,۹$ بدست آمد که نشان دهنده ایستا بودن باقیمانده‌ها است، بنابراین رابطه هم‌انباشتگی بین چهار سری وجود دارد. آزمون جوهانسون نیز رابطه هم‌انباشتگی را تأیید کرد که بدلیل زیاد بودن مطالب از آوردن آن خوداری می‌کنیم. برای بررسی رابطه هم‌انباشتگی کسری ابتدا درجه انباشتگی هر کدام از سریها را برآورد کرده‌ایم که نتایج آن در جدول (۲) نشان داده شده است. مقادیر نشان داده شده در ستون دوم جدول

جدول ۲: برآورد درجه انباشتگی سری‌های زمانی مورد بررسی. δ درجه انباشتگی سری‌های تفاضل‌گیری شده است و $\delta + ۱$ درجه انباشتگی سری‌های اصلی است

| سری زمانی | δ | $\delta + ۱$ |
|-----------|----------|--------------|
| Lcpi | ۰,۳۹۵ | ۱,۳۰۵ |
| LM | ۰,۳۲۸ | ۱,۳۲۸ |
| LER | ۰,۰۳۹۳ | ۱,۰۳۹۳ |
| LGDP | ۰,۳۲۸ | ۱,۳۲۸ |

(۲) مربوط درجه انباشتگی سری‌های تفاضل‌گیری شده است که آن را با δ نشان داده‌ایم و مقادیر نشان داده شده در ستون سوم جدول درجه انباشتگی سری‌های اصلی است در واقع به δ عدد یک اضافه کرده‌ایم.

پس از برآورد درجه انباشتگی سری‌ها و انجام آزمون فرض معلوم شد که سه متغیر GDP، M و CPI دارای درجه انباشتگی مشابه هستند ولی درجه انباشتگی ER با بقیه متفاوت است. بنابراین رابطه رگرسیونی را بین متغیرهای زیر برآورد کرده‌ایم

$$LCPI = -۲,۶ - ۰,۳LGDP + ۰,۸LM + \varepsilon_t \quad (۳۸)$$

حال برای اینکه رابطه هم‌انباشتگی بین این متغیرها وجود داشته باشد درجه انباشتگی باقیمانده‌های مدل بالا باید کوچکتر از درجه انباشتگی متغیرها باشد. پس از برآورد درجه انباشتگی باقیمانده‌ها مقدار آن ۰,۳۷ بدست آمد که رابطه هم‌انباشتگی کسری بین این سه متغیر وجود دارد.

۶ نتایج

پس از تعاریف رابطه هم‌انباشتگی، انباشتگی کسری و رابطه هم‌انباشتگی کسری نشان دادیم که تعداد بردارهای هم‌انباشتگی با توجه به تعریف آنها می‌تواند متفاوت باشد.

کار اصلی در بررسی رابطه هم‌انباشتگی کسری برآورد درجه انباشتگی هر کدام از سری‌ها است. ما در اینجا نشان دادیم که برآوردهای نیمه پارامتری در این زمینه کارا تر از برآوردهای پارامتری هستند. پس از برآورد درجه انباشتگی سری‌ها برآوردهای هم‌انباشتگی اهمیت پیدا می‌کنند. ما همچنین نشان دادیم که یک روش کارا برای برآورد بردارهای هم‌انباشتگی روش FDLs است که به مراتب کارا تر از روش OLS است.

ما کاربرد رابطه هم‌انباشتگی و هم‌انباشتگی کسری را برای بررسی رابطه بین لگاریتم سریهای شاخص قیمت کالاها LCPI، نرخ ارز غیر رسمی LER، تولید ناخالص داخلی GDP و حجم پول LM در ایران را بررسی کردیم. در رابطه هم‌انباشتگی طبق تعریف انگل - گرنجر و همچنین تعریف جوهانسون معلوم شد یک رابطه هم‌انباشتگی بین این سری‌ها وجود دارد که آن را در معادله (۳۷) نشان دادیم. این معادله نشان می‌دهد نرخ تورم در بلند مدت با حجم پول و نرخ ارز غیر رسمی رابطه مستقیم و با تولید ناخالص داخلی رابطه عکس دارد. بنابراین دولت ایران می‌تواند با کاهش حجم پول در دست مردم و کاهش واردات و ارز و همچنین با بالا بردن تولید ناخالص داخلی یعنی تولید کالاهای داخلی و کاهش صادرات نفت نرخ تورم را کاهش دهد.

مراجع

- [1] Chan, N.H , trrin , N.(1995), Infrence for unstable log-memory processes whit applications to fractional unit root autoregressions. Anals of Statistics 23 ,1662-1683 .

- [2] Chueng ,Y.W.,Lai K,S, 1993 A fractional cointegration anlysis of purchasing power parity. *Journal of Bussines and Economic Statistics* 11 ,103-112.
- [3] Dicky , D.A.,Fuller,W.A, 1981 Likelihood ratio statistics for autoregressive time series whit unit root .*Econometrica* 49, 1057-1072.
- [4] Engle ,R.F Granger ,C.W.G. ,1987 Cointegration and error correction : representation , estimation and testing . *Econometrica* 57 , 251-276.
- [5] Engle ,R.F Granger ,C.W.G. ,1991 Long-Run Economic Relationship , Oxford University Press , Oxford.
- [6] Flores Jr., Szafarz,A., 1996 An enlarge definition of cointegration . *Economic Letters* 50 , 193-195.
- [7] Fox, R.,Taqqu ,M.S., 1986 Large-sample properties of parameter estimate for strongly dependent stationary Gayssian time series .*Annals of Statistics* 14 ,517-532.
- [8] Geweke ,J., Poter-Hudak ,S., 1983 ,The estimstion snd application of long memory time series models .*Journal of Time Series Analysis* 4 ,221-238.
- [9] Hannan E.J., 1963 Regression of time series whit errors of measurement . *Biometrika* 50 ,293-302 .
- [10] Hausman ,J.1978 Misspecification tests in econometrics . *Econometrica* 46 ,1251-1271.
- [11] Hurvich ,C.M ,Deo ,R.,Brodsky ,J.,1998 The mean squer error of Geweck and Poter-hudak's estimate of the memory parameter of long memory time series .*Journal of Time Series Analysis* 19 ,19-46.
- [12] Johansen,S., 1991, Estimation and hypothesis testing of cointegrating vectors in Gaussian vectors in autoregressive models . *Econometrica* 59,1551-1580.
- [13] Johansen,S., 1996,Likelihood-Based infrence in conintegrated vector autoregressive models. 2nd Printing .Oxford University Press ,Oxford. Kunsch ,H., 1987 ,Statistical of aspects of self-similar processes . INn : Prohorov ,Yu,Sazonov ,V.V(Eds.), Proceeding of the first world Congress of the Bernolli Society 1.VNU Science Press Utrechh, pp . 67-74.
- [14] Lobato ,I 1995 Multivariate analysis of long memory seies inthe frequency domain . PH.D Thesis , London School of Economics.
- [15] Lobato ,I., 1999 A semiparametric tow-step estimator in the multivariate long memory models . *Journal of Econometrics* 90 129-153.

- [16] Marinucci ,D., Robinson,P.M., 1999a ,Alternative forms of fractional brownian motion . Journal of Statistical Planning and Inference 90 ,129-153.
- [17] Marinucci ,D., Robinson,P.M., 1999b ,Finite sample improvements in statistical inference whit I(1) processes . Journal of Econometrics , forthcoming.
- [18] Marinucci ,D., Robinson,P.M., 2000 ,Weak convergence of multivariate fractional processes. Stochastic processes and their application 86 ,103-120.
- [19] Marinucci ,D., Robinson,P.M., 2001 ,Semiparametric fractional cointegration analysis. Journal of Econometrics 105 ,225-247
- [20] Philips ,P.C.B.,Ouliaris ,S., 1990 ,Asymptotic properties of residual based cointegration test.Econometrica 58 ,165-193 .
- [21] Robinson ,P.M., 1973,The stochastic difference equation whit non-integral differences . Advance In Applied Probability 6 ,524-545
- [22] Robinson ,P.M.,1988 The stochastic difference between econometric statistics . Econometrica 56 ,531-148
- [23] Robinson ,P.M.1994a ,Semiparametric analysis of long memory time series .Annals of Statistics 22 ,519-539
- [24] Robinson ,P.M.1994b , Efficient test of nonstationary hypothesis .Journal of the American Statistical Association 89, 1420-1437
- [25] Robinson ,P.M.1995a ,Log-periodogram regression of time series whit long range dependence.Annals of Statistics 23,1048-1072
- [26] Robinson ,P.M.1995b ,Gaussian semiparametric estimation of long range dependence Annals of Statistics 23,1630-1661
- [27] Robinson ,P.M.1998 ,comment on real and spurious long memory properties of stock market data ,by I.G.Lobato and N.E Savin .Journal of Business and Economic Statistics 16 ,261-182
- [28] Robinson ,P.M., Marinucci ,D.1997 ,Semiparametric frequency-domain analysis of fractional cointegration ,preprint.
- [29] Robinson ,P.M., Marinucci ,D.1999 , Sem Narrow-band analysis of nonstationary processes ,preprint
- [30] Robinson ,P.M., Marinucci ,D.2000, The averaged periodogram of nonstationary vector time series . Statistical Inference for Stochastic Processes 3 ,149-160
- [31] Robinson ,P.M.,Yajima . 2002 , Determination of cointegration rank in fractional systems . Journal of Econometrics 106,217-241

معیارهای انتخاب متغیر در داده‌های چند متغیره بر اساس ساختار کوواریانس و همبستگی

نادر نعمت‌الهی^۱، قدرت‌اله رحمتی^۲

^۱ گروه آمار دانشگاه علامه طباطبایی
^۲ کارشناسی ارشد آمار دانشگاه علامه طباطبایی

چکیده: مؤلفه‌های اصلی یک ترکیب خطی از متغیرهای اولیه می‌باشند. در عمل هنگامیکه تحلیل مؤلفه‌های اصلی روی تعداد زیادی از متغیرها به کار گرفته می‌شود، ممکن است نتوان مؤلفه‌های اصلی به دست آمده را به آسانی تفسیر کرد. در این حالت می‌توان با انتخاب یک زیر مجموعه از متغیرها که تقریباً همان اطلاعات کلی را داشته باشند، تفسیر مؤلفه‌های اصلی حاصل را راحتتر انجام داد. در این مقاله معیارهای مختلف انتخاب متغیر برای ساختار کوواریانس و همبستگی معرفی و بررسی شده است. چون هر معیار یک زیر مجموعه از متغیرها را انتخاب می‌کند، در نتیجه با استفاده از معیار کارایی توضیح داده شده، مناسب‌ترین زیر مجموعه انتخاب می‌شود.

واژه‌های کلیدی: مؤلفه‌های اصلی، معیارهای انتخاب متغیر، ساختار همبستگی، اندازه کارایی، تحلیل خوشه‌ای

۱ مقدمه

در بسیاری از مباحث علمی از قبیل پزشکی، مدیریت، جامعه‌شناسی و اقتصاد مطالعات بر روی تعداد زیادی از متغیرهای تصادفی انجام می‌پذیرد که همبستگی متقابلی بین این متغیرها وجود دارد. به دلیل همین همبستگی چنین به نظر می‌رسد که در اینگونه مسایل بایستی بعد مساله را با از دست دادن حداقل اطلاعات کاهش داد. یکی از روش‌هایی که برای کاهش بعد وجود دارد روش تجزیه و تحلیل مؤلفه‌های اصلی^۱ می‌باشد که هتالینگ^۲ ارایه داده است. مؤلفه‌های اصلی Y_i ، $i = 1, 2, \dots, p$ یک ترکیب خطی از متغیرهای اولیه X_1, X_2, \dots, X_p به صورت زیر می‌باشند:

$$Y_i = \sum_{j=1}^p e_{ji} X_j, \quad i = 1, 2, \dots, p$$

1) Principal Components 2) Hotelling

که معمولاً بوسیله بارهای ezj مؤلفه‌های اصلی تعبیر و تفسیر می‌شوند. حال اگر تعداد متغیرهای اولیه یعنی p بزرگ باشد آنگاه تفسیر مؤلفه‌های اصلی آسان نیست. برای رفع این مشکل انتخاب یک زیر مجموعه از متغیرهای اولیه که تقریباً همان اطلاعات کلی را داشته باشند راه حل عملی می‌باشد. در اینجا یک سوال پیش می‌آید که کدامیک از متغیرها ضروری هستند و در صورت ضروری نبودن بعضی از متغیرها آیا می‌توان آنها را حذف کرد؟ یک راه حل برای این سوال پیدا کردن معیارهایی برای انتخاب زیر مجموعه‌ای از متغیرها است که اولین بار توسط بیل^۳ و دیگران (۱۹۷۶) بررسی شده است. این معیارها به طور کلی به چهار دسته تقسیم می‌شوند.

- ۱) معیارهای که بر اساس همبستگی چندگانه بدست می‌آیند.
- ۲) معیارهایی که از مؤلفه‌های اصلی به دست می‌آیند،
- ۳) معیارهایی که از متغیرهای اولیه به دست می‌آیند،
- ۴) معیارهایی که از تحلیل خوشه‌ای به دست می‌آیند.

بعد از بیل و دیگران، جولیف^۴ و الکاندری^۵ در مقالات مختلف در سال‌های ۱۹۷۲، ۱۹۷۳، ۲۰۰۱ و ۲۰۰۲ به بررسی معیارهای مختلف انتخاب متغیر پرداخته و مک کیپ^۶ نیز در سال (۱۹۸۴) با استفاده از متغیرهای اصلی معیارهای دیگری را نیز عرضه کرد. کادیم^۷ (۱۹۹۵) نیز به بررسی ضعف استفاده از روش معمول انتخاب متغیر بر اساس مؤلفه‌های اصلی پرداخت. در این مقاله سعی شده است بر اساس اینکه متغیرها هم واحد باشند و یا نه، معیارهای انتخاب متغیر بر اساس ساختار کوواریانس و ساختار همبستگی معرفی و بررسی شوند. بر این اساس در بخش دوم معیارهای انتخاب متغیر برای ساختار کوواریانس ارائه می‌شوند. در بخش سوم نیز معیارهای انتخاب متغیر برای ساختار همبستگی معرفی می‌شوند و کاربرد معیارهای انتخاب متغیر در یک مثال بررسی می‌شود.

۲ معیارهای انتخاب متغیر برای ساختار کوواریانس

در تحلیل چند متغیره بر حسب اینکه متغیرها از لحاظ واحد اندازه‌گیری هم واحد باشند یا نه از ساختار کوواریانس یا همبستگی استفاده می‌شود. در این بخش برای ساختار کوواریانس معیارهای انتخاب متغیر را بررسی می‌کنیم. برای این منظور ابتدا سه خاصیت مؤلفه‌های اصلی که توسط مک کیپ (۱۹۸۴) اثبات شده است را بیان می‌کنیم. اکثر معیارهایی که در این بخش و بخش بعد آمده‌اند بر اساس خواص زیر می‌باشند.

خاصیت ۱: به ازای هر عدد صحیح q, p $1 \leq q \leq p$ تبدیل زیر را در نظر بگیرید

$$\tilde{Z} = \tilde{B}'\tilde{X}$$

3) Beale 4) Jolliffe 5) AL-Kandari 6) Mc Cabe 7) Cadima

که Z یک بردار q -بعدی و B' یک ماتریس $q \times p$ می‌باشند و $B'B = I_q$. در این صورت اثر ماتریس $\sum Z$ یعنی $tr(\sum Z)$ ماکزیمم است اگر $B = P_q$ باشد، که در آن P_q شامل q ستون اول ماتریس بردارهای ویژه ماتریس واریانس - کوواریانس متغیرهای اولیه است $\sum Z = B' \sum B$ و ماتریس واریانس - کوواریانس Z است.

خاصیت ۲: با در نظر گرفتن همان ترکیب خطی متعامد $Z = B'X$ در خاصیت ۱، $tr(\sum Z)$ مینیمم است اگر $B = P_q^*$ که P_q^* شامل q ستون آخر ماتریس بردارهای ویژه ماتریس واریانس - کوواریانس متغیرهای اولیه است.

خاصیت ۳: همانند خواص ۱ و ۲، تبدیل $Z = B'X$ را در نظر بگیرید، در این صورت دترمینان ماتریس $\sum Z$ یعنی $\det(\sum Z)$ ماکزیمم است اگر $B = P_q$ باشد.

حال معیارهای انتخاب متغیر را برای ساختار کوواریانس بیان می‌کنیم.

معیار ترکیب واریانس (VC): این معیار زیر مجموعه‌ای از متغیرها را انتخاب می‌کند که دارای واریانس‌های بزرگی باشند. یعنی اگر واریانس متغیرها را به ترتیب نزولی مرتب کنیم آنگاه q متغیری که دارای q واریانس اول باشند، انتخاب می‌شوند.

معیار ترکیب بار (LC): در این معیار ابتدا بایستی تحلیل مؤلفه‌های اصلی را روی داده‌ها انجام داده و سپس k مؤلفه اصلی اول که دارای بیشترین سهم در کل واریانس باشند را در نظر می‌گیریم و بار p متغیر اولیه را در این k مؤلفه اصلی اول به طور نزولی با هم مرتب می‌کنیم. سپس آن q متغیری که (بدون تکرار) دارای بزرگترین بار باشند را نگه می‌داریم و بقیه متغیرها را حذف می‌کنیم. **معیار B_1 :** در این معیار ابتدا تحلیل مؤلفه‌های اصلی را روی p متغیر اولیه انجام داده، سپس آن تعداد از مؤلفه‌های اصلی که بیشتر از α درصد از درصد تجمعی تغییرات کل را توضیح دهند را در نظر گرفته و از هر مؤلفه مورد نظر، متغیری که دارای بیشترین بار باشد را حذف می‌کنیم. تعداد متغیرهای حذف شده برابر است با تعداد مؤلفه‌های اصلی که بیشتر از α درصد تجمعی را توضیح دهند. جولیف (۱۹۷۲) با شبیه‌سازی مقدار مناسب $\alpha = 0.8$ را پیشنهاد کرد. اما جولیف در عمل $\alpha = 0.8$ را مناسب نمی‌داند و باید با توجه به تعداد متغیرهای مورد نظر مقدار α را تعیین کرد.

معیار B_2 : در این معیار ابتدا تحلیل مؤلفه‌های اصلی را روی p متغیر اولیه انجام داده، سپس آن تعداد از مؤلفه‌های اصلی که کمتر از α درصد از درصد تجمعی تغییرات کل را توضیح دهند را در نظر می‌گیریم. در آخر از هر مؤلفه مورد نظر متغیری که دارای بیشترین بار باشد را نگه می‌داریم. برای تعیین مقدار مناسب α همانند معیار B_1 عمل می‌شود.

معیار همبستگی چندگانه (R): این معیار بوسیله یک روش گام به گام به دست می‌آید. بدین صورت که ابتدا متغیری که با $(p - 1)$ متغیر دیگر دارای همبستگی چندگانه ماکزیمم باشد را حذف می‌کنیم، زیرا به خوبی توسط $(p - 1)$ متغیر دیگر توضیح داده شده است. در گام دوم متغیری که با $(p - 2)$ متغیر دیگر دارای همبستگی چندگانه ماکزیمم باشد را حذف می‌کنیم. پس در هر گام هنگامی که m متغیر باقی می‌ماند، متغیری که دارای ماکزیمم همبستگی چندگانه با

($m - 1$) متغیر دیگر باشد را حذف می‌کنیم و این فرآیند را تا رسیدن به q متغیر ادامه می‌دهیم. شرط توقف برای این روش آن است که تمام همبستگی‌های چندگانه بین متغیرها باقیمانده کمتر از یک مقدار R_0 باشد. جولیف با استفاده از شبیه‌سازی مقدار مناسب R_0 را برابر 0.15 پیشنهاد کرده است.

معیار دسته‌بندی همبستگی (CC): در این معیار متغیرهایی که بیشترین ارتباط را با q مؤلفه اصلی اول داشته باشند را نگه می‌داریم. برای انجام این فرآیند به این صورت عمل می‌کنیم که ابتدا مقدار همبستگی بین متغیرها و q مؤلفه اصلی اول را با هم به ترتیب نزولی مرتب می‌کنیم، سپس متغیرهایی که دارای مقدار همبستگی بزرگی باشند را نگه می‌داریم.

تعمیم معیار دسته‌بندی بارها (NLC): در این معیار میانگین قدر مطلق بار متغیرها را در q مؤلفه اصلی اول در نظر می‌گیریم و متغیرهایی که دارای بیشترین میانگین قدر مطلق بار باشند را نگه می‌داریم.

تعمیم معیار دسته‌بندی همبستگی (NCC): در این معیار آن q متغیر را نگه می‌داریم که دارای بیشترین میانگین قدر مطلق مقادیر همبستگی با مؤلفه‌های اصلی مورد نظر باشند. **معیار M_1 :** در این معیار آن q متغیری نگه داشته می‌شوند که دارای کوچکترین دترمینان ماتریس کوواریانس شرطی متغیرهای حذف شده به شرط متغیرهای نگه داشته شده باشند. یعنی q متغیری که کوچکترین دترمینان ماتریس $S_{22}^{-1} S_{12} S_{22}^{-1} S_{12}$ را به دست می‌دهند، که در آن S_{22} ماتریس کوواریانس متغیرهای حذف شده و $S_{22}^{-1} S_{12}$ ماتریس کوواریانس متغیرهای نگه داشته شده می‌باشند.

معیار M_2 : در این معیار آن q متغیری نگه داشته می‌شوند که دارای کوچکترین اثر ماتریس کوواریانس شرطی متغیرهای حذف شده به شرط متغیرهای نگه داشته شده باشند.

معیار M_3 : در این معیار آن q متغیری نگه داشته می‌شوند که دارای کوچکترین نرم ماتریس کوواریانس شرطی متغیرهای حذف شده به شرط متغیرهای نگه داشته شده باشند.

معیار تحلیل خوشه‌ای (C): این معیار بر اساس متناسب کردن p متغیر اولیه به q خوشه از متغیرها عمل می‌کند. در این روش متغیرها بر اساس تحلیل خوشه‌ای سلسله مراتبی به q خوشه تقسیم می‌شوند و سپس از هر خوشه یک متغیر را انتخاب می‌کنیم. انتخاب متغیر از هر خوشه را می‌توان به سه روش زیر انجام داد.

- ۱- آخرین متغیر انتخاب شده در هر خوشه،
- ۲- انتخاب یک متغیر از متغیرهای اولیه یا متغیر داخلی هر خوشه،
- ۳- انتخاب تصادفی یک متغیر از هر خوشه.

الگوریتم تحلیل خوشه‌ای سلسله مراتبی به صورت زیر است.

۱- تعریف یک اندازه تشابه r_{XY} ، بین هر دو گروه از متغیرهای X و Y (در مرحله اول هر متغیر یک گروه محسوب می‌شود)،

۲- محاسبه r_{XY} برای تمام زوج متغیرهای گروهها که تعداد آنها در مرحله اول $(p - 1)p$

می باشد،

۳- ترکیب دو گروه از متغیرها که دارای ماکزیمم r_{XY} باشند،

۴- محاسبه r_{XY} برای گروههای جدید و برگشت به گام سوم.

این فرآیند تا هنگامی که تمام r_{XY} بین گروهها کمتر از یک مقدار r باشند ادامه پیدا می کند. روشهای متفاوتی برای در نظر گرفتن اندازه تشابه r_{XY} وجود دارد. در این مقاله از دو روش اتصال متوسط^۸ و اتصال تکی^۹ استفاده می کنیم که در آن r_{XY} و r'_{XY} به صورت زیر به دست می آیند

$$r_{XY} = \frac{\sum_{i \in X} \sum_{j \in Y} r_{ij}}{n_1 n_2}$$

$$r'_{XY} = \max_{i \in X, j \in Y} r_{ij}$$

که n_1 و n_2 تعداد متغیرهای موجود در هر گروه X و Y ، r_{ij} ضریب همبستگی بین متغیر i ام و متغیر j ام می باشند. جولیف (۱۹۷۲) با استفاده از شبیه سازی مقدار r را برای روشهای اتصال متوسط و اتصال تکی به ترتیب برابر 0.45 و 0.55 در نظر گرفته است. اندازه کارایی: بعد از اینکه هر معیار انتخاب متغیر یک زیر مجموعه از متغیرها را معرفی نمود، بایستی یکی از این زیر مجموعه ها را انتخاب کنیم. در اینجا یک سوال به وجود می آید که کدام یک از این زیر مجموعه ها مناسب است؟ برای پاسخ به این سوال، نیاز به معرفی یک اندازه کارایی داریم که نزدیک بودن نتایج مؤلفه های اصلی روی مجموعه کامل و زیر مجموعه ای از آنها را بسنجد. در این مقاله دو اندازه کارایی را در نظر می گیریم. اندازه کارایی اول عبارت است از

$$E_1 = \frac{\sum_{j=1}^k \lambda_j^n}{tr(\mathbf{S})} 100$$

$$E_2 = \frac{\sum_{j=1}^k \lambda_j}{tr(\mathbf{S})} 100$$

که در آن λ_j ، λ_j^n و \mathbf{S} به ترتیب واریانس زامین مؤلفه اصلی اولیه، واریانس زامین مؤلفه اصلی جدید و ماتریس کوواریانس نمونه می باشند که برای k مؤلفه اصلی اول مقایسه می شوند. برای به دست آوردن اندازه کارایی دوم، ابتدا درصدی از تغییرات کل که توسط زیر مجموعه ای از متغیرها توضیح داده می شود را در نظر می گیریم، که به صورت زیر است

$$\frac{\sum_{i=1}^q S_{i,r}^2 + \sum_{h=1}^{p-q} (R_{h,r}^2 \times S_{h,d}^2)}{tr(\mathbf{S})}$$

که در آن $S_{i,r}^{\gamma}$ واریانس نمونه i امین متغیر در متغیرهای نگه داشته شده، $S_{h,d}^{\gamma}$ واریانس نمونه h امین متغیر در متغیرهای حذف شده و $R_{h,r}^{\gamma}$ نیز مربع همبستگی چندگانه h امین متغیر حذف شده با q متغیر نگه داشته شده می‌باشند. سپس رابطه بالا را با رابطه زیر مقایسه می‌کنیم

$$\frac{\sum_{i=1}^p (S_i^{\gamma} \times R_{i,y_1^n, y_2^n, \dots, y_k^n}^{\gamma})}{tr(\mathbf{S})}$$

که در آن مربع همبستگی چندگانه i امین متغیر با k مؤلفه اصلی اولیه می‌باشد. یعنی

$$E_{\gamma} = \frac{\frac{\sum_{i=1}^q S_{i,r}^{\gamma} + \sum_{h=1}^{p-q} (R_{h,r}^{\gamma} \times S_{h,d}^{\gamma})}{tr(\mathbf{S})} \cdot 100}{\frac{\sum_{i=1}^p (S_i^{\gamma} \times R_{i,y_1^n, y_2^n, \dots, y_k^n}^{\gamma})}{tr(\mathbf{S})} \cdot 100}$$

۳ معیارهای انتخاب متغیر مربوط به ساختار همبستگی

هنگامیکه داده‌های چند متغیره هم واحد نباشند، تحلیل مؤلفه اصلی بر اساس ماتریس همبستگی که بر اساس متغیرهای استاندارد شده $\tilde{x}_i = \frac{x_i - \bar{x}}{s}$ بدست می‌آید، انجام می‌شود. در این حالت برای انتخاب یک زیر مجموعه از متغیرها بایستی معیارهای جدیدی تعریف شوند. البته بعضی از معیارها انتخاب متغیر از قبیل معیار همبستگی چندگانه (R) ، معیارهای M_1 ، M_2 ، M_3 و معیار تحلیل خوشه‌ای (C) برای ساختار کوواریانس و همبستگی مشترک می‌باشند. در ادامه معیارهایی که فقط مربوط به ساختار همبستگی هستند ارائه می‌شود.

معیار B_1^* : در این معیار، ابتدا تحلیل مؤلفه‌های اصلی را روی تمام p متغیر اولیه انجام داده و سپس آن تعداد از مؤلفه‌های اصلی که مقادیر ویژه آنها کمتر از λ_0 باشد را در نظر گرفته و از هر مؤلفه مورد نظر، متغیری که دارای بیشترین بار باشد را حذف می‌کنیم. تعداد متغیرهای حذف شده برابر تعداد مقادیر ویژه کوچکتر از λ_0 می‌باشد. جولیف (۱۹۷۲) با شبیه‌سازی مقدار مناسب $\lambda_0 = 0.7$ را پیشنهاد کرده است.

معیار B_2^* : در این معیار نیز ابتدا تحلیل مؤلفه‌های اصلی را روی تمام p متغیر اولیه انجام داده و سپس آن تعداد از مؤلفه‌های اصلی که مقادیر ویژه آنها بزرگتر از λ_0 باشد را در نظر گرفته و از هر مؤلفه مورد نظر، متغیری که دارای بیشترین بار باشد را نگه می‌داریم. تعداد متغیرهای نگه داشته شده برابر تعداد مقادیر ویژه بزرگتر از λ_0 می‌باشد. جولیف (۱۹۷۲) با شبیه‌سازی مقدار مناسب $\lambda_0 = 0.7$ را پیشنهاد کرده است.

معیار M_1^* : در این معیار آن q متغیری را نگه می‌داریم که دارای کوچکترین دترمینان ماتریس

همبستگی شرطی متغیرهای حذف شده به شرط متغیرهای نگه داشته شده باشند. یعنی q متغیری که کوچکترین دترمینان ماتریس $R_{22}^{-1} R_{21} R_{11}^{-1} R_{21}$ را به دست می‌دهند، که در آن ماتریس همبستگی متغیرهای حذف شده و R_{11} ماتریس همبستگی متغیرهای نگه داشته شده می‌باشند.

معیار M_2^* : در این معیار آن q متغیری را نگه می‌داریم که دارای کوچکترین اثر ماتریس همبستگی شرطی متغیرهای حذف شده به شرط متغیرهای نگه داشته شده باشند.

معیار M_3^* : در این معیار آن q متغیری را نگه می‌داریم که دارای کوچکترین نرم ماتریس همبستگی شرطی متغیرهای حذف شده به شرط متغیرهای نگه داشته شده باشند.

اندازه کارایی: در این حالت برای انتخاب زیر مجموعه مناسبی از متغیرها از اندازه کارایی زیر استفاده می‌کنیم.

$$E_1^* = \frac{\frac{\sum_{j=1}^k \lambda_j^n}{100}}{\frac{\sum_{j=1}^k \lambda_j}{p} 100}$$

که این رابطه برای k مؤلفه اصلی اول محاسبه می‌شود. کاربرد مباحث فوق را روی داده‌های مثال زیر که از کتاب «آشنایی با روش‌های آماری چند متغیره» تالیف مانلی^{۱۰}، ترجمه دکتر محمد مقدم و دیگران (۱۳۷۳) گرفته شده است، به کار می‌بریم.

مثال: مقدار پروتئین موجود در ۹ گروه غذایی برای ۲۵ کشور اروپایی اندازه گرفته شده، که این ۹ گروه غذایی عبارتند از گوشت قرمز (X_1)، گوشت سفید (X_2)، تخم مرغ (X_3)، شیر (X_4)، ماهی (X_5)، غلات (X_6)، غذاهای نشاسته‌ای (X_7)، آجیل و دانه‌های روغنی (X_8) و میوه‌جات و سبزیجات (X_9). نتایج مربوط به استفاده از معیارهای انتخاب متغیر برای انتخاب زیر مجموعه‌ای از متغیرها بر اساس ساختار کوواریانس در جدول ۱ ارائه شده است. با توجه به جدول ۱، بر اساس اندازه کارایی E_1 و E_2 ، مناسب‌ترین زیر مجموعه عبارت است از

$$X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$$

نتایج مربوط به استفاده از معیارهای انتخاب متغیر برای انتخاب زیر مجموعه‌ای از متغیرها بر اساس ساختار همبستگی در جدول ۲ ارائه شده است. با توجه به اندازه کارایی جدول ۲، مناسب‌ترین زیر مجموعه عبارت است از

$$X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$$

توجه کنید که برای تعیین زیر مجموعه‌ای از متغیرها توسط معیارهای M_1 ، M_2 و M_3 (و یا معادلاً M_1^* ، M_2^* و M_3^*) بایستی تمام افزای‌های ممکن ماتریس واریانس کوواریانس (یا همبستگی)

جدول ۱: معیارهای انتخاب متغیر و زیر مجموعه های انتخاب شده بر اساس ساختار کوواریانس

| معیار انتخاب متغیر | متغیرهای انتخاب شده | E_1 | E_2 |
|--------------------|--|-------|-------|
| VC | X_2, X_4, X_6 | ۰٫۹۶ | ۰٫۹۴ |
| LC | X_1, X_2, X_4, X_6 | ۰٫۹۷ | ۰٫۹۹ |
| B_1 | X_1, X_2, X_4, X_6 | ۰٫۹۷ | ۰٫۹۹ |
| B_2 | X_1, X_2, X_4, X_6 | ۰٫۹۷ | ۰٫۹۹ |
| R | X_1, X_2, X_5, X_6, X_9 | ۰٫۷۲ | ۰٫۸۹ |
| CC | X_1, X_2, X_4, X_6 | ۰٫۹۷ | ۰٫۹۹ |
| NLC | X_1, X_2, X_4, X_6 | ۰٫۹۷ | ۰٫۹۹ |
| NLC | X_1, X_2, X_4, X_8 | ۰٫۴۱ | ۰٫۵۷ |
| M_1 | $X_1, X_2, X_4, X_5, X_6, X_7, X_9$ | ۰٫۹۸ | ۱٫۰۵ |
| M_2 | $X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$ | ۱ | ۱٫۰۶ |
| M_3 | $X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$ | ۱ | ۱٫۰۶ |
| C | $X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$ | ۱ | ۱٫۰۶ |

جدول ۲: معیارهای انتخاب متغیر و زیر مجموعه های انتخاب شده بر اساس ساختار همبستگی

| معیار انتخاب متغیر | زیرمجموعه های انتخاب شده | E_1^* |
|--------------------|--|---------|
| B_1^* | X_2, X_3, X_6, X_7, X_9 | ۰٫۶۴ |
| B_2^* | X_1, X_5, X_6, X_7, X_8 | ۰٫۶۴ |
| R | X_1, X_2, X_5, X_6, X_9 | ۰٫۶۴ |
| M_1 | $X_1, X_2, X_4, X_5, X_6, X_7, X_9$ | ۰٫۸ |
| M_2 | $X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$ | ۰٫۹۰۵ |
| M_3 | $X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$ | ۰٫۹۰۵ |
| M_1^* | X_7 | ۰٫۱۲ |
| M_2^* | $X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_9$ | ۰٫۹ |
| M_3^* | $X_1, X_2, X_4, X_5, X_6, X_7, X_8, X_9$ | ۰٫۹۰۵ |
| C | $X_1, X_2, X_5, X_6, X_7, X_8, X_9$ | ۰٫۸۱ |

نمونه را بدست آورده و سپس دترمینان، اثر و یا نرم ماتریس $S_{22,1}$ یا $(R_{22,1})$ را برای هر کدام حساب کرده و مقدار مینیمم را بدست آوریم. برای انجام این روش برنامه‌ای با استفاده از MATLAB تهیه شده است که نتایج این برنامه در جداول ۱ و ۲ مشاهده می‌گردد.

نتیجه: در تحلیل داده‌های چند متغیره، اگر تعداد متغیرها زیاد باشد آنگاه بوسیله معیارهای انتخاب متغیر گفته شده در فوق زیر مجموعه‌هایی از این متغیرها انتخاب می‌شود. سپس با استفاده از معیارهای کارایی توضیح داده شده مناسب‌ترین معیار و در نتیجه مناسب‌ترین زیر مجموعه از متغیرها جهت تحلیل مؤلفه‌های اصلی انتخاب می‌شود.

مراجع

- [۱] مانلی، آشنایی با روش‌های آماری چند متغیره، ترجمه دکتر محمد مقدم و دیگران، ۱۳۷۳، انتشارات پيشتاز علم.
- [2] AL-Kandari, N. (2002). Correlation structures and variable selection in principal component analysis. unpublished manuscript.
- [3] AL-Kandari, N. and Jolliffe, I. (2001). Variable selection and interpretation of covariance principal components. *Communication in statistics: simulation and computation*, 30(2), 339-354.
- [4] Beale, E.M.L. Kendall, M.G. and Mann, D.W. (1967). The discarding of variable in multivariate analysis. *Appl. Statist.*, 16, 223-236.
- [5] Cadima, J.F.C.L., Jolliffe, I.T. (1995). Loading and correlation in the interpretation of principal components. *Journal of applied statistics*, 22(2), 203-314.
- [6] Jolliffe, I.T. (1972). Discarding variables in a principal component analysis : artificial data. *Appl. Statist.* 21, 160-173.
- [7] Jolliffe, I.T. (1973). Discarding variables in a principal component analysis : artificial data. *Appl. Statist.* 22, 21-31.
- [8] Mc Cabe, G.P. (1984). Principal variables. *Tecnometrics*, 26(2), 137-144.

پیشگویی در نمونه‌گیری از جامعه متناهی تحت مدل‌های خطای اندازه‌گیری و کاربرد آن در برآورد هزینه و درآمد خانوار

نادر نعمت‌الهی^۱، پوریا رضاسلطانی^۲

^۱ گروه آمار دانشگاه علامه طباطبائی
^۲ سازمان مدیریت و برنامه‌ریزی کشور

چکیده: در برخی از بررسی‌های نمونه‌ای عملی ممکن است مقادیر واقعی مشخصه‌ها (متغیرها) مشاهده نشوند، بلکه مقادیر دیگری که با خطاهای اندازه‌گیری آمیخته شده‌اند، مشاهده شوند. در این مقاله با فرض مدل خطای اندازه‌گیری ساده‌ای بین مقدار واقعی و مقدار مشاهده شده مشخصه، پارامترهای جامعه متناهی به روش‌های کلاسیک، بیز و مینیماکس پیشگویی می‌شوند. سپس با استفاده از داده‌های طرح هزینه و درآمد خانوار، میانگین هزینه سالیانه خانوارها با فرض مدل خطای اندازه‌گیری پیشگویی و با برآورد جاری که بدون در نظر گرفتن خطای اندازه‌گیری محاسبه شده، مقایسه می‌شود.

واژه‌های کلیدی: خطاهای اندازه‌گیری، پیشگویی، جامعه متناهی، نمونه‌گیری دومرحله‌ای، پیشگویی بیز و پیشگوی مینیماکس

۱ مقدمه

مبحث نمونه‌گیری یکی از مباحث مهم علم آمار و دارای کاربرد فراوان است. در یک بررسی نمونه‌ای (نمونه‌گیری) با دو نوع خطا، خطای نمونه‌گیری^۱ و خطای غیرنمونه‌گیری^۲ سروکار داریم. در یک بررسی نمونه‌ای، اطلاعات به دست آمده از نمونه مبنای نتیجه‌گیری دربارهٔ جامعه را تشکیل می‌دهد. از آنجایی که نمونه صرفاً بخشی از جامعه است، این اطلاعات کامل نبوده و بنابراین هرگونه نتیجه‌گیری همراه با خطا خواهد بود. این نوع خطا که به دلیل استفاده از یک نمونه به جای سرشماری ایجاد می‌شود، اصطلاحاً خطای نمونه‌گیری گفته می‌شود. بدیهی است که تنها راه حذف این خطا، انجام یک سرشماری است. در مراحل اجرای یک نمونه‌گیری امکان بروز خطاهایی از قبیل خطای پوشش، خطای اندازه‌گیری، خطا در تکمیل پرسشنامه و خطای پردازش داده‌ها وجود دارد. تأثیر این گونه خطاها بر نتایج بررسی‌ها، اصطلاحاً خطای غیرنمونه‌گیری گفته می‌شود.

1) Sampling Error 2) Nonsampling Error

در یک بررسی نمونه‌ای معمولاً کار بران برای ارزیابی دقت روش به کار گرفته شده به خطاهای استاندارد^۳ برآورد شده متکی هستند. خطاهای استاندارد برآورد شده تنها منعکس‌کننده خطاهای نمونه‌گیری هستند و برای استنباط‌های آماری از قبیل بازه‌های اطمینان و آزمون فرض‌ها استفاده می‌شوند. نکته مهمی که بایستی متذکر شد این است که در روش‌های متفاوت نمونه‌گیری که معمولاً مطالعه می‌شوند، فرمول‌هایی که ارائه می‌شوند مبتنی بر این فرض هستند که وقتی واحدهای نمونه مشخص می‌شوند و به اندازه‌گیری مشخصه آنها پرداخته می‌شود و براساس این اندازه‌ها، فرمول‌ها به صورت ریاضی و احتمالی بیان می‌شوند، خطای اندازه‌گیری وجود ندارد و تنها خطای موجود، خطای نمونه‌گیری است. در حالی که در برخی از بررسی‌های نمونه‌ای عملی ممکن است مقدار واقعی مشخصه، مشاهده نشود بلکه مقادیر دیگری که با خطاهای اندازه‌گیری آمیخته شده‌اند، مشاهده شوند.

خطای اندازه‌گیری به مشاهده مشخصه‌هایی که در یک بررسی اندازه‌گیری می‌شوند مربوط است و بعضی وقت‌ها به عنوان «خطای مشاهده‌ای» نسبت داده می‌شود. خطای اندازه‌گیری به عنوان قسمتی از خطای گردآوری داده‌ها، در مقابل خطای نمونه‌گیری، بی‌پاسخی، پوشش یا پردازش داده‌ها پدید می‌آید. خطای اندازه‌گیری می‌تواند از چهار منبع به وجود آید: پرسشنامه، روش گردآوری داده‌ها، مصاحبه‌گر و پاسخگو.

در یک بررسی نمونه‌ای، اگر یک مدل خطای اندازه‌گیری بین مقدار واقعی و مقدار مشاهده شده مشخصه، در نظر گرفته شود آن‌گاه می‌توان پارامترهای جامعه متناهی را به طور واقعی‌تری پیشگویی نمود.

در خصوص به کار بردن مدل‌های خطای اندازه‌گیری، ابتدا اسپرنت^۴ (۱۹۶۶) مدل رگرسیونی که در آن متغیرهای تبیینی با خطای اندازه‌گیری آمیخته شده‌اند را معرفی کرد و روشی براساس رهیافت کمترین توان‌های دوم تعمیم‌یافته، برای برآورد ضریب‌های رگرسیونی خطی پیشنهاد کرد. بولفارین^۵ (۱۹۹۱)؛ موخوپادیای^۶ (۱۹۹۲ و ۱۹۹۴)؛ باتاچاریا^۷ (۱۹۹۷)؛ پیشگوی کلاسیک میانگین و واریانس جامعه متناهی را تحت مدل مکان خطای اندازه‌گیری در نظر گرفتند. در مورد پیشگویی بیز و مینیماکس بولفارین و زاکس^۸ (۱۹۹۱) مقاله‌ای ارائه کردند.

در این مقاله با فرض مدل خطای اندازه‌گیری ساده‌ای بین مقدار واقعی و مقدار مشاهده شده مشخصه، پارامترهای جامعه متناهی به روشهای کلاسیک، بیز و مینیماکس پیشگویی می‌شوند. سپس با استفاده از داده‌های طرح هزینه و درآمد خانوار، میانگین هزینه سالیانه خانوارها پیشگویی و با برآورد جاری مقایسه می‌شود.

بر این اساس در بخش ۲ بطور خلاصه به نمونه‌گیری مبتنی بر طرح و مبتنی بر مدل اشاره شده و در بخش ۳ پیشگویی تحت مدل خطای اندازه‌گیری جمعی مطرح می‌گردد. در بخش ۴ پیشگویی بیز و مینیماکس میانگین جامعه مورد بررسی قرار گرفته و پیشگویی میانگین هزینه سالیانه خانوارها در بخش ۵ ارائه می‌شود.

3) Standard Errors 4) Sprent 5) Bolfarine 6) Mukhopadhyay
7) Bhattacharyya 8) Zacks

۲ نمونه‌گیری مبتنی بر طرح و نمونه‌گیری مبتنی بر مدل

در رهیافت مبتنی بر طرح در بررسی نمونه‌ای، مقادیر متغیر (مشخصه) مورد نظر (Y - مقادیر) جامعه به عنوان کمیت‌های ثابت در نظر گرفته می‌شوند. در این رهیافت احتمال‌های انتخاب نمونه‌های مرتبط با طرح نمونه‌گیری مطرح هستند که در به دست آوردن امید ریاضی، واریانس، اربیبی و ویژگی‌های دیگر برآوردگرها استفاده می‌شوند.

از طرف دیگر، در رهیافت مبتنی بر مدل، مقادیر متغیرهای مورد نظر جامعه به عنوان متغیرهای تصادفی در نظر گرفته می‌شوند و ویژگی‌های برآوردگرها به توزیع توأم این متغیرهای تصادفی وابسته هستند. به طور مثال وقتی فشار خون اندازه‌گیری می‌شود، مقدار مشاهده شده برآوردی از مقدار واقعی است. بنابراین مقدار واقعی فشار خون واحدها به طور صحیح مشاهده نمی‌شوند بلکه مقادیر مشاهده شده، برآوردهایی از فشار خون واقعی واحدها هستند.

فرض کنید $\mathbf{Y} = (Y_1, Y_2, \dots, Y_N)$ به برداری از Y مقادیر مربوط به N واحد جامعه اشاره کند. از دیدگاه طرح، این Y - مقادیر، مقادیر ثابتی هستند. از دیدگاه طرح، این Y - مقادیر، متغیرهای تصادفی با توزیع توأم ξ هستند. فرض کنید $p(s)$ به احتمال انتخاب نمونه s تحت یک طرح نمونه‌گیری اشاره کند. s دنباله‌ای یا زیرمجموعه‌ای از واحدهای جامعه است. به وسیله این نمونه می‌توان برخی از کمیت‌های جامعه (Y) را برآورد یا پیشگویی کرد. Y می‌تواند میانگین جامعه یا مقدار کل جامعه باشد. پیشگو یا برآوردگر \hat{Y} ، تابعی از Y - مقادیر نمونه است. برآوردگر یا پیشگوی \hat{Y} برای Y طرح - ناریب^۹ گفته می‌شود اگر امید ریاضی شرطی آن به شرط تحقق N ، Y مقادیر جامعه، برابر با مقدار تحقق یافته Y باشد. یعنی

$$E(\hat{Y} | \mathbf{Y}) = Y.$$

که در آن $\mathbf{Y} = (Y_1, Y_2, \dots, Y_N)$ است. برآوردگر یا پیشگوی \hat{Y} برای Y مدل - ناریب^{۱۰} گفته می‌شود اگر به شرط هر نمونه s ، امید ریاضی شرطی \hat{Y} برابر امید ریاضی Y باشد. یعنی

$$E(\hat{Y} | s) = E(Y).$$

پیشگوی مدل - ناریب \hat{Y} برای Y تحت توزیع توأم ξ جامعه و هر طرح نمونه‌گیری انتخاب نمونه s به شرط نمونه s ، ناریب است. در نتیجه این ناریبی به طرح نمونه‌گیری که با استفاده از آن نمونه s انتخاب شده است، بستگی ندارد.

برآوردگر یا پیشگوی \hat{Y} برای Y ناریب است (یعنی ناریب غیرشرطی)، اگر امید ریاضی \hat{Y} برابر امید ریاضی Y باشد ($E(\hat{Y}) = E(Y)$).

۳ پیشگویی تحت مدل خطای اندازه‌گیری جمعی

فرض کنید جامعه متناهی P با تعداد مشخص N واحد با برچسب‌های $1, \dots, i, \dots, N$ نمایش داده شده باشد. در ارتباط با واحد i - ام جامعه، y_i مقدار واقعی مشخصه مورد مطالعه است که نمی‌توان آن را به طور صحیح مشاهده کرد اما یک مقدار متفاوت Y_i که با خطای اندازه‌گیری آمیخته شده است، مشاهده می‌شود. فولر^{۱۱} (۱۹۸۷، ۱۹۸۹) رفتار کلی برای استنباط تحت مدل‌های خطای اندازه‌گیری جمعی را مورد مطالعه قرار داده است. بولفارین (۱۹۹۱)، موخوادیای (۱۹۹۲) و (۱۹۹۴) و باتاچاریا (۱۹۹۷) و دیگران مدل ذیل را برای پیشگویی میانگین و واریانس جامعه متناهی تحت مدل خطای اندازه‌گیری در نظر گرفتند.

$$\begin{aligned} y_i &= \mu + e_i, & e_i &\sim N(0, \sigma_{ee}), \\ Y_i &= y_i + u_i, & u_i &\sim N(0, \sigma_{uu}), \end{aligned} \quad (1)$$

که در آن μ میانگین مقدار واقعی مشخصه مورد نظر در جامعه، $(\sigma_{ee} > 0)$ و $(\sigma_{uu} > 0)$ مقادیر ثابت هستند و نماد $e_i \sim N(0, \sigma_{ee})$ به متغیرهای تصادفی نرمال مستقل و دارای توزیع یکسان $(i.i.d.)$ ، با میانگین صفر و واریانس σ_{ee} اشاره می‌کند و e_i و u_j برای $i, j = 1, 2, \dots, N$ از یکدیگر مستقل هستند.

برای پیشگویی میانگین جامعه، رده پیشگوهای خطی ذیل در نظر گرفته می‌شود

$$g(s, \mathbf{Y}_s) = b_s + \sum_{i \in s} b_{is} Y_i$$

که در آن b_s و b_{is} مقادیر ثابتی هستند که به Y - مقادیر بستگی ندارند و $\mathbf{Y}_s = (Y_1, Y_2, \dots, Y_n)$ به مجموعه مقادیر مشاهده شده واحدهای نمونه s اشاره می‌کند.

پیشگوی $g(s, \mathbf{Y}_s)$ برای میانگین جامعه (\bar{y}) ناریب خواهد بود اگر و تنها اگر

$$E(g(s, \mathbf{Y}_s)) = E(b_s + \sum_{i \in s} b_{is} Y_i) = E(\bar{y}) = \mu$$

با توجه به این که میانگین نمونه $(\bar{Y}_s = \frac{1}{n} \sum_{i \in s} Y_i)$ در شرط فوق صدق می‌کند از این رو پیشگوی ناریب میانگین جامعه (\bar{y}) خواهد بود. همچنین با توجه به این که Y_i - ها از یکدیگر مستقل و دارای توزیع نرمال هستند، با استفاده از قضیه برآورد (پیشگو) 13 MVU رانو^{۱۴} (۱۹۷۳)، میانگین نمونه (\bar{Y}_s) پیشگوی ناریب با مینیمم واریانس جامعه (\bar{y}) است.

نتیجه حاصل از بخش قبل این است که در هر طرح نمونه‌گیری ناآگاهی بخش با اندازه مؤثر ثابت نمونه (n) میانگین نمونه پیشگوی بهینه میانگین جامعه است. بنابراین \bar{Y}_s برای هر طرح

11) Fuller 12) independent identically distributed 13) Minimum Variance Unbiased 14) Rao

نمونه‌گیری ناآگاهی بخش با اندازه مؤثر ثابت نمونه (n) و تحت مدل (۱) پیشگوی بهینه برای میانگین جامعه (\bar{y}) است.

از طرف دیگر با معرفی این که V به واریانس متناظر با عملیات توأم مدل (۱) و طرح نمونه‌گیری p اشاره می‌کند، واریانس پیشگوی میانگین جامعه به صورت ذیل حاصل می‌شود

$$V(\bar{Y}_s - \bar{y}) = \left(\frac{1}{n} - \frac{1}{N}\right)\sigma_{ee} + \frac{1}{n}\sigma_{uu} = \frac{r_0}{N^2} \quad (2)$$

که در آن $r_0 = N^2 \left\{ \left(\frac{1}{n} - \frac{1}{N}\right)\sigma_{ee} + \frac{1}{n}\sigma_{uu} \right\}$ است.

مدل‌هایی که با استفاده از مدل (۱) توصیف شدند به پارامترهای σ_{uu} و σ_{ee} بستگی دارند که آنها نامعلوم هستند، از این رو ساختاری از برآوردگرهای سازگار برای پارامترهای مدل به اطلاعات اضافی از جامعه نیاز دارد. اکنون دو حالت ذیل برای برآورد این پارامترها در نظر گرفته می‌شود. حالت اول فرض می‌شود واریانس خطای اندازه‌گیری، σ_{uu} ، معلوم است. همان‌گونه که در فولر (۱۹۸۷) اشاره شده است، با اندازه‌گیری‌های تکراری می‌توان مقادیر معقولی را برای σ_{uu} به دست آورد. با استفاده از روش برآوردگر گشتاوری، برآوردگر σ_{ee} به صورت $s_Y^2 - \sigma_{uu}$ حاصل می‌شود، که در آن $s_Y^2 = \frac{1}{(n-1)} \sum_{i \in s} (Y_i - \bar{Y}_s)^2$ است. حالت دوم که ممکن است در موارد خاصی روی دهد این است که نسبت $\delta = \frac{\sigma_{ee}}{\sigma_{uu}}$ معلوم باشد. در این حالت برآوردگر روش گشتاوری σ_{uu} برابر با $\frac{1}{1+\delta} s_Y^2$ است.

برای به دست آوردن پیشگوی بهینه واریانس جامعه $(V(y) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2)$ همانند پیشگوی میانگین جامعه عمل می‌شود و رده پیشگوهای ناریب درجه دوم ذیل در نظر گرفته می‌شود.

$$g(s, \mathbf{Y}_s) = b_s + \sum_{i \in s} b_{is} Y_i^2 + \sum_{i \neq i' \in s} b_{ii'} Y_i Y_{i'}$$

پیشگوی $g(s, \mathbf{Y}_s)$ برای واریانس جامعه $(V(y))$ ناریب خواهد بود اگر و تنها اگر

$$E(g(s, Y_s)) = V(y)$$

با توجه به این که $s_Y^2 = \frac{1}{(n-1)} \sum_{i \in s} (Y_i - \bar{Y}_s)^2$ در شرط فوق صدق می‌کند بنابراین پیشگوی ناریب واریانس جامعه $(V(y))$ خواهد بود.

همچنین Y_i ها از یکدیگر مستقل و دارای توزیع نرمال هستند و همانند پیشگوی ناریب با مینیمم واریانس میانگین جامعه، با استفاده از قضیه برآورد (پیشگو) MVU راتو (۱۹۷۳) نتیجه می‌شود که واریانس نمونه (s_Y^2) ، پیشگوی ناریب با مینیمم واریانس $V(y)$ است.

۴ پیشگویی بیز و مینیماکس

در بخش قبل در مدل (۱)، μ کمیته ثابت در نظر گرفته شده بود. در روش بیزی μ متغیری تصادفی در نظر گرفته می‌شود و تغییرات آن توسط یک توزیع احتمال بیان می‌گردد. در این بخش پیشگوی بیز و مینیماکس مقدار کل جامعه در نمونه‌گیری تصادفی ساده و دو مرحله‌ای مورد بررسی قرار می‌گیرد.

۱.۴ پیشگویی بیز و مینیماکس مقدار کل جامعه در نمونه‌گیری تصادفی ساده

موخوپادیای (۲۰۰۰) با در نظر گرفتن مدل (۱) و نظر به این که در نمونه‌های بزرگ، توزیع اغلب متغیرهای وابسته به موضوع‌های اقتصادی و اجتماعی (حداقل به طور تقریبی) نرمال هستند، یک توزیع پیشین نرمال برای μ در نظر گرفته است.

$$\mu \sim N(\tau, \theta^2) \quad (۳)$$

در قضیه ذیل پیشگوی بیز مقدار کل جامعه تحت مدل (۱) و توزیع پیشین (۳) به دست آورده می‌شود.

قضیه ۱.۴ (موخوپادیای (۲۰۰۰) تحت مدل (۱) و توزیع پیشین (۳) برای μ ، اگر $Y_s = (Y_1, Y_2, \dots, Y_n)$ یک نمونه تصادفی ساده به اندازه n انتخاب شده از جامعه P باشد، آنگاه پیشگوی بیز مقدار کل جامعه $(T(y) = \sum_{i=1}^n y_i)$ تحت تابع زبان درجه دوم خطا، عبارت است از

$$\hat{T}_B = \frac{n\bar{Y}_s \sigma_{ee}}{\sigma^2} + \left[\frac{n\sigma_{uu} + (N-n)\sigma^2}{\sigma^2} \cdot \frac{\tau\sigma^2 + n\bar{Y}_s\theta^2}{\sigma^2 + n\theta^2} \right] \quad (۴)$$

که در آن $\sigma^2 = \sigma_{ee} + \sigma_{uu}$ و مخاطره بیز آن عبارت است از

$$r_\theta(\hat{T}_B) = n\sigma_e^2 + (N-n)\sigma_{ee} + (N - \frac{n\sigma_{ee}}{\sigma^2})^2 \cdot \frac{\theta^2\sigma^2}{\sigma^2 + n\theta^2} \quad (۵)$$

برای به دست آوردن پیشگوی مینیماکس مقدار کل جامعه از روش بیز حدی استفاده می‌شود. چون

$$\lim_{\theta \rightarrow \infty} r_\theta(\hat{T}_B) = \frac{N(N-n)}{n} \sigma^2 + N\sigma_{uu} = r.$$

و همچنین مخاطره پیشگوی $N\bar{Y}_s$ با توجه به رابطه (۲) برابر با r است لذا پیشگوی $N\bar{Y}_s$ پیشگوی مینیماکس مقدار کل جامعه است.

پیشگویی بیز واریانس جامعه در نمونه‌گیری تصادفی ساده

باتاچاریا (۱۹۹۷) پیشگوی بیز واریانس جامعه $(V(y) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2)$ را تحت مدل (۱) و با در نظر گرفتن توزیع پیشین (۳) به صورت ذیل به دست آورد.

قضیه ۲.۴ (باتاچاریا ۱۹۹۷) تحت مدل (۱) و توزیع پیشین (۳) برای μ ، پیشگوی بیز واریانس جامعه تحت تابع زیان درجه دوم خطا عبارت است از

$$\begin{aligned} E(V(y)|Y_s) &= E \left[\left(\frac{n}{N} s_y^2 + \frac{(N-n)}{N} s_{yr}^2 + \frac{n(N-n)}{N^2} (\bar{y}_r - \bar{y}_s)^2 \right) | Y_s \right] \\ &= \frac{\sigma_e^2}{N} (n-1 + 2\lambda) + \frac{(N-n-1)}{N} \sigma_{ee} + \frac{n(N-n)}{N^2} \times \\ &\quad \left\{ \sigma_y^2 + \frac{\sigma_{ee}^2}{\sigma^2} \left[\bar{Y}_s^2 - 2\bar{Y}_s \left(\frac{\tau\sigma^2 + n\bar{Y}_s\theta^2}{\sigma^2 + n\theta^2} \right) + \right. \right. \\ &\quad \left. \left. \left(\frac{\tau\sigma^2 + n\bar{Y}_s\theta^2}{\sigma^2 + n\theta^2} \right)^2 + \frac{\theta^2\sigma^2}{\sigma^2 + n\theta^2} \right] \right\} \end{aligned}$$

که در آن

$$\begin{aligned} \lambda &= \frac{(n-1)\sigma_{ee}^2 \cdot s_Y^2}{\sigma^2 \sigma_y^2}, \quad s_Y^2 = \frac{1}{(n-1)} \sum_{i \in S} (Y_i - \bar{Y}_s)^2, \\ \sigma_e^2 &= \frac{\sigma_{uu} \cdot \sigma_{ee}}{\sigma^2}, \quad \sigma^2 = \sigma_{ee} + \sigma_{uu}, \quad Y_s = (Y_1, Y_2, \dots, Y_n). \end{aligned}$$

۲.۴ پیشگویی بیز در نمونه‌گیری دومرحله‌ای

اخیراً در جامعه‌های انسانی، بررسی‌های نمونه‌ای دومرحله‌ای کاربرد فراوان پیدا کرده‌اند. بولفارین (۱۹۹۱) و موخوپادیای (۱۹۹۵) پیشگویی بیز در نمونه‌گیری دومرحله‌ای از جامعه متناهی تحت مدل خطای اندازه‌گیری را مورد مطالعه قرار داده‌اند.

در این قسمت جامعه متناهی P که شامل K زیرجامعه (خوشه) دوبه‌دو مجزا به اندازه‌های M_h ، $h = 1, 2, \dots, K$ است، در نظر گرفته می‌شود. همچنین y_{hj} به مشخصه مورد نظر مربوط به واحد j - ام در خوشه h - ام جامعه، $h = 1, 2, \dots, K$ و $j = 1, 2, \dots, M_h$ اشاره می‌کند. برای انتخاب نمونه از جامعه متناهی P ، در مرحله نخست یک نمونه n خوشه‌ای از K خوشه جامعه انتخاب می‌شود، این نمونه n خوشه‌ای با s نمایش داده می‌شود. در مرحله دوم از هر یک از خوشه‌ها که در مرحله نخست انتخاب شده‌اند، نمونه s_h ای به اندازه m_h از آن خوشه انتخاب می‌شود. پس از این که نمونه از جامعه انتخاب شد، مقدار کل جامعه را می‌توان

به صورت ذیل نوشت

$$T(y) = \sum_{h \in s} \sum_{j \in s_h} y_{hj} + \sum_{h \in s} \sum_{j \notin s_h} y_{hj} + \sum_{h \notin s} \sum_{j=1}^{M_h} y_{hj} \quad (6)$$

برای پیشگویی مقدار کل جامعه مانند قسمت‌های قبلی مدل ذیل در نظر گرفته می‌شود

$$\begin{aligned} y_{hj} &= \mu_h + e_{hj}, \\ \mu_h &= \mu + v_h, \\ Y_{hj} &= y_{hj} + u_{hj} \end{aligned} \quad (7)$$

که در آن μ ، میانگین مقدار واقعی مشخصه مورد نظر در جامعه، و μ_h ، میانگین مقدار واقعی مشخصه در خوشه h - ام جامعه است و e_{hj} ، v_h و u_{hj} دویه‌دو از یکدیگر مستقل هستند، $e_{hj} \sim N(0, \sigma_{eeh})$ ، $v_h \sim N(0, \delta^2)$ و $u_{hj} \sim N(0, \sigma_{uuh})$ ، $h = 1, 2, \dots, K$ و $j = 1, 2, \dots, M_h$ در این مدل Y_{hj} به مقدار مشاهده شده مشخصه مربوط به واحد j - ام در خوشه h - ام اشاره می‌کند. مقادیر نمونه که به روش دومرحله‌ای از جامعه متناهی انتخاب شده است به صورت ذیل نمایش داده می‌شود

$$Y_s = \{Y_{hj}; j \in s_h, h \in s\}$$

قضیه ۳.۴ (موخوپادایای ۱۹۹۵) تحت مدل (۷) و با توجه به این که در روش نمونه‌گیری دومرحله‌ای μ_h دارای توزیع $N(\mu, \delta^2)$ است و با فرض یک توزیع پیشین ناآگاهی‌بخش برای μ (یعنی، $P(\mu) \propto const.$)، اگر $\{Y_{hj}; j \in s_h, h \in s\}$ یک نمونه تصادفی باشد که به روش دومرحله‌ای از جامعه متناهی انتخاب شده باشد آن‌گاه پیشگوی بیز مقدار کل جامعه $(T(y) = \sum_{h=1}^K \sum_{j=1}^{M_h} y_{hj})$ تحت تابع زیان درجه دوم خطا عبارت است از

$$\begin{aligned} \hat{T}_B = E(T|Y_s) &= \sum_{h \in s} m_h \frac{\sigma_{eeh}}{\sigma_{eeh} + \sigma_{uuh}} \bar{Y}_{sh} \\ &+ \sum_{h \in s} m_h \frac{\sigma_{uuh}}{\sigma_{eeh} + \sigma_{uuh}} \{\lambda_h \bar{Y}_{sh} + (1 - \lambda_h) \bar{Y}_s\} \\ &+ \sum_{h \notin s} (M_h - m_h) \{\lambda_h \bar{Y}_{sh} + (1 - \lambda_h) \bar{Y}_s\} + \bar{Y}_s \cdot \sum_{h \notin s} M_h \end{aligned} \quad (8)$$

و مخاطره بیز آن عبارت است از

$$V(T|Y_s) = \sum_{h \in s} B_h^2 c_{hh} + \sum_{h \notin s} M_h^2 c_{hh} + \sum_{h \in s} \sum_{q \neq h \in s} B_h B_q c_{hq}$$

$$\begin{aligned}
 & + \sum_{h \notin s} \sum_{q \neq h \notin s} M_h M_q c_{hq} + \sum_{h \in s} \sum_{q \neq h \notin s} B_h M_q c_{hq} \\
 & + \sum_{h \in s} m_h \frac{\sigma_{uuh} \cdot \sigma_{eeh}}{\sigma_{eeh} + \sigma_{uuh}} + \sum_{h \in s} (M_h - m_h) \sigma_{eeh} + \sum_{h \notin s} M_h \sigma_{eeh}. \quad (9)
 \end{aligned}$$

که در آن $B_h = m_h \frac{\sigma_{uuh}}{\sigma_{eeh} + \sigma_{uuh}} + (M_h - m_h)$ است.

۵ پیشگویی هزینه سالیانه خانوارها

یکی از شیوه‌های ارزیابی سطح زندگی افراد یک جامعه، دسترسی به الگوهای متفاوت هزینه و درآمد خانوارها با ویژگی‌های متفاوت اجتماعی و اقتصادی است که این هدف با اجرای طرح «آمارگیری از هزینه و درآمد خانوار» تحقق می‌یابد.

با توجه به این که نتایج آمارگیری از هزینه و درآمد خانوار می‌تواند بیان‌کننده میزان و نحوه توزیع هزینه‌ها و درآمدها و چگونگی استفاده خانوارها از تسهیلات زندگی در سطح ملی و منطقه‌ای باشد، لذا این طرح کاربرد بسیار وسیعی در ارزیابی عملکرد و نتایج تصمیم‌گیری‌ها و سیاست‌های گذشته و برنامه‌ریزی و سیاست‌گذاری اقتصادی، اجتماعی و فرهنگی کشور دارد.

اهمیت روزافزون کاربرد نتایج طرح هزینه و درآمد خانوار در کشور ایجاب می‌کند که مطالعاتی برای بهبود روش‌های کار در مراحل متفاوت تهیه، اجرا و استخراج نتایج این طرح انجام شود و در صورت لزوم برای افزایش درستی و دقت نتایج، برخی از این روش‌ها مورد تجدید نظر واقع شوند. با توجه به اهمیت موضوع به ویژه در مرکز آمار ایران که متولی تولید آمار رسمی کشور است، مدیریت پژوهش‌های آماری مرکز آمار ایران در سال ۱۳۷۷ یک طرح تحقیقاتی تحت عنوان «مقایسه دو شیوه اطلاع‌گیری از هزینه و درآمد خانوار» انجام داده است. در این طرح تحقیقاتی از تعدادی از خانوارهای معمولی ساکن شهر تهران، اطلاعات یکسانی به دو روش، روش مصاحبه حضوری (روش فعلی اجرای طرح) و روش ثبت اطلاعات توسط خانوار که انتظار می‌رود دارای دقت و درستی بیشتری باشد، در دو مرحله گردآوری شده است. با مقایسه نتایج حاصل از این دو مرحله و انجام تجزیه و تحلیل‌های لازم، روش فعلی گردآوری اطلاعات طرح «هزینه و درآمد خانوار» مورد ارزیابی قرار گرفته است.

در این بخش با استفاده از اطلاعات طرح تحقیقاتی فوق‌الذکر، مدلی برای خطای اندازه‌گیری هزینه خانوارها در نظر گرفته و تحت این مدل میانگین هزینه سالیانه خانوارها پیشگویی می‌شوند. لازم به ذکر است که در آمارگیری جاری هزینه و درآمد خانوار، برآوردها با فرض این که خطای اندازه‌گیری صفر است، محاسبه می‌شوند.

مشخصات کلی

هدف اصلی این بخش پیشگویی کلاسیک میانگین هزینه سالیانه خانوارهای معمولی ساکن شهر تهران تحت مدل خطای اندازه‌گیری و مقایسه آن با برآوردهای شیوه جاری است. جامعه آماری کلیه خانوارهای معمولی ساکن شهر تهران است و هر یک از این خانوارها به عنوان یک واحد آماری هستند. چارچوب آماری مورد استفاده، فهرست کلیه بلوکهای شهر تهران در سرشماری عمومی نفوس و مسکن سال ۱۳۷۵ به انضمام تعداد خانوار هر یک از این بلوکهاست و زمان آمارگیری برای سؤال‌های متفاوت، بهمن ماه سال ۱۳۷۷ و سال مورد نظر (از اول اسفند ماه ۱۳۷۶ لغایت روز سی‌ام بهمن ماه ۱۳۷۷) است.

روش نمونه‌گیری و اندازه نمونه

در این طرح، روش نمونه‌گیری دومرحله‌ای است. مرحله اول مربوط به انتخاب بلوک‌های نمونه و مرحله دوم مربوط به انتخاب خانوارهای نمونه در هر یک از این بلوکهاست. با توجه به محدودیت‌های اجرایی، امکانات طرح و اهداف مورد نظر، تعداد خانوار نمونه این طرح 50° خانوار در نظر گرفته شده است. در مرحله اول نمونه‌گیری، 50° بلوک انتخاب شده است و در مرحله دوم نمونه‌گیری از هر یک از بلوک‌ها که در مرحله اول انتخاب شده‌اند، 10° خانوار انتخاب می‌شود.

یادآور می‌شود همان‌طور که در فوق مطرح شد، هدف اصلی پیشگویی میانگین هزینه خانوارهای معمولی ساکن شهر تهران است. ولیکن برای پیشگویی پارامترها در رهیافت مبتنی بر مدل، پارامترها با توجه به مدل (۱) پیشگویی می‌شوند که در این مدل فرض شده است مانده‌ها (متغیرها) مستقل و دارای توزیع یکسان نرمال هستند، لذا باید ابتدا این فرض‌ها آزمون شوند. با آزمودن مانده‌ها (متغیرها) مشخص شد که فرض نرمال بودن برای آنها صادق نیست ولیکن با تغییر متغیر ریشه چهارم هزینه خوراکی و لگاریتم هزینه غیرخوراکی سالیانه خانوارها فرض نرمال بودن صادق شد. در این قسمت به همین دلیل ابتدا برآوردها و پیشگوهای میانگین ریشه چهارم هزینه خوراکی سالیانه خانوارها و لگاریتم هزینه غیرخوراکی سالیانه خانوارها محاسبه می‌شوند و سپس با تبدیل عکس، میانگین هزینه خوراکی سالیانه خانوارها و میانگین هزینه غیرخوراکی سالیانه خانوارها پیشگویی خواهند شد.

مقادیر مربوط به برآورد (بدون خطای اندازه‌گیری) میانگین هزینه‌های خوراکی و غیرخوراکی به همراه خطای استاندارد آنها و همچنین پیشگوی کلاسیک (با در نظر گرفتن خطای اندازه‌گیری) میانگین هزینه‌های خوراکی و غیرخوراکی به همراه خطای استاندارد آنها در جداول ذیل آمده‌اند. برای محاسبه مقادیر این جداول از روابط بخش ۳ استفاده شده است.

مقادیر مربوط به هزینه خوراکی سالیانه خانوارها

| خطای استاندارد | برآورد واریانس | میانگین | روش |
|----------------|------------------------------|-------------|--|
| ۳۰۰۷۵۴,۶۲۲۳ | $۹,۰۴۵۳۳۴۲۸۳ \times ۱۰^{۱۰}$ | ۶۵۰۸۴۲۷,۶۲۴ | برآورد معمولی (بدون خطای اندازه‌گیری) |
| ۱۶۹۳۹۵,۲۰۵۵ | $۲,۸۶۹۴۷۳۵۶۴ \times ۱۰^{۱۰}$ | ۶۸۹۳۴۶۷,۸۴۶ | پیشگوی کلاسیک (با خطای اندازه‌گیری) |

مقادیر مربوط به هزینه غیرخوراکی سالیانه خانوارها

| خطای استاندارد | برآورد واریانس | میانگین | روش |
|----------------|------------------------------|-------------|--|
| ۸۶۸۱۲۴,۹۷۸۵ | $۷,۵۳۶۴۰۹۷۸۳ \times ۱۰^{۱۱}$ | ۱۵۸۹۴۷۳۴,۶ | برآورد معمولی (بدون خطای اندازه‌گیری) |
| ۲۳۲۱۵۷,۳۷۹ | $۵,۳۸۹۷۰۴۸۶۴ \times ۱۰^{۱۰}$ | ۱۷۱۳۹۴۵۸,۶۳ | پیشگوی کلاسیک (با خطای اندازه‌گیری) |

همان‌طور که از جدول‌های فوق ملاحظه می‌شود برآورد واریانس برآوردگر میانگین هزینه سالیانه خانوارها بدون در نظر گرفتن خطای اندازه‌گیری از برآورد واریانس پیشگوی میانگین هزینه سالیانه خانوارها با در نظر گرفتن خطای اندازه‌گیری بیشتر است. این به دلیل استفاده از نمونه‌ای است که به شیوه نمونه‌گیری دومرحله‌ای از جامعه انتخاب شده است و با ملاحظه فرمول برآورد واریانس برآوردگر میانگین جامعه مشخص می‌شود که آن تابعی از واریانس درون بلوک‌ها و واریانس بین بلوک‌هاست و به دلیل ناهمگنی بین بلوک‌ها، برآورد واریانس برآوردگر میانگین جامعه نسبت به برآورد واریانس پیشگوی میانگین جامعه بزرگتر شده است.

در انتها پیشنهاد می‌شود در مواردی که خطای اندازه‌گیری مشخصه‌ها قابل اغماض نیستند برای دسترسی به اطلاعات واقعی‌تر از جامعه مورد نظر با در نظر گرفتن یک مدل خطای اندازه‌گیری، پارامترهای جامعه با استفاده از این رهیافت پیشگویی شوند.

مراجع

- [۱] گزارش طرح تحقیقاتی مقایسه دو شیوه اطلاع‌گیری از هزینه و درآمد خانوار (۱۳۷۸)، مدیریت پژوهش‌های آماری، مرکز آمار ایران.
- [2] Bhattacharyya, S. (1997). Some studies on estimation of mean and variances in finite population sampling. Unpublished Ph.D. Thesis submitted to Indian Statistical Institute, Calcutta.
- [3] Bolfarine, H. (1991). Finite-population prediction under error-in-variables super-population models. *Canadian Journal of Statistics*, 19(2), 191-207.
- [4] Bolfarine, H., and Zacks, S. (1991). Bayes and minimax prediction in finite population. *Journal of Statistical Planning and Inference*, 28, 139-151.
- [5] Fuller, W. A. (1987). *Measurement Error Models*. Wiley, New York.
- [6] Fuller, W. A. (1989). Prediction of the true values for the measurement error models, In P. J. Brown and W. A. Fuller, eds., *Statistical Analysis of Measurement Error Models and Applications* Am. Math. Soc., Providence, RI.
- [7] Mukhopadhyay, P. (1992). On prediction in finite population under error-in-variables superpopulation models., Indian Statistical Institute Technical Report No. ASC/92/11.
- [8] Mukhopadhyay, P. (1994). Prediction in finite population under error-in-variables superpopulation models. *Journal of Statistical Planning and Inference*, 41, 151-161.
- [9] Mukhopadhyay, P. (1995). Bayes and minimax estimator for two-stage sampling from a finite population under measurement error models. *Communications in statistics, Theory and Methods*, 24(3), 663-674.
- [10] Mukhopadhyay, P. (2000). *Topics in Survey Sampling*. Springer-Verlag, New York.
- [11] Rao, C. R. (1973). *Linear Statistical Inference*. Second edition, John Wiley and Sons, New York.
- [12] Sprent, P. (1966). A generalized least squares approach to linear function relationships. *J. Roy. Statist. Soc. Ser. B*, 28, 279-297.

استفاده از برآوردگرهای عادی و رگرسیونی نمونه‌گیری مجموعه رتبه‌دار در برآورد مقدار کل محصول گندم ایران

نادر نعمت‌الهی^۱، لطیف سعادتی^۲

^۱ گروه آمار دانشگاه علامه طباطبایی

^۲ کارشناسی ارشد آمار دانشگاه علامه طباطبایی

چکیده: در بعضی از جوامع مورد بررسی ساختار جامعه معلوم نیست و معمولاً از روش نمونه‌گیری تصادفی ساده برای برآورد میانگین متغیر مورد نظر استفاده می‌شود. برآوردگر حاصل از روش نمونه‌گیری تصادفی ساده به علت اینکه کنترلی روی واحدهای انتخاب شده نیست از درجه اعتماد بالایی برخوردار نمی‌باشد. در این مواقع می‌توان از روش نمونه‌گیری مجموعه رتبه‌دار استفاده کرد. روش‌های نمونه‌گیری مجموعه رتبه‌دار یک برآوردگر ناریب برای میانگین جامعه مورد بررسی بدست می‌دهد که این برآوردگر کاراتر از برآوردگر نمونه‌گیری تصادفی ساده با اندازه نمونه برابر می‌باشد.

در این مقاله برآوردگرهای عادی و رگرسیونی روش‌های نمونه‌گیری مجموعه رتبه‌دار ساده، میانه‌ای و حدی مورد بررسی قرار می‌گیرند. سپس مقدار کل محصول گندم ایران را توسط این روشها برآورد کرده و نشان می‌دهیم که این برآوردگرها نسبت به برآوردگر تصادفی ساده از دقت بیشتری برخوردار می‌باشند.

واژه‌های کلیدی: آماره مرتب، نمونه‌گیری مجموعه رتبه‌دار، برآوردگر رگرسیونی

۱ مقدمه

روش‌های نمونه‌گیری مجموعه رتبه‌دار یک برآوردگر ناریب برای میانگین جامعه مورد بررسی بدست می‌دهند که این برآوردگر دارای واریانس کوچکتر از برآوردگر نمونه تصادفی ساده با اندازه نمونه برابر می‌باشد. علاوه بر آن در بسیاری از کارهای عملی اندازه‌گیری تک تک واحدهای نمونه مشکل یا پرهزینه می‌باشد، اما واحدها می‌توانند به آسانی و بی هیچ هزینه‌ای رتبه‌بندی شوند و واحدهای رتبه‌بندی شده خاصی اندازه‌گیری گردند. در مسائل زیست محیطی و کشاورزی رتبه‌بندی واحدها بدون اینکه مقدار واقعی آنها اندازه‌گیری شود امکان‌پذیر است. در بعضی موارد (مثل اندازه‌گیری میزان آلودگی خاک یک منطقه) بخش اعظم هزینه مربوط به تجزیه و تحلیل و کارهای آزمایشگاهی میباشد، درحالیکه تشخیص قابلیت (پتانسیل) واحدهای نمونه کار ساده‌ای است. در این حالت می‌توانیم از روش نمونه‌گیری مجموعه رتبه‌دار استفاده کنیم، یعنی می‌توانیم

تعداد زیادی از واحدهای جامعه را انتخاب کنیم و فقط زیر نمونه‌ای از واحدهای انتخاب شده را از نظر کمیت مورد نظر، اندازه‌گیری کنیم. همچنین استفاده از روش نمونه‌گیری مجموعه رتبه‌دار باعث طبقه‌بندی جامعه به طبقات همگن می‌شود. پس با روش نمونه‌گیری مجموعه رتبه‌دار علاوه بر افزایش دقت و کارایی می‌توان جامعه را بدون اینکه از ساختار جامعه اطلاعی در دست باشد، طبقه‌بندی کرد (و از تمام طبقات نمونه انتخاب کرد) که همین امر باعث افزایش دقت برآوردگر میانگین جامعه می‌شود.

روش نمونه‌گیری مجموعه رتبه‌دار اولین بار توسط مک‌انتایر^۱ (۱۹۵۲) ارائه شد. در سال ۱۹۶۷ تاکاهاسی^۲ و واکیما^۳ نشان دادند که میانگین حسابی حاصل از روش نمونه‌گیری مجموعه رتبه‌دار یک برآوردگر ناریب برای میانگین جامعه و کاراتر از روش نمونه‌گیری تصادفی ساده با اندازه نمونه برابر می‌باشد. در سال ۱۹۷۷ استوکس^۴ کاربرد متغیر همراه در روش نمونه‌گیری مجموعه رتبه‌دار را مطرح نمود. در سال ۱۹۹۳ پتیل^۵، سینها^۶ و تایل^۷ برآوردگر رگرسیونی نمونه‌گیری تصادفی ساده براساس متغیر همراه را با برآوردگر عادی نمونه‌گیری مجموعه رتبه‌دار مورد مقایسه قرار دادند. در سال ۱۹۹۷ یو^۸ و لام^۹ برآوردگر رگرسیونی براساس نمونه‌گیری مجموعه رتبه‌دار را در نظر گرفته و آن را با برآوردگر عادی نمونه‌گیری مجموعه رتبه‌دار و برآوردگر رگرسیونی نمونه‌گیری تصادفی ساده مقایسه نمودند. در سال ۱۹۹۶ سماوی^{۱۰}، ابودایه^{۱۱} و احمد^{۱۲} روش نمونه‌گیری مجموعه رتبه‌دار حدی را معرفی کردند. در سال ۱۹۹۷ و ۱۹۹۸ مطلق^{۱۳} روش نمونه‌گیری مجموعه رتبه‌دار میانه‌ای را با و بدون استفاده از متغیر همراه ارائه کرد و در سال ۲۰۰۱ برآوردگرهای رگرسیونی چندین روش نمونه‌گیری مجموعه رتبه‌دار را برای برآورد میانگین جامعه بکار برد و آنها را با برآوردگرهای عادی روش نمونه‌گیری مجموعه رتبه‌دار مورد مقایسه قرار داد.

در این مقاله ضمن بیان روش‌های نمونه‌گیری مجموعه رتبه‌دار ساده، میانه‌ای و حدی و ارایه برآوردگرهای عادی و رگرسیونی در این روشها، مقدار کل محصول گندم ایران را با استفاده از آنها برآورد و با یکدیگر مقایسه می‌کنیم. بر این اساس در بخش ۲ برآوردگرهای عادی و در بخش ۳ برآوردگرهای رگرسیونی را در روش‌های نمونه‌گیری مجموعه رتبه‌دار ساده، میانه‌ای و حدی را مورد بررسی قرار داده و خواص آنها را بیان می‌کنیم (برای اثبات این خواص به مراجع فوق مراجعه شود). در بخش ۴ مقدار کل محصول گندم ایران را با استفاده از برآوردگرهای فوق برآورد و با یکدیگر مقایسه می‌کنیم.

1) McIntyre 2) Takahasi 3) Wakimoto 4) Stokes 5) Patil 6) Sinha
7) Taili 8) Yu 9) Lam 10) Samawi 11) Abu Dayyeh 12) Ahmad
13) Muttalak

۲ برآوردگرهای عادی در روش‌های نمونه‌گیری مجموعه رتبه‌دار

۱.۲ روش نمونه‌گیری مجموعه رتبه‌دار^{۱۴} ساده

در روش نمونه‌گیری مجموعه رتبه‌دار ساده برای انتخاب یک نمونه n تایی، ابتدا یک نمونه $n^۲$ انتخاب شده و به n زیر مجموعه n تایی تقسیم می‌شوند و اعضای هر مجموعه بر اساس صفت مورد نظر رتبه‌بندی می‌شوند. سپس از نمونه n تایی اول واحد دارای کوچکترین رتبه، از نمونه n تایی دوم واحد دارای دومین رتبه، ... و از نمونه n تایی n ام واحد دارای بزرگترین رتبه اندازه‌گیری می‌شود.

فرض کنید $Y_{۱۱}, Y_{۱۲}, \dots, Y_{۱n}; Y_{۲۱}, Y_{۲۲}, \dots, Y_{۲n}; \dots; Y_{n۱}, Y_{n۲}, \dots, Y_{nn}$ واحد انتخاب شده از جامعه باشند و فرض می‌شود که متغیرهای تصادفی مستقل با تابع توزیع یکسان $F(Y)$ و تابع چگالی یکسان $f(y)$ بوده و $Y_{i[۱]}, Y_{i[۲]}, \dots, Y_{i[n]}$ آماره‌های مرتب نمونه $Y_{i۱}, Y_{i۲}, \dots, Y_{in}$ برای $i = ۱, ۲, \dots, n$ باشند. نمونه اندازه‌گیری شده n تایی از این $n^۲$ واحد انتخابی عبارتند از $Y_{۱[۱]}, Y_{۲[۲]}, \dots, Y_{n[n]}$ که متغیرهای تصادفی مستقل می‌باشند. میانگین واحدهای اندازه‌گیری شده عبارت است از

$$\bar{Y}_{RSS} = \bar{Y}_{[n]} = \frac{1}{n} \sum_{i=1}^n Y_{i[i]} \quad (۱)$$

هنگامیکه تعداد نمونه زیاد باشد n (بزرگتر از ۳ یا ۴) رتبه‌بندی عناصر به راحتی امکان‌پذیر نیست. در چنین حالتی برای رتبه‌بندی راحتتر و ایجاد خطای کمتر، عموماً روش نمونه‌گیری مجموعه رتبه‌دار با اندازه نمونه کوچکتر انجام گرفته و سپس روش نمونه‌گیری چند بار تکرار می‌شود. مثلاً برای انتخاب یک نمونه n تایی که $n = nm$ ، روش نمونه‌گیری مجموعه رتبه‌دار n تایی، m بار تکرار می‌شود، که در آن n تعداد کل نمونه انتخابی مورد نظر می‌باشد. اگر m تکرار وجود داشته باشد در این صورت m مشاهده برای \bar{Y} بدست می‌آید و میانگین کل بصورت زیر خواهد بود.

$$\bar{Y}_{RSS} = \bar{Y}_{[n]} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m Y_{i[i]j} \quad (۲)$$

که به آن برآوردگر روش نمونه‌گیری مجموعه رتبه‌دار ساده گویند. خواص این برآوردگر در قضیه زیر بیان شده است.

قضیه ۱.۲ برآوردگر \bar{Y}_{RSS} یک برآوردگر ناریب برای میانگین جامعه است و واریانس آن با m بار تکرار عبارت است از

$$\sigma^۲(\bar{Y}_{[n]m}) = Var(\bar{Y}_{[n]m}) = \frac{\sigma_{[n]}^۲}{nm} = \frac{\sigma^۲}{nm} - \frac{1}{n^۲m} \sum_{k=1}^n (\mu_{n,k[k]} - \mu)^۲$$

و اگر \bar{Y}_{SRS} میانگین نمونه تصادفی ساده n تایی باشد آنگاه همواره $Var(\bar{Y}_{RSS}) \leq Var(\bar{Y}_{SRS})$ است.

دلیل اینکه میانگین مجموعه رتبه‌دار با میانگین نمونه تصادفی ساده حاصل از یک نمونه n تایی و نه n^2 تایی مقایسه می‌شود این است که در وضعیت مورد مطالعه ما، هزینه رتبه‌بندی عناصر انتخابی مورد نظر نیست و تنها هزینه عناصر اندازه‌گیری شده مورد نظر است.

۲.۲ روش نمونه‌گیری مجموعه رتبه‌دار میانه‌ای

در روش نمونه‌گیری مجموعه رتبه‌دار میانه‌ای n نمونه n تایی (نمونه n^2 تایی) به‌طور تصادفی از جامعه انتخاب می‌شود. واحدهای داخل هر نمونه بر اساس متغیر مورد نظر رتبه‌بندی می‌شود. اگر تعداد نمونه فرد باشد، از هر نمونه n تایی رتبه‌بندی شده، $\frac{n+1}{2}$ امین واحد (میانه نمونه) انتخاب می‌شود. اگر اندازه نمونه زوج باشد از $\frac{n}{2}$ نمونه اول $\frac{n}{2}$ امین واحد رتبه‌بندی شده و از $\frac{n}{2}$ نمونه بعدی $\frac{n+2}{2}$ امین واحد رتبه‌بندی شده، انتخاب می‌شود.

فرض کنید $Y_{11}, Y_{12}, \dots, Y_{1n}; Y_{21}, Y_{22}, \dots, Y_{2n}; \dots; Y_{n1}, Y_{n2}, \dots, Y_{nn}$ متغیرهای تصادفی مستقل باشند که دارای تابع توزیع $F(y)$ و تابع چگالی $f(y)$ باشند و اگر $Y_{i[1]}, Y_{i[2]}, \dots, Y_{i[n]}$ آماره‌های مرتب $Y_{i1}, Y_{i2}, \dots, Y_{in}$ باشند. اگر تعداد نمونه فرد باشد از هر نمونه $Y_{i[\frac{n+1}{2}]}$ یعنی میانه i امین نمونه $(i = 1, 2, \dots, n)$ انتخاب می‌شود، به عبارت دیگر $\frac{n+1}{2}$ امین آماره مرتب از i امین نمونه انتخاب می‌شود. اگر اندازه نمونه زوج باشد از $\frac{n}{2}$ نمونه‌ها، $\frac{n}{2}$ امین آماره مرتب یعنی $Y_{i[\frac{n}{2}]}$ به ازای $(i = 1, 2, \dots, \frac{n}{2})$ و از $\frac{n}{2}$ نمونه بعدی $\frac{n+2}{2}$ امین آماره مرتب یعنی $Y_{i[\frac{n+2}{2}]}$ به ازای $i = \frac{n}{2} + 1, \dots, n$ انتخاب می‌شود. به عبارت دیگر $Y_{i[\frac{n}{2}]}$ ، $i = 1, 2, \dots, \frac{n}{2}$ و $Y_{i[\frac{n+2}{2}]}$ ، $i = \frac{n}{2} + 1, \dots, n$ نشان دهنده نمونه‌های اندازه‌گیری شده روش نمونه‌گیری مجموعه رتبه‌دار میانه‌ای خواهند بود. اگر تعداد نمونه فرد باشد میانگین مجموعه رتبه‌دار میانه‌ای به صورت زیر خواهد بود

$$\bar{Y}_{MRSS1} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m Y_{i[\frac{n+1}{2}]} \quad (3)$$

و در حالتی که تعداد نمونه زوج باشد میانگین عبارت است از

$$\bar{Y}_{MRSS2} = \frac{1}{n} \left(\sum_{i=1}^{\frac{n}{2}} Y_{i[\frac{n}{2}]} + \sum_{i=\frac{n}{2}+1}^n Y_{i[\frac{n+2}{2}]} \right) \quad (4)$$

در قضیه زیر خواص این دو برآوردگر بیان شده‌اند.

قضیه ۲.۲ اگر \bar{Y}_{MRSS2} و \bar{Y}_{MRSS1} میانگینهای مجموعه رتبه‌دار میانه‌ای باشند آنگاه الف) \bar{Y}_{MRSS2} و \bar{Y}_{MRSS1} یک برآوردگر ناریب برای میانگین جامعه (μ) می‌باشند اگر توزیع جامعه حول μ متقارن باشد. ب) $Var(\bar{Y}_{MRSS1})$ و $Var(\bar{Y}_{MRSS2})$ از واریانس نمونه‌گیری تصادفی ساده $Var(\bar{Y})$ کوچکتر است. ج) اگر توزیع حول میانگین μ متقارن نباشد، میانگین مربع خطا (MSE)، \bar{Y}_{MRSS2} و \bar{Y}_{MRSS1} از واریانس نمونه‌گیری تصادف ($Var(\bar{Y}_{SRS})$) کوچکتر خواهد بود.

۳.۲ روش نمونه‌گیری مجموعه رتبه‌دار حدی

روش نمونه‌گیری مجموعه رتبه‌دار حدی به این صورت است که برای انتخاب یک نمونه n تایی، n نمونه n تایی (نمونه n^2 تایی) از جامعه‌ای که باید میانگین آن برآورد شود، انتخاب می‌شود. فرض می‌شود که بزرگترین و کوچکترین واحدهای این مجموعه می‌تواند با یک نگاه اجمالی تعیین شوند. این روش ساده و یک فرآیند عملی می‌باشد. از اولین مجموع n تایی کوچکترین آماره مرتب واحدها، از دومین مجموع n تایی دوم بزرگترین واحدها رتبه‌بندی شده و از مجموع n تایی سوم، کوچکترین واحد رتبه‌بندی شده و به همین ترتیب از $(n-1)$ امین مجموعه بزرگترین واحد رتبه‌بندی شده اندازه‌گیری می‌شوند. برای انتخاب n امین واحد از n امین مجموعه (مجموعه آخری) بسته به این که تعداد نمونه فرد باشد یا زوج به صورت زیر انتخاب می‌شود. الف) اگر تعداد نمونه n زوج باشد، بزرگترین واحد رتبه‌بندی شده، اندازه‌گیری می‌شود. این روش نمونه‌گیری را با $ERSS_a$ نشان می‌دهند. اگر تعداد نمونه فرد باشد به دو صورت زیر می‌توان نمونه را انتخاب کرد. ب) برای اندازه‌گیری n امین واحد، میانگین اندازه‌گیری شده کوچکترین و بزرگترین واحد n امین مجموعه را در نظر می‌گیرند. این روش نمونه‌گیری را با $ERSS_b$ نشان می‌دهند. ج) برای اندازه‌گیری n امین واحد، مقدار میانه n امین مجموعه در نظر گرفته می‌شود. این روش نمونه‌گیری را با $ERSS_c$ نشان می‌دهند. برآوردگرهای حاصل از روش نمونه‌گیری مجموعه رتبه‌دار حدی به صورت زیر خواهند بود.

$$\bar{Y}_{ERSS_a} = \bar{Y}_a = \frac{1}{\frac{n}{2}} \left(\frac{1}{\frac{n}{2}} \sum_{i=1}^{\frac{n}{2}} Y_{r_{i-1}[n]} + \frac{1}{\frac{n}{2}} \sum_{i=1}^{\frac{n}{2}} Y_{r_i[n]} \right) \quad (5)$$

$$\bar{Y}_{ERSS_b} = \bar{Y}_b = \frac{Y_{1[n]} + Y_{2[n]} + \dots + Y_{n-1[n]} + \frac{1}{2} [Y_{n[n]} + Y_{n[n]}]}{n} \quad (6)$$

$$\bar{Y}_{ERSS_c} = \bar{Y}_c = \frac{Y_{1[n]} + Y_{2[n]} + \dots + Y_{n-1[n]} + Y_{n[\frac{n+1}{2}]}}{n} \quad (7)$$

قضیه ۳.۲ وقتی تعداد نمونه n زوج و توزیع جامعه متقارن باشد، \bar{Y}_{ERSS_a} یک برآوردگر ناریب برای میانگین جامعه با واریانس زیر است.

$$Var(\bar{Y}_{ERSS_a}) = \frac{1}{2n}(\sigma_{n, \lfloor n/2 \rfloor}^2 + \sigma_{n, n/2}^2)$$

که در آن $\sigma_{n, \lfloor n/2 \rfloor}^2 = Var(\bar{Y}_{n[\lfloor n/2 \rfloor]})$ و $\sigma_{n, n/2}^2 = Var(\bar{Y}_{n[n/2]})$. همچنین اگر تعداد نمونه فرد باشد $ERSS_b$ و $ERSS_c$ یک برآوردگر ناریب برای میانگین جامعه است و واریانس این برآوردگرهای به ترتیب عبارت است از

$$Var(\bar{Y}_{ERSS_b}) = \frac{1}{4n^2}[(2n-1)(\sigma_{n, \lfloor n/2 \rfloor}^2 + \sigma_{n, n/2}^2) + 2\sigma_{(1, n)}^2] \quad (8)$$

$$Var(\bar{Y}_{ERSS_c}) = \frac{n-1}{2n^2}[(\sigma_{n, \lfloor n/2 \rfloor}^2 + \sigma_{n, n/2}^2)] + \frac{1}{n^2}\sigma_{n, n[\frac{n+1}{2}]}^2 \quad (9)$$

که $\sigma_{(1, n)} = Cov(Y_{n[\lfloor n/2 \rfloor]}, Y_{n[n/2]})$ است.

۳ برآوردگرهای رگرسیونی در روش‌های نمونه‌گیری مجموعه رتبه‌دار

هنگامیکه رتبه‌بندی متغیر اصلی Y واحدهای نمونه مشکل یا پر هزینه است اما یک متغیر همراه مانند X وجود دارد که با متغیر اصلی رابطه قوی دارد و به آسانی قابل رتبه‌بندی است، می‌توان از متغیر همراه X برای رتبه‌بندی استفاده کرد. در این روش همه عناصر انتخاب شده بر اساس متغیر همراه رتبه‌بندی می‌شوند و فقط متغیر اصلی Y متناظر با بعضی از آنها اندازه‌گیری خواهند شد. عناصر متناظر اندازه‌گیری شده، نمونه نهایی را تشکیل می‌دهند. حال به وسیله این عناصر اندازه‌گیری شده می‌توان برآوردگر عادی و یا رگرسیونی را بدست آورد که در این مقاله تنها به بررسی برآوردگرهای رگرسیونی در دو حالت معلوم یا نامعلوم بودن میانگین جامعه متغیر همراه خواهیم پرداخت.

۱.۳ برآوردگرهای رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار با معلوم بودن میانگین جامعه متغیر همراه

برای انتخاب یک نمونه n تایی از متغیر اصل Y با استفاده از متغیر همراه X به صورت زیر عمل می‌کنیم. n نمونه دو متغیره (زوجی) n تایی (نمونه n^2 تایی) از جامعه مورد مطالعه انتخاب شده، نمونه‌های زوجی بر اساس متغیر همراه X رتبه‌بندی می‌شوند. سپس بوسیله یکی

از روش‌های نمونه‌گیری مجموعه رتبه‌دار متغیر اصلی متناظر با متغیر همراه انتخاب و اندازه‌گیری می‌شوند. این عمل m بار تکرار می‌گردد تا یک نمونه nm تایی حاصل گردد.

فرض می‌شود که رگرسیون Y روی X خطی بوده $\frac{X-\mu_X}{\sigma_X}$ و $\frac{Y-\mu_Y}{\sigma_Y}$ دارای توزیع یکسان باشند. برای مثال هر دو دارای توزیع نرمال دو متغیر باشند. در این صورت رگرسیون Y روی X بصورت

$$Y = \mu_Y + \frac{\rho\sigma_Y}{\sigma_X}(X - \mu_X) + \epsilon \quad (10)$$

است که در آن X و ϵ مستقل ϵ دارای میانگین صفر و واریانس $(1 - \rho^2)\sigma_Y^2$ می‌باشد. فرض کنید $X_{k[k]j}$ و $Y_{k(k)j}$ نمونه‌های انتخابی مجموعه رتبه‌دار ساده با استفاده از متغیر همراه باشند که در آن k امین واحد رتبه‌بندی شده متغیر همراه در k امین نمونه و در j امین تکرار و $Y_{k(k)j}$ متغیر اصلی اندازه‌گیری شده متناظر با آن است. همچنین فرض کنید رگرسیون Y روی X خطی باشد. در این صورت رابطه (۱۰) را می‌توان بصورت زیر نوشت

$$Y_{k(k)j} = \mu_Y + \frac{\rho\sigma_Y}{\sigma_X}(X_{k[k]j} - \mu_X) + \epsilon_{kj} \quad k = 1, 2, \dots, n, j = 1, 2, \dots, m \quad (11)$$

میانگین نمونه متغیرهای اصلی انتخابی با nm تکرار برابر خواهد بود با $\bar{Y}_{RSSC} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m Y_{k(k)j}$

هنگامیکه میانگین متغیر همراه μ_X معلوم باشد، برآوردگر تفاضلی میانگین جامعه متغیر اصلی Y با به کارگیری روش مجموعه رتبه‌دار (ساده) بصورت زیر خواهد بود

$$\bar{Y}_{RSSD} = \bar{Y}_{RSSC} + \beta_1(\mu_X - \bar{X}_{RSS}) \quad (12)$$

که در آن β_1 نامعلوم است و یک برآوردگر برای آن عبارت است از

$$\hat{\beta}_1 = \frac{\sum_{k=1}^n \sum_{j=1}^m (X_{k[k]j} - \bar{X}_{RSS})(Y_{k(k)j} - \bar{Y}_{RSSC})}{\sum_{k=1}^n \sum_{j=1}^m (X_{k[k]j} - \bar{X}_{RSS})^2} \quad (13)$$

در نتیجه برآوردگر رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار ساده برای میانگین جامعه متغیر اصلی (μ_Y) به صورت زیر خواهد بود

$$\bar{Y}_{RSSLR} = \bar{Y}_{RSSC} + \hat{\beta}_1(\mu_X - \bar{X}_{RSS}) \quad (14)$$

که خواص آن در قضیه زیر بیان شده است.

قضیه ۱.۳ اگر رگرسیون Y روی X خطی و بفرم (۱۱) باشد آنگاه برآوردگر رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار (ساده) یک برآوردگر نااریب برای میانگین جامعه است و واریانس آن عبارت است از

$$Var(\bar{Y}_{RSSLR}) = \frac{\sigma_y^2}{nm} (1 - \rho^2) [1 + E(\frac{\bar{Z}_{RSS}^2}{s_{z_1}^2})]$$

که در آن $\bar{Z}_{RSS} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^m Z_{k[k]j}$ ، $Z_{k[k]j} = \frac{X_{k[k]j} - \mu_x}{\sigma_x}$

$$S_{z_1}^2 = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^m ((Z_{k[k]j} - \bar{Z}_{RSS})^2)$$

حال فرض کنید $X_{k[m]j}$ و $Y_{k(m)j}$ به ترتیب نشان دهنده میانه X و متغیر اصلی متناظر بدست آمده از k امین نمونه n ، $k = 1, 2, \dots, n$ در j امین تکرار در صورت فرد بودن اندازه نمونه باشد و همچنین $\frac{n}{2}$ امین آماره مرتب از k امین نمونه به ازای $k = 1, \dots, \frac{n}{2}$ و $\frac{n+1}{2}$ امین آماره مرتب از k امین نمونه به ازای $k = \frac{n}{2} + 1, \dots, n$ در j امین تکرار در صورت زوج بودن اندازه نمونه باشد. اگر رگرسیون Y روی X خطی باشد در این صورت رگرسیون خطی $Y_{k(m)j}$ روی $X_{k[m]j}$ با توجه به رابطه (۱۰) عبارت است از

$$Y_{k(m)j} = \mu_Y + \frac{\rho\sigma_Y}{\sigma_X}(\mu_X - X_{k[m]j}) + \epsilon_{kj} \quad (15)$$

$k = 1, 2, \dots, n \quad j = 1, 2, \dots, m$

اگر میانگین متغیر همراه معلوم باشد، برآوردگر تفاضلی میانگین جامعه متغیر اصلی Y با به کارگیری روش نمونه‌گیری مجموعه رتبه‌دار میانه‌ای به صورت زیر خواهد بود

$$\bar{Y}_{MRSSD} = \bar{Y}_{MRSSC} + \beta_2(\bar{X}_{MRSS} - \mu_X) \quad (16)$$

که در آن

$$\bar{Y}_{MRSSC} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^m Y_{k(m)j}$$

$$\bar{X}_{MRSS} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^m X_{k[m]j}$$

و β_2 یک مقدار ثابت و نامعلوم است و یک برآوردگر برای آن عبارت است از

$$\hat{\beta}_2 = \frac{\sum_{k=1}^n \sum_{j=1}^m (X_{k[m]j} - \bar{X}_{MRSS})(Y_{k(m)j} - \bar{Y}_{MRSSC})}{\sum_{k=1}^n \sum_{j=1}^m (X_{k[m]j} - \bar{X}_{MRSS})^2} \quad (17)$$

بنابراین برآوردگر رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار میانه‌ای به صورت زیر خواهد بود

$$\bar{Y}_{MRSSD} = \bar{Y}_{MRSSC} + \hat{\beta}_2(\mu_X - \bar{X}_{MRSS}) \quad (18)$$

که خواص آن در قضیه زیر بیان شده است.

قضیه ۲.۳ اگر رگرسیون Y روی X بفرم (۱۱) باشد برآوردگر رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار میانه‌ای یک برآوردگر ناریب برای میانگین جامعه و دارای واریانس بصورت زیر است

$$Var(\bar{Y}_{MRSSLR}) = \frac{\sigma_y^2}{nm} (1 - \rho^2) [1 + E(\frac{\bar{Z}_{MRSS}^2}{s_{zr}^2})]$$

که در آن $\bar{Z}_{MRSS} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^n Z_{k[m]j}$ ، $Z_{k[m]j} = \frac{X_{k[m]j} - \mu_x}{\sigma_x}$

$$s_{zr}^2 = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^n ((Z_{k[k]j} - \bar{Z}_{MRSS})^2)$$

حال فرض کنید $Y_{k(e)j}$ و $X_{k[e]j}$ به ترتیب نشان دهنده کوچکترین مقدار X و متغیر اصلی متناظر بدست آمده از k امین نمونه $k = 1, 2, \dots, \frac{n}{p}$ و بزرگترین واحد رتبه‌بندی شده X و متغیر اصلی متناظر بدست آمده از $k = \frac{n}{p} + 1, \dots, n$ در k ژامین تکرار در صورت زوج بودن اندازه نمونه باشد. همچنین این نماد را برای نشان دادن کوچکترین واحد رتبه‌بندی شده و مقدار متغیر اصلی متناظر با آن در k امین نمونه $k = 1, 2, \dots, \frac{n-1}{p}$ و میان X و متغیر اصلی متناظر با آن در k امین نمونه $k = \frac{n+1}{p}$ و بزرگترین واحد رتبه‌بندی شده X و متغیر اصلی متناظر با آن در k امین نمونه $k = \frac{n+1}{p} + 1, \dots, n$ و ژامین تکرار، در صورتی که اندازه نمونه فرد باشد را بکار می‌بریم.

در این صورت برآوردگر رگرسیونی نمونه‌گیری مجموعه رتبه‌دار حدی برای μ_Y عبارت است از

$$\bar{Y}_{ERRSSLR} = \bar{Y}_{ERSSC} + \hat{\beta}_3 (\mu_X - \bar{X}_{ERSS}) \quad (19)$$

که $\bar{X}_{ERSS} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^m X_{k[e]j}$ ، $\bar{Y}_{ERSSC} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^m Y_{k(e)j}$ و $\hat{\beta}_3$ برابر است با

$$\hat{\beta}_3 = \frac{\sum_{k=1}^n \sum_{j=1}^m (X_{k[e]j} - \bar{X}_{ERSS})(Y_{k(e)j} - \bar{Y}_{ERSSC})}{\sum_{k=1}^n \sum_{j=1}^m (X_{k[e]j} - \bar{X}_{ERSS})^2} \quad (20)$$

خواص $\bar{Y}_{ERRSSLR}$ در قضیه زیر بیان شده است.

قضیه ۳.۳ اگر رگرسیون Y روی X خطی باشد آنگاه برآوردگر رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار حدی یک برآوردگر ناریب برای میانگین جامعه است و واریانس آن عبارت است از

$$Var(\bar{Y}_{ERRSSLR}) = \frac{\sigma_y^2}{nm} (1 - \rho^2) [1 + E(\frac{\bar{Z}_{ERSS}^2}{s_{zr}^2})]$$

که در آن $\bar{Z}_{ERSS} = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^n Z_{k[e]j}$ ، $Z_{k[e]j} = \frac{X_{k[e]j} - \mu_x}{\sigma_x}$ و

$$S_{z_r}^2 = \frac{1}{nm} \sum_{k=1}^n \sum_{j=1}^n ((Z_{k[e]j} - \bar{Z}_{ERSS})^2)$$

۲.۳ برآوردگرهای رگرسیونی در نمونه‌گیری مجموعه رتبه‌دار با نامعلوم بودن میانگین متغیر کمکی

در بیشتر کاربردها میانگین جامعه متغیر کمکی نامعلوم است و با استفاده از روش نمونه‌گیری دوگانه یا دو فازی^{۱۵} برآورد می‌شود. برای این کار در فاز اول یک نمونه n^2m تایی به صورت تصادفی انتخاب می‌شود و از آن برای برآورد میانگین جامعه متغیر کمکی μ_X استفاده می‌شود. سپس از این نمونه n^2m تایی یک زیر نمونه nm تایی در فاز دوم به صورت تصادفی ساده یا نمونه‌گیری مجموعه رتبه‌دار انتخاب می‌شود که برای مطالعه خصیصه Y استفاده می‌شود. فرض کنید \bar{X}_d میانگین نمونه X بر اساس مشاهدات n^2m تایی X در فاز اول باشد. در این صورت \bar{X}_d یک برآوردگر ناریب برای μ_X است. اگر نمونه‌گیری تصادفی ساده در فاز دوم نمونه‌گیری استفاده شود، برآوردگر رگرسیونی با استفاده از نمونه‌گیری دو گانه به صورت زیر خواهد بود

$$\bar{Y}_{LRD} = \bar{y} + \hat{\beta}(\bar{X}_d - \bar{x}) \quad (21)$$

قضیه ۴.۳ \bar{Y}_{LRD} یک برآوردگر ناریب برای میانگین جامعه است و اگر (X, Y) دارای توزیع نرمال دو متغیره باشند آنگاه واریانس آن عبارت است از

$$Var(\bar{Y}_{LRD}) = \frac{\sigma_y^2}{nm} (1 - \rho^2) \left[1 + \frac{n-1}{n} \frac{1}{nm-3} \right] + \frac{\rho^2 \sigma_y^2}{n^2 m} \quad (22)$$

فرض کنید نمونه‌گیری مجموعه رتبه‌دار ساده در دومین فاز از نمونه‌گیری دوگانه استفاده شود. یو ولام (۱۹۹۷) برآوردگر رگرسیونی میانگین جامعه μ_Y را به صورت زیر ارائه کردند

$$\bar{Y}_{RSSLRD} = \bar{Y}_{RSSC} + \hat{\beta}_1 (\bar{X}_d - \bar{X}_{RSS}) \quad (23)$$

که $\hat{\beta}_1$ توسط رابطه (۱۳) بدست می‌آید.

قضیه ۵.۳ برآوردگر \bar{Y}_{RSSLRD} یک برآوردگر ناریب برای میانگین جامعه μ_Y است و واریانس آن عبارت است از

$$Var(\bar{Y}_{RSSLRD}) = \frac{\sigma_y^2}{nm} (1 - \rho^2) \left[1 + E \left(\frac{(\bar{Z}_{RSS} - \bar{Z})^2}{s_{z_1}^2} \right) \right] + \frac{\rho^2 \sigma_y^2}{n^2 m} \quad (24)$$

که $\bar{Z} = \frac{\bar{X}_d - \mu_X}{\sigma_X}$ ، \bar{Z}_{RSS} و $S_{Z_1}^2$ در قضیه ۱.۳ معرفی شده‌اند. اگر نمونه‌گیری مجموعه رتبه‌دار میانه‌ای در دومین فاز از نمونه‌گیری دوگانه استفاده شود، برآوردگر رگرسیونی میانگین جامع μ_Y به صورت زیر خواهد بود.

$$\bar{Y}_{MRSSLRD} = \bar{Y}_{MRSSC} + \hat{\beta}_2(\bar{X}_d - \bar{X}_{MRSS}) \quad (25)$$

که $\hat{\beta}_2$ توسط رابطه (۱۷) به دست می‌آید. همانند قبل برآوردگر $\bar{Y}_{MRSSLRD}$ یک برآوردگر ناریب برای میانگین جامعه μ_Y با واریانس زیر است

$$Var(\bar{Y}_{MRSSLRD}) = \frac{\sigma_y^2}{nm}(\lambda - \rho^2)[\lambda + E(\frac{(\bar{Z}_{MRSS} - \bar{Z})^2}{s_{z_1}^2})] + \frac{\rho^2 \sigma_y^2}{n^2 m} \quad (26)$$

که \bar{Z}_{MRSS} و $S_{Z_1}^2$ در قضیه ۲.۳ معرفی شده‌اند. اگر نمونه‌گیری مجموعه رتبه‌دار حدی در دومین فاز از نمونه‌گیری دوگانه استفاده شود، برآوردگر رگرسیونی میانگین جامعه μ_Y به صورت زیر خواهد بود

$$\bar{Y}_{ERSSLRD} = \bar{Y}_{ERSSC} + \hat{\beta}_3(\bar{X}_d - \bar{X}_{ERSS}) \quad (27)$$

که $\hat{\beta}_3$ توسط رابطه (۲۰) به دست می‌آید. همانند قبل برآوردگر $\bar{Y}_{ERSSLRD}$ یک برآوردگر ناریب برای میانگین جامعه μ_Y با واریانس زیر است

$$Var(\bar{Y}_{ERSSLRD}) = \frac{\sigma_y^2}{nm}(\lambda - \rho^2)[\lambda + E(\frac{(\bar{Z}_{ERSS} - \bar{Z})^2}{s_{z_3}^2})] + \frac{\rho^2 \sigma_y^2}{n^2 m} \quad (28)$$

که \bar{Z}_{ERSS} و $S_{Z_3}^2$ در قضیه ۳.۳ معرفی شده‌اند.

۴ کاربرد روش‌های مختلف نمونه‌گیری مجموعه رتبه‌دار در برآورد مقدار کل محصول گندم ایران

برای نشان دادن چگونگی اجرای روش‌های نمونه‌گیری مجموعه رتبه‌دار، به عنوان کاربرد، میزان تولید گندم کشور را برآورد می‌کنیم. با استفاده از اطلاعات موجود در سالنامه آماری ۱۳۸۰ مرکز آمار، ابتدا کشور ایران را به شهرستانهای مختلف تقسیم می‌کنیم. این شهرستانها به عنوان واحدهای جامعه در نظر گرفته می‌شوند. متغیر مورد نظر در هر واحد (شهرستان) میزان تولید محصول گندم در آن واحد (شهرستان) است و در جاهایی که متغیر همراه لازم است سطح زیر کشت گندم آن شهرستان به عنوان متغیر همراه در نظر گرفته می‌شود. میزان تولید محصول گندم بر

حساب تن و سطح زیر کشت آن بر حسب هکتار می‌باشد. جامعه آماری مورد نظر شهرستانهای ایران است و حجم جامعه ۲۷۶ می‌باشد. همچنین تعداد نمونه انتخابی برای برآورد مقدار کل گندم ایران ۲۰ نمونه می‌باشد.

با به کارگیری روش‌های مختلف نمونه‌گیری مجموعه رتبه‌دار برآوردهای مقدار کل محصول گندم ایران بدست آمده که نتایج آن در جداول زیر ارائه شده است.

جدول ۱: دقت نسبی برآوردهای عادی روش نمونه‌گیری مجموعه رتبه‌دار نسبت به نمونه‌گیری تصادفی ساده

| روش نمونه‌گیری | خطای استاندارد برآوردها | برآورد مقدار کل گندم ایران | دقت نسبی |
|-------------------|-------------------------|----------------------------|----------|
| رتبه‌دار ساده | ۱,۶۹۷,۰۸۱ | ۸,۸۶۴,۵۲۷ | ۵,۲۸ |
| رتبه‌دار میانه‌ای | ۱,۱۳۶,۰۰۴ | ۷,۳۸۶,۹۲۱ | ۱۱,۷۹ |
| رتبه‌دار حدی | ۱,۰۵۸,۲۵۲ | ۸,۸۹۱,۲۵۹ | ۱۳,۵۹ |

جدول ۲: دقت نسبی برآوردهای رگرسیونی مجموعه رتبه‌دار با μ_X معلوم نسبت به نمونه‌گیری تصادفی ساده

| برآورد رگرسیونی | خطای استاندارد برآوردها | برآورد مقدار کل گندم | دقت نسبی |
|---------------------------------|-------------------------|----------------------|----------|
| رتبه‌دار ساده μ_X معلوم | ۱,۵۱۰,۴۸۷ | ۸,۴۴۷,۴۹۸ | ۶,۶۷ |
| رتبه‌دار میانه‌ای μ_X معلوم | ۱,۶۷۵,۷۴۸ | ۱۰,۴۷۹,۰۰۳ | ۵,۴۲ |
| رتبه‌دار حدی μ_X معلوم | ۱,۵۵۱,۱۳۳ | ۸,۶۲۷,۹۹۹ | ۶,۳۲ |
| ساده μ_X معلوم | ۳,۷۷۲,۹۳۱ | ۱۳,۴۵۳,۲۶۷ | ۱,۰۷ |

از جدول ۱ معلوم می‌شود که برآورد روش نمونه‌گیری مجموعه رتبه‌دار حدی در بین برآوردهای عادی روش نمونه‌گیری مجموعه رتبه‌دار کاراتر از بقیه برآوردهای دیگر است.

از جدول ۲ معلوم می‌شود که برآورد رگرسیونی مجموعه رتبه‌دار ساده با معلوم بودن میانگین متغیر کمکی در بین برآوردهای رگرسیونی روش نمونه‌گیری مجموعه رتبه‌دار با معلوم بودن میانگین متغیر کمکی کاراتر از بقیه برآوردهای دیگر است.

از جدول ۳ معلوم می‌شود که برآورد رگرسیونی مجموعه رتبه‌دار ساده با معلوم نبودن میانگین متغیر کمکی در بین برآوردهای رگرسیونی روش نمونه‌گیری مجموعه رتبه‌دار با معلوم نبودن میانگین متغیر کمکی کاراتر از بقیه برآوردهای دیگر است.

جدول ۳: دقت نسبی برآوردهای رگرسیونی مجموعه رتبه‌دار با معلوم نبودن μ_X نسبت به نمونه‌گیری تصادفی ساده

| دقت نسبی | برآورد مقدار کل گندم ایران | خطای استاندارد برآوردها | برآورد رگرسیونی |
|----------|----------------------------|-------------------------|-----------------------------------|
| ۶/۵۸ | ۸,۶۲۸,۸۵۶ | ۱,۵۲۰,۱۴۴ | رتبه‌دار ساده μ_X نامعلوم |
| ۵/۴۲ | ۱۰,۲۸۱,۱۱۶ | ۱,۶۷۵,۷۴۸ | رتبه‌دار میانه‌ای μ_X نامعلوم |
| ۶/۲۲ | ۹,۱۰۲,۷۱۰ | ۱,۵۶۳,۶۷۱ | رتبه‌دار حدی μ_X نامعلوم |
| ۱/۰۷ | ۱۴,۴۳۵,۶۴۹ | ۳,۷۷۹,۴۲۸ | تصادفی ساده μ_X نامعلوم |

جدول ۴: مقایسه دقت نسبی برآوردهای کارای روش‌های نمونه‌گیری مجموعه رتبه‌دار

| دقت نسبی | برآورد مقدار کل گندم ایران | خطای استاندارد برآوردها | روش نمونه‌گیری |
|----------|----------------------------|-------------------------|--|
| ۵۹/۱۳ | ۸,۸۹۱,۲۵۹ | ۱,۰۵۸,۲۵۲ | رتبه‌دار حدی |
| ۵۴/۵ | ۹,۸۷۵,۴۲۹ | ۱,۶۵۷,۱۲۷ | رتبه‌دار میان‌نمای با استفاده از متغیر همراه |
| ۶۷/۶ | ۸,۴۴۷,۴۹۸ | ۱,۵۱۰,۴۸۷ | رتبه‌دار ساده $UX/معلوم$ |
| ۵۸/۶ | ۸,۶۲۸,۸۵۶ | ۱,۵۲۰,۱۴۴ | رتبه‌دار ساده $UX/نامعلوم$ |

در آخر به مقایسه برآوردهای کارای بدست آمده از جداول بالا می‌پردازیم.
 از جدول ۴ معلوم می‌شود که برآورد روش نمونه‌گیری مجموعه رتبه‌دار حدی کاراتر از بقیه
 برآوردهای کارا می‌باشد.
 از جداول ۱ تا ۴ معلوم می‌شود که روش‌های مختلف نمونه‌گیری مجموعه رتبه‌دار همیشه
 کاراتر از روش نمونه‌گیری تصادفی ساده با تعداد نمونه برابر می‌باشند.

مراجع

- [1] McIntyre, G. A. (1952). A method of unbiased selective sampling, using ranked set. *Australian Journal of Agricultural Research*, 3, 385-390.
- [2] Muttlak, H. A. (1997). Median ranked set sampling, *Journal of Applied Statistical Sciences*, 6, 245-255.
- [3] Muttlak, H. A. (1998). Median ranked set sampling with concomitant variables and a comparison with ranked set sampling and regression estimators, *Environmetrics*, 9, 225-267.
- [4] Muttlak, H. A. (2001). Regression estimators in extreme and median ranked set samples, *Journal of Applied Statistical Sciences*, 28(8) 1003-1017.
- [5] Patil, G. P., Sinha, A. K. & Taillie, C. (1993). Relative precision of ranked set sampling: a comparison with ranked set sampling and regression estimators, *Environmetrics*, 4, 399-412.
- [6] Samawi, H., Abu-Dayyeh, W. & Ahmad, S. (1996). Extreme ranked set sampling, *The Biometrical Journal*, 30, 557-586.
- [7] Stokes, S. L. (1977). Ranked set sampling with concomitant variables, *Communication in Statistics*, AG, 1207-1211.
- [8] Takahasi, K. & Wakimoto, K. (1967). On unbiased estimates of the population mean based on the sample stratified by mean of ordering, *Annals of Statistical Mathematics*, 20, 1-31.
- [9] Yu, P.L. H. & Lam, K. (1997). Regression estimator in ranked set sampling, *Biometrics*, 53, 1070-1080.

انتشار داده‌های آماری و کنترل افشای اطلاعات فردی

حمیدرضا نواب‌پور^۱، محمد بردبار عشرت‌آبادی^۲

^۱ گروه آمار دانشگاه علامه طباطبائی

^۲ محقق پژوهشکده آمار

چکیده: با افزایش تقاضا برای آمارها در عرصه‌های جدید، مانند آمارهای زیست محیطی، آمارهای فقر، آمارهای جنسیتی، آمارهای فرهنگی و . . . ، سازمانهای ملی آمار اقدام به برنامه‌ریزی برای تولید و انتشار اینگونه آمارها علاوه بر آمارهایی که به صورت ادواری تولید و منتشر می‌کنند، کرده‌اند. برنامه‌ریزان، اقتصاددانان، و پژوهشگران همواره به داشتن داده‌های خرد تمایل داشته‌اند، زیرا این داده‌ها مانور آنها را در تحلیل پدیده‌های اجتماعی، اقتصادی و فرهنگی بیشتر می‌کند. اما طبق قانون، انتشار داده‌های جمع‌آوری شده نباید منجر به افشای هویت واحدهای اطلاع‌گیری شده، گردد. این امر همواره بهانه‌ای برای عدم انتشار داده‌های حاصل از آمارگیریها به نحوی که مورد تقاضای کاربران آنهاست، بوده است و یا در برخی موارد داده‌ها در جدولهای انتشاراتی ایمن نبوده و امکان افشای هویت واحدهای اطلاع‌گیری شده وجود داشته است. هدف این مقاله آرایه روشهایی است که به توسط آنها می‌توان هم تقاضای کاربران آمارهای تولید شده را اجابت کرد و هم از افشای اطلاعات فردی که خلاف قوانین عمومی آمار کشورها است جلوگیری کرد.

واژه‌های کلیدی: آمارهای رسمی، محرمانگی، داده‌های خرد، رکوردهای نایمن، متغیر سلسله مراتبی، جدول پیشابندی، فراداده

۱ مقدمه

امروزه مدیریت جامعه‌ها مبتنی بر اطلاعات است، لذا برنامه‌ریزان، اقتصاددانان، و پژوهشگران برای تصمیم‌سازی نیازمند اطلاعات آماری دقیق، روزآمد، و به موقع هستند. از این رو کشورها با تصویب قانون عمومی آمار، اقدام به تأسیس سازمانهای ملی آمار با وظیفه اصلی تولید و انتشار آمارهای رسمی نموده‌اند. معمولاً در قانونهای عمومی آمار احاد ملت ملزم به آرایه اطلاعات صحیح به پرسشگران هستند. در همین حال این قانونها به آنها اطمینان می‌دهد که اطلاعات آنها حفاظت شده و جز در تهیه آمارهای کلی به‌کار برده نمی‌شود. به عبارت دیگر این قانونها به پاسخگویان اطمینان می‌دهد که نحوه انتشار اطلاعات مربوط به افراد موجب افشای هویت آنها نخواهد شد. اصل محرمانگی اطلاعات فردی در همه نظامهای آماری پذیرفته شده است، لذا سازمانهای ملی آمار برای رعایت این اصل محدودیتهایی را برای انتشار داده‌های حاصل از

آمارگیریها اعمال می‌کنند به طوری که بعضا نیازهای پژوهشگران به برخی اطلاعات آماری تامین نمی‌شود و یا گاهی امکان افشای هویتها از جدولهای انتشاراتی سازمانهای ملی آمار وجود دارد. بنابراین در طول سه دهه گذشته ایجاد تعادل بین تأمین نیازهای آماری کاربران و جلوگیری از افشای هویت واحدهای اطلاع‌گیری شده، یکی از دغدغه‌های مهم سازمانهای ملی آمار کشورها بوده است.

یکی از هدفهای مهم یک سازمان ملی آمار، تأمین نیازهای اطلاعات آماری جامعه تا حد امکان است. این هدف، سازمان ملی آمار را موظف می‌کند تا اطلاعات آماری مورد نیاز کاربران را با کمک ابزارهای قانونی جمع‌آوری کرده و به صورت مناسبی در اختیار آنان قرار دهد. اما نکته مهم در این فرآیند، وجود خطر افشای هویت فردی توسط فرد یا افراد متخلف است. شخص متخلف فرضی ممکن است علاقمند به دستیابی به اطلاعات در مورد افرادی خاص باشد، یا اینکه تلاش کند تا با افشای اطلاعات، گردآورنده‌ی آنها را بی‌اعتبار کند.

عموما دو نوع افشای اطلاعات در نظر گرفته می‌شود: افشای هویت و افشای صفت. در نوع اول که در واقع مهمترین نوع افشا است، ابتدا فرد مورد نظر شناسائی شده و سپس بر اساس هویت فرد، اطلاعات مربوط به او از داده‌ها استخراج می‌گردد. اما همواره تعیین هویت شرط لازم برای افشای اطلاعات حساس در مورد یک پاسخگو نیست. در برخی موارد تنها دانستن این که یک پاسخگو عضوی از یک گروه است، بدون این که معلوم شود که کدام یک از آنها است، برای افشای اطلاعات در مورد وی کفایت می‌کند، این نوع افشاء را افشای صفت می‌نامند. عمل افشاء به هر دلیل که صورت پذیرد موجب بی‌اعتمادی عمومی نسبت به سازمان ملی آمار به عنوان حافظ اطلاعات شخصی پاسخگویان شده و کاهش همکاری یا عدم تمایل به همکاری پاسخگویان را در طرحهای آمارگیری به همراه خواهد داشت. نتیجه این امر عدم تولید آمار دقیق، روزآمد، به موقع، ارزان و گسترده خواهد بود که طبیعتاً امکان تأمین آمارهای مورد نیاز کاربران به صورت مناسب و مطلوب توسط سازمان ملی آمار وجود نخواهد داشت.

۲ مرور نوشتگان

موضوع کنترل افشای آماری^۱ اولین بار توسط دالینوس در سال ۱۹۷۷ مطرح شد و سپس سایر آمارشناسان در کشورهای امریکایی و اروپایی بحث و تحقیق در این موضوع را شروع کردند. کاکس (۱۹۸۰) روشهای پنهان‌سازی مقدماتی و مکمل را بحث کرده است. روشهای پرشیده^۲ شامل روشهای گرد کردن تصادفی توسط کاکس (۱۹۸۷)، روش پاسخ پس تصادفی شده^۳ توسط گوویلیو و همکاران (۱۹۸۸) بررسی شده‌اند. ویلنبرگ و دوال (۱۹۹۸) روی روشهای بازکدگذاری

1) Statistical Disclosure Control 2) Perturbative Methods 3) Post Randomized Response Method (PRAM) 4) Global Recoding

عام^۴ و پنهان سازی موضعی^۵ کار کرده‌اند. رهیافت‌های پرشیده نیز در شکل‌های مختلف توسط فاینبرگ و همکاران (۱۹۹۸) ارایه شده‌اند. پیشرفت‌های حاصل در کنترل افشای داده‌های آماری را می‌توان در دو گزارش کمیته فدرال روش شناختی آماری، در سالهای ۱۹۷۸ و ۱۹۹۴ و نیز در مقاله‌های شماره‌های ویژه سالهای ۱۹۹۳ و ۱۹۹۸ مجله آمار رسمی^۶ مشاهده نمود. در ایران تا بحال هیچگونه مطالعه و پژوهشی در موضوع کنترل افشای داده‌های آماری به جز پایان‌نامه کارشناسی ارشد آقای محمد بردبار عشرت آبادی (۱۳۸۲) صورت نگرفته است و نیز هیچ تجربه‌ای در ایمن سازی اطلاعات آماری منتشر شده وجود ندارد.

۳ روشهای محدودسازی افشای خرد داده‌های خرد

هرگاه یک سازمان آماری بخواهد یک مجموعه داده‌های خرد را منتشر کند، متغیرهای شناسایی مستقیم را از فایل داده‌ها حذف می‌کند. بنابراین شیوه مستقیمی برای تعیین اینکه یک رکورد معین مربوط به کدام پاسخگو است وجود ندارد. ظاهراً می‌توان گفت که در صورت انتشار چنین مجموعه داده‌های خردی خطر افشایی وجود ندارد اما این مطلب درست نیست. به مثال زیر توجه کنید. فرض کنید یک سازمان آماری بخواهد مجموعه داده‌های خردی شامل اطلاعاتی در مورد محل سکونت، شغل و سابقه جنایی پاسخگویان را منتشر نماید. علاوه بر این فرض کنید یک رکورد با ترکیب مقادیر زیر در مجموعه داده‌های خرد وجود داشته باشد: «محل سکونت: تهران، شغل: شهردار، سابقه جنایی: یک مورد».

اگر چه فرض بر این است که نام یا آدرس پاسخگو منتشر نمی‌شود، اما افراد زیادی در ایران پی می‌برند که پاسخگو چه کسی است. به ویژه می‌توانند نتیجه بگیرند که این پاسخگو یعنی شهردار تهران دارای یک سابقه جنایی است. به طور کلی هر پاسخگویی که در یک ترکیب از مقادیری که به تعداد کم در جامعه اتفاق می‌افتد قرار بگیرد، ممکن است در معرض خطر شناسایی باشد.

تشخیص اینکه در فایل داده‌های خرد کدام رکوردها و ترکیبهای متغیرها در معرض خطر افشای هستند، بستگی به چگونگی تعریف سناریوی افشای دارد. یک چنین سناریویی به تشریح مدلی می‌پردازد که یک شخص متخلف در جهت دستیابی به اطلاعات محرمانه یک یا چند پاسخگو در یک مجموعه داده‌های خرد، به چه صورتی ممکن است عمل کند و چه اطلاعات پیشینی در مورد یک جمعیت خاص می‌تواند داشته باشد. پس از مشخص شدن رکوردها و ترکیبهای ناایمن برای انتشار در بخش ارزیابی خطر افشای، مرحله حفاظت از فایل داده‌های خرد آغاز می‌شود. در این مرحله توسط روشهای محدودسازی خطر افشای مربوط به داده‌های خرد، یک فایل داده‌های خرد ناایمن به یک فایل داده‌های خرد ایمن قابل انتشار با خطر افشایی که در حد قابل قبولی پائین است، تبدیل می‌شود. این روشها به دو گروه عمده تقسیم می‌شوند.

5) Local Suppression 6) Journal of Official Statistics

گروه اول شامل روشهای ناپریشیده^۷ مانند روشهای باز کدگذاری عام و پنهان سازی موضعی است که به سبب اعمال آنها روی فایل داده‌های خرد، خطائی به داده‌ها افزوده نمی‌شود و فقط اندکی از میزان اطلاعات موجود در داده‌ها کاسته می‌شود. گروه دوم روشهای پریشیده نام دارند، مانند روش پس تصادفی سازی که به واسطه اعمال این روشها، مقادیر اصلی مشاهده شده، با مقادیر دیگری تعویض می‌شوند. این عمل به گونه‌ای صورت می‌گیرد که امکان انجام تحلیل‌های آماری رایج روی فایل پریشیده شده امکانپذیر است.

۱.۳ روش باز کدگذاری عام و روش پنهان سازی موضعی

این دو روش از رایجترین روشهای محدودسازی افشاء در داده‌های خرد می‌باشند و در گروه روشهای ناپریشیده قرار می‌گیرند. از این دو روش اغلب به صورت توأم استفاده می‌شود. هر دو روش باز کدگذاری عام و پنهان سازی موضعی برای متغیرهای رسته‌ای به کار می‌روند. در روش باز کدگذاری عام چندین رسته از یک متغیر مانند A با یکدیگر ترکیب شده و رسته جدیدی را پدید می‌آورند. مقادیر A متناظر با رسته‌های جدید، مجدداً کد گذاری می‌شوند. تعویض رسته‌ها تنها مختص بخش نایمن فایل داده‌های خرد نیست بلکه به منظور دستیابی به یک گروه بندی یکنواخت، روی تمام فایل داده‌های خرد اعمال می‌شود.

در روش پنهان سازی موضعی، مقدار یک متغیر مانند A در یک رکورد K با یک مقدار گمشده تعویض می‌گردد. در حالی که باز کدگذاری عام روی کل فایل داده‌های خرد تأثیر می‌گذارد، پنهان سازی موضعی تنها برای یک مقدار به خصوص در یک رکورد نایمن به کار می‌رود. هر دو روش موجب از دست رفتن مقداری از اطلاعات موجود در داده‌های خرد می‌شوند. علاوه بر این پنهان سازی موضعی ممکن است در صورتی که مقادیر گمشده نادیده گرفته شوند موجب برآوردگرهای اریب گردد. به همین دلیل سعی می‌شود تا از آن در مقیاسی کوچک استفاده شود. هنگامی که باید تعداد زیادی ترکیبهای نایمن حذف شوند، روش باز کدگذاری عام ترجیح داده می‌شود. وقتی که از تعداد اندکی پنهان سازی موضعی استفاده شود، غالباً اریبی تولید شده در برآوردگرها ناچیز است. در عمل یک توازن بین استفاده از این دو روش به وجود می‌آید. اغلب در ابتدا بعضی از متغیرها باز کدگذاری عام می‌شوند و در مورد باقیمانده رکوردهای نایمن از پنهان سازی موضعی برای حداقل مقادیر ممکن، استفاده می‌شود. مینیمم سازی تعداد پنهان سازی‌ها یک موضوع اساسی است که باید مورد توجه قرار گیرد.

۲.۳ روش پس تصادفی شده (PRAM)

پرام (PRAM) یک روش پریشیده است که برای کنترل افشای متغیرهای رسته‌ای در فایل داده‌های خرد به کار می‌رود. در این روش برای هر رکورد در فایل داده‌های خرد، امتیاز تعدادی از متغیرها متناظر با یک مکانیسم احتمالاتی معین، تغییر می‌کند، وقتی که فایل داده‌های خرد

7) Non-Perturbative Methods

پریشیده شد، شناسایی رکوردهای متناظر با اشخاصی معین در جامعه، دشوار می‌گردد. بنابراین رکوردها در فایل اصلی در برابر خطر شناسائی محافظت شده‌اند که این امر، هدف پرام است. از طرف دیگر تا وقتی که مکانیسم احتمالاتی استفاده شده معلوم باشد، مشخصه‌های داده‌های اصلی که پنهان شده‌اند می‌تواند از فایل داده‌های پریشیده شده برآورد شوند. بنابراین به کار بردن تمام تحلیل‌های آماری امکانپذیر است. فرض کنید ξ متغیری رسته‌ای در فایل اصلی باشد، هدف به کار بردن پرام روی این متغیر است. همچنین فرض کنید X همان متغیر رسته‌ای در فایل پریشیده شده است و ξ و X دارای k رسته، $k = 1, 2, \dots, K$ باشند. اگر $P_{kl} = P(X = l | \xi = k)$ احتمال این باشد که امتیاز $\xi = k$ در فایل اصلی به $X = l$ در فایل پریشیده شده تغییر یابد، آنگاه $P = p_{kl}$ ماتریسی $K \times K$ با اعضای p_{kl} است. P یک ماتریس مارکوف است یعنی در رابطه $PI = I$ که در آن I برداری $K \times 1$ از یک‌ها است، صدق می‌کند.

در اینجا فرض کلی بر این است که P معکوس‌پذیر است. این ویژگی برای به کار بردن پرام ضروری نیست ولی P^{-1} می‌تواند در برآورد توزیع فراوانی ξ در فایل اصلی، به خوبی برای برآورد واریانس تولید شده با استفاده از روش پرام، به کار رود.

پرام را می‌توان به صورت مستقل برای متغیرهای مختلف و همچنین هم زمان روی بیش از یک متغیر نیز به کار برد. مثال زیر تأثیر پرام بر محدودسازی افشاء را تشریح می‌کند.

مثال: فرض کنید فایل داده‌های خرد شامل n رکورد است. این فایل یک نمونه تصادفی ساده به اندازه n از جامعه‌ای با اندازه N است. فایل دقیقاً شامل یک زن جراح است. پرام، روی متغیر جنسیت و برای هر رکورد مستقل از رکوردهای دیگر به کار رفته است. امتیاز جنسیت با احتمال 0.9 بدون تغییر باقی می‌ماند و با احتمال 0.1 تغییر می‌کند. یعنی:

الف) ξ : متغیر جنسیت دارای دو رده $1 = \text{مرد}$ و $2 = \text{زن}$ (در فایل اصلی)

ب) X : متغیر جنسیت دارای دو رده $1 = \text{مرد}$ و $2 = \text{زن}$ (در فایل پریشیده شده)

ج) $P_{11} = 0.9$, $P_{12} = 0.1$, $P_{21} = 0.1$, $P_{22} = 0.9$

حال فرض کنید که شخص متخلف بداند که جامعه دارای 1 زن جراح و 99 مرد جراح است. احتمال این که وی پی ببرد، زن جراح در فایل پریشیده شده، در واقع زن جراح در جامعه می‌باشد برابر است با:

$$P(\xi = 2 | X = 2) = \frac{P_{22}P(\xi = 2)}{P_{12}P(\xi = 1) + P_{22}P(\xi = 2)}$$

$$= \frac{0.9 \times 0.1}{0.1 \times 0.99 + 0.9 \times 0.1} \approx 0.08$$

مشاهده می‌شود که مقدار این احتمال (0.08)، بسیار اندک است. بنابراین داده‌های پریشیده شده به اندازه کافی ایمن هستند. حال فرض کنید اگر متغیر جنسیت با احتمال 0.9999 بدون تغییر

بماند و با احتمال 0.0001 تغییر کند، آنگاه احتمال اینکه زن جراح در فایل پرشیده شده متناظر با زن جراح در جامعه باشد برابر 0.99 است که احتمال بالایی است و داده‌های پرشیده شده ایمن به نظر نمی‌رسند.

این مثال مشخص می‌کند که به یک کمیت مشخص نیاز است تا معلوم شود که آیا فایل نتیجه شده از اعمال پرام در واقع ایمن شده است؟ این ایده منجر به معرفی کمیتی با عنوان نسبت مورد انتظار به عنوان معیاری برای میزان عدم اطمینان ایجاد شده توسط پرام می‌گردد.

فرض کنید $\xi^{(r)}$ ($X^{(r)}$) امتیاز ξ (X) برای r امین رکورد در فایل داده‌های خرد باشد. به کار بردن پرام، بدین معنا است که با $\xi^{(r)} = k$ مفروض، امتیاز $X^{(r)}$ از توزیع احتمال $p_{k1}, p_{k2}, \dots, p_{kK}$ استخراج شود.

نسبت مورد انتظار امتیاز k ($k = 1, \dots, K$) یعنی $ER(k)$ به صورت زیر معرفی می‌شود:

$$k = 1, \dots, K$$

$$ER(k) = \frac{P_{kk} T_{\xi}(k)}{\sum P_{lk} T_{\xi}(l)}$$

که در آن $T_{\xi}(k)$ تعداد رکوردهای فایل اصلی برای $\xi^{(k)} = k$ است. اگر k در اصل یک امتیاز کمیاب باشد، نسبت مورد انتظار یعنی $ER(k)$ ، نسبت بین تعداد متوسط رکوردهایی که واقعا متعلق به امتیاز کمیاب k هستند و تعداد متوسط رکوردهایی است که در نتیجه به کار بردن پرام امتیاز k را کسب می‌کنند. مقدار کوچک ($ER(k)$) نشان دهنده این است که بیشتر احتمال می‌رود که یک رکورد برای $X^{(r)} = k$ در اصل متعلق به این امتیاز نبوده و بنابراین ایمن کننده فایل پرشیده شده است.

هنگامی که پرام برای یک فایل داده‌های خرد به کار می‌رود، این سؤال پیش می‌آید که: پرام چه تأثیری روی انواع تحلیل‌های آماری دارد؟ روشن است که تحلیلی که روی فایل پرشیده شده صورت گیرد، نتایجی متفاوت از تحلیل روی فایل اصلی دارد. این امکان وجود دارد تا نتایجی که از تحلیل‌های معینتی روی فایل پوشیده شده حاصل می‌شود، بتواند به نتایجی که می‌توانست از فایل اصلی حاصل گردد، تبدیل شوند.

برای مثال فرض کنید بخواهیم یک تحلیل رگرسیونی روی متغیر عددی وابسته y و متغیر رسته‌ای مستقل ξ انجام دهیم، و فرض کنید که پرام برای ξ با K رسته استفاده شده است. برای هر رکورد در فایل داده‌ها، متغیرهای ظاهری $\delta_1, \delta_2, \dots, \delta_k$ به صورت زیر تعریف می‌شوند:

$$\delta_k = \begin{cases} 1 & \longrightarrow \xi^{(r)} = k \\ 0 & \longrightarrow \xi^{(r)} \neq k \end{cases} \quad k = 1, \dots, K$$

علاوه بر این فرض کنید که:

$$T_{\xi}^y(k) = \sum y^{(r)} I_{\{\xi^{(r)}=k\}}$$

$$T_{\xi}^y = (T_{\xi}^y(1), \dots, T_{\xi}^y(K))$$

که در آن $y^{(r)}$ بر مقدار y برای r امین رکورد در فایل I بر تابع نشانگر^۸ دلالت دارد و T_{ξ}^y مجموع مقادیر y برای رکوردهایی است که دارای امتیاز k می‌باشند. کویمان (۱۹۹۷) نشان داد که T_{ξ}^y به طور نااریب توسط \bar{T} برآورد می‌شود.

اکنون رگرسیون y روی ξ با رگرسیون y روی $\delta_1, \dots, \delta_K$ یکسان است. ضرایب رگرسیونی توسط $(D^t D)^{-1} D^t y$ داده می‌شوند که در آن D ماتریسی $n \times K$ است که عضو (r, j) آن برابر مقدار δ_j برای r امین رکورد در فایل داده‌های اصلی است. (با توجه به اینکه δ_j فقط اعداد ۰ و ۱ را اختیار می‌کند، اعضای ماتریس D اعداد ۰ و ۱ خواهند بود.)

$(D^t D)$ می‌تواند به صورت نااریب توسط \bar{T} (توجه کنید که $(D^t D)$ ماتریسی قطری است که عضو (k, k) آن $T_{\xi}(k)$ است) و $D^t y$ می‌تواند به صورت نااریب توسط \bar{T} برآورد شوند. بعلاوه اینکه برآوردها نااریبند، برآوردگر رگرسیونی نتیجه شده سازگار است. این نتیجه برای تمام فنون تحلیل آماری بر پایه گشتاورهای مرتبه دوم داده‌ها از قبیل تحلیل تشخیصی^۹ و تحلیل واریانس^{۱۰} صادق است. در عمل ابتدا باید تصمیم گرفت که برای کدام متغیرها باید از پرام استفاده کرد. سپس، برای هر یک از متغیرهایی که پرام برای آنها به کار خواهد رفت، باید تصمیم گرفت که کدام رسته می‌تواند به کدام رسته و با چه احتمالی تغییر کند.

۳.۳ روش ریزانبوهش^{۱۱}

روش ریزانبوهش در گروه روشهای پرشیده جای می‌گیرد و برای متغیرهای کمی به کار می‌رود. این روش در ساده‌ترین حالت بر روی یک متغیر و در حالت پیچیده‌تر بر روی بیش از یک متغیر اعمال می‌شود. ایده اصلی در این روش این است که قواعد محرمانگی این اجازه را می‌دهد که اگر در یک مجموعه داده‌ها، پاسخگویان در گروههای k عضوی یا بیشتر قرار گیرند به طوری که k یک مقدار آغازین باشد و هیچ پاسخگویی از روی مقدارش شناسایی نشود، آنگاه آن مجموعه داده‌ها امکان انتشار دارد.

مسئله افزاری که در ریزانبوهش پیش می‌آید با مسئله خوشه‌بندی کلاسیکی که هدف آن تقسیم کردن جامعه به تعداد ثابتی گروههای ناسازگار، بدون توجه به اندازه گروهها است، متفاوت اما در ارتباط است. قیدی که در ریزانبوهش وجود دارد در مورد اندازه گروهها است. افزازهای نتیجه شده در ریزانبوهش نمی‌توانند شامل گروههایی با اندازه کوچکتر از k باشند. چنین افزازی، k -افراز نام دارد. طبق تعریف k -افراز بهینه، آن افزازی است که همگنی درون گروه، ماکسیمم باشد، همگنی بیشتر درون گروه، فقدان اطلاعات کمتری را به همراه دارد، زیرا در ریزانبوهش مقادیر

8) Indicator Function 9) Discriminant Analysis 10) Analysis of Variance
11) Microaggregation

جدول ۱: میزان فروش مؤسسات تجاری به تفکیک نوع و ناحیه فعالیت

| ناحیه فعالیت | A | B | C | کل |
|-----------------|----|----|-----|-----|
| ۱ | ۱۱ | ۴۷ | ۵۸ | ۱۱۶ |
| ۲ | ۱ | ۱۵ | ۳۳ | ۴۹ |
| ۳ | ۲ | ۳۱ | ۲۰ | ۵۳ |
| کل | ۱۴ | ۹۳ | ۱۱۱ | ۲۱۸ |

یک گروه با مقدار متوسط گروه، تعویض می‌شوند. در بحث خوشه‌بندی، استفاده از معیار مجموع توانهای دوم برای اندازه‌گیری همگنی، مساله‌ای رایج است که از آن در ریزانپوهش نیز استفاده می‌شود.

۴ روشهای محدودسازی افشاء در جدول‌های انتشاراتی

۱.۴ داده‌های جدولی

جدول‌ها رایجترین محصولات سازمان‌های آماری می‌باشند که شامل داده‌های جمع‌بندی شده به عنوان مقادیر خانه‌های جدول می‌باشند و به دو صورت جدول‌های مقداری و جدول‌های فراوانی منتشر می‌شوند. به عنوان مثال میزان سرمایه‌گذاری کارخانجات بر اساس ناحیه و نوع فعالیت، یک جدول مقداری و تعداد کارخانجات برحسب ناحیه و نوع فعالیت یک جدول فراوانی را تشکیل می‌دهند.

۲.۴ تشخیص خانه‌های حساس

به علت اینکه داده‌ها مربوط به افراد پاسخگو به صورت تکی نیستند، به نظر می‌رسد که خطر افشای اطلاعات فردی وجود نداشته باشد، در صورتی که همیشه اینگونه نیست. برای روشن شدن مطلب، جدول زیر را که نشان دهنده‌ی میزان فروش تعدادی از مؤسسات تجاری به تفکیک نوع و ناحیه فعالیت است، در نظر بگیرید. این جدول در نگاه اول برای انتشار، پذیرفتنی به نظر می‌رسد زیرا تنها شامل داده‌های خلاصه شده است و داده‌های از پاسخگویان تکی ندارد. به هر حال این نتیجه‌گیری عجولانه است.

برای مثال، فرض کنید که این جدول از یک مشاهده کامل حاصل شده است و تنها یک مؤسسه در ناحیه B با فعالیت ۲ وجود دارد. در نتیجه می‌توان فهمید که میزان فروش این مؤسسه تجاری برابر ۱۵ است. پس اگر قرار باشد که میزان فروش این مؤسسه تجاری حفاظت شود، جدول نباید منتشر شود. بنابراین جدولهایی که دارای خانه‌های شامل داده‌های مربوط به تنها یک پاسخگو هستند، نباید منتشر شوند. ولی همیشه موضوع به این سادگی نیست. فرض

کنید در ناحیه B و فعالیت ۲ به جای یک مؤسسه، دو مؤسسه تجاری وجود داشته باشد. در این حالت هر کدام از این دو مؤسسه به راحتی می‌تواند به میزان فروش دیگری پی ببرد. ممکن است منتشرکننده‌ی داده‌ها به جای وجود حداقل دو عضو در هر خانه‌ی جدول، خواستار وجود سه عضو یا بیشتر در هر خانه‌ی جدول باشد. متأسفانه این حالت نیز همیشه رضایتبخش نیست. فرض کنید در ناحیه B ، ده مؤسسه تجاری دارای فعالیت ۲ باشند و یکی از آنها ۹۵ درصد مقدار کل خانه را ارایه نماید. در این حالت اگر شخص متخلف بداند که میزان فروش این مؤسسه خیلی بالا است، می‌تواند یک برآورد تقریباً دقیق از میزان فروش آن به‌دست آورد. در چنین حالتی گفته می‌شود که این خانه‌ی جدول، تحت تسلط این عضو است.

بنابراین برای کنترل خطر افشاء در جدول‌ها، ابتدا باید دنبال راهی گشت که توسط آن بتوان خانه‌هایی را که مناسب انتشار نیستند و امکان افشای اطلاعات تکی پاسخگویان توسط آنها وجود دارد (خانه‌های حساس) شناسایی نمود که برای این موضوع معمولاً از قواعدی استفاده می‌شود که در ادامه بیان می‌شوند.

در مورد جدول‌های مقداری، عمومی‌ترین ملاکی که به‌کار می‌رود، قاعده‌ی تسلط (n, k) است که به صورت زیر تعریف می‌شود: یک خانه جدول حساس است اگر مجموع مقادیر عرضه شده توسط n عدد از بزرگترین پاسخگویان به مقدار کل خانه، بیشتر از k درصد مقدار کل خانه را شامل شود.

معمولاً مقدار کوچکی (ماکسیمم ۵) برای n و مقدار بزرگی (مثلاً ۸۰) برای k منظور می‌شود. ایده‌ی اصلی این قاعده این است که، حفاظت‌کننده‌ی داده‌ها از این که $n - ۱$ پاسخگو بتوانند به وسیله‌ی جمع کردن مقادیرشان و مقایسه‌ی آن با مقدار کل خانه، یک برآورد دقیق از مقدار پاسخگوی m به‌دست آورند، جلوگیری می‌کند.

قاعده‌ی دیگری که در مورد جدول‌های مقداری به‌کار می‌رود قاعده‌ی پیشین-پسین $((p, q))$ است. فرض می‌کنیم که هر پاسخگو قبل از انتشار جدول می‌تواند مقدار عرضه شده توسط پاسخگوی دیگر به مقدار یک خانه‌ی جدول را با اختلاف کمتر یا برابر q درصد از مقدار واقعی آن، برآورد نماید. اگر پس از انتشار جدول، این امکان برای برخی پاسخگویان به وجود آید که مقدار پاسخگویی دیگر به خانه‌ی جدول را با اختلاف کمتر یا برابر p درصد ($p < q$) مقدار اصلی آن، برآورد نماید، این خانه از جدول، حساس شناخته می‌شود.

تعیین حساسیت یک خانه می‌تواند از راههای گوناگون تعیین شود. در این قسمت نشان می‌دهیم که قاعده‌ی تسلط (n, k) و قاعده‌ی پیشین-پسین $((p, q))$ در حقیقت اعضایی از یک رده‌ی بزرگتر از اندازه‌های حساسیت‌اند که رده‌ی معیارهای حساسیت خطی^{۱۴} نامیده می‌شوند.

در مورد جدول‌های فراوانی که در خانه‌های آن، تعداد پاسخگویان دارای شرایط آن خانه‌ها آمده است، قاعده‌ای که غالباً برای تشخیص خانه‌های حساس به کار می‌رود بدین صورت است که نباید تعداد پاسخگویان در هر خانه‌ی جدول از عدد مشخصی (مقدار آغازین) مثلاً سه کمتر

12) (n,k) Dominance Rule 13) (p,q) Prior- Posterior rule 14) Linear Sensitivity Measures

جدول ۲: قسمتی از یک جدول بزرگ در مورد تخلف‌های زیست محیطی

| تخلف | بله | خیر | کل |
|--------------------|-----|-----|----|
| ناحیه و فعالیت | - | - | - |
| ناحیه B و فعالیت ۳ | ۴ | ۱ | ۵ |
| - | - | - | - |

باشد. در غیر اینصورت، امکان شناسایی پاسخگویان این خانه‌ها به صورت تکی وجود دارد. به هر حال حالتیابی وجود دارد که چنین قاعده‌ی تسلطی برای تعیین خانه‌های حساس مناسب نیست. به‌عنوان مثال جدول ۲ که قسمتی از یک جدول بزرگتر در مورد تخلفهای زیست محیطی مؤسسات است را در نظر بگیرید. این جدول مشخص می‌کند که تنها یک مؤسسه دارای تخلف نبوده است. یک شخص بیرونی می‌تواند با احتمال ۸۰ درصد چهار مؤسسه دارای تخلف را تشخیص دهد در حالی که هر کدام از مؤسسات با احتمال ۷۵ درصد می‌تواند تخلف را تشخیص دهد. پس از تشخیص خانه‌های حساس، باید قبل از انتشار جدول، روش‌های محدودسازی افشاء را روی آنها به کار برد تا در برابر خطر افشای اطلاعات مربوط به پاسخگویان تکی، ایمن گردند.

۳.۴ روش‌های باز طراحی جدول و پنهان‌سازی خانه‌ای

این روش‌ها که از گروه روش‌های غیر پرشیده هستند، زیاد مورد استفاده قرار می‌گیرند و از قدیمی‌ترین روش‌های محدودسازی خطر افشاء در جدول‌ها می‌باشند. ممکن است در یک جدول، در یک سطر یا ستون، خانه‌های زیادی، حساس تشخیص داده شوند. در این حالت توصیه می‌شود ابتدا جدول دوباره طراحی شود. یعنی در طبقه‌بندی متغیرهای تبیینی^{۱۵} تجدیدنظر شود. در این صورت جزئیات اطلاعات آماری ارایه شده در جدول کاهش می‌یابد اما از تعداد خانه‌های حساس به مقدار زیادی کاسته خواهد شد. اگر هنوز تعداد خانه‌های حساس زیادی باقی مانده باشد می‌توان باز هم جزئیات جدول را کاهش داد. هنگامی که تعداد خانه‌های حساس اندک باشد، می‌توان از روش‌های کنترل افشای دیگری که به صورت موضعی عمل می‌کنند استفاده نمود. نمونه‌ای از این روشها، روش پنهان‌سازی خانه‌ای است. یکی از رایج‌ترین روش‌هایی که برای کنترل خطر افشاء در جدول‌ها به کار می‌رود، روش پنهان‌سازی خانه‌ای است. در این روش مقدار خانه‌ی حساس از جدول حذف شده و به جای آن از یک علامت مثلا حرف X استفاده می‌شود. برای توضیح بیشتر، جدول ۳ را در نظر بگیرید. فرض کنید که مقدار خانه متناظر با فعالیت دو و ناحیه C طبق قاعده تسلط به کار رفته، حساس تشخیص داده شده و نمی‌تواند منتشر شود. مقدار این خانه پنهان می‌شود. به طور کلی، پنهان کردن خانه حساس، به تنهایی کافی نیست. همانطور که در جدول ۳ مشخص است، خانه پنهان

15) Explanatory Variables

جدول ۳: میزان سرمایه‌گذاری کارخانجات به تفکیک ناحیه و نوع فعالیت (پس از پنهان‌سازی مقدماتی)

| ناحیه | A | B | C | کل |
|----------|----|-----|----|-----|
| فعالیت ۱ | ۲۰ | ۵۰ | ۱۰ | ۸۰ |
| ۲ | ۸ | ۱۹ | X | ۴۹ |
| ۳ | ۱۷ | ۳۲ | ۱۲ | ۶۱ |
| کل | ۴۵ | ۱۰۱ | ۴۴ | ۱۹۰ |

شده می‌تواند با استفاده از جمع‌های حاشیه‌ای به آسانی محاسبه شود. گزینه دیگر، پنهان کردن خانه‌های اضافه‌ی دیگری در جدول که حساس نیستند، است. این عمل، پنهان‌سازی خانه‌ای مکمل و خانه‌ها، پنهان شده‌های مکمل نامیده می‌شوند. یک الگوی پنهان‌سازی خانه‌ای مکمل در جدول ۳، پنهان کردن مقادیر متناظر با خانه‌های فعالیت دو و سه در ناحیه A و فعالیت سه در ناحیه C است.

۴.۴ افزودن نوفه به داده‌های خرد قبل از جدول‌سازی

یک روش دیگر حفاظت پاسخگویان تکی افزودن نوفه به داده‌های آنهاست. این شیوه برای داده‌های جدول‌های مقداری به کار می‌رود. فرض می‌شود که به داده‌ی متناظر با یک مؤسسه، عدم اطمینانی به میزان اندک مثلاً ده درصد افزوده شود (درصد استفاده شده، نزد حفاظت‌کننده‌ی داده‌ها محرمانه باقی می‌ماند). در نتیجه اگر یک خانه شامل تنها یک مؤسسه باشد یا اینکه اگر تنها یک مؤسسه، خانه را تحت تسلط داشته باشد، آنگاه مقدار این خانه نمی‌تواند به عنوان تقریبی نزدیک از مقدار مؤسسه مسلط استفاده شود، زیرا این مقدار دارای نوفه‌ی اضافه شده است (در این حالت با ۱۰ درصد تغییر می‌کند) با افزودن نوفه‌ی، از افشای مقدار صحیح مؤسسه مسلط جلوگیری می‌شود.

۵.۴ روش گرد کردن

در این روش مقادیر خانه‌های جدول به مضرب‌های صحیح یک عدد پایه‌ی ثابت، گرد می‌شوند. گرد کردن، یک حالت ویژه‌ی داده‌های پرشده است که در آن برای حفظ اطلاعات، داده‌ها تغییر می‌یابند.

چندین روش گرد کردن وجود دارد. ساده‌ترین روش، روش گرد کردن استاندارد است که در آن هر مقدار از خانه‌های جدول به نزدیکترین مضرب پایه‌ی گرد کردن، گرد می‌شود. گرد کردن استاندارد دو محدودیت عمده دارد. اول اینکه، معمولاً ساختار جمع‌پذیری جدول حفظ نمی‌شود، برای مثال، جدول ۴ را که جدولی یک بعدی و نشان دهنده‌ی تعداد مؤسسات تجاری دارای فعالیت

جدول ۴: تعداد مؤسسات تجاری با فعالیت دو به تفکیک ناحیه

| کل | ناحیه B | ناحیه A | ناحیه |
|----|---------|---------|----------|
| ۴ | ۲ | ۲ | فعالیت |
| | | | فعالیت ۲ |

دو به تفکیک ناحیه است را در نظر بگیرید. فرض کنید بخواهیم این جدول را با پایه‌ی گردکردن، ۵، گرد کنیم. هر دو مقدار خانه‌های داخلی جدول به عدد صفر گرد می‌شوند، در حالی که جمع حاشیه‌ای جدول به عدد ۵ گرد می‌شود. این امر موجب ناسازگاری در جدول می‌گردد. دومین محدودیت که از نظر کنترل افشاء مهم است، این است که گاهی اوقات امکان دارد که مقادیر خانه‌های اصلی آشکار گردند یا اینکه به تعداد معدودی از مقادیر اصلی محدود گردد. این نقص گردکردن استاندارد، انگیزه‌ی جستجو برای روشهای دیگر گردکردن که جمع‌پذیری را حفظ نمایند، به وجود آورد. ایده‌ی این روشها این است که هر مقدار (به جز مقادیری که مضر بی از مقدار پایه‌اند) به یکی از دو نزدیکترین مضرب مقدار پایه به طوری که جدول جمع‌پذیر باشد، گرد می‌شوند.

روشهای گردکردن قطعی نیستند بلکه احتمالاتی‌اند. علاوه بر این دارای این خاصیت هستند که خطای گردکردن دارای امید ریاضی صفر است. این امر در صورتی تحقق می‌یابد که یک عدد مانند a که می‌تواند به صورت $a = kb + r$ نوشته شود که در آن $0 \leq r < b$ و k عددی صحیح و b پایه‌ی صحیح گردکردن است، به a^* که مقادیر kb و $(k + 1)b$ را با احتمالهای زیر اختیار می‌کند گرد شود.

$$p(a^* = (k + 1)b) = \frac{r}{b} p(a^* = kb) = 1 - \frac{r}{b}$$

به سادگی دیده می‌شود که:

$$E(a^*) = (k + 1)b \times \frac{r}{b} + kb \times \left(1 - \frac{r}{b}\right) = kr + r + kb - kr = kb + r = a$$

می‌توان نشان داد تنها a ی که در ۱ تعریف شده است در ۲ صدق می‌کند. در جدولهای یک و دو بعدی مسأله گردکردن با مشکل خاصی روبرو نمی‌شود اما در ابعاد بالاترگاهی اوقات غیرممکن است. ذکر این نکته ضروری است که در حال حاضر در کشور ما با وجود ارائه اطلاعات آماری به صورت بسیار کلی باز هم می‌توان موارد زیادی از افشای دقیق یا تقریبی اطلاعات شخصی را به دست آورد. در حالی که در صورت اعمال روشهای ذکر شده در این مقاله بر روی داده‌های آماری می‌توان در عین ارائه مفصلتر اطلاعات آماری، جنبه‌ی محرمانگی اطلاعات شخصی نیز را رعایت نمود.

۵ کاربرد

بخش صنعت در برنامه‌های توسعه اقتصادی، اجتماعی و فرهنگی کشور به عنوان تأمین کننده بسیاری از نیازهای برنامه‌ی توسعه، جایگاه ویژه‌ای را به خود اختصاص داده است. بنابراین ضرورت شناخت دقیق، رفع مشکلات و نارسائی‌ها و انجام برنامه‌ریزی‌های آتی این بخش موجب شده است تا اطلاعات تفصیلی و به هنگام کارگاه‌های صنعتی در قالب زیر جامعه‌های ویژه‌ی جمع‌آوری و منتشر شوند.

طرح آمارگیری از کارگاه‌های صنعتی ده نفر کارکن و بیشتر^{۱۶} هر ساله توسط مرکز آمار ایران و به صورت سرشماری در قالب همین سیاست صورت می‌پذیرد. جمع‌آوری اطلاعات دقیق در این بخش از اهمیت ویژه‌ای برخوردار است، لذا پاسخگویان باید مطمئن باشند که اطلاعاتی که توسط آنها ارائه می‌شود محفوظ بوده و امکان افشاء و استفاده‌های غیرمجاز از آنها وجود ندارد. تاکنون اطلاعات مربوط به کارگاه‌های صنعتی ده نفر کارکن و بیشتر برای حفظ محرمانگی اطلاعات فقط در سطح استان منتشر شده‌اند، در حالی که می‌توان هم زمان با انتشار اطلاعات در سطح جغرافیائی کوچکتر، جنبه‌ی محرمانگی آنها را نیز رعایت کرد. در این بخش نشان داده خواهد شد که برخی اطلاعات در جدول‌های انتشاراتی این آمارگیری قابل افشاء است. سپس با استفاده از نرم‌افزارهای ذیربط چگونگی ساختن جدول‌های ایمن را شرح خواهیم داد.

۱.۵ چند نمونه از اطلاعات قابل افشاء از جدول‌های انتشاراتی

در این قسمت نمونه‌هایی از جدول‌های انتشاراتی مرکز آمار ایران به‌عنوان بزرگترین تولید کننده آمار در کشور که دارای خطر افشای اطلاعات فردی هستند را به عنوان مثال ذکر می‌کنیم. در ایران کمتر انتشار اطلاعات آماری به صورت داده‌های خرد صورت گرفته است. انتشارات جدولی نیز به صورت بسیار کلی صورت می‌گیرند به طوری که به عنوان مثال در انتشار اطلاعات مربوط به طرح آمارگیری از کارگاه‌های صنعتی ده نفر کارکن و بیشتر، متغیر ناحیه تنها در سطح کشور و استان ارائه می‌شود. با این حال، در همین سطح کلی نیز می‌توان موارد قابل افشایی را مشاهده کرد. در زیر دو نمونه از جدول‌های نایمن انتشار یافته توسط مرکز آمار ایران آمده است. در مطالعه مثالهای زیر توجه به دو نکته ضروری است:

الف) برای حفظ محرمانگی از بردن نام واحدی که هویت آن قابل افشاء است، خودداری شده است، و

ب) برای اختصار از انعکاس جدول‌های انتشاراتی اجتناب شده و تنها به شماره جدولها از نشریه «نتایج آمارگیری از کارگاه‌های صنعتی ده نفر کارکن و بیشتر ۱۳۷۹» ارجاع داده شده است.

۱۶) مشخصات این طرح را می‌توان در نشریه «نتایج آمارگیری از کارگاه‌های صنعتی ۱۰ نفر کارکن و بیشتر ۱۳۷۹» از انتشارات مرکز آمار ایران یافت.

جدول ۵: قسمتی از جدول شماره ۱ نشریه نتایج آمارگیری از کارگاههای صنعتی ده نفر کارکن و بیشتر ۱۳۷۹

| ۸ | ۷ | ۶ | ۵ | ۴ | ۳ | ۲ | ۱ | جمع | تعداد شاغلان |
|---|---|---|---|---|---|---|---|-----|--|
| | | | | | | | | | فعالیت |
| : | : | : | : | : | : | : | : | : | : |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | تولید محصولات از توتون، تنباکو و سیگار |
| : | : | : | : | : | : | : | : | : | : |
| 5 | 4 | 4 | 0 | 1 | 0 | 0 | 0 | 14 | تولید وسایل نقلیه موتوری |
| : | : | : | : | : | : | : | : | : | : |

مثال ۱. قسمتی از جدول شماره ۵ از نشریه نتایج آمارگیری از کارگاههای صنعتی ده نفر کارکن و بیشتر سال ۱۳۷۹ جمع‌آوری شده در سال ۱۳۸۰ که با عنوان «تعداد کارگاههای صنعتی ده نفر کارکن و بیشتر برحسب تعداد شاغلان و نوع فعالیت ۱۳۷۹:» آمده است، به صورت زیر است: در این جدول کدهای متغیر تعداد شاغلان به صورت زیر تعریف شده است: ۱-۱۹: ۱۰-۲۹: ۲۰-۳۹: ۳۰-۴۹: ۴۰-۴۹: ۵۰-۴۹۹: ۶۰-۴۹۹: ۷۰-۹۹۹: ۸۰-۱۰۰۰: و بیشتر.

این جدول یک جدول فراوانی است و همانطور که از آن مشخص است تنها یک کارگاه تولید محصولات از توتون و تنباکو و سیگار در کل کشور وجود دارد که تعداد کارکنانش بیشتر از ۱۰۰۰ نفر است. به علت یکتا بودن این کارگاه در جامعه، به راحتی می‌توان از جدول شماره (۲) نشریه، عمومی بودن نحوه مدیریت و دولتی بودن وضعیت حقوقی کارگاه، از جدول شماره (۵) همان نشریه، تعداد دقیق شاغلان به تفکیک تولیدی یا غیرتولیدی بودن و مرد یا زن بودن آنها، از جدولهای شماره (۶) و (۷)، تعداد دقیق شاغلان با مزد و حقوق و بدون مزد و حقوق کارگاه به تفکیک تولیدی یا غیرتولیدی بودن و مرد یا زن بودن آنها، از جدول شماره (۸)، تعداد دقیق شاغلان به تفکیک سطح مهارت آنها (کارگر ساده، ماهر، تکنسین یا مهندس)، از جدول شماره (۱۱)، میزان مزد و حقوق و سایر پرداختیها به شاغلان، از جدول شماره (۱۷)، تعداد دقیق شاغلان به تفکیک سطح سواد (بیسواد، کمتر از دیپلم، دیپلم، لیسانس، فوق لیسانس یا دکترا)، از جدول شماره (۲۳)، میزان پرداختیهای کارگاه، از جدول شماره (۲۴)، ارزش داده‌های فعالیت صنعتی، از جدول شماره (۲۵)، ارزش ستانده‌های فعالیت صنعتی، از جدول شماره (۲۸)، ارزش سرمایه‌گذاربهای کارگاه، از جدول شماره (۲۹)، ارزش انواع موجودی انبار کارگاه و خلاصه تمام اطلاعات مربوط به این کارگاه خاص قابل دسترسی است. اگر رضایت این کارگاه برای انتشار اطلاعات خاص آن کسب نشده باشد، یک نمونه از افشای کامل اطلاعات فردی اتفاق افتاده است.

مثال ۲. این مثال در مورد کارگاههای تولید وسایل نقلیه موتوری است. در نگاه اول خطر افشانی احساس نمی‌شود، در حالی که این چنین نیست. به عنوان مثال در جدول شماره (۴)، تعداد شاغلان این ۱۴ کارگاه تولید وسایل نقلیه موتوری ۳۲۵۱۰ نفر ذکر شده است. منتشرکننده اطلاعات می‌داند که تعداد شاغلان بزرگترین کارگاه ۱۴۷۱۷ نفر، دومین کارگاه بزرگ ۵۱۱۳ نفر

و سومین کارگاه بزرگ ۲۷۶۱ نفر است که جمعا حدود ۷۷ درصد کل شاغلان این بخش را به خود اختصاص داده‌اند. در صورت اتحاد، این سه کارگاه می‌توانند تعداد شاغلان دیگر کارگاهها را برآورد نمایند. برای توضیح بیشتر، فرض کنید x_1, \dots, x_{14} مقادیر عرضه شده توسط ۱۴ کارگاه باشد. طبق قاعده پیشین-پسین (p, q) این خانه از جدول (۴) حساس است اگر دومین کارگاه بزرگ بتواند پس از انتشار جدول، مقدار صفت بزرگترین کارگاه را با اختلاف کمتری یا برابر p درصد از مقدار واقعی آن برآورد کند ($p < q$). فرض کنید دومین پاسخگوی بزرگ قبل از انتشار جدول بتواند هر یک از مقادیر دیگر پاسخگویان را با اختلاف کمتری یا برابر 5° درصد مقدار واقعی آنها برآورد کند ($q = 5^\circ$). در این صورت دومین کارگاه بزرگ پس از انتشار جدول می‌تواند پی ببرد که x_1 حداکثر برابر

$$x_1 + \left(\frac{q}{100}\right) \sum_{i=3}^{14} x_i = 17417 + \left(\frac{5^\circ}{100}\right) \times 9980 = 22407$$

و حداقل برابر

$$x_1 + \left(\frac{q}{100}\right) \sum_{i=3}^{14} x_i = 17417 + \left(\frac{5^\circ}{100}\right) \times 9980 = 12427$$

است یعنی تعداد شاغلان بزرگترین کارگاه را در بازه $\{12427 \text{ و } 22407\}$ برآورد نماید که اختلاف آن با مقدار واقعی کمتر از ۲۹ درصد است ($p = 29$).

به همین صورت در جدول شماره (۲۴) این نشریه، ارزش کل مواد خام و اولیه، لوازم بسته‌بندی، ابزار و وسایل کار کم‌دوام مصرف شده توسط این کارگاهها در خانه مربوط به آنها 13204176 میلیون ریال ذکر شده است. با توجه به این که مقدار این متغیر برای بزرگترین کارگاه 7080000 میلیون ریال، برای دومین کارگاه بزرگ 3300000 میلیون ریال و برای سومین کارگاه بزرگ برابر 626000 میلیون ریال است، این سه کارگاه جمعا حدود ۸۳ درصد مقدار کل خانه را ارایه کرده‌اند. حال فرض کنید دومین کارگاه بزرگ قبل از انتشار جدول، میزان این متغیر برای هر کارگاه را بتواند با اختلاف کمتری یا برابر 5° درصد مقدار واقعی آنها برآورد کند، در اینصورت پس از انتشار جدول طبق مطالب مطرح شده در بالا، این کارگاه مقدار این متغیر برای بزرگترین کارگاه را می‌تواند در بازه $\{5877912 \text{ و } 8282088\}$ یعنی با تفاوتی کمتر از ۱۷ درصد مقدار واقعی برآورد نماید.

۲.۵ ایمن‌سازی اطلاعات آماری انتشاراتی

سازمان‌های ملی آمار کشورها داده‌های حاصل از آمارگیرها را یا به صورت داده‌های خرد و یا به صورت داده‌های جدولی منتشر می‌کنند. همانگونه که اشاره شد در هر صورت امکان افشاء وجود دارد. به منظور ایمن‌سازی داده‌های خرد از نرم‌افزار $\mu ARGUS$ و برای ایمن‌سازی داده‌های

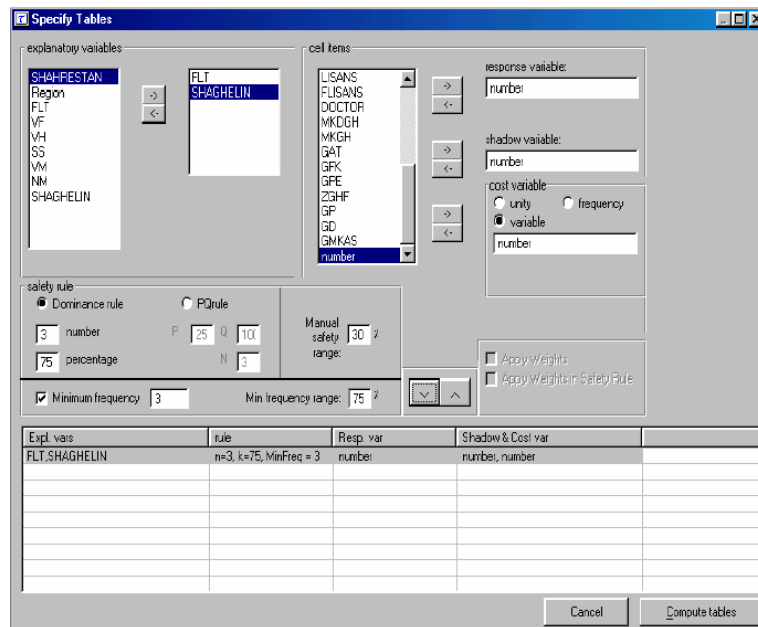
جدولی از نرم افزار $\tau ARGUS$ که حاصل پروژه ${}^{14}CASC$ است استفاده می شود. این پروژه با هدف تهیه ابزارهای عملی برای کنترل افشای آماری توسط تعدادی از کشورهای اتحادیه اروپایی، سازماندهی و اجرا شده است.

با توجه به اینکه داده های سرشماری از کارگاه های صنعتی دارای ده نفر کارکن و بیشتر به صورت جدول های پیشابندی منتشر می شوند، لذا در این بخش به نحوه ایمن سازی جدول های انتشاراتی طرح مذکور در استان تهران در سال ۱۳۷۹ با استفاده از نرم افزار $\tau ARGUS$ پرداخته می شود.

رایجترین شکل انتشار داده های آماری، جدول های پیشابندی هستند، که خود به دو نوع جدول های فراوانی و جدول های مقداری تقسیم می شوند. اگر چه اطلاعات جدول ها، خلاصه شده هستند ولی با توجه به نوع جدول و اعداد خانه ها، امکان افشاء وجود دارد. به منظور ساختن جدول های انتشاراتی ایمن، به فایل داده های خرد و فایل فراداده حاوی مشخصات متغیرهای موجود در فایل داده های خرد نیاز است. هر دو فایل به صورت متن هستند. در فایل داده های خرد هر سطر اصلی شامل نام یک متغیر، نقطه شروع، طول و مقدار گمشده می باشد. در سطرهای بعد مشخصات متغیر وجود دارند. در صورتی که فایل فراداده وجود نداشته باشد، نرم افزار امکان ایجاد آن را دارد. پس از آن که فایل فراداده توسط نرم افزار $\tau ARGUS$ خوانده شد، می توان جدول هایی را که قصد حفاظت آنها در برابر خطر افشاء وجود دارد، مشخص کرد. بدین منظور نرم افزار پنجره زیر را برای مشخص کردن جدول های درخواستی باز می کند. در قسمت بالای پنجره شکل ۱، دو قسمت قابل مشاهده است. در یک قسمت، متغیرهای تبیینی (explanatory variables) که برای ساختن جدول به کار می روند، مشخص می شوند و در قسمت دیگر متغیر پاسخ (response variable) که برای محاسبه جمع های خانه های مورد استفاده قرار می گیرد، تعیین می شود. گزینه ی Shadow Variable متغیری است که برای به کار بردن قاعده ی تسلط استفاده می شود. این متغیر در بیشتر اوقات همان متغیر پاسخ است اما می توان متغیر دیگری را نیز انتخاب کرد.

گزینه ی Cost Variable نشان دهنده متغیر مورد استفاده برای محاسبه ی مینیمم فقدان اطلاعات ناشی از اعمال روش های کنترل افشاء است. این متغیر لزوماً همان متغیر پاسخ نیست بلکه می تواند متغیر تعداد پاسخگویان عرضه کننده ی مقدار خانه های جدول باشد. در قسمت پائین پنجره پارامترهای قاعده ی تسلط مورد نظر را می توان وارد نمود. در اینجا دو قاعده ی تسلط آورده شده است. قاعده ی تسلط (n, k) و قاعده ی پیشین-پسین (p, q) که یکی از آن دو را می توان انتخاب کرد. باید توجه کرد که قاعده ی p درصد همان قاعده ی پیشین-پسین (p, q) درصد با در نظر گرفتن $q = 100$ است. در هر صورت نرم افزار بازه ی حفاظت را تعیین می کند. در حالتی که جدول های انتشاراتی فراوانی هستند، حفاظت کننده ی جدول ممکن است به سادگی خانه هایی که دارای فراوانی ای کمتر از یک مقدار آغازین معین هستند را به عنوان خانه های ناایمن در نظر بگیرد.

در اینجا به عنوان نمونه جدول تعداد کارگاهها به تفکیک نوع فعالیت و تعداد شاغلان درخواست



شکل ۱: پنجره انتخاب جدول

| Variable | dim 1 | dim 2 | Code | Label | Freq | dim 1 | dim 2 |
|-----------|-------|-------|------|-------|------|-------|-------|
| FLT | 42 | 4323 | | Total | 2716 | 285 | 0 |
| SHAGHELIN | 285 | 4323 | 15 | | 272 | 0 | 86 |
| | | | 151 | | 71 | 0 | 51 |
| | | | 1512 | | 3 | 1 | 3 |
| | | | 1514 | | 7 | 0 | 6 |
| | | | 1515 | | 18 | 0 | 16 |
| | | | 1516 | | 16 | 0 | 16 |
| | | | 1517 | | 1 | 1 | 1 |
| | | | 1518 | | 2 | 1 | 2 |
| | | | 1519 | | 24 | 0 | 22 |
| | | | 152 | | 21 | 0 | 18 |
| | | | 1520 | | 21 | 0 | 18 |
| | | | 153 | | 43 | 0 | 32 |
| | | | 1531 | | 27 | 0 | 22 |
| | | | 1532 | | 3 | 1 | 3 |
| | | | 1533 | | 13 | 0 | 12 |
| | | | 154 | | 127 | 0 | 47 |
| | | | 1542 | | 5 | 0 | 4 |
| | | | 1543 | | 7 | 0 | 7 |
| | | | 1544 | | 18 | 0 | 14 |
| | | | 1545 | | 8 | 0 | 8 |
| | | | 1546 | | 56 | 0 | 28 |
| | | | 1547 | | 4 | 1 | 4 |
| | | | 1548 | | 29 | 0 | 22 |
| | | | 155 | | 10 | 0 | 10 |
| | | | 1551 | | 3 | 1 | 3 |
| | | | 1553 | | 1 | 1 | 1 |
| | | | 1555 | | 4 | 1 | 4 |
| | | | 1556 | | 2 | 1 | 2 |
| | | | 16 | | 1 | 1 | 1 |
| | | | 160 | | 1 | 1 | 1 |

شکل ۲: پنجره نشان دهنده‌ی تعداد خانه‌های نایمن به تفکیک رسته‌های متغیرهای تبیینی

شده است. خانه‌های دارای چهار پاسخگو و کمتر نایمن در نظر گرفته شده‌اند. پس از اینکه جدول و قواعد حفاظت آن مشخص شدند، نرم‌افزار مشخصات جدول مورد نظر را نشان می‌دهد که در آن، تعداد خانه‌های نایمن ساخته شده توسط رسته‌های متغیرهای تبیینی انتخاب شده در ابعاد جدول و ابعاد پائینتر نمایش داده می‌شود. (شکل ۲) در نمونه‌ی بالا، جدول درخواستی دارای ۴۳۲۸ خانه‌ی نایمن است (این تعداد از خانه‌ها دارای فراوانی کمتر از چهار هستند). تعداد زیادی از این خانه‌های نایمن به دلیل تفصیل زیاد متغیرها حاصل شده‌اند. با توجه به تعداد خانه‌های نایمن هر متغیر در ابعاد مختلف می‌توان در مورد متغیر نامزد برای انجام بازکدگذاری تصمیم گرفت. بازکدگذاری یک متغیر موجب جزئیات کمتر و در عین حال تعداد ترکیبها نایمن کمتر می‌گردد. نرم‌افزار جدول درخواستی را به صورت زیر نمایش می‌دهد. (شکل ۳) در سمت چپ پنجره، جدول درخواست شده مشاهده می‌شود. خانه‌های نایمن با رنگ قرمز، خانه‌های ایمن با رنگ سیاه، خانه‌های پنهان‌سازی مکمل با رنگ آبی و خانه‌های خالی با علامت - مشخص شده‌اند (در اینجا به دلیل سیاه و سفید بودن چاپگر، رنگها مشخص نیستند). در پائین پنجره چند گزینه مشاهده می‌شود. با انتخاب گزینه‌ی Change View می‌توان جای

Table: FLT x SHAGHELIN | number

| | tot | 10 | 11 | 12 | 13 | 14 |
|------|-------|----|----|----|-----|-----|
| tot | 2,716 | 61 | 77 | 97 | 103 | 108 |
| 15 | 272 | 8 | 4 | 7 | 7 | 12 |
| 151 | 71 | 1 | - | 2 | 1 | 2 |
| 1512 | 3 | - | - | - | - | - |
| 1514 | 7 | - | - | - | - | - |
| 1515 | 18 | 1 | - | - | - | - |
| 1516 | 16 | - | - | 1 | - | 1 |
| 1517 | 1 | - | - | - | - | - |
| 1518 | 2 | - | - | - | - | - |
| 1519 | 24 | - | - | 1 | 1 | 1 |
| 152 | 21 | - | - | - | - | - |
| 1520 | 21 | - | - | - | - | - |
| 153 | 43 | 1 | - | 1 | - | 1 |
| 1531 | 27 | - | - | - | - | - |
| 1532 | 3 | - | - | - | - | - |
| 1533 | 13 | 1 | - | 1 | - | 1 |
| 154 | 127 | 6 | 4 | 4 | 6 | 9 |
| 1542 | 5 | - | - | 1 | - | - |
| 1543 | 7 | - | - | 1 | - | - |
| 1544 | 18 | 1 | 2 | - | 1 | 1 |
| 1545 | 8 | - | - | 1 | 1 | 1 |
| 1546 | 56 | 5 | 2 | 1 | 4 | 3 |
| 1547 | 4 | - | - | - | - | 1 |

Cell Information:

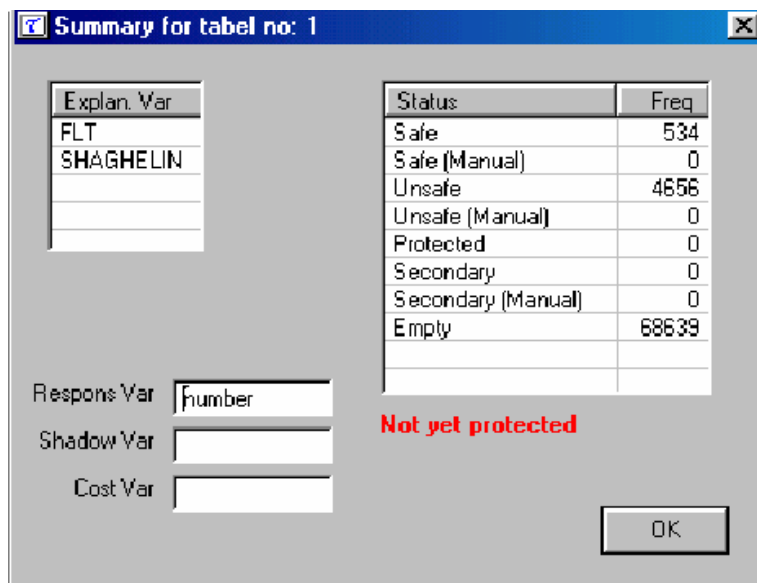
Value: 2,716
 Status: Safe
 Cost: 2,716
 Shadow: 2,716
 # contributions: 2716
 Top n of shadow: 1, 1, 1

Change status:

Suppress:
 HyperCube
 Opt/XPRESS
 Opt/CPLEX

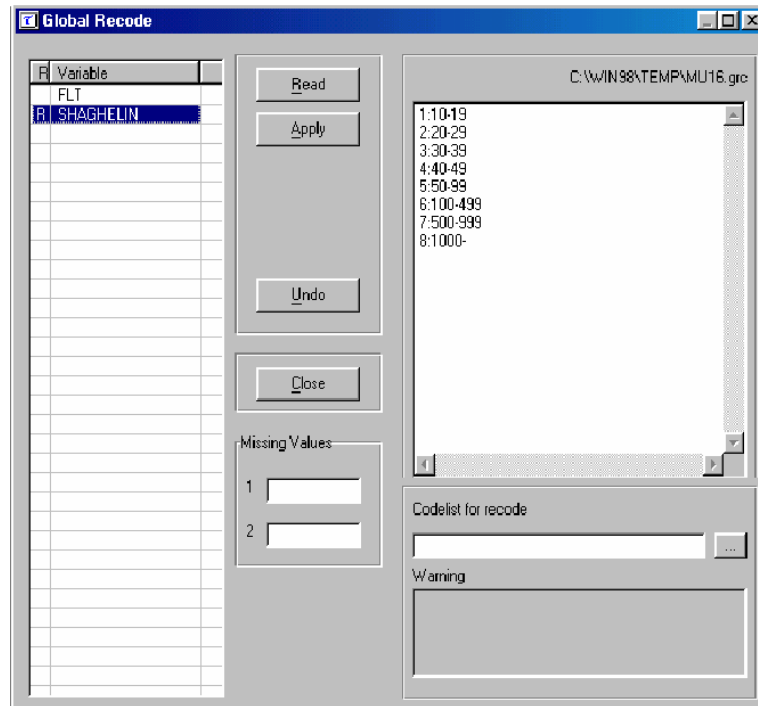
3 dig. separator
 Output View

شکل ۳: پنجره‌ی نشان دهنده‌ی جدول درخواست شده



شکل ۴: پنجره‌ی نشان دهنده‌ی وضعیت جدول درخواست شده (قبل از تعدیل)

متغیرهای سطری و ستونی جدول را عوض نمود. اگر جدول دارای بیشتر از دو بعد باشد، متغیر سوم به صورت لایه‌ای نشان داده می‌شود. گزینه‌ی Summary Table چشم‌اندازی از تعداد خانه‌های ایمن و نایمن بدست می‌دهد. در این حالت تعداد ۵۳۴ خانه ایمن، ۴۶۵۶ خانه نایمن اولیه و ۶۸۶۳۹ خانه خالی وجود دارد (شکل ۴). در قسمت بالای سمت راست شکل ۴ می‌توان اطلاعات مربوط به خانه‌ی انتخاب شده‌ی جدول را مشاهده نمود. در قسمت پائین تر گزینه‌هایی برای تغییر وضعیت خانه یا خانه‌هایی معین وجود دارد. اگر بنا به دلایلی مایل باشیم علاوه بر قاعده‌ی ایمنی به کار رفته، خانه‌ای را از حالت ایمن به نایمن یا برعکس تبدیل کنیم، می‌توان از این قسمت استفاده کرد. اکنون می‌توان این جدول را به لحاظ محرمانگی حفاظت کرد. یکی از روشهای قابل به کارگیری توسط این نرم‌افزار، روش بازکدگذاری است. با انتخاب گزینه‌ی Recode پنجره‌ای مانند پنجره‌ی زیر باز خواهد شد (شکل ۵) در سمت چپ نام متغیرهای تمیینی آمده است که می‌توان متغیر مورد نظر را انتخاب کرد. به عنوان مثال در شکل زیر متغیر تعداد شاغلان انتخاب شده است. اگر متغیر مورد نظر مانند این متغیر انتخاب شده، ناسلسله مراتبی باشد، قسمت سمت راست پنجره حالت ویرایش کدها را خواهد داشت. بدین صورت که رسته‌های جدید متغیر که شامل یک یا چند رسته قبلی‌اند، نوشته می‌شوند. می‌توان این بازکدگذاری را در فایل جداگانه‌ای نوشته و در اینجا توسط گزینه‌ی Read آن را بازخوانی نمود. با انتخاب گزینه‌ی Apply کدهای جدید روی متغیر به کار می‌روند. به عنوان مثال



شکل ۵: پنجره‌ی انتخاب روش بازکدگذاری برای متغیر تعداد شاغلان

The screenshot shows a software window titled "Table: FLT x SHAGHELIN | number". It contains a data table with columns labeled 'tot', '1', '2', '3', '4', '5', '6', '7', '8' and rows numbered from 15 to 1555. The 'tot' column has a value of 2,715. A 'Cell Information' panel on the right displays details for the selected cell (2,715), including its status (Safe), cost (2,715), shadow (2,715), and number of contributions (2,715). The panel also includes buttons for 'Change status' (Set to safe, Set to Unsafe, Set to Protected), 'Elicode', 'Suppress' (HyperCube, Dpt/Press, Dpt/CPlex), 'Write table', and 'Close'.

شکل ۶: پنجره‌ی نشاندهنده‌ی جدول درخواست شده (پس از بازکردگذاری متغیر تعداد شاغلان)

نتیجه‌ی بازکردگذاری روی متغیر تعداد شاغلان را در شکل بعد می‌بینیم. (شکل ۶) با این عمل تعداد خانه‌های ایمن به تعداد ۶۱۱ خانه افزایش می‌یابد و تعداد خانه‌های نایمن اولیه به ۶۲۴ و خانه‌های خالی به ۷۵۵ خانه کاهش می‌یابند.

اگر متغیر انتخابی از نوع سلسله مراتبی باشد در پنجره باز شده، کدها به صورت درختی در سمت راست نشان داده می‌شوند.

در این پنجره توسط گزینه‌ی Maximum level می‌توان میزان ماکسیمم جزئیات این متغیر را تعیین نمود. مثلاً برای متغیر نوع فعالیت ابتدا سطح دو در نظر گرفته می‌شود. نتیجه، در جدولی ظاهر می‌شود. در این حالت تعداد خانه‌های ایمن به تعداد ۳۴۰، خانه‌های نایمن اولیه به ۲۰۵ و تعداد خانه‌های خالی به ۲۳۵ خانه تغییر می‌یابند.

در صورتیکه نتیجه رضایتبخش نباشد، می‌توان دوباره به مرحله‌ی بازکردگذاری برگشت و کردگذاری جدیدی را انتخاب نمود. به عنوان مثال دوباره روی متغیر نوع فعالیت بازکردگذاری انجام می‌دهیم و یک رقم دیگر از کد آنرا حذف می‌کنیم که در نتیجه تعداد خانه‌های ایمن، خانه‌های نایمن اولیه و خانه‌های خالی به ترتیب به تعداد ۱۳۹، ۴۶ و ۶۵ خانه تغییر می‌یابند. به علت اینکه هنوز تعداد خانه‌های نایمن زیاد است در صورت لزوم می‌توان باز هم از مقدار جزئیات

کاست. به عنوان مثال اگر رسته‌های ۷ و ۸ متغیر تعداد شاغلان با هم ترکیب می‌شوند، تعداد خانه‌های ایمن به ۱۳۶، خانه‌های نایمن اولیه به ۳۵ و خانه‌های خالی به ۵۴ خانه تغییر می‌یابند. جدول ۶ نتیجه اعمال مراحل بالا است. در این جدول، با توجه به قاعده‌ی تشخیصی خانه‌های نایمن استفاده شده، خانه‌های نایمن اولیه به صورت خطدار مشخص شده‌اند. پس از پایان بازگذرداری متغیرها، باقیمانده‌ی خانه‌های نایمن را می‌توان توسط روش پنهان‌سازی محافظت نمود. در این مرحله به منظور حفاظت خانه‌های نایمن باید تعدادی خانه‌ی اضافه به عنوان پنهان‌سازی مکمل انتخاب شوند.

۶ نتیجه‌گیری

یکی از مهمترین وظایف نظام‌های آماری کشورها تولید آمارهای مورد نیاز کاربران آنها با مشخصات دقیق بودن، به موقع بودن، پوشش کامل به لحاظ موضوعی و جغرافیایی داشتن، ارزان بودن، و به نحو مناسب اطلاع رسانی شدن، است. به منظور حفظ محرمانگی اطلاعات فردی، نظام‌های آماری موظفند که به شیوه‌های اطلاع رسانی آماری کنند که هویت افراد افشاء نشود. این امر که غالباً جنبه قانونی نیز دارد، بهانه‌ای شده است که برخی از سازمانهای ملی آمار که وظیفه تولید و اطلاع رسانی آماری دارند از انتشار داده‌های خرد جلوگیری کنند و یا داده‌های حاصل از آمارگیریها را به گونه‌ای منتشر کنند که از آنها هویت افراد قابل افشاء باشد. به منظور انتشار ایمن داده‌ها چه به صورت خرد و چه به صورت جدول می‌توان از روش‌های کنترل افشاء داده‌های آماری استفاده نمود. این روش‌ها هم از افشای هویتها جلوگیری می‌کند و هم به سازمان‌های ملی آمار اجازه می‌دهد تا به صورت گسترده به اطلاع رسانی آماری مبادرت ورزند. برای این منظور لازم است از نرم‌افزارهای موجود همانند $\mu - ARGUS$ و $\tau - ARGUS$ استفاده شود.

مراجع

- [۱] محمد بردبار عشرت‌آبادی (۱۳۸۲)، رساله‌ی پایان‌نامه کارشناسی ارشد، فنون محدودسازی افشای داده‌های آماری و کاربرد آنها، دانشگاه علامه طباطبایی.
- [2] Cox, L. H. (1980), "Suppression Methodology and Statistical Disclosure Control," *Journal of the American Statistical Association*, Vol. 75, pp. 337-385
- [3] Cox L. H. (1987), "A Constructive Procedure for Unbiased Controlled Rounding," *Journal of the American Statistical Association*, Vol. 82, pp. 520-524.
- [4] Fienberg, S. E., Markov, U. G., and Steel, R. J. (1998), "Disclosure Limitation using Perturbation and related Methods for Categorical Data," *Journal of official statistics*, Vol. 14, pp. 485-502.
- [5] Gouweleeuw, J. M., Kooiman, p., Willenborg, L. C. R. J., and De Wolf, P. P. (1998), "Post Randomisation for Statistical Disclosure Control: Theory and Implementation," *Journal of Official Statistics*, Vol. 14, pp. 463-474.
- [6] Kooiman, P., Willenborg, L.C.R.J. and Gouweleeuw, J.M. (1997), PRAM: A Method for Disclosure Limitation of Microdata, *Statistics Netherlands*, The Netherlands, Research Paper no. 97 o 5
- [7] Willenborg, L. C. R. J. and De Waal, T. (1998), *Statistical Disclosure Control in Practice*, Springer- Verlag, New York.

برآورد واریانس به روش جک نایف^۱ در آمارگیریهای چندچارچوبی^۲

مرجان نورینی

مرکز آمار ایران

چکیده: روشهای معمول برآورد واریانس از قبیل روش بسط سری تیلور (روش دلتا^۳) در آمارگیریهای چندچارچوبی عموماً مستلزم محاسبه مشتقات جزئی بوده و این محاسبات با افزایش تعداد چارچوبها پیچیدهتر می شود. برآورد واریانس به روش جک نایف روش دیگری است که ضمن سهولت در محاسبه، موجب کاهش چشمگیری در اریبی می شود. در این مقاله به معرفی برآوردگرهای چندچارچوبی، استفاده از روش جک نایف در برآورد واریانس برآوردگرهای چندچارچوبی و مقایسه آن با روش بسط سری تیلور طی یک مطالعه شبیه سازی می پردازیم.

واژه های کلیدی: بررسی های پیچیده^۴، برآورد خطی واریانس^۵، روشهای باز نمونه گیری^۶

۱ مقدمه

چارچوب آماری یا فهرست واحدهای آمارگیری، اساس و مبنای یک طرح آمارگیری نمونه ای را تشکیل می دهد. گاه ممکن است چارچوبی که کلیه واحدهای جامعه مورد مطالعه را پوشش دهد در دسترس نباشد، اما امکان دستیابی به پوشش کامل، با تلفیقی از دو یا چند چارچوب فراهم شود. در چنین حالتی به منظور دسترسی به پوشش مناسب، از دو چارچوب یا بیشتر بطور همزمان استفاده می شود. گاهی نیز ممکن است یک چارچوب، پوشش کامل را برای جامعه مورد مطالعه فراهم کند، اما چارچوب ناقص دیگری موجود باشد که هزینه آمارگیری از آن کمتر از هزینه آمارگیری از چارچوب کامل باشد. در این شرایط به دلیل پایین تر بودن هزینه آمارگیری از این چارچوب، می توان با هزینه ای مشخص و ثابت، از دو چارچوب استفاده کرده و حجم نمونه را بزرگتر و کارایی را افزایش داد. گاهی نیز ممکن است یک چارچوب فهرستی کامل در دسترس باشد اما عملاً با گذشت زمانی نسبتاً طولانی بدلیل بروز تغییرات فراوان در آن، منبعی برای بروز خطاهای غیر نمونه گیری شود. از آنجا که یک فهرست ناحیه ای کمتر در معرض تغییرات می باشد ترکیب آن با یک چارچوب از اعضای جامعه که احتمالاً ناقص باشد می تواند نتایج مفیدی را حاصل نماید. چنین آمارگیریهایی تحت عنوان آمارگیریهای چندچارچوبی بکار می روند.

1) Jackknife 2) Multiple Frame Surveys 3) Delta Method 4) Complex Surveys 5) Linearization Varince Estimation 6) Resampling Mthods

شاید بتوان گفت اولین شالوده آمارگیریهای چندچارچوبی در سال (۱۹۴۹) با آمارگیری از فروشگاههای خرده فروشی که مجری آن دفتر سرشماری آمریکا بود گذاشته شد. سپس هارتلی^۷ (۱۹۶۲) نظریه مقدماتی چارچوبهای چندگانه را توسعه داد. او با فرض اینکه اجتماع چارچوبها، جامعه را پوشش می دهد، واحدهای جامعه را به زیرمجموعه های دو به دو ناسازگار شامل اجتماعها و اشتراکهای چارچوبهای مختلف تقسیم کرد. بعد از وی لاند^۸ (۱۹۶۸)، فولر و بورمیستر^۹ (۱۹۷۲)، وگل^{۱۰} (۱۹۷۵)، فورد و بوسکر^{۱۱} (۱۹۷۶)، بانکیر^{۱۲} (۱۹۸۶)، اسکینر^{۱۳} (۱۹۹۱)، اسکینر و راتو^{۱۴} (۱۹۹۶) و کالتون و اندرسون^{۱۵} (۱۹۸۶) مقالاتی را در این زمینه ارائه نموده اند. در آمارگیریهای چندچارچوبی، برآوردهای مختلفی برای برآورد مجموع جامعه پیشنهاد می شود. این برآوردها عموماً تابعی غیر خطی از مقادیر نمونه ای می باشند. برای برآورد واریانس این گونه برآوردها، روشهای مختلفی از جمله بسط سری تیلور (خطی سازی) پیشنهاد شده است. اسکینر و راتو (۱۹۹۶) روش بسط سری تیلور را پیشنهاد دادند. استفاده از این روش مستلزم محاسبه مشتقات جزئی بوده و این محاسبات با افزایش تعداد چارچوبها پیچیده تر می شود. لیکن برآورد واریانس به روش جک نایف ضمن سهولت در محاسبه، موجب کاهش چشمگیری در آریبی می شود. در این مقاله ضمن استفاده از روش جک نایف در برآورد واریانس به مقایسه آن با روش بسط سری تیلور می پردازیم. بدین منظور در بخش ۲ برآوردهای چندچارچوبی را معرفی می کنیم. در بخش ۳ ابتدا شرح مختصری از روش جک نایف را ارائه داده و سپس به استفاده از آن در آمارگیریهای چندچارچوبی می پردازیم. در بخش ۴ نتایج حاصل از مقایسه و بررسی برآوردها و برآوردهای واریانس از طریق شبیه سازی ارائه می گردد.

۲ برآوردهای مجموع جامعه

برای سادگی فرض کنید دو چارچوب A و B موجود است که هر دو ناقص و دارای واحدهای مشترک با یکدیگر هستند بطوریکه مجموع آنها رویهم، کل جامعه مورد مطالعه را پوشش می دهد. از چارچوبهای A و B ، سه $(2^2 - 1)$ حوزه^{۱۶} دو به دو ناسازگار به دست می آید.

حوزه a : شامل واحدهایی است که فقط در چارچوب A می باشند. $a = A \cap B^c$

حوزه b : شامل واحدهایی است که فقط در چارچوب B می باشند. $b = A^c \cap B$

حوزه ab : شامل واحدهایی است که در هر دو چارچوب می باشند. $ab = A \cap B$

(c نشان دهنده مکمل مجموعه می باشد)

N_A و N_B تعداد واحدها در چارچوبهای A و B ، N_a ، N_b و N_{ab} تعداد واحدهای موجود در حوزه های a ، b و ab می باشند.

7) Hartley 8) Lund 9) Fuller and Burmeister 10) Vogel 11) Ford and Boseker 12) Bankier 13) skinner 14) Rao 15) Kalton and Anderson 16) domain

دو نمونه مستقل s_B و s_A بر اساس طرحهای نمونه‌گیری احتمالی $p_A(s_A)$ و $p_B(s_B)$ به اندازه n_B و n_A از دو چارچوب فوق گرفته می‌شود. احتمال شمول نمونه حاصل از چارچوب A ، $\pi_i^A = p\{i \in s_A\}$ و احتمال شمول نمونه حاصل از چارچوب B ، $\pi_i^B = p\{i \in s_B\}$ می‌باشد. براساس نمونه‌های مستقل فوق، n_a و n_{ab}^A ، تعداد واحدهای نمونه‌گیری حاصل از چارچوب A می‌باشند که به ترتیب در حوزه‌های a و ab ، قرار دارند. به همین ترتیب n_b و n_{ab}^B نیز تعداد واحدهای نمونه‌گیری حاصل از چارچوب B می‌باشند که به ترتیب در حوزه‌های b و ab هستند. با فرض اینکه $Y_a, Y_b, Y_{ab}, \mu_a, \mu_b, \mu_{ab}$ به ترتیب مجموع و میانگین جامعه در حوزه‌های a, b, ab باشند، داریم:

$$Y = Y_a + Y_{ab} + Y_b \quad (۱)$$

هدف برآورد Y است.

چندین برآوردگر نقطه‌ای تحت عنوان برآوردگر دو چارچوبی برای برآورد Y پیشنهاد شده که همگی آنها به فرم $\hat{Y} = \hat{Y}_a + \hat{Y}_{ab} + \hat{Y}_b$ می‌باشند. هر یک از این برآوردگرها بسته به اینکه اطلاعات حاصل از دو نمونه برای برآورد Y چگونه با هم ترکیب می‌شوند با هم تفاوتی دارند. با فرض اینکه N_B و N_A معلوم بوده و $N_a > 0$ و $N_b > 0$ ، وزنهای نمونه‌گیری w_i^A و w_i^B عبارتند از:

$$w_i^A = N_A \left[\pi_i^A \sum_{j \in s_A} \left(\frac{1}{\pi_j^A} \right) \right]^{-1}$$

$$w_i^B = N_B \left[\pi_i^B \sum_{j \in s_B} \left(\frac{1}{\pi_j^B} \right) \right]^{-1}$$

دو متغیر نشانگر زیر را برای هر یک از چارچوبهای A و B در نظر می‌گیریم:

$$\delta_A(i) = \begin{cases} 1 & \text{واحد } i \text{ متعلق به چارچوب } A \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

$$\delta_B(i) = \begin{cases} 1 & \text{واحد } i \text{ متعلق به چارچوب } B \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

در نتیجه برآوردگرهای مجموع در سه حوزه a و b و ab بقرار زیر می‌باشند:

$$\hat{Y}_a^A = \sum_{i \in s_A} w_i^A (1 - \delta_B(i)) y_i$$

$$\hat{Y}_{ab}^A = \sum_{i \in s_A} w_i^A \delta_B(i) y_i$$

$$\hat{Y}_b^B = \sum_{i \in s_B} w_i^B (1 - \delta_A(i)) y_i$$

$$\hat{Y}_{ab}^B = \sum_{i \in s_B} w_i^B \delta_A(i) y_i$$

برآوردگر اندازه هر یک از حوزه‌ها نیز بصورتی مشابه با قرارداد $y_i = 1$ در تعاریف \hat{Y}_{ab}^A و \hat{Y}_b^B بدست می‌آیند. همچنین دو برآوردگر دیگر را نیز بصورت زیر تعریف می‌کنیم:

$$\begin{aligned}\hat{Y}_{ab}(\theta) &= \theta \hat{Y}_{ab}^A + (1 - \theta) \hat{Y}_{ab}^B \\ \hat{N}_{ab}(\theta) &= \theta \hat{N}_{ab}^A + (1 - \theta) \hat{N}_{ab}^B\end{aligned}$$

۱.۲ برآوردگرهای هارتلی و فولر - بورمیستر

هارتلی در سال (۱۹۶۲) برآوردگر دوچارچوبی زیر را برای برآورد مجموع جامعه پیشنهاد کرد:

$$\hat{Y}_H(\theta) = \hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}(\theta) \quad (۲)$$

فولرو بورمیستر هم در سال (۱۹۷۲) برآوردگر زیر را پیشنهاد دادند:

$$\hat{Y}_{FB}(\beta_1, \beta_2) = \hat{Y}_a^A + \hat{Y}_b^B + \hat{Y}_{ab}(\beta_1) + \beta_2(\hat{N}_{ab}^A - \hat{N}_{ab}^B) \quad (۳)$$

مقادیر بهینه پارامترهای θ ، β_1 و β_2 به ترتیب واریانس برآوردگرهای $\hat{Y}_H(\theta)$ و $\hat{Y}_{FB}(\beta_1, \beta_2)$ را می‌نیم کرده و بنابراین به کواریانسهای \hat{Y}_{ab}^A و \hat{Y}_{ab}^B بستگی دارند. اما در عمل این کواریانسها مجهول بوده و باید از روی نمونه برآورد شوند. در نتیجه $\hat{Y}_H(\hat{\theta}_H)$ و $\hat{Y}_{FB}(\hat{\beta}_{FB})$ بطور کلی توابعی خطی از y با وزنهای یکسان برای همه متغیرها نخواهند بود و برای هر متغیر پاسخ وزنها باید جداگانه محاسبه شوند. برای مثال اگر \hat{Y}_{H_1} ، \hat{Y}_{H_2} و \hat{Y}_{H_3} برآوردگرهای هارتلی برای تعداد کل بیمارانی که دچار تنگی نفس در گروههای سنی $16-17$ ، $17-45$ و بالای 45 سال باشد، آن گاه $\hat{Y}_{H_1} + \hat{Y}_{H_2} + \hat{Y}_{H_3}$ لزوماً مساوی برآورد هارتلی تعداد کل بیماران دچار تنگی نفس در جامعه نخواهد بود. در نتیجه فولرو بورمیستر برآن شدند برآوردگری ارائه دهند که از یک مجموعه وزنهای یکسان برای همه متغیرهای پاسخ استفاده کند. اما برآوردگری که ارائه دادند تحت طرح نمونه‌گیری تصادفی ساده طراحی شده بود و مستقیماً با طرحهای نمونه‌گیری پیچیده بکار نمی‌رفت. چون عموماً بصورتی ناسازگار با این نوع نمونه‌گیریها طراحی شده بود. اسکینر (۱۹۹۱) از این برآوردگر تفسیری بصورت یک برآوردگر حداکثر درستنمایی داشت.

۲.۲ برآوردگر شبه حداکثر درستنمایی

اسکینر وراثو (۱۹۹۶) برآوردگری ارائه دادند که علاوه بر استفاده از یک مجموعه یکسانی از وزنها، اصلاح شده یک برآوردگر حداکثر درستنمایی تحت طرح نمونه‌گیری تصادفی ساده برای رسیدن به برآوردگری سازگار با طرحهای نمونه‌گیری پیچیده بود. آنها برآوردگر شبه حداکثر درستنمایی به

صورت زیر را ارائه دادند:

$$\hat{Y}_{PML}(\theta) = \frac{N_A - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_a^A} \hat{Y}_a^A + \hat{N}_{ab}^{PML}(\theta) \frac{\hat{Y}_{ab}(\theta)}{\hat{N}_{ab}(\theta)} + \frac{N_B - \hat{N}_{ab}^{PML}(\theta)}{\hat{N}_b^B} \hat{Y}_b^B \quad (4)$$

بطوریکه $\hat{N}_{ab}^{PML}(\theta)$ تابعی از \hat{N}_{ab}^A و \hat{N}_{ab}^B و θ بوده و کوچکترین ریشه معادله درجه دو به فرم زیر است:

$$\left(\frac{\theta}{N_B} + \frac{1-\theta}{N_A}\right)x^2 - \left(1 + \frac{\theta}{N_B}\hat{N}_{ab}^A + \frac{1-\theta}{N_A}\hat{N}_{ab}^B\right)x + \hat{N}_{ab}(\theta) = 0$$

این برآوردگر بر خلاف برآوردگرهای هارتلی و فولر - بورمیستر، تابعی خطی از y است. اسکینر و راثو (۱۹۹۶) انتخاب $\theta = \theta_P$ را برای می نیم کردن واریانس مجانبی $\hat{N}_{ab}^{PML}(\theta)$ با

$$\theta_P = \frac{N_a N_B V(\hat{N}_{ab}^B)}{N_a N_B V(\hat{N}_{ab}^B) + N_b N_A V(\hat{N}_{ab}^A)} \quad (5)$$

پیشنهاد کردند.

در عمل N_a و N_b و واریانسها در رابطه (۵) مجهول بوده و باید از روی داده ها برآورد شود. در نتیجه برآوردگر حاصل عبارت از $\hat{Y}_{PML}(\hat{\theta}_P)$ است.

۳.۲ برآوردگرهای چارچوب منفرد

بانکیر (۱۹۸۶)، کالتون و اندرسون (۱۹۸۶) و اسکینر (۱۹۹۱) برای برآورد مجموع جامعه با در نظر گرفتن مشاهدات بگونه ای که آنها از یک چارچوب منفرد با وزنه های اصلاح شده برای مشاهدات در حوزه متداخل ab نمونه گیری شده اند، برآوردی ارائه دادند. وزنه های اصلاح شده برای برآوردگرهای تک چارچوبی اسکینر (۱۹۹۱) و کالتون و اندرسون (۱۹۸۶) به شناسایی واحدهای مشابه در نمونه ها نیازی ندارد. این وزنه ها عبارتند از:

$$w_i = \begin{cases} (\pi_i^A + \pi_i^B)^{-1} & i \in ab \\ (\pi_i^A)^{-1} & i \in a \\ (\pi_i^B)^{-1} & i \in b \end{cases}$$

آنها برآوردگر زیر را برای برآورد مجموع جامعه پیشنهاد کردند:

$$\hat{Y}_{SF} = \sum_{i \in s_A} w_i y_i + \sum_{i \in s_B} w_i y_i \quad (6)$$

همانطور که بانکیر (۱۹۸۶) اشاره کرد، این برآوردگر قابل تعمیم به بیش از دو چارچوب می باشد. این برآوردگر از هیچ گونه اطلاعات کمکی در ارتباط با N_A و N_B استفاده نمی کند. از آنجائیکه استفاده از اطلاعات کمکی باعث کاهش واریانس و افزایش دقت برآوردگر [بانکیر (۱۹۸۶)] می شود، اسکینر و راثو (۱۹۹۶) دو روش برای تصحیح این برآوردگر ارائه دادند که عبارتند از: الف) روش تعدیل نسبی^{۱۷} [بانکیر (۱۹۸۶)]

ب) روش رگرسیونی اسکینر و راثو پیرو اثبات قضیه (۱) اسکینر (۱۹۹۱) نشان دادند که فرآیند تعدیل به برآوردگر تجربی زیر همگراست:

$$\hat{Y}_{SFrake} = \frac{N_A - \hat{N}_{ab}^{rake}}{\hat{N}_a} \hat{Y}_a + \frac{\hat{N}_{ab}^{rake}}{\hat{N}_{abS}} \hat{Y}_{abS} + \frac{N_B - \hat{N}_{ab}^{rake}}{\hat{N}_b} \hat{Y}_b \quad (۷)$$

که در آن:

$$\hat{Y}_{abS} = \sum_{s_A} w_i \delta_i^B y_i + \sum_{s_B} w_i \delta_i^A y_i$$

$$\hat{N}_{abS} = \sum_{s_A} w_i \delta_i^B + \sum_{s_B} w_i \delta_i^A$$

و \hat{N}_{ab}^{rake} کوچکترین ریشه معادله درجه دو بفرم زیر می باشد:

$$\hat{N}_{abS} x^2 - [\hat{N}_{abS}(N_A + N_B) + \hat{N}_{aS}^A \hat{N}_{bS}^B] x + \hat{N}_{abS} N_A N_B = 0 \quad (۸)$$

تصحیح با N_A و N_B معلوم از طریق روش رگرسیونی برآوردگر زیر را نتیجه می دهد:

$$\hat{Y}_{SFreg} = \hat{Y}_S + \hat{\beta}^T [N_A - \hat{N}_S^A, N_B - \hat{N}_S^B] \quad (۹)$$

بطوریکه مقدار بهینه $\hat{\beta}$ عبارتست از:

$$\hat{\beta}_S^T = -Cov\{[\hat{N}_{AS}/V(\hat{N}_{AS}), \hat{N}_{BS}/V(\hat{N}_{BS})], \hat{Y}_S\}$$

۳ برآورد واریانس

اسکینر و راثو (۱۹۹۶) با استفاده از بسط سری تیلور، روشی را برای برآورد واریانس \hat{Y}_{PML} ارائه دادند. در این بخش ابتدا برآورد واریانس به روش جک نایف برای برآوردگرهای دو چارچوبی را ارائه داده و سپس نشان می دهیم که برآوردگر جک نایف واریانس بطور مجانبی هم ارز با برآوردگر خطی واریانس می باشد. برای سادگی حالت دوچارچوبی را در نظر می گیریم اما این نتایج براحتی قابل تعمیم به آمارگیریهای چندچارچوبی می باشد.

17) Raking Ratio

نخستین بار کوئنولی^{۱۸} (۱۹۵۹) جک نایف را به عنوان روشی بر اساس حذف هر بار یک مشاهده از مجموعه داده‌های اولیه و محاسبه مجدد برآوردگر با استفاده از بقیه داده‌ها برای کاهش اریبی برآوردگر ضریب همبستگی پیاپی^{۱۹} معرفی نمود. در مقاله‌ای در سال (۱۹۵۶) این روش را تعمیم داده و خواص کلی کاهش اریبی آن را در یک جامعه متناهی بیان کرد. سپس توکی^{۲۰} (۱۹۵۸) پیشنهاد کرد که این روش علاوه بر کاهش اریبی، می‌تواند برای برآورد واریانس نیز استفاده شود.

فرض کنید چارچوب‌های A و B به ترتیب دارای H و L طبقه باشند بطوریکه طبقات h و l هر یک از آنها به ترتیب شامل N_h^A و N_l^B واحد بوده و در مجموع \tilde{N}_h^A و \tilde{N}_l^B واحد نمونه‌گیری اولیه^{۲۱} (psu) دارند که به ترتیب \tilde{n}_h^A و \tilde{n}_l^B واحد از آنها نمونه‌گیری می‌شود. بطوریکه انتخاب می‌شود. وزنهای طبقات h و l در چارچوبهای A و B را به ترتیب با $W_h^A = N_h^A/N_A$ و $W_l^B = N_l^B/N_B$ نشان می‌دهیم. لازم به یادآوری است که اغلب در نمونه‌گیریهایی چندمرحله‌ای، نمونه بدون جایگذاری از واحدهای نمونه‌گیری اولیه با احتمال متناسب با اندازه (به دلیل کارایی بالای نمونه‌گیری بدون جایگذاری نسبت به نمونه‌گیری با جایگذاری) بعمل آمده و سپس در مرحله برآورد واریانس برای صرفه‌جویی در زمان و هزینه و اجتناب از محاسبات بالقوه پر زحمت احتمالات شمول توام، نمونه‌گیری را بگونه‌ای در نظر می‌گیرند که گویی با جایگذاری انتخاب شده است. این تقریب عموماً منجر به بیش برآورد واریانس برآوردگر مجموع می‌شود اما در صورتی که کسرهای نمونه‌گیری مرحله اول بزرگ نباشد، اریبی نسبی بزرگ نخواهد بود. برآوردگر واریانس جک نایف و دیگر برآوردگرهای واریانس بازنمونه‌گیری از این روش در برآورد واریانس استفاده می‌کنند.

حال فرض می‌کنیم که نمونه‌هایی به اندازه (≥ 2) \tilde{n}_h^A و \tilde{n}_l^B به ترتیب با احتمالات شمول متناسب با اندازه، $\tilde{\pi}_{hi}^A = \tilde{n}_h^A p_{hi}^A$ و $\tilde{\pi}_{lj}^B = \tilde{n}_l^B p_{lj}^B$ از چارچوبهای A و B انتخاب می‌شود. بطوریکه p_{hi}^A و p_{lj}^B به ترتیب برابر با احتمال انتخاب واحد i ام و j ام در طبقات h و l متناسب با اندازه psu ها بوده و $\sum_i p_{hi}^A = 1$ و $\sum_j p_{lj}^B = 1$ و $p_{hi}^A = u_i / \sum_A u_i$ و $p_{lj}^B = u_j / \sum_B u_j$ که u_i اندازه تقریبی i امین psu می باشد. با در نظر گرفتن دو بردار $\mathbf{A} = (Y_a, Y_{ab}, N_{ab}, N_A)^T$ و $\mathbf{B} = (Y_b, Y_{ab}, N_{ab}, N_B)^T$ برآورد \mathbf{A} عبارت خواهد بود

18) Quenouille 19) Serial Correlation Coefficient 20) Tukey 21) Primary Sampling Units

از:

$$\begin{aligned} \hat{A} &= \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} \hat{A}_{hi} / \tilde{\pi}_{hi}^A = \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} \hat{A}_{hi} / \tilde{n}_h^A p_{hi}^A = \sum_{h=1}^H \sum_{i=1}^{\tilde{n}_h^A} N_h^A a_{hi} / \tilde{n}_h^A \\ &= \sum_{h=1}^H N_h^A \bar{a}_h \end{aligned}$$

که \hat{A}_{hi} برآورد نارایب بردار مجموع جوامع در i امین psu نمونه از طبقه h ام بوده و $a_{hi} = \hat{A}_{hi} / (N_h^A p_{hi}^A)$ می باشد. برآوردگر $\hat{B} = \sum_{l=1}^L \sum_{j=1}^{\tilde{n}_l^B} N_l^B b_{lj} / \tilde{n}_l^B$ و b_{lj} نیز بصورتی مشابه تعریف می شوند. تحت فرضیات نمونه گیری با جایگذاری، a_{hi} ها برآوردگرهای نارایب مستقل و هم توزیع میانگین جامعه در طبقه h ام چارچوب A می باشند در حالیکه برای $h \neq h'$ ، a_{hi} و $a_{h'i'}$ مستقل بوده اما لزوماً هم توزیع نیستند. بطور مشابه b_{lj} ها هم برآوردگرهای نارایب مستقل و هم توزیع بردار میانگین جامعه در طبقه l ام چارچوب B می باشند.

پارامتری به فرم $\tau = g(\bar{A}, \bar{B})$ که تابعی از میانگینهای جامعه $\bar{A} = A/N$ و $\bar{B} = B/N$ است را در نظر می گیریم. میانگینهای $\bar{A} = A/N$ و $\bar{B} = B/N$ بصورت $\hat{A} = \sum_{l=1}^L W_l^B \bar{b}_l$ و $\hat{A} = \sum_{h=1}^H W_h^A \bar{a}_h$ برآورد شده و τ نیز بصورت $\hat{\tau} = g(\hat{A}, \hat{B})$ برآورد می شود.

برآوردگرهای دوچارچوبی می توانند در قالب $\hat{\tau}$ با در نظر گرفتن مجموع جامعه به صورت $Y = NY = Ng(\bar{A}, \bar{B})$ بیان شوند که میانگین \bar{Y} عبارتست از:

$$g(\bar{A}, \bar{B}) = \bar{Y} = (N_A/N)(\bar{A}_1 + \theta \bar{A}_2) + (N_B/N)(\bar{B}_1 + (1 - \theta)\bar{B}_2) \quad (10)$$

در این صورت برآوردگر هارتلی عبارت از $\hat{Y}_H(\theta_H) = Ng(\hat{A}, \hat{B})$ خواهد بود. برآوردگر PML نیز به ازای θ ثابت، می تواند بصورتی مشابه بعنوان تابعی از \hat{A} و \hat{B} و با توجه به اینکه $\hat{N}_{ab}^{PML}(\theta)$ تابعی از \hat{A} و \hat{B} است، و سپس با استفاده از رابطه (۴) بیان شود.

ماتریس واریانس - کواریانس بردار \hat{A} برابر با $\sum_{h=1}^H (W_h^A)^2 \sum_{i=1}^{\tilde{n}_h^A} a_{hi}^2 / \tilde{n}_h^A$ بوده که با $S_A = \sum_{h=1}^H (W_h^A)^2 S_h^A / \tilde{n}_h^A$ برآورد می شود. بطور مشابه واریانس \hat{B} یعنی \sum_B نیز بصورت $S_B = \sum_{l=1}^L (W_l^B)^2 S_l^B / \tilde{n}_l^B$ برآورد شده که S_h^A و S_l^B به ترتیب برآوردگرهای

واریانس $\sqrt{\tilde{n}_l^B \bar{b}_l} \sqrt{\tilde{n}_h^A \bar{a}_h}$ می‌باشند و عبارتند از:

$$S_h^A = (\tilde{n}_h^A - 1) \sum_{i=1}^{\tilde{n}_h^A} (a_{hi} - \bar{a}_h)(a_{hi} - \bar{a}_h)^T$$

$$S_l^B = (\tilde{n}_l^B - 1) \sum_{j=1}^{\tilde{n}_l^B} (b_{lj} - \bar{b}_l)(b_{lj} - \bar{b}_l)^T$$

برای اثبات سازگاری برآوردگر جک نایف واریانس و همچنین هم‌ارزی مجانبی آن با برآوردگر خطی واریانس به شروط زیر احتیاج داریم .

الف) فرض می‌کنیم به ازای هر h و l ، $W_h^A \tilde{n}_A / \tilde{n}_h^A = O(1)$ و $W_l^B \tilde{n}_B / \tilde{n}_l^B = O(1)$ باشد. در مجموع فرض می‌کنیم $\sum_h W_h^A \tilde{n}_h^A = O(1)$ و $\sum_l W_l^B \tilde{n}_l^B = O(1)$ باشد.

ب) فرض کنید $g_A(\tilde{\mathbf{a}}, \tilde{\mathbf{b}})$ برداری q بعدی از مشتقهای اول تابع $g(\cdot)$ با توجه به مولفه‌های \mathbf{a} و $g_B(\tilde{\mathbf{a}}, \tilde{\mathbf{b}})$ نیز برداری r بعدی از مشتقهای اول $g(\cdot)$ با توجه به مولفه‌های \mathbf{b} باشد که مقدار هر یک در $\tilde{\mathbf{a}}$ و $\tilde{\mathbf{b}}$ محاسبه شده‌اند. نیز فرض کنید $g_A''(\tilde{\mathbf{a}}, \tilde{\mathbf{b}})$ یک ماتریس $q \times q$ بعدی از مشتقهای دوم $\partial^2 g / (\partial a_j \partial a_k)$ و $g_B''(\tilde{\mathbf{a}}, \tilde{\mathbf{b}})$ یک ماتریس $r \times r$ بعدی از مشتقهای دوم $\partial^2 g / (\partial b_j \partial b_k)$ در $\tilde{\mathbf{a}}$ و $\tilde{\mathbf{b}}$ باشند. g_A'' و g_B'' در یک همسایگی از $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ پیوسته و کراندار باشند.

ج) $\tilde{n} = \tilde{n}_A + \tilde{n}_B$ بطوریکه $\tilde{n}_A / \tilde{n} \rightarrow k \in (0, 1)$

شروط فوق توسط راتو و وو (۱۹۸۵) برای تحقیق و بررسی خواص برآوردگرهای واریانس در طرح نمونه‌گیری طبقه‌بندی چندمرحله‌ای مورد استفاده قرار گرفته است. شرط (ج) نیز تضمین می‌کند که نمونه حاصل از یک چارچوب بطور مجانبی روی نمونه دیگر هیچ تأثیری نمی‌گذارد.

قضیه ۱ فرض می‌کنیم شروط (الف) -- (ج) برقرار باشند. آن‌گاه:

$$Var(\hat{\tau}) = g_A^T(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) \sum_{AgA} (\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) + g_B^T(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) \sum_{BgB} (\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) + o(\tilde{n}^{-1}) \quad (11)$$

و برآوردگر خطی واریانس فوق عبارتست از:

$$v_L(\hat{\tau}) = g_A^T(\hat{\tilde{\mathbf{A}}}, \hat{\tilde{\mathbf{B}}}) S_{AgA}(\hat{\tilde{\mathbf{A}}}, \hat{\tilde{\mathbf{B}}}) + g_B^T(\hat{\tilde{\mathbf{A}}}, \hat{\tilde{\mathbf{B}}}) S_{BgB}(\hat{\tilde{\mathbf{A}}}, \hat{\tilde{\mathbf{B}}}) \quad (12)$$

فرض می‌کنیم $\hat{\tau}_{(hi)}^A$ برآوردگری به فرم $\hat{\tau}$ باشد که بعد از حذف مشاهدات i امین psu نمونه طبقه h ام به دست آمده و $\hat{\tau}_{(hi)}^A = g(\hat{\mathbf{A}}_{(hi)}, \hat{\mathbf{B}})$ که $\hat{\mathbf{A}}_{(hi)}$ نیز برآوردگر $\tilde{\mathbf{A}}$ بعد از حذف i امین psu نمونه طبقه h ام در چارچوب A می‌باشد. بطور مشابه $\hat{\tau}_{(lj)}^B = g(\hat{\mathbf{A}}, \hat{\mathbf{B}}_{(lj)})$ که

$\hat{B}_{(lj)}$ برآورد \bar{B} بعد از حذف psu نمونه طبقه l ام چارچوب B است. برآوردگر جک نایف واریانس $\hat{\tau}$ (با توجه به استقلال نمونه‌ها) عبارتست از:

$$v_J(\hat{\tau}) = \sum_{h=1}^H \frac{\tilde{n}_h^A - 1}{\tilde{n}_h^A} \sum_{i=1}^{\tilde{n}_h^A} (\hat{\tau}_{(hi)}^A - \hat{\tau})^2 + \sum_{l=1}^L \frac{\tilde{n}_l^B - 1}{\tilde{n}_l^B} \sum_{j=1}^{\tilde{n}_l^B} (\hat{\tau}_{(lj)}^B - \hat{\tau})^2 \quad (13)$$

قضیه ۲ فرض می‌کنیم شروط (الف) -- (ج) برقرار باشند. آن گاه:

$$v_J(\hat{\tau}) = v_L(\hat{\tau}) + o_p(\tilde{n}^{-1}) \quad (14)$$

برآوردگر جک نایف واریانس فوق در حالتی که $\hat{\tau}$ تابع همواری از \hat{A} و \hat{B} باشد مورد بررسی قرار گرفته و طی قضایای فوق سازگاریش و نیز هم ارزیش با برآوردگر خطی واریانس ثابت می‌شود. در بین برآوردگرهای دو چارچوبی، برآوردگرهای SF و $SFrake$ بصورت توابع همواری از میانگین جوامع قابل بیان هستند. دیگر برآوردگرها نیز مادامی که پارامترهای $\theta_H, \theta_p, \beta_{SB}$ و β_{FB} ثابت بوده و از روی داده‌ها برآورد نشوند، نیز توابع همواری از میانگین جوامع خواهند بود. بنابراین قضیه (۲) بیان می‌دارد که برآوردگر جک نایف واریانس برای شکل بهینه هر برآوردگر سازگار است.

۱.۳ جک نایف کامل و اصلاح شده

برآوردگرهای $\hat{\theta}_H, \hat{\theta}_p, \hat{\beta}_{SB}$ و $\hat{\beta}_{FB}$ همگی توابعی از واریانس - کواریانس میانگین جوامع، S_A و S_B هستند و بطور کلی نمی‌توانند بصورت توابع قابل همواری از میانگین در نمونه‌گیری طبقه‌بندی بیان شوند (واریانسها معمولاً نامعلوم بوده و باید از روی نمونه برآورد شوند) بنابراین قضیه (۲) همیشه نمی‌تواند بطور مستقیم با برآوردگرهایی که توابعی از S_A و S_B هستند بکار رود. مشکل دیگر که در استفاده از جک نایف در خیلی از نمونه‌های طبقه بندی رخ می‌دهد، هنگام محاسبه $S_{(lj)}^B$ و $S_{(hi)}^A$ می‌باشد، چون

$$S_{(hi)}^A = S_A + \frac{(W_h^A)^2}{\tilde{n}_h^A - 2} \left[\frac{2S_h^A}{\tilde{n}_h^A} - \frac{\tilde{n}_h^A}{(\tilde{n}_h^A - 1)^2} (a_{hi} - \bar{a}_h)(a_{hi} - \bar{a}_h)^T \right] \quad (15)$$

در صورتی که $\tilde{n}_h^A = 2$ باشد قابل محاسبه نیست. بطور مشابه $S_{(lj)}^B$ نیز وقتی که $\tilde{n}_l^B = 2$ قابل محاسبه نیست. برای رفع این مشکل راتو و لهر (۱۹۹۷) یک برآورد جک نایف اصلاح شده در حالتی که $\tilde{n}_h^A = 2$ یا $\tilde{n}_l^B = 2$ ارائه دادند. جک نایف (جک نایف اصلاح شده برای طرح نمونه‌گیری با ۲ واحد نمونه‌گیری اولیه در هر طبقه که بعداً به معرفی آن می‌پردازیم) حتی زمانی که برآوردگرها به S_A و S_B بستگی دارند، برآورد سازگاری از واریانس می‌دهد.

در $\hat{\tau}$ فرض بر این است که پارامترهای $\theta_H, \theta_P, \beta_{FB}$ و $\beta_{S\psi}$ ثابت بوده و از روی نمونه برآورد نمی‌شوند. اما اگر این پارامتر ثابت نبوده و به S_A و S_B بستگی داشته باشند باید اثر این برآوردگرها در $\hat{\tau}$ منظور شود. فرض می‌کنیم که

$$\hat{\tau} = g(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = f(\hat{\mathbf{A}}, \hat{\mathbf{B}}, \beta)$$

$$\beta = [\mathbf{E}_A + \mathbf{E}_B]^{-1}[e_A + e_B]$$

که هر عنصر ماتریس $p \times p$ بعدی \mathbf{E}_A و هر عنصر بردار p بعدی e_A ترکیب خطی از عناصر \sum_A هستند و به همین ترتیب برای \mathbf{E}_B و e_B . آن گاه برآوردگرهای $\hat{Y}_{PML}(\hat{\theta}_p)$ و $\hat{Y}_H(\hat{\theta}_H)$ ، $\hat{Y}_{FB}(\hat{\beta}_{FB})$ و $\hat{Y}_{S\psi}(\hat{\beta}_{S\psi})$ همگی به فرم $N\hat{\zeta} = Nf(\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\beta})$ می‌باشند بطوریکه $\hat{\beta}$ برآورد کواریانسهای جامعه یعنی S_A و S_B را در β جایگزین می‌کند. برای مثال اگر بردارهای \mathbf{A} و \mathbf{B} را به صورت زیر در نظر بگیریم:

$$\mathbf{A} = (Y_a, Y_{ab}, N_{ab}, N_A)^T \quad \mathbf{B} = (Y_b, Y_{ab}, N_{ab}, N_B)^T$$

آن گاه برآوردگر هارتلی $\hat{Y}_H(\hat{\theta}_H)$ را می‌توان به صورت زیر نوشت:

$$\hat{Y}_H(\hat{\theta}_H) = \hat{Y}_H(\theta_H) + (\hat{\beta}_H - \beta_H)[N_A(\hat{\mathbf{A}}_{\psi} - \bar{\mathbf{A}}_{\psi}) - N_B(\hat{\mathbf{B}}_{\psi} - \bar{\mathbf{B}}_{\psi})]$$

بطوریکه مقدار بهینه β_H برابر است با:

$$\beta_H = -\frac{k^2 \sum_A(1, \psi) - \sum_B(1, \psi) - \sum_B(\psi, \psi)}{k^2 \sum_A(\psi, \psi) + \sum_B(\psi, \psi)}, \quad k = N_A/N_B$$

$\sum_A(1, \psi)$ عبارت از $(1, \psi)$ امین عنصر \sum_A می‌باشد. برای برآوردگر شبه حداکثر درست‌نمایی،

پارامتر مناسب عبارت از $\sum_B(3, 3)/\sum_A(3, 3)$ می‌باشد. $\beta_p = k^3 \sum_A(3, 3)/\sum_B(3, 3)$

فرض کنید $\hat{\zeta}_{(hi)}^A$ برآوردگری به فرم $\hat{\zeta}$ باشد که بعد از حذف مشاهدات i امین psu نمونه طبقه h حاصل شده است. آن گاه

$$\hat{\zeta}_{(hi)}^A = f(\hat{\mathbf{A}}_{(hi)}, \hat{\mathbf{B}}, \hat{\beta}_{(hi)}^A)$$

که $(\hat{\beta}_{(hi)}^A)$ نیز برآوردگر β با استفاده از S_B و $S_{(hi)}^A$ با فرض $(\tilde{n}_h^A > 2)$ می‌باشد. به طور مشابه

$$\hat{\zeta}_{(lj)}^B = f(\hat{\mathbf{A}}, \hat{\mathbf{B}}_{(lj)}, \hat{\beta}_{(lj)}^B)$$

به شرط اینکه $\tilde{n}_l^B > 2$ باشد. آن گاه برآوردگر جک نایف واریانس $\hat{\zeta}$ عبارتست از:

$$v_J(\hat{\zeta}) = \sum_{h=1}^H \frac{\tilde{n}_h^A - 1}{\tilde{n}_h^A} \sum_{i=1}^{\tilde{n}_h^A} (\hat{\zeta}_{(hi)} - \hat{\zeta})^2 + \sum_{l=1}^L \frac{\tilde{n}_l^B - 1}{\tilde{n}_l^B} \sum_{j=1}^{\tilde{n}_l^B} (\hat{\zeta}_{(lj)} - \hat{\zeta})^2 \quad (16)$$

برآوردگر جک نایف اصلاح شده واریانس، $\hat{V}_{MJ}(\hat{\zeta})$ نیز به همان فرم $\hat{V}_J(\hat{\zeta})$ بوده با این تفاوت که به جای $\hat{\beta}_{(hi)}^A$ یا $\hat{\beta}_{(lj)}^B$ از $\hat{\beta}$ استفاده کرده و قابل کاربرد با $\tilde{n}_h^A = 2$ یا $\tilde{n}_l^B = 2$ می‌باشد. برآوردگر خطی واریانس $\hat{\zeta}$ عبارتست از:

$$v_L(\hat{\zeta}) = \hat{g}_A^T(\hat{A}, \hat{B})S_A\hat{g}_A(\hat{A}, \hat{B}) + \hat{g}_B^T(\hat{A}, \hat{B})S_B\hat{g}_B(\hat{A}, \hat{B}) \quad (17)$$

که \hat{g}_A و \hat{g}_B از $\hat{\beta}$ بجای β استفاده می‌کنند. هر سه برآوردگر واریانس $v_J(\hat{\zeta})$ ، $v_{MJ}(\hat{\zeta})$ و $v_L(\hat{\zeta})$ بطور مجانبی در حالتی که اختلاف بین هر جفت از آنها از مرتبه $o_p(\tilde{n}^{-1})$ باشد هم‌ارزند. [رائو و لهر (۱۹۹۷)]

برآوردهای جک نایف واریانس مذکور برای حالت نمونه‌گیری طبقه‌بندی بدون جایگذاری از واحدهای نمونه‌گیری اولیه با فرض ناچیز بودن کسر نمونه‌گیری مطرح شده و قابل کاربرد با طرحهای نمونه‌گیری طبقه‌بندی چندمرحله‌ای بدون جایگذاری نیز می‌باشند. در حالتی که نمونه‌گیری با جایگذاری از واحدهای نمونه‌گیری اولیه به عمل آید نیز همین فرمولهای واریانس صادق هستند. یک حالت خاص ممکن است زمانی رخ دهد که در یکی از چارچوبها مثلاً چارچوب فهرستی B ، نمونه‌گیری تصادفی ساده طبقه‌بندی انجام شده و کسرهای نمونه‌گیری ناچیز باشد. در این حالت b_{lj} برداری از مقادیر مرتبط با z امین واحد در طبقه l ام چارچوب B و $\tilde{n}_l^B = n_l^B$ عبارت از تعداد واحدهای نمونه‌گیری شده از N_l^B واحد در طبقه l ام چارچوب B می‌باشد. در صورتی که کسر نمونه‌گیری n_l^B/N_l^B ناچیز نباشد، بطور مشابه با روش جک نایف در حالت چارچوب منفرد (ولتر (۱۹۸۵))، $(n_l^B - 1)/n_l^B$ با $(n_l^B - 1)/N_l^B$ جایگزین می‌شود.

۲.۳ محاسبه برآورد جک نایف واریانس برآوردهای دوچارچوبی

تمام برآوردهایی دو چارچوبی می‌توانند به فرم زیر بیان شوند:

$$\hat{Y} = \sum_{t \in s_A} \tilde{w}_t^A y_t + \sum_{t \in s_B} \tilde{w}_t^B y_t \quad (18)$$

که از وزنه‌های اصلاح شده \tilde{w}_t^A و \tilde{w}_t^B استفاده می‌کند. برای مثال $\hat{Y}_{PML}(\theta)$ از وزنه‌های زیر استفاده می‌کند:

$$\tilde{w}_t^A = \begin{cases} w_t^A [N_A - \hat{N}_{ab}^{PML}(\theta)] / \hat{N}_a^A & t \in a \\ w_t^A \theta \hat{N}_{ab}^{PML}(\theta) / \hat{N}_{ab}(\theta) & t \in ab \end{cases}$$

$$\tilde{w}_t^B = \begin{cases} w_t^B [N_B - \hat{N}_{ab}^{PML}(\theta)] / \hat{N}_b^B & t \in b \\ w_t^B (1 - \theta) \hat{N}_{ab}^{PML}(\theta) / \hat{N}_{ab}(\theta) & t \in ab \end{cases}$$

اگر $\hat{\theta}_p$ را به عنوان برآوردگر پارامتر θ در نظر بگیریم، کفایت در وزنه‌های فوق بجای θ برآورد آن را جایگزین کنیم.

برای محاسبه‌ی برآورد جک نایف $\hat{Y}_{(hi)}^A$ بسادگی وزنه‌های w_t^A را با $w_{t(hi)}^A$ جایگزین می‌کنیم. اگر واحد t در خوشه‌ی k ام طبقه‌ی g در چارچوب A باشد آن گاه

$$w_{t(hi)}^A = \begin{cases} \circ & hi = gk \\ \frac{\tilde{n}_h^A}{(\tilde{n}_h^A - 1)} w_t^A & h = g \quad i \neq k \\ w_t^A & h \neq g \end{cases}$$

بطور مشابه $\hat{Y}_{(lj)}^B$ را از وزنه‌های $w_{t(lj)}^B$ متناظراً محاسبه کرده و سپس برآورد جک نایف واریانس \hat{Y} از رابطه (۱۴) بدست می‌آید.

اگر وزنه‌های اصلاح شده \tilde{w}_t^A و \tilde{w}_t^B به عناصر S_A و S_B بستگی داشته باشند همانطور که برای $\hat{Y}_{PML}(\hat{\theta}_p)$ اتفاق افتاد برای محاسبه برآورد جک نایف کامل $\hat{Y}_{(hi)}^A$ احتیاج به محاسبه $S_{(hi)}^A$ داریم. ماتریس $S_{(hi)}^A$ به ازای $\tilde{n}_h^A \geq 3$ می‌تواند با استفاده از یک جک نایف جداگانه درون هر تکرار جک نایف، این بار با مجموعه داده‌ای که مشاهدات i امین طبقه‌ی h ام چارچوب A حذف شده است. اگر $\tilde{n}^A - H$ و $\tilde{n}^B - L$ هر دو بزرگ باشند، آن‌گاه یک فاصله اطمینان تقریبی ۹۵ درصد برای τ می‌تواند به صورت زیر بنا شود:

$$\hat{\zeta} \pm 1/96 \sqrt{\hat{V}(\hat{\zeta})} \quad (19)$$

در غیر اینصورت با استفاده از صدک توزیع t با درجه آزادی برابر $\min(\tilde{n}^A - H, \tilde{n}^B - L)$ بجای $1/96$ به یک فاصله اطمینان محافظه کاری دست یابیم.

۴ نتایج شبیه سازی

در این بخش با استفاده از شبیه سازی به مطالعه خواص تجربی برآوردگرهای واریانس می‌پردازیم. فرض می‌شود که جامعه نامتناهی بوده و طرح نمونه‌گیری برای هر چارچوب یک طبقه دارد. یک نمونه خوشه‌ای دو مرحله‌ای با \tilde{n}_A خوشه و m عنصر از هر خوشه به عنوان نمونه حاصل از چارچوب A ($n_A = \tilde{n}_A \cdot m$) و یک نمونه تصادفی ساده با n_B مشاهده به عنوان نمونه حاصل از چارچوب B تولید می‌شود. با فرض نامتناهی بودن جامعه، N_a/N و N_b/N را با γ_a و γ_b جایگزین می‌کنیم. نمونه حاصل از چارچوب A شامل مقادیر $\{(y_{ij}, m_{ai}), i = 1, 2, \dots, \tilde{n}_A \quad j = 1, \dots, m\}$ می‌باشد. عبارت از تعداد عناصر نمونه‌گیری شده از i امین خوشه نمونه‌ای، متعلق به حوزه a و y_{ij} نیز مقدار مرتبط با

زامین عنصر نمونه‌ای در i امین خوشه نمونه می‌باشد. نمونه حاصل از چارچوب B شامل مقادیر نمونه‌ای $\{(y_j, n_b), j = 1, 2, \dots, n_B\}$ است که در آن n_b تعداد عناصر نمونه‌ای متعلق به حوزه b و y_j مقدار مرتبط با زامین عنصر نمونه‌ای است.

برای شبیه‌سازی، نمونه‌گیری را $R=10000$ بار تکرار کرده و از هر نمونه، برآوردهای $\hat{Y}_H(\hat{\theta}_H)$ ، $\hat{Y}_{SFreg}(\hat{\beta}_{SF})$ و $\hat{Y}_{PML}(\hat{\theta}_P)$ و برآوردهای تک چارچوبی \hat{Y}_{SF} و \hat{Y}_{SFrake} و $\hat{Y}_{SFreg}(\hat{\beta}_{SF})$ را با استفاده از مقادیر بهینه $\hat{\theta}_H$ و $\hat{\beta}_{SF}$ و $\hat{\theta}_P$ و نیز سه برآوردهای واریانس خطی (L)، جک نایف کامل (J) و جک نایف اصلاح شده (MJ) را به ازای مقادیر مختلف پارامترهای طرح محاسبه می‌کنیم. در محاسبه جک نایف کامل به دلیل اینکه از مقادیر برآورد شده پارامترها استفاده می‌شود و این مقادیر، خود به ماتریس‌های S_B و S_A وابسته‌اند، $S_{(ij)}^B$ و $S_{(hi)}^A$ را با استفاده از یک جک نایف جداگانه درون هر تکرار جک نایف محاسبه می‌کنیم. در جدول (۱) نتایج میانگین توان دوم خطای تجربی EMSE (متوسط مربع انحراف برآوردگر از مقدار واقعی) برآوردهای دوچارچوبی و همچنین اریبی نسبی 23 (RB) و خطای معیار نسبی 24 (RSD) برآوردهای خطی، جک نایف و جک نایف اصلاح شده ارائه شده که بصورت زیر محاسبه می‌شوند:

$$EMSE = \frac{1}{R} \sum_{r=1}^R (\hat{Y}_r - Y)^2$$

$$RB = 100(EV - EMSE)/EMSE$$

$$RSD = (\text{انحراف معیار } 10000 \text{ برآورد واریانس}) / \sqrt{EMSE}$$

\hat{Y}_r عبارت از مقدار \hat{Y} برای r امین تکرار شبیه سازی و R تعداد تکرارهای شبیه سازی و EV عبارت از میانگین 10000 برآورد واریانس در هر روش است. به دلیل مشابه بودن نتایج اریبی نسبی خطی و جک نایف اصلاح شده فقط نتایج حاصل از روش خطی ارائه شده است. همچنین با استفاده از فاصله اطمینان ۹۵ درصد نرمال

$$(\text{برآورد}) \pm 1/96 SE$$

و نیز صدک توزیع t با $(\tilde{n}_A - 1)$ درجه آزادی، احتمال پوشایی تجربی برای روشهای خطی و جک نایف کامل داده شده است. در این شبیه‌سازی، تمام محاسبات با استفاده از نرم افزار S-PLUS 2000 صورت گرفته است.

نتایج شبیه‌سازی در جدول (۱) نشان می‌دهد که اریبی نسبی هر سه برآوردهای واریانس با افزایش حجم نمونه میل به کاهش دارند. همچنین روش جک نایف از پایداری کمتری نسبت به برآوردهای خطی واریانس برخوردار است. بعنوان مثال وقتی $\tilde{n}_A = 10$ خطای معیار نسبی جک نایف تقریباً دو برابر خطای معیار نسبی خطی است با این حال با افزایش حجم نمونه پایداری برآوردهای جک نایف بهبود می‌یابد.

23) Relative Bias 24) Relative standard deviation

نتایج حاصل از فاصله اطمینان نشان می‌دهد که روش جک نایف کامل به وضوح احتمال پوشایی بالاتری از دو روش دیگر دارد. با استفاده از صدک توزیع t با $n_A - 1$ درجه آزادی بجای ۱/۹۶، جک نایف کامل احتمال پوشایی نزدیک به ۰/۹۵ دارد در حالیکه دو روش دیگر به وضوح احتمال پوشایی کمتری دارند. با افزایش حجم نمونه این احتمال پوشایی افزایش می‌یابد.

۵ نتیجه‌گیری

برآوردگر جک نایف واریانس را بطور نظری توجیه کرده با استفاده از شبیه‌سازی نشان دادیم که دارای اریبی کوچکتری از برآوردگر خطی واریانس است. همچنین برآوردگر جک نایف به وضوح دارای احتمال پوشایی بالاتری از برآوردگر خطی است خصوصاً وقتی از صدک t استفاده می‌شود احتمال پوشایی جک نایف خیلی نزدیک ۰/۹۵ است.

جک نایف معمولاً با توابع غیر خطی نظیر نسبت مجموع دو جامعه بکار می‌رود. در این حالت مشتقات جزئی که در محاسبه‌ی برآورد خطی چنین مقادیر غیر خطی استفاده می‌شود، در بررسی‌های دوچارچوبی پیچیده‌تر از تک چارچوبی است. در صورتی که با استفاده از جک نایف می‌توان از این محاسبات پیچیده اجتناب کرد. برآوردگرهای جک نایف اخیراً در اداره آمار کانادا برای محاسبه برآورد واریانس در آمارگیری وسیع ملی کودکان و جوانان^{۲۵} استفاده می‌شود. روشهای باز نمونه‌گیری دیگر برآورد واریانس نظیر تکرار متعادل پاسخ^{۲۶} (BRR) و بوت استرپ^{۲۷} نیز می‌توانند به موازات جک نایف در این گونه طرحها مورد استفاده قرار گیرند. حسن این روشها این است که برخلاف جک نایف با توابع غیر هموار نظیر میانه نیز قابل کاربرد هستند.

25) National Longitudinal Survey of Children and Youth 26) Balanced Repeated Replication 27) Bootstrap

| ت | | احتمال پوشایی | | اریبی نسبی | | EMSE | برآوردگر |
|-------|-------|---------------|-------|-----------------|------------------|-------|----------|
| J | L | J | L | J | L | | |
| ۰,۹۴۳ | ۰,۹۲۳ | ۰,۹۱۹ | ۰,۸۸۷ | ۱۵,۴۵ (۱,۹۰) | -۱۹,۹۴ (۰,۷۸) | ۷,۱۶ | H |
| ۰,۹۴۶ | ۰,۹۳۱ | ۰,۹۲۱ | ۰,۹۰۱ | ۱۰,۴۶ (۱,۶۴) | -۱۲,۷۵ (۰,۷۷) | ۶,۸۳ | FB |
| ۰,۹۵۰ | ۰,۹۳۹ | ۰,۹۲۴ | ۰,۹۱۲ | ۱۱,۴۶ (۱,۸۷) | -۹,۹۱ (۰,۸۲) | ۶,۵۸ | PML |
| ۰,۹۷۲ | ۰,۹۷۲ | ۰,۹۴۸ | ۰,۹۴۸ | ۴,۱۷ (۰,۴۷) | ۴,۱۷ (۰,۴۷) | ۱۱,۲۶ | SF |
| ۰,۹۴۳ | ۰,۹۲۸ | ۰,۹۱۸ | ۰,۸۹۷ | ۱۱,۴۲ (۱,۷۸) | -۱۳,۷۳ (۰,۸۲) | ۶,۹۷ | SFreg |
| ۰,۹۵۰ | ۰,۹۴۴ | ۰,۹۲۳ | ۰,۹۱۸ | ۴,۸۵ (۱,۱۹) | -۰,۱۱ (۱,۰۸) | ۶,۵۵ | SFrake |
| ۰,۹۴۴ | ۰,۹۲۳ | ۰,۹۲۸ | ۰,۹۰۶ | ۱,۸۰ (۰,۷۴) | -۱۵,۹۷ (۰,۴۱) | ۳,۵۸ | H |
| ۰,۹۴۳ | ۰,۹۳۰ | ۰,۹۲۸ | ۰,۹۱۳ | -۰,۶۳ (۰,۶۷) | -۱۲,۶۹ (۰,۴۰) | ۳,۵۲ | FB |
| ۰,۹۴۵ | ۰,۹۳۵ | ۰,۹۳۰ | ۰,۹۱۹ | -۱,۱۰ (۰,۷۴) | -۱۱,۲۶ (۰,۴۲) | ۳,۴۴ | PML |
| ۰,۹۵۸ | ۰,۹۵۸ | ۰,۹۴۴ | ۰,۹۴۴ | -۰,۶۹ (۰,۲۳) | -۰,۶۹ (۰,۲۳) | ۵,۸۳ | SF |
| ۰,۹۴۲ | ۰,۹۲۸ | ۰,۹۲۸ | ۰,۹۱۳ | ۰,۰۸ (۰,۷۲) | -۱۲,۹۲ (۰,۴۲) | ۳,۵۷ | SFreg |
| ۰,۹۴۴ | ۰,۹۴۰ | ۰,۹۲۹ | ۰,۹۲۵ | -۳,۰۷ (۰,۵۵) | -۵,۳۳ (۰,۵۳) | ۳,۴۸ | SFrake |

جدول (۱). لازم به توجه است که در این جدول نتایج تنها به ازای $\gamma_a = ۰/۱$ و $\gamma_b = ۰/۲$ ارائه شده است. در بخش بالایی، $\tilde{n}_A = ۱۰$ و $n_B = ۱۰۰$ و در بخش پائینی $\tilde{n}_A = ۲۰$ و $n_B = ۲۰۰$ می باشد. عبارت از میانگین مربعات خطای مونت کارلو برای ۱۰۰۰۰ بار تکرار شبیه سازی است. اریبی نسبی روش خطی و جک نایف و همچنین خطای معیار نسبی آنها در پرانتز زیر مقادیر اریبی نسبی داده شده است.

مراجع

- [1] Hartley, H. O.(1962). Multiple Frame Surveys, in Proceedings of the Social Statistics Section,American Statistical Association, 203-206.
- [2] Lund, R. E. (1968) . Estimators in Multiple Frame Surveys, in Proceedings of the Social Statistics Section, American Statistical Association, 282-288.
- [3] Fuller, W. A. and Burmeister, L. F. (1972). Estimators for Samples Selected From two Overlapping Frames,in Proceedings of the Social Statistics Section, American Statistical Association, 245-249.
- [4] Hartley, H. O.(1974). Multiple Frame Methodology and Selected Applications, Sankhya, Ser. C, 36, 203-206.
- [5] Vogel, F. A. (1975). Survey With Overlapping Frame-Problems in Application, in Proceedings of the Social Statistics Section, American Statistical Association, 694-699.
- [6] Ford, B. L. and Bosecker, R. R. (1979) Multiple Frame Estimation With Stratified Overlap Domain ,in Proceeding of the Social Statistics Section, American Statistical Association, 219-224.
- [7] Bankier, M. D. (1986). Estimators Based On Several Stratified Samples With Applications to Multiple Frame Surveys , Journal of the American Statistical Association, 81, 1074-1079.
- [8] Skinner, C. J. (1991). On the Efficiency of Raking Ratio Estimation for Multiple Frame Surveys , Journal of the American Statistical Association, 86, 779-784.
- [9] Skinner, C. J. and Rao, J. N. K. (1996). Estimation in Dual Frame Survey With Complex Designs, Journal of the American Statistical Association , 91, 349-356.
- [10] Kalton, G. and Anderson, D. W. (1986). Sampling Rare Populations, Journal of the Royal Statistical Society ,Ser. A, 149, 65-82.
- [11] Rao, J. N. K. and Lohr, S. L. (1997), Jackknife Variance Estimation in Dual Frame Surveys, Technical Report, Department of Mathematics and Statistics, Carleton University.
- [12] Rao, J. N. K. and Wu, C. F. J. (1985). Inference from Stratified Samples: Second-Order Analysis of Three Methods for Nonlinear Statistics, Journal of the American Statistical Association , 80, 620-630.
- [13] Wolter, K. M. (1985). Introduction to Variance Estimation, New York: Springer-Verlag.

- [14] Rao, J. N. K. and Lohr, S. L. (2000). Inference From Dual Frame Survey Journal of The American Statistical Association, 95, 271-280.

آزمون برازش توزیع لجستیک در تعیین توانایی و رتبه‌بندی در امتحان‌ها

سید مقتدی هاشمی پرست^۱، سید حسن مرتضوی نصیری^۲، کمال عقیق^۳

^۱ مرکز مطالعات، تحقیقات و ارزشیابی آموزشی

^۲ مرکز تحقیقات سیاست علمی کشور

^۳ دانشگاه صنعتی خواجه نصیرالدین طوسی

چکیده: در این مقاله ابتدا با استفاده از توزیع لجستیک سه پارامتری در پاسخگویی به سوالات چهارگزینه‌ای در امتحانات ورودی، میزان توانایی هر یک از داوطلبان برای هر سوال در نتیجه هر امتحان با روش نظریه سوال-پاسخ (IRT) مورد اندازه‌گیری و بر اساس توانایی نمره او تعیین شده و مبتنی بر این اطلاعات رتبه‌بندی داوطلبان انجام پذیرفته است. با استفاده از روش متداول سوالات چهارگزینه‌ای، توزیع‌های تجربی و نظری نیز محاسبه و رتبه‌بندی انجام پذیرفته است. دو نتیجه حاصل مورد قیاس قرار گرفته و با استفاده از آزمون نکویی برازش تناظر بین رتبه‌بندی‌ها مورد مطالعه و بحث قرار گرفته‌اند. در پایان نتایج حاصل در تعیین رتبه‌بندی دانشجویان در امتحان تولیمو ۱۲ به کار گرفته شده و با نتایج حاصل از رتبه‌بندی متداول مورد مقایسه قرار گرفته است.[†]

واژه‌های کلیدی: نظریه سوال-پاسخ، IRT، رتبه‌بندی، توزیع لجستیک، توانایی

۱ مقدمه

در این مقاله ابتدا به تحلیل و بررسی آزمون‌هایی که به صورت تست‌های چهار جوابی است با استفاده از مدل نظریه سوال-پاسخ^۱ (IRT) مبتنی بر مدل لجستیک سه پارامتری می‌پردازیم. در این تحلیل

(۱) برآورد پارامترهای هر یک از سوالات (دشواری b، تشخیص a و حدس c)،
(۲) برآورد احتمال آرایه پاسخ درست در هر یک از سطوح توانایی به آنها توسط نرم‌افزار ترسیم منحنی مشخصه‌های سوالات و نمودار تابع اطلاع آنها،

(۳) طبقه‌بندی سوالات،

(۴) برآورد توانایی هر یک از امتحان دهندگان،

(۵) نمره‌گذاری نتایج بر اساس برآورد توانایی هر امتحان دهنده (نمره‌گذاری مبتنی بر توانایی)،

(۶) مقایسه نتایج بر اساس نمرات خام و نمرات مبتنی بر توانایی با استفاده از آزمون کای اسکور.

†) با تشکر از سازمان سنجش و آموزش کشور که با پشتیبانی و حمایت آن سازمان این تحقیق انجام گرفته است.

1) Theory Item-Response

با وجود اینکه سالیان درازی است که آزمون‌گری کلاسیک از رواج و شهرت فراوانی برخوردار است اما همچنان دارای کاستی‌هایی است که سودمندی آن را به عنوان مبنا و اساسی برای آزمون‌گری مدرن محدود می‌سازد. شکوفایی و گسترش فن‌آوری کامپیوتر و توسعه روزافزون کاربرد آن در روانسنجی سبب شده است تا برخی از محدودیتهای نظریه کلاسیک آشکارتر و برجسته‌تر شود. کاربرد نظریه کلاسیک در موقعیتهای سنتی آزمون‌گری، اعم از گروهی یا فردی بسیار مناسب است. در این موقعیت‌ها که عملاً نفس رتبه‌بندی اهمیت پیدا می‌کند و از تعیین توانایی افراد غفلت می‌شود، آزمون‌هایی یکسان یا دسته‌هایی از سوالات موازی به تمام اعضای جمعیت هدف (مثلاً اشخاصی که درصدد ورود به دانشگاه هستند) داده می‌شود. این مجموعه سوالات را می‌توان به صورت چاپ شده یا از طریق کامپیوتر به امتحان دهندگان ارائه داد. وقتی از توانایی در زمینه‌های علمی صحبت می‌کنیم، می‌توانیم از اصطلاحات توصیف کننده‌ای همچون توانایی خواندن یا توانایی ریاضی استفاده کنیم. هر یک از این اصطلاحات درست به همان چیزی اشاره دارد که متخصصان روان‌سنجی از آن با عنوان صفت نامشهود^۲ یا صفت مکنون^۳ یاد می‌کنند. با وجود اینکه متخصصان روان‌سنجی چنین اصطلاحی را به راحتی توصیف می‌کنند و افراد مطلع می‌توانند خواص آن را برشمارند، اما اندازه‌گیری این مقولات مفهومی، به آسانی اندازه‌گیری خصوصیات فیزیکی و جسمانی نیست. یکی از هدف‌های اولیه اندازه‌گیری‌های روانی و تربیتی این است که تعیین کند افراد از چه مقداری از این صفت‌ها یا توانایی‌ها برخوردارند.

از اوایل دهه ۱۹۵۰ مدل‌هایی از اندازه‌گیری‌های روانی عرضه شدند که بر این تنگناهای نظریه کلاسیک انگشت گذاردند. امروزه رایج‌ترین و پخته‌ترین این مدل‌ها، مجموعه‌ای را تشکیل می‌دهند که به شناخت و تعیین ویژگیها و مشخصه‌های ریاضی پاسخهای امتحان دهندگان به آزمون کمک می‌کند. این مدل‌ها به عنوان مدل‌های نظریه سوال-پاسخ یا مدل‌های IRT شناخته شده‌اند. نظریه سوال-پاسخ (IRT) رابطه مشخصه‌ها یا پارامترهای هر سوال و ویژگیها یا توانایی افراد (صفت مکنون) را با احتمال ارائه پاسخ صحیح بررسی می‌کند. فرض بر آن است که مقدار این صفت (توانایی) همیشه در امتداد پیوستاری تک بعدی^۴ تغییر می‌پذیرد که معمولاً آن را با حرف θ نشان می‌دهند. طبق نظریه سوال-پاسخ، هم سوالات آزمون و هم افراد پاسخ دهنده به آنها، روی مقیاس θ از کمترین تا بیشترین مقدار چیده می‌شوند. جای شخص i ام که روی محور θ با θ_i نشان داده می‌شود معمولاً به توانایی یا شایستگی وی تعبیر می‌شود. جای سوال j روی محور θ که معمولاً با b_j نشان داده می‌شود، به دشواری سوال تعبیر می‌شود. به طور شمی یا شهودی می‌توان انتظار داشت وقتی که $b_j - \theta_i$ افزایش پیدا می‌کند، احتمال وقوع پاسخ صحیح به سوال j ام نیز به طور یکنواخت افزایش یابد. تمام مدل‌های IRT احتمال پاسخ صحیح به سوالی از آزمون‌های چند گزینه‌ای با یک پاسخ درست را به عنوان تابعی از θ مشروط به یک یا چند پارامتر سوال نشان می‌دهند. برای هر سوال می‌توان احتمال پاسخ صحیح دادن یا موافقت با طبقه خاصی از پاسخ را روی نمودار نشان داد. این توابع معرف رگرسیون غیرخطی احتمال وقوع پاسخ صحیح هستند. در این توابع احتمال وقوع پاسخ صحیح برحسب وجود مقدار معینی

2) unobservable trait 3) latent trait 4) unidimensional

از یک توانایی یا یک صفت مکنون نظیر هشیاری یا توانایی کلامی در نزد پاسخ دهنده تعریف می شود (هالین، دراسگو و پارسونز^۵).

احتمال وقوع هر پاسخ به عنوان تابعی از پارامترها و به طور جداگانه برای هر سوال و مشخصه های هر شخص، مدل سازی می شود. پارامترهای سوال معرف خواصی از سوال است که آن را به عنوان میزان دشواری و توانایی تشخیص سوال می شناسیم و مشخصه های شخص معرف میزان توانایی آزمون شونده است. اگر احتمال وقوع پاسخ به عنوان تابعی از مشخصه های (توانایی) فرد معرفی شود، تابع مورد نظر در IRT بعنوان تابع پاسخ شناخته می شود. در این مقاله گزارشی این از تحلیل و بررسی در آزمون تخصصی زبان انگلیسی سازمان سنجش آموزش کشور که با استفاده از مدل سه پارامتری انجام پذیرفته است عرضه می شود.

۲ تابع لجستیک به عنوان مبنای نظریه سوال-پاسخ

اگر بتوان با متغیر تصادفی X مشخصه های رفتار یک سیستم را در مقابل یک عمل یا آزمایش بیان کرد، قانون توزیع X ، توزیعی است که به نام توزیع لجستیک خوانده می شود. مثلاً میزان پاسخ های گروهی از امتحان دهندگان در مقابل یک پرسش، یا اندازه و میزان نحوه رفتار یک جامعه در برخورد با یک پدیده اجتماعی، دارای توزیع لجستیک است. توابع لجستیک اغلب برای توضیح گونه های خاصی از توزیع رشد بکار می رود. این مدلها، مدل های خوبی برای توضیح دادن توزیع و گسترش اطلاعات در جوامع و جمعیتها می باشند. همچنین مدل های خوبی برای تشریح و تبیین توزیع رشد بیولوژیک جمعیت در گونه هایی است که میزان رشد جمعیت شان آن قدر زیاد بوده که اکو سیستم آنها را به نقطه اشباع نزدیک کرده است. بنابراین می توان انتظار داشت که احتمال ارایه پاسخ صحیح به سوالاتی که یک توانایی معین را اندازه می گیرند به موازات افزایش توانایی در آزمودنی، افزایش یابد. منحنی تابع توزیع این توابع همانند تابع توزیع توابع نمایی، در ابتدا افزایش یا رشد سریع دارند اما به دلیل وجود بعضی محدودیت ها که حد و حدودی را بر جمعیت اعمال می کند، این افزایش به آرامی کاهش پیدا کرده و سپس قطع می شود. بنابراین مدل استاندارد ریاضی برای نظریه سوال-پاسخ و برای منحنی مشخصه های سوال، مدل تابع لجستیک^۶ می باشد. نظریه سوال-پاسخ با فرض وجود رابطه ای ریاضی بین توانایی ها (یا دیگر صفت های مفروض) و احتمال جوابگویی به سوالها، به مطالعه و بررسی نمرات سوال و امتحان می پردازد.

5) Hulin, Drasgow & Parsons 6) logistic function

۳ مدل‌های یک، دو و سه پارامتری (بیکر [1])

وقتی در تحلیل هر سوال برآورد دو مشخصه آن یعنی میزان دشواری و تشخیص آن مورد نظر باشد مدل دو پارامتری با تابع توزیع احتمال زیر را انتخاب می‌کنیم.

$$p(\theta) = \frac{1}{1 + e^{-a(\theta-b)}} \quad (1)$$

در رابطه فوق θ نمایشگر متغیر توانایی و a نمایشگر پارامتر تشخیص و b نمایشگر پارامتر دشواری سوال است. واضح است که پارامترهای a و b برای هر سوال ثابت‌اند و دامنه تغییرات آنها از $-\infty$ تا $+\infty$ است ولی در عمل $2/8^\circ < a < 2/8^\circ - 2/8^\circ$ و $3 < b < 3 - 3$ در نظر گرفته می‌شود. b نقطه ای روی محور توانایی است که احتمال پاسخی گویی متناظر با آن 0.5 باشد. باید توجه داشت که بر اساس مفروضات نظریه سوال-پاسخ، متغیر توانایی از پارامترهای تشخیص و دشواری سوال مستقل است.

اگر در رابطه فوق مقدار a را برابر ۱ قرار دهیم، مدل فوق به مدل لجستیک یک پارامتری تبدیل خواهد شد.

برین باوم^۷ [3] با اندکی دستکاری در تابع لجستیک، مدل سه پارامتری را ارائه داد. اما با این وجود مدل وی همواره به عنوان تابع لجستیک سه پارامتری خوانده شده است. وقتی که عامل شانس در پاسخگویی دخیل باشد، مدلی سه پارامتری خواهیم داشت. برای مدل سه پارامتری می‌توان از فرمول زیر استفاده کرد:

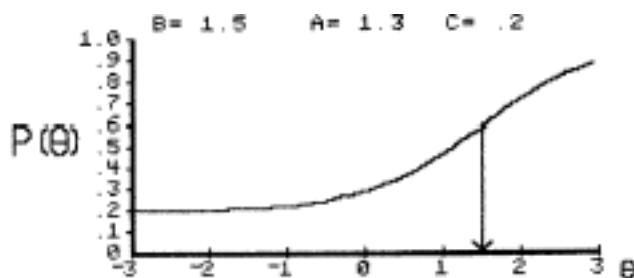
$$p(\theta) = c + (1 - c) \frac{1}{1 + e^{-a(\theta-b)}}$$

یعنی در این مورد تنها ضریب حدس ($0 < c < 1$) وارد می‌شود که معمولاً حدس بالای 0.35 قابل قبول نیست. با توجه به آنچه گفته شد در واقع می‌توان برای توانایی هر امتحان دهنده یک مقدار عددی یا یک نمره منظور کرد که جای او را روی محور توانایی مشخص می‌کند. این مقدار توانایی را با حرف θ و احتمال اینکه امتحان دهنده با این سطح توانایی به این سوال پاسخ درست بدهد را با $p(\theta)$ نشان می‌دهیم. در شکل ۱ محور x نشانگر سطح توانایی و محور y نشان دهنده احتمال وقوع پاسخ درست به یک سوال از امتحان می‌باشد. منحنی S گونه نشان دهنده احتمال وقوع پاسخ درست از سوی افرادی است که مقدار توانایی آنها متفاوت است.

۴ برآورد پارامترهای یک سوال: بیکر [2]

فرض کنیم در آزمونی، یک نمونه مرکب از n تا امتحان شونده به یک سوال موجود در امتحان پاسخ داده‌اند. سطح توانایی این امتحان دهندگان روی محور توانایی توزیع شده است. این امتحان

7) A. Birnbaum



شکل ۱: منحنی مشخصه‌های سوال برای مدل سه پارامتری

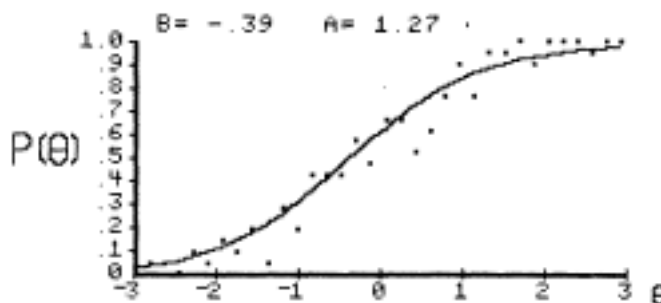
دهندگان را به t گروه توانایی در طول محور توانایی تقسیم می‌کنیم. فرض کنیم n_i امتحان دهنده در گروه $i, i = 1, 2, 3, \dots, t$ ، با سطح توانایی یکسان θ_i قرار دارند و فرض کنیم r_i امتحان دهنده از این گروه به سوال، پاسخ صحیح داده‌اند. بنابراین برآورد احتمال پاسخ‌های صحیح در سطح توانایی θ_i ، با محاسبه نسبت زیر (یعنی نسبت پاسخ‌های صحیح مشاهده شده به کل پاسخها در همین زیرگروه) امکان پذیر است:

$$p(\theta_i) = \frac{r_i}{n_i}$$

کار اساسی ما این است که منحنی مشخصه‌های سوال را که بهترین برازش^۸ برای نسبت پاسخ‌های مشاهده شده است، به دست آوریم. فرآیند برازش منحنی، با استفاده از روش برآورد حداکثر درست‌نمایی^۹ انجام می‌گیرد.

در این روش، برآوردهای اولیه برای مقادیر a و b بدین قرار فرض می‌شود $a = 1$ و $b = 0$. آنگاه با تکیه بر این برآوردها، مقدار $P(\theta_j)$ با استفاده از معادله تابع لجستیک مدل دو پارامتری در سطح توانایی θ_j محاسبه می‌شود. سپس این مقدار، با مقدار مشاهده شده $p(\theta_i)$ در هر زیرگروه توانایی مقایسه می‌شود. بعد برای پارمترهای برآورد شده سوال سرشکنی تعیین می‌شود. نتیجه این کار آن است که توافق بهتری بین منحنی مشخصه‌های سوال (که با استفاده از مقادیر برآورد شده برای پارمترها رسم می‌شود) و نسبت مشاهده شده پاسخ‌های درست، حاصل می‌شود. فرآیند سرشکنی برآوردها تا جایی ادامه می‌یابد که مقدار آن به اندازه‌ای کوچک شود که دیگر امکان بهبود این توافق بسیار ناچیز باشد. در این نقطه، تکرار عملیات مراحل برآورد متوقف می‌شود و آخرین مقادیر بدست آمده a و b برآورد پارمترهای سوال می‌باشند. با قرار دادن این برآوردها در معادله منحنی مشخصه‌های سوال می‌توان احتمال پاسخ درست $P(\theta_i)$ را برای هر سطح توانایی محاسبه کرد. منحنی حاصل، منحنی مشخصه‌های سوال است که بهترین برازش برای

8) fitting 9) maximum likelihood estimation



شکل ۲: منحنی مشخصه‌های سوال برآزش شده با نسبت مشاهده شده پاسخ‌های درست

پاسخ‌های داده شده به آن سوال است. شکل ۲ یک منحنی مشخصه سوال را نشان می‌دهد که با نسبت پاسخ‌های صحیح مشاهده شده برآزش شده است. نتیجه حاصل آن است که در مدل IRT پارامترهای a و b محاسبه شده در هر زیرگروه از توانایی همیشه یکسان خواهد بود. این در حالی است که اندیس دشواری سوال در مدل کلاسیک از یک زیرگروه به زیرگروه دیگر متغیر می‌باشد. به همین دلیل، تعبیر و تفسیر دشواری سوال آن طور که در نظریه IRT تعریف شده، ساده‌تر می‌باشد. در اینجا برآورد پارامترهای سوال را برای مدل لجستیک سه پارامتری توضیح می‌دهیم. فرض کنیم x_j ، تابعی دو مقداری باشد که مقدار ۱ در آن برای پاسخ درست و صفر برای پاسخ نادرست به سوال j ام در نظر گرفته شده است. فرض کنیم X^m ماتریس $k \times 1$ تایی است که مؤلفه‌های آن صفر یا یک می‌باشند یعنی پاسخ‌های دو حالتی شخص m ام به مجموعه‌ای از k سوال باشد. پس فرض کنیم که ماتریس X ماتریس پاسخ‌گویی n تا امتحان دهنده به k تا سوال باشد. احتمال مشاهده X پیش از مشاهده واقعی پاسخ‌های نمونه به قرار زیر است:

$$p(X|\theta, \beta) = \prod_i p(x_i|\theta_i, \beta) = \prod_i \prod_j p_j(\theta_i)^{x_i^j} q_j(\theta_i)^{1-x_i^j} \quad (2)$$

که $q_j(\theta) = 1 - p_j(\theta)$ و بردارهای پارامتری ثابت نامعلوم در مدل هستند. برای مدل لجستیک سه پارامتری، مؤلفه‌های عبارتند از پارامترهای دشواری، تشخیص و حدس برای k تا سوال. معادله (۲) بر این نکته متکی است که احتمال پاسخ‌گویی صحیح از سوالاتی به سوالاتی (هر جفت سوال) در سطح توانایی خاص θ ، مستقل از یکدیگر هستند. این قضیه فرض بنیادی تمام مدل‌های IRT است. بعد از اینکه مشاهده عملی نمونه‌ها انجام شد، دیگر معادله (۲)، یک احتمال نخواهد بود و به جای آن به ازای مقدار x_i معین، معادله را می‌توان به عنوان تابع درست‌نمایی x_i از θ و به شرط x_i تعریف کرد. مقدار تابع درست‌نمایی، به دسته‌ای خاص از مقادیر و پارامترهای واقعی درست‌نمایی بستگی دارد که در نقطه x_i مشاهده بشود. در

مدل لجستیک سه پارامتری، رایج ترین روش برآورد پارامترهای سوال، برآورد حداکثر درست‌نمایی کناری^{۱۰} است. در این روش برآورد پارامترها از طریق انتگرال‌گیری از تابع درست‌نمایی روی توزیع توانایی به دست می‌آید. یعنی

$$L(\beta|X) = \prod_i \int p(x_i|\theta, \beta) f(\theta) d\theta \quad (۳)$$

که $f(\theta)$ تابع چگالی توانایی می‌باشد. برآوردهای حداکثر درست‌نمایی کناری، مقادیری از هستند که معادله (۳) را ماکزیمم می‌کند.

اگر پارامترهای سوال در مدل IRT معلوم باشند آنگاه برآورد توانایی برای نمونه‌ای از امتحان دهندگان با استفاده از روش برآورد حداکثر درست‌نمایی^{۱۱} آسانتر و صریحتر خواهد بود. در همین نوشته خواهیم دید که چگونه می‌توان میزان توانایی یا θ را با استفاده از برآورد پارامترهای سوال که در جریان کالیبره کردن سوالات حاصل می‌شود، برآورد کرد. برآورد همزمان توانایی و پارامترهای سوال در مدل‌های مختلف (اعم از یک، دو یا سه پارامتری) کاری دشوار و پر دردسر است. البته تهیه نرم‌افزارهای کامپیوتری یا نرم‌افزارهای موجود این کار را نسبتاً آسان ساخته است.

۵ برآورد توانایی امتحان دهنده

در نظریه سوال-پاسخ، هدف اولیه آرایه یک آزمون به امتحان دهنده این است که جای او روی مقیاس توانایی مشخص شود. برای اندازه‌گیری توانایی مبتنی بر تعدادی سوال (N) خواهد بود که هر کدام جنبه‌ای از آن توانایی را اندازه می‌گیرد. در بحث پارامترهای سوال و برآورد آنها فرض می‌شود که پارامتر توانایی امتحان دهندگان معلوم است. بر عکس برای برآورد (پارامتر) توانایی نامعلوم یک امتحان دهنده، فرض بر این خواهد بود که مقادیر عددی پارامترهای سوالات آزمون معلوم‌اند. وقتی امتحانی اجرا می‌شود، هر امتحان دهنده به هر یک از سوالات آن پاسخی می‌دهد و به پاسخ‌ها به صورت دو وجهی^{۱۲} نمره داده می‌شود. پس برای هر یک از سوالات آزمون، نتیجه نمره‌ای برابر با یک یا صفر خواهد بود. در نظریه سوال-پاسخ برای برآورد توانایی امتحان دهنده، از روش حداکثر درست‌نمایی استفاده می‌شود. این روش، همان طور که برای برآورد پارامترهای سوال نیز در این حالت وجود داشت، فرآیندی است که به دفعات تکرار می‌شود. انجام آن به این شکل است که با ملحوظ داشتن مقادیر معلوم پارامترهای سوالات، ابتدا عددی به عنوان مقدار پیشین^{۱۳} برای توانایی امتحان دهنده تعیین می‌شود. از این مقادیر برای محاسبه میزان احتمال وقوع پاسخ صحیح به سوال از سوی امتحان دهنده استفاده می‌شود. آنگاه سرشکنی^{۱۴} برای

10) Marginal Maximum Likelihood estimation 11) Maximum Likelihood estimation
12) dichotomously 13) priori value 14) adjustment

برآورد توانایی به دست می‌آید که توافق احتمالات محاسبه شده را با بردار پاسخ‌های امتحان دهنده به سوالات اصلاح کرده و بهبود می‌بخشد. این فرآیند تکرار و تکرار می‌شود تا اندازه میزان حاصل به به قدری کوچک شود که تغییرات در توانایی برآورد شده قابل اغماض باشد. نتیجه عبارت است از برآورد پارامتر توانایی امتحان دهنده. بعد این فرآیند برای هر یک از امتحان دهندگانی که در آزمون شرکت داشته‌اند به طور جداگانه تکرار می‌شود. به هر حال این روش مبتنی بر رویکردی است که با هر امتحان دهنده، جداگانه برخورد می‌کند. از این رو موضوع اصلی آن این است که چطور می‌شود توانایی یک فرد واحد را برآورد کرد. معادله برآورد در فرآیند برآورد حداکثر درست‌نمایی به صورت زیر است:

$$\hat{\theta}_{s+1} = \hat{\theta}_s + \frac{\sum_{i=1}^N a_i [u_i - P_i(\hat{\theta}_s)]}{\sum_{i=1}^N a_i^2 - P_i(\hat{\theta}_s)Q_i(\hat{\theta}_s)}$$

که در آن توانایی برآورد شده برای امتحان دهنده با تکرار s است، a_i پارامتر تشخیص سوال i برای $i = 1, 2, 3, \dots, k$ است، $u_i = 1$ برای پاسخ درست به سوال i و $u_i = 0$ برای پاسخ نادرست به سوال i است، $p_i(\hat{\theta}_s)$ احتمال پاسخ صحیح به سوال i در مدل منحنی مشخصه‌های سوال داده شده در سطح توانایی $\hat{\theta}$ با تکرار s است، $Q_i(\hat{\theta}_s) = 1 - p_i(\hat{\theta}_s)$ احتمال پاسخ نادرست به سوال i در مدل منحنی مشخصه‌های سوال داده شده در سطح توانایی $\hat{\theta}$ با تکرار s است. اینک می‌توانیم خطای استاندارد برآورد را طبق فرمول زیر محاسبه کنیم:

$$SE(\hat{\theta}) = \frac{1}{\sqrt{\sum_{i=1}^N a_i^2 P(\hat{\theta})Q(\hat{\theta})}}$$

قابل توجه است که زیر رادیکال این معادله، همان مخرج معادله برآورد توانایی است. مراحل برآورد حداکثر درست‌نمایی در حالتی که امتحان دهنده به همه سوالات پاسخ غلط و یا در حالتی که امتحان دهنده به همه سوالات پاسخ درست داده باشد قابل محاسبه نیست. یعنی برآورد توانایی در حالت اول $-\infty$ و در حالت دوم $+\infty$ می‌شود.

۶ مربع خی

در تحلیل‌های IRT مسئله مهم این است که آیا یک منحنی مشخصه سوال معین با داده‌های مشاهده شده، خوب برازش ۱۵ شده است یا نه؟ روش‌های مختلفی برای آزمون نکویی برازش پیشنهاد شده است. برای بررسی نکویی برازش از آزمون آماری مجذور خی ۱۶ استفاده شد.

$$\chi^2 = \sum_{i=1}^t n_i \frac{[p(\theta_i) - P(\theta_i)]^2}{P(\theta_i)Q(\theta_i)}$$

15) goodness-of-fit 16) chi-square (chi-square goodness-of-fit)

که t عبارت است از تعداد زیر گروه‌های توانایی، θ_i سطح توانایی گروه i ام، n_i تعداد امتحان دهندگانی که سطح توانایی θ_i دارند، $p(\theta_i)$ احتمال مشاهده پاسخ صحیح گروه i ، و $P(\theta_i)$ احتمال پاسخ صحیح برای گروه i است که با استفاده از برآورد پارامترها در مدل منحنی مشخصه سوال محاسبه شده است و $Q(\theta_i) = 1 - P(\theta_i)$.

اگر مقدار بدست آمده از یک مقدار معین بزرگتر باشد، منحنی مشخصه سوال که با برآورد پارامترهای سوال به دست آمده، برانزده داده‌ها نیست. این امر ممکن است به دو دلیل اتفاق بیافتد. اول اینکه مدل نادرستی برای ترسیم و تحلیل منحنی مشخصه‌های سوال انتخاب شده باشد. دوم اینکه مقدار احتمالات مشاهده شده پاسخ‌های صحیح آن قدر متفرق باشند که توان منحنی برازش شده ای را برای مدل مورد انتظار به دست آورد.

تحت مدل مشخص IRT. فراوانی مورد انتظار انتخاب k توسط پاسخ دهندگان با استفاده از فرمول زیر محاسبه می‌شود.

$$E_i(k) = N \int P(v_i = k | \theta = t) f(t) dt$$

$f(t)$ در این فرمول عبارت است از چگالی توانایی که معمولاً استاندارد نرمال فرض می‌شود زیرا تابع سوال-پاسخ بر مبنای این توزیع مقیاس سازی می‌شود. فون در والنبرگ^{۱۷} نشان داده است که آماره مربع خبی برای هر یک از تک سوال‌ها در بسیاری از موارد به فرض تک بعدی بودن آزمون، حساس نیست. برای اجتناب از این مسئله پیشنهاد شده است که مجذور خبی برای سوالات در دسته‌های دو تایی و سه تایی محاسبه شود. در صورتی که این دسته سوالات از نابرازندگی مشابهی^{۱۸} برخوردار باشند، مقادیر بزرگی از مربع خبی را نشان خواهند داد. فراوانی مورد انتظار برای یک جفت سوال برای خانه‌های (k, k') در جدول دو طرفه برای سوال‌های i و i' به صورت زیر محاسبه می‌شود:

$$E_{i,i'}(k, k') = N \int P(v_i = k | \theta = t) P(v_{i'} = k' | \theta = t) f(t) dt$$

اگر عدم برازشی بین مدل و داده‌ها رخ دهد مربع خبی برای دسته سوالات دو تایی و سه تایی در سطح توانایی یکسان افزایش خواهد یافت. دراسگو و همکارانش ([4]) برای آسان شدن مقایسه مربع خبی بر حسب حجم نمونه پیشنهاد کرده‌اند که مقادیر مربع خبی برای نمونه‌ای سه هزار نفره سرشکن شده و به درجه آزادی تقسیم شود. آنها در پی مطالعات مختلفی در یافته‌اند که نسبت‌های سرشکن شده χ^2 بر درجه آزادی که کمتر از ۳ باشد نشان دهنده نکویی برازش مدل و داده‌هاست. به دست آمدن مقادیر بزرگ برای این نسبت‌ها ممکن است غلط بودن فرض استقلال موضعی یا تک بعدی بودن آزمون را نشان دهد.

17) Van der Wollenberg 18) similar misfits

۷ تحلیل نتایج آزمون تولیمو-الف-تحلیل سؤالات آزمون تولیمو

بر این اساس آنچه تا اینجا گفته شد، برای تمام سؤالات آزمون تولیمو ۱۲، پارامترهای سؤالات (یعنی دشواری d، تشخیص a و حدس c) برآورد شد. چون سؤالات از نوع چند گزینه‌ای بودند، امکان پاسخگویی به سؤالات از طریق حدس وجود داشت، بنابراین پارامتر حدس نیز برای تعیین احتمال پاسخگویی صحیح به سؤال وقتی که دانش یا توانایی درکار نیست، برآورد شد. منحنی مشخصه‌های سؤالات نیز بر اساس مقادیر برآورده شده‌ی پارامترها رسم شد و سرانجام سؤالات بر اساس مشابهت مقدار پارامترهای آنها و نیز بر اساس مقایسه منحنی‌های مشخصه سؤالات، به ۸ گروه زیر تقسیم شدند (جدول ۱).

حال که مشخصات هر یک از گروه‌های سؤالات را بیان کردیم، این مسئله مطرح می‌شود که کدامیک از انواع سؤالات برای سنجش توانایی یا دانش زبان آزمودنی‌ها مناسبتر است. اگر هدف از اجرای آزمون تولیمو سنجش توانایی دانش زبان و تعیین سطوح مختلف توانایی افراد باشد. اولاً دشواری سؤالات آزمون باید به گونه‌ای باشند که احتمال اندکی از پاسخگویی درست به سوال برای افراد برخوردار از سطح توانایی پایین نیز وجود داشته باشد و در تمام موارد با بالا رفتن سطح توانایی، احتمال پاسخگویی درست نیز افزایش یابد. در سؤالات چهارگزینه‌ای مقدار این احتمال برای سطوح پایین توانایی چیزی در حد و حدود احتمال پاسخگویی تصادفی یا پاسخگویی درست از روی شانس است. در واقع سؤالات باید از دشواری متوسط برخوردار باشند. ثانیاً سؤالات باید از پارامتر تشخیص مناسبی نیز برخوردار باشد تا هم بتواند افراد دارای سطح بالای توانایی را از افراد با سطح پایین توانایی صراحتاً متمایز کند و هم سطوح مختلف بین این دو را شناسایی کند. با مفروض دانستن هدف فوق، سؤالات گروه اول و دوم، سؤالات مناسبی نیست، چون مقدار پارامتر دشواری سؤالات این دو گروه خیلی بالا است و احتمال پاسخگویی درست از سوی تواناترین افراد تنها ۰/۵ بوده است و احتمال پاسخگویی درست از سوی افراد با توانایی متوسط و پایین نیز بسیار ناچیز بوده است. در نتیجه این دو گروه برای تعیین سطوح مختلف توانایی دانش زبان مناسب نیستند و با این دو گروه از سؤالات می‌توان تا حدی فقط تواناترین افراد را در زمینه دانش زبان مشخص نمود.

با فرض فوق، سؤالات گروه سوم و چهارم نیز مناسب نیست، چون پارامتر دشواری بالایی دارند و احتمال پاسخگویی درست تنها برای افراد برخوردار از سطح توانایی بالاتر از یک ($\theta > 1$)، به مقدار قابل توجهی بوده و از این سطح به بعد رو به افزایش بوده است. افراد با توانایی متوسط و پایین هم به سختی و با احتمال بسیار اندکی توانسته‌اند به این گروه از سؤالات پاسخ درست بدهند. پس با این گروه سؤالات نیز فقط می‌توان تواناترین افراد را مشخص کرد.

سؤالات گروه پنجم نیز برای تعیین سطوح مختلف توانایی دانش زبان مناسب نیستند چون افرادی که در سطوح متوسط و پایین توانایی قرار گرفته‌اند ($\theta < 1$) با احتمال بسیار ناچیزی توانسته‌اند به این گروه از سؤالات پاسخ درست بدهند. با این گروه از سؤالات تنها می‌توان

جدول ۱: طبقه‌بندی سوالات آزمون تولیمو

| گروه | |
|-------|--|
| اول | این گروه که مشتمل بر ۵ سوال می‌باشد دارای خصوصیات زیر است: پارامتر دشواری ($b = 3$)، پارامتر تشخیص بالاتر از یک ($a > 1$)، و پارامتر حدس کمتر از 0.2 ($c < 0.2$). همان طور که پارامتر دشواری و منحنی مشخصه‌های سوالات نشان می‌دهد، سوالات موجود در این گروه، سوالاتی بسیار دشوار هستند به گونه‌ای که احتمال ارایه پاسخ درست در بالاترین سطح توانایی یعنی $\theta = 3$ فقط 0.5 بوده است و این احتمال برای افراد با توانایی متوسط و پایین هم بسیار ناچیز بوده است. با این سوالات تنها تا حدودی می‌توان تواناترین افراد را از سایرین تشخیص داد و می‌توان گفت که سوالات اصلا قادر به شناسایی افراد با توانایی متوسط و پایین و تمیز آنه از یکدیگر نیست. |
| دوم | این گروه که مشتمل بر ۱۳ سوال می‌باشد دارای خصوصیات زیر است: پارامتر دشواری ($b = 3$)، پارامتر تشخیص بالاتر از یک ($a > 1$) و پارامتر حدس بالاتر از 0.2 ($c > 0.2$). در این گروه ۱۳ سوال جای گرفته است که همانند سوالات گروه قبلی بسیار دشوارند، طوری که احتمال ارایه پاسخ درست به این سوالات از سوی افراد با سطح توانایی بالا $\theta = 3$ تنها 0.5 بوده است و افراد با توانایی متوسط و پایین چندان قادر به پاسخگویی درست به سوالات نبودند. البته با این تفاوت که در این دسته از سوالات احتمال پاسخگویی درست با تکیه بر شانس برای آنها بیشتر بوده یا می‌توان گفت در حد معقولی بوده است. با این گروه از سوالات نیز تنها تا حدودی می‌توان تواناترین افراد را مشخص کرد. |
| سوم | این گروه که مشتمل بر ۱۵ سوال می‌باشد دارای خصوصیات زیر است: پارامتر دشواری ($b = 3$)، پارامتر تشخیص بین 0.5 تا یک ($1 < a < 0.5$)، پارامتر حدس کمتر از 0.15 ($c < 0.15$). این گروه از سوالات شامل ۱۵ سوال می‌شود. با توجه به منحنی مشخصه‌های سوالات می‌توان گفت از سطح توانایی یک ($\theta > 1$) به بالا، به موازات افزایش توانایی احتمال دادن پاسخ درست به این سوالات هم افزایش می‌یابد. اما برای افرادی که روی مقیاس توانایی در سطح متوسط و پایین قرار گرفته‌اند احتمال ارایه پاسخ درست به سوالات بسیار پایین بوده است و با توجه به دشواری سوالات می‌توان گفت که این احتمال حتی برای تواناترین افراد ($\theta = 3$) هم فقط 0.5 بوده است. در نتیجه با این گروه از سوالات هم تنها تا حدودی می‌توان توانایی‌های بالاتر از متوسط و نیز تواناترین افراد را در زمینه دانش زبان شناسایی کرد. |
| چهارم | ۳۴ سوال در این گروه قرار گرفته که ۱۵ سوال دارای پارامتر دشواری $b = 3$ ، ۱۴ سوال پارامتر دشواری بین ۲ تا ۳ ($2 < b < 3$)، و ۵ سوال هم پارامتر دشواری کمتر از ۲ ($b < 2$) دارند. منحنی مشخصه‌های سوالات این گروه شیب بسیار ملایمی دارد زیرا تمام سوالات پارامتر تشخیص پایینی دارند و در نتیجه سوالات از قدرت تمیز کافی برخوردار نیستند و نمی‌تواند افرادی را که در سطح توانایی بالا قرار گرفته‌اند از افرادی که در سطح توانایی پایین قرار دارند به صراحت تفکیک کند. |

| گروه | |
|------|---|
| پنجم | <p>این گروه که مشتمل بر ۱۹ سوال می‌باشد دارای خصوصیات زیر است: دشواری ($۱ < b < ۲/۵$) و تشخیص بالاتر از $۰/۵$ ($a > ۰/۵$) همان گونه که منحنی مشخصه‌های سؤالات نشان می‌دهد از سطح توانایی یک به بالا ($\theta > ۱$)، احتمال پاسخگویی درست به سؤالات، روند افزایشی داشته است، یعنی با بالا رفتن توانایی، احتمال پاسخگویی درست نیز افزایش یافته است. بنابراین این گروه از سؤالات به خوبی توانسته است افرادی را که در زمینه دانش زبان دارای توانایی بالاتر از متوسط بوده‌اند، شناسایی کند. اما در سطوح توانایی متوسط و پایین تر از آن قادر به تفکیک و تشخیص تفاوتها نبوده است.</p> |
| ششم | <p>سه سؤال در این گروه جای دارد با پارامتر دشواری بین $۰/۵$ تا ۱ ($۱ < b < ۰/۵$)، پارامتر تشخیص بزرگتر از $۰/۵$ ($a < ۰/۵$) و پارامتر حدس بزرگتر از $۰/۲$. با توجه به منحنی مشخصه‌های سؤالات از سطح توانایی متوسط به بالا، احتمال پاسخگویی درست به سؤالات رو به افزایش بوده است و با افزایش توانایی، احتمال پاسخگویی درست نیز افزایش یافته است. پارامتر تشخیص سؤالات نیز مطلوب است. با این گروه می‌توان افراد با توانایی متوسط و بالا را در زمینه دانش زبان به خوبی مشخص نمود. اما این سؤالات قادر به تفکیک افراد در سطوح پایینتر از متوسط نمی‌باشد.</p> |
| هفتم | <p>پنج سؤال در این گروه جای دارد با پارامتر دشواری کمتر از $۰/۵$ ($b < ۰/۵$)، پارامتر تشخیص بین $۰/۵$ تا ۱ ($۰/۵ < a < ۱$) است. با توجه به پارامترهای دشواری این سؤالات می‌توان گفت که احتمال پاسخگویی درست به این سؤالات برای کسانی که در سطوح پایین توانایی قرار داشته‌اند کمتر بوده و هر چقدر توانایی افزایش یافته این احتمال نیز افزایش می‌یافته است. البته این موضوع تا سطح توانایی ۲ ($\theta = ۲$) صادق بوده و احتمال پاسخگویی درست برای سطوح مختلف توانایی بالاتر از آن یکسان بوده است. به عبارت دیگر پارامتر تشخیص سؤالات نیز تا حد زیادی مناسب است یعنی این گروه از سؤالات تا سطح توانایی $\theta = ۲$ به خوبی افراد را از یکدیگر متمایز می‌کند، اما از این سطح توانایی به بالا دیگر خوب عمل نمی‌کند. در واقع این گروه از سؤالات نمی‌تواند تواناترین افراد را از افراد توانا تفکیک کند. اما برای سنجش و تمیز سطوح توانایی ضعیف تا بالاتر از متوسط کارایی دارد.</p> |
| هشتم | <p>این گروه فقط شامل یک سؤال می‌شود که دارای پارامتر دشواری $b = -۰/۳۲$، پارامتر تشخیص $a = ۰/۲۶$ و حدس $c = ۰/۳$ است. با توجه به پارامتر دشواری سؤال می‌توان گفت احتمال پاسخگویی درست برای افرادی هم که در سطوح پایین توانایی قرار داشته‌اند، زیاد بوده است. با بالا رفتن توانایی احتمال پاسخ درست به کندی افزایش یافته است طوری که از این لحاظ بین سطوح مختلف توانایی، تفاوتها چندان محسوس نیست. با توجه به پایین بودن مقدار پارامتر تشخیص می‌توان گفت که سؤال نمی‌تواند تفاوتها را در بین توانایی‌های بالاتر از متوسط و در بین سطوح پایینی به صراحت شناسایی کند.</p> |

اشخاصی را که در سطح بالا توانایی دانش زبان قرار گرفته‌اند، شناسایی کرد. سوالات گروه ششم نیز بار نبل به هدف مذکور مناسب نیست زیرا اشخاصی را که در مقیاس توانایی در سطوح پایینتر از متوسط ($\theta < 0$) قرار گرفته‌اند را نمی‌توانند مشخص نمایند با این گروه تنها می‌توان افرادی را که در مقیاس توانایی در سطوح متوسط و بالا هستند مشخص کرد. با این وجود می‌توان ادعا کرد که سوالات این گروه و گروه بعدی مناسب‌ترین سوالات این آزمون بوده‌اند. سوالات گروه هفتم با توجه به داشتن پارامتر تشخیص مناسب و پارامتر دشواری پایین تقریباً تا سطوح توانایی بالا یعنی ($\theta = 2$) خیلی خوب عمل کرده‌اند و به خوبی افراد را تفکیک کرده‌اند، اما قادر به تمایز گذاری بین تواناترین افراد و افراد توانا نبوده‌اند. همان طور که در مورد گروه قبلی گفته شد به رغم این مسایل می‌توان گفت به لحاظ مقایسه‌ای این گروه از سوالات را می‌توان مناسبترین سوالات این آزمون قلمداد کرد.

گروه هشتم هم که تنها شامل یک سؤال می‌شود برای سنجش سطوح مختلف توانایی مناسب نیست. زیرا اگر چه به علت پایین بودن مقدار پارامتر دشواری سؤال، افراد از سطوح پایین‌تر توانایی قادر به پاسخگویی درست به این سؤال بوده‌اند و احتمال پاسخگویی درست نیز با بالا رفتن توانایی افزایش یافته است، اما به دلیل پایین بودن مقدار پارامتر تشخیص سؤال می‌توان گفت که این سؤال نمی‌تواند افراد را در سطوح مختلف توانایی به صراحت تفکیک نماید.

آنچه که تاکنون ذکر شد در مورد مناسب بودن سؤالاها با فرض این موضوع بود که هدف از آزمون برگزار شده، سنجش تمام سطوح مختلف توانایی دانش زبان بوده است.

حال اگر هدف از اجرای آزمون شناسایی و تعیین افراد با توانایی متوسط و بالا در زمینه دانش زبان باشد، سوالات آزمون باید به گونه‌ای باشند که از سطوح توانایی متوسط به بالا افراد قادر به پاسخگویی درست به سوالات باشند و با افزایش سطح توانایی، احتمال پاسخگویی درست نیز افزایش یابد. در واقع دشواری سوالات باید متناسب با سطوح توانایی متوسط و بالا باشد تا بتوان افرادی را که در این سطوح قرار دارند مشخص کرد. علاوه بر این سوالات باید از پارامتر تشخیص مناسبی نیز برخوردار باشند بر این اساس سوالات گروه اول و دوم و سوم مناسب نمی‌باشند، زیرا همان طور که گفته شد، مقدار پارامتر دشواری سوالات این سه گروه بالا است و احتمال پاسخگویی درست برای تواناترین افراد فقط ۵/۰ بوده است و مقدار این احتمال برای افراد با توانایی متوسط هم بسیار اندک بوده است بنابراین با این سوالات تا حدودی می‌توان فقط تواناترین افراد را در زمینه دانش زبان مشخص کرد.

ب- برآورد توانایی آزمودنی‌ها در آزمون تولیمو

متأسفانه هیچ راهی برای از پیش دانستن پارامتر واقعی توانایی وجود ندارد و بهترین راه تنها برآورد آن است. به لحاظ منطقی می‌توان گفت مقدار متوسط برآوردهایی که توسط کامپیوتر محاسبه می‌شود به مقدار پارامتر توانایی آزمودنی‌ها نزدیک است. وقتی پارامتر دشواری سوالات برابر یا نزدیک به پارامتر توانایی آزمودنی باشد، میانگین توانایی برآورد شده نزدیک به مقدار توانایی وی است. نکته حائز اهمیت در برآورد توانایی آزمودنی‌ها در نظریه سؤال-پاسخ این است که توانایی آزمودنی نسبت به سوالات ثابت است یعنی نسبت به سؤال‌هایی که برای سنجش توانایی

مورد نظر به کار می‌روند نامتغیر است. این بدین معنا است که از یک آزمون با هر جایگاهی که در سرتاسر مقیاس توانایی داشته باشد (یعنی با هر میزانی از دشواری) می‌توان برای برآورد توانایی آزمودنی استفاده کرد. بر این اساس توانایی هر یک از امتحان دهندگان آزمون تولیمو با استفاده از نرم‌افزارهای تهیه شده برآورد شد (جدول نتایج پیوست). برای هر آزمودنی مقدار توانایی، نمره براساس توانایی، نمره خام (روش کلاسیک)، رتبه‌بندی براساس نمرات توانایی و بر اساس نمرات خام و اختلاف دو روش رتبه‌بندی ارائه شده است.

در آزمون تولیمو ۷۸۴ آزمودنی شرکت کرده‌اند که بر اساس جدول نتایج پیوست دامنه توانایی آزمودنی‌ها بین ۳/۳۵۸- تا ۴/۶۶۵+ است. توزیع فراوانی توانایی آزمودنیها به قرار زیر است: ۹۰ آزمودنی توانایی بالاتر از ۳+، ۴۰۷ آزمودنی توانایی بین ۲+ تا ۳+، ۲۴۶ آزمودنی توانایی بین ۱+ تا ۲+، ۳۱ آزمودنی توانایی بین صفر تا ۱+ و ۱۰ آزمودنی توانایی زیر صفر دارند.

محاسبه نمره در هر سطح توانایی و مقایسه رتبه‌ها- بعد از برآورد توانایی هر یک از آزمودنی‌ها مرحله بعدی محاسبه نمره هر یک از آنها است. برای محاسبه نمره آزمودنی‌ها در هر سطح از توانایی، احتمال پاسخگویی درست به تک تک سؤالاتی را که آزمودنی به آن پاسخ درست داده است را با هم جمع می‌کنیم. در آزمون تولیمو دامنه نمرات محاسبه شده بر اساس نظریه سؤال-پاسخ بین ۱/۵۴ تا ۸۵/۸۸ است.

پس از محاسبه نمرات هر یک از آزمودنی‌ها، رتبه‌بندی آنها بر اساس نمرات نظریه سؤال-پاسخ و نمرات خام (روش کلاسیک) انجام پذیرفت. آنگاه نمره آزمودنی‌ها بر پایه هر دو روش نمره‌گذاری، رتبه‌بندی شد و نتایج نیز با یکدیگر مقایسه شد. همانطور که در جدول نتایج پیوست مشاهده می‌شود در بعضی موارد نمرات خام آزمودنی‌ها یکسان است، اما نمرات محاسبه شده آنها بر اساس نظریه سؤال-پاسخ یکسان نیست، چون احتمال پاسخگویی درست به سؤالات یکسان و یا مختلف، در سطوح مختلف توانایی متفاوت است. در واقع بر اساس نظریه سؤال-پاسخ هیچ دو آزمودنی نمره مشابه کسب نخواهند کرد و تفاوتی هر چند اندک بین نمرات مشاهده می‌شود.

مقایسه رتبه‌بندی آزمودنی‌ها نشان می‌دهد که فقط در ۳۱ مورد بین دو روش رتبه‌بندی تفاوتی وجود ندارد و این آزمودنی‌ها بر اساس هر دو روش رتبه یکسانی را کسب کرده‌اند. به عنوان مثال از رتبه ۱ تا ۶ نتایج مشابه هستند و هم چنین رتبه‌های ۷۸۱ تا ۷۸۴ نیز یکسان هستند به این معنا که یک آزمودنی بر اساس هر دو روش رتبه ۷۸۴ را کسب کرده است. همچنین آزمودنی دیگری نیز بر اساس هر دو روش نمره‌گذاری رتبه ۳۱۰ را کسب کرده است.

اما در بیشتر موارد بین رتبه‌ها همسانی وجود ندارد. در بعضی از موارد اختلاف رتبه‌ها ناچیز است و در بعضی از موارد این اختلاف خیلی چشمگیر است. به عنوان مثال آزمودنی خاصی، بر اساس نظریه سؤال-پاسخ، رتبه ۴۶۸ را کسب کرده است و بر پایه روش کلاسیک رتبه‌اش ۵۷۳ شده است و اختلاف بین دو روش رتبه‌گذاری برابر ۱۰۵ است. شاید یکی از دلایل این تفاوت‌های چشمگیر بین دو روش رتبه‌گذاری این باشد که این قبیل آزمودنی‌ها به سؤالاتی پاسخ درست داده‌اند که بر طبق تحلیل صورت گرفته مبتنی بر نظریه سؤال-پاسخ، سؤالات مناسبی برای اهداف سنجش بوده‌اند.

اختلاف بین رتبه‌ها، دو الگو دارد. در الگوی اول رتبه محاسبه شده بر اساس نظریه سؤال-پاسخ کمتر از رتبه محاسبه شده بر اساس روش کلاسیک است و در واقع اختلاف بین رتبه‌ها مثبت شده است. به عنوان مثال یک آزمودنی بر اساس نظریه سؤال-پاسخ رتبه 9° کسب کرده است، اما بر اساس نمرات خام (روش کلاسیک) رتبه 10.9° کسب کرده که در اینجا اختلاف رتبه برابر 1.9° است. یا نمونه دیگر آزمودنی دیگری بوده است که بر اساس نظریه سؤال-پاسخ رتبه 35.4° کسب کرده است، اما بر اساس روش کلاسیک رتبه 41.4° را کسب کرده، که اختلاف بین دو رتبه 5.7° است.

در الگوی دوم رتبه کسب شده بر اساس نمرات خام (روش کلاسیک) بالاتر از رتبه‌های مبتنی بر نظریه سؤال-پاسخ است. برای مثال یکی از آزمودنی‌ها بر پایه نمرات خام رتبه‌اش 11.8° شده است، اما بر پایه رویه سؤال-پاسخ رتبه‌اش 13.8° شده است که اختلافی برابر 2° بین دو روش رتبه‌بندی وجود دارد. یا نمونه دیگر آزمودنی است که بر اساس رویه کلاسیک رتبه‌اش 47.6° شده است، اما طبق رویه سؤال-پاسخ رتبه 52.5° را کسب کرده است که اختلافی برابر 4.9° بین دو روش مذکور مشاهده می‌شود.

همان‌طور که قبلاً ذکر شد در 31 مورد بین دو روش رتبه‌گذاری تفاوتی وجود ندارد اما در سایر موارد رتبه‌ها یکسان نیستند. در بعضی موارد اختلاف رتبه‌ها ناچیز و در بعضی موارد بسیار فاحش است. به همین منظور تفاوت بین رتبه‌های دو روش نمره‌گذاری را در 4 طبقه با فاصله طبقاتی 25 گروه‌بندی و فراوانی هر طبقه مشخص شده است (جدول ۲). همان‌طور که در جدول مشخص است اختلاف رتبه زیر 25 بیشترین فراوانی را دارد و با بیشتر شدن اختلاف رتبه‌ها فراوانی کمتر می‌شود. چنانچه فراوانی تفاوت رتبه بالاتر از 76 کمترین مقدار را دارد یعنی همان‌طور که از جدول 3 پیداست بیشترین جابجایی در رتبه‌ها را می‌توان در گروه افزایش رتبه‌های از 1 تا 25 دید که رتبه 28° نفر تغییر کرده است و این میزان تغییر در یک جمعیت حدود 8° نفره تغییر ناچیزی به حساب نمی‌آید. به ویژه که باید توجه داشت که به همین نسبت جابجایی قرینه‌وار یعنی کاهش رتبه را شاهد بوده‌ایم. پس از آن تغییرات را در افزایش رتبه‌ها از 26° تا 5° رتبه شاهد هستیم و کمترین تغییرات را در گروه پایینی تغییرات رتبه‌ها می‌توان سراغ گرفت.

برای مقایسه بیشتر بین نمرات و رتبه‌های دو شیوه سنتی و IRT همبستگی بین سطوح مختلف آن محاسبه شد. جالب توجه است که همبستگی بین نمرات و رتبه‌ها در دو سطح ده درصدی بالا و پایین نمرات بیش از 0.98 بود. همبستگی بین سطوح میانی نمرات نیز بیش از 0.96 بوده است. محدود بودن تغییرات فاحش بین دو روش حاکی از آن است که اندازه‌گیری و نمره‌دهی بر حسب نظریه سؤال-پاسخ به تغییرات فاحش در جابجایی رتبه‌ها نمی‌انجامد، بلکه با توجه به دقت در تعیین پارامترهای سوالات و تکرار محاسبات در برآورد توانایی در چندین دور محاسباتی می‌توان استنتاج کرد که این شیوه اندازه‌گیری تصویری دقیق‌تر از نتایج به دست می‌دهد، علاوه بر این می‌توان ادعا کرد که حتی اگر شاهد همین تغییرات هم نبودیم حداقل نتایج حاصل که مزیت IRT نسبت به رویه سنتی است عبارت است از فراهم ساختن امکان برآورد تواناییها، تعیین احتمال پاسخگویی به هر سؤال روی سطوح مختلف توانایی، ترسیم مقیاس واحدی برای

جدول ۲: توزیع فراوانی جابجایی رتبه‌ها در کل نمرات توانایی صرف نظر از کاهش یا افزایش رتبه

| فرآوانی | اختلاف رتبه |
|---------|-------------|
| ۳۱ | ۰ |
| ۵۷۱ | ۱-۲۵ |
| ۱۳۷ | ۲۶-۵۰ |
| ۳۶ | ۵۱-۷۵ |
| ۹ | ۷۶ به بالا |
| ۷۸۴ | جمع کل |

جدول ۳: توزیع فراوانی اختلاف رتبه در کل نمرات توانایی بر حسب افزایش و کاهش رتبه

| فرآوانی | اختلاف رتبه |
|---------|-------------|
| ۲۸۰ | ۱-۲۵ |
| ۶۸ | ۲۶-۵۰ |
| ۱۸ | ۵۱-۷۵ |
| ۶ | ۷۶ به بالا |
| ۳۱ | ۰ |
| ۲۹۱ | -۱-(-۲۵) |
| ۶۹ | -۲۶-(-۵۰) |
| ۱۸ | -۵۱-(-۷۵) |
| ۳ | ۷۶ به پایین |
| ۷۸۴ | جمع کل |

نشان دادن مشخصات امتحان دهندگان و سوالات، و ترسیم نمودار مشخصه‌های سوالات. اینها امکانات فوق العاده‌ای است که IRT علاوه بر در اختیار گذاشتن نتایج دقیقتر، فراروی ما می‌گذارد. پیش‌بینی می‌شود با حذف سوالاتی که با هدف سنجش همخوانی ندارند، و قدرت تمیزکافی بین توانایی‌های مختلف ندارند، در برآورد توانایی آزمودنی‌ها تفاوت‌هایی مشاهده شود.

۸ برآزش مدل و داده‌ها

همان طور که قبلاً گفته شد مدل انتخابی برای تحلیل آزمون و برآورد توانایی‌ها، مدل سه پارامتری بوده است. مقادیر حاصل از محاسبه مربع خی در اکثر موارد از عدد معیار بزرگتر نبود. با این حال این دو موضوع بایستی مورد بررسی قرار می‌گرفت: یکی موضوع تک بعدی یا چند بعدی بودن مدل و دیگری مسئله نامتغیر بودن پارامترها.

جدول ۴: توزیع فراوانی اختلاف رتبه در ۲۵ درصد بالای نمرات توانایی

| اختلاف رتبه | فراوانی | درصد | درصد تراکمی |
|-------------|---------|------|-------------|
| ۱-۲۵ | ۹۳ | ۴۷,۴ | ۴۷,۴ |
| ۰ | ۱۶ | ۸,۲ | ۵۵,۶ |
| -۱-(-۲۵) | ۸۷ | ۴۴,۴ | ۱۰۰ |
| جمع کل | ۱۹۶ | ۱۰۰ | |

جدول ۵: توزیع فراوانی اختلاف رتبه در ۲۵ درصد پایین نمرات توانایی

| اختلاف رتبه | فراوانی | درصد | درصد تراکمی |
|-------------|---------|------|-------------|
| ۱-۲۵ | ۶۲ | ۳۱,۶ | ۳۱,۶ |
| ۲۶-۵۰ | ۱۳ | ۶,۶ | ۳۸,۲ |
| ۵۱-۷۵ | ۲ | ۱,۰ | ۳۹,۲ |
| ۷۶ به بالا | ۱ | ۰,۵ | ۳۹,۷ |
| ۰ | ۱۲ | ۶,۱ | ۴۵,۸ |
| -۱-(-۲۵) | ۸۱ | ۴۱,۳ | ۸۷,۲ |
| -۲۶-(-۵۰) | ۲۱ | ۱۰,۷ | ۹۸ |
| -۵۱-(-۷۵) | ۳ | ۱,۵ | ۹۹,۵ |
| ۷۶ به بالا | ۱ | ۰,۵ | ۱۰۰ |
| جمع کل | ۱۹۶ | ۱۰۰ | |

جدول ۶: توزیع فراوانی اختلاف رتبه در میانه توزیع نمرات

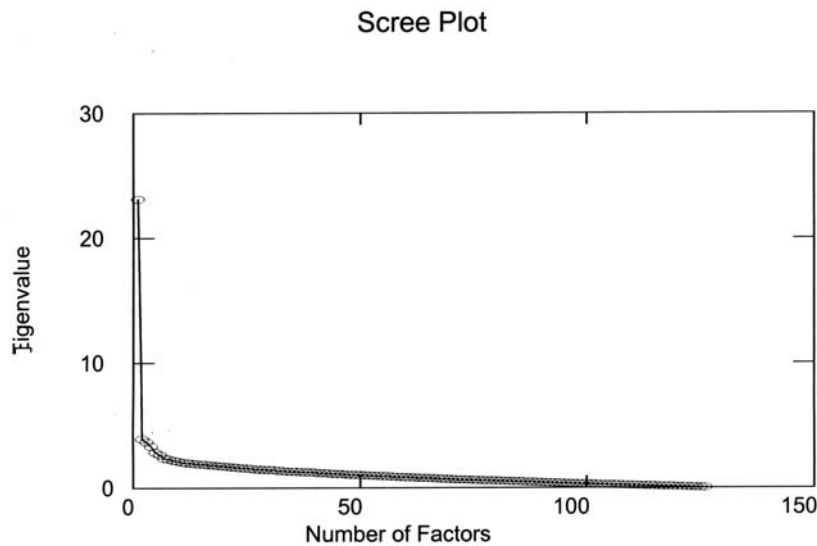
| اختلاف رتبه | فراوانی | درصد | درصد تراکمی |
|-------------|---------|------|-------------|
| ۱-۲۵ | ۱۲۵ | ۳۱,۸ | ۳۱,۸ |
| ۲۶-۵۰ | ۵۵ | ۱۴ | ۴۵,۸ |
| ۵۱-۷۵ | ۱۶ | ۴,۱ | ۴۹,۹ |
| ۷۶ به بالا | ۵ | ۱,۳ | ۵۱,۲ |
| ۰ | ۳ | ۰,۸ | ۵۲ |
| -۱-(-۲۵) | ۱۲۳ | ۳۱,۴ | ۸۳,۲ |
| -۲۶-(-۵۰) | ۴۸ | ۱۲,۲ | ۹۵,۷ |
| -۵۱-(-۷۵) | ۱۰ | ۳,۸ | ۹۹,۵ |
| ۷۶ به بالا | ۲ | ۰,۵ | ۱۰۰ |
| جمع کل | ۳۹۲ | ۱۰۰ | |

۹ ابعاد آزمون

اغلب مدل‌های اندازه‌گیری این فرض را دستمایه کار خود قرار می‌دهند که مفهومی که داریم اندازه می‌گیریم یک بعدی است، یعنی تنها یک عامل برجسته است که مبنای توضیح آن رفتار خاص است. رویکردهای مختلفی برای واریسی و ارزیابی تک بعدی بودن مدل پیشنهاد شده است. (برای اطلاع از آنها می‌توان به کتاب هاتی [5]) مراجعه کرد). روش‌های رایج عبارت است از تعیین تعداد مقادیر ویژه‌ای^{۱۹} که بزرگتر است از ۱، بررسی نمودار اسکری^{۲۰}، و مقایسه نسبت اولین مقدار ویژه به دومین مقدار ویژه. در این تحلیل ما از نمودار اسکری و مقایسه نسبت اولین مقدار ویژه با دومین مقدار ویژه استفاده کرده‌ایم که بر مبنای تکنیک تحلیل عاملی به بررسی ابعاد آزمون می‌پردازد. البته در اینجا برای تعیین تعداد عوامل دخیل در مقیاس، به جای روش تحلیل عناصر اصلی^{۲۱} (PCA)، از روش عامل یابی محور اصلی^{۲۲} (PAF) استفاده کرده‌ایم. زیرا در روش تحلیل عامل یابی محور اصلی فقط واریانس مشترک در محاسبات در نظر گرفته می‌شود در حالی که در روش تحلیل عناصر اصلی هم واریانس اشتراکی و هم واریانس اختصاصی در محاسبات دخیل است.

همان‌طور که قبلاً گفته شد سوالات آزمون تولیمو دو انتخابی هستند یعنی پاسخ آنها از دو حالت صحیح یا غلط بیرون نیست. درباره این گونه داده‌های دو انتخابی (صحیح یا غلط) لازم است بدانیم که تحلیل عاملی بر روی همبستگی تتراکوریک پاسخ‌ها انجام می‌شود. از بین نرم‌افزارهای معروف برای انجام این محاسبه از نرم افزار SYSTAT 10.2 استفاده شد که در محاسبه همزمان همبستگی‌های تتراکوریک و عامل یابی محور اصلی کارایی چشمگیری دارد. در زیر ابتدا مقادیر اولین مقدار ویژه (۲۳/۱۰۸) تا بیستمین مقدار آرایه شده است. با مشاهده نخستین مقدار ویژه می‌توان دریافت که اختلاف آن با مقادیر بعدی بسیار فاحش است. در قسمت بعدی نیز سهم هر یک از عوامل در واریانس همبستگی‌ها دیده می‌شود که سهم عامل نخست نسبت به سایر عوامل بسیار بیشتر است. همچنین نمودار اسکری نیز نشان می‌دهد که عامل‌های دوم به بعد در جایی از نمودار قرار گرفته‌اند که نمودار به حالت تخت نزدیک شده است در حالی که عامل اول با فاصله بسیاری در بالای نمودار قرار گرفته است. بنابراین می‌توان از تک بعدی بودن آزمون و مقیاس مطمئن بود.

19) Eigenvalues 20) scree plot 21) Principal Components Analysis 22) Principle Axis Factoring



مقادیر ویژه

| | | | | |
|--------|-------|-------|-------|-------|
| ۱ | ۲ | ۳ | ۴ | ۵ |
| ۲۳,۱۰۸ | ۳,۹۲۵ | ۳,۶۹۳ | ۳,۳۳۶ | ۲,۸۴۰ |
| ۶ | ۷ | ۸ | ۹ | ۱۰ |
| ۲,۶۵۰ | ۲,۳۵۴ | ۲,۳۰۱ | ۲,۱۹۰ | ۲,۱۳۱ |
| ۱۱ | ۱۲ | ۱۳ | ۱۴ | ۱۵ |
| ۲,۰۳۵ | ۱,۹۸۹ | ۱,۹۶۵ | ۱,۹۰۹ | ۱,۸۸۷ |
| ۱۶ | ۱۷ | ۱۸ | ۱۹ | ۲۰ |
| ۱,۸۷۹ | ۱,۸۲۱ | ۱,۷۹۹ | ۱,۷۵۲ | ۱,۷۴۱ |

از آنجایی که بررسی ابعاد مدل نشان داد مقیاس ما یک بعدی است ناچار باید بررسی شود که آیا پارامترهای سوالات و توانایی آزمودنی‌ها نسبت به هم نامتغیر هستند یا نه. برای این کار گروه چهارم از سوالات که ۳۴ سوال بودند را که با توجه به هدف آزمون سوالات مناسبی نبودند حذف گردید. یعنی این سوالات از قدرت تمیز کافی برخوردار نیستند و نمی‌توانستند در سطوح مختلف توانایی آزمودنی‌ها را به صراحت تفکیک نمایند. پس از حذف این سوالات مجدداً برای هر یک از آزمودنی‌ها بر آورد توانایی بر اساس روش حداکثر درست‌نمایی، محاسبه نمره خام و نمره بر اساس نظریه IRT و رتبه‌بندی آزمودنی‌ها با توجه به ۱۰۶ سوال باقیمانده صورت گرفت. پس از آن نتایج این مرحله و مرحله قبل مورد مقایسه قرار گرفتند تا مشخص شود حذف سوالات چه تاثیری در برآورد توانایی آزمودنی‌ها دارد. بر این اساس میانگین تفاوت توانایی‌های برآورد شده با احتمال ۹۹ درصد در فاصله ۰/۱۰۵۷- تا ۰/۲۷۵۶ قرار گرفته‌اند.

بررسی نتایج نشان می‌دهد که در اکثریت موارد تفاوت بین توانایی‌های برآورد شده بسیار ناچیز است و در فاصله برآورد شده قرار گرفته‌اند و میزان جا به جایی افراد در مقیاس توانایی با حذف سوالات قابل اغماض است. این بیانگر نامتغیر بودن توانایی آزمودنی‌ها نسبت به سوالهایی است که برای برآورد توانایی به کار می‌روند.

برای بررسی استقلال پارامترهای سوالات از نمونه، ۱۰۰ آزمودنی به طور تصادفی از میان ۷۸۴ آزمودنی انتخاب شد. محاسبات پارامترهای سوالات تفاوت معنی‌داری را نشان نمی‌داد.

مراجع

- [1] Baker, F. B. (2001), The basics of item response theory.
- [2] Baker, F. B. (1992), Item response theory: Parameter estimation techniques NewYork Marcel Dekker.
- [3] Birnbaum, A. (1968), Some latent trait models and their use in inferring an examinee's ability. In F.M. Lord & M.R. Novick, Statistical theories of mental test scores. Reading, Mass.: Addison-Wesley Publishing.
- [4] Drasgow, F., Levine M.V., Tsien, S., Williams B.A., & Mead, A.D. (1995), Fitting polytomous item response theory models to multiple-choice tests. Applied Psychological Measurement, 19, 143-165
- [5] Hulin, C.L., Drasgow, F., & Parsons, C.K. (1983), Item response theory Applications to psychological measurement. Homewood, IL: Dow Jones.