

بنام خدا

ششمین کنفرانس بین‌المللی آمار ایران

مجموعه مقالات

۱۳۸۱ تا ۶ شهریور ماه

دانشگاه تربیت مدرس

تهران، ایران

خلاصه مقالات ششمین کنفرانس بین‌المللی آمار ایران

تیراژ : ۱۰۰۰

تنظيم و صفحه‌آرایی : رحیم صفری

چاپ و صحافی : سازمان سنجش آموزش کشور

تاریخ انتشار : تابستان ۱۳۸۱

بنام خدا

پیشگفتار

با یاری خداوند متعال و همکاری انجمن آمار ایران و همت والای همکاران و دانشجویان محترم و حمایت مسئولین محترم دانشگاه و سازمانها و دستگاههای مختلف ششمین کنفرانس بین المللی آمار ایران در روزهای ۴ تا ۶ شهریور ماه ۱۳۸۱ در دانشگاه تربیت مدرس برگزار می‌شود. استقبال چشمگیر اسانید و دانشجویان دانشگاه و همچنین محققین و دست اندکاران فعالیتهای آماری در دستگاهها و سازمانهای مختلف در ارسال مقالات و گزارشات علمی در زمینه‌های مختلف علوم آماری موجب گردید این کنفرانس از سطح علمی مطلوبی برخوردار شود و زمینه‌های ارتباط بین دانشگاه و کاربران آمار در دستگاههای مختلف و واحدهای صنعتی بیش از گذشته فراهم گردد.

متقاضیان شرکت در کنفرانس تعداد ۳۴۰ مقاله به دبیرخانه ارسال نمودند که کمیسیونهای تخصصی شش‌گانه کمیته علمی کنفرانس تعداد ۱۵۷ مورد را بعنوان سخنرانی، تعداد ۷۳ مقاله را برای ارائه بصورت پوستر و ۵ مورد را بصورت کارگاه آموزشی پذیرش نمودند و بدلیل محدودیت مدت زمان برگزاری کنفرانس ۶۹ مقاله را واحد اولویت لازم برای ارائه در کنفرانس ندانستند. بعلاوه از دانشمندان صاحب نام علم آمار داخلی و خارج کشور برای ارائه سخنرانی‌های عمومی و تخصصی دعوت بعمل آمد که تعدادی از آنها به دلایلی از شرکت در کنفرانس عذرخواهی نمودند و نهایتاً بالغ بر ۳۶ مقاله توسط مدعوین داخلی و خارجی در کنفرانس ارائه شد.

برگزاری ششمین کنفرانس آمار ایران بدون همکاری صمیمانه تعدادی از اسانید دانشگاههای سراسر بعنوان هیأت داوران و تلاشهاي خستگی ناپذیر همکاران گرامی دانشگاه علی‌الخصوص همکاران دانشکده علوم پایه و دانشجویان و دانش آموختگان محترم دانشگاه امکان پذیر نبود و جا دارد از تمامی عزیزانی که در برگزاری کنفرانس زحماتی را متحمل شدند قدردانی و تشکر شود.

انجام کلیه امور کامپیوتری و تهیه و تنظیم مجموعه‌های خلاصه و اصل مقالات با همکاری و پشت کاربی دریغ و خستگی ناپذیر آقایان رحیم صفری، امیر کاووسی، شهاب جولانی و سرکار خانم ثریا خفری صورت پذیرفت که از این عزیزان نهایت تشکر و امتنان را دارد.

حمایت کنندگان:

این کنفرانس با حمایتها و پشتیبانی انجمن آمار و سازمانها و موسسات زیر برگزار می‌گردد که بدینوسیله نهایت قدردانی خود را از آنانی که تا زمان چاپ این مجموعه کنفرانس را مورد حمایت‌های مختلف قرار دادند ابراز می‌داریم.

- | | |
|--|-------------------------------|
| ۱ - وزارت علوم، تحقیقات و فناوری | ۲ - وزارت نفت |
| ۳ - وزارت پست و تلگراف و تلفن | ۴ - وزارت جهاد کشاورزی |
| ۵ - وزارت نیرو و برق منطقه تهران | ۶ - وزارت صنایع و معادن |
| ۷ - سازمان سنجش آموزش کشور | ۸ - نهاد ریاست جمهوری |
| ۹ - موسسه پژوهش و برنامه ریزی آموزش عالی | ۱۰ - مرکز آمار ایران |
| ۱۱ - پژوهشکده آمار | ۱۲ - پژوهشکده تربیت بدنی |
| ۱۳ - پژوهشکده هوا فضا | ۱۴ - پژوهشکده صنایع رنگ |
| ۱۵ - دانشگاه امام حسین | ۱۶ - انسیتو پاستور |
| ۱۷ - بانک رفاه کارگران | ۱۸ - بانک کشاورزی |
| ۱۹ - شرکت مهندسی سیستم پاس | ۲۰ - سازمان انتقال خون |
| ۲۱ - شرکت هواپیمایی جمهوری اسلامی ایران | ۲۲ - کمیسیون یونسکو |
| ۲۳ - دانشگاه پیام نور | ۲۴ - وزارت کار و امور اجتماعی |

دبير ششمین کنفرانس بين المللی آمار ايران

شهریور ۱۳۸۱

اعضای کمیته برگزاری:

- ۱- دکتر محسن محمدزاده (مسئول کمیته و دبیر کنفرانس)
- ۲- دکتر محمدقاسم وحیدی اصل (مسئول کمیته علمی)
- ۳- دکتر سقراط فقیه زاده (مسئول کمیته برنامه ریزی)
- ۴- دکتر انوشیروان کاظم نژاد (مسئول کمیته اجرایی)
- ۵- دکتر غلامرضا بابایی (مسئول کمیته تبلیغات)
- ۶- دکتر سید محمد حسینی

اعضای کمیته علمی کنفرانس:

- | | |
|-------------------------------|---|
| دانشگاه شهید بهشتی | ۱- دکتر ابوالقاسم وحیدی اصل (مسئول کمیته) |
| دانشگاه تربیت مدرس | ۲- دکتر محسن محمد زاده (دبیر کنفرانس) |
| دانشگاه شهید بهشتی | ۳- دکتر محمد رضا مشکانی |
| دانشگاه صنعتی اصفهان | ۴- دکتر احمد پارسیان |
| دانشگاه تربیت مدرس | ۵- دکتر سقراط فقیه زاده |
| دانشگاه تربیت مدرس | ۶- دکتر انوشیروان کاظم نژاد |
| دانشگاه علامه طباطبائی | ۷- دکتر حمید رضا نواب پور |
| دانشگاه علوم پزشکی شهید بهشتی | ۸- دکتر یدا... محرابی |
| پژوهشکده آمار | ۹- دکتر عباس گرامی |

فهرست مندرجات

عنوان		صفحة
اثرات طبقه‌بندی نادرست در تحلیل داده‌های گستته	۱	
علیرضا ابدی، کاظم محمد، محمدرضا مشکانی، انوشیروان کاظم‌نژاد		
رکوردها و حجم نمونه	۱۴	
مهدی دوست پرست، جعفر احمدی		
اطلاع در توزیع بر نوع (BURR XII)	۲۶	
جعفر احمدی، مصطفی رزمخواه		
برآورد رگرسیون عام و کاربرد آن در بررسی هزینه	۴۱	
حمید بیدرام، محمد صالحی مرزیجرانی		
برآورد پارامترهای مدل رگرسیونی با رهیافت ماکزیمم آنتروپی	۵۰	
عین‌الله پاشا، محسن محمدزاده، علی آقامحمدی		
مقدمه‌ای بر داده‌کاوی	۵۸	
حمیدرضا نوابپور، علی‌اصغر حائری		
کاهش اربیی انتخاب در آزمایشهای دنباله‌ای	۷۰	
آرزو جبیبی‌راد		
بررسی بعضی روش‌های رگرسیون جایگزین LS و مقایسه آنها با یکدیگر	۷۸	
محمد رضا ریبعی، مجتبی گنجعلی		
تقریب خطی اثرات متغیرهای کمکی تغییرپذیر در مدل کاکس توسط مدل‌های پویا	۹۲	
علیرضا رضائی، سیامک نوربلوچی		

۱۰۷..... روش‌های احتمالاتی در حل مسائل دترمینیستیک	بیژن زنگنه
۱۲۱..... توسعه روش سری زمانی فازی وارائه یک مورد کاوی	عباس سقایی، رسول نورالسناء
?? پیش‌بینی سریهای زمانی به کمک شبکه‌های عصبی	جواد شاکر آرانی
۱۳۳..... رگرسیون ژرفا	حمید شریف، خلیل شفیعی
۱۵۲..... آزمایش‌های آمیزه‌ای و کاربرد آنها در فرمولاسیون مواد غذایی (تهیه نوشیدنی ایزوله پروتئین سوبا)	مجید صیوتی، فتح میکاییلی
۱۶۳.. MCDM با استفاده از DEWMA و EWMA تخمین بهینه وزنها در نمودارهای	رضا کاظم‌زاده، عزیزا... معماریانی، مهدی کرباسیان
۱۷۰..... مقایسه مدل رگرسیون پواسن و مدل رگرسیون دوچمله‌ای منفی در تعیین عوامل موثر بر حاملگی ناخواسته	شبنم کریمی، انوشیروان کاظم‌زاده، سید مهدی سادات هاشمی
۱۸۲..... تعیین پارامتر پنجره در برآوردتابع چگالی توسط B-اسپلاین	محسن محمدزاده، رضا صالحی
۱۹۲..... استفاده از کریگینگ عام در همه‌گیری‌شناسی جغرافیایی بیماریها	محسن محمدزاده، انوشیروان کاظم‌زاده، سفراط فقیه‌زاده، یدالله واقعی
۲۰۵..... شناسایی داده‌های دورافتاده فضایی	محسن محمدزاده، علی محمدیان مصمم
۲۱۵..... سرعت همگرایی الگوریتم EM و تسریع آن	سید محسن میرحسینی، عین الله پاشا

برآورد مینیماکس مجاز پارامتر خانواده کای— دوتبدیل یافته در فضای پارامتر کراندار ۲۲۶	
نادر نعمت الهی، محمد جعفری جوزانی	
مقدمه‌ای بر افسای داده‌ها حمید رضا نواب‌پور، محمد بردبار	۲۳۸
توسعه کنترل فرآیند آماری با استفاده از رویکرد بیز تجربی رسول نورالسناء، محمدرضا لطفی	۲۵۷
تصمیم آماری بیزی بوسیله برآوردگرهای MLM در فضاهای غیرحقیقی حسین نوری امامزاده، مهدی دوست پرست	۲۶۴
انتخاب مدل برای خوشبندی احتمالاتی با استفاده از معیار BIC محمد قاسم وحیدی‌اصل، محسن محمدزاده، محمد قربانی	۲۷۵
کاربرد گرافهای تصادفی بازه‌ای در تحلیل خوشبای محمد قاسم وحیدی‌اصل، محسن محمدزاده، فرانک گودرزی	??
نمایه ??	

اثرات طبقه‌بندی نادرست در تحلیل داده‌های گستته

علیرضا ابدی^۱، کاظم محمد^۲، محمدرضا مشکانی^۳، انوشیروان کاظم نژاد^۴

P ۱۱۰۷۰

^۱ دانشگاه تربیت مدرس

^۲ دانشگاه علوم پزشکی تهران

^۳ دانشگاه شهید بهشتی

^۴ دانشگاه تربیت مدرس

چکیده: طبقه‌بندی نادرست از منابع اریبی است که باعث کاهش کارایی در تحلیل داده‌های گستته می‌شود. در این مقاله با هدف معرفی خطای طبقه‌بندی نادرست ابتدا تاریخچه این موضوع و مباحثت مهم مطرح شده پیرامون آن ارائه و سپس انواع منابع خطا در طبقه‌بندی نادرست بررسی شده است. به منظور نشان دادن چگونگی تاثیر این خطا در تحلیل نتایج مطالعات، شاخصهای مهم قابل محاسبه در یک جدول دو در دو که شامل میزانهای وقوع، خطر نسبی، تعداد موردنیاز برای بیمار و نسبت بخت می‌باشد معرفی و آنگاه در قالب یک مثال فرضی به بررسی تاثیر طبقه‌بندی غلط در نتایج حاصله پرداخته شده است. تاثیر طبقه‌بندی نادرست وقتی که هر یک از متغیرهای وضعیت بیماری یا در معرض خطر بودن یا هر دو متغیر توأم^۱ غلط طبقه‌بندی شده باشند مورد بررسی قرار گرفته و مهمترین روش‌های تعديل تاثیر طبقه‌بندی نادرست براساس روش‌های ماتریسی، مدل سازی و روش‌های بیزی تشریح شده است.

واژه‌های کلیدی: خطای طبقه‌بندی نادرست، تحلیل داده‌های گستته.

۱ مقدمه

خطای اندازه‌گیری متغیرهای با مقیاس اسمی یا رتبه‌ای عموماً خطای طبقه‌بندی نادرست نامیده می‌شود. این خطا در هر نوع مطالعه‌ای ممکن است واقع شود.

در مطالعات گذشته نگراگر تعریف دقیق و روشنی از مواجهه با خطای نداشته باشیم انتساب افراد به گروههای موردنظر با خطا همراه خواهد بود.

^۱ Misclassification Error

همچنین در مواردی که در تشخیص بیماری از تستهای آزمایشگاهی استفاده می‌نمائیم و این روشها از حساسیت^۲ و ویژگی^۳ صدرصد برخوردار نباشد انتساب افراد به گروههای بیمار و غیربیمار نیز با خطای همراه خواهد بود که این امر بیشتر در مطالعات آینده‌نگر واقع می‌شود.

در برخی مطالعات نیز بیماری و مواجهه با خطر توامًا با خطای اندازه‌گیری همراه است که بیشتر در مطالعات مقطعی ممکن است واقع شود.

اندازه‌گیری غیردقیق مواجهه با خطر، تشخیص غلط، گزارشهای فردی غیر صحیح یا هر عامل دیگری ممکن است باعث شود تا فردی بطور غلط دریکی از خانه‌های جدول توافقی قرار گیرد. خطای فوق الذکر عامل انحراف در تعیین شاخصها و تحلیل داده‌ها می‌گردد.

در بسیاری از تحقیقها متأسفانه بدون توجه به احتمالات طبقه‌بندی نادرست به تحلیل نتایج پرداخته می‌شود که می‌تواند نتایج گمراه کننده‌ای به همراه داشته باشد. در نتیجه توجه به ساختار احتمالات طبقه‌بندی نادرست برای دستیابی به نتایج واقع بینانه ضروری است.

موضوع طبقه‌بندی نادرست ابتدا توسط براس (۱۹۵۴) مورد توجه قرار گرفته است. توسط چن (۱۹۸۹) مطالعات مروری در این خصوص صورت پذیرفته و تنبیبن (۱۹۷۲) روش برآورده براساس نمونه‌گیری دو مرحله‌ای را با وجود طبقه‌بندی نادرست ارائه نموده و توسط هیتوسکی و راپ (۱۹۹۲) به دوبار نمونه‌گیری شرطی تعیین داده شده است. اسپلندوهوبی (۱۹۸۷) مدل‌های عمومی خطی را برای داده‌های ناکامل براساس اطلاعات کمکی بکار برده‌اند. لیو ولیانگ (۱۹۹۱) تصحیح خطای طبقه‌بندی نادرست را براساس مدل‌های عمومی خطی ارائه نموده‌اند. هیرونوری و شیز (۲۰۰۰) استنباط در مورد احتمالات طبقه‌بندی نادرست براساس پاسخهای تکرار شده دو تایی را ارائه نموده‌اند.

۲ برخی شاخصها و تحلیلهای در جداول دو طرفه

تحلیل داده‌ها در بسیاری از مطالعات با محاسبه برخی شاخصها یا میزانها در جداول توافقی همراه می‌باشد. جدول فرضی زیر را در نظر بگیرید:

^۲ Sensitivity

^۳ Specificity

		بیمار	غیربیمار	جمع
		A	B	$A + B$
		C	D	$C + D$
جمع		$A + C$	$B + D$	

برخی شاخصهای قابل محاسبه در این جدول را داؤسنون و تراپ به شرح زیر را به نموده‌اند:

- ۱) میزان وقوع در گروه تجربه

$$EER = \frac{A}{A + B} \quad (\text{Experimental Event rate})$$

۲) میزان وقوع در گروه کنترل

$$CER = \frac{C}{C + D} \quad (\text{Control Event Rate})$$

۳) قدر مطلق کاهش خطر

$$ARR = |EER - CER| \quad (\text{Absolute Risk Reduction})$$

۴) تعداد مورد نیاز تیمار برای جلوگیری از یک مورد بیماری

$$NNT = \frac{1}{ARR} \quad (\text{Number Needed to Treat.})$$

۵) خطر نسبی

$$RR = \frac{EER}{CER} \quad (\text{Relative Risk})$$

۶) کاهش خطر نسبی

$$RRR = \frac{ARR}{CER} \quad (\text{Relative Risk Reduction})$$

(۷) نسبت بخت

$$OR = \frac{[A/(A+C)]/[C/(A+C)]}{[B/(B+D)]/[D/(B+D)]} = \frac{A/C}{B/D} = \frac{AD}{BC}$$

براساس علامت گذاریهای فوق الذکر میزان وقوع ER (Event Rate) برای محاسبه اندازه خطر استفاده می‌شود. EER میزان وقوع در گروه تجربه است که براساس تعداد بیمار در گروه مواجهه یافته با عامل خطر محاسبه می‌شود. CER همین میزان در گروه کنترل است که مواجهه با عامل خطر نداشته‌اند. ARR نشانگر قدر مطلق کاهش خطر در مقایسه با خطر مبنا می‌باشد. NNT نمایشگر تعداد افرادی است که باید تحت این تیمار یا درمان قرار بگیرند تا انتظار داشته باشیم یکنفر به بیماری مبتلا نشود. به عبارتی این شاخص فایله درمان را مشخص می‌سازد.

RR نشانگر نسبت بروز در افراد مواجهه یافته با عامل خطر به بروز در افراد مواجهه نیافته می‌باشد. RRR نیز اندازه کاهش خطر نسبی براساس خطر مبنا را مشخص می‌کند. OR نیز شاخص برای بررسی خطر در مطالعات مورد شاهدی است که نسبت حاصل ضرب مقاطع نیز گفته شده است. odds یا بخت، نسبت احتمال وقوع واقعه موردنظر به عدم وقوع آن است که هم می‌تواند برای بروز و هم می‌تواند برای شیوع تعریف شود. بخت و نسبت می‌توانند برای بیان فراوانی بیماری بکار روند. وقتی که نسبت مقدار کوچکی باشد بخت تقریبی از نسبت می‌باشد. برای مثال اگر نسبت سیگاریها در جامعه‌ای ۲/۰ باشد بخت آن عبارت است از:

$$Odds = \frac{۰/۲}{۰/۸} = ۰/۲۵$$

یعنی در جامعه مورد نظر برای هر یکنفر سیگاری ۰ نفر غیرسیگاری وجود دارد. نسبت بخت نیز حاصل تقسیم بخت در گروه مورد به بخت در گروه شاهد است.

۳ تاثیر طبقه‌بندی نادرست در تحلیل اطلاعات

با افزایش تعداد متغیرهایی که مواجه به طبقه‌بندی نادرست می‌باشد تحلیل اطلاعات پیچیده‌تر خواهد شد. در دائره المعرف آمارزیستی موضوع براساس آنالیز تک متغیر و دو متغیره و چندمتغیره بررسی شده است. در حالت تک متغیره A را متغیر گروه‌بندی شده واقعی و A^* را متغیر گروه بندی شده نادرست در نظر می‌گیریم. با فرض استقلال بین واحد‌ها احتمال طبقه‌بندی نادرست بصورت

$$P_r = (A^* = j | A = k) = \theta_{j,k} \quad j, k = ۱, ۲, \dots, m$$

است که پارامتر $\theta_{j,k}$ ساختار طبقه‌بندی غط را نشان می‌دهد و می‌تواند به صورت ماتریس $\phi = [\theta_{j,k}]_{m \times m}$ با عضوهای غیر منفی و جمع ستونهای یک باشد. برای متغیرهای پاسخ دوتایی با فرض $A=2$ (وجود بیماری) و $A=1$ (عدم وجود بیماری) این ماتریس به صورت

$$\begin{pmatrix} \beta & 1-\alpha \\ 1-\beta & \alpha \end{pmatrix} \quad (1)$$

می‌باشد که α نشان دهنده حساسیت و β نشان دهنده ویژگی می‌باشد. اثر بکارگیری طبقه‌بندی جایگزین A^* می‌تواند بصورت زیر خلاصه شود:

$$\Pi_{A^*} = \phi \pi_A \quad (2)$$

در حالیکه

$$\Pi_{A^*} = (\Pi_{A^*}(1), \dots, \Pi_{A^*}(m))'$$

و

$$\Pi_A = (\Pi_A(1), \dots, \Pi_A(m))'$$

نسبتهای جامعه در گروه متغیر جانشین و متغیر اصلی می‌باشند. بنابراین نسبتهای نمونه A^* برآوردهای اribی از π_A خواهند بود: ماهیت این اribی در حالت دوتایی براساس رابطه (2) بصورت زیر خواهد بود:

$$\Pi_{A^*}(1) = \beta \Pi_A(1) + (1-\alpha) \Pi_A(2) \quad (3)$$

$$\Pi_{A^*}(2) = (1-\beta) \Pi_A(1) + \alpha \Pi_A(2) \quad (4)$$

در آنالیز دو متغیره اگر A معرف وضعیت بیماری و B معرف وضعیت مواجهه با خطر باشد ابتدا حالتی را در نظر می‌گیریم که طبقه‌بندی غلط در رابطه با متغیر پاسخ باشد. در این حالت متغیر پاسخ A^* بوسیله A اندازه گیری می‌شود و $\Pi_{A|B}(j|l)$ نشانگر نسبت واحدهایی از جامعه که متغیر A در وضعیت j برای گروه در معرض عامل خطر با $B=l$ باشد و همچنین $\Pi_{A^*|B}(j|l)$ نشانگر نسبت متناظر برای A^* باشد. اگر توجه به تفاوت نسبت پاسخ در دو گروه در معرض نمائیم آنگاه برآوردگر نااریب

$$\Pi_{A^*|B}(2|2) - \Pi_{A^*|B}(2|1)$$

در حالت کلی برای تفاضل واقعی

$$\Pi_{A|B}(2|2) - \Pi_{A|B}(2|1)$$

اریب خواهد بود. در حالتی که هر دو گروه در معرض خطر حساسیت α و ویژگی β یکسان داشته باشند مکانیزم طبقه‌بندی نادرست، غیرافترaci^۱ نامیده می‌شود. در این حالت براساس رابطه (۴) خواهیم داشت.

$$\Pi_{A^*|B}(2|2) - \Pi_{A^*|B}(2|1) = (\alpha + \beta - 1)[\Pi_{A|B}(2|2) - \Pi_{A|B}(2|1)] \quad (5)$$

اندازه‌گیری شده بوسیله A^* با مقدار واقعی A برابر خواهد بود. در حالت طبقه‌بندی نادرست غیرافترaci نسبت

$$\frac{\Pi_{A|B}(2|2)}{\Pi_{A|B}(2|1)}$$

به سمت مقدار تحت فرض صفر کاهش خواهد یافت. در حالت دوم اگر متغیر پاسخ درست گروه‌بندی شده باشد و متغیر در معرض عامل خطر B غلط طبقه‌بندی شده باشد آنگاه A^* گروه‌بندی غلط متغیر B خواهد بود. اگر طبقه‌بندی نادرست B غیرافترaci با توجه به متغیر A باشد آنگاه

$$\begin{aligned} \Pi_{A|B^*}(2|2) - \Pi_{A|B^*}(2|1) &= \frac{(\alpha_B - \beta_B - 1)\Pi_B(2)\Pi_B(1)}{\Pi_{B^*}(2)\Pi_{B^*}(1)} \\ &\times [\Pi_{A|B}(2|2) - \Pi_{A|B}(2|1)] \end{aligned}$$

که

$$\Pi_B = (\Pi_B(1), \Pi_B(2))' = (\Pi_{B^*}(1), \Pi_{B^*}(2))'$$

که به ترتیب نسبتهاي جامعه B و B^* بوده و α_B و β_B حساسیت و ویژگی طبقه‌بندی B هستند. فاکتور $[(\Pi_{A|B}(2|2) - \Pi_{A|B}(2|1)) / (\Pi_{A|B}(2|2) - \Pi_{A|B}(2|1))]$ در (۶) بین صفر و یک است و اثر آن نیز مانند حالت قبل کاهشی می‌باشد.

^۱ Nondifferential

حالت سوم این است که هر دو متغیر A و B در معرض طبقه‌بندی نادرست باشند.
در این حالت زوج جانشین (A^*, B^*) تواماً توسط زوج (A, B) با احتمال طبقه‌بندی
نادرست

$$P_r(A^* = j, B^* = k | A = j, B = k)$$

بیان می‌شود. طبقه‌بندی نادرست A و B مستقل گفته می‌شود اگر:

$$\begin{aligned} P_r(A^* = j^*, B^* = k^* | A = j, B = k) &= P_r(A^* = j^* | A = j, B = k) \\ &\quad \times P_r(B^* = k^* | A = j, B = k) \end{aligned}$$

و همچنین طبقه‌بندی نادرست A و B غیر افتراقی گفته می‌شود اگر:

$$\begin{aligned} P_r(A^* = j^* | A = j, B = k) &= P_r(A^* = j^* | A = j) \\ P_r(B^* = k^* | A = j, B = k) &= P_r(B^* = k^* | B = k) \end{aligned}$$

تحت شرایط استقلال و طبقه‌بندی نادرست غیر افتراقی در A و B گالن و همکاران، نشان داده‌اند که برآورده ناواریب $\Pi_{A^*|B^*}(2|2) - \Pi_{A^*|B^*}(1|1)$ مجددًا تفاوتی کمتر از مقدار واقعی را نشان می‌دهد. از سوی دیگر اگر A و B یا هر دو تحت طبقه‌بندی نادرست افتراقی باشند گلدبرگ، نشان داده که اربیی حاصل لزوماً کاهشی نسبت به مقدار تفاوت واقعی نیست و هر دو جهت را ممکن است بگیرد. در رابطه با آزمون فرض در جداول دو طرفه نیز اگر میان A و B وابستگی وجود نداشته باشد با وجود احتمالات طبقه‌بندی نادرست غیر افتراقی، میان جانشینهای A^* و B^* وابستگی وجود نخواهد داشت. مارشال و همکاران، نشان داده‌اند که آزمون عدم وابستگی میان A و B براساس A^* و B^* دارای سطح معنی داری صحیح اما با توان کمتر خواهد بود.

به منظور روشنتر شدن موضوع با مثالی فرضی اثرات طبقه‌بندی نادرست در تحلیل اطلاعات جداولی دو طرفه وقتی که تشخیص بیماری با احتمالات طبقه‌بندی نادرست مواجه باشد را بررسی می‌کنیم: با در نظر گرفتن ساختار جدول (۱) داده‌های فرضی زیر را به عنوان داده‌های واقعی در نظر می‌گیریم:

اگر روش تشخیصی بکار گرفته شده برای تشخیص بیمار از غیر بیمار حساسیت ۹/۰ و ویژگی ۸/۰ داشته باشد آنگاه براساس احتمالات فوق الذکر امید ریاضی مقادیر مشاهده شده به صورت جدول ۳ خواهد بود:

		مجموع	غیر بیمار	بیمار	
		مواجهه	۷۰	۹۰	۱۵۰
		غیر مواجهه	۲۰	۱۳۰	۱۵۰
		مجموع	۸۰	۲۲۰	۳۰۰

جدول ۱: طبقه‌بندی واقعی بیمار و غیر بیمار در دو گروه

		مجموع	غیر بیمار	بیمار	
		مواجهه	۷۲	۷۸	۱۵۰
		غیر مواجهه	۴۴	۱۰۶	۱۵۰
		مجموع	۱۱۶	۱۸۴	۳۰۰

جدول ۲: طبقه‌بندی بیمار و غیر بیمار بر اساس روش‌های جاری (همراه با خطای)

که فراوانیهای مشاهده شده براساس حساسیت و ویژگی روش‌های تشخیص بیماری بر اساس جدول ۴ بدست آمده است.

براساس جداول ۲ و ۳ برخی شاخصهای قابل محاسبه در جدول ۵ ارائه شده است.

براساس نتایج فوق تأثیر طبقه‌بندی نادرست در برآورد شاخصهای ملاحظه می‌شود. بدیهی است خطای طبقه‌بندی نادرست متاثر از نسبت بیمار به غیر بیمار، حجم مشاهدات در هر گروه، خطای طبقه‌بندی در تشخیص بیماری و تعیین مواجهه و ... می‌باشد و اندازه خطای براساس عوامل فوق الذکر متفاوت خواهد بود. به منظور نشان دادن این اصل که "هر مقدار خطای طبقه‌بندی بیشتر شود برآورد نسبت بخت به سمت مقدار تحت فرض صفر میل خواهد نمود" محاسبات فوق برای تعیین نسبت بخت در حالت‌های مختلف حساسیت و ویژگی در جدول ۶ آورده شده است:

به روشنی ملاحظه می‌شود که نسبت بخت داده‌های واقعی $\frac{4}{3}$ است و با کم شدن حساسیت و ویژگی روش تشخیصی مقدار نسبت بخت به سمت عدد ۱ که همان مقدار تحت فرض صفر است میل می‌کند.

۴ راههای تعديل خطای طبقه‌بندی نادرست

به منظور اصلاح اریبی حاصل از طبقه‌بندی نادرست لازم است اطلاعاتی در مورد ساختار طبقه‌بندی نادرست داشته باشیم. در حالت کلی چون این ساختار نامشخص است باید

	بیمار	غیر بیمار
موجهه	$TP + FP$	$TN + FN$
غیر موجهه	$TP + FP$	$TN + FN$

جدول ۳:

	بر اساس طبقه‌بندی نادرست	بر اساس طبقه‌بندی واقعی اشخاص
<i>EER</i>	۰/۴	۰/۴۸
<i>CER</i>	۰/۱۳	۰/۲۹
<i>ARR</i>	۰/۲۷	۰/۱۹
<i>NNT</i>	۳/۷	۵/۳
<i>RR</i>	۳/۱	۱/۶۵
<i>RRR</i>	۲/۱	۰/۶۵
<i>OR</i>	۴/۳	۲/۲

جدول ۴:

بر اساس داده‌های کمکی به برآورد احتمالات طبقه‌بندی نادرست پردازیم. برای بدست آوردن داده‌های کمکی به برآورد احتمالات طبقه‌بندی نادرست پردازیم. برای بدست آوردن داده‌های کمکی در روش پیشنهادی عبارتند از:

۱ – نمونه‌های معتبر

بر اساس ابزاری که به آن استاندارد طلائی^۱ گفته می‌شود امکان تعیین مقدار واقعی یا دقیقتر A^* وجود دارد. البته در برخی مطالعات ممکن است این روش بسیار گران باشد. احتمال خطأ در استاندارد طلائی صفر نیست اما این مقدار خطأ باید قابل اغماس باشد. مناسب است که نمونه معتبر زیر نمونه‌ای از داده‌های اولیه باشد تا نسبتها را گروه‌بندی A در نمونه معتبر برآوردگر ناریب نسبتها را متناظر در جامعه باشد. روش‌های دیگری در ارتباط با انتخاب نمونه‌های معتبر برای شرایط متفاوت وجود دارد.

۲ – اندازه‌گیریهای مکرر

در صورت عدم وجود استاندارد طلائی با روش اندازه‌گیریهای مکرر به برآورد پارامترهای طبقه‌بندی نادرست می‌پردازیم. تعداد تکرار بستگی به مدل خواهد داشت و این اندازه‌گیریها باید استقلال شرطی داشته باشند. عموماً سه بار اندازه‌گیری برای اکثر مدل‌ها کفایت دارد. در صورت وجود اطلاعات لازم از ساختار طبقه‌بندی غلط با کمک روش‌های زیر می‌توان خطای طبقه‌بندی نادرست را تعدیل نمود.

^۱ Gold Standard

		غیربیمار	بیمار
$Sen = 0/9$	مواجهه	۵۴	۹۶
$Spe = 1$	غیرمواجهه	۱۸	۱۳۲
$OR = 4/1$	مواجهه	۶۶	۸۴
$Sen = 0/8$	غیرمواجهه	۴۲	۱۰۸
$Spe = 0/8$	مواجهه	۷۵	۷۵
$OR = 2$	غیرمواجهه	۵۵	۹۵
$Sen = 0/7$	مواجهه	۸۷	۶۳
$Spe = 0/5$	غیرمواجهه	۷۹	۷۱
$OR = 1/2$	مواجهه	۷۵	۷۵
$Sen = 0/5$	غیرمواجهه	۷۵	۷۵
$Spe = 0/5$	مواجهه		
$OR = 1$	غیرمواجهه		

جدول ۵

الف - روش‌های ماتریسی

با استفاده از روش برآورده بردار نسبتهاي Π_A را برای متغير با گروه‌بندی نادرست A بدست می‌آوریم. فرض کنیم که مجموعه داده اولیه به حجم n_p و یک مجموعه داده معتبر به حجم n_v در دسترس باشد. داده‌های معتبر که برآوردها را مشخص می‌سازند بصورت

$$\hat{\Theta}(A^*|A)$$

در نظر گرفته می‌شوند. برای ماتریس احتمالات طبقه‌بندی نادرست داریم

$$\Theta_{jk} = P_r(A^* = j|A = k)$$

و برآورده ماتریس Π_A بصورت

$$\hat{\Pi}_A^M = \left[\hat{\phi}(A^*|A) \right]^{-1} \hat{\Pi}_{A^*}$$

این برآورده‌گر در مطالعات مختلف کاربرد دارد. با روش دلتا نیز امکان برآورده ماتریس واریانس وجود دارد. این روش برای مسائل ساده تر مناسب است و از آن جمله در مواردی

که تعداد گروه‌بندی متغیرها یا واریانس کم باشد. این روش در مواردی که گروه‌بندی متغیرها زیاد باشد جداول پراکنده و برآوردهای غیر دقیق برای برخی پارامترها را به همراه خواهد داشت.

ب - روش مدل سازی

اگر A را یک مجموعه از متغیرهای با احتمال طبقه‌بندی نادرست و A^* نیز مجموعه متغیرهای جانشین و C مجموعه متغیرهای بدون طبقه‌بندی نادرست باشند و مجموعه متغیر L وابستگی داده‌ها را به گروهی مشخص نماید آنگاه توزیع توام (A, A^*, C, L) می‌تواند تحت دو مدل زیر مشخص شود:

۱) مدلی برای متغیرهای صحیح (A, C, L) که اثر متقابل بین L و (A, C) نشان‌گر تفاوت میان توزیع متغیرهای صحیح در میان نمونه‌ها می‌باشد.

۲) مدلی برای احتمالات طبقه‌بندی غلط بواسیله اثرات متقابل درون A^* و بین A^* با (A, C) مشخص می‌شود. مدل با توجه به متغیرهای صحیح (A, C) اشباع شده می‌باشد غالباً هر دو مدل لگ خطی سلسله مراتبی در نظر گرفته می‌شود اما مدل توأم تولید شده توسط آنها در حالت کلی ممکن است لگ خطی نباشد.

ج - روشهای بیزی

با فرض اینکه A^* نشان‌گر وضعیت متغیر با طبقه‌بندی غلط و A نشان‌گر وضعیت واقعی آن متغیر باشد در حالتی که احتمالات طبقه‌بندی نادرست نامشخص باشند تابع درستنمایی قابل تشخیص نبوده و استنباطهای کلاسیک را غیر ممکن می‌سازد. و از سوی دیگر اگر حساسیت و ویژگی مشخص باشند آنگاه استنباط در مورد نسبت‌ها در جداول توافقی آسان و روشن می‌باشد. اما در حالت واقع گرایانه ما اطلاعات تقریبی از حساسیت و ویژگی روشهای داریم و این اطلاعات را به عنوان بهترین حدسهای برای سایر برآوردها به کار خواهیم برد. براساس روشهای بیزی عدم اطمینان در مورد حساسیت و ویژگی روشهای را با یک توزیع پیشین نشان داده و آنرا با اطلاع پیشین ترکیب کرده و برآوردهای بهتری برای احتمالات طبقه‌بندی نادرست بدست خواهیم آورد.

۵ بحث و نتیجه‌گیری

طبقه‌بندی غلط موجب انحراف در برآوردهای کمیتهای موردنظر براساس مشاهدات می‌گردد. اربیی حاصل از طبقه‌بندی نادرست تابعی از حساسیت و ویژگی روشهای بکار گرفته شده برای طبقه‌بندی است. طبقه‌بندی غلط در برخی حالات موجب تضعیف وابستگی میان

متغیرها می‌گردد. در برخی موارد نیز ممکن است انحرافهای ناشی از خطای طبقه‌بندی غلط در هر جهتی واقع شود. لازم است با انتخاب داده‌های کمکی همچون داده‌های معتبر یا اندازه‌گیریهای مکرر برآورده از ساختار احتمالات طبقه‌بندی نادرست بدست آوریم تا بتوانیم برآوردهای تعدیل شده‌ای بدست آوریم. در حالت طبقه‌بندی نادرست تلاش براین است که اربیی حاصله را تعدیل کنیم و در جهت اصلاح اربیی پیشنهادهایی بدھیم در این موارد در صورت امکان از ملاک MSE به جای واریانس استفاده خواهیم نمود که معمولاً اطلاع واقع بینانه‌تری در اختیار قرار می‌دهد. در این مقاله با یک مثال فرضی اثرات طبقه‌بندی نادرست در برآورد برخی شاخصهای قابل محاسبه در جداول توافقی 2×2 نشان داده شده است. براساس نتایج حاصله مشخص گردیده برآورد شاخصهایی چون نسبت بخت به مقدار زیاد تحت تاثیر اندازه خطای طبقه‌بندی نادرست بوده و بالافراش خطای طبقه‌بندی نادرست اندازه برآورد نسبت بخت به مقدار تحت فرض صفر میل می‌کند. نتایج مورد اشاره با نتایج ارائه شده در دیگر تحقیقها هماهنگی دارد.

مراجع

- Armitage, Peter and Colton (1991) *Encyclopedia of Biostatistics*, John Wiley, 2615 - 2621.
- Bross, I. (1954) *Misclassification in 2×2 tables*, Biometrics 10, 488 - 495.
- Chen, T.T (1989) *A review of methods for misclassified categorical data in epidemiology*, *Statistics in medicine* 8, 1095 - 1106.
- Dawson, Beth and Trapp, Robert (2001) *Basic and Clinical Biostatistics*, MC Graw - Hill, 51 - 54
- Espeland, M. A. & Hui, S.L. (1987) *A general approach to analyzing epidemiologic data that contain misclassification errors*, *Biometrics* 43, 1001-1012
- Haitovsky, Y & Rapp, J. (1992) *Conditional resampling for misclassified multinomial data with applications to sampling inspection*. *Technometrics*

34, 473 - 483.

Hironori, F. and Shizu, I. (2000) *Inference about misclassification probabilities from repeated binary response*, *Biometrics* 56, 706 - 711.

Liv, X. and liang, K. Y. (1991) *Adjustment for nondifferential Misclassification error in the generalized linear model*, *statistics in Medicine* 10, 1197 - 1211

Marshal, J.R.,Priore, R.,Graham, S. (1981) *On the distortion of risk estimates in multiple exposure level case - control studies*, *American Journal of Epidemiology* 113, 464 - 473

Tenenbein, A. (1972) *A double sampling scheme for estimating from misclassified multinomial data with applications to sampling inspection*, *Technometrics* 14, 187 - 202

رکوردها و حجم نمونه

جعفر احمدی، مهدی دوست پرست

P ۱۳۰۵۴

دانشکده علوم ریاضی دانشگاه فردوسی مشهد

چکیده: یکی از مسائل مهم در استنباط و تحلیل آماری برآورد پارامترهای مجھول جامعه می‌باشد. موارد متعددی وجود دارد که بدلاً لیل نامعلومی تعدادی از مشاهدات ناپدید شده‌اند و پژوهشگر سعی در برآورد تعداد آنها دارد، بنابراین یکی از پارامترهای مجھول در این حالت حجم نمونه (n) می‌باشد. روش‌های زیادی وجود دارد که می‌توانند جهت تخمین n بکار گرفته شوند و بکاربردن هر روش نیازمند وجود اطلاعاتی درباره شاخصهایی از جامعه است.

فرض کنید X_1, X_2, \dots, X_r دنباله‌ای از مشاهدات باشد، r را یک رکورد پائین (بالا) می‌گوئیم هرگاه از تمام مشاهدات ما قبل خود کوچکتر (بزرگتر) باشد. در این مقاله، با درنظر گرفتن $N^{(n)}$ بعنوان تعداد رکوردها در نمونه‌ای به حجم n , ابتدا به بررسی خواص توزیعی $N^{(n)}$ می‌پردازیم، سپس براساس $N^{(n)}$, برآوردگرهای نااریب، درستنمایی ماکزیمم و گشتاوری برای n بدست می‌آوریم. در پایان، این برآوردگرهای را مورد مقایسه قرار می‌دهیم.

واژه‌های کلیدی: رکوردهای پائین (بالا)، تقریر (تحدّب)، برآوردگر نااریب، برآوردگر درستنمایی ماکزیمم، برآوردگر گشتاوری، اعداد استرلینگ نوع اول.

۱ مقدمه

نظریه رکوردها شاخه‌ای نسبتاً جدید می‌باشد و رشد اصلی خود را در چند دهه اخیر نموده است. این نظریه علاوه بر ویژگیهای مهم تئوری دارایی کاربردهای خاص عملی است. تغییرات جوی، بعضی از مسائل ترافیک، پیشامدهای طبیعی مانند باد و مسائل مربوط به تعیین مقاومت مصالح و محاسبه احتمال کارافت و ... از جمله کاربردهای مهم آن است.

وقتی ما با مشاهدات متوالی سرکار داریم، مشاهدات تصادفی مربوط به آنها مورد توجه قرار می‌گیرند. مثلاً ممکن است طی سالهای متمادی میزان بارندگی روزانه یا هفتگی را در منطقه‌ای ثبت کرده و بخواهیم در مورد تعداد روزها یا هفته‌هایی که ممکن است میزان بارندگی بیش از گذشته باشد را مطالعه کنیم. یا طراح یک سازه بخواهد بداند که

در طول عمر مفید سازه مورد طرح وی احتمال وقوع زلزله‌ای بزرگتر یا مخربتر از آنچه که در گذشته اتفاق افتاده، چقدر است؟ بعبارت دیگر اگر سازه‌ای براساس بزرگترین زلزله در تاریخ منطقه‌ای مورد نظر محاسبه و طراحی شود، احتمال تنزل کارائی آن در اثر بار زلزله در ۲۰ یا ۵۰ سال آینده چقدر است؟ یا فرض کنید دنباله از تولیداتی را که ممکن است در اثر فشار متلاشی شوند در نظر بگیریم و مایل به تعیین مینیمم فشار متلاشی شدن برای این دنباله از تولیدات باشیم، ما اولین محصول را آزمایش می‌کنیم تا اینکه با حداقل فشار متلاشی شود، X_1 را عنوان فشار گسیختگی اولین محصول ثبت می‌کنیم، سپس محصولات بعدی را در نظر می‌گیریم و X_m را عنوان فشار گسیختگی محصول m ام نامیده و آنرا در صورتی ثبت می‌کنیم که

$$X_m < \min(X_1, \dots, X_{m-1})$$

بنابراین در صورتی فشار گسیختگی را ثبت خواهیم کرد که کمترین مقدار در بین مقادیر ثبت شده قبلی را داشته باشد.

عنوان مثال دیگر X_{ij} را بالاترین سطح آب رودخانه‌ای در روز زام در مکان i ام باشد. فرض کنید علاقمند به مطالعه ماکریم موضعی z_{ij} در هر مکان باشیم. بنابراین مشاهده می‌کنیم گاهی اوقات ماکریم موضعی و گاهی اوقات مینیمم آن برای ما اهمیت دارد. مزیت این نظریه را می‌توان در عدم نیاز به مشاهدات غیر رکوردي، سادگي و غير تقريري بودن برخی نتایج آن بویژه آن دسته که مستقل از توزيع مشاهدات هستند، دانست. تحقیق براساس رکوردها تاکنون مورد توجه متخصصین بسیاری واقع شده است. شاید بتوان گفت چاندلر بطور برجسته‌ای مطالعه مقادیر رکورد را شروع کرد و رزنيک (۱۹۷۳) و شورک (۱۹۷۳)، نظریه مجانبی رکوردها را تکمیل کردند. بعدها، نقطه نظرات جالبی برای نسبت دادن مقادیر رکورد با فرایندهای فرین معین معرفی شدند. ند گلیک (۱۹۷۸) با استفاده از عنوان فریبنده "شکستن رکوردها و شکستن مرزها" تحقیقات بیست و پنج ساله اول را جمع آوری کرد. استنباط آماری براساس رکوردها در سالهای اخیر نیز مورد توجه قرار گرفته است می‌توان به کارلين و گیلفاند (۱۹۹۳)، فیوروبرگ و هال (۱۹۹۸)، گالاتی و پادجت (۱۹۹۴)، احمدی و ارقامی (۱۳۷۵) (۱۳۷۳ ش) (۲۰۰۲) (۲۰۰۱)، (۲۰۰۲ ب) و احمدی (۱۳۷۳ ش) (۲۰۰۰) مراجعه کرد. احمدی و ارقامی (۱۳۷۳ ش) (۲۰۰۱) با ارائه نظریه اطلاع فیشر براساس رکوردها نشان دادند که در بعضی موارد، رکوردها بسیار مناسب تراز مشاهدات و یا آماره‌های ترتیبی عمل می‌کنند. مطالب جامع درباره رکوردها را می‌توان در کتابی با عنوان "رکوردها" نوشته آرنولد، بالاکریشنان و ناگاراجا (۱۹۹۸) یافت.

در این مقاله، بخش ۲ شامل تعاریف پایه‌ای ازداده‌های رکوردی می‌باشد. در بخش ۳

دربارهٔ چگالی و خصوصیات $N^{(n)}$ بحث خواهیم کرد و با توجه به قضایای ارائه شده در بارهٔ آمارهٔ تعداد رکوردها، در بخش ۴ برآوردهای نقطه‌ای ناواریب، درستنمایی ماکزیمم و گشتاوری برای حجم نمونه معرفی می‌کنیم. بخش (۵) شامل مقایسه این برآوردهای می‌باشد.

۲ آماره‌های رکوردی

در این بخش چند تعریف پایه‌ای که در این مقاله از آنها استفاده می‌کنیم، بیان می‌کنیم.

۱.۲ رکورد

فرض کنید X_1, X_2, \dots, X_n از متغیرهای تصادفی هم توزیع، مستقل و از توزیع مشترک F باشند. مشاهده X_i را یک رکورد پائین (بالا) می‌گوئیم هرگاه از همه مشاهدات ما قبل خود کوچکتر (بزرگتر) باشد.

دردامه، خود را به رکوردهای پائین محدود می‌کنیم. نتایج حاصل را می‌توان برای رکوردهای بالا با اندکی تغییر مورد استفاده قرار داد.

۲.۱ زمان رکورد

طبق تعریف ۱.۲، زمانی که در آن رکورد رخ می‌دهد یک متغیر تصادفی می‌باشد که بصورت زیر تعریف می‌شود

$$L(1) = 1, L(k) = \min[j : j > L(k-1), X_j < X_{L(k-1)}]$$

که $L(k)$ زمان رخ دادن k امین رکورد پائین می‌باشد.

۲.۲ تعداد رکوردها

قرار دهید $I_i = \sum_{j=1}^n I_{ij}$ آنگاه $N^{(n)}$ متغیر تصادفی تعداد رکوردهای پائین در بین X_n, \dots, X_1 می‌باشد که متغیرهای نشانگر I_{ij} را بصورت زیر تعریف می‌کنیم

$$I_{ij} = \begin{cases} 1 & X_i = \min(X_1, \dots, X_i) \\ 0 & \text{در غیر این صورت} \end{cases}$$

۴.۲ مثال

فرض کنید دنباله‌ای از مشاهدات با $n = 10$ با مقادیر زیر داریم

$20, 22, 18, 25, 30, 12, 40, 32, 10, 11$

درین صورت طبق تعریف ۱.۲، رکوردهای پائین عبارتند از:

$20, 18, 12, 10$

و طبق تعریف ۲.۲، زمان رخدادن رکوردهای پائین عبارتند از:

$$L(1) = 1, L(2) = 3, L(3) = 6, L(4) = 4$$

و طبق تعریف ۳.۲، تعداد رکوردهای پائین

$$N^{(10)} = 4$$

می‌باشد.

لازم به ذکر است که دنباله رکوردهای بالای نمونه فوق عبارتند از

$20, 22, 25, 30, 40$

۳ خواص توزیعی تعداد رکوردها

درین بخش با ارائه قضایا و لمهایی، رفتار $N^{(n)}$ را بررسی می‌کنیم و با ایده‌هایی که از این بخش می‌گیریم، زمینه لازم برای رسیدن به هدف خود یعنی برآورد $N^{(n)}$ فراهم می‌کنیم.
لم ۱: متغیرهای تصادفی $I_i = 1, \dots, n$ (دارای توزیع $\binom{1}{i}$) هستند و از هم مستقل‌اند.
برهان: برای هر $j < i$

$$\begin{aligned} P(I_i = 1, I_j = 1) &= P(Y_i = \min(Y_1, \dots, Y_i), Y_j = \min(Y_1, \dots, Y_j)) \\ &= P(Y_i = \min(Y_1, \dots, Y_i) > Y_j = \min(Y_{i+1}, \dots, Y_j)) \\ &= P(Y_i = \min(Y_1, \dots, Y_i)) \times P(\min(Y_1, \dots, Y_i) > \min(Y_{i+1}, \dots, Y_j)) \times P(Y_j = \min(Y_{i+1}, \dots, Y_j)) = \frac{1}{i} \frac{j-i}{j} = \frac{1}{ij} = P(I_i = 1)P(I_j = 1) \end{aligned}$$

و به همین ترتیب سایر موارد بطور مشابه ثابت می‌شود.

لم ۲: بدیهی است که

$$E(N^{(n)}) = E(\sum_{i=1}^n I_i) = \sum_{i=1}^n \frac{1}{n}$$

$$Var(N^{(n)}) = \sum_{i=1}^n \frac{1}{i}(1 - \frac{1}{i}) = \sum_{i=1}^n \frac{1}{i} - \sum_{i=1}^n \frac{1}{i^2}$$

نتیجه ۱: از روابط فوق می‌توان یک کران بالا برای احتمال وقوع حداقل تعدادی از رکوردها را با توجه به نامساوی یک طرفه چیشید ارائه داد:

$$P(N^{(n)} \geq r) \leq \frac{Var(N^{(n)})}{Var(N^{(n)}) + [r - E(N^{(n)})]^2}$$

عنوان مثال

$$P(N^{(1000)} \geq 18) \leq 0.05$$

قضیه ۱: تابع جرم احتمال $N^{(n)}$ عبارت است از:

$$g(n; k) = P(N^{(n)} = k) =$$

$$\sum_{i_1=1}^{n-k+2} \sum_{i_2=i_1+1}^{n-k+3} \dots \sum_{i_k=i_{k-1}+1}^n \frac{1}{(i_1-1)(i_2-2)\dots(i_k-1)n}$$

برهان: با توجه به تعریف، اولین مشاهده را رکورد قرار داده ایم پس

$$P(L(2) = n) = P(\text{رکورد پائین دوم در زمان } n \text{ رخ دهد})$$

$$= P(X_1 = \min(X_1, \dots, X_{n-1}), X_n < X_1)$$

$$= \frac{(n-1)! * 1 * 1}{n!} = \frac{1}{n(n-1)}$$

اما برای اعداد صحیح $1 < n_1 < n_2 < \dots < n_k$ داریم

$$P(L(1) = 1, L(2) = n_2, L(3) = n_3, \dots, L(k) = n_k)$$

$$P(I_1 = 1, \dots, I_{n_1-1} = 1, I_1 = 1, I_{n_1+1} = 1, \dots, I_{n_k} = 1)$$

(از لم ۱)

$$= \frac{1}{2} \frac{2}{3} \frac{3}{4} \dots \frac{n_2-2}{n_2-1} \frac{1}{n_2} \frac{n_2}{n_2+1} \dots \frac{1}{n_k} = [(n_2-1)(n_3-1)\dots(n_k-1)n_k]^{-1}$$

و نتیجه حاصل می‌شود.

در بعضی از مراجع، تابع جرم احتمال $N^{(n)}$ که با استفاده از اعداد استرلینگ نوع اول بیان می‌شود بصورت زیر بیان می‌کند:

$$g(n; k) = P(N^{(n)}) = \frac{|S_n^{k-1}|}{(k-1)!} \approx \frac{[ln(n)]^{k-1}}{n(k-1)!}$$

تقریب برای n به اندازه کافی بزرگ برقرار است. برای اثبات می‌توان به ند گلیک (۱۹۷۸) رجوع کرد.

مورنو و دیگران (۱۹۹۶)، تابع مولد احتمال $N^{(n)}$ ، بصورت زیر ارائه شده است :

$$E(S^{N^{(n)}}) = P_n(s) = \frac{\Gamma(n+s)}{\Gamma(n+1)\Gamma(s)} \quad \forall s \in R$$

$$k = 1, \dots, n+1 \text{ و برای } g(1, 1) = 1 : 3 \text{ لم}$$

$$g(n+1; k) = \frac{n}{n+1}g(n; k) + \frac{1}{n+1}g(n; k)$$

برهان :

$$\begin{aligned} g(n+1; k) &= P(N^{(n+1)} = k) \\ &= \sum_{i=1}^k P(N^{(n+1)} = k | N^{(n)} = i)P(N^{(n)} = i) \\ &= P(N^{(n+1)} = k | N^{(n)} = k) \times P(N^{(n)} = k) + \\ &P(N^{(n+1)} = k | N^{(n)} = k-1) \times P(N^{(n)} = k-1) + \dots + \dots \\ &= P(I_{n+1} = 1 | N^{(n)} = k) \times g(n; k) + P(I_{n+1} = 1 | N^{(n)} = k-1) \times g(n; k-1) \\ &= \frac{n}{n+1}g(n; k) + \frac{1}{n+1}g(n; k-1) \end{aligned}$$

لم ۴ : تابع جرم احتمال $N^{(n)}$ ، بازای n ثابت، تابعی مقعر است.

برهان : $\Delta^2(n; k)$ و $\Delta^1(n; k)$ را به صورت زیر تعریف می‌کنیم:

$$\Delta(n; k) = g(n; k+1) - g(n; k)$$

$$\Delta^1(n; k) = \Delta(n; k+1) - \Delta(n; k)$$

کافی است ثابت کنیم :

$$\Delta^1(n; k) \leq 0 \quad \forall k$$

برای $n = 1, 2, 3$ بارسم شکل دیده می‌شود که $g(n; k)$ تابعی مقعر است. فرض

کنیم برای n برقرار باشد. از لم ۳ داریم :

$$\Delta^1(n+1; k) = \frac{1}{n+1}[n\Delta^1(n; k) + \Delta^1(n; k-1)]$$

باتوجه به فرض استقرار $\Delta^1(n; k) \leq 0$ و $\Delta^1(n; k-1) \leq 0$ می‌باشد. لذا $\Delta^1(n+1; k) \leq 0$.

لم ۵: اگر $N^{(n)}$ براساس $g(n; k)$ باشد.

الف) $Mo(N^{(n)}) \leq Mo(N^{(n+1)}) \leq Mo(N^{(n)} + 1)$

ب) برای $n > 2$ اگر و فقط اگر $Mo(N^{(n+1)}) = Mo(N^{(n)}) + 1$

$$n\Delta(n; Mo(N^{(n)})) + \Delta(n; Mo(N^{(n)})) - 1 \geq 0$$

برهان:

الف) از لم ۳ نتیجه می‌شود که

$$\Delta(n+1; k) = \frac{1}{n+1} [n\Delta(n; k) + \Delta(n; k-1)]$$

و با توجه به تقریر $g(n; k)$

$\Delta(n; k-1) > 0$ و $\Delta(n; k) > 0$ آنگاه $k \leq Mo(N^{(n)}) - 1$ اگر

$$\Delta(n+1; k) > 0 \quad (1)$$

آنگاه $\Delta(n; k-1) < 0$ و $\Delta(n; k) < 0$ $k \geq Mo(N^{(n)}) + 1$ اگر

$$\Delta(n+1; k) < 0 \quad (2)$$

فرض کنیم در (۱)، $k = Mo(N^{(n)}) - 1$ لذا

$$g(n+1; Mo(N^{(n)})) > g(n+1; Mo(N^{(n)})-1) > \dots > g(n+1; 1) \quad (3)$$

و همچنین در (۲)، $k = Mo(N^{(n)}) + 1$ لذا

$$g(n+1; Mo(N^{(n)})+1) > g(n+1; Mo(N^{(n)})+2) > \dots > g(n+1; n+1) \quad (4)$$

$Mo(N^{(n)}) \leq Mo(N^{(n+1)}) \leq Mo(N^{(n)}) + 1$ پس

ب) اگر $\Delta(n+1; Mo(N^{(n)})) > 0$ آنگاه

$$g(n+1; Mo(N^{(n)})) + 1 > g(n+1; Mo(N^{(n)}))$$

و با توجه به رابطه‌های (۳) و (۴)

$$\forall k \neq Mo(N^{(n)}) + 1 \quad g(n; k) < g(n+1; Mo(N^{(n)})) + 1$$

پس ۱. بر عکس اگر $Mo(N^{(n+1)}) = Mo(N^{(n)}) + 1$

$$Mo(N^{(n+1)}) = Mo(N^{(n)}) + 1$$

آنگاه با فرض $k = Mo(N^{(n)})$

$$\Delta(n+1; Mo(N^{(n)})) > 0$$

لم ۶ : بازای k ثابت، $n \geq k$ برای $g(n; k)$ دارای خصوصیات زیر است :

الف) برای $3 \leq k \leq n$ $g(k; k) < g(k+1; k)$

$$\lim_{n \rightarrow \infty} g(n; k) = 0$$

ج) $g(n; k)$ تابع تک مدی است یعنی

اگر $g(n; k) > g(n-1; k)$ آنگاه $g(n+1; k) > g(n; k)$

اگر $g(n+1; k) < g(n; k)$ $g(n; k) < g(n-1; k)$

برهان : الف) از تابع جرم احتمال $g(n; k)$ به آسانی بدست می‌آید.

ب) برای n های بزرگ در پیفر (۱۹۹۱) ثابت کرد که

$$g(n; k) = e^{-\lambda_n} \frac{\lambda_n^{k-1}}{(k-1)!} + O\left(\frac{1}{Ln(n)}\right)$$

که $\lambda_n = Ln(n) + \gamma$ و γ ثابت اویلر است و (ب) بدست می‌آید.

ج) از لم ۳ داریم:

$$g(n+1; k) - g(n; k) = \frac{g(n; k-1) - g(n; k)}{n+1} = \frac{-\Delta g(n; k-1)}{n+1}$$

اگر $\Delta g(n; k-1) < 0$ آنگاه $g(n+1; k) - g(n; k) > 0$ ثابت

$Mo(N^{(n-1)}) \leq k-1$ الف) $Mo(N^{(n)}) \leq k$ واژلم ۵ مقرر است لذا

مجددًا از تقریر $\Delta g(n-1; k-1) < 0$ داریم لذا

$$g(n; k) - g(n-1; k) = \frac{-\Delta g(n-1; k)}{n} > 0 \quad (5)$$

حالت دوم بطور مشابه ثابت می‌شود.

برآورد ۴

در این بخش، برآوردهای نااریب، درستنمایی ماکزیمم و گشتاوری با توجه به قضایا و لمهای بخش ۳ بیان می‌کنیم.

قضیه ۲ : ۱ $T_1 = 2^{N^{(n)}} - 1$ یک برآوردهای نااریب n با واریانس

$$Var(T_1) = \frac{(n+3)(n+2)(n+1)}{6} - (n+1)^2$$

برهان : با توجه به $P_n(s)$ تابع مولد احتمال $N^{(n)}$

$$E(T_1) = E(2^{N^{(n)}}) - 1 = P_n(2) - 1 = n$$

$$Var(T_1) = Var(2^{N^{(n)}}) = E(2^{2N^{(n)}}) - [E(2^{N^{(n)}})]^2 =$$

$$P_n(4) - (P_n(2))^2 = \frac{\Gamma(n+4)}{\Gamma(n+1)\Gamma(3)} - (n+1)^2$$

اگر k تعداد رکوردهای پائین مشاهده باشد آنگاه $\hat{n} = T_2$ (برآوردگر درستنمایی ماکزیمم n) آن n ای است که $g(n; k_0)$ را ماکزیمم کند یعنی

$$g(\hat{n}; k_0) = \max_n g(n; k_0)$$

لذا می‌توان برآوردگر درستنمایی n را بصورت زیر بیان کرد.

قضیه ۳ : فرض کنید k تعداد رکوردهای پائین مشاهده شده در نمونه‌ای به حجم n (مجهول) باشد. در این صورت برآوردگر درستنمایی n عبارت است از:

$$\hat{n} = \min\{n | Mo(N_{(n)}) = k_0\}$$

برهان : اگر \hat{n} را ماکزیمم کند، در این صورت $g(\hat{n}; k_0) \leq g(\hat{n}-1; k_0)$ و $\Delta(\hat{n}; k_0 - 1) > \Delta(\hat{n}-1; k_0 - 1) < 0$ از (۵) و از $Mo(N^{(\hat{n})}) \leq g(\hat{n}; k_0)$ تقریباً جرم احتمال $N^{(\hat{n})}$

$$Mo(N^{(\hat{n}-1)}) \leq k_0 - 1 , \quad Mo(N^{(\hat{n})}) \geq k_0. \quad (6)$$

$$Mo(N^{(\hat{n}-1)}) \leq k_0 - 1 < k_0 \leq Mo(N^{(\hat{n})}) \Rightarrow Mo(N^{(\hat{n}-1)}) < Mo(N^{(\hat{n})})$$

اما از لم ۵

$$Mo(N^{(\hat{n}-1)}) \leq Mo(N^{(\hat{n})}) \leq Mo(N^{(\hat{n}-1)}) + 1 \Rightarrow$$

$$Mo(N^{(\hat{n})}) = Mo(N^{(\hat{n}-1)}) + 1$$

از (۶) داریم :

$$Mo(N^{(\hat{n})}) - 1 \leq k_0 - 1 , \quad Mo(N^{(\hat{n})}) \geq k_0 \Rightarrow$$

$$Mo(N^{(\hat{n})}) = k_0 , \quad Mo(N^{(\hat{n}-1)}) = k_0 - 1$$

یعنی \hat{n} جایی قرارداد که با یک واحد کاستن از آن، یک واحد از مد $g(n; k_0)$ کاسته می‌شود. یعنی

$$\hat{n} = \min\{n | Mo(N^{(n)}) = k_0\}$$

واثبات کامل می‌شود.

می‌دانیم که $E(N^{(n)}) = \sum_{k=1}^n \frac{1}{k} = \ln(n) + \gamma$. پس $N^{(n)} = \sum_{k=1}^n I_k$ برای n بصورت زیر ارائه داد:

$$T_2 = \exp(N^{(n)} - \gamma)$$

بازوجه به تابع مولّد احتمال $N^{(n)}$:

$$E(T_2^k) = E(\exp(N^{(n)} - \gamma))^k = P_n(\exp(k))\exp(-k\gamma)$$

چون برآورده n باید مقادیر صحیح بگیرد لذا T_2^* را نزدیکترین عدد صحیح به $\exp(N^{(n)} - \gamma)$ تعریف می‌کنیم.

ارهارد کرامر ثابت کرد هنگامی که $\rightarrow \infty$, $N^{(n)}$, بااحتمال یک

$$\frac{T_2(N^{(n)})}{\exp(N^{(n)} - \gamma)} \rightarrow 1 \quad (4)$$

ممکن است پیشنهاد شود برای n , فاصله $[T_2(k), T_2(k+1)]$ که

$$T_2 = \max\{m \geq k \mid Mo(N^{(m)}) = k\}$$

مناسب بنظر برسد. ولی طول این فاصله بطورنمایی افزایش می‌پابد زیرا

$$T_2 = T_2(k+1) - 1$$

و طول فاصله برابر

$$T_2(k) - T_2(k+1) = T_2(k+1) - T_2(k) - 1$$

$$= \exp(k)\exp(-\gamma)[e - 1]$$

است.

۵ مقایسه

مورنو و دیگران (۱۹۹۶) پیش بینی کردند که برای $3 \leq N^{(n)}$

$$T_1 < T_2 < T_2^*$$

اما مدتی بعد، ارهارد کرامر همین رابطه را با استفاده از اعداد استرلینگ نوع اول ثابت کرد. مناسب بنظر می‌رسد که برای چند مقدار از $N^{(n)}$ برآوردهای مختلف ارائه شود.

$N^{(n)}$	T_1	T_2	T_2^*
۱	۱	۱	۲
۲	۳	۲	۴
۳	۷	۸	۱۱
۴	۱۵	۲۵	۳۱
۵	۳۱	۷۳	۸۳
۶	۶۳	۲۰۴	۲۲۷
۷	۱۲۷	۵۶۵	۶۱۶
۸	۲۵۵	۱۵۵۶	۱۶۷۴
۹	۵۱۱	۴۲۷۲	۴۵۵۰
۱۰	۱۰۲۳	۱۱۶۹۸	۱۲۳۶۷

جدول ۱: مقایسه برآوردها

مراجع

۱. احمدی، ج (۱۳۷۳) «رکوردها در توزیعهای پیوسته» رساله کارشناسی ارشد (آمار ریاضی)، دانشکده علوم ۲، دانشگاه فردوسی مشهد.
۲. احمدی و ارقامی (۱۳۷۵) «رکوردها در توزیع بتای نوع اول» اندیشه آماری ، سال اول، شماره ۲، ۹-۱۳.

- [3] Ahmadi, J. and Arghami, N. A.(2002a)*Nonparametric confidence intervals based record values datas*, accepted for publication in *Statistical papers*.
- [4] Ahmadi, J. and Arghami, N. A.(2002b) *Comparing the Fisher information In record values and IID observations*, to apaer in *Statistics*.
- [5] Ahmadi, J. and Arghami, N. A.(2001) *On the Fisher information in record values*, *Metrika* 53, 3, 195-206.
- [6] Ahmadi, J (2000) Record values - Theory and Applications, Ph.D thsies, Ferdowsi University of Mashhad, Iran.
- [7] Arnold, B. C., Balakrishnan, N. and Nagaraja, H. N. (1992) *A first course in order statistics*, John Wiley, New York.
- [8] Arnold, B. C., Balakrishnan, N. and Nagaraja, H. N. (1998) *Records*, John Wiley, New York.

- [9] Carlin, P. B. and Gelfand, A. E. (1993) *Parametric likelihood inference for record breaking problems*, *Biometrika*, 80, 3, 507-515.
- [10] Erhard Gramer. *Asymptotic properties of estimators of the sample size in a record model*.
- [11] Feuerverger, A. and Hall, P. (1998) On statistical inference based on record values, *Extremes*, 12, 169-190.
- [12] Glick,N.(1978).*Breaking records and breaking boards*,Am.math. Monthly 85,2-26.
- [13] Gulati, S. and Padgett, W. J. (1994) Smooth nonparametric estimation of the distribution and density functions from record-breaking data, *Comm. Stat.-Theory Methods*, 23(5), 1259-1274.
- [14] Moreno-Rebollo, J.L, Barranco-Chamorro,I.,L. pez-Bl. squeez, F.and G..mez-G. mez,T. (1996). *On the estimation of the unknown sample size from the records*. statist. Probab. letters 31, 7-12.
- [15] Resnick,S.I. (1973)*Records values and maxima*, Ann. Probab.1,650-662
- [16] Shorrock, R. W. (1973) *Record values and inter-record times*, J. Appl. Probab. 10, 543-555.

اطلاع در توزیع بر نوع ۱۲ (BURR XII)

جعفر احمدی، مصطفی رزمخواه

P ۱۴۲۸۳

گروه آمار، دانشگاه فردوسی مشهد

چکیده: خانواده توزیعهای بر در سال ۱۹۴۲ توسط بر معرفی شده است. یکی از اعضای این خانواده توزیع بر نوع ۱۲ می‌باشد که دارای کاربردهای وسیعی بویژه در کنترل کیفیت، تحلیل داده‌های نقص فنی، کارآیی مسکن‌ها در آزمایشات پزشکی و مدل‌های قابلیت اطمینان در سیستمهای مهندسی می‌باشد.

این توزیع در مقالات بصورت دوپارامتری معرفی شده است. برای بدست آوردن ماتریس اطلاع فیشر به محاسبه انگرالهایی مواجه می‌شویم که به سادگی قابل حل نیستند. در این مقاله ابتدا امید ریاضی تعدادی توابع پایه‌ای از این توزیع را بدست می‌آوریم و آنگاه نشان خواهیم داد که تمام امید ریاضی‌های مربوط به ماتریس اطلاع فیشر را می‌توان از نتایج بدست آمده محاسبه نمود. بعلاوه با روش پارامترسازی، یک پارامتر مقیاس برای این توزیع معرفی می‌کنیم و ماتریس اطلاع فیشر را برای حالت سه‌پارامتری بدست می‌آوریم.

واژه‌های کلیدی: تابع درستنمایی، برآوردگر درستنمایی ماکزیمم، ماتریس اطلاع فیشر.

۱ مقدمه

منشأ کار روی این خانواده از توزیعها مقاله‌ای بود که در سال ۱۹۴۲ توسط بر(BURR) ارائه شد که در آن وی قصد دارد تا توزیع‌هایی را که می‌توانند فرم‌های بسیار متنوعی را بگیرند و اعطاف پذیری زیادی دارند ایجاد کند.

خانواده بر توسط جواب‌هایی برای یک معادله تفاضلی در تابع توزیع تجمعی معرفی شده است. یک جواب ویژه که باعث بوجود آمدن توزیع B12 شده، در عمل بسیار پر کاربرد واقع گشته است؛ می‌توان به برخی از این کاربردها در بخش‌هایی از کنترل کیفیت، طول عمر یا مدل بندی زمان نقص، مدل بندی توزیع درآمد، bio-assay و آزمون فرض اشاره نمود.

بر(۱۹۴۲) تعدادی از فرم‌های توابع توزیع تجمعی را توصیف کرد که در عمل معلوم

شده است که برای برازش داده‌ها مفید می‌باشد. هدف از استفاده از این توابع توزیع تجمعی، آسان کردن تحلیل‌های ریاضی تا رسیدن به یک برازش منطقی به داده‌ها با استفاده از روش گشتاورها می‌باشد. این سیستم توزیع‌ها با در نظر گرفتن توابع توزیعی که برای معادله تفاضلی زیر مناسب هستند معرفی می‌شوند:

$$dF = F(1 - F)g(x)dx$$

که $1 \leq F \leq 0$ و $g(x)$ یک تابع مناسب و نامنفی برای هر محدوده x است. جواب این معادله تفاضلی به ازای $g(x)$ داده شده عبارتست از:

$$F(x) = \frac{1}{1 + e^{-G(x)}}$$

که در آن

$$G(x) = \int_{-\infty}^x g(u)du$$

بر، دوازده جواب برای این معادله تفاضلی (با در نظر گرفتن انتخاب (g)) ارائه کرده است که دوازدهمین آنها توزیع بر نوع ۱۲ نام گرفته است.

اگر X بصورت توزیع بر نوع دوازده توزیع شده باشد، از نماد ۱۲ B استفاده می‌کنیم. برای این توزیع، تابع توزیع احتمال بصورت زیر معرفی می‌گردد

$$F(x) = 1 - (1 + x^c)^{-k}, \quad x > 0$$

و تابع چگالی احتمال عبارتست از

$$f(x) = kcx^{c-1}(1 + x^c)^{-(k+1)}, \quad x > 0 \quad (1)$$

این تابع چگالی تک مدلی با مد $\left(\frac{c-1}{kc+1}\right)^{\frac{1}{c}}$ برای $c > 1$ می‌باشد. و اگر $c \leq 1$ این چگالی فرم L به خود می‌گیرد.

در بسیاری از پژوهش‌های آماری، متغیرهایی که مدل می‌شوند توابعی از تعدادی متغیرهای دیگر هستند. بنابراین بیان تعدادی از روابط بین توزیع ۱۲ B و تعدادی دیگر از توزیع‌های معروف برای تسهیل چنین مدل‌بندی‌هایی مفید است. اگر $X \sim B(12, p)$ باشد.

$$U = (1 + X^c)^{-1} - 1$$

U دارای توزیع بتای نوع اول با پارامترهای k و 1 می‌باشد. $V = X^c - 2$ دارای توزیع بتای نوع دوم (پیرسون نوع ۶) با پارامترهای k و 1 می‌باشد.

$W = kX^c - 3$ دارای توزیع فیشر با درجات آزادی (۲, ۲k) می‌باشد. شایان توجه است که Timr.L.Fry.(1993) خواص خانواده توزیع بر را در حالت‌های یک متغیری و دو متغیری بررسی نموده است، همچنین احمدی (۱۳۷۳) بر اساس رکوردها این توزیع را توصیف کرد.

هدف اصلی این مقاله، محاسبه ماتریس اطلاع فیشر می‌باشد. می‌دانیم که اطلاع فیشر در نامساوی کرامر- رائو برای بدست آوردن کران پایین واریانس برآوردهای ناارب و همچنین در واریانس توزیع حدی برآوردهای درستنمایی ماکزیمم ظاهر می‌شود.

اخیراً بررسی اطلاع فیشر مورد توجه متخصصین زیادی قرار گرفته است، از جمله می‌توان به (Ahmadi j and Arghami NR (1996)، Zheng G and Park S (1996) و Gastwirth JL (2000) اشاره نمود.

در بخش دوم این مقاله ابتدا به تعریف اطلاع فیشر می‌پردازیم و سپس لمحهای را اثبات می‌کنیم که در محاسبه ماتریس اطلاع فیشر مورد استفاده قرار می‌گیرند. در بخش سوم لگاریتم درستنمایی و مشتقات آن را محاسبه می‌کنیم. در بخش چهارم امید ریاضی توابعی پایه‌ای از این توزیع را که در تعیین ماتریس اطلاع فیشر مفید واقع می‌گردد محاسبه می‌کنیم. در بخش پنجم ماتریس اطلاع فیشر را معرفی می‌کنیم. و بالاخره در بخش ششم با پارامترسازی یک پارامتر مقیاس برای این توزیع تعریف کرده و ماتریس اطلاع فیشر متناظر را برای حالت سه‌پارامتری بدست می‌آوریم.

۲ نتایج مقدماتی

تعریف اطلاع فیشر: مقدار اطلاع فیشر درباره پارامتر θ در متغیر تصادفی X ، که دارای چگالی $f(x; \theta)$, $\theta \in A$ می‌باشد، برابر است با:

$$I_X(\theta) = E\left(\left[\frac{\partial}{\partial \theta} \log f(X; \theta)\right]^2\right) \geq 0$$

امید ریاضی فوق را نسبت به چگالی $f(x; \theta)$ حساب می‌کنیم و فرض می‌کنیم که وجود داشته باشد.

ماتریس اطلاع فیشر: اطلاع فیشر را می‌توان برای θ موقعی که چند بعدی باشد تعیین داد و یک ماتریس به نام ماتریس اطلاع فیشر بدست آورد. مثلاً برای $(\theta_1, \theta_2) = \theta$ موضوع را شرح می‌دهیم.

فرض کنید X دارای چگالی $f(x; \theta_1, \theta_2)$ باشد. همچنین فرض می‌کنیم تابع f از

نظر مشتق‌گیری و انتگرال‌گیری دارای شرایط لازم باشد. ماتریس 2×2 زیر را ماتریس اطلاع برای $(\theta_1, \theta_2) = \theta$ می‌نامیم.

$$I_X(\theta_1, \theta_2) = I_X(\theta) = \begin{pmatrix} I_{11}(\theta) & I_{12}(\theta) \\ I_{21}(\theta) & I_{22}(\theta) \end{pmatrix}$$

لم ۱. فرض کنید X داری توزیع B ۱۲ باشد، آنگاه

$$E(X^r) = \frac{\Gamma(\frac{r}{c} + 1)\Gamma(k - \frac{r}{c})}{\Gamma(k)}$$

که در آن $\Gamma(\cdot)$ تابع گاما می‌باشد.

اثبات: با استفاده از (۱) داریم

$$E(X^r) = ck \int_0^\infty x^{c+r-1} (1+x^c)^{-(k+1)} dx$$

که با استفاده از تغییر متغیر

$$\begin{aligned} E(X^r) &= k \int_1^\infty (t-1)^{\frac{r}{c}} t^{-(k+1)} dt = k \int_0^1 z^{k-\frac{r}{c}-1} (1-z)^{r/c} dz \\ &= kB(k - \frac{r}{c}, \frac{r}{c} + 1) = \frac{\Gamma(k - \frac{r}{c})\Gamma(\frac{r}{c} + 1)}{\Gamma(k)} \end{aligned}$$

لم ۲. فرض کنید X داری توزیع B ۱۲ باشد، آنگاه

$$E_k\left(\frac{X^c \log X}{1+X^c}\right) = \frac{k}{k+1} E_{k+1}(X^c \log X)$$

اثبات: با استفاده از (۱)

$$\begin{aligned} E_k\left(\frac{X^c \log X}{1+X^c}\right) &= \int_0^\infty \left(\frac{x^c \log x}{1+x^c}\right) c k x^{c-1} (1+x^c)^{-(k+1)} dx \\ &= \int_0^\infty (x^c \log x) c \frac{k}{k+1} (k+1) x^{c-1} (1+x^c)^{-(k+1)+1} dx \\ &= \frac{k}{k+1} E_{k+1}(X^c \log X) \end{aligned}$$

لم ۳. فرض کنید X داری توزیع B ۱۲ باشد، آنگاه

$$E_k\left[\frac{X^c (\log X)^\gamma}{(1+X^c)^\gamma}\right] = \frac{k}{k+\gamma} E_{k+\gamma}[X^c (\log X)^\gamma]$$

اثبات:

$$\begin{aligned} E_k \left[\frac{X^c (\log X)^\gamma}{(1 + X^c)^\gamma} \right] &= \int_0^\infty \frac{x^c (\log x)^\gamma}{(1 + x^c)^\gamma} c k x^{c-1} (1 + x^c)^{-(k+1)} dx \\ &= \int_0^\infty x^c (\log x)^\gamma c \frac{k}{k+\gamma} (k+\gamma) x^{c-1} (1 + x^c)^{-(k+1)+\gamma} dx \\ &= \frac{k}{k+\gamma} E_{k+\gamma} [X^c (\log X)^\gamma] \end{aligned}$$

لم ۴. فرض کنید X داری توزیع B ۱۲ باشد، آنگاه

$$E_k \left(\frac{\log X}{1 + X^c} \right) = \frac{k}{k+1} E_{k+1} (\log X)$$

لم ۵. فرض کنید X داری توزیع B ۱۲ باشد، آنگاه

$$E_k \left[\frac{X^c \log X}{(1 + X^c)^\gamma} \right] = \frac{k}{k+\gamma} E_{k+\gamma} (X^c \log X)$$

لم ۶. برای تابع چگالی احتمال $f(x; \theta_1, \dots, \theta_k)$ تحت شرایط نظم داریم
الف) امید ریاضی متغیر نمره صفر است، یعنی

$$E \left(\frac{\partial \ell(\theta_1, \theta_2, \dots, \theta_k)}{\partial \theta_i} \right) = 0 \quad i = 1, 2, \dots, k$$

(ب)

$$Var \left(\frac{\partial \ell}{\partial \theta_i} \right) = -E \left(\frac{\partial^2 \ell}{\partial \theta_i^2} \right) \quad i = 1, 2, \dots, k$$

(ج)

$$Cov \left(\frac{\partial \ell}{\partial \theta_i}, \frac{\partial \ell}{\partial \theta_j} \right) = -E \left(\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right) \quad i, j = 1, 2, \dots, k$$

که در آن

$$\ell(\theta_1, \theta_2, \dots, \theta_k) = \sum_{i=1}^n \log f(x_i; \theta_1, \theta_2, \dots, \theta_k)$$

۳ لگاریتم درستنماهی و مشتقهای آن

برای بدست آوردن ماتریس اطلاع فیشر این توزیع، ابتدا لازم است تابع درستنماهی مربوطه را تشکیل داده و مشتقهای لگاریتم درستنماهی را محاسبه کنیم. برای این منظور فرض می‌کنیم X_1, X_2, \dots, X_n یک نمونه تصادفی از توزیع (۱) باشد. می‌دانیم که تابع درستنماهی یک نمونه تصادفی عبارتست از

$$\ell = \sum_{i=1}^n \log f(x_i; c, k)$$

آنگاه از (۱) داریم

$$\ell = n \log(ck) + (c - 1) \sum_{i=1}^n \log x_i - (k + 1) \sum_{i=1}^n \log(1 + x_i^c)$$

مشتق جزئی ℓ نسبت به c عبارتست از

$$\frac{\partial \ell}{\partial c} = \frac{n}{c} + \sum_{i=1}^n \log x_i - (k + 1) \sum_{i=1}^n \frac{x_i^c \log x_i}{1 + x_i^c} \quad (2)$$

همینطور داریم

$$\frac{\partial \ell}{\partial k} = \frac{n}{k} - \sum_{i=1}^n \log(1 + x_i^c) \quad (3)$$

از طرفی بسادگی می‌توان نشان داد که

$$\frac{d}{dc} \sum_{i=1}^n \frac{x_i^c \log x_i}{1 + x_i^c} = \sum_{i=1}^n \frac{x_i^c (\log x_i)'}{(1 + x_i^c)^2}$$

ولذا مشتقهای جزئی مرتبه دوم لگاریتم درستنماهی، به شرح ذیل می‌باشد

$$\frac{\partial^2 \ell}{\partial k^2} = -\frac{n}{k^2} \quad (4)$$

$$\frac{\partial^2 \ell}{\partial c^2} = -\frac{n}{c^2} - (k + 1) \sum_i^n \frac{x_i^c (\log x_i)'}{(1 + x_i^c)^2} \quad (5)$$

$$\frac{\partial^2 \ell}{\partial c \partial k} = -\sum_i^n \frac{x_i^c (\log x_i)'}{(1 + x_i^c)^2} \quad (6)$$

۴ امید ریاضی توابع پایه‌ای

بدست آوردن امیدهای مشتقات جزئی مرتبه اول و دوم لگاریتم درستنمایی، مستلزم محاسبه امید ریاضی توابع پایه‌ای زیر می‌باشد

$$\log X, \log(1 + X^c), \frac{X^c(\log X)}{1 + X^c}, \frac{X^c(\log X)^2}{(1 + X^c)^2}$$

فرض کنید X دارای توزیع B ۱۲ باشد، در اینصورت با محاسبه امید توابع یاد شده نتایج زیر را داریم.

نتیجه ۱.

$$E(\log X) = -\frac{\gamma + \psi(k)}{c} \quad (\textcircled{V})$$

که ψ و γ ثابت اویلر است.

اثبات : بنا به $\lim_{r \rightarrow 0} E(X^r)$ داریم

$$E(X^r) = \frac{\Gamma(\frac{r}{c} + 1)\Gamma(k - \frac{r}{c})}{\Gamma(k)} \quad (\textcircled{A})$$

با مشتق گیری از طرفین (A) نسبت به r داریم

$$E(X^r \log X) = \frac{\Gamma'(\frac{r}{c} + 1)\Gamma(k - \frac{r}{c}) - \Gamma(\frac{r}{c} + 1)\Gamma'(k - \frac{r}{c})}{c\Gamma(k)} \quad (\textcircled{B})$$

با فرض $r = 0$ در (B)

$$\begin{aligned} E(\log X) &= \frac{\Gamma'(1)\Gamma(k) - \Gamma(1)\Gamma'(k)}{c\Gamma(k)} \\ &= \frac{\Gamma'(1) - \frac{\Gamma'(k)}{\Gamma(k)}}{c} = -\frac{\gamma + \psi(k)}{c} \end{aligned}$$

نتیجه ۲.

$$E[\log(1 + X^c)] = \frac{1}{k} \quad (\textcircled{C})$$

اثبات : بدیهی است هر گاه X دارای توزیع (۱) باشد، آنگاه $Y = 1 + X^c$ از یک توزیع پارتو با تابع چگالی احتمال $1 - f(y) = ky^{-(k+1)}$ ، $y \geq 1$ پیروی می‌کند. همینطور $Z = \log Y$ ، توزیع نمایی منفی با میانگین k^{-1} دارد. بنابراین نتیجه می‌شود که

$$E[\log(1 + X^c)] = \frac{1}{k} \quad (۴)$$

نتیجه ۳.

$$E[(\log X)^r] = \frac{\frac{\pi^r}{1} + \gamma^r + 2\gamma\psi(k) + (\psi(k))^r + \psi'(k)}{c^r} \quad (۵)$$

اثبات : با مشتقگیری از طرفین (۹) بدست می‌آوریم

$$E[X^r(\log X)^r] = \frac{\Gamma''(\frac{r}{c} + 1)\Gamma(k - \frac{r}{c}) - 2\Gamma'(\frac{r}{c} + 1)\Gamma'(k - \frac{r}{c}) + \Gamma(\frac{r}{c} + 1)\Gamma''(k - \frac{r}{c})}{c^r \Gamma(k)}$$

با فرض $r = 0$ در عبارت فوق داریم

$$\begin{aligned} E[(\log X)^r] &= \frac{\Gamma''(1)\Gamma(k) - 2\Gamma'(1)\Gamma'(k) + \Gamma(1)\Gamma''(k)}{c^r \Gamma(k)} \\ &= \frac{\frac{\pi^r}{1} + \gamma^r + 2\gamma\psi(k) + (\psi(k))^r + \psi'(k)}{c^r} \end{aligned}$$

که ψ' مشتق تابع ψ می‌باشد.

نتیجه ۴.

$$E_k\left(\frac{X^c \log X}{1 + X^c}\right) = \frac{k}{k+1} E_{k+1}(X^c \log X) \quad (۶)$$

(نوشتن امید بصورت E_k برای تاکید بر نقش k در تابع چگالی احتمال توزیع مورد نظر می‌باشد).

اثبات : با توجه به لم ۲ داریم

$$E_k\left(\frac{X^c \log X}{1 + X^c}\right) = \frac{k}{k+1} E_{k+1}(X^c \log X)$$

حال با استفاده از معادله (۹) و جایگذاری $r = c$ و همچنین جایگزین کردن $1 + k$ به k داریم

$$E_k\left(\frac{X^c \log X}{1 + X^c}\right) = \frac{k}{k+1} \left[\frac{\Gamma'(2)\Gamma(k) - \Gamma(2)\Gamma'(k)}{c\Gamma(k+1)} \right]$$

$$= \frac{1 - \gamma - \psi(k)}{c(k + 1)}$$

نتیجه ۵. برای تابع ψ رابطه بازگشتی زیر را داریم

$$\psi(k + 1) = \psi(k) + \frac{1}{k} \quad (12)$$

اثبات: می‌توانیم بنویسیم

$$E_k\left(\frac{X^c \log X}{1 + X^c}\right) = E_k(\log X) - E_k\left(\frac{\log X}{1 + X^c}\right)$$

با استفاده از لم ۴ داریم

$$E_k\left(\frac{X^c \log X}{1 + X^c}\right) = E_k(\log X) - \frac{k}{k + 1} E_{k+1}(\log X)$$

حال با استفاده از معادلات (۷) و (۱۲) داریم

$$\frac{1 - \gamma - \psi(k)}{c(k + 1)} = -\frac{\gamma + \psi(k)}{c} + \frac{k}{k + 1} \left(\frac{\gamma + \psi(k + 1)}{c} \right)$$

$$\frac{1 - \gamma - \psi(k)}{c(k + 1)} = \frac{-\gamma k - \gamma - k\psi(k) - \psi(k) + k\gamma + k\psi(k + 1)}{c(k + 1)}$$

$$1 = -k\psi(k) + k\psi(k + 1)$$

نهایتاً داریم

$$\psi(k + 1) = \psi(k) + \frac{1}{k}$$

نتیجه ۶.

$$E_k\left[\frac{X^c (\log X)^r}{(1 + X^c)^r}\right] = \frac{k}{c^r (k + 1)(k + 2)} \left[\frac{\pi^r}{\Gamma} + \gamma^r - 2\gamma \right.$$

$$\left. + 2(\gamma - 1)\psi(k + 1) + (\psi(k + 1))^r + \psi'(k + 1) \right] \quad (14)$$

اثبات: با توجه به لم ۳ داریم

$$E_k\left[\frac{X^c(\log X)^\gamma}{(1+X^c)^\gamma}\right] = \frac{k}{k+\gamma} E_{k+1}(X^c(\log X)^\gamma)$$

با بصورت زیر داریم

$$E_k\left[\frac{X^c(\log X)^\gamma}{(1+X^c)^\gamma}\right] = E_k\left(\frac{(\log X)^\gamma}{1+X^c}\right) - E_k\left(\frac{(\log X)^\gamma}{(1+X^c)^\gamma}\right)$$

$$= \frac{k}{k+1} E_{k+1}[(\log X)^\gamma] - \frac{k}{k+\gamma} E_{k+2}[(\log X)^\gamma] \quad (15)$$

با قرار دادن (۱۱) در (۱۵) و با استفاده از نتیجه ۵ و همچنین رابطه بازگشتی $\psi'(k+1) = \psi'(k+1) - \frac{1}{(k+1)^\gamma}$

$$\begin{aligned} E_k\left[\frac{X^c(\log X)^\gamma}{(1+X^c)^\gamma}\right] &= \frac{k}{k+1} E_{k+1}[(\log X)^\gamma] - \frac{k}{k+\gamma} E_{k+2}[(\log X)^\gamma] \\ &= \frac{k}{k+1} \left\{ \frac{\frac{\pi^\gamma}{\gamma} + \gamma^\gamma + 2\gamma\psi(k+1) + (\psi(k+1))^\gamma + \psi'(k+1)}{c^\gamma} \right\} \\ &\quad - \frac{k}{k+\gamma} \left\{ \frac{\frac{\pi^\gamma}{\gamma} + \gamma^\gamma + 2\gamma\psi(k+2) + (\psi(k+2))^\gamma + \psi'(k+2)}{c^\gamma} \right\} \\ &= \frac{k}{c^\gamma(k+1)(k+2)} \left[\frac{\pi^\gamma}{\gamma} + \gamma^\gamma + 2\gamma k \psi(k+1) + 4\gamma \psi(k+1) \right. \\ &\quad + (k+2)(\psi(k+1))^\gamma + (k+2)\psi'(k+1) - 2\gamma k \psi(k+2) - 2\gamma \psi(k+2) \\ &\quad \left. - (k+1)(\psi(k+2))^\gamma - (k+1)\psi'(k+2) \right] \\ &= \frac{k}{c^\gamma(k+1)(k+2)} \left[\frac{\pi^\gamma}{\gamma} + \gamma^\gamma - 2\gamma + 2(\gamma-1)\psi(k+1) + (\psi(k+1))^\gamma + \psi'(k+1) \right] \end{aligned}$$

۵ اطلاع فیشر در حالت دو پارامتری

اکون که امید ریاضی توابعی پایه‌ای از توزیع $B(12, c)$ را محاسبه کرده‌ایم، بسادگی می‌توانیم امید ریاضی مشتقه جزئی دومتابع درستنمایی را به منظور تعیین درایه‌های ماتریس اطلاع فیشر این توزیع، معرفی کنیم.

۱- میزان اطلاع فیشر توزیع $B(12, c)$ نسبت به پارامتر c

از معادله (۵) داریم

$$E\left(\frac{\partial \ell}{\partial c}\right) = -\frac{n}{c} - n(k+1)E\left[\frac{X^c(\log X)}{(1+X^c)^2}\right]$$

با استفاده از (۱۴)

$$\begin{aligned} E\left(\frac{\partial \ell}{\partial c}\right) &= -\frac{n}{c}\left\{1 + \frac{k}{k+2}\left[\frac{\pi^2}{4} + \gamma^2 - 2\gamma\right.\right. \\ &\quad \left.\left.+ 2(\gamma-1)\psi(k+1) + (\psi(k+1))^2 + \psi'(k+1)\right]\right\} \end{aligned}$$

بنا به لم ۶. ب می‌دانیم میزان اطلاع نسبت به پارامتر c عبارتست از

$$-E\left(\frac{\partial \ell}{\partial c}\right)$$

۲- میزان اطلاع فیشر توزیع χ^2 نسبت به پارامتر k

با استفاده از (۴) امید کمیت معلوم $\frac{\partial \ell}{\partial k}$ عبارتست از

$$E\left(\frac{\partial \ell}{\partial k}\right) = E\left(\frac{-n}{k^2}\right) = \frac{-n}{k^2}$$

بنا به لم ۶. ب در مورد میزان اطلاع نسبت به پارامتر k داریم

$$-E\left(\frac{\partial \ell}{\partial k}\right) = \frac{n}{k^2}$$

۳- میزان اطلاع فیشر توزیع χ^2 نسبت به پارامترهای c و k

از معادله (۶) و با استفاده از معادله (۱۲) داریم:

$$E\left(\frac{\partial \ell}{\partial c \partial k}\right) = -nE\left(\frac{X^c \log X}{1+X^c}\right) = -n\left[\frac{1-\gamma-\psi(k)}{c(k+1)}\right]$$

و بنا به لم ۶. پ در مورد میزان اطلاع نسبت به پارامترهای c و k داریم:

$$-E\left(\frac{\partial \ell}{\partial c \partial k}\right) = n\left[\frac{1-\gamma-\psi(k)}{c(k+1)}\right]$$

توجه داریم که علیرغم شکل انعطاف ناپذیر نسبی عناصر ماتریس اطلاع فیشر بدست آمده، به راحتی می‌توان آنها را با استفاده از روش‌های عددی و با به عدد درآوردن ψ و ψ' محاسبه کرد.

۶ - پارامترسازی

چندین روش برای معرفی یک پارامتر مقیاس در معادله (۱) وجود دارد که عادی‌ترین آن تعریف $Y = \theta X$ به ازای $\theta > 0$ می‌باشد. در اینصورت Y دارای تابع توزیع تجمعی $F(\frac{y}{\theta}; c, k)$ و تابع چگالی احتمال $f(\frac{y}{\theta}; c, k)$ می‌باشد. این نحوه پارامترگذاری باعث بوجود آمدن رابطه $x_i = \frac{y_i}{\theta}$ می‌شود. لذا برای داده‌های Y_1, Y_2, \dots, Y_n داریم

$$f(y; c, k) = \frac{1}{\theta} ck \left(\frac{y}{\theta} \right)^{c-1} \left[\left(1 + \frac{y}{\theta} \right)^c \right]^{-(k+1)}$$

در این حالت، لگاریتم درستنمایی ℓ^y عبارتست از

$$\ell^y = n \log ck - nc \log \theta + (c-1) \sum_{i=1}^n \log y_i - (k+1) \sum_{i=1}^n \log \left(1 + \left(\frac{y_i}{\theta} \right)^c \right)$$

مشتقات اول $\frac{\partial \ell^y}{\partial \theta}$ و $\frac{\partial \ell^y}{\partial k}$ را برای حالت دو پارامتری بدست آوردیم و اما در مورد داریم

$$\begin{aligned} \frac{\partial \ell^y}{\partial \theta} &= -\frac{nc}{\theta} - (k+1) \sum_{i=1}^n \frac{c \left(\frac{y_i}{\theta} \right)^{c-1} \left(-\frac{y_i}{\theta^2} \right)}{1 + \left(\frac{y_i}{\theta} \right)^c} \\ &= -\frac{nc}{\theta} - \frac{c(k+1)}{\theta} \sum_{i=1}^n \frac{x_i^c}{1 + x_i^c} \end{aligned} \quad (16)$$

همچنین مشتقات جزئی دوم $\frac{\partial^2 \ell^y}{\partial c \partial k}, \frac{\partial^2 \ell^y}{\partial k^2}, \frac{\partial^2 \ell^y}{\partial c^2}$ در حالت دو پارامتری محاسبه شدند و در مورد بقیه مشتقات دوم روابط زیر برقرار است

$$\frac{\partial^2 \ell^y}{\partial k \partial \theta} = \frac{c}{\theta} \sum_{i=1}^n \frac{x_i^c}{1 + x_i^c} \quad (17)$$

$$\frac{\partial^2 \ell^y}{\partial c \partial \theta} = -\frac{n}{\theta} + \frac{(k+1)}{\theta} \sum_{i=1}^n \frac{x_i^c}{1 + x_i^c} + \frac{c(k+1)}{\theta} \sum_{i=1}^n \frac{x_i^c \log x_i}{(1 + x_i^c)^2} \quad (18)$$

$$\frac{\partial^2 \ell^y}{\partial \theta^2} = \frac{nc}{\theta^2} - \frac{c(k+1)}{\theta^2} \sum_{i=1}^n \frac{x_i^c}{1+x_i^c} - \frac{c^2(k+1)}{\theta^2} \sum_{i=1}^n \frac{x_i^c}{(1+x_i^c)^2} \quad (19)$$

برای بدست آوردن امیدهای مشتقات لگاریتم درستنمایی علاوه بر امیدهای بدست آمده در حالت دوپارامتری، اکنون به امید عبارات زیر نیز احتیاج داریم

$$\frac{X^c}{1+X^c}, \frac{X^c}{(1+X^c)^2}, \frac{X^c \log X}{(1+X^c)^2}$$

این امیدها می‌توانند بر حسب امیدهای موجود و با بکارگیری معادله (۱) بصورت زیر بیان شوند

$$E\left(\frac{X^c}{1+X^c}\right) = E_k\left(\frac{X^c}{1+X^c}\right) = \frac{k}{k+1} E_{k+1}(X^c) \\ = \frac{k}{k+1} \frac{\Gamma(k)}{\Gamma(k+1)} = \frac{1}{k+1} \quad (20)$$

$$E\left[\frac{X^c}{(1+X^c)^2}\right] = E_k\left[\frac{X^c}{(1+X^c)^2}\right] = \frac{k}{k+2} E_{k+2}(X^c) \\ = \frac{k}{(k+1)(k+2)} \quad (21)$$

تذکر: در این رهیافت از $r = c$ در معادله (۸) استفاده شده است. همچنین برای بدست آوردن معادله (۱۹) k را با 1 و برای بدست آوردن معادله (۲۰) k را با 2 جایگزین کردہ‌ایم.
۳. با توجه به لم ۵

$$E\left[\frac{X^c \log X}{(1+X^c)^2}\right] = E_k\left[\frac{X^c \log X}{(1+X^c)^2}\right] = \frac{k}{k+2} E_{k+2}(X^c \log X)$$

حال با استفاده از (۹) وقتی که $r = c$ داریم

$$E\left[\frac{X^c \log X}{(1+X^c)^2}\right] = \frac{k}{k+2} \left[\frac{1-\gamma-\psi(k)}{c(k+1)} \right] \quad (22)$$

این نتیجه همچنین می‌تواند با نوشتن امید بصورت

$$E_k\left[\frac{X^c \log X}{(1+X^c)^2}\right] = E_k\left(\frac{\log X}{1+X^c}\right) - E_k\left(\frac{\log X}{(1+X^c)^2}\right)$$

$$= \frac{k}{k+1} E_{k+1}(\log X) - \frac{k}{k+2} E_{k+2}(\log X)$$

و با استفاده از معادله (۷) و با نوشتن امید بصورت

$$E_k\left[\frac{X^c \log X}{(1+X^c)^2}\right] = \frac{k}{k+1} E_{k+1}\left[\frac{X^c \log X}{1+X^c}\right]$$

و استفاده از معادله (۱۲) بدست آید.

برای بدست آوردن امیدهای مشتقات لگاریتم درستنماهی در حالت سه پارامتری، توجه به این نکته که در بین مشتقات اول، $\frac{\partial \ell^y}{\partial c}$ و $\frac{\partial \ell^y}{\partial k}$ و در بین مشتقات دوم، $\frac{\partial^2 \ell^y}{\partial c \partial k}$ ، $\frac{\partial^2 \ell^y}{\partial k^2}$ از حالت دوپارامتری تبعیت می‌کنند، حائز اهمیت است. اما در مورد امید سایر مشتقات در حالت سهپارامتری به شرح ذیل داریم

۱. از معادله (۱۶)، امید $\frac{\partial \ell^y}{\partial \theta}$ را با استفاده از معادله (۲۰) و با ساده کردن بدست می‌آوریم

$$E\left(\frac{\partial \ell^y}{\partial \theta}\right) = -\frac{nc}{\theta} + \frac{nc(k+1)}{\theta} E\left(\frac{X^c}{1+X^c}\right) = 0$$

۲. از معادله (۱۷) و با استفاده از معادله (۲۰) داریم

$$E\left(\frac{\partial^2 \ell^y}{\partial \theta \partial k}\right) = \frac{nc}{\theta} E\left(\frac{X^c}{1+X^c}\right) = \frac{nc}{\theta(k+1)}$$

۳. از معادله (۱۸) و با استفاده از معادلات (۲۰) و (۲۲) داریم

$$\begin{aligned} E\left(\frac{\partial^2 \ell^y}{\partial c \partial k}\right) &= \frac{n}{\theta} \left\{ -1 + (k+1)E\left(\frac{X^c}{1+X^c}\right) + c(k+1)E\left[\frac{X^c \log X}{(1+X^c)^2}\right] \right\} \\ &= \frac{nk[1-\gamma-\psi(k+1)]}{\theta(k+2)} \end{aligned}$$

۴. از معادله (۱۹) و با استفاده از معادلات (۲۰) و (۲۱) داریم

$$\begin{aligned} E\left(\frac{\partial^2 \ell^y}{\partial \theta^2}\right) &= \frac{nc}{\theta^2} \left\{ 1 - (k+1)E\left(\frac{X^c}{1+X^c}\right) - c(k+1)E\left[\frac{X^c}{(1+X^c)^2}\right] \right\} \\ &= -\frac{nkc^2}{\theta^2(k+2)} \end{aligned}$$

مراجع

احمدی، ج، (۱۳۷۲)، رکوردها در توزیعهای پیوسته، کارشناسی ارشد، دانشگاه فردوسی مشهد.

Ahmadi, J. and Arghami, NR. (2001), *On the Fisher information in record values*. Metrika 53: 195-206

Hofmann, G. and Nagaraja, N. (2002), *Fisher information in Record Data*. to appear in Metrika.

Park, S. (1996), *Fisher information in order statistics*. J. Amer. Stat. As.91, 385-390.

Timr, L. Fry. (1993), *univariate and multivariate Burr distributions*. PJS, A,1-24.

Zheng, G. and Gastwirth, JL. (2000), *Where is the Fisher information in an ordered sample?*. Statistica Sinica 10, 1267-1280.

برآوردگر رگرسیونی عام و کاربرد آن در بررسی هزینه خانوار

حمید بیدرام^۱، محمد صالحی مرزیجرانی^۲

P11128

^۱ دانشکده ریاضی و کامپیوتر خوانسار، وابسته به دانشگاه اصفهان

^۲ دانشکده ریاضی، دانشگاه صنعتی اصفهان

چکیده: در این مقاله، برآوردگر رگرسیونی عام را که برپایه برآوردگر هارویتز- تامپسون^۱ بنای شده است، معرفی می‌کنیم. علاوه‌نماییم در هنگام مطالعه یک بررسی علاوه بر متغیرهای کمکی از متغیرهای مشترکی که در هر دو بررسی وجود دارد جهت افزایش دقت استفاده نماییم. زیچانگ (۱۹۹۰) در بررسی هزینه خانوار ایالات متحده به خوبی این کار را به انجام رسانیده است. در این مقاله ضمن معرفی روش رنسن و نیونبروک (۱۹۹۷) که روش زیچانگ را تا حد زیادی ساده و کاربردی تر نموده است با مقایسه این دو روش یک مدل رگرسیونی ارائه خواهیم داد که همزمان شامل هر دو متغیرهای کمکی و مشترک بوده و خواهیم دید با این روش نحوه اجرای کار به مراتب ساده‌تر خواهد شد.

واژه‌های کلیدی: برآوردگر رگرسیونی، برآوردگر رگرسیونی عام، برآوردگر هارویتز- تامپسون، متغیر کمکی، نمونه‌گیری با احتمالات نابرابر.

۱ مقدمه

برآوردگر رگرسیونی عام از جمله برآوردگرهایی است که پتانسیل عملی بالایی دارد. برآورد هزینه خانوار هم اینک در تمام ملل دنیا صورت می‌گیرد. که در آن نمونه گیری چند مرحله‌ای انجام می‌شود و در هر مرحله اولاً باید وزنهای متفاوتی برای انتخاب زیردامنه‌ها و واحدهای مختلف در نظر گرفت ثانیاً به علت تفاوت سطح فرهنگ در خانوارها روش جمع آوری اطلاعات نیز باید متفاوت باشد. برآوردگر رگرسیونی عام قادر است از متغیرهای مشترک اندازه‌گیری شده در دو یا چند بررسی علاوه بر متغیرهای کمکی استفاده نموده و دقت را افزایش دهد. در این مقاله به کاربرد عملی این برآوردگر در بررسی هزینه خانوار اشاره خواهیم نمود. در اینجا حجم جامعه را با N و اعضای آن را Y_1, Y_2, \dots, Y_N و پارامتر مجموع کل که معمولاً نامعلوم و قابل برآورد است را با $\sum_{i=1}^N Y_i = t$ نمایش می‌دهیم. به دلیل اهمیت برآوردگر هارویتز- تامپسون در برآوردگر رگرسیونی عام این برآوردگر را معرفی می‌کنیم.

۲ برآوردهارویتزر-تامپسون

جامعه‌ای با حجم N با اعضای Y_1, Y_2, \dots, Y_N را در نظر بگیرید. اگر احتمال انتخاب اعضاء را با π_k نمایش دهیم مجموع کل یعنی $t = \sum_{i=1}^N Y_i$ را با

$$\hat{t}_\pi = \sum_{i=1}^n \frac{Y_i}{\pi_i} \cdot I_i \quad (1)$$

برآورد می‌کنیم که آن را برآوردهارویتزر-تامپسون می‌نامند و در آن

$$I_k = \begin{cases} 1 & \text{اگر عضو } k \text{ ام در نمونه باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

واضح است که اگر در بررسی از نمونه‌گیری تصادفی ساده استفاده کنیم یعنی به همه اعضاء احتمال انتخاب یکسانی را تخصیص دهیم به ازای هر k در جامعه، $\pi_k = \frac{n}{N}$ و بنابراین:

$$\begin{aligned} \hat{t}_\pi &= \sum_{i=1}^n \frac{Y_i}{\frac{n}{N}} \cdot I_i \\ &= N \sum_{i=1}^n \frac{Y_i}{n} \cdot I_i \\ &= N \bar{Y}_s \end{aligned} \quad (2)$$

که \bar{Y}_s میانگین اعضاء نمونه است و (2) همان برآورد مجموع کل در نمونه‌گیری تصادفی ساده است. این برآوردهارویتزر-تامپسون می‌باشد که در [1] مفصلابیان شده است.

۳ برآوردهارویتزر-تامپسونی عام

فرض کنیم بردار متغیرهای کمکی برای k -امین عضو به صورت

$$\underline{x}_k = (x_{1k}, x_{2k}, \dots, x_{jk})' \quad , \quad k \in S$$

می‌باشد که S مجموعه اعضای نمونه است. برآوردهارویتزر-تامپسونی عام $t = \sum_{i=1}^N Y_i$ (مجموع کل نامعلوم جامعه) بصورت زیر فرموله می‌شود.

$$\hat{t}_r = \hat{t}_\pi + \sum_{j=1}^J \hat{\beta}_j (t_{x_j} - \widehat{t_{x_j \pi}}) \quad (3)$$

که در آن $\widehat{t}_\pi = \widehat{\sum_{k=1}^N \frac{x_{j_k}}{\pi_k} I_k}$ براوردگر هارویتز-تامپسون t و $t_{x_j} = \sum_{k=1}^N x_{j_k}$ (مجموع کل متغیرهای کمکی) است و همچنین با در نظر گرفتن مدل رگرسیون

$$Y = X'\beta + \epsilon$$

که در آن $(Y_1, Y_2, \dots, Y_N)' = (Y_1, Y_2, \dots, Y_N)$ یک بردار $N \times 1$ ، X یک ماتریس $J \times N$ از متغیرهای کمکی، β یک بردار $J \times 1$ با صورت $(\beta_1, \beta_2, \dots, \beta_J)$ و ϵ یک بردار $N \times 1$ است. با روش کمترین مربعات خطا بصورت زیر براورد می‌شود:

$$\hat{\beta} = (X'X)^{-1}X'\Sigma^{-1}Y \quad (4)$$

که در آن

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \sigma_N^2 \end{bmatrix}$$

ماتریس واریانس-کوواریانس Y است.

۴ استفاده از براوردگر رگرسیونی عام بعنوان یک تغییل کننده وزنی در بررسی هزینه خانوار

نمونه‌ای شامل n خانوار را در نظر بگیرید. بطوريکه بردار $(s_1, s_2, \dots, s_n)'$ که در آن $s_k = \pi_k^{-1}$ و π_k احتمال انتخاب خانوار k -ام می‌باشد، معلوم و مشخص باشد. به بردار S اصطلاحاً بردار وزنهای آغازین (ابتداي) می‌گويند. همچنین فرض کنیم J حجره طبقه‌بندی پسین وجود داشته باشد بطوريکه تعداد افراد جامعه (N_j) برای هر حجره معلوم باشد. برای مثال در بررسی هزینه مصرف کننده ایالات متحده تعداد $J = 48$ حجره مطابق با ترکیبات ۲ جنس، ۲ نژاد (سیاه و غیرسیاه) و ۱۲ رده سنی وجود دارد.

در واقع بردار $(N_1, N_2, \dots, N_j)'$ بردار متغیرهای کمکی است و N_j -ها تعداد افراد در حجره j -ام را نشان می‌دهد. ترکیب خانوارهای نمونه توسط ماتریس $X = (x_{kj})$ مشخص می‌گردد که x_{kj} تعداد اشخاص نمونه در j -امین حجره از خانوار k -ام می‌باشد. بطوريکه ماتریس X را بصورت زیر نمایش

می‌دهیم:

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1J} \\ x_{21} & x_{22} & \cdots & x_{2J} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nJ} \end{bmatrix}_{n \times J}$$

واضح است که اگر برداری از ماتریس فوق بصورت $(\circ, 1, 0, \dots, 0)$ باشد بیانگر آن است که خانوار مذکور شامل ۲ نفر در حجره اول و یک نفر در حجره دوم خواهد بود. با استفاده از وزنهای آغازین می‌توان تعداد افراد در حجره j -ام را از رابطه

$$\widehat{N}_j = \Sigma_k x_{kj} s_k$$

با بطور کلی از رابطه

$$\widehat{\underline{N}} = X' \underline{S} \quad (5)$$

برآورد کرد.

مسلماً $\widehat{\underline{N}} \neq \underline{N}$. بنابراین دنبال بردار تعديل کننده وزنی هستیم تا بواسطه آن $\widehat{\underline{N}} = \underline{N}$. زیچانگ (۱۹۹۰ و ۱۹۸۶) به روش (General Least Square) GLS بروز آورد بطوریکه $\underline{N} = X' \underline{W}$ بددست $(w_1, w_2, \dots, w_n)'$. روش به اینصورت است که با در نظر گرفتن قید $\underline{N} = X' \underline{W}$ معیار

$$\Sigma_k (w_k - s_k) / s_k = (\underline{W} - \underline{S})' A^{-1} (\underline{W} - \underline{S}) \quad (6)$$

را که در آن $A_{n \times n} = diag(\underline{S})$ ، کمینه می‌کنند بنابراین بردار \underline{W} از کمینه کردن معیار (۶)، بصورت زیر بددست می‌آید:

$$\underline{W} = \underline{S} + A \underline{X} (\underline{X}' A \underline{X})^{-1} (\underline{N} - \underline{X}' \underline{S}) \quad (7)$$

حال اگر فرض کنیم Y یک بردار $1 \times n$ از صفت مورد مطالعه در نمونه باشد با داشتن بردار \underline{W} بسادگی برآوردگر رگرسیونی عالم مجموع کل صفت مورد مطالعه یعنی $t = \Sigma_k Y_k$ بصورت زیر بددست می‌آید:

$$\begin{aligned} \widehat{t}_r &= \underline{Y}' \underline{W} \\ &= \underline{Y}' (\underline{S} + A \underline{X} (\underline{X}' A \underline{X})^{-1} (\underline{N} - \underline{X}' \underline{S})) \\ &= \underline{Y}' \underline{S} + \underline{Y}' A \underline{X} (\underline{X}' A \underline{X})^{-1} (\underline{N} - \underline{X}' \underline{S}) \end{aligned} \quad (8)$$

که با مقایسه با رابطه (۳) برای برآوردگر رگرسیونی عام ملاحظه می‌کنیم:

$$\begin{aligned} \underline{N} &= \underline{X}' \underline{W} = t_x \\ \widehat{\beta}' &= \underline{Y}' \Lambda^{-1} \underline{X} (\underline{X}' \Lambda \underline{X})^{-1} \\ \widehat{t}_\pi &= \underline{Y}' \underline{S} \\ \underline{X}' \underline{S} &= \widehat{t}_{x_\pi} \end{aligned}$$

بنابراین بنظر می‌رسد با داشتن وزنهای \underline{W} براحتی با ضرب بردار نمونه \underline{Y} در بردار وزنهای \underline{W} ، برآوردگر رگرسیونی عام بدست آید.

۵ وجود متغیرهای مشترک دو برسی

زیچانگ (۱۹۸۶ و ۱۹۹۰) در بررسی هزینه خانوار جهت افزایش دقت جمع آوری اطلاعات نمونه را به دو طریق اعمال کرد. یک نمونه را با روش مصاحبه‌ای (پرسش و پاسخ بین مامور و خانوار) و نمونه دیگر را با روش یادداشتی (ارسال پرسشنامه به خانوارها و تکمیل آنها توسط خانوارها) انجام داد. ایشان به این نتیجه مهم دست یافت که جهت افزایش دقت برآورد بهتر است علاوه بر استفاده از متغیرهای کمکی دو برسی از متغیرهای مشترکی (مانند ناحیه نمونه‌گیری، مالکیت، نوع خانوار...) که در هر دو برسی وجود دارد استفاده کرد. انتظار داریم که برآورد مجموع کل توسط هر دو طرح یکسان باشد. زیچانگ با مشخص کردن برسی اول و دوم بترتیب با اندیسهای ۱ و ۲ رابطه (۸) را با بکار بردن نمادهای زیر، تعمیم داد.

$$\begin{aligned} \underline{S} &= \begin{bmatrix} S_1 \\ S_2 \end{bmatrix}, \quad \Lambda = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \\ N_X &= \begin{bmatrix} N_{X^0} \\ N_{X^0} \\ \vdots \end{bmatrix} \\ X &= \begin{bmatrix} X_1^0 & \circ & X_1^c \\ \circ & X_2^0 & -X_2^c \end{bmatrix} \end{aligned} \tag{۹}$$

که در آنها:

S_1 : بردار وزنهای آغازین برسی اول.

S_2 : بردار وزنهای آغازین برسی دوم.

Λ_{12} و Λ_{22} ماتریسهای قطری بصورت $\Lambda_{11} = diag(S_1)$ و $\Lambda_{22} = diag(S_2)$

$$\Lambda_{21} = 0$$

X_1^c و X_2^c : ماتریسهای کمکی بررسی اول و دوم.

X_3^c و X_4^c : ماتریسهای مشترک بررسی اول و دوم.

N_X^c : مجموع عناصر ستونهای ماتریس X در جامعه است.

بنابراین بردار \underline{W} در رابطه (۷) با جایگذاری ماتریسها و بردارهای (۹) محاسبه می‌شود. مثال کاربردی از این روش در [۳] بیان شده است. همانطور که ملاحظه می‌کنیم استفاده از متغیرهای مشترک به فرم ماتریس مشکل و بعضاً با اشتباه همراه است. روشی را در زیر جهت استفاده بهتر متغیرهای مشترک بیان می‌کنیم.

۶ استفاده از متغیرهای مشترک بعنوان متغیرهای رگرسور اضافی

فرض کنید دو بررسی در نمونه دارای متغیرهای مشترکی باشند اگر مجموع کل این متغیرها معلوم باشند آنگاه متغیرهای مذکور را بعنوان متغیرهای کمکی در برآورد رگرسیونی هر کدام از بررسیها در نظر می‌گیریم ولی اگر مجموع این متغیرها نامعلوم باشد ممکن است برآورد مجموع کل آنها را همراه با متغیرهای کمکی به عنوان متغیرهای توضیحی اضافی در برآورد رگرسیونی عام هر دو بررسی بکار ببریم. با مشخص شدن وزنهای برآورد بوجود آمده، مجموع کل متغیرهای کمکی را دوباره برآورد می‌کنیم. در اینجا برآورد رگرسیونی جدیدی بدست می‌آوریم که به آن برآورد رگرسیونی تعديل شده می‌گوییم و براساس دو نمونه S_1 و S_2 به معنی آن خواهیم پرداخت. در واقع با معروفی این روش وجود متغیرهای مشترک را بعنوان یک جمله اضافی در مدل رگرسیون بکار می‌بریم.

فرض کنید در اینجا برای هر عضو u_k یک بردار n -تایی \underline{z}_k از متغیرهای مشترک داشته باشیم و

$t_z = Z' l$ مجموع کل متغیرهای مشترک در جامعه باشد که در آن

$$l_{N \times 1} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \quad Z_{N \times q} = \begin{bmatrix} z'_1 \\ z'_2 \\ \vdots \\ z'_N \end{bmatrix}$$

برآوردگر هارویتز-تمپسون برای t_z را بصورت $t_{z\pi} = Z'_s A_s^{-1} l_s$ نشان می‌دهیم که در آن Z_s مقادیر مشاهده شده متغیرهای مشترک نمونه است و A_s و l_s قبل از معرفی شده‌اند. رنسن و نیونبروک (۱۹۹۷) برآوردگر رگرسیونی تعديل شده را با توجه به نمادهای فوق بصورت زیر معرفی کردند:

$$\widehat{t_{Ar}} = \widehat{t_\pi} + \widehat{B'_A}(t_x - \widehat{t_{x\pi}}) + \widehat{D'_A}(\widehat{t_z} - \widehat{t_{z\pi}}) \quad (10)$$

که $\widehat{D'_A}$ و $\widehat{B'_A}$ از رابطه زیر بدست می‌آید.

$$\begin{pmatrix} \widehat{B_A} \\ \widehat{D_A} \end{pmatrix} = ((X_s Z_s)' A_s (X_s Z_s))^{-1} (X_s Z_s)' A_s Y_s$$

و $\widehat{t_z}$ یک بردار q -تایی از برآوردهای t_z است. در واقع رابطه (10) را با توجه به تعریف برآوردگر رگرسیونی عام به صورت زیر نیز می‌توان نوشت:

$$\begin{pmatrix} t_x - \widehat{t_{x\pi}} \\ \widehat{t_z} - \widehat{t_{z\pi}} \end{pmatrix} \widehat{t_{Ar}} = \widehat{t_\pi} + \begin{pmatrix} \widehat{B_A} \\ \widehat{D_A} \end{pmatrix} \quad (11)$$

حال به خاصیت اساسی برآوردگر رگرسیونی تعديل شده می‌پردازیم. در واقع این خاصیت تکمیل کننده روش رنس و نیونبروک می‌باشد. فرض کنید:

$$\begin{aligned} \widehat{t_r} &= \widehat{t_\pi} + \widehat{\beta}'(t_x - \widehat{t_{x\pi}}) \\ \widehat{\beta}' &= (X_s' A_s X_s)^{-1} X_s' A_s Y_s' \\ \widehat{t_{zr}} &= \widehat{t_{z\pi}} + \widehat{L}'(t_x - \widehat{t_{x\pi}}) \\ \widehat{L} &= (X_s' A_s X_s)^{-1} X_s' A_s Z_s \end{aligned}$$

برآوردهای رگرسیونی عام مجموع کل t_Y و t_Z باشند. فرض می‌کنیم X_s پر رتبه ۱ باشد اگر $J \geq 1$ باشد برآوردگر رگرسیونی تعديل شده (10) یعنی $\widehat{t_{Ar}}$ بصورت زیر نوشه می‌شود:

$$\widehat{t_{Ar}} = \widehat{t_r} + \widehat{D'_A}(\widehat{t_z} - \widehat{t_{zr}}) \quad (12)$$

که در آن $\widehat{t_{zr}}$ برآوردگر رگرسیونی عام متغیرهای مشترک بررسی اول است و ضریب رگرسیون جرئی $\widehat{D_A}$ بصورت زیر نوشه می‌شود:

$$\widehat{D_A} = (Z_s' R_s Z_s)^{-1} Z_s' R_s Y_s$$

$$R_s = A_s - A_s X_s (X_s' A_s X_s)^{-1} X_s' A_s$$

اثبات (۱۲) در [۴] آمده است.

از (۱۲) ملاحظه می‌کنیم که برآورده رگرسیون تعديل شده برابر است با برآورده رگرسیونی عام باضافه یک جمله تعديل کننده. این جمله تعديل کننده یک عامل اصلاح کننده برآورده رگرسیونی عام و در عین حال عاملی برای ایجاد سارگاری بین برآورد متغیرهای مشترک از دو بررسی می‌باشد. برآورده رگرسیونی تعديل شد. در (۱۲) بردارهای وزنی n -تایی زیر را نتیجه می‌دهد.

$$W_{As} = W_s + R_s Z_s (Z_s' R_s Z_s)^{-1} (\hat{t}_z - \hat{t}_{zr}) \quad (13)$$

$$W_s = \underline{S} + A_s X_s (X_s' A_s X_s)^{-1} (t_x - \hat{t}_{x\pi}) \quad \text{که در آن:}$$

یک بردار n -تایی از وزنهایی است که مطابق با برآورد رگرسیون عام \hat{t}_r بدست می‌آمد. دقت کنیم در اینجا $\circ = X_s' R_s$ و بنابراین

$$\begin{aligned} X_s' W_{As} &= X_s' W_s = t_x \\ Z_s' W_{As} &= \hat{t}_z \end{aligned}$$

یعنی بردار وزنی W_{As} —ای پیدا کردیم که علاوه بر اینکه با حاصلضرب آن در X_s' مجموع کل متغیرهای کمکی در جامعه را می‌دهد، با حاصلضرب در Z_s' (که از نمونه مشاهده می‌شود) مجموع کل متغیرهای مشترک در جامعه را برآورد می‌کند. (چون \hat{t}_z برآورد مجموع کل متغیرهای مشترک در جامعه بود). مثال کلی از هر دو روش زیچانگ و روش مطرح شده فوق در [۴] آمده است.

مراجع

- C.E.Sandral ,B.Swenson and J.Wretman, Model Assisted Survey sampling,
New York :Springer-Verlag(1990)
K.D. Zieschang,"Generalized Least Squares: An Alternative to Principal
person weighting" Washington, DC. Bureau of Labor Statistics (1986).

K.D. Zieschang,"Sample Weighting Methods and Estimation of totals in the consumer Expenditure Surveys", Journal of the American Statistical Association, No:85, PP. 986-1001,(1990).

حمید بیدرام، : برآوردگر رگرسیونی عام و کاربرد آن در بررسی هزینه خانوار: پایان نامه کارشناسی ارشد، (دانشگاه صنعتی اصفهان - ۱۳۷۹).

برآورد پارامترهای مدل رگرسیونی با رهیافت ماکسیمم آتروپی

عین‌اله پاشا^۱، محسن محمدزاده^۲، علی آقامحمدی^۲

^۱ گروه آمار، دانشگاه تربیت معلم

^۲ گروه آمار، دانشگاه تربیت مدرس

چکیده: در استفاده از مدل رگرسیون که در اکثر مطالعات اقتصادی و اجتماعی به کار می‌رود، معمولاً با دو مشکل مواجه می‌شویم. نخست اینکه در مسئله تعداد پارامترهای مجهول که باید برآورد شوند بیش از تعداد مشاهدات است و دوم فقدان صحّت شرایط لازم جهت کاربرد مدل. این مشکلات اغلب هنگامی که فرایند جمع‌آوری داده‌ها خوب طراحی نمی‌شوند، اتفاق می‌افتد. یک مسئله در این حالت همخطی داده‌ها است. یعنی دو یا چند متغیر همبستگی خیلی زیادی با هم دارند. البته همخطی کامل به ندرت رخ می‌دهد، ولی همخطی غیر کامل نیز منجر به افزایش واریانس در برآورد پارامترها به روش LS شده، و این مسئله باعث افت توان آزمون در آزمون فرض پارامترهای مدل می‌شود.

اخيراً رهیافتهاي متعددی برای مواجه با اين گونه مسائل پیشنهاد شده که يکی از آنها رهیافت ماکسیمم آتروپی تعمیم یافته است که در سال ۱۹۹۳ توسط گالن و جودجی ارائه شده است. در این مقاله ضمن مطالعه این روش کارائی آن را با یک شبیه سازی بررسی می کنیم.

۱ مقدمه

هدف از علم آمار حصول اطلاعات از داده‌های نمونه‌ای می‌باشد. در این راستا یکی از پر کاربردترین ابزارهای آماری رگرسیون می‌باشد، که هدف از آن تبیین تغییرات اندازه‌های آزمایشی یک متغیر از روی متغیرهای وابسته دیگری که مقادیر آنها در جریان آزمایش تغییر می‌کند، می‌باشد. با شرکت دادن صریح داده‌های این متغیرهای نافذ در تجزیه و تحلیل آماری، اغلب امکان دارد که از طبیعت رابطه بین این متغیرها مطلع شده، آنگاه از این اطلاع برای توصیف و استنباطهای مربوط به متغیر اولیه مورد علاقه بهره‌گیری کرد. این واقعیت در درجه اول اهمیّت است که بدانیم تجزیه و تحلیل رگرسیون، با برآش یک

مدل بوسیله کوچکترین توانهای دوم خط، به دست آوردن فاصله های اطمینان و بالاخره با آزمون فرضهای متعدد کامل نمی شود. این گامها فقط بیان نیمی از این مطلب است که استنباطهای آماری را وقی می توان انجام داد که مدل مفروض، مناسب باشد. در پیشتر مطالعات مربوط به علوم اجتماعی و علوم اقتصادی، روابط بین متغیرها شکل تجربی دارند، به نحوی که هیچ وقت نمی توانیم مطمئن باشیم که مدل خاصی صحیح است. در اکثر اوقات فرضها و شرایط لازم برای کاربرد رگرسیون موجود نیست. این موضوع به نوع جمع آوری داده ها و... بر می گردد. اگرچه روشهای نمونه گیری نیز بسط و توسعه یافته اند، اما بیشتر داده ها که توسط آماردانان به کار گرفته می شوند، محدود به خطاهای نمونه ای و غیر نمونه ای نظری خود جامعه، تبدیل داده ها، کدبندی و استفاده از تقریب می باشد. استفاده از مدل رگرسیون که در اکثر مطالعات اقتصادی و اجتماعی به کار می رود، معمولاً با دو مشکل مواجه هست. نخست اینکه در مسئله، تعداد پارامترهای مجھول که باید برآورد شوند بیش از تعداد مشاهدات است و دوم اینکه پارامترها در مسئله ثابت نیستند. این مشکلات اغلب هنگامی که فرایند جمع آوری داده ها خوب طراحی نمی شوند، اتفاق می افتد. یک مسئله در این حالت همخطی داده ها است. یعنی دو یا چند متغیر همبستگی خیلی زیادی با هم دارند. در این حالت ماتریس طرح ستونی رتبه کامل نبوده، در نتیجه معکوس ماتریس، حاصل ضرب ماتریس طرح در تراشه آن، موجود نیست. البته باید توجه داشت که همخطی کامل بندرت رخ می دهد، ولی همخطی غیر کامل نیز منجر به افزایش واریانس در برآورد پارامترهای *LS* شده، و این مسئله باعث افت توان آزمون در آزمون فرض پارامترهای مدل می شود.

اخیراً رهیافتهای متفاوتی برای مواجه با این گونه مسائل به وجود آمده که یکی از آنها رهیافت ماکسیمم آنتروپی تعمیم یافته است. اگر چه ایده اصلی ماکسیمم آنتروپی به سالهای ۱۹۵۷ بر می گردد، ولی مشکلاتی در کاربرد این مفهوم به شکل اصلی آن برای برآورد پارامترهای یک مدل آماری موجود بوده است. مثلاً کمیتهای مجھول باید خاصیت توزیع احتمال داشته باشند، ولی در اکثر اوقات پارامترهای مدلها این گونه نیستند. اگر برای مثال هدف، برآورد پارامترهای مدل رگرسیونی $\underline{Y} = \underline{X}\beta + \underline{\epsilon}$ باشد، از آنجاییکه β بردار ثابتی است، توزیع احتمال برای آن بی معنی است. اصل ماکسیمم آنتروپی طوری تعمیم داده شده که برای این منظور نیز بتوان از آن استفاده کرد.

۲ ماکسیمم آنتروپی تعمیم یافته

با در نظر گرفتن محدودیتهای فوق جودجی و گلن (۱۹۹۴) ماکسیمم آنتروپی را به وسیله تعریف پارامترهای مجھول بر حسب توزیعهای احتمال گستته تعمیم دادند. برای

توصیف کاربردهای ماکسیمم آنتروپی تعمیم یافته^۱ (GME) مدل رگرسیونی استاندارد زیر را در نظر بگیرید.

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{E} \quad (1)$$

که در آن \underline{Y} ، یک بردار $(T \times 1)$ بعدی، \underline{X} ماتریس طرح $(T \times K)$ بعدی، $\underline{\beta}$ یک بردار $(K \times 1)$ بعدی از پارامترهای ناشناخته و \underline{E} نیز یک بردار $(T \times 1)$ بعدی از خطاهای می‌باشد. جودجی و گلن (۱۹۹۴) فرض کردند که $\underline{\beta}_k$ ، مؤلفه‌های بردار $\underline{\beta}$ متغیرهای تصادفی گسته‌ای باشند، که M مقدار، $\infty < M \leq \infty$ را از یک بازه بسته با احتمالهای متفاوت اختیار می‌کنند. اگر z_{kl} و z_{km} به ترتیب کران پایین و بالای این بازه باشند، آنگاه β_k را می‌توان به عنوان یک ترکیب خطی محدب از کرانها نوشت. به طور مشابه، مؤلفه‌های بردار \underline{E} را هم می‌توان به عنوان یک متغیر تصادفی گسته که J مقدار، $\infty < J \leq \infty$ را از یک بازه بسته با احتمالهای متفاوت اختیار می‌کند، در نظر گرفت. در این حالت رابطه (1) را می‌توان به صورت زیر بازنویسی کرد.

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{E} = \underline{X}\underline{Z}\underline{P} + \underline{V}\underline{W}$$

که در آن $\underline{E} = \underline{Z}\underline{P}$ و $\underline{Z}\underline{P} = \underline{V}\underline{W}$ ، به طوری که Z یک ماتریس $(K \times KM)$ بعدی، از مقادیری که $\underline{\beta}$ اختیار می‌کند و \underline{P} یک بردار $(KM \times 1)$ بعدی، از احتمالهای ناشناخته است، به طوری که

$$\underline{Z}\underline{P} = \begin{pmatrix} \underline{Z}'_1 & \circ & \cdots & \circ \\ \circ & \underline{Z}'_2 & \cdots & \circ \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \cdots & \underline{Z}'_k \end{pmatrix} \begin{pmatrix} \underline{P}'_1 \\ \underline{P}'_2 \\ \vdots \\ \underline{P}'_k \end{pmatrix}$$

$$\begin{aligned} \underline{Z}'_1 &= (z_{k1}, z_{k2}, \dots, z_{km}), & \underline{P}'_k &= (p_{k1}, p_{k2}, \dots, p_{km}) \\ \underline{Z}'_k \underline{P}_k &= \sum_{m=1}^M z_{km} p_{km} = \beta_k, & k &= 1, 2, \dots, K \end{aligned}$$

$$\begin{aligned} p_{km} &> 0, & k &= 1, 2, \dots, K, & m &= 1, 2, \dots, M \\ \sum_{m=1}^M p_{km} &= 1, & k &= 1, 2, \dots, K \end{aligned}$$

^۱ Generalized Maximum Entropy

به طور مشابه V یک ماتریس $(T \times TJ)$ بعدی، از مقادیری که \underline{E} اختیار می‌کند و \underline{W} نیز یک بردار $(1 \times TJ)$ بعدی، از وزنهای احتمالی است، به طوری که

$$V\underline{W} = \begin{pmatrix} \underline{V}_1' & \circ & \cdots & \circ \\ \circ & \underline{V}_2' & \cdots & \circ \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \cdots & \underline{V}_T' \end{pmatrix} \begin{pmatrix} \frac{\underline{W}_1}{\underline{W}_2} \\ \vdots \\ \vdots \\ \frac{\underline{W}_T}{\underline{W}_T} \end{pmatrix}$$

$$\underline{V}'_t = (v_{t1}, v_{t2}, \dots, v_{tJ}), \quad \underline{W}'_t = (w_{t1}, w_{t2}, \dots, w_{tJ}), \quad t = 1, 2, \dots, T$$

$$\underline{V}'_t \underline{W}'_t = \sum_{j=1}^J v_{tj} w_{tj} = e_t, \quad t = 1, 2, \dots, T$$

$$w_{tj} > 0, \quad t = 1, 2, \dots, T, \quad j = 1, 2, \dots, J$$

$$\sum_{j=1}^J w_{tj} = 1, \quad t = 1, 2, \dots, T$$

با فرضیات بالا می‌توان از مدل ماکسیمم آنتروپی جهت برآورد p_k ها و w_t ها به شکل زیر استفاده کرد.

$$MAX H(\underline{P}, \underline{W}) = - \sum_{k=1}^K \sum_{m=1}^M p_{km} \ln p_{km} - \sum_{t=1}^T \sum_{j=1}^J w_{tj} \ln w_{tj} \quad (2)$$

تحت شرایط

$$\sum_{m=1}^M p_{km} = 1, \quad k = 1, 2, \dots, K \quad (3)$$

$$\sum_{j=1}^J w_{tj} = 1, \quad t = 1, 2, \dots, T \quad (4)$$

$$\underline{Y} = X\underline{\beta} + \underline{E} = XZ\underline{P} + V\underline{W} \quad (5)$$

معادله رابطه (2) تابع هدف و رابطه (5) قید مدلی و روابط (3) و (4) قیدهای نرمال ساز هستند. در چنین ساختاری معادله هدف با مجموع آنتروپی شanon روی پارامترها و خطاهای متناسب است. باید توجه داشت که برآورد $\underline{\beta}$ طوری به دست می‌آید که شامل اطلاعات داده‌ها، مدل و مجموعه قیود باشد. کرانهای پایین و بالای بازه‌های شامل پارامترها

بسته به نوع مسئله می‌توانند متقارن باشند. همچنین هنگامی که توزیع خطاهای نرمال باشد مقادیر ماتریس V را می‌توان از یک توزیع یکتواخت گستته که حول صفر متقارن است، در نظر گرفت. محدوده بازه مربوط به خطای را می‌توان از روی مقادیر y_t ها نیز تعیین کرد. پالکشیم (۱۹۹۴) نشان داد، هنگامی که مقادیر V ها بصورت $\{-\alpha, 0, \alpha\}$ یعنی $J = 3$ ، در نظر گرفته شوند، آنگاه α را می‌توان به عنوان واریانس نمونه‌ای y_t ها در نظر گرفت. در تعیین بازه‌های مربوط به مقادیر Z نیز بازه باید طوری انتخاب شود که مجموعه شدنی، یعنی مجموعه شامل P ها و W هایی که در شرایط (۳) (۵) صادق هستند، تهی نباشد تا در داخل این مجموعه شدنی بتوان آنتروپی را بر حسب P و W ماکسیمم کرد. همچنین تعیین تعداد نقاط یعنی مقادیر M و J نیز مشکلاتی ایجاد می‌کنند اگر چه جودجی و گلن پیشنهاد می‌کنند که در یک بازه مقادیر را به فاصله مساوی از یکدیگر انتخاب کنیم ولی کیتامورا و همکارانش (۱۹۹۸) روشی پیشنهاد می‌کنند که در آن توزیع احتمال‌ها به صورت پیوسته در نظر گرفته شده‌اند.

تابع لاغرانژ را می‌توان به صورت زیر تشکیل داد.

$$\begin{aligned} L = & - \sum_{k=1}^K \sum_{m=1}^M p_{km} \ln p_{km} - \sum_{t=1}^T \sum_{j=1}^J w_{tj} \ln w_{tj} - \sum_{k=1}^K \alpha_k \left(\sum_{m=1}^M p_{km} - 1 \right) \\ & - \sum_{t=1}^T \gamma_t \left(\sum_{j=1}^J w_{tj} - 1 \right) - (\lambda_1, \lambda_2, \dots, \lambda_T) (Y - XZP - VW). \end{aligned}$$

حال از L نسبت به p_{km} و w_{tj} مشتق گرفته، مساوی صفر قرار می‌دهیم

$$\frac{\partial L}{\partial p_{km}} = -1 - \ln p_{km} - \alpha_k - z_{km} \left(\sum_{t=1}^T \lambda_t x_{tk} \right) = 0$$

اگر قرار دهیم $X_k = (x_{1k}, x_{2k}, \dots, x_{Tk})$ و $\underline{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_T)$ از طرفی چون $\sum_{m=1}^M p_{km} = 1$ خواهیم داشت:

$$p_{km}^* = \frac{\exp(z_{km} \underline{\lambda}' X_k)}{\sum_{m=1}^M \exp(z_{km} \underline{\lambda}' X_k)}, \quad k = 1, 2, \dots, K, m = 1, 2, \dots, M$$

و همچنین

$$\frac{\partial L}{\partial w_{tj}} = -1 - \ln w_{tj} - \gamma_t - v_{tj} = 0$$

بنابراین

$$w_{tj} = \exp(1 + \gamma_t + \lambda_t v_{tj})$$

از طرفی چون $\sum_{j=1}^J w_{tj} = 1$ پس

$$w_{tj}^* = \frac{\exp(\lambda_t v_{tj})}{\sum_{j=1}^J \exp(\lambda_t v_{tj})} \quad t = 1, 2, \dots, T, \quad j = 1, 2, \dots, J$$

در نتیجه مقادیر e_t و β_k را می‌توان به صورت زیر برآورد کرد:

$$\begin{aligned} \beta_k^* &= \sum_{m=1}^M z_{kmp} p_{km}^*, & k &= 1, 2, \dots, K \\ e_t^* &= \sum_{j=1}^J v_{tj} w_{tj}^*, & t &= 1, 2, \dots, T \end{aligned}$$

البته در عمل برای برآورد پارامترهای مدل از روش تحلیلی به دلیل پیچیده بودن معادلات نمی‌توان استفاده کرد. به همین خاطر از روش‌های بهینه سازی که برای کمینه و بیشینه کردن توابع هدف تحت قیود مشخص بکار می‌روند، استفاده می‌شود. بیشتر عملیات بهینه سازی در این گونه مسائل با استفاده از نرم‌افزار GAMS انجام می‌شود.

۳ نتایج یک شبیه سازی

برای بررسی کارائی روش گفته شده پارامترهای مدل رگرسیونی زیر را با استفاده از روش ماکسیمم آنتروپی برآورد می‌کنیم.

فرض کنید که یک مدل رگرسیونی به صورت

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + E_i \quad (1)$$

را می‌خواهیم به داده‌ها بیزارانیم. برای تولید داده‌ها فرض می‌کنیم

$$\underline{\beta} = (\beta_0, \beta_1, \beta_2) = (0/2, 0/3, 0/6)$$

باشد. مقادیر E_i را از توزیع نرمال با میانگین صفر و واریانس $9/00$ تولید می‌کنیم. مقادیر x_{i1} و x_{i2} که به ترتیب، $x_{i1}:x_{i2}$: تعداد گرافیتهاي شبکه‌ای موجود در هر میلیمتر مربع از سطح قالب آهن و $x_{i2}:$ درصد ترکیب آهن و کربن موجود در قالب آهن است، را از کتاب باتاچاریا جلد دوم صفحه ۴۳۲ در نظر می‌گیریم. حال برای تولید مقادیر y_t ، ابتدا e_t را از توزیع نرمال به طور تصادفی انتخاب کرده و از طریق معادله (۱) مقادیر y_t را به ازای x_{i1} و x_{i2} و $\underline{\beta}$ مشخص شده، به دست می‌آوریم.

ابتدا ۱۰۰ نمونه به حجم‌های ۶، ۱۰ و ۲۰ تولید کرده و برای هر نمونه مقادیر β را از طریق ماقسیم آتروبی به روش زیربرآورد می‌کنیم

$$\begin{aligned} \underline{Z}_1 &= (\circ, \circ/\Delta, 1)', & \underline{P}_1 &= (p_{11}, p_{12}, p_{13})', & \beta_1 &= \sum_{d=1}^3 z_{1d} p_{1d}^*, \\ \underline{Z}_2 &= (\circ, \circ/\Delta, 1)', & \underline{P}_2 &= (p_{21}, p_{22}, p_{23})', & \beta_2 &= \sum_{d=1}^3 z_{2d} p_{2d}^*, \\ \underline{Z}_3 &= (\circ, \circ/\Delta, 1)', & \underline{P}_3 &= (p_{31}, p_{32}, p_{33}), & \beta_3 &= \sum_{d=1}^3 z_{3d} p_{3d}^* \end{aligned}$$

مقادیر V_t را برای کلیه $t = 1, 2, \dots, T$ حجم نمونه است، برابر با $V_t = (w_{t1}, w_{t2}, w_{t3})'$ در نظر گرفته، در نتیجه $e_t = (w_{t1}, w_{t2}, w_{t3})' - (\circ/\Delta, \circ, \circ)$ برابر است:
با:

$$e_t = \sum_{j=1}^3 v_{tj} w_{tj}^*, \quad t = 1, 2, \dots, T$$

حال برای یک نمونه T تابع هدف به صورت زیر است،

$$-\sum_{k=1}^3 \sum_{d=1}^3 p_{kd} \ln p_{kd} - \sum_{t=1}^T \sum_{j=1}^3 w_{tj} \ln w_{tj}$$

که باید تحت قیود

$$\sum_{j=1}^3 w_{tj} = 1, \quad t = 1, 2, \dots, T$$

$$\sum_{d=1}^3 p_{kd} = 1, \quad k = 1, 2, 3$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{T1} & x_{T2} \end{pmatrix} \begin{pmatrix} \underline{Z}'_1 & \circ & \circ \\ \circ & \underline{Z}'_2 & \circ \\ \circ & \circ & \underline{Z}'_3 \end{pmatrix} \begin{pmatrix} \underline{P}_1 \\ \underline{P}_2 \\ \underline{P}_3 \end{pmatrix}$$

$$+ \begin{pmatrix} V'_1 & \circ & \cdots & \circ \\ \circ & V'_2 & \cdots & \circ \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \cdots & V'_T \end{pmatrix} \begin{pmatrix} W_1 \\ W_2 \\ \vdots \\ W_T \end{pmatrix}$$

ماکسیمم شود، تا مقادیر بردارهای p_k و w_t و $k = 1, 2, \dots, T$ مشخص شوند. در نتیجه مقادیر β_k از روابط (۷) به دست می‌آیند.

مقادیر میانگین ضرایب رگرسیونی β_1, β_2 و β_3 روی ۱۰۰ نمونه و همچنین واریانس‌های آنها طبق جدول ۱ می‌باشد.

میانگین	واریانس	$N = 7$	$N = 10$	$N = 20$			
$\sum_{n=1}^{100} \frac{\beta_1^* n}{100}$	$Var(\beta_1^*)$	۰/۴۲	۰/۰۰۳	۰/۲۵	۰/۰۰۹	۰/۲۸	۰/۰۱۰
$\sum_{n=1}^{100} \frac{\beta_2^* n}{100}$	$Var(\beta_2^*)$	۰/۳۵	۰/۰۰۴	۰/۲۸	۰/۰۰۵	۰/۳۷	۰/۰۴۰
$\sum_{n=1}^{100} \frac{\beta_3^* n}{100}$	$Var(\beta_3^*)$	۰/۵۲	۰/۰۰۸	۰/۴۹	۰/۰۰۹	۰/۵۴	۰/۰۱۰

جدول ۱: میانگین و واریانس ضرایب رگرسیون در ۱۰۰ تکرار

مراجع

- Golan, A., Judge, G. G. and Karp, L. (1996a), A Maximum Entropy Approach to Estimation and Inference in Dynamic Models or Counting Fish in the Sea Using Maximum Entropy, *J. of Economic Dynamics and Control*, 20, 559-582.
- Golan, A., Judge, G. G. and Robinson, S. (1994), Recovering Information from Incomplete or Partial Multisectoral Economic Data, *Review of Economic and Statistics*, 76, 541-549.
- Judge, G. G. and Golan, A. (1993), Recovering Information in the Case of Ill-Posed Inverse Problem with Noise, *Mimeo Department of Agricultural and Natural Resources, University of California, Berkeley, Ca..*
- Kitamura, Y. and Stutzer, M. (1998), An Information-Theoretic Alternative to GME Estimation, *Econometrica*, 65, 861-874.
- Pukelsheim, F. (1994), The Three Sigma Rule, *The American Statistician*, 48, 88-91.
- Shannon, C. E. (1948), A Mathematical Theory of Communication, *Bell System Technical Journal*, 27, 379-423.

آشنایی با داده‌کاوی

علی‌اصغر حائری، دکتر حمیدرضا نواب‌پور

P ۱۳۰ ۱۵

گروه آمار، دانشگاه علامه طباطبائی

چکیده: با گسترش روز افزون فناوری رایانه‌ای و مجهز شدن سازمانها، شرکتهای دولتی و خصوصی و ... به این ابزار، امکان جمع آوری اطلاعات دقیق و جدید با قیمت ارزان برای آنها فراهم آمده، به طوری که این امکان منجر به ایجاد مجموعه‌های بسیار وسیعی از داده‌ها شده است، که شامل مقادیر زیادی از سوابق ثبت شده هستند. این سوابق مدنظر مدیران، برنامه‌ریزان، پژوهشگران و ... بوده و از آنها برای تهیه گزارش‌های هفتگی، ماهانه، پرسنلی، خرد، فروش، موجودی انبار و ... استفاده می‌شود. اخیراً داده‌کاوی^۱ موضوع تعداد زیادی از مقالات در تجارت و مجله‌های آماری و نرم‌افزاری شده است. ولی تا چند سال قبل فقط تعداد کمی از مردم عبارت داده‌کاوی را شنیده بودند، اگرچه داده‌کاوی تکامل یافته رشته‌ای با تاریخ قدیمی است ولی عبارت داده‌کاوی اخیراً و در دهه ۱۹۹۰ مطرح شده است. داده‌کاوی به استخراج ساختار مورد علاقه از پایگاه داده‌ها به ویژه انبار وسیع داده‌ها مربوط می‌شود. در این مقاله سعی شده است مفهوم کلی داده‌کاوی، ضرورت استفاده از داده‌کاوی، روش‌های تحلیل داده‌ها و کاربرد آن مدنظر قرار گیرد.

واژه‌های کلیدی: داده انبار^۲، پردازش تحلیلی آنی (OLAP)^۳، ردیابی^۴، پردازش معاملات آنی (OLTP)^۵، پایگاه داده‌ها^۶ ساختار داده‌ها^۷.

۱ مقدمه

گسترش روزافزون فناوری رایانه‌ای امکان جمع آوری اطلاعات دقیق، جدید و لحظه به لحظه با قیمت ارزان را فراهم آورده است به طوری که این امکان منجر به افزایش انفحار آمیز منابع اطلاعاتی در دهه‌های اخیر شده و ضرورت گردآوری و پردازش منابع اطلاعاتی

^۱ Data Mining

^۲ Data Warehouse

^۳ On-Line Analytic Processing

^۴ Classification

^۵ On-Line Transactional Processing

^۶ Database

^۷ Data Structure

را آشکارتر نموده است. استخراج و کشف سریع و دقیق اطلاعات با ارزش از این منابع اطلاعاتی از جمله اموری است که هر سازمان، شرکت و کشوری به منظور توسعه علمی، فنی و اقتصادی خود به آن نیاز دارد. هم‌اکنون در هر کشوری اغلب سازمانها و شرکتها برای پرسنل، بازرگانی و حتی مراجعه کنندگان، پایگاه داده‌ها ایجاد یا خردباری نموده‌اند به طوری که برای مدیران، برنامه‌ریزان و پژوهشگران و ... جهت گرفتن تصمیمهای راهبردی، گزارش‌های مختلف و توصیف وضعیت کنونی خود و ... می‌تواند حیاتی باشد. بعضی مواقع اطلاعات موردنیاز در بیش از یک پایگاه داده‌ها می‌باشد. اغلب این پایگاه‌های داده‌ها همگن نیستند؛ بنابراین به داده انبار نیاز خواهیم داشت که به تفصیل به آن خواهیم پرداخت. پس از تشکیل داده انبار برای استخراج الگوها و اطلاعات با ارزش به داده‌کاوی خواهیم پرداخت. در ادامه به مفهوم داده انبار و داده‌کاوی، روش‌های داده‌کاوی، مراحل، کاربرد و ضرورت داده‌کاوی خواهیم پرداخت.

۲ مفهوم داده انبار

داده انبار یکی از ابزارهای فناوری مدیریت داده‌ها برای حمایت از داده‌کاوی است و خیلی شبیه به یک یا ترکیبی از چند پایگاه داده‌ها می‌باشد و با زمان تغییر می‌کند. اغلب پایگاه‌های داده‌های مختلف همگن نبوده و تشکیل داده انبار مشکل می‌باشد. به نظر اینمان داده انبارها موضوعات جهت داده شده‌ای هستند که تا حد زیادی به کاربرد داده انبار بستگی دارد.^۱ این داده انبارها ممکن است از منابع داده‌های ناهمگن و مختلف درست شوند. ایده‌ای که در پشت مفهوم داده انبار وجود دارد این است که چندین مدل پردازش برای کمک به یک درخواست اطلاعات از پایگاه‌های داده‌ها در محیط‌های ناهمگن با هر فرایندی لازم است. بنابراین داده انبار به دلیل وجود پایگاه داده‌های ناهمگن ایجاد می‌شود. اطلاعات، درون این پایگاه‌های داده‌ها ترکیب شده و در یک داده انبار قرار می‌گیرند. راههای مختلفی برای ترکیب اطلاعات وجود دارد:

یک روش، تکرار ساده پایگاه داده‌ها می‌باشد. این روش هیچ مزیتی نسبت به تعیین پایگاه داده‌های ناهمگن ندارد.

دومین روش، تکرار اطلاعات می‌باشد؛ اما هر اطلاعات ناسازگار و اضافه‌ای حذف می‌شود. این روش چند مزیت دارد. مثل اهمیت به دست آوردن یک تصویر سازگار از پایگاه داده‌ها.

^۱ Thuraisingham, Bahavan (1999)

سومین روش، انتخاب زیرمجموعه‌ای از اطلاعات پایگاه داده‌ها و انتقال آن به یک انبار می‌باشد. چندین پیامد در اینجا وجود دارد. چگونه زیرمجموعه‌ها انتخاب می‌شوند؟ زیرمجموعه‌های انتخابی تصادفی هستند یا روش‌هایی برای انتخاب داده‌ها به کار برده می‌شود؟ (برای مثال، فرد هر سطری را در یک رابطه می‌تواند انتخاب کند و این سطرها را در انبار ذخیره کند («با فرض اینکه پایگاه داده‌ها رابطه‌ای است»)).

چهارمین روش که تغییر جزئی روش سوم می‌باشد، انواع سؤالاتی را که کاربران ممکن است مطرح کنند تعیین نموده و سپس به تجزیه و تحلیل داده‌ها و صرفاً به ذخیره‌سازی داده‌هایی که موردنیاز کاربر است می‌پردازد. این روش «پردازش تحلیلی آنی» (OLAP) نامیده می‌شود که در مقابل «پردازش معاملات آنی» (OLTP) که عمل تولید داده را انجام می‌دهد، قرار دارد.

داده‌ها اغلب ممکن است به صورت متفاوت و برای کاربردهای مختلف منظور شوند که آنها را داده‌های چندبعدی گویند. بنابراین می‌توان گفت «داده انبار» سیستمی برای ذخیره و ارائه توده‌ای از مقادیر داده‌ها است. روش شناختی ویژه‌ای که برای داده انبار طرح شده است باید برای استفاده در کشوربررسی شود. اگرچه تشابهاتی وجود دارد لیکن برخی از امور در روش شناختی داده انبار با گسترش کاربرد امور مرسوم مغایرت دارد.

روش شناختی داده انبار در اموری که در پژوهش‌های داده انبار گنجانده شده یکپارچگی ایجاد می‌کند. کشف دانش در پایگاه داده‌ها^۱ (KDD) در جهت کشف اطلاعات مفید از مجموعه بزرگ داده‌ها می‌باشد. دانش کشف شده می‌تواند قاعده‌های باشد تا ویژگی‌های داده‌ها، الگوهایی که به طور متناسب رخ می‌دهند، خواص بندی موضوعات در درون پایگاه داده‌ها و غیره را توصیف کند.

یک کاربرسیستم KDD بایستی درک بالایی از قلمرو داده‌ها به منظور انتخاب زیرمجموعه صحیحی از داده‌ها، رده مناسبی از الگوها و معیار خوبی برای الگوهای جالب داشته باشد. بنابراین سیستم KDD باید ابزارهایی با اثر تعاملی داشته باشد نه سیستمهای تجزیه و تحلیل خودکار. لذا کشف دانش از پایگاه داده‌ها باید مثل یک فرایند شامل گامهای زیر باشد:

- ۱) درک قلمرو
- ۲) آماده کردن مجموعه داده‌ها
- ۳) کشف الگوها (داده‌کاوی)
- ۴) پردازش بعد از کشف الگو

^۱ Knowledge Discovery of Database (KDD)

(۵) استفاده از نتایج

۴ مفهوم داده‌کاوی

عبارت داده‌کاوی مترادف با یکی از عبارتهای استخراج دانش، برداشت اطلاعات، بررسی داده‌ها و حتی لایروبی کردن داده‌ها^۱ است که در حقیقت کشف دانش در پایگاه داده‌ها را توصیف می‌کند. بنابراین این ایده که مبنای داده‌کاوی است، یک فرایند با اهمیتی از شناخت الگوهای بالقوه مفید، تازه و درنهایت قابل درک در داده‌ها می‌باشد. واژه کشف دانش در پایگاه داده‌ها در سال ۱۹۸۹ در مراجعه به مفهوم کلی، گستردگی و سطح بالا و به دنبال جستجوی دانش در اطلاعات شکل گرفته است. عبارت داده‌کاوی فنون کاربردی ابزارهای به کار برده شده برای بیان و تجزیه و تحلیل داده‌ها و برای نشانه‌های داوری می‌باشد. این اصطلاح (داده‌کاوی) را آمارشناسان، تحلیلگران داده‌ها و انجمن سیستم‌های اطلاعات مدیریت به کار برده‌اند. در حالیکه محققین یادگیری ماشین و هوش مصنوعی از KDD ها بیشتر استفاده می‌کنند. بنابراین داده‌کاوی چیست؟ در زیر چند تعریف ارائه می‌شود.

«داده‌کاوی در حقیقت کشف ساختارهای جالب توجه، غیرمنتظره و با ارزش از روی مجموعه وسیعی از داده‌ها می‌باشد و فعلیتی است که اساساً با آمار و تحلیل دقیق داده‌ها منطبق است.» David Hand (1999)

«داده‌کاوی فرایند کشف رابطه‌ها، الگوها و روندهای جدید معنی‌داری است که به بررسی حجم وسیعی از اطلاعات ذخیره شده در ابزارها با فناوریهای تشخیص الگو (مانند ریاضی و آمار) می‌پردازد.» Gartner Group

«داده‌کاوی یا به تعبیر دیگر کشف دانش در پایگاه داده‌ها، استخراج غیربدیهی اطلاعات بالقوه مفید از روی داده‌هاست که قبل ناشناخته مانده‌اند. این مطلب برخی از روش‌های فنی مانند خوشبندی، خلاصه‌سازی داده‌ها، فراگیری قاعده‌های رده‌بندی، یافتن ارتباط شبکه‌ها، تحلیل تغییرات و کشف بی‌قاعده‌گی‌ها را شامل می‌شود.» William. J. Frawley, Gregory Piatetsky, Shapiro and Christopher. J. Matheus

بنابراین معمولی‌ترین فنون به کار رفته در داده‌کاوی عبارتند از:

۱) درختهای تصمیم^۲

^۱ data dredging

^۲ Decision Trees

- ۲) شبکه‌های عصبی^۲
- ۳) قاعده‌های استقراء^۱
- ۴) نزدیکترین همسایگی^۲
- ۵) الگوریتم‌های ژنتیک^۳
- و

۴ اختلافها و تشابهات داده‌کاوی با آمار

داده‌کاوی که کشف دانش در پایگاه داده‌ها (KDD) نیز نامیده می‌شود، همپوششی زیادی با آمار دارد به طوری که ممکن است برخی از مردم KDD را زیرمجموعه‌ای از آمار تلقی کنند. این رویکرد برای برخی از آمارشناسان محافظه‌کار مشترک است ولی واقع‌گرایانه نیست. داده‌کاوی از ایده‌ها، ابزارها و روش‌هایی از دیگر سطوح به ویژه یادگیری ماشین، نظریه و فناوری پایگاه داده‌ها استفاده می‌کند و ضرورتی ندارد که به برخی از سطوحی که مدنظر آمار است مربوط شود. در حالت کلی، داده‌کاوی با مجموعه‌های داده‌های بزرگ سروکار دارد. برای آمارشناسان مجموعه داده‌ها با ۱۰۰۰ رکورد ممکن است خیلی بزرگ تلقی شود اما به نظر می‌رسد که در مقایسه با پایگاه داده‌های با یک میلیون بایت کاملاً کوچک باشد. واضح است که برای بررسی مقادیر داده‌های این چنینی شخص بایستی از رایانه استفاده کند، چرا که داده‌کاوی پیش‌زمینه مهمی در علوم رایانه‌ای به ویژه در نظریه پایگاه داده‌ها دارد.

با این حال عبارت داده‌کاوی برای فرد آماری جدید نیست و اغلب در یک مفهوم جرئی مانند فرایند تولید مدل‌های زیادی به منظور به دست آوردن برازش خوب به کار برده شده است.

در آمار تمرکز روی استنباط است و وظیفه اصلی آمار وقتی نمونه‌ای از یک جامعه مشاهده می‌شود استنباط درباره جامعه است و فرضیات آزمون می‌شوند ولی در داده‌کاوی با کل داده‌ها (جامعه) سروکار داریم. فرضیات از روی داده‌ها ساخته می‌شوند در چنین مواردی به آزمون معنی‌داری احتیاجی نیست و مقدار مشاهده شده آماره همان مقدار پارامتر است.

در آمار مدل سیمای دیگری داشته و نقش مهمی را ایفا می‌کند، هر چند که ممکن

^۲ Neural Networks

^۱ Induction Rules

^۲ Nearest Neighbour

^۳ Genetic Algorithms

است معانی مختلفی داشته باشد، به طوری که همه مکانیسم‌ها در آمار به مدلسازی ختم می‌شوند. محاسبات، معیار انتخاب و نمونه‌گیری وغیره جزئیات کمکی برای ساختن یک مدل خوب هستند. ضرورتی ندارد که مدل برای داده‌کاوی درست باشد از این رو الگوریتم‌ها برای جستجوی الگو خیلی مهم هستند هر چند بعضی از این الگوریتم‌ها در آمار خلق شده‌اند ولی برای مجموعه داده‌های بزرگ به تدبیر جدیدتر، سریعتر، مؤثرتر و همچنین وجود KDD احتیاج است.

داده‌کاوی کشف اطلاعات با ارزش ولی ناشناخته است و فرایند آن به طور ضروری برخلاف تجزیه و تحلیل آماری، صرفاً اکتشافی است.

۵ مراحل داده‌کاوی

در اینجا فرایندی برای استخراج دانش نهفته از داده‌انبار، فایل اطلاعات مشتری یا هرپایگاه داده‌های شرکتی‌های دیگر وجود دارد.

۱) تعیین موضوع: قبل از شروع باید روش کرد که با تجزیه و تحلیل خود انتظار داریم به چه چیزی دست یابیم، از قبل بایستی هدف از داده‌کاوی را تعیین کرده و برای خود ثابت کنیم که این هدف قابل اندازه‌گیری است یا نه؟ برخی از اهداف ممکن عبارتند از:

- یافتن رابطه فروش بین محصولات ویژه یا خدمات
- تعیین الگوهای خرید خاص و فوق العاده
- تعیین انواع پتانسیلهای مشتریها
- یافتن روند فروش محصولات

۲) انتخاب داده‌ها: در مرحله قبل هدف تعیین گردید. گام بعدی، انتخاب داده برای رسیدن به این هدف می‌باشد. این امر می‌تواند زیر مجموعه‌ای از داده‌انباری که در اختیار داریم یا از فروشگاه داده‌ها^۱ که حاوی اطلاعات خاص یا فایل اطلاعات مشتری است، باشد. برای این کار تا حد امکان گستره داده‌هایی که بایستی کاویده شوند را بخش‌بندی کنیم.

چند نتیجه کلیدی وجود دارد.

— آیا داده‌ها مناسب برای توصیف پدیده‌هایی هستند که تحلیل داده‌کاوی تلاش

^۱ data mart

- می‌کند برای آنها مدلسازی کند؟
- آیا می‌توانیم سوابق مشتری داخلی را با شیوه زندگی بیرونی داده‌های جمعیتی مطابقت دهیم؟
- آیا داده‌ها ثبات دارند – آیا بعد از تحلیل ویژگیها ثابت می‌ماند؟
- اگر پایگاه داده‌ها را ادغام نماییم آیا می‌توانیم زمینه مشترکی برای پیوند دادن آنها بیابیم؟
- داده‌ها چه رابطه‌ای با هدف کاری دارند؟
- ۳) آماده کردن داده‌ها: در مرحله قبل داده‌ها گردآوری شد. حال باید تصمیم بگیریم کدامیک از ویژگیهای موجود در قالبها قابل استفاده هستند. سهم متخصصین، تولید کنندگان و کاربران را از قلمرو داده‌ها ملاحظه کنیم.
- اتخاذ تدابیری برای بررسی داده‌های گمشده، نوفه خارجی و داده‌های دورافتاده
- تعیین متغیرهای اضافی در مجموعه داده‌ها و تصمیم‌گیری درباره اینکه کدام فیلد باید خارج شود
- در صورت لزوم تصمیم‌گیری درباره تبدیل لگاریتمی، توان دوم، یا ...
- بازرسی شهودی مجموعه داده‌ها برای درک پایگاه داده‌ها
- تعیین توزیع فراوانی داده‌ها
- می‌توانیم برخی از این تصمیم‌ها را تا انتخاب ابزار داده‌کاوی به تعویق بیندازیم. مثلاً اگر ما به یک شبکه عصبی یا شبکه چندبعدی نیاز داشته باشیم ممکن است مجبور شویم برخی از فیلدها را تغییر دهیم.
- ۴) بازرسی داده‌ها: ارزیابی ساختار داده‌ها به منظور تعیین ابزارهای مناسب
- نسبت صفات رسته‌ای یا دو حالتی در پایگاه داده‌ها چیست؟
- طبیعت و ساختار پایگاه داده‌ها چیست؟
- توزیع مجموعه داده‌ها چیست؟
- بایستی ارزیابی عینی ساختار داده‌ها را در مقابل نیاز کاربران برای درک یافته‌ها توانیم بخشیم. برای مثال شبکه‌های عصبی نتایج را نمی‌توانند توصیف کنند.
- ۵) انتخاب ابزارها: اهداف کاری و ساختار داده‌ها دو رهنمای مهم برای انتخاب ابزار داده‌کاوی مناسب می‌باشند. هر دو بایستی ما را به ابزار یکسان هدایت کنند. موقعی که یک مجموعه از ابزارهای بالقوه را ارزیابی می‌کنیم، این پرسش‌ها را ملاحظه کنیم.
- آیا مجموعه داده‌ها صرفاً رسته‌ای است؟
- چه برنامه‌ای برای حمایت از ابزارهای منتخب خود داریم؟
- آیا ابزارهای انتخابی سازگار با ODBC^۱ یک رابط برنامه‌ای از شرکت مایکروسافت

^۱ Open Database Connectivity (ODBC)

است که ارتباط بین برنامه‌های کاربردی ویندوز را برای دسترسی به پایگاههای داده‌های روی شبکه فراهم می‌سازد. هستند؟

— چه قالبی از داده‌ها می‌تواند برای ابزارها مفهوم داشته باشد؟

ابزارهای متعددی وجود دارد که به طور مشابه برای پروژه داده‌کاوی تهیه می‌کنند. تعدادی از ابزارها چندین فناوری را در یک مجموعه از برنامه‌های تحلیل آماری، شبکه عصبی و رده‌بندی کنندۀ نمادین جمع می‌کنند.

۶) قالب پاسخ: به همراه بازرسی داده‌ها، هدف کاری و انتخاب ابزار، قالب پاسخ را تعیین می‌کنند. سوالات کلیدی عبارتند از:

— قالب بهینه حل درخت تصمیم— کد C، دستور SQL— چیست؟

— قالب مناسب گزینه چیست؟

— هدف از حل چیست؟

— کاربر به چه چیزی احتیاج دارد— نمودارها، گزارشها، برنامه—؟

۷) طرح ریزی مدل: در این هنگام است که فرایند داده‌کاوی آغاز می‌شود. عموماً اولین گام، استفاده از چند رده تصادفی برای جدا کردن داده‌ها در مجموعه تولیدی، مجموعه آزمون و ایجاد و ارزیابی یک مدل می‌باشد. تولید قاعده‌های رده‌بندی، درختهای تصمیم، خوشه‌بندی زیرگروه‌ها، امتیازها، برنامه‌ها، وزن‌ها و ارزیابی داده‌ها یا نسبت خطاهای در این مرحله قرار دارند. تفکیک کردن این پیامدها:

— آیا نسبت خطاهادر سطح قابل قبولی است؟ آیا می‌توانیم آنها را بهبود دهیم؟

— چه صفات خارجی را یافته‌ایم؟ آیا می‌توانیم آنها را حذف کنیم؟

— آیا داده‌های اضافی یا روش‌شناختی مختلفی لازم است؟

— آیا مجموعه داده‌های جدیدی تولید و آزمون خواهیم کرد؟

۸) ارزیابی یافته‌ها: نتایج تجزیه و تحلیل را با مشتری کاری یا متخصص قلمرو به بحث و تبادل نظر گذاشته و مطمئن می‌شویم که یافته‌های برای اهداف کاری صحیح و مناسب هستند.

— آیا یافته‌ها حساسیت‌ساز هستند؟

— آیا برای هر گام قبلی و بهبود نتایج برمی‌گردیم؟

— آیا می‌توانیم ابزارهای دیگر داده‌کاوی را برای تکرار یافته‌ها به کار ببریم؟

۹) ارایه یافته‌ها: یک گزارش نهایی برای واحد کاری یا مشتری کاری یا متشتمی تهیه می‌کنیم. این گزارش باید مستند به فرایند داده‌کاوی بی‌عیب و نقص شامل آماده‌سازی داده‌ها، ابزارهای به کار رفته، نتایج آزمون، کد منابع و قاعده‌ها باشد.

بعضی از این نتایج عبارتند از:

– داده‌های اضافی تجزیه و تحلیل را بهبود خواهد داد؟

– چه بیانش اساسی‌ای پوشش دهیم و چگونه قابل کاربرد است؟

– از تجزیه و تحلیل داده‌کاوی چه نتایجی می‌تواند پیشنهاد شود؟

– آیا یافته‌ها مناسب هدف کاری هستند؟

۱۰) هماهنگ کردن پاسخ‌ها: یافته‌ها را با کلیه عالیق کاربران در واحدهای کاری مناسب تقسیم می‌کیم. ما می‌توانیم به طور کامل نتایج تجزیه و تحلیل در روشهای کارشکن را به هم پیوند دهیم، بعضی از حل‌های داده‌کاوی ممکن است موارد زیر را در برگیرد.

– دستورهای SQL برای توزیع مصرف کنندگان نهایی

– تولید سیستم با همکاری کدهای C

– تجمیع قاعده‌ها در نظام تصمیم‌گیری

اگر چه ابزارهای داده‌کاوی، پایگاه داده‌ها را به صورت خودکار تجزیه و تحلیل می‌کنند، ولی اگر مراقب نباشیم ممکن است به یافته‌های نادرست و نتایج غلط منجر شود. در نظر داشته باشید که داده‌کاوی یک فرایند کاری با هدف مشخص است، که برای استخراج دیدگاه‌های رقابتی از سوابق تاریخی در پایگاه داده‌ها می‌باشد.

۶ روشهای داده‌کاوی

الف) گزارش‌دهی و پردازش تحلیل آنی

موقعی که یک تحلیلگر در زمان تجزیه و تحلیل با OLAP کار می‌کند در مورد چیزهایی که به دنبال آن است پیش‌زمینه‌ای دارد و با فرضیاتی که مورد بررسی قرار گرفته شروع می‌کند . در حالیکه در مورد داده‌کاوی، تحلیلگر هیچ شناخت قبلی در رابطه با نتایج احتمالی ندارد. کاربران از طریق OLAP به دنبال پاسخ به پرسش‌های ایشان هستند. هر جستجو برای دستیابی به پرسش خاص ما را با مرحلهٔ پیچیده‌تری رویرو می‌کند، و کاربر به دانش قبلی از نتایج مورد انتظار نیاز دارد. فرایند در داده‌کاوی کاملاً متفاوت است. در حالیکه OLAP برای تجزیه و تحلیل گذشته و کسب بینش‌های تازه به کاربر کمک می‌کند، داده‌کاوی برای پیش‌گویی آینده به کاربر کمک می‌کند.

ب) مدلسازی نظریه رهنمون

همبستگی

آزمون t

تحلیل واریانس

رگرسیون خطی
رگرسیون لوزیستیک
تحلیل تشخیصی
مدلهای پیش‌بینی

ج) مدلسازی داده رهنمون

تحلیل خوش‌های
درختهای تصمیم^۱
تصویر کردن داده‌ها^۲
شبکه‌های عصبی
قاعده‌های پیوند^۳
قاعده استقراء^۴

۷ کاربردهای داده‌کاوی

داده‌کاوی در رشته‌های مختلفی کاربرد دارد که در زیر به برخی از این موارد اشاره می‌کنیم:

- ۱) بازاریابی / خرد فروشی
 - تشخیص الگوهای خرید مشتریها
 - یافتن پیوند میان مشخصه‌های جمعیت‌شناسنامه مشتری
 - پیش‌گویی پاسخ برای (پست کردن) اوراق مبارزه انتخاباتی
 - تجزیه و تحلیل سبد بازار
- ۲) بانکداری
 - کشف الگوهای استفاده از کارتهای اعتباری قلابی
 - تشخیص مشتریهای ثابت (وظیفه‌شناس)
 - پیش‌گویی احتمالی مشتریها برای تغییر کارتهای اعتباری پیوسته آنها
 - تعیین پرداخت کارتهای اعتباری برای گروههای مشتریان
 - یافتن ارتباط مخفی بین نشانگرهای مالی مختلف
 - تشخیص قواعد بازرگانی موجود از داده‌های بازاریابی

^۱ Decision Trees

^۲ Data Visualization

^۳ Association Rules

^۴ Rule Induction

(۳) بیمه و مراقبت بهداشتی

- تجزیه و تحلیل ادعاهای - یعنی اینکه کدامیک از روش‌های پزشکی با هم ادعا شده‌اند
- پیش‌گویی اینکه مشتریها بیمه‌نامه جدید بخوبند
- تشخیص الگوهای رفتاری مشتریان خطری
- تشخیص رفتارهای قلابی

(۴) حمل و نقل (انتقال)

- تعیین توزیع فهرست مقدار فروش
- تجزیه و تحلیل الگوهای بارگیری

(۵) پزشکی

- تعیین الگوی رفتار بیماران برای پیش‌بینی مراجعات به مطب (ویژتهایی که در مطب صورت می‌گیرد)
- تشخیص معالجات موفق پزشکی برای بیماری‌های مختلف

(۶) متن‌کاوی^۱

(۷) شبکه‌کاوی^۲

(۸) صوت‌کاوی^۳

(۹) تصویرکاوی^۴

(۱۰) سیستم هوشمناسی

(۱۱) سیستم ثبت احوال کشور

(۱۲) صنعت

مراجع

Thuraisingham, Bahavan (1999), *Data Mining Technologies, Techniques, Tools, and Trend*, CRC press LLC, Florida.

Michael. J. A. Berry Gordon Linoff (1997), *Data Mining Techniques*, John Wiley, New York.

^۱ Text Mining

^۲ Web Mining

^۳ Audio Mining

^۴ Video Mining

David Hand (1999), *Why data mining is more than statistics writ large*, Imperial College of Science, Technology, and Medicine, Department of Mathematics.

Gartner Group, <http://www.spss.com>

Han Jiawei, Kamber Micheline (2001), *Data Mining concepts and techniques*, Morgan Kaufmann.

Heikki Mannila (1996), *Data Mining: machine learning, statistics, and databases*, 8th International Conference on Scientific and Statistical Database Management, Stockholm, June 18-20, 1996, pp 1-8.

Paulraj Ponniah(2001), *Data Warehousing Fundamentals*, John Wiley, New York.

Rud, Olivia Parr (2001), *Data Mining Cookbook*,John Wiley, New York.
Data Mining vs. Statistics (1997), <http://www.pmsi.fr/dminita.htm>, Article
and figure originally published in 1997 in Science Tribune.

William. J. Frawly, Gregory Piatetsky, Shapiro and Christopher. J. Matheus,
http://pcc.qub.ac.uk/tec/courses/datamining/stu_notes/dm_book_6.html

<http://www.sas.com>

<http://www.kdnuggets.com/fad/data-mining.html>

کاهش اربی انتخاب در آزمایشهای دنباله‌ای

آرزو حبیبی راد

P ۱۲۰۳۱

گروه آمار، دانشگاه فردوسی مشهد

چکیده برای مقایسه اثرات دو یا چند تیمار در یک طرح تحقیقی، در بسیاری از موارد این مشکل وجود دارد که به افراد (واحدهایی) که آزمایشها روی آنها باید انجام شود، در یک زمان دسترسی نداریم بلکه به این افراد (واحدها) در طول زمان دسترسی پیدا می‌کنیم، به طوری که پس از دسترسی به هر واحد لازم است تیمار مناسب بالافاصله به آنها تخصیص داده شود. حال اگر آزمایشگر آگاه باشد و یا حدس بزند که چه تیماری توسط آن آماردان به واحد تخصیص داده می‌شود، می‌تواند آگاهانه یا ناآگاهانه با اعمال نفوذ در انتخاب واحدها برای آزمایش، باعث ایجاد اربی از نتایج آزمایش شود، که اربی حاصل از این اعمال نفوذ را اربی انتخاب می‌نامیم. از این رو در این مقاله به دنبال ارائه طرحهایی هستیم که در آن با توجه به استراتژی‌های استفاده شده توسط آزمایشگر متوجه تعداد حدسهای درست او (آزمایشگر) را تا حدامکان کاهش دهیم.

واژه‌های کلیدی: اربی انتخاب، استراتژی همگرا، استراتژی احتمال متناسب، طرح سکه اربی Efron، طرح کاسه‌ای، طرح سکه اربی تطبیقی.

۱ مقدمه

فرض کنید آزمایشگر (E) علاقمند به مقایسه اثرات دو تیمار A و B بر روی یک موضوعی در جامعه باشد، اما ممکن است این مشکل وجود داشته باشد که وی به واحدها (افراد) در یک زمان دسترسی نداشته باشد بلکه به این واحدها (افراد) در طول زمان دسترسی پیدا کند بطوریکه پس از دسترسی، آزمایشگر (E) باید تصمیم بگیرد که آیا این واحدها برای آزمایش مناسب است یا خیر، و در صورت مثبت بودن پاسخ وی، آنگاه آمارداران (S) تصمیم می‌گیرد که تیمار A یا تیمار B را به واحد مورد نظر انتساب دهد و این روند تا انتخاب نمونه ای به حجم n ادامه پیدا می‌کند.

اما اگر آزمایشگر (E) تمایل به برتری تیمار A نسبت به B داشته باشد می‌تواند آگاهانه یا ناآگاهانه براساس حدس خود واحدهایی را برای تیمار A انتخاب کند که متوسط پاسخ مورد انتظار آن $\Delta + \mu$ (اربی در هر آزمایش می‌باشد) و بر عکس اگر حدس بزند که

تیمار B توسط آماردان انتخاب می شود و احدهایی را انتخاب کند که متوسط پاسخ مورد انتظار آنها $\Delta - \mu$ است.

تعريف ۱: روشی را که طبق آن در هر مرحله آزمایشگر حدس می زند که آماردان چه تیماری را به مریض بعدی اختصاص می دهد استراتژی آزمایشگر می نامیم.

تعريف ۲: روشی را که در آن آماردان برای انتخاب تیمارها اتخاذ می کند طرح آماردان نامیده می شود.

تعريف ۳: کمیت $E(G) - n\Delta$ اربیی انتخاب نامیده می شود که در آن G تعداد حدهای درست آزمایشگر است.

۲ طرحهای آماردان (S)

۱.۲ طرح سکه اریب (EFRON)

فرض کنید r آزمایش از $2n$ (علوم) آزمایش انجام شده، برای آزمایش $(1+r)$ ام اگر تیمار A ، D_r نا با گروه تیمار B اختلاف داشته باشد، انتخاب تیمار برای آزمایش $(1+r)$ ام توسط آماردان به قرار زیر صورت می گیرد

۱- اگر $0 < D_r$ باشد، با احتمال A تیمار B را انتخاب می کند.

۲- اگر $0 = D_r$ باشد، با احتمال $\frac{1}{2}$ تیمار A و با احتمال $\frac{1}{2}$ تیمار B را انتخاب می کند.

۳- اگر $0 < D_r$ باشد، با احتمال p تیمار A و با احتمال q تیمار B را انتخاب می کند.

بطوریکه همواره $p + q = 1, p \geq q$.

طرح فوق را طرح سکه اریب با اربیی p می نامیم و با علامت اختصاری $BCD(p)$ نمایش می دهیم.

در حالیکه $\frac{2}{3} = p$ فرض شود اربیی انتخاب و متعادل بودن طرح (برابری حجم نمونه برای دو تیمار A, B) نسبت به سایر حالتها در وضعیت مطلوبتری قرار دارد. از این رو در این مقاله مقدار $\frac{2}{3} = p$ فرض می شود.

۲.۲ طرح کاسه‌ای (The Urn Design)

ظرفی را در نظر بگیرید که از ابتدا α توپ قرمز و β توپ سفید را در خود جای داده بیک توپ به طور کاملاً تصادفی از ظرف بیرون کشیده می‌شود، اگر توپ سفید باشد، آماردان تیمار A را برای مریض بعدی انتخاب می‌کند و اگر قرمز بود تیمار B را انتخاب می‌کند و هر توپ که بیرون آورده شود نه تنها به ظرف بر می‌گردد بلکه β توپ از رنگ مخالف آن نیز به ظرف برگردانده می‌شود، این طرح به طرح کاسه‌ای معروف است و با علامت اختصاری $UD(\alpha, \beta)$ نمایش داده می‌شود. در حالت خاص $\beta > \alpha = 0$ با احتمال $\frac{1}{2}$ یکی از دو تیمار را انتخاب می‌کند و β توپ رنگ مخالف را به ظرف اضافه می‌کند و فرایند انتخاب تیمارها مشابه بالا ادامه می‌دهد.

اما دیده می‌شود که $UD(0, \beta)$ دارای خواص جالبی است از جمله اینکه $UD(0, \beta) = UD(r, \beta)$ از مقدار β مستقل است و دیگراینکه به راحتی قابل اجرا است زیرا اگر تا مرحله r ام، مریض به تیمار A و j مریض به تیمار B نسبت داده باشیم آنگاه تیمار $A(B)$ را با احتمال $\frac{j}{i+j}$ (یا $\frac{i}{i+j}$) به مریض $r+1$ ام نسبت می‌دهیم.

۲.۲ طرح سکه اربیب تطبیقی (The AdaPtive Biased Coin Design)

فرض کنید که در مرحله r ام $(r = 0, 1, 2, \dots, 2n)$ یک آزمایش دنباله‌ای قرار داریم، بطوریکه تا این مرحله i واحد به تیمار A و j واحد را به تیمار B نسبت داده باشیم. در این طرح آماردان برای انجام مرحله $(r+1)$ ام آزمایش با احتمال $p = p(\frac{D_r}{r})$ تیمار A و با احتمال $q = q(\frac{D_r}{r})$ تیمار B را انتخاب می‌کند. (که $D_r = i - j$). تابع غیر صعودی از $(\frac{D_r}{r})$ است بطوریکه مقادیر p در بازه $[0, 1]$ تغییر می‌کند) از طرفی $p + q = 1$ فرض شده است و همچنین $p(x) = q(-x)$ برای هر $x \in [-1, 1]$. توجه کنید که طرح سکه اربیب $BCD(p)$ حالت خاصی از طرح فوق است اگر $p(x)$ مقدار ثابت (p) بزرگتر از $\frac{1}{2}$ فرض شود.

و اگر $\frac{1-x}{2} = p(x)$ فرض شود طرح فوق معادل طرح $UD(0, \beta)$ است. (Wei ۱۹۷۷a)

۳ استراتژیهای آزمایشگر (E)

۱.۳ استراتژی همگرا (Convergent Strategy)

متداولترین استراتژی برای آزمایشگر در مقابل طرحهای ارائه شده توسط آماردان استراتژی همگرا (θ_1) می‌باشد. بطوریکه تحت استراتژی همگرا (θ_1) آزمایشگر تیماری را حدس

می زند که تا آن مرحله کمتر انتخاب شده است و در نتیجه واحدی (فردی) را انتخاب می کند تا به تیمار مورد نظر وی پاسخ مطلوبتری را بدهد در صورتیکه تعداد تیمارهای انتخاب شده تا آن مرحله برابر باشند $\frac{1}{2}$ حدس می زند که چه تیماری انتخاب شود.

۲.۳ استراتژی احتمال متناسب (Probability Convergent Strategy)

تحت استراتژی احتمال متناسب (θ_2) آزمایشگر (E) در هر مرحله با احتمال γ تیمار A و با احتمال $1 - \gamma$ تیمار B را حدس می زند، سپس متناسب با حدس خود واحدها (افراد) را انتخاب می کند (یعنی اگر نتیجه حدس وی تیمار A باشد واحد (فردی) با متوسط پاسخ $\Delta + \mu$ و اگر نتیجه حدس وی تیمار B باشد واحد (فردی) با متوسط پاسخ $\Delta - \mu$ انتخاب می کند). مقدار $\frac{1}{\gamma} - \frac{1}{1-\gamma} = \gamma$ فرض می شود (به ترتیب تعداد تیمارهای A, B هستند تا مرحله t ام)

۴ خاصیت مجانبی تعادلی طرح

تعريف ۴: طرحی را به طور مجانبی تعادل گوییم که در آن $\frac{D_{2n}}{n}$ در احتمال به صفر همگرا باشد وقتی که در آن $D_{2n} = N_A - N_B$ و $n \rightarrow \infty$.
 که N_B, N_A به ترتیب تعداد تیمارهای B, A در پایان آزمایش هستند (Efron ۱۹۷۱) و (Wei ۱۹۷۷b) نشان دادند که دو طرح $BCD(\frac{2}{3})$ و طرح $UD(0, \beta)$ که هر دو حالت خاصی از طرح سکه اریب تطبیقی هستند به طور مجانبی تعادل هستند.

۵ اریبی انتخاب طرح آماردان (S)

۱.۵ تحت استراتژی همگرا (θ_1)

اریبی انتخاب برای دو طرح $UD(\alpha, \beta)$ و $BCD(p)$ تحت استراتژی (θ_1) توسط Wei (۱۹۷۷a) محاسبه شده است او همچنین نشان داده که به ازاء $n > 4$ ، طرح $UD_{\theta_1}(0, \beta)$ اریبی انتخاب کمتری از طرح $BCD_{\theta_1}(\frac{2}{3})$ دارد (شکل ۱).

۲.۵ تحت استراتژی احتمال متناسب (θ_2)۱.۲.۵ اریبی انتخاب طرح سکه اریب EFRON تحت استراتژی احتمال متناسب (θ_2)

می دانیم اریبی انتخاب برابر است با $E(G) - n$ ، پس $E(G) - n = 2\Delta(E(G))$ فرض کنید در مرحله i ام اختلاف بین دو تیمار d باشد $(D_r = d)$ ، اگر G_r پیشامد حدس درست آزمایشگر برای مرحله $(r+1)$ ام باشد و $X_{r+1} = A$ پیشامد انتخاب تیمار A در مرحله $(r+1)$ باشد داریم ،

$$P(G_r = 1 | D_r = d) = P(G_r = 1 | D_r = d, X_{r+1} = A)P(X_{r+1} = A)$$

$$+ P(G_r = 1 | D_r = d, X_{r+1} = B)P(X_{r+1} = B)$$

تحت استراتژی θ_2 ، احتمال آنکه آزمایشگر درست حدس بزند، به شرط آنکه تیمار A (B) انتخاب شده برابر $\gamma(1-\gamma)$ است پس می توان نوشت ،

$$P(G_r = 1 | D_r = d) = \gamma P(X_{r+1} = A) + (1-\gamma)P(X_{r+1} = B)$$

حال اگر آماردان از طرح سکه اریب EFRON استفاده کند، احتمال انتخاب تیمار A در مرحله $(r+1)$ ام بستگی به $(i-j=d)$ دارد، به طوریکه اگر $d > 0$ باشد این احتمال برابر p است . و اگر $d < 0$ باشد ، با احتمال q و اگر $d = 0$ باشد با احتمال $\frac{1}{2}$ تیمار A (B) را انتخاب می کند . پس می توان نوشت ،

$$P(G_r = 1) = \sum_{d=-r}^r P(G_r | D_r = d)P(D_r = d)$$

در نتیجه ،

$$E(G) = \sum_{r=0}^{n-1} P(G_r)$$

$$= \frac{1}{2} + \sum_{r=1}^{n-1} \sum_{d=-r}^{-1} (p\gamma + (1-\gamma)q)P(D_r = d) + \frac{1}{2} \sum_{r=1}^{n-1} P(D_r = d)$$

$$+ \sum_{r=1}^{n-1} \sum_{d=1}^r (q\gamma + (1-\gamma)p)P(D_r = d)$$

بنابراین $\gamma = \frac{1}{\gamma} + \frac{d}{\gamma r}$ و $1 - \gamma = \frac{1}{\gamma} - \frac{d}{\gamma r}$ پس

$$\Rightarrow E(G) = \frac{1}{\gamma} + \sum_{r=1}^{2n-1} \sum_{d=-r}^r \left(\frac{1}{\gamma} - \frac{|d|}{\gamma r} (q-p) \right) P(D_r = d)$$

$$\frac{1}{\gamma} + \sum_{r=1}^{2n-1} \left(\frac{1}{\gamma} - \frac{(q-p)E|D_r|}{\gamma r} \right)$$

اگر در طرح کاملاً تصادفی (فرض شود، دارای $p = q = \frac{1}{\gamma}$ EFRON)

$$E(G) = \frac{1}{\gamma} + \sum_{r=1}^{2n-1} \frac{1}{\gamma} = \frac{1}{\gamma} + \frac{2n-1}{\gamma} = n$$

پس اریبی انتخاب طرح فوق، صفر است. ولی می‌دانیم که این طرح در حجم نمونه ای کم نامتعادل است.

در حالتیکه $\gamma = p$ است، اریبی انتخاب و متعادل بودن طرح، نسبت به بقیه حالتها در وضعیت مطلوبتری قرار دارد. در نتیجه تحت استراتژی احتمال متناسب، مقدار $E(G)$ را برای حالت $\gamma = p$ محاسبه می‌کنیم پس با فرض $\gamma = p$ و $q = \frac{1}{\gamma}$ داریم،

$$E(G) = \frac{1}{\gamma} + \sum_{r=1}^{2n-1} \left(\frac{1}{\gamma} + \frac{E|D_r|}{\gamma r} \right)$$

۲.۲.۵ اریبی انتخاب طرح سکه اریب تطبیقی وقتی $p(x) = \frac{1-x}{2}$ تحت استراتژی احتمال متناسب

اگر آماردان از طرح سکه اریب تطبیقی با تابع $p(x) = \frac{1-x}{2}$ استفاده کند اریبی انتخاب طرح فوق تحت استراتژی θ_2 به صورت زیر محاسبه می‌شود

$$P(G_r) + \sum_{d=-r}^r \left(p\left(\frac{d}{r}\right)\gamma + (1-\gamma)q\left(\frac{d}{r}\right) \right) P(D_r = d)$$

اگر $q\left(\frac{d}{r}\right) = \frac{1}{\gamma} + \frac{d}{\gamma r}$ و $p\left(\frac{d}{r}\right) = \frac{1}{\gamma} - \frac{d}{\gamma r}$ فرض شود، پس $p(x) = \frac{1-x}{2}$ در نتیجه

$$E(G) = \frac{1}{\gamma} + \sum_{r=1}^{2n-1} \sum_{d=-r}^r \left(\left(\frac{1}{\gamma} - \frac{d}{\gamma r} \right)^2 + \left(\frac{1}{\gamma} + \frac{d}{\gamma r} \right)^2 \right) P(D_r = d)$$

$$\begin{aligned}
 &= \frac{1}{\gamma} + \sum_{r=1}^{\gamma n-1} \sum_{d=-r}^r \left(\frac{1}{\gamma} + \frac{d\gamma}{\gamma r^\gamma} \right) P(D_r = d) \\
 &= \frac{1}{\gamma} + \sum_{r=1}^{\gamma n-1} \left(\frac{1}{\gamma} + \frac{E(D_r)}{\gamma r^\gamma} \right)
 \end{aligned}$$

شکل (۱) اریبی انتخاب طرحهای $BCD(\frac{\gamma}{\gamma}, UD(0, \beta), \theta_1, \theta_2)$ را تحت استراتژیهای برای مقادیر مختلف n نشان می‌دهد.

شکل ۱: آزمایش $2n$ در $(E(G) - n)$ اریبی انتخاب

بنا به شکل واضح است که تحت استراتژی θ_1 ، طرح $UD(0, \beta)$ در مقایسه با طرح $BCD(\frac{\gamma}{\gamma})$ به ازاء $4 > n$ ، اریبی انتخاب کمتری دارد، اما تحت استراتژی θ_2 ، اریبی انتخاب $BCD(\frac{\gamma}{\gamma}, UD(0, \beta))$ به ازاء هر n بیشتر از اریبی انتخاب طرح $BCD_{\theta_2}(\frac{\gamma}{\gamma})$ است، از طرفی

مشاهده می شود که $BCD_{\theta_1}(\frac{3}{4}) > n$ ، اریبی انتخاب کمتری نسبت به طرح $BCD_{\theta_2}(\frac{3}{4})$ دارد.

مراجع

- BLACKWELL, D. & HODGES,J.L.(1957).Design for the control of selection bias. Ann.Math.statist.28,449-60.
- STIGLER ,STEPHEN(1969).The use of random allocation for the control of selection bias.Biometrika .56,553-60.
- EFRON, BRADLEY(1971). Forcing a sequential experiment to be balanced. Biometrika. 58, 403-17.
- WEI,L.J.(1977a),A class of designs for sequential clinical trials. J.Amer. Statist. Assoc. Vol. 72, No.358, 382-386.
- WEI,L.J.(1977b),The adaptiv biased coin design for sequential experiments .The Annals of Statistics .Vol;No.1,92-100.

بررسی بعضی روش‌های رگرسیون جایگزین LS و مقایسه آنها با یکدیگر

محمد رضا ریعی^۱، مجتبی گنجعلی^۲

P11209

^۱ دانشگاه علوم پزشکی و خدمات درمانی گلستان

^۲ دانشگاه شهید بهشتی

چکیده: تعدادی شیوه‌های رگرسیون به عنوان روش‌های جایگزین برای کمترین توان‌های دوم وقتی نقاط نافذ^۱ یا دورافتاده^۲ در داده‌ها وجود دارند، مورد مطالعه قرار گرفته است. هیچ رگرسیون استواری نمی‌تواند در حالت کلی به عنوان بهترین، پذیرفته شود. روش‌های مختلف نقاط ضعف و قوت متفاوتی را با توجه به درصد داده‌های دورافتاده و موقعیت متغیرهای توضیحی و پاسخ، از خود نشان می‌دهند. هدف این مقاله مقایسه روش‌های استوار با سه خاصیت گوناگون الف) نقطه فروپاش ب) کارایی ج) کراندار بودن نقاط موثر است. برآورد ضرایب به دست آمده از روش‌های گوناگون استوار با روش کمترین توان های دوم در دو مثال مقایسه شده‌اند.

واژه‌های کلیدی: کمترین توان‌های دوم، رگرسیون استوار^۳، کارایی، کراندار بودن نقاط موثر^۴، نقطه فروپاش^۵.

۱ مقدمه

وقتی رگرسیون کمترین توان‌های دوم را با n مشاهده برای مدل p پارامتری $Y = X\beta + \epsilon$ به کار می‌بریم فرض‌های بخصوصی را در مورد بردار خط‌ها (ϵ) داریم، یکی از این فرض‌ها آن است که بردار ϵ دارای توزیع $N(o, I\sigma^2)$ است. به‌حال در عمل انحرافاتی از این فرضها وجود دارد. اگر این انحرافات جدی باشند امیدواریم آن‌ها را در رفتار باقیمانده‌ها تشخیص دهیم. بنابراین این انحرافات منجر می‌شود که تعدیلات مناسبی را در مدل و یا در ماتریس متغیرهای کمکی وارد کنیم. مثلاً ممکن است تبدیلی روی یک متغیر یا بیش از یک متغیر وارد کرده یا مدل را با اضافه کردن عبارتهای مرتبه بالاتر تعديل کیم.

^۱ Leverage Pointes

^۲ Outlier

^۳ Robust Regression

^۴ Bounded Influence

^۵ Breakdown Point

اگر در تحلیل این نکته به نظر برسد که خطاهای توزیع غیرنرمال دارند، مخصوصاً در موضعی که توزیع خطای دمها پنهانتری نسبت به توزیع نرمال داشته باشد، یعنی از احتمالهای بیشتری در دمها نسبت به توزیع نرمال برخوردار باشد، بایستی، یک روش رگرسیون استوار را در نظر بگیریم. چنین توزیعهایی دم پهن محتمل هستند که خطاهای بزرگتری را نسبت به حالت نرمال تولید کنند. روش کمترین توانهای دوم در به دست آوردن برآورد پارامترها به هر مشاهده وزن یکسان می‌دهد. روش‌های استوار قادرند به مشاهدات وزن‌های نابرابر اختصاص دهند. به طور کلی مشاهداتی که با قیماندهای بزرگی را تولید می‌کنند به وسیله برآورد استوار کم ورزندند. تعدادی از این روش‌های استوار می‌توانند معرفی و به کار روند.

۲ روش رگرسیون کمترین توانهای دوم (LS)

فرض کنید مدل مورد نظر به صورت زیر باشد:

$$Y = X\beta + \varepsilon \quad (1.2)$$

که در آن Y بردار $(n \times 1)$ مشاهدات و X ماتریس $(n \times (p+1))$ به صورت معلوم است. β بردار $(p+1 \times 1)$ پارامترها و ε بردار $(n \times 1)$ خطاهاست. به قسمی که $\varepsilon = E(\varepsilon) + I\sigma^2$ بنابراین مولفه‌های ε ناهمبسته‌اند. برآورد کمترین توانهای دوم β بردار b است که مقدار ε' را مینیمم می‌کند:

$$b = (X'X)^{-1} X'Y \quad (2.2)$$

اگر خطاهای مستقل باشند $\varepsilon_i \sim N(0, \sigma^2)$ آنگاه b برآورد حداقل درستنمایی β است. وضعیت مقادیر نقش مهمی در برازش کمترین توانهای دوم بازی می‌کند. در حالتی که همه نقاط در تعیین ارتفاع خط، وزن یکسانی دارند، شب خط قویاً تحت تاثیر مقادیر دور افتاده است. مشاهدات موثر را می‌توان به دور افتاده‌ها و نقاط نافذ تقسیم کرد.

دور افتاده‌ها: یک دور افتاده در بین مانده‌ها مانده‌ای است که از نظر قدر مطلق خیلی بزرگتر از بقیه است، و شاید به فاصله بیشتر از سه یا چهار برابر انحراف معیار از میانگین مانده‌ها قرار دارد. به عبارت دیگر عنصر دور افتاده معمولاً به داده‌ای اطلاق می‌شود که دارای مقدار منتهی‌الیهی برای متغیر پاسخ باشد.

نقطه نافذ: نقطه داده‌ای که مقدار منتهی‌الیهی برای یکی از متغیرهای توضیحی را داراست یک نقطه نافذ نامیده می‌شود. نقاط نافذ اثری نامناسب در برآوردهای ضرایب

رگرسیونی دارد، زیرا که زمانی نقاط نافذ وجود دارند این نقاط هیچ تبعیتی از بقیه داده‌ها در مدل نمی‌کنند و اثر نامطلوبی در مدل باقی می‌گذارند.

نقطه فروریزش: فرض کنید X عبارت است از یک نمونه n تایی (x_i, y_i) و همچنین برآورده رگرسیون، T باشد. فرض کنید $(m; T, X)$ ، سوپریمم $\|T(X') - T(X)\|$ برای تمام نمونه‌های آشفته X' باشد. وقتی این نمونه‌های آشفته، با جایگزین کردن هر m تا از نقاط اصلی نمونه با مقادیر دلخواه به دست آیند، آنگاه نقطه فروریزش T برای X عبارتست از:

$$\varepsilon^*(T, X) = \min\left\{\frac{m}{n}; \beta(m; T, X)\right\} \quad (3.2)$$

به بیان دیگر، کوچکترین مقدار اختشاش که سبب می‌شود برآورده، فاصله دلخواهی $T(X)$ داشته باشد نقطه فروریزش نام دارد. توجه کنید که این تعریف شامل هیچ توزیع احتمالی نیست. برای کمترین توان‌های دوم $\frac{1}{n} \varepsilon^*(T, X)$ است، چون یک مشاهده بد می‌تواند سبب فروریزش شود.

کارایی: نسبت میانگین مربع خطابرا کمترین توان‌های دوم به میانگین توان‌های دوم خطابرا هر رگرسیون استوار دیگر را کارایی گویند.

خاصیت کراندار بودن نقاط موثر: روش کمترین توان‌های دوم به هر مشاهده وزن یکسان می‌دهد در صورتی که روش‌های استوار قادرند به مشاهدات وزن‌های نابرابر اختصاص داده و اثر آن‌ها را در مدل کم کنند.

ما در اینجا به مقایسه روش‌های رگرسیون استوار با رگرسیون کمترین توان‌های دوم با سه خاصیت نقطه فروریزش، کارایی، خاصیت کراندار بودن نقاط موثر می‌پردازیم.

۳ رگرسیون حداقل قدر مطلق انحرافات (LAD)

برآوردهای کمترین قدر مطلق انحرافات^۱ a و b ای هستند که عبارت زیر را کمینه می‌کنند:

$$\sum_i |y_i - (a + bx_i)| \quad (1.3)$$

تفاضل $y_i - (a + bx_i)$ از خط $\hat{y} = a + bX$ نامیده می‌شود.

هرگاه هر یک از موارد زیر رخ دهد روش رگرسیون LAD نسبت به روش رگرسیون LS ترجیح داده می‌شود:

^۱ Least Absolute Deviations

الف) خطاهای دارای توزیع متقارنی باشد که میانه آن برآورده کارتری از نمونه باشد. ب) خطاهای دارای توزیع دم پهن باشند. پ) داده دورافتاده در متغیر پاسخ داشته باشیم. ت) هم خطی چندگانه در بین متغیرهای مستقل وجود داشته باشد. ث)تابع زیان قدر مطلق خطاب مناسبتر از تابع زیان درجه دوم باشد.

خواص سه گانه

۱- چون در مقابل نقاط نافذ عاجز است لذا نقطه فروبریش برای نمونه‌های متناهی همچنان برابر $\frac{1}{n}$ است که با افزایش نمونه به سمت صفر می‌کند. ۲- با توجه به (۱) این روش از خاصیت کرانداری نقاط نافذ برخوردار نیست. ۳- اگر این خطاهای دارای یک توزیع لاپلاس باشند، آنگاه کارایی برابر $1 < \frac{\sigma}{\sigma} = ۰$ درصد است.

۴- برآوردهای M

فرض کنید که (u) تابعی تعریف شده از u است و s برآورد مقیاس باشد (نه لزوماً به روش LS). یک برآوردهای استوار گویند اگر عبارت:

$$\sum_{i=1}^n \rho\left(\frac{e_i}{s}\right) = \sum_{i=1}^n \rho\left\{\frac{Y_i - x'_i \beta}{s}\right\} \quad (1.4)$$

را کمینه کند. برای کمینه کردن رابطه فوق با مشتق گیری از این رابطه نسبت به پارامترهای $p = k + ۰, ۱, ۲, \dots, k, \beta_j$ معادله به شکل:

$$\sum_{i=1}^n x_{ij} \psi\left\{\frac{Y_i - x'_i \beta}{s}\right\} = ۰ \quad j = ۰, ۱, ۲, \dots, k \quad (2.4)$$

می‌رسیم. به دلیل عدم جواب صریح این معادلات، لازم است جوابها از روش تکراری محاسبه شوند. پس از محاسبه با استفاده از شکل ماتریسی داریم:

$$X'W_\beta X\beta = X'W_\beta Y \quad (3.4)$$

که در آن $W_\beta = \text{diagonal}(w_{1\beta}, w_{2\beta}, \dots, w_{n\beta})$ است، که

$$W_{i\beta} = \frac{\psi\{(Y_i - x'_i \beta)/s\}}{(Y_i - x'_i \beta)/s} \quad i = ۱, ۲, \dots, n \quad (4.4)$$

و در صورتی که دقیقاً $Y_i - x'_i \beta = ۰$ باشد آن را یک تعریف می‌کنیم. می‌توان جواب تکراری را به صورت زیر نوشت:

$$\hat{\beta}_{q+1} = (X'W_q X)^{-1} X'W_q Y \quad q = ۰, ۱, ۲, \dots \quad (5.4)$$

این روش تکراری را کمترین توان‌های دوم با تکرار تجدید وزن^۱ (IRLS) می‌نامند.

بزای استراتژی سایقمه مل ماعی ارde راوتسا دروآرد کی یا:

$$s = \text{median}|e_i - \text{median}(e_i)| / ۰/۶۷۴۵ \quad (۶.۴)$$

که وقتی n بزرگ و خطاهای به طور نرمال توزیع شده باشند تقریباً برآوردگری نااریب برای $sd(y_i) = \sigma$ است.

۵ رگرسیون کمترین میانه توان‌های دوم (LMS)

رسو^۱ (۱۹۸۴) در این روش مجموع توان‌های دوم مانده‌ها را با میانه توان‌های دوم مانده‌ها جایگزین کرد، در نتیجه این برآورد می‌تواند در مقابل اغتشاش ۵۰ درصدی داده‌ها مقاومت کند. برآورد کمترین میانه توان‌های دوم^۲ (LMS) به وسیله مینیمم کردن میانه^۳ e_i^2 ها یعنی

$$\underset{\beta}{\text{minimize}} \text{med}_i e_i^2 \quad (۱.۵)$$

به دست می‌آید. این پیشنهاد اساس ایده همپل^۴ (۱۹۷۸، صفحه ۳۸۰) است.

برای برآورد پارامتر مقیاس، σ ، آن را به صورت زیر برآورد می‌کنیم:

$$s = ۱/۴۸۲۶(۱ + \frac{۵}{n-p-۱}) \sqrt{\text{med}_i e_i^2} \quad (۲.۵)$$

که $۱/۴۸۲۶$ یک عمل تصحیح مجانبی برای حالت خطاهای نرمال است.

تعريف: مشاهده^۱ $(x_i, y_i) = (x_i, y_i, x_{i1}, \dots, x_{ip})$ متعلق به فضای خطی بردارهای سطری $1 + p$ هستند و پارامتر مجهول β یک بردار ستونی $1 + p$ بعدی $(\beta_0, \beta_1, \dots, \beta_p)$ است. گوییم مشاهدات در موقعیت کلی هستند، هرگاه هر $1 + p$ تای آن‌ها یک β منحصر بفردی را تضمین کنند. مثلاً در حالتی که $2 = p$ است جمله بالا به این معنی است که هر جفت مشاهده (x_{i1}, x_{i2}, y_i) و (x_{j1}, x_{j2}, y_j) نمی‌توانند هم خط باشند.

قضیه: اگر $1 > p$ و مشاهدات در موقعیت کلی باشند آنگاه نقطهٔ فروریزش روش LMS برابر است با

$$([\frac{n}{2}] - p + ۲)/n \quad (۳.۵)$$

^۱ Iteratively Reweighted Least Squares

^۲ Rousseeuw

^۳ Least Median of Squares

^۴ Hample

اغلب حد عبارت بالا را وقتی n به سمت بینهایت میل می‌کند در نظر می‌گیرند. بنابراین می‌توان گفت که (با ثابت) روش LS دارای نقطهٔ فروریزش صفر است. در حالی که نقطهٔ فروریزش روش LMS برابر با 5° درصد است.

خواص سه گانه

- ۱ - دارای نقطهٔ فروریزش 50% است.
- ۲ - از کارایی کمی برخوردار است.
- ۳ - دارای خاصیت کراندار بودن نقاط موثر نیست.

۶ رگرسیون کمترین توان‌های دوم پیراسته (LTS)

رسو(۱۹۸۴) روش کمترین توان‌های دوم پیراسته^۱ (LTS) را به صورت زیر معرفی می‌کند :

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^h (e^{\top})_{i:n} \quad (1.7)$$

که $(e^{\top})_{1:n} \leq \dots \leq (e^{\top})_{n:n}$ توان دوم مانده‌های مرتب شده است.

قضیه: نقطهٔ فروریزش روش LTS با $h = [n/2] + [(p+1)/2]$ برابر است با:

$$\varepsilon_n^* = ([n-p]/2) + 1/n \quad (2.6)$$

با انتخاب h تقریباً برابر $n/2$ نقطهٔ فروریزش به 50% می‌رسد. در روش LTS یک قاعده برای $\hat{\sigma}$ به صورت :

$$\hat{\sigma}_{LTS} = C \sqrt{\frac{1}{n} \sum_{i=1}^h (e^{\top})_{i:n}} \quad (3.6)$$

است که C یک عامل تصحیح نمونه متناهی و بزرگتر از یک است. در هر یک از موارد فوق زمانی یک مشاهده دورافتاده است که فقط و فقط نسبت $|e_i/\hat{\sigma}|$ بزرگ باشد.

خواص سه گانه

- ۱ - دارای نقطهٔ فروریزش 50% است.
- ۲ - کارایی LTS از LMS بیشتر و نرخ همگرایی آن سریعتر است.
- ۳ - دارای خاصیت کراندار بودن نقاط موثر است.

^۱ Least Trimmed of Squares

S-برآوردها

S-برآوردها را در سال ۱۹۸۴^۲ و ۱۹۸۷^۳ توسط رسو و یوهای^۲ و رسو و لروی^۳ به صورت زیر تعریف شده است:

$$\underset{\beta}{\text{minimize}} \quad S(e_1(\beta), \dots, e_n(\beta)) \quad (1.7)$$

و برآورده می‌باشد که $s(e_1(\beta), \dots, e_n(\beta)) = s(e_1(\hat{\beta}), \dots, e_n(\hat{\beta}))$ است. انحراف معیار $\hat{\sigma}_{S-est} = \sqrt{s(e_1(\hat{\beta}), \dots, e_n(\hat{\beta}))}$ به عنوان جوابی از:

$$\frac{1}{n} \sum_{i=1}^n \rho\left(\frac{e_i}{s}\right) = K \quad (2.7)$$

تعریف می‌شود و K اغلب برابر $E_\Phi[\rho]$ است که Φ تابع توزیع نرمال استاندارد است. تابع ρ با شروط زیر مشخص می‌شود: ۱) ρ متقارن و به طور پیوسته مشتق پذیر است. ۲) وجود دارد به طوری که ρ روی فاصله $[c, \infty)$ اکیداً افزیشی و روی فاصله $(-\infty, c]$ ثابت است. ۳)

$$\frac{K}{\rho(c)} = \frac{1}{2}$$

اگر بیش از یک جواب برای معادله (۲.۷) وجود داشته باشد آنگاه:

$$S(e_1, \dots, e_n) = \sup\{s; (1/n) \sum \rho(e_i/s) = K\} \quad (3.7)$$

اگر هیچ جوابی برای معادله (۲.۷) وجود نداشت آنگاه $S(e_1, \dots, e_n)$ را برابر صفر قرار می‌دهیم.

با توجه به مطالب فوق ($S(\beta) = S(e_1(\beta), \dots, e_n(\beta))$) یک نوع S-برآوردها استوار از انحراف معیار مانده هاست. به علاوه هدف S-برآوردها همانند M-برآوردها کمینه کردن انحراف معیار مانده‌ها است.

قضیه: یک S-برآوردها ساخته شده از تابع ρ که در شرایط (۱) تا (۳) صدق می‌کند دارای نقطهٔ فروریزش:

$$\varepsilon_n^* = ([n/2] - p + 2)/n \quad (4.7)$$

است، اگر برای هر نمونه $\{(x_i, y_i); i = 1, \dots, n\}$ در موقعیت کلی باشد.

همان طور که مشاهده می‌شود در حد نقطهٔ فروریزش به سمت $50/5 = 50$ درصد میل می‌کند.

^۲ Yohai

^۳ Leroy

خواص سه گانه

۱- دارای نقطهٔ فروپاش ۵۰٪ است. ۲- کارایی S-برآورد از LTS و LMS بیشتر است. ۳- دارای خاصیت کراندار بودن نقاط موثر نیست.

M-برآوردهای تعمیم یافته (GM) \wedge

مالوز^۱ (۱۹۷۵) توابع وزنی w را در روش هوبربه صورت زیر به کار می برد:

$$\sum_{i=1}^n w(x_i)\psi(e_i/\hat{\sigma})x_i = \circ \quad (1.8)$$

: متفرگ رژیم ریز تر وصفه^۴ با قوبل و مرف^۲ پوشش و رزان اوته مایه

$$\sum_{i=1}^n w(x_i)\psi(e_i/w(x_i)\hat{\sigma})x_i = \circ \quad (2.8)$$

که این M-برآوردهای تعمیم یافته^۳ تک تک مشاهدات موثر را کراندار می کنند.

قضیه: نقطهٔ فروپاش روش GM برابر

$$\varepsilon_n^* = \frac{1}{p+1} \quad (3.8)$$

که تابعی نزولی از تعداد متغیرهای توضیحی (p) است (مارونا، باستانس و یوهای ۱۹۷۹^۴). با توجه به تعریف GM در معادله (۲.۸) چند انتخاب ممکن را در جدول ۱ خلاصه کرده و در زیر یک مورد از آنها را به فرم الگوریتم در می آوریم. در ضمن خواص سه گانه در جدول خلاصه شده است.

جدول ۱: چند GM-برآورد خاص و خواص سه گانه آنها

^۱ Mallows

^۲ Schweppe

^۳ Generalized M-estimators

^۴ Maronna, Bustos and Yohai

اجزا	GM1	GM2	GM3	GM4	GM5	GM6
فرمول	شُووب	شُووب	شُووب	شُووب	شُووب	شُووب
برآورد اولیه	LTS	—M	—S	—S	—S	LS
برآورد مفیاس	$\hat{\sigma}_{LMS}$	$\hat{\tau}_k$	$\hat{\sigma}_{S-est}$	$\hat{\sigma}_{S-est}$	$\hat{\sigma}_{S-est}$	$\hat{\tau}_k$
w تابع وزن	$\min(1, b/RD^2)$	$1/ z $	$med z / z $	$med z / z $	$med z / z $	$w_{i,k}$
تابع ψ	هوبر	هوبر	هوبر	هوبر	توکی	هوبر
ثابت‌های تنظیم کننده	۱/۳۴۵	$1/0.5p^{1/2}$	۱/۳۴۵	۱/۳۴۵	۴/۶۸۵	$2\sqrt{(p+1)/n}$
خواص سه‌گانه						
کاریابی بالا	دارد	دارد	دارد	دارد	دارد	دارد
فروزیش بالا	دارد	ندارد	ندارد	دارد	دارد	دارد
خاصیت کرانداری	دارد	دارد	دارد	دارد	دارد	دارد

الگوریتم GM1 در $K = ۰$ امین تکرار برآورد LTS پارامترها را محاسبه کرده، سپس به صورت زیر عمل می‌کنیم. ۱- مانده‌های $e_{i,k}$ را یافته و با استفاده از رابطه (۲.۶) مقدار $\hat{\sigma}_{LMS}$ را بیابید.
۲- وزن‌های :

$$w_{i,k} = \min\left\{1, \left[\frac{b}{(x_i - m_x)'C^{-1}(x_i - m_x)}\right]^{a/2}\right\} \quad (4.8)$$

را به دست آورید، که در آن $a = ۲, b = \chi_{0.95}^2(p+1)$ و $m_x.a = ۲, C = MVE$ برآورد کننده‌های مکانی و کواریانس همراه با متغیرهای کمکی دریبضی گون با کمترین حجم (MVE) هستند که می‌توانند با تابع $Splus$ $scov.mve()$ محاسبه شوند. توجه کنید که در جدول ۱ داریم $.z_i = ۱ - h_i$ و $RD^2 = (x_i - m_x)'C^{-1}(x_i - m_x)$ ۳- وزن‌ها را استفاده کرده و برآوردهای جدید را با استفاده از رابطه زیر محاسبه کنید:

$$\hat{\beta}_{k+1} = \hat{\beta}_k + (X'BX)^{-1} X'W \psi\left(\frac{e_{i,k}}{w_{i,k} \hat{\sigma}_{LMS}}\right) \hat{\sigma}_{LMS} \quad (5.8)$$

که در آن $B = diag(w_i), W = diag(w_i)$ است. ۴- گام‌های ۱، ۲ و ۳ را تا همگرایی تکرار کنید (آنقدر ادامه دهید تا پارامترها در دو تکرار متوالی به هم نزدیک باشند).

خلاصه: جدول ۲ مقایسه‌ای از روش‌های رگرسیونی استوار را با سه مشخصه ذکر شده خلاصه می‌کند.

جدول ۲: خواص سه‌گانه روش‌های رگرسیون استوار

شیوه	ε^*	روش محاسبه	معیار
LS	۰٪	ساده	بهترین برآوردهای ناریب خطی
برآورد -M	۰٪	تکراری	کارا
GM	p با افزایش به سمت صفر میل می کند	بالا ولی با افزایش به توابعی وزنی از xها	تکراری کراندار کردن نقطه نافذ و کارا
LMS LTS برآورد -S	ثابت و ۵۰٪ به سمت میل می کند	تکنیک های تعاقب تصویری	نقطه فروبریزش بالا

۹ کاربردهای عملی

مثال ۱ (وزن مخصوص چوب): بیست نمونه از ورقه های نازک مقاطع عرضی چوب کاج به ضخامت 30 cm تهیه شده اند. این مقاطع در مایع کلروزول رنگ شده و مقادیر زیر تعیین شده اند.

جدول ۳: داده های مربوط به عاملهای آناتومی و وزن مخصوص برشهای چوب کاج (درایر و اسمیت (۱۹۸۷))

وزن مخصوص چوب	تاریخ تابستانه %	جذب نور چوب	جذب نور چوب	جذب نور چوب	تعداد رشته ها در یک میلیمتر مربع	تعداد رشته ها در یک میلیمتر مربع	یک میلیمتر مربع از چوب بهاره %
۰/۵۲۴	۰/۵۲۴	۸۴/۱	۵۲/۸	۴۶/۵	۱۰۵۹	۵۲	۵۲/۲
۰/۵۳۵	۰/۵۳۵	۸۸/۷	۵۴/۵	۵۲/۷	۱۲۵۶	۶۵۱	۶۵۱
۰/۵۷۰	۰/۵۷۰	۹۲/۰	۵۲/۱	۴۹/۴	۱۲۷۳	۶۰۶	۶۰۶
۰/۵۲۸	۰/۵۲۸	۸۷/۹	۵۰/۳	۴۸/۹	۱۱۵۱	۶۲۰	۶۲۰
۰/۵۴۸	۰/۵۴۸	۹۱/۵	۵۱/۹	۵۳/۱	۱۱۳۵	۵۴۷	۵۴۷
۰/۵۵۵	۰/۵۵۵	۹۱/۴	۵۵/۲	۵۴/۹	۱۲۳۶	۵۵۷	۵۵۷
۰/۴۸۱	۰/۴۸۱	۸۲/۴	۴۵/۵	۵۶/۲	۱۲۳۱	۴۸۹	۴۸۹
۰/۵۱۶	۰/۵۱۶	۹۱/۳	۴۴/۳	۵۶/۶	۱۵۶۴	۶۸۵	۶۸۵
۰/۴۷۵	۰/۴۷۵	۸۵/۴	۴۶/۴	۵۹/۲	۱۱۸۲	۵۳۶	۵۳۶
۰/۴۸۶	۰/۴۸۶	۹۱/۴	۵۶/۴	۶۳/۱	۱۵۶۴	۶۸۵	۶۸۵
۰/۵۵۴	۰/۵۵۴	۸۶/۷	۴۸/۱	۵۰/۶	۱۵۸۸	۶۶۴	۶۶۴
۰/۵۱۹	۰/۵۱۹	۸۱/۲	۴۸/۴	۵۱/۹	۱۳۳۵	۷۰۳	۷۰۳
۰/۴۹۲	۰/۴۹۲	۸۹/۱	۵۱/۹	۶۲/۰	۱۳۹۵	۶۵۳	۶۵۳
۰/۵۱۷	۰/۵۱۷	۸۸/۹	۵۶/۵	۵۰/۵	۱۱۱۴	۵۸۶	۵۸۶
۰/۵۰۲	۰/۵۰۲	۸۸/۹	۵۷/۰	۵۲/۱	۱۱۴۳	۵۲۴	۵۲۴
۰/۵۰۸	۰/۵۰۸	۹۱/۹	۶۱/۲	۵۰/۵	۱۳۲۰	۵۲۲	۵۲۲
۰/۵۲۰	۰/۵۲۰	۹۰/۴	۶۰/۸	۵۴/۶	۱۲۴۹	۵۸۰	۵۸۰
۰/۵۰۶	۰/۵۰۶	۹۱/۸	۵۳/۴	۵۲/۲	۱۰۲۸	۴۴۸	۴۴۸
۰/۵۹۵	۰/۵۹۵	۹۲/۹	۵۳/۲	۴۲/۹	۱۰۵۷	۴۷۶	۴۷۶
۰/۵۶۸	۰/۵۶۸	۹۰/۰	۵۶/۱	۴۲/۴	۱۰۵۷	۵۲۸	۵۲۸

با مشاهدهٔ خروجی Splus مدل رگرسیونی به روش LS به صورت زیر به دست می‌آید:

$$\hat{Y} = 0/4421 + 0/0053X_3 - 0/0018X_4 + 0/0044X_5 \quad (1.9)$$

جهت تشخیص نقاط موثر می‌توان از ماتریس کلاه‌دار H، آماره کوک و همچنین باقیمانده‌های استاندارد شده می‌توان استفاده کرد. مقادیر فوق در جدول زیر خلاصه شده است. رسو وون‌زومرن^۱ (۱۹۹۰) استفاده از برآوردکننده‌های LMS را برای تشخیص دورافتاده‌ها پیشنهاد کرده‌اند. با استفاده از روش‌های فوق مشاهدات ۴، ۶، ۸ و ۱۸ را به عنوان یک نقطهٔ موثر نشان می‌دهد.

جدول ۴: مقادیر مانده‌های استاندارد شده و استیوینت شده و آماره کوک برای مشاهدات مشکوک

مشاهدات	h_{ii}	cook's	std.res	stud.res
۴	۰/۲۵۶۴	۰/۰۱۳	-۱/۳۲۷۸	-۱/۳۶۸۵
۶	۰/۰۹۴۲	۰/۰۹۱۸	۲/۳۰۰۴	۲/۸۱۰۶
۸	۰/۵۰۹۷	۰/۷۲۱۲	-۲/۰۴۰۲	-۲/۳۴۵۴
۱۸	۰/۲۲۴۴	۰/۰۷۳۱	-۱/۰۷۷۰	-۱/۰۸۳۷
min	۰/۰۹۴۲	۰/۰۰۰۳	-۲/۰۴۰۲	-۲/۳۴۵۵
max	۰/۵۵۷۴	۰/۷۲۱۲	۲/۳۰۰۴	۲/۸۱۰۶

جدول ۵: برآورد ضرایب مدل با استفاده روش‌های مختلف

روش	β_0	β_1	β_2	β_3	β_4	β_5
LS	۰/۴۴۲۱	۰/۰۰۰۱	۰/۰۰۰۰	-۰/۰۰۵۳	-۰/۰۰۱۸	۰/۰۰۴۴
LAD	۰/۳۸۷۱	۰/۰۰۰۲	۰/۰۰۰۰	-۰/۰۰۵۵	-۰/۰۰۳۸	۰/۰۰۶۱
$S - est$	۰/۳۳۶۶	۰/۰۰۰۱	۰/۰۰۰۰	-۰/۰۰۵۶	-۰/۰۰۵۱	۰/۰۰۷۵
LMS	۰/۳۵۳۷	۰/۰۰۰۱	۰/۰۰۰۰	-۰/۰۰۵۴	-۰/۰۰۴۹	۰/۰۰۷۲
LTS	۰/۴۱۱۶	۰/۰۰۰۲	۰/۰۰۰۰	-۰/۰۰۵۸	-۰/۰۰۳۷	۰/۰۰۵۷
GM1	۰/۴۱۲۶	۰/۰۰۰۲	۰/۰۰۰۰	-۰/۰۰۵۸	-۰/۰۰۳۲	۰/۰۰۵۵
GM6	۰/۳۹۹۹	۰/۰۰۰۲	۰/۰۰۰۰	-۰/۰۰۵۷	-۰/۰۰۳۶	۰/۰۰۵۸

انتخاب بهترین مدل

رسو (۱۹۸۴) پیشنهاد می‌کند در مسائلی که با متغیرهای متعددی مواجه هستیم یک برآوردگر سیار استوار شبه LMS و LTS را می‌توان برای پیدا کردن داده‌های دورافتاده به کار گرفت. در حالت کلی در مسائل عملی به نظر می‌رسد که بهتر است هم روش LMS و هم روش LS را در نظر بگیریم. اگر آن‌ها با هم توافق داشتند و نتایجشان با هم نسبتاً یکی بود آنگاه به نتیجه LS می‌توان اعتماد کرد. اما اگر بین آن‌ها اختلاف معنی داری

^۱ Van Zomeren

وجود داشت آنگاه باید به مطالعه باقیمانده هایی که از روش LMS به دست می آیند، داده هایی که بحران ایجاد کرده اند را پیدا کرد. پس از یافتن این داده های مؤثر، پیشنهاد می کنیم در M -برآورد گرها یا GM-برآورد گرها مختلف محاسبه شده، روشی مناسبتر و کارتر است که به این داده های مؤثر و دور افتاده وزن کمتری را نسبت دهد. به نظر میرسد از M -برآورد گرها توکی و GM1-برآورد گرها GM6 و GM1 بهترین باشند.

مثال ۲ (شبیه سازی): مدل زیر را در نظر گیرید:

$$y = 1 + x + \varepsilon \quad (2.9)$$

که در آن $(1, 1) \sim N(0, 1)$ و $x \sim N(0, 1)$ است. پس از محاسبه y با استفاده از مقادیر تولید شده x و ε در اولین و دومین مطالعه شبیه سازی ($n = 50, n = 100$) برای آزمودنی های دوم و سی و پنجم مقادیر x را به ترتیب با ۶ و ۶- تعویض می کنیم تا عملکرد GM_6 و GM_1 ، LS را مورد بررسی قرار دهیم. در مطالعه سوم شبیه سازی ($n = 15$) مقادیر x برای آزمودنی های ۲ و ۱۳ را به ترتیب با ۶ و ۶- تغییر می دهیم

جدول ۶: خلاصه برآورد پارامترهای مدل با استفاده از روش های LS و GM1 و GM6

روش	$n = 100$		$n = 50$		$n = 15$	
	β_0	β_1	β_0	β_1	β_0	β_1
LS	۰/۹۴۶	۰/۵۴۰	۱/۱۳۶	۰/۲۶۸	۱/۰۲۰	-۰/۱۶۴
GM1	۱/۰۲۳	۱/۰۹۸	۰/۹۸۹	۰/۸۱۰	۱/۱۷۸	۰/۹۵۹
GM6	۱/۰۳۶	۱/۰۵۱	۱/۱۲۰	۰/۶۶۳	۱/۶۷۷	۰/۰۳۷

همان طور که جدول نشان می دهد برای نمونه های بزرگ از آنجا که مقدار واقعی $\beta_0 = ۱$ و $\beta_1 = ۱$ است، GM1 و GM6 بهتر از LS هستند ولی برای نمونه های کوچک تنها GM1 مقادیری نزدیک به مقادیر واقعی پارامترها تولید می کند.

۱۰ تئیجه گیری

با توجه به مطالب بیان شده:

- ۱- اگر در داده ها نقاط موثر (دورافتاده و نافذ) وجود نداشته باشد روش LS از کارایی لازم برخوردار است.
- ۲- بهتر است ابتدا نتایج روش LS را با یکی از روش های LTS یا LMS مقایسه کرده اگر منطبق نبودند به دنبال روش های دیگر استوار برای به دست آوردن مدل می گردیم.
- ۳- اگر تنها در مشاهدات، داده دورافتاده وجود داشته باشد از روش M-برآوردها استفاده می کنیم.

۴- اگر در مشاهدات، علاوه بر دورافتاده‌ها نقاط نافذ نیز مشاهده شود—برآوردها کارترین روش هستند.

مراجع

- Draper, N. and H. Smith (1987). *"Applied Regression Analysis"*, John Wiley and Sons.
- Draper, N. and H. Smith (1998). *"Applied Regression Analysis"*, Second Edition, John Wiley and Sons.
- Hampel, F.R. (1978). *"Optimally Bounding the Gross-Error Sensitivity and the Influence of Position in Factor Space"*, Proceedings of The Statistical Computing Section of The American Statistical Association, Washington, D.c., PP.59-64.
- Huber, P.J.(1981). *"Robust Statistics"*, New York: John Wiley.
- Mallows, C.L. (1975). *"On Some Topics in Robustness"*, Unpublished Memorandum, Bell Telephone Laboratories: Murray Hill, NJ.
- Marronna,R. A., O. Bustos and V. Yohai(1979). *"Bias and Efficiency-Robustness of General M-estimators for Regression with Random Carriers"*, in Smoothing Techniques for Curve Estimation, eds. T.gasser and M.Rosenblatt, New York: Springer Verlag, pp. 91-116.
- Rousseeuw, P.J. (1984). *"Least Median of Squares Regression"*, Journal of The American Statistical Association December, Vol.79, PP.388,871-880.
- Rousseeuw, P.J., and A.J. Leroy (1987). *"Robust Regression and Outlier Detection"*. Wiley,New York.
- Rousseeuw, P.J., and V.J. Yohai (1984). *"Robust Regression by Means of S-estimators"*, Robust and Nonlinear Time Series Analysis, eds. J. Franke, W. H.Hardle, and D.Martin, Springer-Verlag: Heidelberg,Germany, pp. 256-272.

Rousseeuw, P.J., and B.C. Van Zomeren (1990). "Unmasking Multivariate Outliers and Leverage Points", Journal of the American Statistical Association, Vol. 85, pp. 633-651.

تقریب خطی اثرات متغیرهای کمکی تغییر پذیر در مدل کاکس توسط مدلهای پویا

علیرضا رضایی^۱، سیامک نور بلوجی^۲

A15226

^۱ مرکز آمار ایران

^۲ دانشگاه مینسوتا

چکیده: در مدل کاکس^۱ یک فرض این است که اثر متغیرهای کمکی^۲ در طول زمان ثابت است. در بسیاری از مسائل واقعی به نظر می‌رسد این فرض ضعیف است و بهتر است که این اثرات به عنوان تابعی از زمان در نظر گرفته شوند. مساله اصلی، برآورد یکتابع است که نمی‌تواند با استفاده از روش‌های معمولی با تعداد محدودی مشاهده برآورد شود. گامرمان^۳ (۱۹۹۱) طول کل زمان مطالعه را به N بازه افزار نمود سپس تقریب ساده $\beta(t)$ را در هر کدام از این بازه‌ها در نظر گرفت و آنها را با استفاده از مدل‌های پویا برآورد نمود. یک راه برای بهبود تقریب $\beta(t)$ استفاده از تقریب خطی به جای تقریب ساده در هر بازه می‌باشد. ما در این مقاله از تقریب خطی $\beta(t)$ در هر بازه استفاده کردیم و با استفاده از مدل‌های پویا این تقریبها را برآورد نموده‌ایم.

واژه‌های کلیدی: تحلیل بقا، تحلیل بیزی، فرآیند دیریکله، مدل خطر مناسب، مدل‌های پویا.

۱ مقدمه

در طول چهاردهه اخیر مطالعات زیادی روی آنالیز داده‌های بقا انجام گرفته است. از توابع بسیار مهم در آنالیز بقا تابع خطر است. که عبارت است از:

$$\lambda(t) = \lim_{\delta \rightarrow 0} \frac{P(t \leq T < t + \delta | T \geq t)}{\delta}$$

در بسیاری از مسائل اتفاق می‌افتد که در کنار داده بقای بیمار اطلاعات اضافی دیگری از قبیل سن، جنس، میزان فشار خون، سطح هموگلوبین خون و... وجود دارد که می‌توانند

^۱ Cox

^۲ Covariates variables

^۳ Gamerman

به تنهایی یا در کنار همیگر بر روی تابع خطر بیمار تاثیر داشته باشند. به این اطلاعات اضافی متغیرهای کمکی گفته می شود. کاکس (۱۹۷۲) مدل زیر را برای تابع خطر ارائه داد که وابستگی اندازه تابع خطر بیمار به مقادیر متغیرهای کمکی را نشان می دهد.

$$\lambda(t) = \lambda_0(t) \exp(Z' \beta) \quad (1)$$

$Z' = (Z_1, \dots, Z_p)$ بردار p بعدی متغیرهای کمکی، $\beta' = (\beta_1, \dots, \beta_p)$ بردار p بعدی اثرات متغیرهای کمکی و $\lambda_0(t)$ تابع خطر پایه^۱ است.

از زمانیکه کاکس مدل بالا را ارائه داد تا کنون کوششهای زیادی نسبت به بهبود این مدل و قویتر کردن فرضهایی که در آن به کار رفته انجام گرفته است. از فرضهایی که در این مدل شده این است که اثرات متغیرهای کمکی در طول زمان ثابت هستند اما در بسیاری از حالات به نظر می رسد این فرض بسیار ضعیف است و واقع بینانه تر این است که آنرا به عنوان یک تابع از زمان: $\beta(t)$ در نظر گیریم. بنابراین رابطه ۱ به صورت زیر تغییر می کند.

$$\lambda(t) = \lambda_0(t) \exp(Z' \beta(t)) \quad (2)$$

که $\beta(t) = (\beta_1(t), \dots, \beta_p(t))'$ بردار p بعدی اثرات متغیرهای کمکی است. حال در اینجا مسئله اصلی برآورد $\beta(t)$ که یک تابع است، می باشد که با روش‌های معمولی نمی توان آن را با تعداد محدود مشاهده برآورد نمود. در اینجا دو روش عمده برای برآورد وجود دارد. روش اول این است که تغییرات در طول زمان، وابسته به تعدادی پارامتر در نظر گرفته شود یا اصطلاحاً پارامتری شوند، آنگاه با برآورد پارامترها یک برآورد برای $\beta(t)$ حاصل می شود. روش دوم این است که زمان کل مطالعه به تعداد محدودی بازه افزایش شود، آنگاه $\beta(t)$ در این بازه‌ها تقریب شود و با بدست آوردن برآورده برای این تقریبها، برآورده برای $\beta(t)$ بدست آید. در روش اول فرضهای زیادی برای شکل تغییرات $\beta(t)$ مورد نیاز است، اما روش دوم با ساخت یک تقریب «خوب» مواجه است. گامرمان (۱۹۹۱) از تقریب ثابت $\beta(t)$ در هر کدام از این بازه‌ها استفاده نمود و این تقریب ثابت در هر بازه را با استفاده از مدل‌های پویا برآورد نمود. برای بهبود تقریب می توان از تقریب خطی استفاده نمود. ما در این مقاله از تقریب خطی $\beta(t)$ استفاده کردیم و آنها را در هر بازه برآورد نموده‌ایم.

^۱ Baseline hazard function

۲ تقریب خطی اثرات تغییر پذیر متغیرهای کمکی

در مدل کاکس می‌توان $\lambda_0(t) = \ln(\lambda_0(t))$ را به صورت $\beta_0(t)$ نوشت و Z و $\beta(t)$ را به صورت $\beta(t) = (\beta_0(t), \beta_1(t), \dots, \beta_p(t))^T$ و $Z = (1, Z_1, \dots, Z_p)$ تغییر داد. آنگاه رابطه ۲ به صورت زیر ساده‌تر می‌شود.

$$\lambda(t) = \exp(Z\beta(t)) \quad (2)$$

فرض کنید n داده بقا، t_1, t_2, \dots, t_n مشاهده شده‌اند که بعضی از آنها سانسور^۱ شده‌اند و کل زمان مطالعه با استفاده از نقاط τ_1, \dots, τ_N به صورت زیر افزایشده است: $(\tau_N \geq t_n)$

$$\begin{cases} I_1 = [\circ, \tau_1] \\ I_2 = (\tau_1, \tau_2] \\ \vdots \\ I_N = (\tau_{N-1}, \tau_N] \end{cases}$$

با این فرض که تغییرات $\beta(t)$ داخل هر کدام از این بازه‌ها تقریباً به صورت خطی است، داریم:

$$\beta(t) \sim \begin{cases} \underline{\theta}_{10} + \underline{\theta}_{11}t, & \tau_0 \leq t \leq \tau_1 \\ \underline{\theta}_{20} + \underline{\theta}_{21}t, & \tau_1 < t \leq \tau_2 \\ \vdots \\ \underline{\theta}_{N0} + \underline{\theta}_{N1}t, & \tau_{N-1} < t \leq \tau_N \end{cases} \quad (4)$$

که در آن $(\underline{\theta}_{i0}, \underline{\theta}_{i1}, \dots, \underline{\theta}_{ip})$ و $\underline{\theta}_{i0} = (\theta_{i01}, \theta_{i02}, \dots, \theta_{i0p})$ با $i = 1, \dots, N$ ، $i = 1, \dots, N$ ، $j = 1, \dots, n$ ، $j = 1, \dots, n$ توجه به روابط ۲ و ۳ تابع خطر زامین بیمار در i امین بازه برابر است با

$$\lambda_i^{(j)}(t) = \exp\{Z_j(\underline{\theta}_{i0} + \underline{\theta}_{i1})\}, \quad i = 1, \dots, N \quad j = 1, \dots, n \quad (5)$$

^۱ Censored

با تغییر متغیر $\gamma_{ij} = \exp\{Z_j \theta_i\}$ و $\rho_{ij} = \exp\{Z_j \theta_i\}$ رابطه ۵ به شکل زیر ساده تر می‌شود.

$$\lambda_i^{(j)}(t) = \rho_{ij} \gamma_{ij}^t, \quad i = 1, \dots, N \quad j = 1, \dots, n \quad (6)$$

بدون کاستن از کلیت مساله و برای سادگی در نمادگذاری می‌توان از اندیس j صرفنظر نمود و شکل کلی تابع خطر تجمعی برای بیماران را به صورت زیر بدست آورد.

$$\Lambda(t) = \begin{cases} \frac{\rho_1}{\ln \gamma_1} [\gamma_1^t - \gamma_1^{\tau_0}] & , \quad \tau_0 \leq t \leq \tau_1 \\ \vdots \\ \sum_{k=1}^i \left\{ \frac{\rho_k}{\ln \gamma_k} [\gamma_k^{\tau_k} - \gamma_k^{\tau_{k-1}}] \right\} + \frac{\rho_i}{\ln \gamma_i} [\gamma_i^t - \gamma_i^{\tau_{i-1}}] & , \quad \tau_{i-1} < t \leq \tau_i \quad (7) \\ \vdots \\ \sum_{k=1}^N \left\{ \frac{\rho_k}{\ln \gamma_k} [\gamma_k^{\tau_k} - \gamma_k^{\tau_{k-1}}] \right\} & , \quad t \geq \tau_N \end{cases}$$

با استفاده از روابطی که تابع بقا و خطر با یکدیگر دارند، شکل تابع بقای هر بیمار در بازه i ام برابر است با

$$S(t) = \exp\left\{-\sum_{k=1}^i \left\{ \frac{\rho_k}{\ln \gamma_k} [\gamma_k^{\tau_k} - \gamma_k^{\tau_{k-1}}] \right\} - \frac{\rho_i}{\ln \gamma_i} [\gamma_i^t - \gamma_i^{\tau_{i-1}}]\right\} \quad (8)$$

بنابراین اگر B_{i-1} پیشامد زنده ماندن تا زمان t_{i-1} باشد، آنگاه براحتی با استفاده از تعریف احتمال شرطی و رابطه ۶) می‌توان نتیجه گرفت که $S(t|B_{i-1})$ و $f(t|B_{i-1})$ برابرند با

$$\begin{cases} S(t|B_{i-1}) = \exp\left\{-\frac{\rho_i}{\ln \gamma_i} [\gamma^t - \gamma^{\tau_{i-1}}]\right\} \\ f(t|B_{i-1}) = \rho_i \gamma_i^t \exp\left\{-\frac{\rho_i}{\ln \gamma_i} [\gamma^t - \gamma^{\tau_{i-1}}]\right\} \end{cases} \quad (9)$$

۳ تجزیه زمانی تابع درستنمایی

با توجه به روابط بالا پارامترهای مجھولی که در هر بازه i ، $i = 1, \dots, N$ وجود دارند عبارتند از: $\theta_1, \theta_2, \dots, \theta_i$ واضح است که تابع خطر هر بیمار در هر بازه به مقدار آن در بازه های قبلی وابسته است، از اینرو مقادیر پارامترهای $\theta_1, \theta_2, \dots, \theta_i$ نیز به مقادیر پارامترهایی که

در بازه های قبل وجود دارند وابسته‌اند. در استفاده از مدل‌های پویا باید این وابستگی با یک رابطه بیان شود که اصطلاحاً به آن معادله سیستم گفته می‌شود. استنباط با استفاده از مدل‌های پویا شامل تلفیق اطلاعات نهفته از مشاهدات تحت عنوان تابع درستنمایی و اطلاعات موجود در معادله سیستم به عنوان پیشین و تشکیل تابع پسین می‌باشد، به این صورت که در هر مرحله به این طریق استنباطی برای پارامترهای $\rho_1, \rho_2, \dots, \rho_N$ بدست می‌آید، با استفاده از این اطلاعات به عنوان پیشین برای مرحله بعد و همچنین مشاهدات مرحله بعد و معادله سیستم، این استنباطها بهنگام می‌شود. حال مساله اصلی در این قسمت، بدست آوردن تابع درستنمایی در هر مرحله است زیرا تابع درستنمایی در کل زمان مطالعه شامل تمام پارامترهای مورد بررسی در هر بازه می‌باشد و نمی‌توان از آن برای انجام استنباط در n امین بازه $N = 1, \dots, n$ استفاده نمود. بنابراین باید تابع درستنمایی را در هر بازه تجزیه کرد به طوریکه تابع درستنمایی تجزیه شده در هر بازه شامل پارامترهای مجهول همان بازه باشد تا بتوان بر اساس آن استنباطی برای پارامترهای آن بازه بدست آورد و از آن برای مرحله بعد استفاده کرد. به این علت این نوع تجزیه، تجزیه زمانی نام گرفته است که تابع درستنمایی بین بازه‌هایی تجزیه شده است که زمان مطالعه با استفاده از آنها افزایش شده است. گام‌مان (۱۹۹۱) با در نظر گرفتن تقریب ساده $(t)^{\beta}$ درون هر بازه، تابع درستنمایی را درون آنها تجزیه نمود. ما برای تجزیه تابع درستنمایی به صورت زیر عمل نموده‌ایم که تقریباً مشابه روش تجزیه گام‌مان است.

لم ۱ : فرض کنید T_1, T_2, \dots, T_n متغیر تصادفی مستقل زمان بقای بیماران باشند که تابع بقای آنها در n امین بازه به صورت رابطه ۷ است. به علت سهولت در نوشتن، بدون کاستن از کلیات مساله فعلاً پارامترهای توزیع برای همه آنها یکسان و برابر با $i = \rho_i, \theta_i$ درنظر می‌گیریم. همچنین قرار دهید: $(\rho_1, \gamma_1, \dots, \rho_N, \gamma_N) = \mathcal{P}$. تابع درستنمایی $N = 1, \dots, n$ که از کل مشاهدات بدست می‌آید را می‌توان به صورت زیر بین بازه‌ها تجزیه نمود.

$$\begin{aligned} l(T|\mathcal{P}, D_{\circ}) &= \prod_{i=1}^n l_i(T^{(i)}|\mathcal{P}, D_{i-1}) \\ &= l_i(T^{(i)}|A, D_{i-1}) = \prod_{w \in F(I_i)} f(t_{iw}|\mathcal{P}, D_{i-1}) \\ &= \times \prod_{q \in R(I_i)} S(t_{iq}|\mathcal{P}, D_{i-1}), \quad i = 1, \dots, N \end{aligned}$$

که $w \in F(I_i)$ و $t_{iw} = t_w$ اگر $T^{(i)} = \{min(T_k, t_i), k \in F(I_i) \cup R(I_i)\}$ و $q \in R(I_i)$ و $t_{iq} = min\{t_q, t_i\}$ عبارت از کلیه اطلاعات موجود تا زمان i .

اثبات: روابط فوق با استفاده از تعاریف احتمال شرطی و قضیه بیز و روابط موجود بین تابعهای بقا، خطر و چگالی در آنالیز بقا بدست می‌آیند. برای اثبات کامل رجوع کنید به [۲۱].

ملاحظه می‌شود که تابع درستنمایی تجزیه شده در هر بازه شامل پارامترهای همان بازه می‌باشد.

۴ معادلات تکامل زمانی پارامترها (معادله سیستم)

همانطور که اشاره شد آنالیز یک مسئله با استفاده از مدل‌های پویا نیاز به مشخص بودن روابط زیر دارد:

- ۱) مدل نمونه‌گیری در هر بازه: عبارت است از تابع چگالی مشاهدات بدست آمده در هر بازه که همان عوامل تجزیه شده تابع درستنمایی داخل بازه‌های زمانی می‌باشد.
- ۲) معادله سیستم:

$$\begin{aligned}\underline{\theta}_0 i &= \underline{\theta}_{0,i-1} + \underline{\varepsilon}_i \\ \underline{\theta}_1 i &= \underline{\theta}_{1,i-1} + \underline{\varepsilon}_i\end{aligned}\quad (10)$$

که در آن فرض براین است که $\underline{\varepsilon}_i \sim N(\underline{\theta}_0, b_i \sigma_i I_p)$ و b_i عبارت از طول بازه i ام باشد.

ما برای بیان تکامل در طول زمان از یک مدل گام برداری تصادفی استفاده کردیم. زیرا در مسائل آنالیز بقا که در آنها اثرات متغیرهای کمکی در طول زمان ثابت نیست، مقدار یک اثر در زمان t به شدت به مقدار آن در زمان قبل از آن وابسته است. هر چه از زمان t جلوتر و یا عقب تر برویم این وابستگی کمتر می‌شود. از این‌رو به نظر می‌رسد که مدل گام برداری تصادفی برای در طول زمان مناسب باشد. هر چه طول یک بازه کوتاه‌تر باشد، پارامترهایی که در بازه بعد هستند، بیشتر به پارامترهایی که در این بازه هستند وابسته‌اند و اختلاف کمی بین آنهاست. در نتیجه واریانس ϵ_i کاهش می‌یابد و به همین ترتیب برای حالتی که طول بازه زیاد است، وابستگی پارامترها در بازه‌ها کمتر می‌شود و واریانس ϵ_i افزایش می‌یابد. برای نشان دادن این موضوع از ضربی b_i در توزیع $\underline{\theta}_0$ و $\underline{\theta}_1$ ها استفاده شده است.

همچنین در ابتدای مطالعه توزیع پیشینی که برای $(\underline{\theta}_0, \underline{\theta}_1) = \underline{\theta}_0$ در نظر می‌گیریم یک توزیع نرمال p بعدی با میانگین صفر و واریانس I_p است. واریانس

زیاد تقریباً حالت بی اطلاعی را می‌رساند، اما علت انتخاب صفر برای میانگین این است که در مسائل آنالیز بقا به خصوص در آزمون اثرات متغیرهای کمکی، همیشه فرض اولیه یا فرض صفر این است که متغیر کمکی بدون اثر است مگر این که خلاف آن مشاهده شود. از اینرو از این فرض برای انتخاب میانگین $\underline{\theta}$ استفاده کردۀ ایم. بنابراین می‌توان دریافت که مدل و فرضهایی که در بالا برای $\underline{\theta}$ شده است، منطقی می‌باشند.

۵ تحلیل روش پیشنهادی

با توجه به روابط (۹) به نظر می‌رسد که می‌توان با استفاده از روشی که وست و هریسون و میگان در سال ۱۹۸۵ ارائه دادند برای این حالت نیز استفاده نمود. در روش پیشنهادی آنها اولین فرض این بود که مدل نمونه گیری نمایی باشد، یعنی مشاهداتی که بدست آمده‌اند از متغیر تصادفی باشند که توزیع آن متعلق به خانواده نمایی باشد. مدل نمونه گیری در روش پیشنهادی ما نمایی نیست زیرا شکل کلی خانواده نمایی به صورت زیر است:

$$p(Y_t|\eta_t, \phi) = \exp[\phi\{Y_t\eta_t - a(\eta_t)\}]b(Y_t, \phi), \quad (11)$$

که در آن η_t پارامتر طبیعی توزیع است به طوریکه داریم: $E[Y_t|\eta_t, \phi] = \mu_t = a'(\eta_t)$ و ϕ پارامتر مقیاسی توزیع است که $\phi/a = a'/\phi$ مطابق رابطه (۱۱) در خانواده نمایی فقط یک حاصلضرب جدایی ناپذیر $Y_t\eta_t$ بین مشاهده و پارامتر طبیعی وجود دارد که به هیچ وجه نمی‌توان آنها را از هم جدا کرد، اما در رابطه (۸) دو حاصلضرب وجود دارد که آنها را نمی‌توان از هم جدا نمود و به همین دلیل نمی‌توان تابع (Y_t, ϕ) را تشکیل داد، در نتیجه نمی‌توان از روش وست و هریسون و میگان در این مساله استفاده نمود.

از طرفی اگر بخواهیم از روش‌های معمول بیز استفاده کنیم حجم بسیار زیادی از محاسبات پیچیده در مساله وارد می‌شود. ما این مساله را با استفاده از روش‌های تولید داده حل کردۀ ایم.

مولر^۱ (۱۹۹۹) روشی را برای انجام تحلیل در اینگونه مسائل ارائه داد. او فرض کرد مدل نمونه گیری مشاهدات و معادله سیستم به ترتیب به صورت زیر باشد:

$$y_t \sim p(y_t|\theta_t) : \text{مدل نمونه گیری مشاهدات}$$

$$\text{معادله سیستم} : \theta_t \sim g_t(\theta_{t-1}) + \omega_t$$

^۱ Muller

که در آن $(.)_t$ معلوم و لزومی ندارد نرمال باشد و $(.)_t$ یک تابع معلوم و ω_t یک برداری متغیر تصادفی است که لزوماً دارای توزیع نرمال نیست.

روش او مبتنی بر انتگرال مونت کارلو می‌باشد، به این صورت که یک نمونه مونت کارلو $\{\theta_1, \dots, \theta_n\} = A_1$ از توزیع پیشین که در مرحله اول برای پارامتر سیستم در نظر گرفته می‌شود $\pi(\theta_1|D_0)$ تولید می‌شود، آنگاه با استفاده از روش متروپلیس نمونه اولیه A_1 به یک نمونه مونت کارلو $\{\eta_1, \dots, \eta_n\} = B_1$ از توزیع پیشین $\pi(\theta_1|D_1)$ تبدیل می‌شود. با استفادهٔ مستقیم از معادلهٔ سیستم، B_1 تبدیل به یک نمونه مونت کارلو A_2 از توزیع پیشین $\pi(\theta_2|D_1)$ می‌شود، به این طریق که ابتدا یک نمونه به اندازه n از توزیعی که برای ω در نظر گرفته شده است تولید می‌شود $\{\omega_1, \dots, \omega_n\}$ ، آنگاه قرار داده می‌شود:

$$\theta_i = g(\eta_i) + \omega_i, \quad i = 1, \dots, n$$

و این عملیات آنقدر انجام می‌شود تا فرآید به نقطهٔ پایانی برسد. برای جزئیات بیشتر رجوع کنید به [۱۶]

در روش پیشنهادی مولاز الگوریتم متروپلیس برای تبدیل نمونه مونت کارلو پیشین به یک نمونه مونت کارلوی پسین استفاده شده است. در بخش زیر اشاره کوتاهی به این الگوریتم شده است.

۱.۵ الگوریتم متروپلیس

فرض کنید A_t یک نمونه مونت کارلو از توزیع پیشین در اختیار است، الگوریتم متروپلیس با این نمونه، یک زنجیر مارکف x_1, \dots, x_M تولید می‌کند به طوری که این زنجیر مارکف دارای توزیع تقریبی پسین است.

فرض کنید $(.)_t$ توزیع پسین $p(\theta|data)$ باشد و $\theta_j \in A_i$ ، $j = 1, \dots, n$ ، $i = 1, \dots, n$ باشند، y_1 امین نمونه تولید شده در i امین بازه از توزیع پیشین $\pi(\theta_i|D_{i-1})$ باشد. به ترتیب زیر یک زنجیر مارکف تشکیل می‌شود.

$$x_1 = \theta_i \quad (1) \text{ قرار دهید: } x_1 = \theta_i$$

$$y_1 \text{ را از توزیع کاندید } (y_1|x_1) \text{ تولید کنید.} \quad (2)$$

$$x_2 = x_1 \text{ با احتمال } \frac{p(y_1)}{p(x_1)} \text{ قرار دهید: } \theta = \min\{1, \frac{p(y_1)}{p(x_1)}\} \quad (3)$$

$x_2 = y_1$ با x_2 تولید شده به مرحله ۲ بروید و به همین ترتیب M بار ادامه دهید. x_M را به عنوان یک نمونه تولید شده از توزیع پسین در نظر گیرید.

توزیع کاندید توزیعی است که در شرط $g(y|x) = g(x|y)$ صدق کند و به این علت نام توزیع کاندید برای آن انتخاب شده است که چگالی بیشتری روی مقادیری از دامنه توزیع پسین که احتمال آنها بیشتر است دارد. در عمل می‌توان این توزیع را نرمال با میانگین x_i و واریانس توزیع پسین در نظر گرفت.

تنها مشکل در کاربرد این روش در مدل‌های پویا این است که شکل توزیع پسین برای انجام مرحله سوم الگوریتم فوق مشخص نیست، اما می‌توان آن را به صورت زیر برطرف نمود.

می‌دانیم $\pi(\theta) \cdot l(data|\theta)$ ، بنابراین اگر بتوان $\pi(\theta)$ را با استفاده از نمونه مونت کارلویی که از آن در اختیار است برآورد نمود، آنگاه استفاده از روش فوق برای تحلیل مدل‌های پویایی که برآوردهای پارامترهای مدل شکل بسته‌ای ندارند، مقدور می‌باشد. دقت کنید که ضریب تناسب در رابطه بیز در صورت و مخرج رابطه‌ای که در مرحله

۳ الگوریتم است ازین می‌رود.

برای برآوردتابع توزیع می‌توان از فرآیند دیریکله استفاده کرد.

۲.۵ فرآیند دیریکله

برآورد تابع توزیع با استفاده از فرآیند دیریکله به صورت زیر است:

اگر $P \in \mathcal{D}(\theta)$ آنگاه بازای هر $t \in Be(\alpha(-\infty, t], \alpha((t, \infty)))$ $F(t) \in \mathcal{D}(\theta)$ مخاطره بیز (مربع خطأ) در حالت بدون نمونه وقتی مینیمم می‌شود که

$$\hat{F}(t) = E(F(t)) = F_{\circ}(t) = \alpha((- \infty, t]) / \alpha(R)$$

واضح است که $F_{\circ}(t)$ حدس اولیه ما برای شکل تابع توزیع را نشان می‌دهد.

برای وقتی که یک نمونه به اندازه n داشته باشیم، قاعده بیز عبارت است از:

$$\begin{aligned} \hat{F}_n(t|x_1, \dots, x_n) &= \frac{\alpha((- \infty, t]) + \sum_{i=1}^n \delta_{x_i}((- \infty, t])}{\alpha(R) + n} \\ &= p_n F_{\circ}(t) + (1 - p_n) F_n(t|x_1, \dots, x_n) \end{aligned} \quad (12)$$

که $F_n(t|x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}((- \infty, t])$ و $p_n = \alpha(R) / (\alpha(R) + n)$ تابع توزیع شهودی نمونه است.

یک مسئله بسیار مهم تعیین مقدار $\alpha(R)$ است که اغلب در مسائل مجھول است. واضح است که نمی‌توان از $F_{\circ}(t)$ برآوردهای برآورد آورد. در مسئله برآورد تابع

توزیع در فرایند دیریکله، باید (t) و $\alpha(R)$ به طور مجزا انتخاب شوند. (رجوع کنید به فرگوسن [۱۰])

فرض کنید در تجزیه و تحلیل مدلی که برای تجزیه و تحلیل در نظر گرفتیم، در مرحله n ام هستیم و می خواهیمتابع توزیع نمونه مونت کارلویی را که از توزیع پیشین در این مرحله داریم برآورد نماییم. حدس اولیه و یا پیشین ما برای شکل تابع توزیع مشاهده ای که در این مرحله قرار دارد، یک توزیع نرمال با میانگین و واریانس توزیع پسینی است که در مرحله قبل بدست آمده است. اما برای انتخاب $\alpha(R)$ به ترتیب زیر عمل می کنیم.

تعریف می کنیم Z_n برایر است با تعداد مشاهده های مجزایی که در یک نمونه به اندازه n در دسترس است. آنگاه $E(Z_n) = M \times [\log((n + \alpha(R)) / \alpha(R))]$ هر چه n افزایش یابد $E(Z_n)$ نیز افزایش می یابد. در مسائلی که $\alpha(R)$ مجهول است از طریق Z_n استنباطهایی برای این مقادیر مجهول بدست می آورند. با توجه به این موضوع که تعداد مشاهدات یکسان در هر نمونه مونت کارلو A_t ، تقریباً است بنابراین مقدار Z_n برابر با n است و از طریق رابطه بالا و حل آن به روش عددی، برآورده برای M بدست می آید که از آن برای برآورد $\alpha(R)$ استفاده شده است. حال برآورد تابع توزیع در هر مرحله با توجه به رابطه (۱۲) حاصل می شود.

۶ مثال عددی

داده های جدول زیر را در نظر گیرید. که در آن تیمارها دو نوع روش برای برداشتن غده

تیمار	زمانهای بقا (روز)
روش ۱	۱۲۲، ۱۰۸، ۱۰۳، ۹۵، ۷۴، ۷۲، ۶۰، ۴۸، ۴۴، ۴۲، ۱۷ ، ۱۹۷، ۱۹۵، ۱۹۳، ۱۸۵، ۱۸۳، ۱۷۰، ۱۶۷، ۱۴۴ ، ۴۴۵، ۴۰۱، ۳۱۵، ۳۰۷، ۲۵۴، ۲۳۵، ۲۳۴، ۲۰۸ ، ۷۹۵، ۵۸۰، ۵۷۷، ۵۶۷، ۵۴۲، ۵۲۸، ۴۸۴، ۴۶۴ ، ۱۴۵۵، ۱۳۶۶، ۱۲۳۲، ۱۲۱۴، ۱۱۷۴، ۸۵۵ ، ۱۷۳۶، ۱۶۲۶، ۱۶۲۲، ۱۵۸۵
روش ۲	۳۰۱، ۲۶۲، ۲۵۰، ۲۱۶، ۱۸۲، ۱۲۵، ۱۰۵، ۶۳، ۱ ، ۳۸۲، ۳۸۲، ۳۸۰، ۳۵۸، ۳۵۶، ۳۵۴، ۳۴۲، ۳۰۱ ، ۵۲۴، ۵۲۲، ۴۹۹، ۴۸۹، ۴۶۰، ۴۰۸، ۳۹۴، ۳۸۸ ، ۷۸۶، ۷۷۸، ۷۴۸، ۶۷۶، ۵۷۵، ۵۶۹، ۵۶۲، ۵۳۵ ، ۱۴۲۰، ۱۲۷۱، ۱۲۴۵، ۹۷۷، ۹۶۸، ۹۵۵، ۷۹۷ ، ۱۶۹۴، ۱۶۹۰، ۱۵۵۱، ۱۵۱۶، ۱۴۶۰، ۱۴۲۰

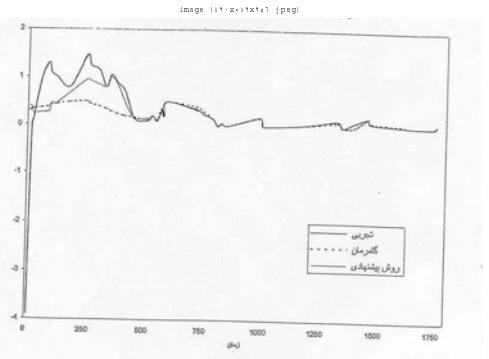
سرطانی می‌باشند. هدف مقایسه اثرات دو تیمار فوق روى زمان بقای بیماران سرطانی است. این کار را می‌توان با استفاده از مدل کاکس با یک متغیر کمکی که یکتابع نشانگر θ_0 است (۱ برای روش اول و ۰ برای روش دوم) انجام داد. دلایلی وجود دارد که اثر این متغیر کمکی در طول زمان ثابت نیست و بنابراین نمی‌توان از مدل معمولی کاکس استفاده نمود. کارترا^۱ (۱۹۸۳) فرض کرد که تغییرات این اثر در طول زمان به صورت $\theta_0 + \theta_1 t$ است و مقادیر θ_0, θ_1 را برآورد نمود. گامرمان (۱۹۹۱) تقریب ساده اثرات تیمارها را در نظر گرفت و با استفاده از مدل‌های پویا آنها را برآورد نمود. (رجوع کنید به گامرمان [۱۲]) ما با استفاده از روشی که ارائه دادیم این اثرات را به صورت خطی در هر بازه برآورد می‌نماییم.

ابتدا در شروع مطالعه یعنی در زمان صفر، توزیع پیشینی که برای $\theta_0 = (\theta_{011}, \theta_{012})$ و $\theta_1 = (\theta_{111}, \theta_{112})$ در نظر می‌گیریم، یک توزیع نرمال با میانگین صفر و واریانس I_4^{1000} است. زیرا فرض کرده ایم که $\theta_{111}, \theta_{112} = (\theta_{011}, \theta_{012})$ از یکدیگر مستقل هستند. البته این فرضی است که برای همه $\theta_i = (\theta_{0i}, \theta_{1i})$ ، $i = 1, \dots, N$ ها در نظر گرفته‌ایم.

اطلاعاتی که از مشاهدات درون هر بازه بدست می‌آید، یعنی همان عوامل تجزیه شده تابع درستنمایی در هر بازه را نیز به تک تک مشاهدات تقسیم کرده‌ایم، به این صورت که با توجه به استقلال متغیرهای تصادفی زمانهای بقا در بازه‌های زمانی، در بازه τ_i ام برای آنالیز بیزی از یک مشاهده استفاده می‌نمائیم، بنابراین باید مشاهدات درون هر بازه را مرتب نمود و این کار به این صورت انجام می‌گیرد که ابتدا زمانهای مرگ بیمارانی که در بازه τ_i ام زنده هستند، مقدار τ_i را قرار می‌دهیم. لازم است که نماد گذاری مناسبتری به کار ببریم. قرار می‌دهیم $D_{i-1, j}$ عبارت است از مجموع اطلاعات موجود تا قبل از زمان τ_i و اطلاعات حاصل از اولین مشاهده مرتب شده تا زامین مشاهده مرتب شده در بازه τ_i ام. همچنین p_{ij} و m_{ij} به ترتیب برابرند با میانگین توزیع پیشین و پسینی که از تابع درستنمایی زامین بیمار در بازه τ_i ام بدست می‌آید.

همانگونه که قبلاً اشاره گردید از میانگین پسین m_{ij} در هر مرحله برای برآورد توزیع پیشین در مرحله بعد استفاده می‌نمائیم تا از تمام توابع درستنمایی مشاهدات مرتب شده در بازه τ_i ام استفاده کنیم. آنگاه با توجه به معادله سیستم وارد مرحله (بازه τ_i+1) می‌شویم و همین‌طور ادامه می‌دهیم تا به آخرین بازه برسیم. برآوردهایی که برای θ_i ها بدست می‌آیند، میانگین توزیع پسینی هستند که از آخرین مشاهده در بازه τ_i ام بدست می‌آیند.

^۱ Carter



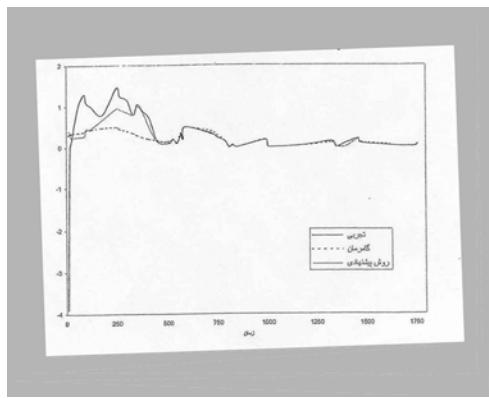
شکل ۱ : نمودار مقدار اثر متغیر کمکی در طول زمان و مقایسه روش‌های تقریب

در شکل صفحه بعد مقادیر برآورد پیشنهادی در این رساله و مقادیر برآورد گامرمان و مقادیری را که از قبل و از تجربه در اختیار است، مشاهده می نمایید.
همانطور که در شکل ملاحظه می نمایید مقادیر پیشنهادی نزدیک به مقادیر تجربی هستند و به وضوح می توان دید که تقریب اثر متغیر کمکی بهتر شده است.

مراجع

Antoniak, C. E. "Mixtures of dirichlet process with applications to Bayesian nonparametric problems". (1974) *Ann. Statist*, 2, 1152-1174.

Chang, Myron. N. and Chung Daehyun. "Isotonic window estimation of the baseline hazard function in Cox's regression model under order Restrict-



شکل ۲: نمودار مقدار اثر تغییر کمکی در طول زمان با استفاده از روش پیشنهادی

- tion". (1998) *Scand. J. Statist.* 25, 151-161.
- Chap, T. L. " Applied Survival analysis". (1997) New York: Wiely.
- Chatfield, Christopher. "Time series analysis". (1985) Boston: PWS-Kent Publishing.
- Collet, D. " Modelling survival data in medical research". (1994) London: Chapman and Hall.
- Cox, D. R and Oakes, D. " Analysis of survival data". (1984) London: Chapman and Hall.
- Dellaportas, P. and Smith, A. F. M. "Bayesian inference for generalized and proportional hazard Model via Gibbs sampling". (1993) *Appl. Statist.* 42, 443-459.

- Escobar, D. and West, M. "Bayesian density estimation and inference using mixture". (1995) *J. Am. Statist. Ass*, 90, 577-588.
- Escobar, D. "Estimating normal means with a dirichlet process prior". (1994) *J. Am. Statist Ass*, 89, 268-277.
- Ferguson, T, S. "A Bayesian analysis of some nonparametric problems". (1973) *Ann. Statist*, 1, 209-230.
- Ferguson, T, S. "Bayesian density estimation by mixtures of normal distributions". (1994) New York: Academic Press, pp. 287-302.
- Gamerman, D. "Dynamic Bayesian models for survival data". (1991) *Appl. Statist.*, 40, 63-79.
- Ghorai, J. K. "Nonparametric Bayesian estimation of a survival function under the proportional hazard model". (1989) *Commun. Statist*, 18, 1831-1842.
- Gilks, W. R. and Wild. P. "Adaptive rejection sampling for Gibbs sampling". (1992) *Appl. Statist*, 41, 337-348.
- Meinhold, R. J and Singpurwalla, N. D. "Understanding the Kalman filter". (1983) *Am. Statistn*, 37, 123-127.
- Muller, P. "Posterior integration in dynamic models". (1999) Technical Report, Duke University, Institut of Statistics and Decision Science.
- Smith, A. F. M. and Geleffand, A. E. "Bayesian statistics whitout tears: A sampling-resampling perspective". (1992) *Am. Statistn*, 46. 84-88.
- Tsai, W. Y. "Estimation of survival curves from dependent censorship models via a generalized self-consistent property with nonparametric Bayesian estimation application". (1986) *Ann. Statist*, 14, 238-249.
- West, M. and Harrison, P. J. and Migon, H. S. "Dynamic generalized linear models and Bayesian forecasting (with discussion)". (1985) *J. Am. Statist. Ass*, 80, 73-97.

West, M. "Bayesian kernel density estimation". (1990) Technical Report, 90-A02, Duke University, Institute of Statistics and Decision Science.

علیرضا رضایی، برآورد تقریب خطی اثرات تغییر پذیر در مدل کاکس با استفاده از مدل‌های پویا، رساله کارشناسی ارشد، دانشکده علوم ریاضی—گروه آمار، دانشگاه شهید بهشتی (۱۳۷۹).

روشهای احتمالاتی در حل مسائل دترمینیستیک

بیژن ظهوری زنگنه

P ۱۱۹۰۰

دانشکده علوم ریاضی، دانشگاه صنعتی شریف

چکیده: اثبات قضیه‌های احتمال بر اساس تکنیکهای آنالیز ریاضی را در اغلب قضیه‌های نظریه احتمال دیده‌ایم. در این سخنرانی قصد داریم که جریان معکوس این پدیده را یعنی کاربرد روشهای تصادفی را در حمله به مسائل آنالیز کلاسیک بررسی کنیم. یکی از ابتدایی‌ترین مثالهای این روشهای اثبات قضیه تقریب وایراشتراس به وسیله احتمالات است. این روشهای در حل مسائل نظریه پتانسیل، مسئله دیریشله و مسئله شرایط مرزی مارتین نیز کاربرد دارد. در این سخنرانی عمومی سعی خواهیم کرد با زبان شهودی و غیررسمی به بعضی از این کاربردها پردازیم.

واژه‌های کلیدی: روشهای احتمالاتی

احتمال به روایت تاریخ

نظریه احتمال در ابتدا و نیز پس از آن برای مدتی طولانی، عبارت بود از صورت آرمانی و تحلیل برخی پدیده‌های زندگی واقعی در خارج از حیطه ریاضیات، اما اندک در نیمه نخست این قرن احتمال ریاضی، بخشی معمولی از ریاضیات شد. تا قرن پانزدهم هیچگونه بررسی علمی در مورد پیشامدهای تصادفی انجام نشد. دانش پژوهان ایتالیایی لوکا با چولی (۱۴۴۵-۱۵۱۴)، نیکولا تارتالگیا (۱۴۹۹-۱۵۵۷)، چرولامو کارданو (۱۵۰۰-۱۵۷۱)، از جمله پیشکسوتان دانش ریاضی هستند که احتمالهای مربوط به بسیاری از بازیهای تصادفی را محاسبه کرده‌اند. به قول دیود مامفرد در مقاله طلوع عصر روشهای تصادفی «اگر جلوتر بیاییم، می‌بینیم که در عصر رنسان، کارданو شخصیت بی‌نظیری است. او به خاطر کتاب فن کبیرش (۱۵۴۵) غالباً مبدع خوانده می‌شود. ظاهراً وی یکی از خبره‌ترین افراد در زمینه عملیات صوری جبر بود به‌طوری که تبعات قواعد منطقی جبر را یک گام فراتر از اسلام خویش برد. ولی او در عین حال، معتقد به قمار هم بود و در کتاب «بازیهای شانسی» خود نخستین تحلیل را از قوانین شانس ارائه کرد، اما خیجالت می‌کشید آن را انتشار دهد و این کتاب تا ۱۶۶۳ به چاپ نرسید، یعنی تقریباً مقارن با زمانی که یاکوب برنولی کار خود را آغاز کرد.» (مامفرد، نشر ریاضی ۱۳). به هر حال پیشرفت واقعی

در فرانسه از سال ۱۶۵۴ آغاز شد، از وقتی بلر پاسکال (۱۶۲۳-۱۶۶۲) و پیردو فرما (۱۶۰۱-۱۶۶۵) دو ریاضیدان نامی نامه‌هایی به یکدیگر رد و بدل کردند که در این در مورد نامه‌ها روش‌های کلی محاسبه احتمالها بحث کرده‌اند. در سال ۱۶۵۵ داشمند معروف آلمانی کریستین هویگنس (۱۶۹۵-۱۶۲۹) به آنها پیوست، و این همکاری بسیار پر ثمر بود. در سال ۱۶۵۷ هویگنس اولین کتاب درباره احتمال را تحت عنوان «درباره محاسبات بازیهای شانسی» نوشت. این کتاب به منزله تولد واقعی احتمال محسوب می‌شود.

بعد از این تمام غول‌های ریاضی مانند برنولی، لاپلاس، پواسن و گاووس که استاد دقت در رشته‌های دیگر ریاضی بودند، قضیه‌های غلط و یا کم دقیقی را در احتمال ثابت کردند تا اینکه در سال ۱۹۰۰ در کنگره بین‌المللی ریاضیدانها در پاریس، دیوید هیلبرت (۱۸۶۲-۱۹۴۳) مساله را که به عقیده او حل آنها در پیشرفت ریاضیات موثر بود پیشنهاد کرد. یکی از این مسائل، بحث اصول موضوعی نظریه احتمال بود. در راستای رسیدن به این هدف کارهایی به وسیله امیل بورل (۱۸۷۱-۱۹۵۶)، و برنشتاین (۱۸۸۰-۱۹۶۸) انجام شد، تا اینکه در سال ۱۹۳۳ اندری کلموگروف (۱۹۰۳-۱۹۸۷) به صورتی موفق آمیز نظریه احتمال را اصل موضوعی کرد.

نظریه احتمال کلموگرف بر اساس نظریه اندازه لیگ (۱۹۰۲) بنیان گذاشته شد. کلموگرف در نخستین صفحات تکنگاشت مشهور خود درباره نظریه احتمال صراحتاً می‌گوید که متغیرهای تصادفی حقیقی مقدار، همان توابع اندازه‌پذیرند و امیدهای ریاضی انتگرال‌های آنها مع‌هذا، اندازه‌پذیری یکتابع حقیقی مقدار را تعریف می‌کند و وقتی که می‌خواهد امید ریاضی یک متغیر تصادفی را تعریف کند صاف و ساده نمی‌گوید که این مقدار برابر است با انتگرال متغیر تصادفی نسبت به اندازه احتمال مفروض، بلکه انتگرال را نیز تعریف می‌کند (دوب، نشر ریاضی ۱۲) این سنت در اغلب کتب نظریه احتمال باقی مانده است (به عنوان مثال رجوع شود به Chung 1977).

کلموگروف در تکنگاشت خود مطالب بسیار مهمی را مطرح کرد. او فضای احتمال، ساختن فرآیند تصادفی روی فضای بینهایت بعدی احتمال و امید شرطی نسبت به یک میدان سیگما‌بی را معرفی کرد. فهم قضیه توسعی کلموگرف روی فضای بینهایت بعدی مدت‌ها برای بسیاری از ریاضیدانها مشکل بود. جوزف دوب از بنیان‌گذاران احتمال در آمریکا می‌گوید:

(نویسنده به یاد می‌آورد که منظور کلموگروف از اندازه روی فضای تابعی را تا زمانی دراز پس از آنکه تکنگاشت وی را خوانده بود، در نمی‌یافته است) (دوب، نشر ریاضی ۱۲)

این رویکرد مورد توجه و توافق اغلب احتمال‌دانها قرار گرفت. تا جایی که خیلی از احتمال‌دانها احتمال را جزئی از آنالیز می‌دانند. دوب می‌گوید «برخی ریاضیدانان بر آنند که هرگاه با خواص تحلیلی احتمال و امید ریاضی سروکار داشته باشیم، موضوع بخشی از آنالیز است، ولی اگر با دنباله‌های نمونه‌ای و توابع نمونه‌ای سروکار داشته باشیم، موضوع عبارت است از احتمال، نه آنالیز. این مولفان در موقعیت جالب توجهی هستند از این رو که در نظر کردن به تابع دو متغیره $x(t, \omega) \rightarrow x(t, \omega)$ مثلاً در فرآیندهای تصادفی اگر خانواده توابع (\cdot, t, ω) هنگامی که t تغییر می‌کند مورد مطالعه باشد آن را آنالیز می‌خوانند، ولی اگر خانواده توابع (ω, \cdot, t) هنگامی که ω تغییر می‌کند مورد بررسی باشد آن را احتمال می‌نامند و قطعاً آنالیز به حساب نمی‌آورند. دقیقتر بگوییم، ایشان بحث پیرامون توزیعها و پرسشهای مربوطه را آنالیز می‌دانند، اما بحث‌های به زبان توابع نمونه‌ای را آنالیزنمی‌دانند. این دیدگاه در قول ذیل بیان شده است.

پروتر: ایتو (Ito) در سال ۱۹۴۴، با ارائه انتگرال‌ش که در آن فرآیندهای تصادفی انتگرال بودند، توانست بخش چند بعدی را با تکنیک‌های احتمالاتی محض مورد بررسی قرار دهد، که نسبت به روش‌های آنالیزو فلر بهتر است» (دوب، نشر ریاضی ۱۲).

در هر حال دوب این ایده را که احتمال مستقل از آنالیز باشد قبول ندارد و در آخر این مقاله با آوردن استدلالی درباره توابع رادماچر (Rademacher)، و قضیه پل لوی می‌نویسد: «دیگر بر عهده خواننده است که داوری کند کدامیک از این نتایج نظریه اندازه‌ای است و کدام یک احتمالاتی، و آیا اصلاً بیرون راندن احتمال ریاضی از قلمرو آنالیز فایده‌ای دارد، و اگر دارد، آیا نظریه اندازه را هم نباید بیرون راند؟» (دوب، نشر ریاضی ۱۲).

در مقابل رهیافت دوب، رهیافت دیگری است که به وسیله دیوید مامفرد مطرح می‌گردد که معتقد است که «رهیافت دیگر آن است که مفهوم «متغیر تصادفی» در مرکز توجه قرار بگیرد و همه کارها با انواع و اقسام دستکاری در متغیرهای تصادفی انجام شود». (مامفرد، نشر ریاضی ۱۳).

دیوید مامفرد در مقاله «طلوع عصر روش‌های تصادفی» می‌خواهد ریاضی جدیدی به وجود آورد او می‌نویسد «در رهیافت تقلیل‌گرا، متغیر تصادفی بر حسب اندازه تعريف می‌شود که خود بر حسب نظریه اعداد حقیقی تعريف می‌شود، و این را هم نظریه مجموعه‌ها تعريف می‌کند که خودش بر اساس حساب محمولات تعريف می‌شود. در عوض من می‌خواهم بگوییم که باید علی‌الاصول بتوان متغیرهای تصادفی را در مبانی منطق و ریاضیات ادغام کرد و به صورت‌بندی شفافتر و کاملتری از دیدگاه تصادفی رسید. من خودم هنوز صورت‌بندی کامل و قطعی از این قضیه ندارم».

در هر حال چه ما مانند «دوب» نظریه احتمال را جزیی از آنالیز بدانیم چه مانند «دیوید مامفرد» قصد داشته باشیم ریاضی جدیدی با تغییر در مبانی آن به وجود آوریم و یا مانند پروتر بخشی را آنالیز و بخشی را احتمال بنامیم، این مسئله مسلم است که نظریه احتمال کاملاً آغشته به آنالیز است و براساس آن توسعه یافته حال چه بخواهد در درون آنالیز باقی بماند و چه از آن خارج شود. روش‌های تصادفی دارای شهود ویژه خود است که در بقیه قسمت‌های آنالیز وجود ندارد. این شهود باعث مطرح شدن مسائل زیادی در احتمال گشته و منبع گسترش آن است. در اغلب مسائل برای اثبات قضیه‌های احتمال از آنالیز کمک می‌گیریم. در این سخنرانی قصد داریم برای اثبات مسائل کلاسیک آنالیز از روش‌های تصادفی کمک بگیریم و آنها را اثبات کنیم. بنابراین با نگاه یک احتمال دان به اشیاء آنالیز ریاضی می‌پردازیم.

نگاه احتمالاتی به آنالیز کلاسیک

بازه بسته $[1, 0]$ را در نظر بگیریم. می‌خواهیم به این بازه دترمینیستیگ با تعبیر احتمالی جان تازه‌ای بدیمیم.

امیل بورل در سال ۱۹۰۹ هر $x \in [0, 1]$ را به صورت بسط دودوئی

$$x = {}^{\circ}.x_1, x_2, \dots$$

نوشت که رقم x_j یا صفر است یا یک: این رقمها توابعی از x هستند. (اگر بازه $[1, 0]$ را با اندازه لبگ در نظر بگیریم که یک اندازه احتمال بر این بازه است، این تابعها به شکل معجزه آسایی متغیرهایی تصادفی می‌شوند که دقیقاً همان توزیعهایی را دارند که در محاسبه احتمالات پرتاب سکه به کار می‌روند. یعنی 2^{-n} برابر است با احتمال منسوب به این رویداد که در یک آزمایش پرتاب سکه نخستین n پرتاب دنباله معینی از شیر و خط به دست بدهد، و 2^{-n} همچنین طول کل (=اندازه لبگ) تعدادی متناهی بازه است که نقاط متعلق به آنها بسطهایی دودوئی با دنباله‌ای مشخص از صفرها و یکها در n جایگاه خاص دارند.) (دوب، نشر ریاضی ۱۲). با این دید بازه $[1, 0]$ فضای نمونه‌ای آزمایش برنولی است (رجوع شود به Adams, Guiliemin).

از بسط دودوئی بالا و آزمایش برنولی یعنی پرتاب مستقل بینهایت بار سکه، می‌توان نشان داد که $\frac{x_1 + x_2 + \dots + x_n}{n}$ به $\frac{1}{2}$ میل می‌کند. (ولی بیان قویتری از قانون اعداد بزرگ حکمی بود که بورل به دست آوردید طی یک برهان اشتباه و غیرقابل تصحیح- که این دنباله از میانگینها به ازای تقریباً هر x به $\frac{1}{2}$ میل می‌کند (با احتمال ۱). یک سال بعد فیر

(Faber) برهان درستی برای این حکم ارائه کرد و از آن هنگام برهانهای بسیار ساده‌تری هم به دست آمدند، [Billingsley, Breiman, Chung] و زنگنه [فرشه، حرمت بورل را نگه داشت «برهان بورل بیش از حد کوتاه است. در آن چندین استدلال میانی حذف شده است و نیز احکامی بدون برهان فرض شده‌اند» (دوب، نشر ریاضی ۱۲).]

با استفاده از قضیه قانون قوی اعداد بزرگ برای دنباله‌های برنولی می‌توان به اثبات ساده‌برنستین از قضیه واپاشتراس برای تقریب توابع پیوسته با چند جمله‌ای‌ها دست یافت.

گیریم $f(x)$ یک تابع پیوسته روی بازه $[0, 1]$ باشد. ثابت می‌کنیم f حد یکنواخت چندجمله‌ای‌های برنیشتین به صورت

$$B_n(x) = \sum_{k=0}^n f(k/n) C_n^k x^k (1-x)^{n-k}$$

است.

برهان. فرض کنیم X_1, X_2, \dots, X_n دنباله‌ای از متغیرهای تصادفی (i.i.d.) برنولی با

$$P\{X_i = 1\} = x, P\{X_i = 0\} = 1 - x$$

است و آنگاه $S_n = X_1 + X_2 + \dots + X_n$

$$\begin{aligned} E(f(\frac{S_n}{n})) &= \sum_{k=0}^n f(\frac{k}{n}) P\{S_n = k\} \\ &= \sum_{k=0}^n f(k/n) C_n^k x^k (1-x)^{n-k} \\ &= B_n(x). \end{aligned}$$

چون تابع f روی $[0, 1]$ پیوسته یکنواخت است و بنابراین برای هر $t > 0$ وجود دارد طوری که

$$|x - y| \leq \delta \Rightarrow |f(x) - f(y)| \leq t$$

اما تابع f چون پیوسته است بنابراین کراندار است در نتیجه یک M چنان وجود دارد که برای هر x

$$|f(x)| \leq M < \infty$$

با به کار بردن این نامساوی داریم

$$\begin{aligned}
 |f(x) - B_n(x)| &= \left| \sum_{k=0}^n f(x) C_n^k (\lambda - x)^{n-k} \right. \\
 &\quad \left. - \sum_{k=0}^n f(k/n) C_n^k x^k (\lambda - x)^{n-k} \right| \\
 &= \left| \sum_{k=0}^n (f(x) - f(k/n)) C_n^k x^k (\lambda - x)^{n-k} \right| \\
 &\leq \sum_{k=0}^n |f(x) - f(k/n)| C_n^k x^k (\lambda - x)^{n-k} \\
 &= \sum_{\{k:|k/n-x|\leq\delta\}} |f(x) - f(k/n)| C_n^k x^k (\lambda - x)^{n-k} \\
 &\quad + \sum_{\{k:|k/n-x|>\delta\}} |f(x) - f(k/n)| C_n^k x^k (\lambda - x)^{n-k} \\
 &\leq \epsilon \sum_{k=0}^n C_n^k x^k (\lambda - x)^{n-k} \\
 &\quad + 2M \sum_{\{k:|k/n-x|>\delta\}} C_n^k x^k (\lambda - x)^{n-k}
 \end{aligned}$$

حال چون برای $0 \leq k \leq n, n \geq 1$ ، جرم احتمال در نقطه k احتمال دو جمله‌ای

$$P_n(k) = C_n^k x^k (\lambda - x)^{n-k}$$

است، بنابراین

$$\sum_{\{k:|k/n-x|>\delta\}} P_n(k) = P\left\{\left|\frac{S_n}{n} - x\right| > \delta\right\}$$

در نتیجه

$$|f(x) - B_n(x)| \leq \epsilon + 2M P\left\{\left|\frac{S_n}{n} - x\right| > \delta\right\} \quad (*)$$

اما چون برای هر متغیر تصادفی ξ داریم

$$P\{|\xi - E(\xi)| \geq \epsilon\} \leq \frac{V(\xi)}{\epsilon^2}$$

و در حالت $E(\xi) = x, \xi = \frac{S_n}{n}$

$$Var\left(\frac{S_n}{n}\right) = \frac{Var(S_n)}{n^2} = \frac{nx(\lambda - x)}{n^2} = \frac{x(\lambda - x)}{n}$$

بنابراین

$$P\left\{\left|\frac{S_n}{n} - x\right| > \delta\right\} \leq \frac{Var\left(\frac{S_n}{n}\right)}{\delta^2} = \frac{x(\lambda - x)}{n\delta^2} \leq \frac{1}{4n\delta^2}$$

در نتیجه داریم

$$\begin{aligned}|f(x) - B - n(x)| &\leq \epsilon + 2MP\left\{\frac{S_n}{n} - x\right\} > \delta\} \\&\leq \epsilon + 2M\frac{1}{\gamma n \delta^2} \\&= \epsilon + \frac{M}{\gamma n \delta^2}\end{aligned}$$

و بنابراین

$$\lim_{n \rightarrow \infty} \max_{0 \leq x \leq 1} |f(x) - B_n(x)| = 0$$

که نتیجه قضیه وایرشتراس است.

یکی از کاربردهای جالب احتمالات در آنالیز قضیه (Szegő) است برای اثبات آن می‌توان به (Adams, Guillemin صفحه ۱۶۷) رجوع کرد.

یکی از کاربردهای مهم احتمالات در آنالیز کلاسیک اثبات قضیه رادن-نیکودیم به وسیله مارتینگل‌ها است. چون ما معمولاً احتمال و امید شرطی را به وسیله قضیه رادن-نیکودیم تعریف می‌کنیم این نظریه مارتینگل‌ها است که بر اساس این قضیه استوار می‌گردد. اما اثبات قضیه رادن-نیکودیم به وسیله مارتینگل‌ها چنان‌هم غیرمنتظره نیست، برای دیدن این قضیه رجوع کنید به ظهوری زنگنه، نظریه فرآیندهای تصادفی.

آخرین مسئله دترمینیستیک که ما می‌خواهیم در این سخنرانی به آن پردازیم مسئله دیریشله است. این مسئله یکی از مسئله‌های مهم نظریه معادلات دیفرانسیل پاره‌ای بیضوی است و نظریه پتانسیل دیریشله به حساب می‌آید.

مسئله دیریشله

فرض کنیم D یک میدان در \mathbb{R}^1 باشد. فرض کنیم تابع f روی مرز D تعریف شده باشد می‌خواهیم ثابت کنیم یک تابع هارمونیک u روی میدان D چنان موجود است که $u|_{\partial D} = f$.

اثبات وجود این مسئله با روش‌های دترمینیستیک پیچیده است و تنها وجود جواب می‌توان را ثابت کرد. در صورتی که در راه حل احتمالاتی که ارائه می‌دهیم نه تنها وجود جواب را ثابت می‌کنیم بلکه فرمولی برای جواب نیز ارائه می‌دهیم. اثبات با روش‌های احتمالاتی این قضیه بر اساس آنالیز تصادفی و خواص حرکت برونی است. در بند بعد ما به طور شهودی و توصیفی درباره حرکت برونی و خواص آن صحبت می‌کنیم و سپس مسئله دیریشله را ثابت می‌کنیم. برای مطالعه عمیق حرکت برونی و خواص آن رجوع کنید به Revuz & Yor, Karatzas.

حرکت بروني و مسئله ديريشله

حرکت بروني نامی است که به حرکت نامنظم گرده گیاهان که در آب معلق هستند داده شده است. رابرت براون گیاهشناس معروف انگلیسی برای اولین بار در ۱۸۲۸ با مشاهده این حرکت، متوجه اهمیت آن در مطالعه ذرات معلق میکروسکوپی شد. پس از آن دامنه کاربرد حرکت بروني از مطالعه ذرات معلق میکروسکوپی بسیار فرازره است و شامل مدل‌سازی قیمت‌های سهام، نوافه حرارتی در مدارهای الکتریکی، برخی حالت‌های حدی در سیستم‌های صفحه و موجودی و اختلالات تصادفی در انواع دیگر از سیستم‌های فیزیکی، زیستی، اقتصادی و مدیریت شده است.

آنچه براون در ابتدا مشاهده نمود این بود که گرده‌های گیاهان درون مایع دارای حرکت‌اند و علاقه‌مند شد تا قانون و علت این حرکت را بیابد، اما از عهده این کاربرنیامد و مساله بدون پاسخ باقی ماند. سپس در سال ۱۹۰۶ میلادی، اینشتین موفق به حل مسئله شد و علت حرکت را بمباران دانه‌های گرده توسط ملکولهای مایع معرفی نمود. با این حال اولین مدل ریاضی حرکت بروني در تزدکتری ریاضی بشلیه (Bachelier) در سال ۱۹۰۰ میلادی در دانشگاه پاریس و برای مدل اقتصادی مطرح شد.

بشلیه توزیعهای مهم متعددی استخراج کرده بود که همگی به فرآیند حرکت براونی در R مربوط بودند، «از جمله توزیع مربوط به تغییر پیشینه در طول یک بازه زمانی. بدین منظور وی توزیعهای متناظر با یک قدم زدن تصادفی گستته را پیدا می‌کرد و سپس حد را هنگامی که طول قدمها به سمت صفر میل می‌کرد به دست می‌آورد. دقیقت بگویم، آنچه بشلیه استخراج نمود توزیعهایی بودند که برای فرآیند حرکت بروني کارایی داشتند، به فرض آنکه اصلاً چیزی تحت عنوان حرکت بروني وجود داشته باشد، و به فرض اینکه بشود آن را با آن قدم زدن‌های تصادفی تقریب زد.» (دوب، نشر ریاضی ۱۲)

پس از آن نوربرت وینر ریاضیدان برجمسته و نابغه قرن بیستم در سال ۱۹۱۸ مدل ریاضی این حرکت را به طور کامل بررسی کرد. «توجه کنید که شکی در وجود حرکت بروني نیست: حرکت براونی را می‌شود زیر میکروسکوپ نظاره کرد. ولی هنوز برهانی برای وجود یک فرآیند تصادفی، یک حساب ریاضی، با خواص مطلوب در دست نبود. وینر (۱۹۲۳) فرآیند مطلوب حرکت بروني را که امروزه گاه فرآیند وینر نامیده می‌شود ساخت. بدین منظور وی از رهیافت دانیل به نظریه اندازه استفاده کرد تا اندازه‌ای با خواص دیل برفضای S از توابع پیوسته به دست آورد: اگر (x, t) متغیری تصادفی باشد که با مقداریک تابع در S در زمان t تعریف شده باشد، فرآیند تصادفی این متغیرهای تصادفی فرآیندی تصادفی است با اعضای S به عنوان توابع نمونه‌ای، و با توزیعهای توأمی که

برای فرآیند حرکت برونوی داشتیم به عنوان توزیعهای توازن متغیر تصادفی» (دوب، نشر ریاضی ۱۲)

ریاضیدانان رحمت برای تولید یک نظریه عمیق ریاضی می‌کشید اما متأسفانه آنچه باقی می‌ماند نتیجه این نظریه به صورت یک مقاله است، همانطور که ۷۵ ریاضیدان در سال ۱۹۶۲ در بیانیه‌ای در مورد تدریس ریاضی اعلام کردند «تفکر ریاضی تنها استدلال استنتاجی نیست، همچنین اثبات صرف هم نمی‌باشد. فرآیندهای ذهنی و فکری که اثبات و چگونگی اثبات را ارائه می‌کنند همانند خود اثبات که نتیجه تفکر ریاضی است بخشی از تفکر ریاضی محسوب می‌شود. استخراج مفاهیم درست از وضعیت‌های محسوس و ملموس، تعمیم از حالات شهود، استدلال استقرایی، استدلال از طریق تمثیل، زمینه‌های شهودی که برای آشکار کردن یک حدسیه به کار می‌روند همگی سبک و طریقه ریاضی گونه تفکر است.» (دوب، نشر ریاضی ۱۲) خوشبختانه نوربرت وینر با نوشتن کتابهای توصیفی مانند «من یک ریاضیدان هستم» پشت پرده تفکر و فرآیند به وجود آوردن یک نظریه را تا اندازه‌ای بررسی کرده است، بنابراین داستان حرکت برونوی را از زبان نوربرت وینر می‌شنویم. وینر در M.I.T. استخدام شده بود «ساختمان M.I.T. در ساحل رودخانه چارلز ساخته شده بود و طوری قرار داشت که می‌شد مستقیماً واز پنجره‌های آن، از چشم انداز گسترش سرزمین زیبای دور و بر آن لذت برد، به خصوص وجود رودخانه، موجب شادی بود. به نظر می‌رسید که می‌توان از بام تا شام به تماشای ناز و کرشمه‌های عجیب و غریب آب نشست. ولی آن چه در میان این همه زیبایی‌ها به طرف خود می‌کشید، ریاضیات و فیزیک بود. آن قانون‌مندی‌های ریاضی، که همه این توده بی‌نظم و ناآرام آب را هدایت می‌کند، کدام است؟ مگر اهمیت اصلی ریاضیات در این نیست که می‌تواند نظم و ترتیبی را که زیر این هرج و مرج و نابسامانی ظاهر دور و بر ما پنهان شده است، پیدا کند؟ رودخانه چارلز، گاهی ناگهان از موج‌های بلند، با شانه‌های بلندکف، پوشیده می‌شود و گاه چنان چین خودگی ملایمی دارد که به زحمت می‌توان موج‌های کوتاه آن را دید. طول موج‌های آن، گاه از دو پا سه بند انگشت تجاوز نمی‌کند و گاه به چند متر می‌رسد. چگونه می‌توان بیان جزئیات این منظره غرق شویم؟ برایم روشن بود که باید استفاده کنیم تا در تنوع بی‌پایان جزئیات این منظره غرق شویم؟ برایم روشن بود که این مسئله، با مسئله میانگین آماری بستگی دارد که با انتگرال لیگ خویشاوند است.» (وینر، انتشارات فاطمی، صفحه ۴۱).

وینر در کتاب خود به آشنازی به آثار ویلارد «ویلارد گیبس، اشاره می‌کند: یکی از بزرگ‌ترین دانشمندان آمریکایی است که در واقع رشته تازه‌ای از دانش را پایه گذاشت، رشته‌ای که در حد فاصل فیزیک و ریاضیات قرار دارد» (وینر، انتشارات فاطمی، صفحه

(۴۲). «گیبس آثار بسیار جالبی هم در فیزیک و هم در ریاضیات دارد، ولی کارهای اساسی او در زمینه مکانیک آماری؛ بیش از هر چیز دیگری، برایم جالب بود، همین کارهای او بود که، تا حد زیادی، مسیر خاص زندگی مرا مشخص کرد.»

دیدگاه سنتی در فیزیک، که از نیوتون بزرگ سرچشمه می‌گیرد، بستگی خلل ناپذیری با تصویرهای دترمینیستیکی دارد و، بر طبق آن، معرفت دقیق چگونگی جهان و یا هر قسمت بسته‌ای از آن در یک لحظه معین، شامل معرفت دقیق آن در زمان‌های بعدی هم می‌باشد. بنابر تصویر اصلی نیوتون، اگر موقعیت و سرعت ذره‌ها را در موج‌های سطح رودخانه چارلز بدانیم، می‌توان حرکت این موج‌ها را در همه سده‌های آینده محاسبه کنیم. متأسفانه با وسیله‌های اندازه‌گیری که در اختیار داریم، و همه آن با دستهای آدمی ساخته شده‌اند، نمی‌توانیم مقادیر مطلقاً دقیق و سرعت همه ذره‌ها را، در لحظه اولین زمان، به دست آوریم. به این ترتیب، فیزیک، که باید عملأً به درد پدیده‌های طبیعت بخورد، ناگزیر مواجه با مشکلی می‌شود: چگونه می‌توان با تکیه بر این داده‌های تقریبی مربوط به وضع اولیه، که به کمک وسیله‌های موجود به دست می‌آید، درباره واقع امر قضاوت کرد؟» (وینر، انتشارات فاطمی، صفحه ۴۳).

با این ایده وینر فضای نمونه‌ای را فضای توابع پیوسته یا در واقع فضایی که هر عضو آن یک موج باشد در نظر گرفت. متغیر تصادفی از فضای توابع پیوسته $\Omega = C[0, T]$ به اعداد حقیقی تعریف شده بود. در اینجا وینر مفهوم مهمی را کشف کرده بود، توابعی که روی فضایی تعریف می‌شوند که دامنه آن نقطه (در فضای چندبعدی) نیست بلکه دامنه آن یک فضای بینهایت بعدی است. خود می‌گوید «به تعمیم مفهوم احتمال، در مواردی مربوط می‌شود که «حالات‌ای ممکن» را نمی‌توان به صورت نقاطهای یک صفحه یا حوزه‌ای از فضای دانظر گفت، ولی خصلت منحنی‌های را دارند که معرف اشیای متحرکی هستند.» (وینر، انتشارات فاطمی، صفحه ۴۴).

به این ترتیب، حرکت برونوی موقعیتی را در برابر ما قرار می‌دهد که در آن، ذره‌ها به رسم منحنی‌هایی مشغول‌اند، و این منحنی‌ها، به مجموعه‌ای آماری از منحنی‌ها تعلق دارند. «این حرکت، بهترین زمینه برای اندیشه‌های من در مورد به کار بردن انتگرال‌گیری لیگ در فضای منحنی‌ها بود. و ضمناً دارای این خصوصیت بود که موضوع آن، از لحاظ فیزیکی، به دنیای واقع مربوط می‌شد و دقیقاً به اندیشه‌های گیبس بستگی داشت. در واقع، در این جا بود که توانستم با به کار بردن نظرات خود در تعمیم نظریه انتگرال‌گیری، به موفقیت بزرگی برسم. خود حرکت برونوی، موضوعی نبود که در فیزیک، بدون بررسی باقی مانده باشد. ولی در کارهای اساسی و عمیقی که اینشتین و سمولوخوفسکی در این زمینه کرده‌اند، یا به رفتار یک ذره در یک لحظه زمانی ثابت پرداخته‌اند و یا به خصلت‌های

آماری مجموعه بزرگی از ذره‌ها در جریان زمان: ولی خاصیت‌های ریاضی خط سیر ذره‌های جداگانه، هیچ‌گاه مورد مطالعه قرار نگرفته بود). در مورد موضوع اخیر، تقریباً هیچ چیز روش نبود، البته اگر اظهار نظر عمیق یه رن فیزیک دان فرانسوی، را که در کتاب خود به نام «اتم‌ها» آورده است، به حساب نیاوریم. او می‌گوید) خط سیر به کلی بی‌نظم ذره‌ها، که در اثر حرکت بروندی به وجود می‌آید، آدمی را به یاد منحنی‌های پیوسته ریاضی دانان می‌اندازد که در هیچ نقطه خود مشتق نداشته باشند» (وینر، انتشارات فاطمی، صفحه ۴۸ و ۴۹).

«با کمال تعجب، و در عین حال لذت، دریافتم که با چنین درکی از حرکت براونی، می‌توان نظریه صوری آن را در حد بالایی از کمال و ظرافت تنظیم کرد. در چارچوب این نظریه توانستم ملاحظه یه رن را ثابت کنم که، به استثنای چند موردی که احتمال آن را در مجموع برابر صفر است، مسیرهای حرکت براونی، منحنی‌های پیوسته‌اند، که در هیچ کجا مشتق ندارند» (وینر، انتشارات فاطمی، صفحه ۴۹).

بدین ترتیب وینر در سال (۱۹۲۳) از خاصیت گاوی بودن حرکت براونی استفاده می‌کند و حرکت براونی را به صورت یک سری (فوریه) با پایه دنباله‌های توابع شاودر، و ضرایب متغیرهای تصادفی گوسی می‌سازد. حرکت براونی بعدها با اصلاحاتی به وسیله پل لوی در سال (۱۹۴۸) و چیشلسکی (۱۹۶۱) به صورت ساده تری ساخته می‌شود (برای این روش رجوع کنید به ۱۹۶۶ Lamperti). حرکت براونی یک فرآیند تصادفی است که دارای خاصیت مارکفی، گوسی، مارتینگلی، با نموهای ایستا و مستقل است. دارای مسیرهای پیوسته‌ای است که در هیچ نقطه‌ای مشتق ندارد.

حال به حل مسئله دیریشله می‌پردازیم. با توجه به این واقعیت که حرکت براونی دارای خاصیت قوی مارکف نیز می‌باشد، یعنی اگر B_t یک حرکت براونی و T یک زمان توقف باشد آنگاه فرآیند جدید $B_{T+t} - B_T$ یک حرکت براونی جدید است. حرکت براونی یک مارتینگل موضعی است یعنی اگر T یک زمان توقف باشد آنگاه $B_{T \wedge t}$ (که در اینجا t می‌نیم و T است) یک مارتینگل است.

فرض کنیم D یک مجموعه فشرده نسبی باشد، یعنی \bar{D} فشرده باشد. قرار می‌دهیم

$$\tau_D = \inf\{t : B_t \in D^c\}$$

یعنی زمان برخورد حرکت براونی با مرز مجموعه D یا اولین زمانی که حرکت براونی می‌خواهد از مجموعه D خارج شود. می‌توان ثابت کرد که τ_D یک زمان توقف است، توجه شود که τ_D متناهی است. یعنی حرکت براونی با احتمال ۱ از مجموعه D خارج می‌شود. اگر f تابعی باشد که روی D تعریف شده باشد آنگاه $f(B_t)$ فقط برای $t < \tau_D$

شده است.

قضیه: فرض کنیم f یک تابع هارمونیک در D باشد (یعنی $\Delta f = 0$ در D) و فرض کنیم $B_0 = x \in D$ (یعنی حرکت براونی در لحظه صفر در نقطه x باشد) آنگاه $f(B_t)$ یک مارتینگل موضعی است.

برهان: با استفاده از آنالیز تصادفی و فرمول ایتو & Karatzas, Oksendal, Revuz & Yor و ظهوری زنگنه] داریم

$$f(B_t) = f(x) + \int_0^t \nabla f(B_s) dB_s + \frac{1}{2} \int_0^t \Delta f(B_s) ds$$

چون f هارمونیک است، $\Delta f = 0$ و از آنجا

$$f(B_t) = f(x) + \int_0^t \nabla f(B_s) dB_s$$

چون انتگرال تصادفی یک مارتینگل است، نتیجه برقرار است. Q.E.D.

قضیه (مسئله دیریشله)

فرض کنیم D یک میدان باشد که در آن برای تمام x ها

$$P^x(\tau_D < +\infty) = 1$$

و گیریم f یک تابع مثبت بدل اندازه‌پذیر باشد. تابع $(h(x) := E^x(f(B_{\tau_D}))$ را در نظر می‌گیریم، اگر x موجود باشد که $h(x) < +\infty$ آنگاه $h(x)$ روی D هارمونیک است.

تبصره ۱. توجه کنید P^x و E^x ، احتمال و امید شرطی به شرط $\{B_0 = x\}$ (یعنی حرکت براونی در لحظه صفر در نقطه x است) باشد.

تبصره ۲. در اینجا فرمول صریح $h(x) := E^x(f(B_{\tau_D}))$ حل مسئله دیریشله را به ما می‌دهد.

برهان. گیریم $x \in D$ و گیریم V گویی به مرکز x باشد که $D \cap V \subset \bar{V}$ قرار می‌دهیم

$$\tau_v = \inf\{t : B_t \in \partial V\}$$

اگر $x, \tau_v < \tau_D$. لذا

$$\begin{aligned} E^x(f(B_{\tau_D})) &= E^x(E(f(B_{\tau_D}) | \mathcal{F}_{\tau_v})) \\ &= E^x(E(f(B_{\tau_D}) | B_{\tau_v})) \\ &= h(B_{\tau_v}) \end{aligned}$$

بنابراین $(h(x) = E^x(h(B_{\tau_v}))$. اما حرکت براونی نسبت به دوران ناوردا است، لذا روی ∂V یکنواخت است. حال

$$h(x) = \frac{1}{|\partial V|} \int_{\partial V} h(y) dy$$

در نتیجه برای هر گویی به مرکز x که در D واقع شود، اگر h متناهی یا نامتناهی باشد

$$h(x) = \frac{1}{|V|} \int_V h(y) dy$$

پس $(h(x)$ متناهی است اگر و فقط اگر h در یک همسایگی x متناهی باشد، که نشان می‌دهد مجموعه $\{x : h(x) < +\infty\}$ باز است. از همبند بودن D نتیجه می‌شود که اگر x یی در D موجود باشد که آنگاه برای هر y , $h(x) < +\infty$, $h(y) < +\infty$ ، لذا h هارمونیک است.

مراجع

- Adams, M., Guillemin, V., Measure Theory and Probability.
- Bass, R.F. Probabilistic Techniques in Analysis, Springer-Verlag 1995.
- Billingsley, P., Probability and Measure, New York: Wiley.
- Breiman, Probability, Siam 1992.
- Chung, K.L., A Course in Probability Theory, Academic Press, 1977.
- Kac, M., Statistical Independence in Probability, Analysis and Number Theory. Mathematical Association of America (Carus Mathematical Monograph, no 12) 1959.
- Karatzas, Shreve E., Brownian Motion and Stochastic Calculus Springer-Verlag New York 1988.
- Lamperti, J., probability W.A. Benjamin Inc 1966.
- OKsendal, B., Stochastic Differential Equations Springer-Verlag, Berlin.

Revuz, D., Yor, M., Continuous Martingales and Brownian Motion Springer
Berlin 1992.

در باب برنامه ریزی درسی دبیرستان، ترجمه جواد حاجی‌بابایی، رشد آموزش ریاضی
شماره ۴۶، ۱۳۷۴.

دوب، جوزف (ترجمه عطاءالله تقاء) سیر پیدایش دقت در احتمال ریاضی (۱۹۵۰-۱۹۰۰)،
نشر ریاضی، سال ۱۲، شماره ۱ و ۲.

ظهوری زنگنه، بیژن، نظریه فرآیندهای تصادفی، جزوایات درسی، دانشگاه صنعتی شریف.

ظهوری زنگنه، بیژن، آنالیز تصادفی، جزوایات درسی، دانشگاه صنعتی شریف.

مامفرد، دیوبد، ترجمه شاپور اعتماد، طلوع عصر روش‌های تصادفی، نشر ریاضی سال
۱۳، شماره ۱.

وینر، نوربرت، ترجمه پرویز شهریاری، من یک ریاضیدانم، انتشارات فاطمی، خرداد
۱۳۶۸.

توسعه روش سری زمانی فازی و ارائه یک مورد کاوی

رسول نورالسناء، عباس سقائی

P11018

دانشکده صنایع، دانشگاه علم و صنعت

چکیده: پس از ارائه مجموعه های فازی توسط دکتر زاده، کاربردهای متفاوتی برای آن تعریف گردیده است. پیش بینی داده هایی که ذاتاً قطعی نبوده باعث گردید تا اصول سری زمانی فازی برای پیش بینی فرایندهای پویایی که مشاهداتی زبانی دارند، توسط محققین مطرح گردد. روش ارائه شده مورد تحلیل آماردانان قرار گرفت و ایراداتی بر آن وارد شد. در این مقاله سعی شده است تا با توجه به ابهاماتی که روش مذکور داشته است، الگویی براساس تلفیق مدل های سری زمانی وابسته به زمان و مستقل از زمان برای مشاهدات زبانی مطرح گردد که بتواند ابهامات مربوطه را مرتفع نماید. در این مدل از شبکه های عصبی مصنوعی نیز استفاده گردیده است. بعلاوه جایگاه کاربردی مناسبی برای این تکنیک مشخص گردیده که در آن حیطه یک مورد کاوی نیز انجام پذیرفته است.

واژه های کلیدی: سری زمانی، مجموعه های فازی، داده های زبانی، شبکه های عصبی.

۱ مقدمه

تئوری فازی در سال ۱۹۶۵ برای اولین بار توسط زاده (۱۹۶۵) طی مقاله ای به نام «مجموعه های فازی» به عنوان ابزاری نوین برای مدل سازی عدم قطعیت معرفی شد. ایده مجموعه های فازی پاسخ ذهن خلاق زاده به ناتوانی روش های موجود در تحلیل رفتار سیستمهای پیچیده و خصوصاً سیستمهای جاندار و طبیعی بود. در سال ۱۹۶۲ چیزی را بدین مضمون برای سیستم های بیولوژیک نوشت «مالا ساساً به نوع جدیدی ریاضیات نیازمندیم، ریاضیات مقادیر مبهم یا فازی که توسط توزیع احتمالات قابل توصیف نیستند» پس از آن وی ایده اش را در مقاله «مجموعه های فازی تجسم بخشدید. با پیدایش تئوری فازی، بحث و جدلها پیرامون آن نیز آغاز گردید. بعضی ها آنرا تایید نموده و کار روی این زمینه جدید را شروع کردند و برخی دیگر آن را برخلاف اصول علمی می دانستند. بسیاری از مفاهیم بنیادین تئوری فازی به وسیله زاده در اواخر دهه ۶۰ و اوایل دهه ۷۰ میلادی مطرح گردید. نظریه مجموعه های فازی که شاید بتوان گفت به انگیزه ایجاد ابزاری برای مدل سازی عدم قطعیت در سیستم های پیچیده معرفی شده به تبیین نوعی دیگر از عدم قطعیت پرداخت

که شاید بهترین لغت برای توصیف آن «ابهام» باشد.

یکی از بخش‌های آمار پیش‌بینی است. وجود هرگونه عدم قطعیت موجب شد تا تکنیک‌های زیادی برای پیش‌بینی تدوین گردد. این علم به عنوان یکی از شاخه‌های اصلی در آمار کاربردی بطور وسیعی در سایر علوم مورد استفاده قرار گیرد. در این بین تجزیه و تحلیل سری‌های زمانی نقش بسیار مهمی را ایفا می‌کند. یک سری زمانی را مجموعه‌ای از مشاهدات می‌دانیم که هر یک متعلق به زمانی مشخص بوده و اغلب فاصله بین زمان مشاهدات متولی برابر می‌باشد. که این مقادیر با $t_1, Y_1, t_2, Y_2, \dots$ (از متغیر تصادفی Y) در زمانهای t_1, t_2, \dots نشان داده می‌شوند. لذا Y تابعی از t بوده و آنرا به صورت $Y=F(t)$ مشخص می‌کنند. آنچه در این مقاله مورد توجه قرار گرفته «سری زمانی فازی» است. بطوری که مشاهده می‌شود این عبارت در ارائه مدل‌ها و روش‌های مختلف مورد استفاده قرار گرفته است. در این مقاله سعی داریم تا تحقیقات موجود را که با عنوان «سری زمانی فازی» توان بوده‌اند را مشخص نموده و بر حسب موضوع مربوطه تقسیم بندی نماییم. و به تشریح هریک پردازیم. در بخش دوم مقاله سری زمانی فازی با داده‌های قطعی را بررسی می‌کنیم. و سری زمانی فازی با داده‌های زبانی^۱ را در بخش سوم مورد بررسی قرار می‌دهیم در این بخش روشی را توسعه داده و الگوریتمی برای آن پیشنهاد نموده‌ایم. در بخش چهارم مثالی را برای تشریح بیشتر روش پیشنهادی مطرح نموده و در خاتمه نتیجه گیری نموده‌ایم.

۲ سری زمانی فازی با داده‌های قطعی

روشهایی که در این بخش مورد توجه قرار می‌دهیم، روش‌هایی است که داده‌های سری زمانی، مقادیری قطعی هستند که به توسط مدل‌سازی فازی و سری زمانی مورد تحلیل قرار گرفته و پیش‌بینی مقادیر را ممکن نموده است. باید توجه کرد که مفاهیم سری زمانی فازی پس از توسعه رگرسیون فازی مطرح شده است. در این بین تاباکا(۱۹۸۷، ۱۹۸۲، ۱۹۹۲) رگرسیون فازی را مورد بررسی قرارداد و روشی را توسعه داد که محیط فازی را جایگزین خطای مدل نمود. این مدل پیش‌بینی را به توسط فاصله‌ای نشان می‌داد. یکی از ضعفهای این روش زمانی است که نقاط دورافتاده در مشاهدات وجود داشته باشند، و موجب وسیع شدن بیش از حد فاصله ارائه شده گردد. یکی از اولین روش‌های کاربردی پیش‌بینی فازی تحت عنوان رگرسیون فازی توسط واتادا(۱۹۹۲) مطرح شد. در این روش فاصله‌ای، امکان رخداد را نشان می‌دهد. این مدل براساس مدل رگرسیون امکان تدوین شده و یکی از ضعفهای آن این است که وزن توابع هدف بطور ذهنی تعیین می‌گردد. چن(۱۹۹۶)،

^۱ Linguistic Data

روش سری زمانی فازی را براساس مطالب سونگ و چیسوم (که در بخش بعد به تفضیل مطرح شده است) بنا نهاد. داده های مورد تحلیل قطعی بوده و از رابطهٔ فازی برای تشریح رابطهٔ بین داده ها استفاده گردید. این روش قادر بود پیش بینی را براساس یک مقدار (نه بصورت یک فاصله) ارائه کند. یکی از جدیدترین روشهای سری زمانی فازی را تسنگ و همکارانش (۱۱۰۰۲) ارائه نمودند. سریهای زمانی ARIMA که توسط باکس و جنکینز (۱۹۷۶) ارایه شده اند به عنوان تکنیکهای کاربردی در علوم اجتماعی، اقتصادی، مهندسی و... از جایگاه ویژه‌ای برخوردار است. تحقیقات تسنگ و همکارانش بر سری زمانی فازی که داده های قطعی را در مدلهای ARIMA مورد بررسی قرار می‌دهند معطوف گردیده است. علاوه براین در این روش دو مقدار «بهترین امکان» و «بدترین امکان» برای پیش بینی ارائه می‌گردد. و درنتیجه فاصله‌ای برای امکان رخداد مشخص می‌شود. ایشان این روش را FARIMA^۲ نامیدند.

یکی از دلایل اصلی که باعث شده است تا مدلهای سری زمانی فازی با داده های قطعی مورد توجه قرار گیرند، آن است که مدلهای سری زمانی (مانند ARIMA) جهت بررسی و پیش بینی به ۵۰ یا ۱۰۰ مشاهده برای تعیین مناسب مدل و پیش بینی مقادیر آن نیاز دارند. ولی در سری زمانی فازی، به تعداد کمتر مشاهده برای دقیقی یکسان نیاز داریم. خلاصه روشهای مورد بحث در جدول شماره^۱ آمده است.

جدول ۱ : مقایسه روشهای سری زمانی با داده های قطعی

Watada	Chen	FARIMA	ARIMA	Compare
ندارد	ندارد	دارد	دارد	تطابق‌بامدل باکس و جنکینز
تابع فازی	رابطهٔ فازی	تابع فازی	تابعی براساس مشاهدات قبل	ارتباط‌پروردی و خروجی مدل
فاصلهٔ امکان	برآوردهای نقطه‌ای	فاصلهٔ امکان	فاصلهٔ اطمینان	چگونگی پیش بینی
کمتر از ARIMA	کمتر از ARIMA	کمتر از ARIMA	حداقل ۵ نمونه برآورد خوب	تعداد نمونه‌های موردنیاز

^۱ Fuzzy ARIMA

۳ سری زمانی فازی با داده‌های زبانی

در زندگی روزمره اغلب کلماتی وجود دارند که برای توصیف متغیرها استفاده می‌شوند. به عنوان مثال هنگامی که می‌گوییم «امروز هوا سرد است.» ما از واژه سرد استفاده می‌کنیم. واضح است که متغیر «دماهی هوا امروز» می‌تواند مقادیری مانند ۱۰، ۵، ۷— یا ۲۰— را اختیار کند. هنگامی که متغیر، واژه‌ها را به عنوان مقدار می‌گیرد، در آنصورت چهارچوب مشخصی برای فرموله کردن آن در ریاضیات کلاسیک نداریم. برای ایجاد چنین چهارچوبی، مفهوم متغیرهای زبانی تدوین گردیده است.

۱.۳ نظریه موضوع

یک متغیر زبانی بوسیله چهار پارامتر (X, T, U, M) مشخص می‌گردد.

X : نام متغیر زبانی / دماهی هوا امروز

U : دامنه‌ای که متغیر در آن مقدار می‌گیرد / $U = [-20, 50]$

T : مجموعه نامهایی متغیر زبانی است. / خیلی سرد، سرد، متعادل، گرم، خیلی گرم

$T =$

M : یک قاعده است که هر مقدار زبانی در T را به یک مجموعه فازی در U مرتبط می‌سازد و غالباً به وسیله توابع تعلق نشان داده می‌شود.

پیش‌بینی مقادیر زبانی به کمک مشاهدات گذشته می‌تواند کاربرد وسیعی داشته باشد. این موضوع توسط محققینی مورد توجه قرار گرفته، که از مهمترین آنها می‌توان به تحقیقات سونگ و چیسوم (۱۹۹۳) اشاره نمود. ایشان برای اولین بار سری زمانی فازی با داده‌های زبانی را مطرح نمودند و شرایط استفاده از این روش را؛ پویا بودن فرایند، فازی بودن مشاهدات، قابل تعریف بودن دامنه مقادیر فازی در R^1 و غیر قابل حل بودن توسط مدل‌های سنتی سری زمانی در این روش دو نوع مدل مطرح شده است که عبارتند از مدل پایدار در زمان و مدل ناپایدار در زمان برشمردند. و تعاریف و قضایایی نیز برای هریک ارائه شده که برخی از آن در ذیل آورده شده است (برای اطلاعات بیشتر، رجوع شود به سونگ و چیسوم (۱۹۹۳-۱۹۹۴))

تعریف ۱: فرض کنید $(Y(t), t = \dots, 0, 1, 2, \dots)$ زیرمجموعه‌ای از R^1 است، به عنوان دامنه مجموعه‌های فازی $f_i(t)$ ($i = 1, 2, \dots$) تعریف گردد و دسته‌ای از $F(t)$ ($i = 1, 2, \dots$)، $f_i(t)$, باشد، لذا $F(t)$ را سری زمانی بر حسب $(Y(t), t = \dots, 0, 1, 2, \dots)$ می‌نامند.

تعریف ۲: فرض کنید $F(t)$ فقط توسط $F(t-1)$ ایجاد می‌گردد. و به صورت $-F(t)$

۱ و یا $F(t) = F(t-1) \circ R(t,t-1)$ نشان می‌دهیم. بطوری که $R(t,t-1)$ نشان‌هندۀ رابطه فازی بین $F(t)$ و $F(t-1)$ می‌باشد و \circ عملگر $\max\min$ است.

تعریف ۳: فرض کنید $F(t)$ فقط توسط $F(t-1)$ یا توسط $F(t-2)$ یا \dots $F(t-m)$ ایجاد می‌گردد. این رابطه فازی را می‌توان بصورت تساوی زیر نشان داد.

$$F(t) = (F(t-1) \cup F(t-2) \cup \dots \cup F(t-m)) \circ R(t,t-1)$$

تعریف ۴: فرض کنید $R(t,t-1)$ یک مدل مرتبۀ اول از $F(t)$ باشد. اگر به ازای هر t , $R(t,t-1)$ مستقل از t باشد، یعنی برای هر t تساوی $R(t,t-1) = R(t-1,t-2)$ برقرار باشد، آنگاه $F(t)$ یک سری زمانی مستقل از زمان و درغیر اینصورت یک سری زمانی وابسته به زمان نامیده می‌شود.

قضیه ۱: $F(t) = F(t-1)$ یک سری زمانی فازی است. اگر به ازای هر t تساوی $F(t) = F(t-1)$ برقرار باشد و همچنین $F(t)$ دارای تعداد عناصر محدود باشد آنگاه $F(t)$ یک سری زمانی مستقل از زمان نامیده می‌شود.

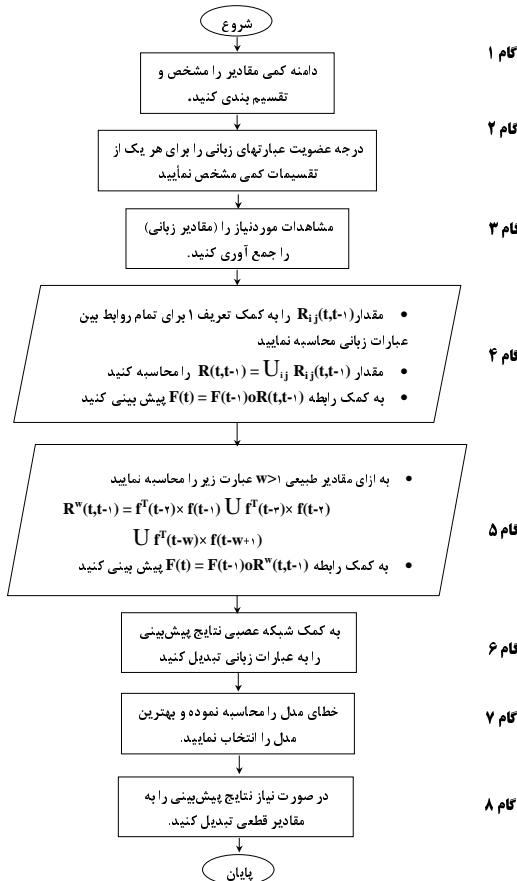
قضیه ۲: اگر $F(t)$ یک سری زمانی فازی باشد، $F(t) = F(t-1)$ و برای هر t و عناصر محدود $(f_i(t), i=1,2,\dots)$ وجود داشته باشد، خواهیم داشت:

$$R(t, t-1) = \dots f_{i1}(t-1) * f_{i0}(t) U f_{i2}(t-2) * f_{i1}(t-1) U \dots U f_{im}(t-m) * f_{im-1}(t-m+1) \dots$$

بطوری که $m > 0$ بوده و تمام زوجهای مجموعه‌های فازی متفاوت باشند.

این سری مقالات مورد بحث و توجه خاصی قرار گرفت که می‌توان به مقاله سولیوان و وودال (۱۹۹۴) اشاره کرد که ایراداتی رانیزمنعکس نموده بودند. مباحثی که مطرح شده باعث شد تا اصلاحاتی را در روش سونگ و چیسوم پیشنهاد نماییم و الگوریتمی جدید را که در شکل ۱ آمده است توسعه دهیم. کلیات ایراداتی که بر آن گرفته شده بود و در روش پیشنهادی اصلاح گردیده است به شرح ذیل می‌باشد.

* اصرار محققین برای استفاده از این مدل جهت پیش‌بینی داده‌های قطعی کارایی این مدل را در مقایسه با روش‌های آماری متعدد قرارداده بود. که در روش پیشنهادی ساختار مدل برای داده‌های فازی طراحی گردیده است. دو جنبه کاربردی برای آن پیشنهاد شده است که عبارتند از: الف) هنگامی که داده‌ها ذاتاً مقادیری از متغیرهای زبانی هستند یا بدلیل سهولت بیشتر یا هزینه کمتر اندازه گیری مقادیر زبانی جایگزین مقادیر قطعی می‌گردد. ب) در بسیاری از مواقع بدلیل مشخص نبودن مقادیر دقیق پیشین یک متغیر، نمی‌توان مدلی برای پیش‌بینی ارائه نمود. در صورتی که شاید بتوان مقادیر قبل را بطور



شکل ۱ : الگوریتم پیشنهادی

نادریق بیان کرد. و به کمک مدل پیشنهادی مقادیر آتی را پیش بینی نموده و به تدریج مقادیر دقیق را جمع آوری و تکمیل نمود.

* تعریف توابع عضویت در مقالات یکسان نبوده‌اند. که‌این مورد نیز مرتفع گردیده است.

* برای انتخاب و تشخیص مدل‌های ناپایدار و پایدار در زمان روش مناسبی پیشنهاد نگردیده بود که در روش پیشنهادی معیار انتخاب خطای کمتر مدل تعریف گردیده است.

* شبکه عصبی مصنوعی استفاده شده بطور نامناسب بکار بسته شده است. که‌این مورد نیز مرتفع گردیده است.

۴ مثال

در کارخانه‌ای سوابق تولید و فروش درسنوات قبل ثبت شده ولی سوابق تقاضای محصول ثبت نگردیده است. و مدیریت اکنون به این نتیجه رسیده که برای برنامه‌ریزی تولید و کنترل موجودی نیازمند پیش‌بینی تقاضای سال آتی می‌باشد. مطابق مدل‌های سنتی پیش‌بینی می‌باشد مشاهدات آتی تا چند دهه جمع‌آوری شده و پس از کفایت مشاهدات، مدل مناسب برآش گردد. ولی مطابق مدل سری زمانی فازی پیشنهاد شده در شکل ۱، گام‌های اشاره شده در ذیل برداشته شد.

گام ۱: ابتدا دامنه تقاضای سالانه بصورت $[13000, 20000] = U$ تعریف گردید. و به نواحی زیر تقسیم شد.

$$U_1 = [13000, 14000]$$

$$U_2 = [14000, 15000]$$

$$U_3 = [15000, 16000]$$

$$U_4 = [16000, 17000]$$

$$U_5 = [17000, 18000]$$

$$U_6 = [18000, 19000]$$

$$U_7 = [19000, 20000]$$

گام ۲: عبارتهای های زبانی A با هفت مجموعه فازی به عنوانی:

$(A_1 = \text{بسیار بسیار کم}), (A_2 = \text{بسیار کم}), (A_3 = \text{متوسط}), (A_4 = \text{بسیار زیاد}), (A_5 = \text{بسیار زیاد}), (A_6 = \text{بسیار بسیار زیاد}), (A_7 = \text{بسیار بسیار زیاد})$

در فضای U با درجه عضویت های ذیل تعریف گردیدند.

$$A_1 = \{U_1 / 1, U_2 / 0.5, U_3 / 0, U_4 / 0, U_5 / 0, U_6 / 0, U_7 / 0\}$$

$$A_2 = \{U_1 / 0.5, U_2 / 1, U_3 / 0.5, U_4 / 0, U_5 / 0, U_6 / 0, U_7 / 0\}$$

$$A_3 = \{U_1 / 0, U_2 / 0.5, U_3 / 1, U_4 / 0.5, U_5 / 0, U_6 / 0, U_7 / 0\}$$

$$A_4 = \{U_1 / 0, U_2 / 0, U_3 / 0.5, U_4 / 1, U_5 / 0.5, U_6 / 0, U_7 / 0\}$$

$$A_5 = \{U_1 / 0, U_2 / 0, U_3 / 0, U_4 / 0.5, U_5 / 1, U_6 / 0.5, U_7 / 0\}$$

$$A_6 = \{U_1 / 0, U_2 / 0, U_3 / 0, U_4 / 0, U_5 / 0.5, U_6 / 1, U_7 / 0.5\}$$

$$A_7 = \{U_1/0, U_2/0, U_3/0, U_4/0, U_5/0, U_6/0.5, U_7/1\}$$

گام ۳: در گام سوم براساس نظرات خبرگان کارخانه، میزان تقاضای سالهای قبل بطور ذهنی و به صورت عبارتهای زیانی تعیین گردید که در جدول ۲ آورده شده است. روابط احصا شده از جدول ۲ عبارتند از:

$$A_1 \rightarrow A_1$$

$$A_1 \rightarrow A_2$$

$$A_2 \rightarrow A_3$$

$$A_3 \rightarrow A_3$$

$$A_3 \rightarrow A_4$$

$$A_4 \rightarrow A_3$$

$$A_4 \rightarrow A_6$$

$$A_6 \rightarrow A_1$$

$$A_7 \rightarrow A_7$$

$$A_7 \rightarrow A_6$$

سال	مقادیر زبانی	نتایج پیش‌بینی (درجات عضویت)							مقادیر پیش‌بینی
		U۱	U۲	U۳	U۴	U۵	U۶	U۷	
۱۳۶۰	A _۱
۱۳۶۱	A _۱
۱۳۶۲	A _۱
۱۳۶۳	A _۲
۱۳۶۴	A _۲
۱۳۶۵	A _۲	۰,۵	۰,۵	۰,۵	۰,۵	۰,۵			A _۲
۱۳۶۶	A _۲	۰,۵	۰,۵	۱	۰,۵				A _۲
۱۳۶۷	A _۲	۰,۵	۰,۵	۱	۰,۵				A _۲
۱۳۶۸	A _۴		۰,۵	۱	۰,۵				A _۴
۱۳۶۹	A _۴		۰,۵	۰,۵	۰,۵	۰,۵			A _۴
۱۳۷۰	A _۴		۰,۵	۰,۵	۱	۰,۵			A _۴
۱۳۷۱	A _۲		۰,۵	۰,۵	۱	۰,۵			A _۲
۱۳۷۲	A _۲		۰,۵	۰,۵	۱	۰,۵			A _۲
۱۳۷۳	A _۲		۰,۵	۱	۰,۵	۰,۵			A _۲
۱۳۷۴	A _۲		۰,۵	۱	۰,۵	۰,۵			A _۲
۱۳۷۵	A _۲		۰,۵	۱	۰,۵				A _۴
۱۳۷۶	A _۴		۰,۵	۱	۰,۵				A _۴
۱۳۷۷	A _۶		۰,۵	۰,۵	۰,۵	۰,۵			A _۴
۱۳۷۸	A _۶					۰,۵	۰,۵	۰,۵	A _۶
۱۳۷۹	A _۷					۰,۵	۱	۰,۵	A _۶
۱۳۸۰	A _۷					۰,۵	۰,۵	۰,۵	A _۶
۱۳۸۱	A _۶					۰,۵	۰,۵	۱	A _۶
۱۳۸۲						۰,۵	۰,۵	۱	A _۶

جدول ۲. نتایج پیش‌بینی مدل انتخاب شده

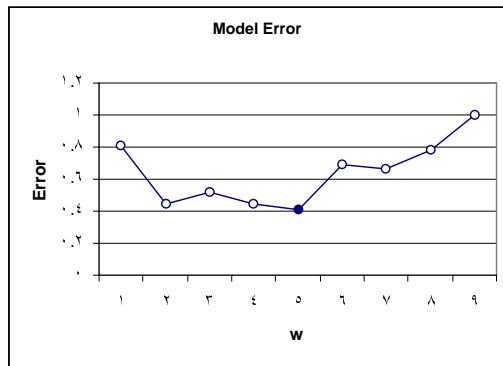
گام ۴: در این مرحله از الگوریتم، برای هر رابطه فوق ماتریس ($R_{ij}(t, t - ۱)$ محاسبه و در نهایت ($R(t, t - ۱)$ بدست می‌آید. و در نهایت مطابق رابطه ارائه شده در شکل ۱، پیش‌بینی مقادیر به صورت درجات عضویت انجام پذیرفت.

گام ۵: مطابق رابطه ارائه شده در شکل ۱، به ازای مقادیر $1 < w$ رابطه ($R^w(t, t - ۱)$ محاسبه گردید و مقادیر پیش‌بینی به صورت درجات عضویت انجام پذیرفت.

گام ۶: نتایج حاصل از هر دو روش گام ۴ و گام ۵ که به صورت مقادیر پیش‌بینی بر حسب درجات عضویت می‌باشند به کمک شبکه عصبی RBF^۱ به مقادیر زبانی تبدیل گردید. ورودی این شبکه عصبی از برداری با ۷ درایه که درجات عضویت تعريف شده در گام ۲ می‌باشند، تشکیل گردیده است و خروجی آنرا رده‌های متغیرهای زبانی تعريف شده در گام ۲ تشکیل می‌دهند. این شبکه به کمک نرم افزار MATLAB (5.3.1) آموزش دیده و جهت تبدیل درجات عضویت مقادیر پیش‌بینی شده به مقادیر زبانی مورد استفاده قرار گرفته است.

گام ۷: در این گام خطای مدل پیش‌بینی محاسبه می‌گردد. بدین ترتیب که ابتدا مقدار

^۱ Radial Basis Function Neural Network



شکل ۲: خطای مدل

رده پیش بینی بدست آمده از شبکه عصبی، به صورت عدد صحیح گرد می شود. سپس قدر مطلق اختلاف بین این مقادیر و مقادیر رده واقعی محاسبه می گردد. در نهایت با تقسیم مجموع کلیه مقادیر قدر مطلق بر تعداد پیش بینی انجام شده، میزان خطای مدل بدست می آید. براساس همین روش، مقدار خطای مدل گام ۴ و مدلها ی گام ۵(به ازای w های مختلف) محاسبه گردیده است. نتایج محاسبات به صورت یک نمودار در شکل ۲ نشان داده شده است. همانطور که در شکل مشاهده می شود کمترین مقدار خطای متعلق به مدل گام ۵(به ازای $w=5$) می باشد. نتایج عددی محاسبات مربوط به مدل مذکور شامل مقادیرزبانی واقعی، مقادیرزبانی پیش بینی شده و مقادیر درجات عضویت در جدول ۲ موجود می باشد.

گام ۸: در صورت نیاز می توان مقادیر پیش بینی را به صورت مقادیر عددی قطعی نیز تبدیل نمود. برای این منظور می توان میانه بخشی از دامنه (که در گام ۱ معرفی گردیده است)، و دارای بالاترین درجه عضویت است را به عنوان مقدار پیش بینی قطعی در نظر گرفت. (به عنوان مثال مقدار 1950° برای سال 1382) و یا در روشی دیگر، میانگین وزنی درجات عضویت و تقسیمات دامنه را محاسبه نماییم. (مقدار 1875° برای سال 1382)

اکنون مدیریت کارخانه مقدار تقاضای سال 1382 را به ۳ شکل؛ فازی، عبارت زبانی و مقدار قطعی پیش بینی نموده است.

۵ تیجه گیری

در این مقاله توانایی مدلسازی فازی را در سریهای زمانی مورد بررسی قرار دادیم و مشاهده نمودیم که در داده های قطعی می توان تعداد مشاهدات مورد نیاز را برای تخمین مدل و انجام پیش بینی در مقایسه با روش های معمول آماری کاهش داد. به علاوه توانستیم با بررسی روش های سری زمانی فازی که برای مشاهدات زبانی تعییه گردیده است روشی را توسعه دهیم که با ترکیب مدلسازی فازی و شبکه های عصبی، قادر است با دریافت مشاهدات زبانی مقادیر پیش بینی را به صورت فازی، عبارت زبانی و مقادیر قطعی ارائه نماید.

مراجع

- Cressie, N. (1993), *Statistics for Spatial Data*, Revised edition, John Wiley, New York.
- David, M. (1977), *Geostatistical Ore Reserve Estimation*, Eevier Scintific Publishing Co., 364 P.
- Box,G.P. Jenkins,G.M.(1976), *Forecasting and Control*, Holden-day Inc.,San Francisco, CA.
- Chen,S.M. (1996), *Forecasting enrollments based on fuzzy time series*, Fuzzy Sets and System, vol.81(3), 311-319 P.
- Tanaka,H. (1987), *Fuzzy data analysis by possibility linear models*, Fuzzy Sets and Systems, vol.24(3), 363-375 P.
- Tanaka,H. Ishibuchi,H. (1992), *Possibility regression analysis based on linear programming*, Fuzzy Regression Analysis, Omnitech Press, Heidelberg, 47-60 P.
- Tanaka,H. Uejima,S. Asai,K. (1982), *Linear regression analysis with fuzzy model*, IEEE Trans. System Man Cybernet.,vol.12(6), 903-907 P.
- Watada,J. (1992), *Fuzzy time series analysis and forecasting of sales volume*, Kacprzyk,J. Fedrizzi,M. (Eds.), Fuzzy Regression Analysis, Omnitech Press, Heidelberg, 211-227 P.

- Song,Q. Chissom,B.S. (1993), *Fuzzy time series and its models*, Fuzzy Sets and Systems, vol.54(3), 269-277 P.
- Song,Q. Chissom,B.S. (1993), *Forecasting enrollments with fuzzy time series-Part I*, Fuzzy Sets and Systems, vol.54(1), 1-9 P.
- Song,Q. Chissom,B.S. (1994), *Forecasting enrollments with fuzzy time series-Part II*, Fuzzy Sets and Systems, vol.62(1), 1-8 P.
- Tseng,Fang-Mei. Tzeng,Gew-Hshiung. Yu,Hsiao-Cheng. Yuan,J.C. (2001), *Fuzzy ARIMA model for forecasting the foreign exchange market*, Fuzzy Sets and Systems, vol.118, 9-19 P.
- Sullivan,J. Woodall,W.H. (1994), *A comparison of fuzzy forecasting and Markov modeling*, Fuzzy Sets and Systems, vol.64, 279-293 P.
- Zadeh,L.A. (1965), *Fuzzy Sets*, Inform and Control, vol.8, 338-353 P.

آزمایش‌های آمیزه‌ای و کاربرد آن در فرمولاسیون مواد غذایی

مجید صیوتی، فتح میکائیلی

P ۱۳۲۲۶

گروه آمار، دانشگاه علامه طباطبایی

چکیده: بسیاری از محصولات مورد مصرف امروزی از مخلوط کردن دو یا بیشتر از دو مولفه بوجود می‌آیند، برای مثال تولید آلیاژهای فولاد و چدن، بتونهای سیمانی، رنگهای نقاشی و با تولید معجون‌های مختلف و ... از این گونه‌اند. همچنین در اکثر موارد خواص محصولات فقط به نسبت مواد تشکیل دهنده محصول بستگی دارد.

طرح آزمایش‌های آمیزه‌ای، شاخه‌ای از طرح آزمایش‌های است که در آن پاسخ، فقط تابعی از نسبت مولفه‌های تشکیل دهنده آمیزه است. تفسیر آزمایش‌های آمیزه‌ای به ایده روش‌شناسی رویه پاسخ وابسته است که یک روش دنباله‌ای است که با استفاده از ابزارهای طرح آزمایش و رگرسیون روابط بین اثرات تیماری را توضیح می‌دهد.

یک گروه کارشناسی متشکل از متخصصین صنایع غذایی و آماراقدام به فرمولاسیون نوعی نوشیدنی به نام ایزووله پروتئین سویا کرداند. در تهیه این نوشیدنی از یک طرح لایپتیم (که به کمک نرم‌افزار SAS تهیه شده است) استفاده شده و هدف بهینه کردن طعم نوشیدنی به ازای ترکیبات نسبت‌های مختلف تشکیل دهنده نوشیدنی بوده است. در این نوشیدنی که برای اولین بار و به صورت آزمایشگاهی تولید شده، در نهایت یک فرمول بهینه برای آن بدست آمده است.

واژه‌های کلیدی: آزمایش آمیزه‌ای، سادک، چندجمله‌ای کانونی، روش‌شناسی رویه پاسخ، طرح‌های بهینه، فرمولاسیون.

۱ مقدمه

بسیاری از محصولات مورد استفاده امروزی از ترکیب دو یا چند جزء از مواد مختلف بدست می‌آید، مثلاً فولاد و چدن آلیاژهای هستند که از ترکیب نسبتهای خاصی (متفاوت) از کربن و آهن بوجود می‌آیند و یا معجونهای مختلف از نسبت‌های مختلفی از عصاره‌های میوه‌های مختلف تشکیل می‌شود.

در مواقعي خواص محصولات یا آمیزه‌ها به نسبتی (بر حسب حجم یا وزن و یا حجم

مولکولی و غیره) اجزاء تشکیل دهنده آن بستگی دارد، در مثالهای بالا خواص فولاد و چدن (از لحاظ شکنندگی، استحکام وغیره) با هم متفاوت است و یا طعم معجون بستگی به نسبت عصارهای تشکیل دهنده آن دارد. بررسی این گونه آزمایشها در آمار به نوع خاصی از طراحی آزمایشها بر می‌گردد که به آزمایش‌های آمیزه‌ای^۱ معروفند.

۲ مساله اصلی

با مقدمه کوتاهی که در مورد آزمایش‌های آمیزه‌ای بیان شد به تعریف کلی و دقیق مساله می‌پردازیم.

تعریف: آزمایش آمیزه‌ای، آزمایشی است که در آن پاسخ فقط تابعی از نسبت اجزاء تشکیل دهنده آمیزه است و به مقدار کل آمیزه بستگی ندارد.

در واقع در این‌گونه از آزمایش‌ها عاملها، نسبت‌های اجزاء تشکیل دهنده می‌باشند. این نسبتها نامنفی هستند و مجموع آنها برابر یک است. فرض کنید که در آمیزه q مولفه‌ای، x_i نسبت مولفه نام مخلوط باشد، داریم که:

$$\sum_{i=1}^q x_i = 1 \quad , x_i \leq 0 \quad , i = 1, 2, \dots, q \quad (1)$$

قید اخیر محدودیت اساسی مربوط به اجزاء متشکله آزمایش آمیزه‌ای است، به این معنی که هر مولفه بوسیله دیگر مولفه‌ها تعیین می‌شود و درجه آزادی یک سیستم آمیزه‌ای با q مولفه برابر با $1 - q$ است.

۳ فضای پارامتری

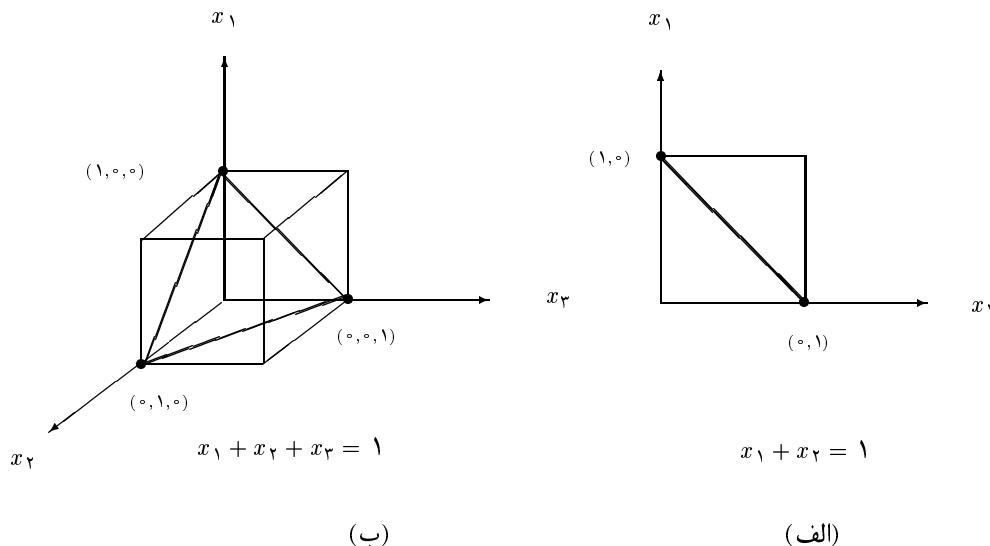
با توجه به قیود (1) فضای پارامتری شامل مجموعه نقاط واقع بر روی یا داخل سادک^۱ است.

تعریف: یک سادک، یک چند وجهی محدب و منتظم n بعدی است که دقیقاً $n + 1$ رأس دارد. در حالت صفر بعدی یک سادک یک نقطه، یک بعدی یک خط، دو بعدی مثلث و در حالت سه بعدی یک چهار وجهی و الی آخر می‌باشد. در شکل (1) فضای پارامتری برای $q = 2$ و $q = 3$ نشان داده شده است.

^۱ Mixtures Experiments

^۱ Simplex

در دستگاه مختصات سادک برای مثال با $q = 3$ مولفه، ناحیه آزمایشی یک مثلث است که رئوس این مثلث بر ترکیب خالص یعنی ترکیب تک مولفه‌ای^۲ منطبق است، یعنی نقطه‌ای که با $x_i = 1, x_j = 0, x_k = 0$ برای $i, j, k = 1, 2, 3$ و $i \neq j \neq k$ مشخص می‌شود. يالهای (ضلعها) مثلث ترکیبات دو مولفه‌ای^۳ (جزء سوم برابر با صفر) را نشان داده و نقاط داخل مثلث ترکیبات سه مولفه‌ای^۴ که تمام مولفه‌ها مخالف با صفر هستند را نشان می‌دهد.



شکل ۱ : (الف) فضای عوامل برای $q = 2$ سادک یک بعدی و (ب) فضای عوامل برای $q = 3$ سادک دو بعدی

۴ روش‌شناسی رویه پاسخ

در مسایل آمیزه‌ای، داده‌های حاصل از آزمایش در یک مقیاس کمی مثل عملکرد یا خصیصه فیزیکی نهایی تعریف می‌شود. هدف برنامه آزمایشی، ارائه الگوی پاسخ آمیزه در قالب یک معادله ریاضی برای پیش‌بینی پاسخ برای ترکیبی از اجزاء می‌باشد و یا به عبارت دیگر باید با آن معادله ریاضی بتوانیم برای اثری که هر جزء به صورت تکی یا به صورت ترکیبی روی آمیزه دارد، یک اندازه بدست آوریم. برای تحلیل مسایل آمیزه‌ای

^۲ Simple Component

^۳ Bineary Components

^۴ Ternary Components

از مجموعه‌ای از فنون ریاضی و آماری سودمند به نام روش‌شناسی رویه پاسخ^۱ استفاده می‌کنیم. در این روش ما قادر خواهیم بود رابطه بین عوامل و پاسخ را به صورت منحنی یا رویه بیان کنیم. به طور کلی این روش بر چهار گام استوار است:

- ۱) داده‌ها با استفاده از یک آزمایش از پیش طراحی شده جمع‌آوری می‌گردد.
- ۲) یک مدل ریاضی (معمولًاً چند جمله‌ای) با استفاده از تکنیکهای برآش منحنی (روشهای رگرسیونی) به داده‌ها برآش می‌شود.
- ۳) به دنبال ناحیه‌ای که مقادیر پاسخ بهینه در آن باشد، جستجو می‌گردد (با استفاده از منحنی‌های تراز پاسخ).
- ۴) نقاط اضافی برای صحّت مدل برآش شده جمع‌آوری می‌گردد.

۵ چند جمله‌ایهای کانونی

قیود (۱) نه تنها فضای پارامتری را تغییر می‌دهد بلکه چند جمله‌ایهای که برای برآورد پاسخ در این گونه آزمایشها به کار می‌رود نیز متفاوت خواهد بود.
فرض کنید چند جمله درجه اول با q مولفه بصورت زیر باشد.

$$\eta = \beta_0 + \sum_{i=1}^q \beta_i x_i$$

حال با ضرب کردن β_0 در x_i داریم:

$$\begin{aligned} \eta &= \beta_0 \sum_{i=1}^q x_i + \sum_{i=1}^q \beta_i x_i \\ &= \sum_{i=1}^q \beta_i^* x_i \end{aligned}$$

که در آن برای هر i داریم:

$$\beta_i^* = \beta_0 + \beta_i, \quad i = 1, 2, \dots, q$$

^۱ Response Surface Methodology

این چند جمله‌ایها به چند جمله‌ایهای کانونی^۱ معروفند که چند جمله‌ای درجه دوم و درجه سوم به صورت زیر با اعمالی مشابه بدست می‌آید.

$$\begin{aligned}\eta &= \sum_{i=1}^q \beta_i x_i + \sum_{1 \leq i < j \leq q} \beta_{ij} x_i x_j \\ \eta &= \sum_{i=1}^q \beta_i x_i + \sum_{1 \leq i < j \leq q} \beta_{ij} x_i x_j + \sum_{1 \leq i < j \leq q} \beta_{ij} x_i x_j (x_i - x_j) \\ &\quad + \sum_{1 \leq i < j < k \leq q} \beta_{ijk} x_i x_j x_k\end{aligned}$$

۶ فرمولاسیون نوشیدنی ایزوله پروتئین سویا

۱.۶ بیان مساله

فرآورده‌های سویا یک منبع مهم و کم هزینه پروتئین، مواد معدنی و ویتامین می‌باشد که نقش مهمی در سلامتی دارند و جهت جلوگیری از بیماریهای قلبی، چشم، سرطان، دیابت، پوکی استخوان و چاقی توصیه می‌شود. ایزوله پروتئین سویا^۱ (SPI)، فرآورده‌ای از آرد لوبیای سویای چربی زدایی شده است که شامل بیش از ۹۰ درصد پروتئین می‌باشد.

اهمیت پروتئین سویا در رژیم غذایی، بعلت ظرفیت بالای اسیدهای آمینه ضروری، میزان بالای آن در دانه سویا و ارزان بودن آن است، پروتئین سویا تمام اسیدهای آمینه ضروری بدن را فراهم می‌کند. ایزوله پروتئین سویا به دلیل خصوصیات عملکردی و ارزش تغذیه‌ای به منظور تولید فرآورده‌های جانشین لبنتیات، فرآورده‌های گوشتی، سوپ و پودر نوشیدنی‌های فوری با طعمهای متفاوت به کار می‌رود. از طرف دیگر در رژیم غذایی ورزشکاران، جهت تامین پروتئین اضافه توصیه می‌شود، چرا که ورزشکاران به مقدار پروتئین بیشتری نسبت به سایر افراد نیازمندند. همچین لازم است با دادن یک رژیم غذایی مناسب افزایش وزن ورزشکار را به افزایش توده عضلانی تبدیل کرد. از آنجایی که نیاز ورزشکاران به مقدار پروتئین بیشتر از سایر افراد است در نتیجه مقدار دریافت آب بدن آنها نیز افزایش می‌یابد تا دفع مواد حاصل از فعل و انفعالات ترکیبات ارتقی به راحتی انجام گیرد، که به کار بردن یک نوشیدنی پروتئینی به حل این مشکل کمک می‌کند.

^۱ Canonical Polynomials

^۱ Soyabeon Protein Isolate

۲.۶ تعیین پاسخ

طعم فرآورده‌های سویا مشکل عمدۀ و عامل محدود کننده گسترش کاربرد پروتئین سویا در غذای انسانی است. با بهمود خواص عملکردی و طعم ایزووله پروتئین سویا، می‌توان از آن جهت تولید یک نوشیدنی پروتئینی استفاده کرد. با تهیه فرمولاسیون مناسب از ایزووله پروتئین سویا در غالب یک نوشیدنی و با تأکید بر طعم و پذیرش کلی آن از طرف مصرف کننده می‌توان رغبت بیشتری در میان مصرف کنندگان ایجاد نمود.

۳.۶ اهداف

به طور کلی هدف‌های این تحقیق بصورت زیر بوده است:

- ۱) برازش رویه پاسخ مناسب.
- ۲) تهیه بهترین فرمولاسیون برای این محصول.

۴.۶ تعیین مولفه‌ها

جهت تهیه پودر نوشیدنی، SPI، شیرین کننده (شکر)، ایزووله پروتئین سویا، اسانس و مالتود کسترنین به صورت خشک مخلوط می‌شود و سپس در آسیاب کاملاً همگن می‌گرددند. در نهایت پودر حاصله به نسبت ۱ به ۵ در شیر حل می‌شود. به دلیل محدودیتهای آزمایشگاهی و همچنین قیدهایی که بر روی مولفه‌های تشکیل دهنده مشابه این محصول وجود دارد چهار مولفهٔ تشکیل دهنده ایزووله پروتئین سویا به صورتی که در جدول ۱ مشاهده می‌شود محدود شده‌اند.

جدول ۱: محدودیتها بر روی مولفه‌های واقعی (بر حسب گرم)

مولفه‌ها	کد حد بالا	حد پایین (درصد)	حد بالا (درصد)	کد حد پایین	حد بالا	حد پایین (درصد)
شکر	C _۱	۵	۱۱	۵	۲۵	۵۵
ایزووله پروتئین سویا	C _۲	۶	۱۱	۶	۳۰	۵۵
اسانس	C _۳	۰۱	۲	۰۱	۰۵	۱۰
مالتودکسترنین	C _۴	۱	۲	۱	۵	۱۰

مولفه‌های تشکیل دهنده در مقیاس گرم اندازه‌گیری شده‌اند و مولفه‌ها در هر آزمایش ۲۰ گرم بوده است که این مقدار تقریباً در ۱۰۰ سی سی شیر حل شده است.

۵.۶ ارزیابی پاسخ

ارزیابی پاسخ یک روش علمی است که در آن برای اندازه‌گیری تعبیر و تفسیر خصوصیات حسی غذا از حواس انسان استفاده می‌شود، اهمیت این روشها به قابلیت ویژه آنها برمی‌گردد. به عنوان مثال طعم غذا، بوی غذا و اینکه غذا چگونه به نظر می‌رسد یا احساس می‌شود توسط هیچ روش فیزیکی و شیمیایی قابل تشخیص و اندازه‌گیری نیست.

برای همین منظور گروهی متشکل از کارشناسان صنایع غذایی و آمارجهت ارزیابی خصوصیات ارگانولیپتکی^۱ شامل طعم، رنگ، مزه و به طور کلی پذیرش کلی تست پانل توسط ۶ داور انجام گرفته و میانگین نمرات آنها به عنوان پاسخ برای آن ترکیب در نظر گرفته شده است. تست پانل در قالب تست هدنیک^۲ پنج نقطه‌ای که ۲- بدترین نمره و ۲+ بهترین نمره بوده است. البته ملاحظاتی از قبیل یکنواختی نمونه، دمای نمونه، ظروف نوشیدنی، کدگذاری، ترتیب ارائه نمونه، زمان آزمایش، محیط آزمایش و ... در نظر گرفته شده است.

۶.۶ انتخاب طرح

همانگونه که در بخش ۴.۶ گفته شد مولفه‌های به کار رفته برای تهیه ایزوله پروتئین سویا در کرانهای بالا و پایین محدود هستند بنابراین طرحهای آمیزه‌ای استاندارد برای این آزمایش مستقیماً نمی‌تواند مورد استفاده قرار گیرد. با استفاده از مولفه‌نامها^۳ که به صورت تبدیل زیر انجام می‌شود آزمایش آمیزه‌ای به کار گرفته می‌شود.

$$X_i = \frac{a_i - L_i}{20 - \sum L_i}$$

که در آن a_i وزن واقعی مولفه نام آمیزه است.

اسنی^۴ [۱۹۷۵] با استفاده از اصولی موفق شد طرحی را برای مدل درجه دوم ارائه کند که این طرح شامل رؤوس فرین، مراکز صفحات مقید، نقطه مرکزی کلی و مراکز بزرگترین یالها است. نقاط نامزد برای انتخاب طرح شامل ۳۱ نقطه رؤوس، نقاط میانی یالها و وجهها را تشکیل می‌دهند. به علت محدودیت‌های زمانی و امکانات از ۳۱ نقطه نامزد، ۱۶ نقطه به عنوان طرح LD_۱ پیتمم با استفاده از ازوال SAS^۵ انتخاب شده است. این نقاط با فرض درست بودن مدل درجه دوم شفه بوده است.

^۱ Organoleptic

^۲ Hedonic

^۳ Pseudocomponents

^۴ Snee

^۵ Optimal Experimental Designs

طرح بهینه شامل ۹ نقطه راس فرین، ۴ نقطه میانی یالها، ۳ نقطه میانی وجه ها است.
۴ نقطه در رؤوس فرین هر کدام هر کدام یک بار تکرار شده تا آزمون نکویی برآش میسر
گردد.

۷.۶ تحلیل داده‌ها

تحلیل داده‌ها یکی از مهمترین مراحل طرح آزمایشی است، برای تحلیل داده‌ها و برآش رویه پاسخ به چهار مولفه‌نما از روال^۱ RSREG و^۲ GLM در نرم افزار SAS استفاده شده است.

جدول ۲: آزمون نکویی برآش و تحلیل واریانس

Regression	df	Sum of Sqre	F-Value	P-Value
Model	۹	۱۳۸۶	۳۳۸	۰۰۳۵
Error	۱۰	۴۵۵		
Lack of fit	۶	۲۸۲	۳۴۹	۰۱۲
Pure error	۴	۰۱۸		
Total	۱۹	۱۸۴۱		

در پیوست برنامه مربوط به آن با جزئیات بیشتری آورده شده است. آماره‌های مهم و مرتبط از خروجی مورد نظر استخراج شده و در ادامه بحث خواهد شد. همانگونه که در جدول ۲ دیده می‌شود آزمون LOF مناسب بودن مدل درجه دوم را نشان می‌دهد.

مقدار R^2 برابر با ۷۵٪ و همچنین مقدار Root MSE برابر با ۶۸٪ است که نشان دهنده خوبی مدل درجه دوم است. همچنین با استفاده از روال GLM برای برآش مدل و برآورد ضرایب استفاده شده است.

جدول ۳: برآورد پارامترهای رویه پاسخ

Variable	Coefficient	S.E	t-Value	P-Value
X_1	-۵۰۲۵	۱۰۸۳	-۲۸۶	۰۰۱۲
X_2	-۱۰۳۶	۱۰۷۰	-۰۸۰	۰۴۴۳
X_3	-۴۰۶۰۴	۱۰۴۰۲۶	-۳۸۹	۰۰۰۱
X_4	۳۲۰۵۰	۹۰۷	۳۵۸	۰۰۰۳
$X_1 X_3$	۵۲۰۰۸۰	۱۲۷۰۸۶	۴۰۷	۰۰۰۱
$X_2 X_3$	۴۶۸۰۹۳	۱۲۸۰۶۴	۳۶۴	۰۰۰۲

^۱ Response Surface Regression

^۲ General Linear Models

همچنین آزمونهای نرمالیتی و همسانی واریانس باقیمانده‌ها نیز انجام گرفته که نتایج نرمال بودن باقیمانده‌ها و همسانی واریانس برای تمام سطوح X را نشان می‌دهد. نسبتهاي بهينه (نسبتهاي از مولفه‌ها که به ازاي آنها پاسخ بهينه، طعم مطلوب، بدست آمده است) به صورت زير(بر حسب درصد) است:

$$\begin{array}{ll} P_1 = 41 & P_2 = 47 \\ P_3 = 3 & P_4 = 9 \end{array}$$

مراجع

فرهوش، ر. (۱۳۷۴). بررسی امکان تولید ایزوله پروتئین سویا، پایان‌نامه کارشناسی ارشد علوم و صنایع غذایی، دانشکده کشاورزی، دانشگاه تربیت مدرس.

جهانیان، ل. (۱۳۸۱). تهیه نوشیدنی ایزوله پروتئین سویا، پایان‌نامه کارشناسی ارشد علوم و صنایع غذایی، دانشکده کشاورزی، دانشگاه تربیت مدرس.

Cornell, J. A. (1990), *Experiments with Mixtures: Designs, Models and the Analysis of Mixture Data*. 2nd edition. John Wiley & Sons, New York, 1990.

Myers, R. H. and Montgomery, D. C. (1995), *Response Surface Methodology: process and product optimization using designed experiments*. John Wiley & Sons, New York, 1995.

Ruguo, Hu. (1999), *Food Product Design: A Computer-Aided Statistical Approach*, Technomic Publishing Company, Inc.

Scheffé, H. (1958), Experiments with Mixtures. *J. R. Stat. Soc, B*, 20, No.2, 344-360.

Snee, R. D. (1979), Experimenting with Mixtures. *CHEMTECH*, 9, 702-710.

Snee, R.D. (1975) Experimental Design for Quadratic Models in Constrained Mixture Spaces. *Techonometrics*, 17, 149-159

پیوست

روال RSREG روش حداقل مربعات را برای برآورد پارامترهای درجه دوم به کار می‌گیرد. مدل‌های رویه پاسخ مخصوصاً هنگامی که هدف برآورد پاسخ بهینه است به کار گرفته می‌شوند. به طور کلی روال RSREG می‌تواند برای اهداف زیر مورد استفاده شود:

- ۱) آزمون برای عدم برازش
- ۲) آزمون اثرات عاملهای تکی
- ۳) تحلیل ساختار کانوئیک رویه پاسخ برآورده شده
- ۴) تحلیل ستیغی برای رسیدن به پاسخ بهینه
- ۵) پیش‌بینی مقادیر جدید پاسخ

فرم کلی استفاده از این روال در SAS به صورت زیر است:

```
PROG RSREG <options>;
  MODEL response=independent <options>;
  RIDGE <options>;
  WEIGHT <options>;
  ID variables;
  BY variables;
```

استفاده از PROC RSREG و عبارت MODEL اجباری است. در ادامه خروجی کامپیوتری با استفاده از ADX و سایر روالهای موجود در SAS با جزئیات بیشتری آورده شده است.

تخمین بهینه وزنها در نمودارهای EWMA و DEWMA با استفاده از MCDM

رضا برادران کاظمزاده^۱، عزیزا... معماریانی^۱، مهدی کرباسیان^۲

P ۱۱۰۲۴

^۱ بخش مهندسی صنایع، دانشگاه تربیت مدرس

^۲ دانشگاه مالک اشتر

چکیده: ضرایب پیش‌بینی کننده و تصحیح کننده در نمودارهای کنترل میانگین متحرک موزون نمایی EWMA و نمودارهای کنترل میانگین متحرک موزون نمایی جفتی DEWMA نقش اساسی در تنظیم عملیات کنترل کیفیت آماری در فرایند ساخت نیمه‌هادی، الکتروسراام و قطعات اپتیکی دارد. تخمین بهینه این ضرایب که معمولاً بر اساس حداقل کردن مریع خطای M.S.E. انجام می‌شود باعث تشخیص سریع هر گونه تغییری در میانگین و برآکندگی مشخصه کیفی فرایند می‌شود.

در این مقاله از روش تصمیم‌گیری با معیارهای چندگانه M.C.D.M و تکنیک برنامه‌ریزی آرمانی G.P برای کاهش اریب و انحراف معیار آماره آن نمودارها به صورت جداگانه و بر اساس نظر و آرمانهای تصمیم‌گیرنده D.M و اولویتهایی که او قائل است استفاده می‌شود. این مساله با توجه به اهمیت نظرات مدیریت و کارشناسان در طراحی، شروع و ادامه نمودارهای کنترلی، حائز اهمیت است.

۱ مقدمه

نمودارهای کنترل برای اولین بار توسط دکتر والتر شیوارت (۱۹۹۳) معرفی شده‌اند. پیج در سال ۱۹۹۵ برای اولین بار پیشنهاد کرد که از حدود کنترلی به نام حدود هشدار، جدا از حدود کنترل شیوارت استفاده شود. (پاگ ۱۹۹۵) بعد از پیج حدود و معیارهای دیگری برای حساس سازی بیشتر نمودارهای کنترل توسط افرادی از قبیل چمپ و دال در سال ۱۹۹۷ ابداع گردید. پیج در سال ۱۹۵۴ استفاده از نمودارهایی با عنوان جمع تجمعی^۱ CUSUM را مطرح نمود (پاگ ۱۹۵۴). در این نمودار از اطلاعات موجود در نمونه‌های قبلی نیز استفاده می‌شود. نمودار میانگین متحرک موزون نمایی^۲ EWMA با خواصی شبیه به نمودارهای CUSUM برای اولین بار توسط روبرتس در سال ۱۹۵۹ توسعه یافت. لوكاس و سکیوس در سال ۱۹۹۰ نتیجه‌گیری کردند که عملکرد نمودارهای CUSUM

^۱ Cumulative Sum Control Charts

^۲ Exponentially Weighted Moving Average

و EWMA بسیار شبیه به هم هست (لوکاس ۱۹۹۰). در مقالات متعددی نیز به کارایی قویتر نمودارها EWMA در مقایسه با نمودارهای CUSUM اشاره شده است (وان براکل و گان) به هر حال زمانی که پی‌بردن به وجود تغییرات کوچک مورد نظر باشد، نمودارهای CUSUM و EWMA بهتر از نمودارهای شیوارت عمل می‌کنند. البته لوکاس در سال ۱۹۷۳ نشان داده است که نمودارهای شیوارت نسبت به تغییرات بزرگ در میانگین، حساس‌تر از نمودارهای EWMA می‌باشند.

نمودارهای EWMA غالباً از مشاهدات انفرادی استفاده می‌کنند. البته آماره این نمودار به علت اینکه میانگین موزونی از مشاهدات انفرادی باشد. همچنین این آماره کاربرد زیادی در مدل‌سازی سریهای زمانی و پیش‌بینی دارد (جنکینز ۱۹۷۶) و (پری ۱۹۷۳).

بالتر و استفانی در فرایندهای متحول و آشوبناک نمودارهایی پیشنهاد کردند که از یک ضریب تصحیح کننده – پیش‌بینی کننده استفاده می‌کند. این نمودارها با نام ^۱DEWMA شناخته می‌شوند.

۲ نمودارهای EWMA و DEWMA

آماره میانگین متحرک موزون نمایی در نمودار DEWMA به صورت زیر تعریف شده است:

$$Z_t = \lambda X_t + (1 - \lambda) Z_{t-1} \quad (1-2)$$

در این رابطه λ دارای مقدار ثابتی بین $0 < \lambda \leq 1$ است. مقدار Z_t برابر μ می‌باشد. در بعضی مواقع نیز مقدار Z_t از میانگین اطلاعات اولیه (\bar{X}) تخمین زده می‌شود. رابطه $(1-2)$ به صورت زیر بازنویسی می‌شود:

$$\begin{aligned} Z_t &= \lambda X_t + (1 - \lambda)[\lambda X_{t-1} + (1 - \lambda)Z_{t-2}] \\ &= \lambda X_t + \lambda(1 - \lambda)X_{t-1} + (1 - \lambda)^2 Z_{t-2} \\ Z_t &= \lambda \sum_{j=0}^{t-1} (1 - \lambda)^j X_{t-j} X_t + (1 - \lambda)^t Z_0. \end{aligned} \quad (2-2)$$

^۱ Double EWMA

وزنهای λ^j به صورت هندسی با افزایش عمر کاهش می‌یابند. همچنین جمع وزنها به سمت یک میل می‌کند:

$$\lambda \sum_{j=0}^{t-1} (1-\lambda)^j = \lambda \left[\frac{1 - (1-\lambda)^t}{1 - (1-\lambda)} \right] = 1 - (1-\lambda)^t$$

اگر مشاهدات X_i متغیرهای تصادفی مستقل با واریانس σ^2 باشد آنگاه:

$$\begin{aligned} \sigma_{Z_t}^2 &= \sigma^2 \lambda^2 \left[\frac{1 - (1-\lambda)^{2t}}{1 - (1-\lambda)^2} \right] \\ &= \sigma^2 \left(\frac{\lambda}{2-\lambda} \right) [1 - (1-\lambda)^{2t}] \end{aligned} \quad (3-2)$$

با افزایش σ_{Z_t} به مقدار حدی خود یعنی:

$$\sigma_{Z_t}^2 = \sigma^2 \left(\frac{\lambda}{2-\lambda} \right) \quad (4-2)$$

میل خواهد کرد.

از طرف دیگر اگر میانگین هر x_t را μ بنامیم، در آن صورت امیدریاضی Z_t به صورت زیر خواهد بود:

$$E(Z_t) = \mu - (1-\lambda)^t (\mu_0 - \mu) \quad (5-2)$$

در نمودارهای DEWMA آماره W_t یک تصحیح کننده است، به صورت زیر محاسبه می‌شود:

$$W_t = \lambda_2 (x_t - z_{t-1}) + (1 - \lambda_2) W_{t-1} \quad (6-2)$$

که در آن $W_0 = 0$ ، $W_0 < \lambda_2$ ، رابطه فوق را می‌توان به صورت زیر بازنویسی کرد:

$$W_t = \lambda_2 \sum_{j=0}^{t-1} (1 - \lambda_2)^j F_{t-j}$$

که در آن $F_{t-j} = x_t - Z_{t-j}$ می‌باشد و میزان آماره DEWMA برابر مقدار خواهد بود. اگر مشاهدات x_t متغیرهای تصادفی مستقل با واریانس σ^2 باشند و با توجه به رابطه (4-2) واستقلال Z_{t-1} و X_t ، آنگاه R_t دارای واریانس زیر خواهد بود:

$$si_{R_t}^2 = \sigma^2 + \left(\frac{\lambda_1}{2 - \lambda_1} \right) \left(\frac{\lambda_1}{2 - \lambda_1} \right) \quad (7-2)$$

از طرف دیگر اگر میانگین X_t برابر μ باشد و با توجه به رابطه (۵-۲) خواهیم داشت:

$$E(R_t) = \mu - (1 - \lambda_1)^t(\mu_0 - \mu) + [1 - (1 - \lambda_1)^t][\mu - (1 - \lambda_1)^{t-1}(\mu_0 - \mu)]$$

۳ طراحی نمودارهای کنترل DEWMA و EWMA

مهتمرين پارامترهایی که در طراحی اين نوع نمودارها مورد استفاده قرار می‌گيرند مقادير λ_1 و λ_2 مهمترین مقادير می‌باشند. اين امكان وجود دارد که بتوان اين پارامترها را به گونه‌اي انتخاب نمود تا هرگونه انحراف از میانگين و يا افزایش پراکندگي با سرعت تشخيص داده شود. دل کاستلو در سال ۱۹۹۹ يك مدل غير خطی برای بدست آوردن مقادير بهينه اين ضرائب بر اساس كمينه کردن ميزان ميانگين مربع خطأ^۱ MSE بدست آورده است.

۴ استفاده از تکنيکهای MCDM برای تخمین بهينه λ_1 و λ_2

توجه محققین دردهه اخير معطوف به مدلهاي چندمعياره^۲ (MCDM) برای تصميم‌گيري پيچيده گردیده است. شايد بتوان گفت برنامه‌ريزي آرمانی GP^۳ يكى از قديميترين مدلهاي موجود در تصميم‌گيري چند معياره است که با كاريدهای وسیع به كارگيري شده است.

چارلز و کوپر اولين مقاله را درباره G.P در سال ۱۹۵۵ منتشر کردند به طوری که آنها مى‌نیمم کردن مجموع قدر مطلق انحرافات از مقاصد مشخصی را مورد بررسی قرار داده‌اند. تلاش در G.P بر آن است که منطق مدلهاي رياضي بهينه تواماً با تمايل تصميم‌گيرنده^۴ در تامين مقاصد مشخص اراده‌اف مختلف مورد توجه قرار بگيرند.

در اين قسمت ميزان اريب و واريанс آماره نمودارهای DEWMA و EWMA و به صورتی آرمانی بهينه مى‌گردد، آرمان در اريب صفر مى‌باشد (بي اربي) و آرمان در مورد واريанс بستگي به نظر تصميم‌گيرنده، دارد.

با اين اوصاف مدل G.P در نمودار EWMA به صورت زير است:

$$\frac{\text{lexMin}\{n_1 + p_1, n_2\}}{(1-4)}$$

^۱ Mean Square Error

^۲ Multi Squer Error

^۳ Goal-Programing

^۴ Decision Maker

$$(1 - \lambda)^t(\mu_0 - \mu) + n_1 + p_1 = 0 \quad (1)$$

$$\sigma^2 \left(\frac{\lambda}{2 - \lambda} \right) + n_2 = S \quad (2)$$

$$0 < \lambda \leq 1 \quad (3)$$

$$n_1, n_2, p \geq 0 \quad (4)$$

باتوجه به اهمیت اریب، از مدل لیکسوسوگراف با اولویت اریب نسبت به واریانس استفاده شده است.

میزان واریانس آماره است که توسط D.M تعیین می‌شود.

مدل G.P. در حالت DEWMA به صورت زیر می‌باشد:

$$\text{let} \min\{n_1 + p_1, n_2\} \quad (2 - 4)$$

$$(1 - \lambda_1)^t(\mu_0 - \mu) - [1 - (1 - \lambda_1)^t][\mu - (1 - \lambda_1)^{t-1}(\mu_0 - \mu)] + n_1 + p_1 = 0 \quad 1$$

$$\left(\frac{\lambda_2}{2 - \lambda_2} \right) + n_2 = S$$

$$0 < \lambda \leq 1$$

$$n_1, n_2, p_1 \geq 0$$

۵ تیجه‌گیری و پیشنهادات

تعیین وزنهای بهینه λ_1 و λ_2 در نمودارهای EWMA و DEWMA، در کاهش اریب و واریانس و در تیجه تشخیص سریع تغییر در میانگین و افزایش پراکندگی و کاهش خطای نوع دوم (β) کاملاً موثر است. همچنین مقادیر بهینه λ_1 و λ_2 باعث این پایداری این نمودارها می‌شود.

می‌توان گفت وزنهای بهینه λ_1 و λ_2 را در حالتی که نمودارهای فوق برای آماره‌هایی مانند تعداد عیب و یا تعداد قطعات معیوب به کار می‌روند محاسبه کرد. همچنین استفاده از برنامه‌ریزی آرمانی موزون و همچنین محاسبه وزنهای بهینه در این نوع برنامه‌ریزی برای فعالیتهای آینده پیشنهاد می‌شود.

مراجع

- Shewart,W.A(1993). "The Economic control of Qualitiy Manufactured Products". Van Nostrand, New York
- Pages, E.S. (1995). "Control Charts With Warning Lines" Biometrika, 41, 243-259.
- Champ, C.W., and Woodall, W.H., (1987). "Exact Result For Shewart Control Charts With Supplementary Run Rules" Technometrics, 29, 393-399.
- Pages, E.S. (1954). "Continuous Inspection Schemes" Biometrika, 41, 100-114.
- Roberts, S.W. "Control Chart Tests Base On Geometric Moving Averages" Technometrics,(1959) 1, pp 239.250
- Lucas,J.M. and Saccucci, M.S.(1990). "Exponentially Moving Averages Control Chart: Properties and Enhancements" Technometrics, 31, 1-12.
- Van Brackle, L.N., III; Rynols, M.R., Jr.EWMA amd CUSUM Control Charts In THe Procence Of Corrolation, Commuinications In Statistics -Simultain" Vol:26, lss: 3, pp.979, 1008
- Gan, F.F." Desigen Of One-and Two-saided exponential EWMA Charts" Jornal of Quality Technology.
- Lucas, J.M.(1973)"A modified V-Mask Scheme." Technometrics, 15,833-847
- Box, G.E.P., and G.M.jenkins(1976) "Time Series Analysis, Forecasting ,and Control,"Holden-Day, Sanfransisco.
- Perry, R.I. (1973). "Skip-Lot sampling Planes." Journal of Quality Technology, Vol. 5.
- Bulter, S. W. and Stfain, J. a. (1994). Supervisory Run-To-Run Control of a Polysilicon Gate etch using in Situ Ellipsometry." IEEE Transaction on Semiconductor Manufacturing. 7,193-201.

Castillo, E. D.(1999) "Long Run and Transient Analysis of a Double Ewma Feedback Controller.", IIE Transaction, 31, 1137-1169.

مقایسه مدل رگرسیون پواسن و مدل رگرسیون دوجمله‌ای منفی در تعیین عوامل مؤثر بر حاملگی ناخواسته

نوشیروان کاظم‌نژاد، سیدمهدی سادات‌هاشمی، شبم کریمی

P۱۳۲۰۸

گروه آمار زیستی، دانشگاه تربیت مدرس

چکیده: موضوع تحقیق حاضر «مقایسه مدل رگرسیون پواسن و مدل رگرسیون دوجمله‌ای منفی در تعیین عوامل مؤثر بر حاملگی ناخواسته» می‌باشد و متغیر وابسته تعداد حاملگی‌های ناخواسته در نظر گرفته شده است.

در آنالیز رگرسیون در حالتی که متغیر وابسته گستته و نامنفی باشد، مدل رگرسیون پواسن مورد استفاده قرار می‌گیرد. مدل رگرسیون پواسن به عنوان نمونه‌ای از مدل‌های تعمیم یافته خطی است که مدل‌های رگرسیونی را به خانواده نمایی توزیع‌هایی بسط می‌دهد که هر دو توزیع نرمال و پواسن را شامل می‌شوند. شرط اصلی استفاده از مدل رگرسیون پواسن معادل بودن میانگین و واریانس متغیر وابسته می‌باشد. وقتی که میانگین و واریانس داده‌ها به طور تقریبی برابر نباشند، مدل پواسن برآوردهای ناصحیحی از واریانس جملات واستنباط‌هایی گمراه کننده درباره رگرسیون ایجاد می‌کند. برای حل این مشکل می‌توان از مدل رگرسیونی دوجمله‌ای منفی استفاده کرد.

به منظور تعیین مهمترین عوامل خطرساز در حاملگی ناخواسته و تعیین برآورد پارامترها از مدل‌های رگرسیونی پواسن و دوجمله‌ای منفی استفاده شده و مقایسه‌ای بین آنها صورت گرفته است. از بین کلیه متغیرهای مورد بررسی، متغیرهای سطح تحصیلات مادر، تعداد فرزندان زنده دختر، تعداد فرزندان زنده پسر و استفاده از روش پیشگیری مؤثر بوده و متغیرهای سن، سطح تحصیلات پدر و آگاهی قبلی از روش‌های تنظیم خانواده از نظر آماری از مدل حذف شدند و اثرات متقابل نیز وارد مدل نشدند.

واژه‌های کلیدی: مدل‌های تعمیم یافته خطی، رگرسیون پواسن، رگرسیون دوجمله‌ای منفی، حاملگی ناخواسته.

۱ مقدمه

۱.۱ مدل‌های تعمیم یافته خطی

بخش وسیعی از آمار ریست سنجی و اقتصاد، به تحلیل داده‌های پاسخ گسسته مربوط می‌شود. در مواردی که متغیر پاسخ گسسته باشد (متغیرهای مستقل از هر نوع که باشند) مدل‌های تعمیم یافته خطی را به کار می‌بریم.

مدل‌های تعمیم یافته خطی دو ویژگی کلی هستند. اول اینکه برای متغیر وابسته y_i با مقدار مورد انتظار μ_i ، توزیع احتمال y_i به شرط μ_i ، عضوی از خانواده نمایی است. دوم اینکه یکتابع پیوند وجود دارد (تابع تبدیل $(\cdot)^g$) که مقدار مورد انتظار y_i را خطی می‌کند [1].

مدل رگرسیون پواسن و مدل رگرسیون دوجمله‌ای منفی نمونه‌هایی از مدل‌های تعمیم یافته خطی‌اند که برای آنالیز داده‌های شمارشی مورد استفاده قرار می‌گیرند [2].

۲.۱ مدل رگرسیون پواسن

ساده‌ترین مدل برای داده‌های شمارشی، مدل رگرسیون پواسن می‌باشد. مدل رگرسیون پواسن به عنوان نمونه‌ای از مدل‌های تعمیم یافته خطی است که مدل‌های رگرسیونی را به خانواده نمایی توزیع‌هایی بسط می‌دهد که هر دو توزیع نرمال و پواسن را شامل می‌شوند.

روش رگرسیون پواسن مابین آنالیز مبتنی بر داده‌های پیوسته و آنالیز مبتنی بر داده‌های دو حالتی قرار گرفته است. در رگرسیون پواسن متغیر وابسته شمارشی بوده و دارای توزیع پواسن می‌باشد که این موضوع اساس استنباط را تشکیل می‌دهد. معادله رگرسیون پواسن، یک متغیر شمارشی یا یک نرخ را به یک سری متغیرهای مستقل مربوط می‌سازد و ساختاری را برای تحلیل آماری ارائه می‌دهد. در مدل رگرسیون پواسن، لگاریتم نرخ بر اساس مجموع وزنی متغیرهای مستقل (متغیرهای خطی) مدل‌بندی می‌شود. در واقع این مدل، میانگین متغیر وابسته را بر حسب متغیرهای مستقل، مدل‌بندی می‌کند [3].

محدوده اصلی مدل رگرسیون پواسن بر این شرط استوار است که واریانس متغیر وابسته با میانگین آن برابر باشد، اما در بسیاری از پدیده‌ها داده‌ها بسیار پراکنده‌اند. به عبارت دیگر ممکن است واریانس به طور معنی داری بزرگتر از میانگین باشد. وقتی که میانگین و واریانس داده‌ها به طور تقریبی مساوی نباشند، واریانس‌های ضرایب برآورد شده در مدل پواسن اریب خواهند بود. در این حالت دیگر برازش رگرسیون پواسن بر روی داده‌ها مناسب نمی‌باشد، ولی می‌توان این محدودیت را با استفاده از توزیع دوجمله‌ای منفی از میان برداشت [4].

در تحقیقی که در سال ۱۳۷۸ توسط آقای مهران آل ابریشم کارزاده با عنوان «کاربرد مدل رگرسیون پواسن در تعیین رابطه کیست تخدمان ناشی از کاشت نوریلانت با ریسک فاکتورهای مربوطه» انجام شد، ایشان به تفصیل در ارتباط با رگرسیون پواسن و کاربرد آن به بحث و بررسی پرداخته‌اند [۵].

۳.۱ رگرسیون دوجمله‌ای منفی

توزیع دوجمله‌ای منفی برای داده‌های گستته و نامنفی مناسب بوده و به برابری میانگین و واریانس حساس نیست. پس به لحاظ مشابه بودن فرم فضای نمونه‌ای، می‌توان از توزیع دوجمله‌ای منفی به جای توزیع پواسن استفاده کرد و برآوردهای دقیق‌تری را به دست آورد.

تقریباً در همه تحقیقات علوم اجتماعی انتظار داریم که تغییرات تبیین نشده‌ای بین موارد وجود داشته باشد، که این تغییرات تبیین نشده، اختلافهای مرتبط با پیش‌بینی کننده‌های مذکور نشده را منعکس می‌کند.

رگرسیون دوجمله‌ای منفی را می‌توان شکل دیگری از رگرسیون پواسن دانست که شامل مؤلفه‌ای تصادفی است که عدم ثبات در نخ‌های حقیقی رویدادهای واقع شده برای اشخاص را منعکس می‌کند. در واقع مسئله اختلافهای بین موردی تبیین نشده، در مدل رگرسیون دوجمله‌ای منفی حل شده است [۲].

اولین کاربردهای توزیع دوجمله‌ای منفی توسط استیونت (۱۹۰۷) در مقاله‌ای که در آن از توزیع پواسن نیز بحث می‌کند آمده است. استیونت از این توزیع به عنوان صورت دیگری از توزیع پواسن در شمارش‌های توصیف شده روی طرح‌های از یک هماسیتومتر استفاده کرده است. همچنین یول (۱۹۱۰) با استفاده از این توزیع، تعداد مرگ و میرها را در بروز یک بیماری معین در نظر گرفته است. گرین وود و یول (۱۹۲۰) این توزیع را به عنوان یک نتیجه از فرضهای ساده خاصی در مدل پدیده‌های تصادفی به کار گرفتند [۶].

۴.۱ حاملگی ناخواسته

چنانچه زوجین یا یکی از آنها در حوالی به وقوع پیوستن بارداری تصمیم به بچه‌دارشدن نداشته باشند و یا یکی از روش‌های پیشگیری از بارداری را استفاده کرده باشند، بارداری به وقوع پیوسته به نام بارداری برنامه ریزی نشده تلقی می‌گردد. حال اگر بعد از اثبات وقوع این بارداری، یکی از زوجین یا هر دو آنها تمایلی به ادامه یافتن بارداری نداشته و خواهان خاتمه آن باشند، بارداری ناخواسته تلقی می‌گردد که در برخی موارد این قبیل بارداری‌ها، حاملگی برنامه ریزی نشده – ناخواسته نیز نامیده می‌شود [۷].

در زمینه بررسی علل وقوع حاملگی ناخواسته تحقیقات مختلفی در ایران صورت گرفته که در اینجا به ذکر نمونه‌هایی از آن می‌پردازم.

در تحقیقی جهت بررسی مشخصات زنان با حاملگی ناخواسته در بین مراجعین به مراکز بهداشتی – درمانی جنوب شهر تهران، خانم امامی افشار در طی ۴۵ روز در زمستان سال ۱۳۶۸ از چهار مرکز بهداشتی – درمانی که به صورت تصادفی از بین یازده مرکز در منطقه جنوب شهر تهران انتخاب شدند، ۳۵۰ پرونده به روش نمونه‌گیری تصادفی انتخاب و از طریق پرسشنامه‌ای شامل ۳۵ سؤال از مادرانی که حداقل تجربه یک بار حاملگی و زایمان را داشتند، اطلاعات مورد نیاز را جمع آوری نمود. متغیرهای مورد بررسی عبارت بودند از: سن مادر، تعداد فرزندان زنده، میزان تحصیلات مادر و همسرش. جهت تجزیه و تحلیل این اطلاعات از تست کای اسکوئر استفاده گردید و مشخص شد که در این داده‌ها بین دو متغیر میزان تحصیلات و حاملگی ناخواسته ارتباطی معکوس وجود دارد و همچنین رابطه معنی داری بین تعداد فرزندان زنده و حاملگی ناخواسته وجود داشت یعنی با افزایش تعداد فرزندان، درصد حاملگی ناخواسته در این مادران افزایش می‌یافتد [۸].

در تحقیق دیگری که در سال ۱۳۷۵ توسط خانم فریفته منصوریان صورت گرفت ۶۱۲۵ خانم تحت مطالعه قرار گرفتند. هدف از این مطالعه بررسی عوامل مؤثر بر حاملگی ناخواسته و اقدام به سقط بوده است. متغیرهای مورد بررسی عبارت بودند از: شاخص‌های جمعیتی پدر و مادر، تعداد فرزندان زنده، فاصله حاملگی از زایمان یا سقط قبلی، مراجعه برای معاینات و وزن نوزاد. از نظر آماری متغیرهای سن پدر و مادر، تحصیلات پدر و مادر، شغل پدر و مادر، قومیت پدر و مادر، تعداد فرزندان زنده و فاصله حاملگی از زایمان یا سقط قبلی به عنوان عوامل مؤثر بر حاملگی ناخواسته شناخته شدند [۹].

۲ مواد و روشها

۱.۲ مواد

داده‌هایی که در این مقاله مورد تجزیه و تحلیل قرار می‌گیرد، از پرسشنامه‌های مربوط به طرح بررسی تاثیر مشاوره تنظیم خانواده پس از زایمان در پیشگیری از حاملگی‌های ناخواسته در دو سال اول پس از زایمان، که توسط مرکز ملی تحقیقات بهداشت و باروری صورت گرفته، جمع آوری شده است.

در این تحقیق جامعه آماری، کلیه خانم‌هایی هستند که در فاصله سالهای ۱۳۷۷ - ۱۳۷۵ برای زایمان به یکی از مراکز دهگانه شهر تهران (بیمارستان‌های مهدیه، میرزا

کوچک خان، شریعتی، لقمان، امام حسین، فیروزگر، امام خمینی، رهنمون، آرش و اکبر آبادی) مراجعه کرده‌اند.

با توجه به ماهیت مطالعه، از بین خانم‌هایی که برای زایمان به یکی از مراکز دهگانه فوق در طی سالهای ۱۳۷۷ - ۱۳۷۵ مراجعه کرده بودند، نمونه تصادفی ساده به حجم ۴۱۷۷ نفر استخراج گردید.

تعداد حاملگی‌های ناخواسته متغیر وابسته (متغیر پاسخ) و سایر عواملی را که در این قسمت معرفی می‌کنیم متغیرهای مستقل مدل (پیش‌بینی کننده‌ها) می‌باشند.

(۱) AGE: این متغیر سن افراد مورد مطالعه را مشخص می‌کند.

(۲) GRADE: این متغیر سطح تحصیلات مادر را مشخص می‌کند و در پنج رده (۱ - بی‌سودا، ۲ - ابتدایی، ۳ - متوسطه، ۴ - دیپلم و ۵ - بالاتر از دیپلم) طبقه‌بندی شده است.

(۳) FGRADE: این متغیر سطح تحصیلات پدر را مشخص می‌کند و در پنج رده (۱ - بی‌سودا، ۲ - ابتدایی، ۳ - متوسطه، ۴ - دیپلم و ۵ - بالاتر از دیپلم) طبقه‌بندی شده است.

(۴) DAUT: این متغیر نشان دهنده تعداد فرزندان زنده دختر می‌باشد.

(۵) SON: این متغیر نشان دهنده تعداد فرزندان زنده پسر می‌باشد.

(۶) CONTRA: این متغیر نشان دهنده استفاده از روش پیشگیری از حاملگی قبل از وقوع حاملگی است و در دو رده (۱ - پیشگیری داشته و ۲ - پیشگیری نداشته) طبقه‌بندی شده است.

(۷) KNOW: این متغیر نشان دهنده آگاهی قبلی مادر در ارتباط با روش‌های تنظیم خانواده است و در دو رده (۱ - آگاهی قبلی داشته و ۲ - آگاهی قبلی نداشته) طبقه‌بندی شده است.

۲.۲ روش‌ها

در مدل رگرسیون پواسن، لگاریتم نرخ براساس مجموع وزنی متغیرهای مستقل مدل‌بندی می‌شود. تحت فرض اینکه y_i (متغیر وابسته) دارای توزیع پواسن با میانگین μ_i باشد، داریم:

$$Pr[Y_i = y_i] = p(y_i) = (\mu_i^{y_i} e^{-\mu_i}) / y_i, i = 1, 2, \dots, n$$

به طوری که مقادیر برآورد شده از رابطه زیر محاسبه می‌شود:

$$\mu_i = E[Y_i] = e^{\beta X_i}$$

که در آن μ_i تعداد وقایع مورد انتظار، X_i بردار متغیرهای پیش بینی کننده و β بردار ضرایب رگرسیون می‌باشد.

مدل $\mu_i = \exp(\beta X_i)$ مدلتابع نرخ نمایی است که همواره نامنفی می‌باشد. این نوع تابع نرخ در پردازش انواع مختلف داده‌های شمارشی بسیار انعطاف پذیر است [4]. رگرسیون پواسن معرفی شده به وسیله رابطه فوق توسط روش‌های حداکثر درستنمایی استاندارد قابل محاسبه است [5].

توزیع پواسن یک توزیع داولطلب اولیه برای متغیرهای تصادفی گستته می‌باشد، ولی اگر میانگین و واریانس داده‌ها به طور تقریبی برابر نباشد، رگرسیون پواسن مدل نامناسبی خواهد بود. یک روش ساده برای حل این مشکل استفاده از توزیع دوجمله‌ای منفی می‌باشد.

مدل دوجمله‌ای منفی از مدل پواسن با اضافه کردن جمله خطای ϵ که به طور مستقل

$$Ln(\mu_i) = \beta X_i + \epsilon$$

جمله خطای ϵ دارای توزیع گاما با میانگین یک و واریانس α می‌باشد [2]. این مدل ارتباط بین میانگین و واریانس را به صورت زیر نتیجه می‌دهد [4]:

$$Var(y_i) = E(y_i)[1 + \alpha E(y_i)]$$

α در رابطه فوق همان واریانس جمله خطاست و یک پارامتر سنجش پراکندگی است که به وسیله روش‌های حداکثر درستنمایی استاندارد قابل برآورد است [10].

اگر $\alpha = 0$ باشد مدل دوجمله‌ای منفی به مدل پواسن تبدیل می‌شود ($Var(y_i) = E(y_i)$). ولی اگر α اختلاف معنی داری از صفر داشته باشد، مدل دوجمله‌ای منفی انتخابی شایسته و مدل پواسن نامناسب می‌باشد.

در مدل دوجمله‌ای منفی با استفاده از ضربه تعدیلی α محدودیت واریانس مدل پواسن برطرف می‌شود، به دلیل اینکه واریانس ϵ یعنی α پارامتری اضافه شده به مدل واریانس می‌باشد، رگرسیون دوجمله‌ای منفی در مدل بندی ارتباط بین مقدار مورد انتظار و واریانس y_i دارای انعطاف پذیری بیشتری نسبت به مدل به شدت محدود پواسن می‌باشد [4].

در قسمت نتایج به کمک مدل رگرسیون پواسن و مدل رگرسیون دوجمله‌ای منفی اثر هر یک از متغیرها را بر متغیر وابسته (تعداد حاملگی‌های ناخواسته) مورد بحث و بررسی قرار داده‌ایم.

۳ نتایج

با توجه به ماهیت متغیر وابسته (تعداد حاملگی‌های ناخواسته) به منظور مدل بندی عوامل مؤثر بر آن از مدل‌های رگرسیونی پواسن و دوجمله‌ای منفی استفاده می‌کنیم. روش مورد استفاده رگرسیون گام به گام پیشرو با $P_E = ۰/۲$ و $P_R = ۰/۰۵$ می‌باشد. P_E مقدار P-Value برای گنجاندن متغیر در مدل و P_R مقدار P-Value برای حذف متغیر از مدل می‌باشد).

ضرایب پراکندگی α اختلاف معنی داری از صفر ندارد ($۰/۰۵e - ۲ = \alpha$). بنابراین در این داده‌ها مدل‌های رگرسیونی پواسن و دوجمله‌ای منفی برازشی مشابه دارند.

متغیرهای سطح تحصیلات مادر، تعداد فرزندان زندهٔ دختر، تعداد فرزندان زندهٔ پسر و استفاده از روش پیشگیری وارد مدل شدند و هیچ یک از اثرات متقابل وارد مدل نشدند.

۱.۳ برآورد پارامترها

ضرایب برآورده شده در مدل رگرسیون پواسن در جدول (۱) و ضرایب برآورده شده در مدل رگرسیون دوجمله‌ای منفی در جدول (۲) نشان داده شده است.

متغیر	ضرایب برآورده شده	ضرایب برآورده استاندارد	خطای استاندارد	آماره والد	P-Value	فاصله اطمینان ۹۵%
GRADE(2)	-۰/۱۹۶۵	۰/۱۲۶۸	-۱/۵۴۹	۰/۱۲۱	(-۰/۴۴۵۲, -۰/۰۵۲۱)	
GRADE(3)	-۰/۱۸۷۶	۰/۱۲۵۵	-۱/۴۹۵	۰/۱۳۵	(-۰/۴۲۳۶, -۰/۰۵۸۴)	
GRADE(4)	-۰/۲۶۹۵	۰/۱۲۲۷	-۲/۰۲۱	۰/۰۴۲	(-۰/۵۲۹۵, -۰/۰۰۹۵)	
GRADE(5)	-۰/۴۳۵۶	۰/۲۰۶۷	-۲/۱۰۷	۰/۰۲۵	(-۰/۸۴۰۷, -۰/۰۳۰۴)	
DAUT	۰/۴۵۲۱	۰/۰۳۴۱	۱۳/۲۵۴	۰/۰۰۰	(۰/۳۸۵۳, ۰/۵۱۹۰)	
SON	۰/۳۰۱۹	۰/۰۲۰۱	۱۵/۰۱۶	۰/۰۰۰	(۰/۲۶۲۵, ۰/۳۴۱۳)	
CONTRA(2)	-۰/۶۵۴۲	۰/۰۲۱۷	-۹/۱۲۹	۰/۰۰۰	(-۰/۷۹۴۷, -۰/۵۱۳۸)	
ثابت	-۱/۴۵۳۶	۰/۱۲۲۵	-۱۰/۹۷۲	۰/۰۰۰	(-۱/۷۱۳۲, -۱/۱۹۴۰)	

جدول ۱: برآوردهای حداکثر درستنمایی برای مدل پذیرفته شده رگرسیون پواسن

با توجه به ستون مقدار آماره «والد» در جدول (۱) به نتایج زیر می‌رسیم:

به نظر می‌رسد که مهمترین و مؤثرترین متغیر تعداد فرزندان زندهٔ پسر بوده است و پس از آن به ترتیب متغیرهای تعداد فرزندان زندهٔ دختر، و استفاده از روش پیشگیری تاثیر نسبتاً زیادی بر متغیر وابسته، یعنی تعداد حاملگی‌های ناخواسته داشته‌اند.

همچنین با توجه به ستون ضرایب برآورده شده به این نتیجه می‌رسیم که بین متغیرهای تعداد فرزندان زندهٔ پسر و تعداد فرزندان زندهٔ دختر و متغیر وابسته یک ارتباط مستقیم برقرار بوده است (برآورده پارامتر برای این دو متغیر مستقل، مثبت است) یعنی با افزایش

متغیر	برآورده شده	ضرایب برآورده شده	خطای استاندارد	آماره والد	P-value	فاصله اطمینان ۹۵%
GRADE(2)	-۰/۱۹۶۵	۰/۱۲۷۷	-۱/۵۳۹	۰/۱۲۴	(-۰/۴۴۶۸, -۰/۰۵۲۷)	
GRADE(3)	-۰/۱۸۷۶	۰/۱۲۶۵	-۱/۴۸۳	۰/۱۲۸	(-۰/۴۲۵۶, -۰/۰۶۰۳)	
GRADE(4)	-۰/۲۶۹۵	۰/۱۳۳۲	-۲/۰۲۳	۰/۰۴۳	(-۰/۵۳۰۵, -۰/۰۰۸۴)	
GRADE(5)	-۰/۴۲۵۶	۰/۲۱۲۹	-۲/۰۴۶	۰/۰۴۱	(-۰/۸۵۲۸, -۰/۰۱۸۴)	
DAUT	۰/۴۵۲۱	۰/۰۳۴۱	۱۲/۲۴۵	۰/۰۰۰	(۰/۲۸۵۲, ۰/۰۵۱۹۰)	
SON	۰/۳۰۱۹	۰/۰۲۰۱	۱۵/۰۲۸	۰/۰۰۰	(۰/۲۶۲۵, ۰/۳۴۱۲)	
CONTRA(2)	-۰/۶۵۴۲	۰/۰۷۰۶	-۹/۲۶۱	۰/۰۰۰	(-۰/۷۹۲۷, -۰/۰۵۱۵۸)	
ثابت	-۱/۴۵۲۶	۰/۱۳۲۷	-۱۰/۸۷۶	۰/۰۰۰	(-۱/۷۱۵۶, -۱/۱۹۱۷)	

جدول ۲: برآوردهای حداکثر درستنمایی برای مدل پذیرفته شده رگرسیون دوجمله‌ای منفی

تعداد فرزندان زنده پسر (یا تعداد فرزندان زنده دختر) خطر وقوع حاملگی ناخواسته بیشتر می‌شود.

برآورد پارامتر منفی برای متغیر سطح تحصیلات مادر و استفاده از روش پیشگیری، نشان دهنده این مطلب است که بین این متغیرها و متغیر واپسیه یک رابطه منفی وجود داشته است، یعنی با افزایش میزان سطح تحصیلات مادر، تعداد حاملگی‌های ناخواسته کاهش یافته است و چون متغیر استفاده از روش پیشگیری از حاملگی متغیری دوحتایی است (۱- پیشگیری داشته و ۲- پیشگیری نداشته)، وجود رابطه منفی به این معنی است که خطر وقوع حاملگی ناخواسته برای افرادی که پیشگیری داشته‌اند بیشتر از کسانی بوده که پیشگیری نداشته‌اند و این نشان دهنده وجود شکافی در آگاهی از کنترل بارداری، حتی در میان استفاده کنندگان از روشهای پیشگیری از بارداری است.

۲.۳ نسبت خطر

در این بخش می‌خواهیم نسبت خطرهای مربوط به متغیرهای مستقل طبقه‌ای را مورد بررسی قرار دهیم. لازم به ذکر است که در هر دو روش رگرسیونی پواسن و دوجمله‌ای منفی اولین سطح هر متغیر طبقه‌ای به عنوان مرجع درنظر گرفته شده و خطرات نسبی بقیه سطوح نسبت به این سطح سنجیده می‌شود. به شرط ثابت نگهداشتن سایر متغیرها نتایج زیر به دست آمده است:

در ابتدا متغیر مستقل سطح تحصیلات مادر را درنظر می‌گیریم. این متغیر دارای پنج رده می‌باشد. نسبت خطر را به صورت زیر درنظر می‌گیریم:

$$\pi = \frac{\text{نسبت حاملگی ناخواسته در یک گروه}}{\text{نسبت حاملگی ناخواسته در گروه مرجع}}$$

پس داریم:

$$\text{نسبت حاملگی ناخواسته در گروه بی سواد} / \text{نسبت حاملگی ناخواسته در گروه با سواد ابتدایی} = \pi$$

که برابر است با $8216/0$ و فاصله اطمینان 95% برای این نسبت خطر برابر است با $(1/0535, 1/0640)$ که نشان دهنده معنی داری این نسبت خطر در سطح خطای 5% است و خواهیم داشت: $1/2171 = 1/\pi$

بنابراین خطر وقوع حاملگی ناخواسته در مادران بی‌سواد نسبت به مادرانی که تحصیلاتی در حد ابتدایی دارند، $2/1$ برابر بیشتر است.

به همین ترتیب نتایج زیر به دست آمده است:

خطر وقوع حاملگی ناخواسته در مادران بی‌سواد نسبت به مادرانی که تحصیلاتی در حد متوسطه دارند، $2/1$ برابر بیشتر است.

خطر وقوع حاملگی ناخواسته در مادران بی‌سواد نسبت به مادرانی که تحصیلاتی در حد دیپلم دارند، $3/1$ برابر بیشتر است.

خطر وقوع حاملگی ناخواسته در مادران بی‌سواد نسبت به مادرانی که تحصیلاتی بالاتر از دیپلم دارند، $5/1$ برابر بیشتر است.

حال متغیر استفاده از روش پیشگیری از حاملگی را در نظر می‌گیریم. این متغیر دارای دو رده (۱ - پیشگیری داشته و ۲ - پیشگیری نداشته) می‌باشد.

$$\frac{\text{نسبت حاملگی ناخواسته در گروهی که پیشگیری داشته}}{\text{نسبت حاملگی ناخواسته در گروهی که پیشگیری نداشته}} = \pi$$

که برابر است با $5198/0$ و فاصله اطمینان 95% برای این نسبت خطر برابر است با $(1/05982, 1/04517)$ و داریم: $1/9238 = 1/\pi$

بنابراین خطر وقوع حاملگی ناخواسته در مادرانی که از روش‌های پیشگیری از بارداری استفاده می‌کنند نسبت به مادرانی که از روش‌های پیشگیری از بارداری استفاده نمی‌کنند، تقریباً دو برابر بیشتر است.

۴.۳ بحث

استیودنت (۱۹۰۷) جزء اولین افرادی بود که از مدل دوجمله‌ای منفی به عنوان صورت دیگری از توزیع پواسن در داده‌های شمارشی استفاده کرد. مک‌کدربیک (۱۹۱۴) و کیونویل (۱۹۴۹) توزیع‌های آمیخته در رابطه با توزیع دوجمله‌ای منفی را شرح دادند [۶].

گاردنر، مالوی و شاو (۱۹۹۵) مدل‌های پواسن و دوجمله‌ای منفی را، بر روی داده مربوط به وقوع خشونت در اشخاص مبتلا به ناخوشی‌های روانی، مورد مقایسه قرار

دادند. این داده‌ها از مطالعه‌ای مربوط به وقوع خشونت از جامعه‌ای شامل ۷۹۷ نفر که دچار ناخوشی‌های روانی بودند، برگرفته شده است. در این بررسی بیماران مراجعه کننده به بخش اورژانس یک بیمارستان روانی، مدت شش ماه پیگیری شده‌اند.

در مثال مذکور، داده‌ها تغییرات قابل توجهی را در وقوع خشونت بین بیماران نشان می‌دهد (وقوع $M = ۷/۳$, $SD = ۳$).

رگرسیون مربوط به داده‌های خشونت شامل سه متغیر پیش‌بینی کننده برای خشونت می‌باشد:

$$(1) \text{ سن بیماران } (X = ۲۸/۶ \text{ و سال } SD = ۱۱/۱)$$

(2) اندازه‌گیری تکنسین‌های بخش اورژانس درباره اطمینان از اینکه بیمار دچار خشونت شده است.

(3) اندازه‌گیری سابقه خشونت در بیماران قبل از اینکه به بخش اورژانس مراجعه کنند.

متغیرهای مذکور به ترتیب AGE، CONCERN و HISTORY نام‌گذاری شده‌اند.

در داده‌های مربوط به خشونت، واضح است که نمی‌توان رگرسیون پواسن با ثابت μ را برازش کرد چون $Var(y) = ۵۲/۶$ بسیار بزرگتر از $y = ۳$ می‌باشد و در تابع ضرایب برآورد شده به وسیله مدل رگرسیون پواسن (که در جدول (۳) نشان داده شده است) بسیار کوچکتر از مقدار حقیقی آنها است و آزمونهای t مربوط به این ضرایب متوجه شده است.

ضرایب برآورد شده به وسیله برازش مدل رگرسیون دوجمله‌ای منفی در جدول (۴) نشان داده شده است و همانطور که ملاحظه می‌شود نتایج حاصله بسیار بهتر از نتایج به دست آمده به وسیله مدل پواسن است.

بنابراین در مدل بندی داده‌های شمارشی ابتدا باید میانگین و واریانس متغیر وابسته را محاسبه کرد. در صورتی که میانگین و واریانس به طور تقریبی برابر باشند، برای آنالیز داده‌ها از مدل رگرسیون پواسن استفاده می‌کنیم و در غیر اینصورت مدل رگرسیون دوجمله‌ای منفی را مورد استفاده قرار می‌دهیم.

به عبارت دیگر برای استفاده از مدل رگرسیون دوجمله‌ای منفی، ابتدا باید مقدار α را برآورد کرد. در صورتی که α اختلاف معنی داری از صفر داشت، رگرسیون دوجمله‌ای منفی برای آنالیز داده‌ها مدل مناسبی خواهد بود و در غیر اینصورت بهتر است که از مدل پواسن برای مدل بندی داده‌ها استفاده کنیم. در ایران مطالعات مختلفی در زمینه حاملگی

ناخواسته صورت گرفته اما در هیچیک از آنها از رگرسیون پواسن و یا رگرسیون دوجمله‌ای منفی برای آنالیز داده‌ها استفاده نشده است.

متغیر	ضرایب برآورد شده	ضرایب استاندارد	خطای استاندارد	t	$Pr(> t)$
AGE	-۰/۰۴۵	۰/۰۰۲۳	۱۹/۶۹	۰/۰۰۰۱	
CONCERN	۰/۰۸۳	۰/۰۰۲۵	۱۱/۲۰	۰/۰۰۰۱	
HISTORY	۰/۰۴۲۰	۰/۰۳۸۰	۱۱/۲۶	۰/۰۰۰۱	
ثابت	-۲/۴۱۰	۰/۰۶۹۰	-۴۹/۲۹	۰/۰۰۰۱	

جدول ۳: آنالیز رگرسیون پواسن برای نرخ‌های خشونت

متغیر	ضرایب برآورد شده	ضرایب استاندارد	خطای استاندارد	t	$Pr(> t)$
AGE	-۰/۰۴۵۹	۰/۰۰۷۹۹	-۵/۷۴	۰/۰۰۰۱	
CONCERN	۰/۰۹۶۲	۰/۰۳۰۹۰	۳/۱۲	۰/۰۰۲۸	
HISTORY	۰/۰۵۳۶۰	۰/۱۵۵۰۰	۳/۴۵	۰/۰۰۰۸	
ثابت	-۲/۵۵۰۰	۰/۲۶۲۰۰	-۱۳/۵۵	۰/۰۰۰۱	

جدول ۴: آنالیز رگرسیون دوجمله‌ای منفی برای نرخ‌های خشونت

مراجع

- Mcculagh, P. and Nelder J.A. (1989), *Generalized Linear Models (2nd ed.)*, Chapman & Hall , London.
- Piegrosh, W.W. (1990), *Maximum Likelihood Estimation for the Negative Binomial Dispersion Parameter* , Biometrics., 49:863 P.
- Denton, B. , Scott, S. and Chase, B. (1994), *Unintend and Unwanted Pregnancy in Halifax: the Rate and Associated Factors* , Con.J.Public Health, 234 p.
- Gardner , W. , Mulvey E.P. and Shaw E.C. (1995), *Regression Analysis of Counts and Rates: Poisson, Overdispersed Poisson and Negative Binomial Models.*, Psychological Bulletin., 118(3):392 P .
- Selvin , S. (1995), *Practical Biostatistical Methods*, Duxbury, USA.

Poch, M. and Mannering, F. (1996), *Negative Binomial Analysis of Intersection Accident Frequencies*, Journal of Transportation Engineering., March/April 105 P.

آل ابریشم کارزاده مهران. کاربرد مدل رگرسیون پواسن در تعیین رابطه کیست تخدمان ناشی از کاشت نورپلانت با ریسک فاکتورهای مربوط. پایان نامه کارشناسی ارشد آمار حیاتی، تهران: دانشکده علوم پزشکی دانشگاه تربیت مدرس، ۱۳۷۸.

اندرسن ارلینگ بی. الگوهای آماری گستته و کاربرد آن در علوم اجتماعی. ترجمه دکتر علی مشکانی و دکتر ابوالقاسم بزرگ نیا. چاپ اول، مشهد: دانشگاه فردوسی، ۱۳۷۲.

امامی افشار نژهت. حاملگی ناخواسته و ارتباط چند متغیر دموگرافی- اجتماعی و روش‌های پیشگیری از بارداری در بین مراجعین به مرکز بهداشتی - درمانی جنوب تهران. محیط‌شناسی، زمستان ۱۳۷۰.

منصوریان فریفته. کاربرد مدل‌های عمومی لگ خطی در تعیین عوامل مؤثر بر حاملگی ناخواسته و پیامدهای آن. پایان نامه کارشناسی ارشد آمار حیاتی، تهران: دانشکده علوم پزشکی دانشگاه تربیت مدرس، ۱۳۷۵.

تعیین پارامتر پنجره در برآورد تابع چگالی توسط B -اسپلاین

دکتر محسن محمدزاده، رضا صالحی

P ۱۳۰ ۱۹

گروه آمار، دانشگاه تربیت مدرس

چکیده: اغلب در مسائل کاربردی تابع چگالی احتمال نامعلوم است و لازم است براساس مشاهدات حاصل از یک نمونه تصادفی برآورد شود. در آمار روش‌های مختلفی برای برآورد تابع چگالی بکار برد همیشه شود که در این مقاله شیوه برآورد B -اسپلاین مورد مطالعه قرار گرفته و روشی برای تعیین پارامتر پنجره، که میزان همواری برآورد را کنترل می‌کند، ارائه خواهد شد.

واژه‌های کلیدی: B -اسپلاین، پارامتر پنجره، برآورد چگالی احتمال.

۱ مقدمه

تابع چگالی احتمال یک مفهوم اساسی در آمار است که توسط آن رفتار متغیر تصادفی، احتمال پیشامدهای متناظر با آن و بسیاری از خواص متغیر تصادفی قابل توصیف می‌باشد. اما در اغلب مطالعات عملی این تابع نامعلوم است، در عوض یک مجموعه از n مشاهده x_1, \dots, x_n داریم، که فرض می‌شود مقادیر n متغیر تصادفی مستقل و هم توزیع و دارای تابع چگالی مجهول f هستند و معمولاً براساس این مشاهدات به یکی از دوروش پارامتری یا ناپارامتری برآورد می‌شود. در این مقاله ابتدا روش ناپارامتری برآورد یک تابع دلخواه توسط اسپلاین‌ها مورود شده سپس روش B -اسپلاین‌ها برای برآورد تابع چگالی احتمال مورد بررسی قرار می‌گیرد. چون یکی از مسائل مهم در روش B -اسپلاین‌ها تعیین پارامتر پنجره است، که میزان همواری تابع برآورد شده را کنترل می‌کند، روشی برای تعیین مقدار بهینه این پارامتر ارائه خواهد شد.

در تجزیه و تحلیل رگرسیون برای تعیین ارتباط بین متغیر پاسخ و متغیرهای مستقل و برآش منحنی به داده‌ها فرض‌های محدود کننده‌ای بمنظور حصول یک پاسخ منحصر بفرد در نظر گرفته می‌شود. در بسیاری از مواقع نوع ارتباط بین متغیرها مشخص نیست و بجای مفروض داشتن یک الگوی خاص پارامتری، از روشی استفاده می‌شود که داده‌ها ماهیت روند خود را بهتر نشان دهند. اسپلاین همواری روشی است برای برآورد منحنی

که در آن تنها فرض همواری منحنی در نظر گرفته می شود. با در نظر گرفتن مدل:

$$Y_i = g(t_i) + \epsilon_i \quad , \quad i = 1, \dots, n \quad (1)$$

که در آن g نامعلوم است، در حالت کلی نمی توان آنرا بطور منحصر بفرد برآورد کرد، مگر آنکه محدودیتی روی مدل اعمال شود. به عنوان مثال در روش‌های رگرسیون خطی فرض می شود g تابعی خطی از پارامترهای t است و به یکی از روش‌های برآوردهای پارامترهای آن بطور منحصر بفرد برآورده شوند. گذاشتند فرض خطی بودن تابع g ممکن است برای برخی از مجموعه داده ها مناسب باشد لیکن دارای کاستی هایی است که برای رهایی از آنها می توان g را تابعی دلخواه در نظر گرفت و با فرض هموار بودن، آن را برآورده نمود. میزان ناهمواری تابع g در بازه $[a, b]$ را می توان توسط عبارت $\int_a^b [g''(t)]^2 dt$ ^۲ اندازه گیری نمود. لذا می توان در برآورده منحنی هموار g ، علاوه بر عامل مجموع مربعات باقیمانده ها، میزان ناهمواری را نیز در نظر گرفته و ملاک مجموع مربعات جریمه ای

$$S(g, \alpha) = \sum_{i=1}^n (Y_i - g(t_i))^2 + \alpha \int_a^b [g''(t)]^2 dt \quad (2)$$

را که در آن $\alpha > 0$ پارامتر همواری نامیده می شود و تقابل بین نکوبی برازش منحنی به داده ها و میزان ناهمواری تابع g را کنترل می کند، برای برآورده g بکار برد. تابعی مانند، g که دارای مشتق مرتبه دوم پیوسته و کمینه کننده مجموع مربعات جریمه ای (۲) باشد، اسپلاین همواری^۱ نامیده می شود. این شیوه برآورده منحنی، رهیافت جریمه ناهمواری نامیده می شود که در اویانک (۱۹۸۸)، هاردل (۱۹۹۰) و راسنبلات (۱۹۹۱) مطالعه مبسوطی در مورد آن آمده است. تعیین مقدار مناسب پارامتر همواری α یک مسئله مهم است که برای مطالعه آن می توان به هاردل، مارون و هال (۱۹۸۸) و گرین و سیلورمن (۱۹۹۴) مراجعه نمود. یکی از روش‌های پیدا کردن پارامتر هموارسازی، روش اعتبار متقابل و اعتبار متقابل تعمیم یافته می باشد که در استون (۱۹۷۴)، کریون و واهبا (۱۹۷۹)، گرین و سیلورمن (۱۹۹۴) آمده است. محمدزاده (۱۹۹۸) الگوریتمی مناسب برای تعیین مقدار پارامتر هموارسازی بر اساس مشاهدات ارائه نموده است. برای مطالعه بیشتر در زمینه اسپلاین ها می توان به اویانک (۱۹۸۴) و شاماکر (۱۹۹۳) مراجعه نمود. وقتی مدل (۸) به روش حداقل مربعات باقیمانده ها برازش داده شود، حالات زیر رخ می دهد. اگر نقاط بوسیله خط راست به یکدیگر متصل شده باشند، مجموع مربعات باقیمانده ها صفر خواهند شد. اگر شرط همواری g را اعمال کنیم و منحنی دارای مشتق مرتبه دوم پیوسته باشد، منحنی از تمامی نقاط (t_i, Y_i) می گذرد. در این حالت نیز مجموع مربعات

^۱ Smoothing Spline

باقیمانده‌ها صفر است، اما منحنی موج و دارای نوسانات زیادی است. مشتق مرتبه دوم اندازهٔ ناهمواری منحنی را در نقاطی که منحنی تغییر جهت می‌دهد، مشخص می‌کند که ممکن است منفی یا مثبت باشد. مجموع توان دوم این مشتق‌ها روی تمام نقاط میزان ناهمواری منحنی را مشخص می‌کند. برای مقدار داده شده α ، کمینه کردن $(S(g, \alpha))^2$ بهترین راه توافق بین میزان همواری و نیکویی برازش مدل است. در روش اسپلاین همواری مقدار α نقش مهمی در برآورد منحنی ایفا می‌کند. اگر α بزرگ باشد، مولفه اصلی در $(S(g, \alpha))^2$ جملهٔ جریمهٔ ناهمواری است و بنابراین کمینه کننده g خیلی کم انجنا خواهد بود. در حد وقتی α به بی نهایت میل می‌کند، سهم عبارت $\int_a^b [g''(t)]^2 dt$ به سمت صفر میل می‌کند و منحنی \hat{g} همان برازش رگرسیون خطی خواهد بود. از طرف دیگر اگر α نسبتاً کوچک باشد، مجموع مربعات باقیمانده‌ها سهم اصلی را در $(S(g, \alpha))^2$ دارد و برآورد منحنی \hat{g} تا حد زیادی روند داده‌ها را دنبال می‌کند، که ممکن است نوسانات زیادی داشته باشد. در وضعیت حدی اگر α به صفر نزدیک شود، برآورد حاصل به اسپلاین درونیاب^۲ میل می‌کند. برآورد حداقل مربعات جریمه‌ای، \hat{g} ، کمینه کننده تابع $(S(g, \alpha))^2$ روی کلاس تمام توابع دوبار مشتق پذیر g تعریف می‌شود. برآورد منحنی \hat{g} که کمینه کننده تابع $(S(g, \alpha))^2$ است، اسپلاین نامیده می‌شود. تابع $S(t)$ را اسپلاین از درجه r با گره‌های ^۳ t_1, \dots, t_n گویند اگر $a = t_0 < t_1 < \dots < t_n < t_{n+1} = b$

۱) برای هر $n, \dots, 0, k = 0, \dots, n$ در فاصله (t_k, t_{k+1}) یک چند جمله‌ای درجه نابیشتر از r باشد.

۲) $S^{(r-1)}(t), S'(t), \dots, S(t)$ توابع پیوسته‌ای روی $[a, b]$ باشند.

هر چند جمله‌ای یک تابع اسپلاین بدون گره می‌باشد و بنابر تعریف بالا r امین مشتق یک اسپلاین از درجه r یک تابع ثابت تکه‌ای با شکستگی در نقاط t_n, \dots, t_1 است. بر عکس r امین تابع اولیه از یک تابع ثابت تکه‌ای یک اسپلاین از درجه r است. تابع g بر فاصله $[a, b]$ اسپلاین درجه سه نامیده می‌شود اگر برای اعداد حقیقی $a < t_1 < t_2 < \dots < t_n < b$ در دو شرط زیر صدق کند:

۱- روی هر یک از فاصله‌های $(a, t_1), (t_1, t_2), \dots, (t_n, b)$ یک چند جمله‌ای درجه سه باشد.

۲- قطعات چند جمله‌ای در نقاط t_i به نحوی برازش یابند که مشتق‌های اول و دوم آن در هر t_i پیوسته باشد که در نتیجه روی تمام بازه $[a, b]$ پیوسته است.

^۲ Interpolating Spline

^۳ Knots

۲—اسپلاین-B

اسپلاین‌ها رهیافت متنوعی برای بررسی منحنی‌ها در گرافهای کامپیوترا فراهم می‌کنند و مرکب از منحنی‌های چند جمله‌ای متصل در گره‌ها هستند و از درجه یکسان هستند، که درجه B-اسپلاین نامیده می‌شود. فرض کنید U_r, \dots, U_1 یک نمونه تصادفی از توزیع یکنواخت بر فاصله $[1/2, 1/2]$ باشد. چگالی مجموع این متغیرهای تصادفی بصورت

$$B^{(r)}(x) = \frac{\partial}{\partial x} P_r(U_1 + U_2 + \dots + U_r < x), \quad x \in R$$

است. هرتابع $B^{(r)}$ عضوی از کلاس $C^{(r-2)}$ یعنی توابع دارای مشتق پیوسته مرتبه $r-2$ با تکیه‌گاه $[-r/2, r/2]$ است. برای هر عدد صحیح i ، تابع $(.)^{(r)}$ که به بازه $(i+r/2, i+1+r/2)$ محدود شده است یک چند جمله‌ای بدیهی در خارج بازه $[-r/2, r/2]$ ، یا از درجه r در داخل آن است. تابع $(.)^{(r)}$ B-اسپلاین مرتبه r نامیده می‌شود. بعضی اسپلاین‌های اساسی بصورت زیر می‌باشند:

برای $r=1$

$$B^{(1)}(x) = \begin{cases} 1 & -\frac{1}{r} \leq x < \frac{1}{r} \\ 0 & \text{o.w.} \end{cases}$$

برای $r=2$ بنابراین $B^{(2)} = B^{(1)} * B^{(1)}$

$$B^{(2)}(x) = \begin{cases} 1+x & -1 \leq x < 0 \\ 1-x & 0 \leq x \leq 1 \\ 0 & \text{o.w.} \end{cases}$$

برای $r=3$ بنابراین $B^{(3)} = B^{(1)} * B^{(1)} * B^{(1)}$

$$B^{(3)}(x) = \begin{cases} \frac{1}{r}(x + \frac{r}{r})^2 & -\frac{r}{r} \leq x < -\frac{1}{r} \\ -x^2 + \frac{r}{r} & -\frac{1}{r} \leq x \leq \frac{1}{r} \\ \frac{1}{r}(x - \frac{r}{r})^2 & \frac{1}{r} < x \leq \frac{r}{r} \\ 0 & \text{o.w.} \end{cases}$$

والی آخر. واضح است که $B^{(1)}$ یک تابع پله‌ای واحد و ثابت است. بنابراین $B^{(2)}$ خطی تکه‌ای، $B^{(3)}$ درجه دو تکه‌ای می‌باشد.

حال بوسیله انتقال و تغییر مقیاس $B^{(r)}$ ، می‌توان B -اسپلاین را برای مجموعه‌ای از گره‌ها روی R که بطوریکسان تقسیم شده‌اند، بدست آورد:

$$B_{s,h}^{(r)}(x) = B^{(r)}\left(\frac{x}{h} - s\right), \quad x \in R, s \in Z, h > 0 \quad (3)$$

برای هر عدد صحیح مثبت دلخواه r ، هستهٔ اسپلاین بصورت

$$Q^{(r)}(x,y) = \sum_{s \in Z} B^{(r)}(x-s)B^{(r)}(y-s), \quad (y,x) \in R^2 \quad (4)$$

است، که برای $h > 0$ و $(x,y) \in R^2$ هستهٔ مقیاس ساز^۱ بصورت

$$Q_h^{(r)}(x,y) = \frac{1}{h} Q^{(r)}\left(\frac{x}{h}, \frac{y}{h}\right) \quad (5)$$

تعریف می‌شود. کرزیکوفسکی (۲۰۰۱) نشان داد هستهٔ اسپلاین (۴) دارای خواص زیر است:

$$\begin{aligned} Q^{(r)}(x,y) &= Q^{(r)}(y,x) = 0 & |y-x| > r \\ \int Q^{(r)}(x,y) dx &= 1 & y \in R \\ \int x Q^{(r)}(x,y) dx &= y & y \in R, r > 1 \\ \int x^r Q^{(r)}(x,y) dx &= y^r + \frac{r}{r+1} & y \in R, r > 2 \end{aligned}$$

۳ برآورد تابع چگالی

فرض کنید X_1, \dots, X_n یک نمونه تصادفی مستقل با مقادیر حقیقی و چگالی مشترک f باشد. برآوردگر هسته‌ای $\hat{f}_h(x)$ برای برآورد تابع چگالی $f(x)$ در نقطه x بصورت زیر تعریف می‌شود:

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)$$

^۱ Scaling kernel

که در آن $K(\cdot)$ تابع هسته‌ای نامیده می‌شود و می‌تواند هسته یکنواخت، مثلثی، درجه دو، گوسی، نمایی وغیره باشد و h پارامتر پهنانی بازه است و میزان همواری برآورده را کنترل می‌کند.

برآورده سری برای تابع f بصورت زیر تعریف می‌شود:

$$\hat{f}(x) = \frac{1}{n} \sum_{k=1}^n K_m(x - x_k)$$

که در آن

$$K_m(y) = \frac{1}{2\pi} \sum_{j=-m}^m e^{ijy}$$

و m تعداد دوره‌های تناوب 2π است یعنی برای $1 = -\sqrt{-1}$ فاصله $(-\pi, \pi]$ و برای $2 = m$ فاصله $[2\pi, -2\pi]$ بازه تغییرات داده‌ها برای برآورد تابع چگالی می‌باشد.

برآورده اسپلاین تابع چگالی بصورت زیر تعریف می‌شود:

$$\hat{f}_{n,h}(x) = \frac{1}{n} \sum_{j=1}^n Q_h^{(r)}(X_j, x), \quad h \in R_+, \quad x \in R, \quad r \geq 1 \quad (1)$$

این برآورد نه از نوع برآوردهای هسته‌ای و نه از نوع سری می‌باشد، اما بعضی خواص این دو برآورده را داراست. تشابه آن با برآوردهای سری بصورت

$$\hat{f}_{n,h}(x) = \sum_{s \in Z} a_{s,n,h} B_{s,h}^{(r)}(x), \quad x \in R, \quad h \in R_+$$

است، که در آن

$$a_{s,n,h} = \frac{1}{n} \sum_{j=1}^n \frac{1}{h} B_{s,h}^{(r)}(X_j)$$

به عبارت دیگر $\hat{f}_{n,h}(\cdot)$ یک ترکیب خطی از B -اسپلاین است. برآورده (1) بوسیله سیسیلیسکی (۱۹۸۷) داده شده و برای بررسی عمیق‌تر این برآورده و همچنین برآورد چگالی اسپلاین چند متغیره می‌توان به سیسیلیسکی (۱۹۹۰) مراجعه نمود.

۴ تعیین پارامتر پنجره

در این بخش راه حلی برای تعیین پارامتر پنجره $h \in R_+$ در رابطه (۶) ارائه می‌شود برای این منظور میانگین نمونه‌ای $\bar{X}_n = \frac{1}{n} \sum_{j=1}^n X_j$, یعنی برآوردگر نااریب با واریانس کمینه برای میانگین $\sigma_n^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2$ و $m = \int x f(x) dx$ با واریانس نااریب با واریانس کمینه برای واریانس $\sigma^2 = \int (x - m)^2 f(x) dx$ را در نظر می‌گیریم. با توجه به خواص هسته اسپلاین برای تمام $\Omega \in \omega$ داریم :

$$\int \hat{f}_{n,h}(x) dx = \frac{1}{n} \sum_{j=1}^n \int \frac{1}{h} Q^{(r)}\left(\frac{X_j}{h}, \frac{x}{h}\right) dx = 1$$

در نتیجه برآوردگر چگالی (۶) یک چگالی احتمال است. میانگین نمونه‌ای \bar{X}_n و واریانس نمونه‌ای σ_n^2 بترتیب میانگین و واریانس چگالی $(.) \hat{f}_{n,h}$ هستند، پس داریم :

$$\begin{cases} \int x \hat{f}_{n,h}(x) dx = \bar{X}_n \\ \int (x - \bar{X}_n)^2 \hat{f}_{n,h}(x) dx = \sigma_n^2 \end{cases}$$

این سیستم پایه‌ای برای تخمین پارامتر پنجره $h \in R_+$ خواهد بود. قضیه زیر روشی برای بدست آوردن پارامتر پنجره‌ای ارائه می‌دهد.

قضیه : اگر $\{Q_h^{(r)} : h \in R_+\}$ خانواده هسته‌های رابطه (۵) و X_1, \dots, X_n یک نمونه تصادفی باشد، آنگاه برای برآوردگر $\hat{f}_{n,h}$ در رابطه (۶) داریم :

$$\int x \hat{f}_{n,h}(x) dx = \bar{X}_n$$

و درست یک $h_n = h_n(X_1, \dots, X_n) > 0$ وجود دارد بطوریکه

$$\int (x - \bar{X}_n)^2 \hat{f}_{n,h}(x) dx = \sigma_n^2$$

و

$$h_n = h_n(X_1, \dots, X_n) = \frac{\sigma_n}{\sqrt{kn}}, \quad k = \frac{r}{\gamma} \quad (7)$$

اثبات : بنا به خواص هسته اسپلاین داریم :

$$\begin{aligned} \int x \hat{f}_{n,h}(x) dx &= \frac{1}{n} \sum_{j=1}^n \int \frac{x}{h} Q^{(r)}\left(\frac{X_j}{h}, \frac{x}{h}\right) dx \\ &= \frac{1}{n} \sum_{j=1}^n \int u h Q^{(r)}\left(\frac{X_j}{h}, u\right) du \end{aligned}$$

$$= h \frac{1}{n} \sum_{j=1}^n \frac{X_j}{h}$$

$$= \bar{X}_n$$

و

$$\begin{aligned} \int (x - \bar{X}_n)^r \hat{f}_{n,h}(x) dx &= \frac{1}{n} \sum_{j=1}^n \int x^r \frac{1}{h} Q^{(r)}\left(\frac{X_j}{h}, \frac{x}{h}\right) dx - (\bar{X}_n)^r \\ &= \frac{1}{n} \sum_{j=1}^n h^r \int u^r Q^{(r)}\left(\frac{X_j}{h}, u\right) du - (\bar{X}_n)^r \\ &= \frac{1}{n} \sum_{j=1}^n h^r \left(\left(\frac{X_j}{h}\right)^r + \frac{r}{h}\right) - (\bar{X}_n)^r \\ &= \frac{rh^r}{h} + S_n^r \end{aligned}$$

که در آن $\frac{rh^r}{h} + S_n^r = \sigma_n^r$ است. چون $S_n^r = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^r$ بنابراین

$$\frac{r}{h} h^r + \frac{n-1}{n} \sigma_n^r = \sigma_n^r$$

در نتیجه

$$h^r = \frac{\sigma_n^r}{rn} = \frac{\sigma_n^r}{n \frac{r}{h}}$$

اگر قرار دهیم $k = \frac{r}{h}$ آنگاه :

$$h = h_n = h_n(X_1, \dots, X_n) = \frac{\sigma_n}{\sqrt{nk}}$$

و حکم ثابت است.

بنابراین برآورده‌گر چگالی براساس نمونه تصادفی X_1, \dots, X_n به صورت

$$f_n(x) = \hat{f}_{n,h}(x) \quad x \in R \quad (\lambda)$$

است، که در آن $h = h_n$ بوسیله رابطه (۷) داده شده است. کرزیکوفسکی (۲۰۰۱) نشان داد که این برآورده‌گر سازگاری قوی در نرم L_1 دارد یعنی اگر X_1, \dots, X_n یک نمونه تصادفی با تابع چگالی f باشد، بطوریکه $E(X_1^r) < \infty$ در اینصورت برای تمام $x \in R$ با احتمال یک داریم

$$\|f_n - f\|_1 \rightarrow 0$$

۵ بحث و نتیجه گیری

برآوردهای ایجاد شده از B -اسپلاین‌ها تعمیمی از برآوردهای چگالی هیستوگرام‌ها می‌باشد. بطور مثال برای $r = 1$ B -اسپلاین درجه صفر، همان هیستوگرام‌ها می‌باشند. از آنجا که برآوردهای چگالی بر مبنای اسپلاین‌هاست، درجه‌ای از همواری بسادگی کنترل می‌شود. معمولاً در عمل $r = 1, 2, 3, 4$ در نظر می‌گیرند که برای $r = 2$ یک چند جمله‌ای خطی، برای $r = 3$ یک چند جمله‌ای درجه دو و برای $r = 4$ یک چند جمله‌ای درجه سه می‌باشد. با توجه به نمونه تصادفی داده شده، برآوردهای چگالی احتمال با استفاده از B -اسپلاین‌های خطی، درجه دو یا درجه سه بسته به اینکه هر کدام خطای L کمتری داشته باشند صورت می‌گیرد.

برآوردهای چگالی احتمال برای یک نمونه تصادفی مستقل X_1, \dots, X_n بصورت ترکیب خطی از B -اسپلاین‌ها می‌باشد و زمانی که تعداد نمونه افزایش پیدا کند یعنی n بقدر کافی بزرگ باشد بر چگالی احتمال واقعی جامعه منطبق می‌شود.

مراجع

- Ciesielski Z. (1987), Local Spline Approximation and Nonparametric Density Estimation, International Conference on Constructive Theory of Function, 25-31 May Varna.
- Ciesielski Z. (1990), Asymptotic Nonparametric Spline Density Estimation, Probab. Math. Statist. 12, pp. 1-24.
- Ciesielski Z. (1990), Asymptotic Nonparametric Spline Density Estimation in Several Variables. International Series of Numerical Mathematics (Birkhauser Verlag Basel) Vol. 94, pp. 25-53.
- Creven, P. and Wahba, G. (1979), Smoothing Noisy Data with Spline Function, Number Math, 31, 377-390.

- Eubank, R. L. (1984), Approximate Models and Spline, COMMUN. STATIS. Theory. Meth., 13(4), 433-484.
- Eubank, R. L. (1988), Splines Smoothing and Nonparametric Regression, New York, Marcel Dekker.
- Green, P. J. and Silverman, B. W. (1994), Nonparametric Regression and Generalized Linear Model, London, Chapman and Hall.
- Hardle, W. (1990), Applied Nonparametric Regression, Cambridge: Cambridge University Press.
- Hardle, W., Marron, J. S. and Hall (1988), How Far Are Automatically chosen regression Smoothing Parameters from their Optimum, Jornal of the American Statistical Association, Vol. 83, No. 401, Theory and methods.
- Krzykowski G. (2001), On Automatic Choice of the Window Parameter in the Nonparametric Density Estimation, University of Gdansk. Poland.
- Mohammadzadeh, M. (1998), An Algorithm to Find the Smoothing Parameter in Smoothing Splines, Proceeding of the 4th Iranian Statistical Conference, Shahid Beheshti University , Tehran.
- Rossernblatt, M. (1991), Stochastic Curve Estimation, NSFCBMS Regional Conference Serries in Probability and Statistics, Vol 3, Hayward, California: Institute of Mathematical Statistics.
- Schumaker, L. L. (1993), Splines Functions: Basic Theory. Florida, Krieger.
- Stone, M. (1974), Cross Validatory Choice and Assessment of Statistical Prediction, Sco B.

استفاده از کریگینگ عام در همه‌گیری‌شناسی جغرافیایی بیماریها

محسن محمدزاده، انوشیروان کاظم‌نژاد، سقراط فقیه‌زاده، یدالله واقعی

دانشگاه تربیت مدرس

چکیده: یکی از ابزارهای مهم برای تجزیه و تحلیل بیماریها در مبحث همه‌گیری‌شناسی جغرافیایی تهیه نقشه بیماری می‌باشد. اما نوسانات یا اغتشاشات ناشی از عواملی مانند خطاهای اندازه‌گیری و تغییرات سریع میزان بیماری در مناطق مختلف، تعبیر و تفسیر نقشه را مشکل می‌سازد. لذا برای حذف نوسانات شدید لازم است حتی‌امکان نقشه‌های هموار بدست آورده شوند. نظر به اینکه میزان بیماری در مناطق مختلف اغلب دارای همبستگی فضایی هستند، در این مقاله روش کریگینگ عام برای تهیه نقشه آماری بیماری همواری‌ها معرفی شده و به کمک آن نقشه میزان بروز بیماری سل ریوی برای سالهای ۱۳۷۷ و ۱۳۷۸ شهرستانهای کشور تهیه و روند جغرافیایی بیماری در این دو سال مقایسه خواهد شد.

کلید واژه: نقشه آماری، کریگینگ عام، نالیستیابی، سل ریوی

۱ مقدمه

معمولًاً در همه‌گیری‌شناسی جغرافیایی از نقشه‌های بروز یا شیوع بیماریها برای نمایش توزیع جغرافیایی و تجزیه و تحلیل بیماری در بهداشت ملی - منطقه‌ای استفاده می‌شود. یکی از قدیمی‌ترین موارد استفاده از نقشه در همه‌گیری‌شناسی، مطالعه جان اسنو (۱۸۵۴) درباره همه‌گیری وبا در شهر لندن می‌باشد، که با استفاده از ترسیم موقعیت مرگ افراد و موقعیت جغرافیایی پمپ‌های تأمین آب، تأثیر منابع تأمین آب در انتشار وبا را مورد بررسی قرار داد. نقشه‌های بیماری برای دستیابی به تغییرات جغرافیایی کمبودهای بهداشتی و در نتیجه اختصاص منابع مالی و انسانی و تنظیم سیاستهای بهداشتی مورد استفاده قرار می‌گیرند. همچنین برای مطالعه ارتباط بین بروز بیماری با متغیرهای توضیحی مانند عوامل اقتصادی، اجتماعی و فرهنگی می‌توانند مفید باشند. این ارتباط توسط لاوسن و کرسی (۲۰۰۰) و لاوسن (۲۰۰۱) از طریق رگرسیون فضایی و تحت عنوان تجزیه و تحلیل اکولوژیکی مورد مطالعه قرار گرفته است. نقشه‌های معمولی میزان بروز یک

بیماری، معمولاً براساس اطلاعات موجود در بعضی مناطق تهیه می‌شوند. عدم وجود میزان بیماری در بعضی مناطق یا وجود تغییرات جغرافیایی خطای اندازه‌گیری ممکن است اینگونه نقشه‌ها را دچار ناهمواری و نوسانات زیادی کند که در اینصورت تعبیر و تفسیر آنها بسیار دشوار خواهد بود. لذا برای رفع این مشکل لازم است، نقشه آماری به گونه‌ای هموار شود تا بتوان از آن برای تخمین میزان بیماری در مناطق فاقد اطلاع استفاده نمود.

در سه دهه گذشته روش‌های مختلفی به منظور تخمین مقدار یک متغیر مانند میزان بیماری در موقعیتی مشخص با استفاده از مشاهدات واقع در موقعیت‌های مجاور توسعه یافته است، که مهمترین آنها عبارتند از: روش عکس فاصله موزون^۱ (ریپلی، ۱۹۸۱) و واتسون و فیلیپ، ۱۹۸۵ که توسط پوکالا (۱۹۸۹) و گلاتر (۱۹۸۹) برای تهیه نقشه آماری بیماری مورد استفاده قرار گرفته است، روش سطوح روند چندجمله‌ای^۲ (واتسون، ۱۹۷۲ و ریپلی، ۱۹۸۱)، روش اسپلاین‌ها^۳ (دیبور، ۱۹۷۸ و لیند و همکاران، ۱۹۹۵) و روش‌های مبتنی بر بیز (کلایتون و کالدر، ۱۹۸۷، کلایتون و برناردینلی، ۱۹۹۲، برناردینلی و همکاران، ۱۹۹۵ و ماگلین و همکاران، ۲۰۰۰). روش کریگینگ از جمله فنون آماری است که برای تخمین فضایی و تهیه نقشه آماری بیماری به کار برده می‌شود و علیرغم روش‌های مذکور، ساختار همبستگی فضایی میزانهای بیماری را در تجزیه و تحلیل آنها مورد استفاده قرار می‌دهد. علاوه با استفاده از این روش می‌توان در هر موقعیت دلخواه انحراف معیار تخمین میزان بیماری را نیز تعیین کرد و در صورت نیاز نقشه آن را تهیه نمود (کارات و والرون، ۱۹۹۲ و کرسی، ۱۹۹۳).

چون میزانهای بروز بیماری در مناطق مختلف اغلب از یک ساختار همبستگی پیروی می‌کنند، که به موقعیت آنها بستگی دارد، نوعی داده فضایی^۴ هستند و این ساختار همبستگی باید در تجزیه و تحلیل آنها لحاظ گردد. در این مقاله نحوه تعیین ساختار همبستگی داده‌ها و استفاده از روش کریگینگ برای تخمین فضایی میزانهای بیماری بیان شده و نقشه هموار بیماری سل ریوی سالهای ۱۳۷۷ و ۱۳۷۸ با استفاده از داده‌های سل ریوی اسمیر مثبت شهرستانهای کشور که توسط اداره کل مبارزه با بیماریها جمع آوری شده است، تهیه می‌شود.

^۱ Inverse Distance-Weighted

^۲ Trend Polynomial Surfaces

^۳ Splines

^۴ Spatial Data

۲ مفاهیم و روش‌های آماری

داده‌هایی که نوعاً بر حسب موقعیت (مکان) قرار گرفتن آنها در فضای مورد مطالعه همیشه باشند و این همبستگی تابعی از فاصله موقعیت آنها باشد، داده‌های فضایی نامیده می‌شوند. میزان بروز یا شیوع یک بیماری در شهرستانهای مختلف و میزان گازهای آلاینده هوا در استگاههای مختلف سنجش آنودگی هوا مجموعه داده‌های فضایی، در حیطه‌های علوم پزشکی و محیط زیست هستند.

دو مشخصه داده‌های فضایی را از داده‌های معمولی تمایز می‌کند، اولاً هر داده فضایی با موقعیت آن در فضای مورد مطالعه همراه است و ثانیاً این داده‌ها به هم وابسته و این وابستگی مرتبط با موقعیت قرار گرفتن آنها است. داده‌های فضایی با نماد $z(\mathbf{t}_1, \dots, \mathbf{t}_n)$ نشان داده می‌شود، که در آن n تعداد داده‌ها و (x_i, y_i) طول و عرض متناظر با موقعیت جغرافیایی مشاهده \mathbf{t}_i (در فضای دو بعدی) است. ساختار همبستگی داده‌های فضایی به وسیله تابع تغییرنگار^۱ مشخص می‌شود. این تابع همبستگی داده‌های فضایی را بر حسب فاصله نشان می‌دهد و معمولاً مقدار آن با زیاد شدن فاصله افزایش می‌یابد، که نشان دهنده کاهش همبستگی بین داده‌ها می‌باشد. ماترون (۱۹۶۳) برآورد تجربی تغییرنگار را به صورت

$$2\hat{\gamma}(\mathbf{h}) = \frac{1}{|N(\mathbf{h})|} \sum_{i=1}^{|N(\mathbf{h})|} (z(\mathbf{t}_i) - z(\mathbf{t}_i + \mathbf{h}))^2 \quad (1)$$

معرفی نمود، که در آن $\{(\mathbf{t}_i, z(\mathbf{t}_i)) : \mathbf{t}_i - \mathbf{t}_j \in N(\mathbf{h})\}$ مجموعه تمام زوج داده‌های در فاصله \mathbf{h} از یکدیگر و $|N(\mathbf{h})|$ تعداد زوجهای تمایز این مجموعه می‌باشد.

در حالت کلی ساختار همبستگی داده‌ها به فاصله و جهت موقعیت آنها در فضای مورد مطالعه بستگی دارد. اگر همبستگی داده‌ها فقط تابعی از اندازه فاصله بین موقعیت آنها باشد، اصطلاحاً همسانگرد^۲ نامیده می‌شوند. برای برآش یک مدل تغییرنگار مناسب به داده‌ها معمولاً فاصله مورد مطالعه به تعدادی زیرفاصله یا لگ^۳ تقسیم می‌شود و با استفاده از زوج داده‌ایی که در هر لگ قرار می‌گیرند مقدار تغییرنگار تجربی محاسبه می‌شود و به یکی از روش‌های حداقل مربعات معمولی، حداقل مربعات وزون یا درستنما می‌باشد. یک مدل معتبر به آنها برآش داده می‌شود. نظر به اینکه انتخاب تعداد لگ در میزان دقت

^۱ Variogram

^۲ Isotropic

^۳ Lag

برازش مدل تغییرنگار مؤثر است، محمدزاده و همکاران (۱۳۸۰) تعداد لگ بهینه را برای برازش مدل تغییرنگار معروفی کردند.

کرسی (۱۹۹۳) نشان داد چنانچه داده‌های فضایی دارای روند باشند، برآورد تغییرنگار اریب است و برای رفع این نقصه برآورد تغییرنگار براساس داده‌های بدون روند شده را توصیه نمود. برای این کار یک منحنی رگرسیون چندجمله‌ای به موقعیت‌های (x_i, y_i) برازش داده می‌شود و از تفاضل مقدار منحنی روند از هر مشاهده، داده‌های بدون روند شده به دست می‌آیند و تغییرنگار تجربی آنها محاسبه می‌شود. لیکن تغییرنگار داده‌های اصلی و داده‌های بدون روند شده یکسان نیست و به روش حذف روند از داده‌ها بستگی دارد. اما محمدزاده (۱۳۷۸) نشان داد استفاده از تابع کوواریانس تعیین یافته^۱ داده‌های بدون روند شده، مشکل تعیین ساختار همبستگی داده‌ها را تا حدی مرتفع می‌سازد.

ماترون (۱۹۶۳) تخمین فضایی در یک موقعیت مشخص $(x_0, y_0) = t_0$ براساس داده‌های $(t_1, z_1), \dots, (t_n, z_n)$ را به نام دی. جی. کریگ، مهندس معدن آفریقای جنوبی، کریگینگ نامید. کریگینگ همان بهترین تخمین کنندهٔ خطی نالریب (*BLUP*) به صورت $\sum_{i=1}^n \lambda_i z(t_i)$ است، که در آن ضرایب $\lambda_1, \dots, \lambda_n$ به گونه‌ای تعیین می‌شوند که تخمین کننده، نالریب و دارای کمترین واریانس باشد. ضرایب λ در کریگینگ به گونه‌ای خواهند بود که به مشاهدات نزدیک وزن بیشتر و به مشاهدات دور وزن کمتر اختصاص داده می‌شود. این امر به تبعیت از ماهیت داده‌های فضایی است که با افزایش فاصله، همبستگی آنها کاهش می‌یابد.

کریگینگ شامل سه نوع متداول ساده، عادی و عام است که ماترون (۱۹۶۹ و ۱۹۷۳) و کرسی (۱۹۹۳) آنها را مورد بررسی قرار داده‌اند. کریگینگ عادی در مواقعي به کار می‌رود که داده‌ها فاقد روند فضایی باشند و در صورت وجود روند از روش کریگینگ عام استفاده می‌شود. جرئیات بیشتر و محاسبات کریگینگ عام را می‌توان در کرسی (۱۹۹۳) و ویستر و الیور (۲۰۰۱) ملاحظه نمود.

۳ تعیین ساختار همبستگی میزان بروز سل ریوی

مختصات جغرافیایی مراکز شهرستانها نسبت به یک مبدأ واقع در جنوب غربی کشور به صورت $(x_i, y_i) = t_i$ بر حسب کیلومتر از نقشه تقسیمات کشوری ۱۳۷۸ استخراج شده است. سپس از تقسیم تعداد موارد جدید سالیانه سل ریوی اسمیر مثبت برآورد جمعیت ۲۶۲ شهرستان کشور میزان بروز بیماری سل ریوی اسمیر مثبت براساس صد هزار نفر

^۱ Generalized Covariance Function

جدول ۱: مجموع مربعات موزون (WSS) برای مدل‌های تغییرنگار برآورد شده به میزانهای بروز سل ریوی

مدل									سال
Cubic	Gneiting	Powerd Exponential	Circular	Sphricaal	Gaussian	Wave	Matern	Exponential	
۴۲,۷۹	۱۲۶/۴۵	۷۴,۰۷	۱۰۱/۰۱	۱۱۲,۴۰	۵۴,۸۳	۴۴,۰۵	۸۰/۱۲	۸۰/۱۲	۱۳۷۷
۵۷,۰۲	۱۵۵/۲۵	۹۴,۲۹	۱۲۶,۵۵	۱۲۳,۸۲	۷۰/۲۳	۵۹,۴۲	۱۰۱,۸۲	۱۰۱/۸۲	۱۳۷۸

برای سالهای ۱۳۷۷ و ۱۳۷۸ محاسبه و با (t_i) نشان داده شده است.

بدلیل وجود روند فضایی در داده‌های مذکور با برآورد چندجمله‌ای درجه دو به صورت $R(t_i) = a_0 + a_1 x_i + a_2 y_i + a_3 x_i^2 + a_4 y_i^2 + a_5 x_i y_i$ ، روند هریک از مجموعه داده‌های مربوط به سالهای ۱۳۷۷ و ۱۳۷۸ تعیین و داده‌ها به صورت $R(t_i) = z(t_i)$ بدون روند شده‌اند. سپس تغییرنگار تجربی همسانگرد داده‌های بدون روند شده در $\hat{\mu}$ (محمدزاده و همکاران، ۱۳۸۰) محاسبه شده است. هشت مدل متداول تغییرنگار (باری و همکاران، ۱۹۹۷) به روش حداقل مربعات موزون، به تغییرنگار تجربی برآش داده و مجموع مربعات موزون (WSS) مدل‌های برآورد شده به تفکیک سال ۱۳۷۷ و ۱۳۷۸ در جدول ۱ نشان داده شده است. معیار مجموع مربعات موزون، اختلاف مدل برآورد شده با تغییرنگار تجربی را اندازه می‌گیرد و هر چه مقدار آن کمتر باشد مدل برآش شده بهتر خواهد بود. با توجه به نتایج مندرج در جدول ۱ مدل شبیه تغییرنگار Cubic همسانگرد

$$\gamma(|h|) = \begin{cases} c_0 + c \left[7\left(\frac{|h|}{a}\right)^2 - 8/75\left(\frac{|h|}{a}\right)^3 + 3/5\left(\frac{|h|}{a}\right)^5 - 5/75\left(\frac{|h|}{a}\right)^7 \right] & , \quad |h| < a \\ c_0 + c & , \quad |h| \geq a \end{cases}$$

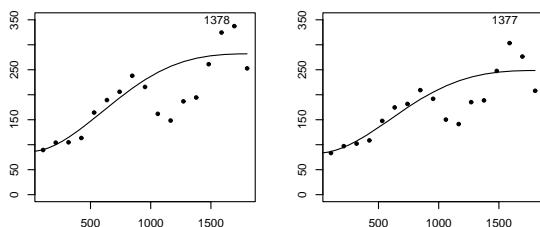
به عنوان بهترین مدل انتخاب شد، که مقادیر پارامترهای c_0 و a به همراه مجموع مربعات موزون (WSS) به تفکیک سال در جدول ۲ نشان داده شده است. شکل ۱ نمودار تغییرنگار تجربی و مدل برآورد شده Cubic همسانگرد را برای سالهای مختلف نشان می‌دهد.

۴ تخمین میزان بروز سل ریوی در چند شهرستان جدید

طبق تقسیمات کشوری سال ۱۳۷۸، جمهوری اسلامی ایران دارای ۲۸۲ شهرستان است که داده‌های سل ریوی ۲۶۲ شهرستان موجود و اطلاعی از میزان بیماری در ۲۰ شهرستان

جدول ۲: برآورد پارامترهای مدل تغییرنگار Cubic برای میزانهای بروز سل ریوی

WSS	\hat{a}	\hat{c}	\hat{c}_0	سال
۴۲,۷۹	۱۹۶۷,۲۵	۸۱,۱۶	۴۲,۸۸	۱۳۷۷
۵۷,۰۲	۱۹۶۹,۸۷	۹۵,۳۲	۴۴,۵۴	۱۳۷۸



شکل ۱: تغییرنگار تجربی و مدل Cubic برازش شده برای سالهای ۱۳۷۷ و ۱۳۷۸

دیگر در دسترس نمی‌باشد. به منظور نمایش قابلیتهای روشن کریگینگ و تکمیل اطلاعات بیماری مطابق با تقسیمات کشوری، میزان بروز بیماری سل در این شهرستانها برای سالهای ۱۳۷۷ و ۱۳۷۸ به روش کریگینگ عام تخمین زده شده و به همراه انحراف معیار آنها در جدول ۳ نشان داده شده است.

هنگامی که در چند موقعیت جدید تخمین انجام می‌شود، چون در این موقعیتها هیچ مشاهده‌ای در دسترس نیست نمی‌توان خطای پیش‌بینی و یا خطای کلی مدل تغییرنگار را ارزیابی کرد. در اینگونه موارد می‌توان به کمک اعتبارسنجی متقابل^۱ (استون، ۱۹۷۴) معیاری را برای دقت تخمینها ارائه داد. در این شیوه یک مشاهده حذف و مقدار آن براساس $n - 1$ مشاهده باقیمانده به روش کریگینگ تخمین زده می‌شود و از اختلاف مقدار واقعی و تخمین آن، خطای کریگینگ در بک موقعیت محاسبه می‌شود. این عمل را برای تمام n موقعیت تکرار و جذر میانگین مربع اختلافهای استاندارد به صورت:

$$RMSP = \left\{ \frac{1}{n} \sum_{j=1}^n \left[\frac{z(t_j) - \hat{z}_{-j}(t_j)}{\sigma_{-j}(t_j)} \right]^2 \right\}^{\frac{1}{2}}$$

محاسبه می‌شود، که در آن (t_j) مقدار واقعی مشاهده در موقعیت t_j و (t_j) \hat{z} تخمین آن براساس بقیه مشاهدات بجز (t_j) می‌باشد و σ_{-j} انحراف معیار تخمین (t_j)

^۱ Cross-Validation

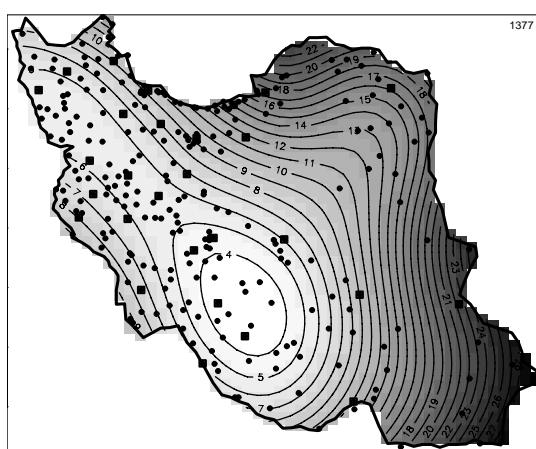
جدول ۳: تخمین فضایی میزان بروز بیماری سل ریوی در ۲۰ شهرستان

نام شهرستان	طول	عرض	سال ۱۳۷۸		سال ۱۳۷۷		میزان انحراف معیار	میزان انحراف معیار
			سال ۱۳۷۸	سال ۱۳۷۷	میزان	انحراف معیار		
اسکو	۲۱۹/۲۰	۱۳۰۸/۸۰	۶/۱۲	۶/۶۶	۴/۸۸	۶/۸۰	۶/۸۰	۵/۷۴
جلفا	۱۷۲/۸۰	۱۴۲۰/۸۰	۷/۲۹	۶/۸۰	۵/۷۴	۶/۹۴	۶/۹۴	۷/۱۸
چالدران	۵۹/۲۰	۱۴۳۲/۰۰	۶/۸۵	۷/۰۲	۴/۶۲	۷/۰۲	۷/۱۸	۷/۵۲
کوثر	۴۰۸/۰۰	۱۲۸۴/۸۰	۷/۷۶	۶/۶۴	۷/۵۲	۷/۷۷	۷/۷۷	۶/۷۶
نیر	۳۷۹/۲۰	۱۳۳۲/۸۰	۸/۰۸	۶/۶۶	۷/۷۷	۷/۷۹	۷/۷۹	۶/۷۶
ابجرود	۴۲۵/۶۰	۱۱۶۰/۰۰	۷/۲۶	۶/۶۳	۶/۳۹	۶/۷۶	۶/۷۶	۶/۷۶
خرمده	۵۰۸/۸۰	۱۱۳۶/۰۰	۶/۸۹	۶/۶۳	۷/۳۷	۷/۷۶	۷/۷۶	۷/۷۶
املش	۵۸۷/۲۰	۱۴۰۴/۸۰	۱۴/۰۳	۶/۸۱	۹/۲۸	۶/۶۵	۶/۷۸	۱۴/۴۲
رضوانشهر	۴۸۸/۰۰	۱۲۸۸/۰۰	۹/۱۱	۶/۶۵	۹/۲۸	۶/۷۸	۶/۷۸	۹/۹۰
ساهکل	۵۵۶/۸۰	۱۲۴۴/۸۰	۹/۵۱	۶/۶۴	۶/۶۴	۶/۶۷	۶/۶۷	۹/۹۰
ماسال	۴۸۸/۰۰	۱۲۶۸/۸۰	۸/۷۳	۶/۶۴	۸/۹۱	۶/۷۷	۶/۷۷	۸/۹۱
جوپیار	۸۲۷/۲۰	۱۲۱۹/۲۰	۱۴/۱۱	۶/۶۶	۱۴/۲۷	۶/۷۹	۶/۷۹	۱۴/۲۷
بندرگز	۹۲۰/۰۰	۱۲۴۴/۸۰	۱۶/۶۹	۶/۶۹	۱۶/۸۲	۶/۸۲	۶/۸۲	۱۶/۸۲
پاکشت	۷۶۸/۰۰	۱۰۸۰/۰۰	۹/۷۸	۶/۶۴	۱۰/۲۵	۶/۷۷	۶/۷۷	۱۰/۲۵
تهران و کرون	۷۲۰/۰۰	۷۶۸/۰۰	۴/۵۱	۶/۶۴	۵/۰۲	۶/۷۷	۶/۷۷	۵/۰۲
ارسنجان	۹۵۸/۴۰	۴۸۶/۶۰	۳/۴۵	۶/۶۸	۲/۷۹	۶/۸۱	۶/۸۱	۲/۷۹
خرم بید	۹۳۷/۶۰	۵۶۰/۰۰	۲/۳۲	۶/۶۸	۲/۸۲	۶/۸۱	۶/۸۱	۲/۸۲
جاجرم	۱۱۳۱/۲۰	۱۲۹۴/۴۰	۱۸/۸۴	۶/۷۶	۱۸/۸۲	۶/۸۹	۶/۸۹	۱۸/۸۲
ابوموسی	۱۱۸۵/۶۰	۵۴/۴۰	۱۱/۱۰	۷/۱۴	۱۰/۱۷	۷/۳۱	۷/۳۱	۱۰/۱۷
راور	۱۲۷۲/۰۰	۶۷۶/۸۰	۱۰/۱۹	۶/۷۸	۸/۵۶	۶/۹۳	۶/۹۳	۸/۵۶

است. وقتی در کریگینگ مدل‌های متفاوتی برای تغییرنگاریه کار گرفته می‌شود، می‌توان میزان دقت تخمین حاصل از کریگینگ‌های مختلف را براساس معیار $RMSD$ با یکدیگر مقایسه نمود. بدیهی است هر چقدر مقدار این معیار به یک نزدیکتر باشد کریگینگ از دقت بیشتری برخوردار می‌باشد. این معیار برای میزانهای بروز سل ریوی ۱۳۷۷ برابر ۱/۲۶ و برای سال ۱۳۷۸ برابر ۱/۲۴ شده است، که به مقدار مورد انتظار یک نزدیک می‌باشند. نظر به اینکه تغییر مدل تغییرنگار سبب تغییر دقت تخمین‌ها می‌شود، معیار $RMSD$ به طور غیرمستقیم دقت مدل تغییرنگار برازش شده را نیز نشان می‌دهد، ممکن است بتوان با یافتن مدل‌های تغییرنگاری که منجر به کاهش $RMSD$ شوند، تخمینهای دقیق‌تری را نیز به دست آورد.

۵ نقشه آماری

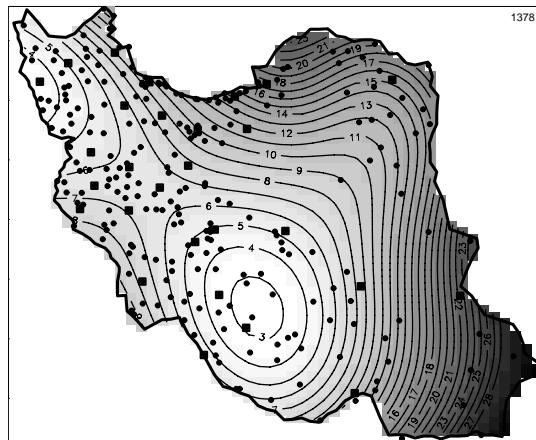
برای تهیه نقشه آماری بروز سل ریوی، کشور در یک شبکه منظم 50×67 با فواصل ۲۸/۸ کیلومتر محاط گردیده است. در تمام گرهای شبکه واقع در داخل کشور میزان بروز بیماری ریوی به روش کریگینگ عam، با روند درجه دوم تخمین زده شده و انحراف معیار آنها محاسبه شده‌اند. شکلهای ۲ و ۳ نقشه‌های تخمین میزان بروز بیماری و شکلهای ۴ و ۵ نقشه‌های انحراف معیار تخمین را در سالهای ۱۳۷۷ و ۱۳۷۸ نشان می‌دهند. در این نقشه‌ها مراکز شهرستانها با نقطه و مراکز استانها با مربع نمایان شده‌اند. افزایش تیرگی رنگ در شکلهای ۲ و ۳ از مناطق مرکزی به سمت مرزهای شرقی کشور نشان‌دهنده کمترین میزان بروز بیماری در مناطق مرکزی است، که به تدریج با نزدیک شدن به سمت مرزهای جنوب شرقی و قسمتی از شمال شرق به بیشترین مقدار می‌رسد.

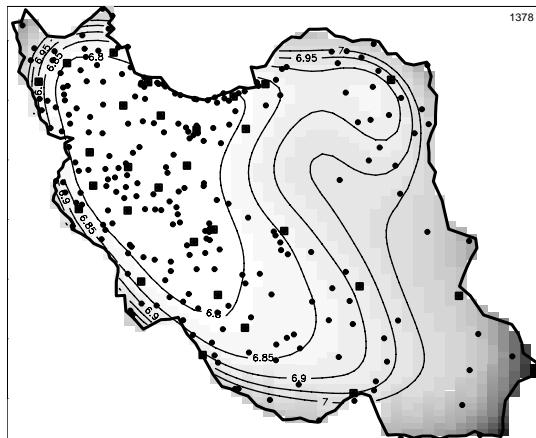


شکل ۲: نقشه میزان بروز سل ریوی، ۱۳۷۷

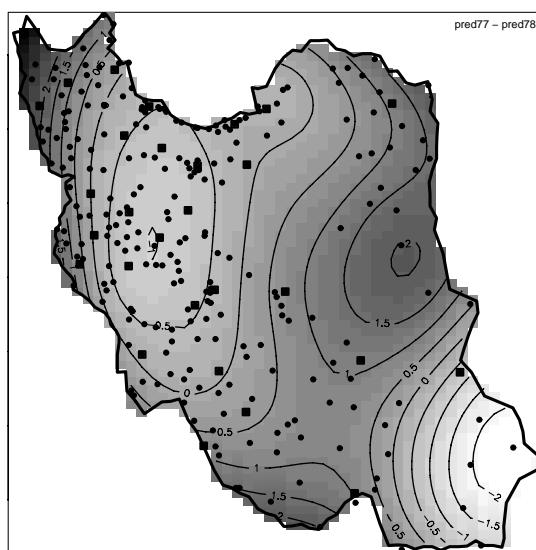
با توجه به شکلهای ۴ و ۵ ملاحظه می‌شود کمترین انحراف معیار مربوط به نیمه غربی کشور است که در آن تجمع شهرستانها بیشتر و در نتیجه داده‌های زیادتری در اختیار هستند. افزایش انحراف معیار به سمت شرق می‌تواند ناشی از پراکندگی شهرستانها و در نتیجه کم بودن تعداد داده‌ها و حتی ناشی از تغییرات زیاد بیماری در این مناطق باشد. بطور کلی انحراف معیار تخمین‌ها در سطح کشور نوسان زیادی ندارد و این نشان‌دهنده آن است که دقت تخمین‌ها تقریباً یکسان است.

با توجه به شکلهای ۲ و ۳ در اغلب مناطق کشور میزان بیماری در سال ۱۳۷۸ نسبت به ۱۳۷۷ کاهش یافته است. به منظور بررسی روند جغرافیایی کاهش یا افزایش بیماری





شکل ۵: نقشه انحراف معیار تخمین میزان بروز بیماری، ۱۳۷۸



شکل ۶: نقشه کاهش میزان بروز بیماری از ۱۳۷۷ به ۱۳۷۸

۶ بحث و نتیجه‌گیری

در این مقاله روش تخمین فضایی کریگینگ عام برای تجزیه و تحلیل داده‌ها در همه‌گیری‌شناسی جغرافیایی بیماریها و تهیه نقشه آماری آنها ارائه شد و داده‌های سل ریوی اسمیر مثبت سالهای ۱۳۷۷ و ۱۳۷۸ شهرستانهای کشور، به عنوان یک مطالعه موردی، مورد بررسی قرار گرفتند، که میزان دقت نتایج تحت تأثیر کمثی بیماری ناشی از عدم

مراجعه بعضی بیماران به شبکه بهداشت یا عدم تشخیص بیماری می‌باشد. صرف نظر از آن نقشه‌های حاصل از روش‌های ارائه شده، روند جغرافیایی بیماری را بطور کامل نشان می‌دهند. روند افزایشی میزان بروز بیماری در شرق ایران می‌تواند دلایل همه‌گیری‌شناسی خاصی داشته باشد، که در اینجا مورد بررسی قرار نگرفته‌اند.

در این مقاله کریگینگ براساس یک مدل تغییرنگار همسانگرد که به روش حداقل مربuat موزون به داده‌ها برازش داده شده صورت پذیرفته است. با توجه به امکان وجود ناهمسانگری تأثیر بسزایی در میزان دقت تخمين‌ها و نقشه بیماریها داشته باشد، که لازم است در مطالعات بعدی مورد توجه قرار گیرد.

کلیه محاسبات و تهییه نمودارها با برنامه‌نویسی در محیط نرم‌افزار SPLUS و به کمک مجموعه توابع GeoS (ربیرو و دیگل، ۲۰۰۰) صورت گرفته است.

از همکاری صمیمانه واحد سل اداره کل مبارزه با بیماریها که داده‌ها و راهنماییهای فنی در مورد نحوه ثبت و جمع آوری آنها را ارائه نموده‌اند تشکر و قدردانی می‌شود.

مراجع

- Barry, J., Crowder, M. and Diggle, P., (1997). Parametric Estimation of the Variogram, *Technical Report ST-97-06*. Dep. Maths and Stats, Lancaster University, Lancaster, UK.
- Bernardinelli, L., Clayton, D., Pascutto, C., Montomoli, C., Ghislandi, M. and Songini, M., (1995). Bayesian Analysis of Space-Time Variation in Disease Risk, *Statistics in Medicine*,, **14**: 2433-2443.
- Carrat, F. and Valleron A., J., (1992). Epidemiologic Mapping Using the Kriging Method: Application to an Influenza-Like Illness epidemic in France, *American Journal of Epidemiology*,, **135**: 1293-1300.
- Clayton, D. and Bernardinelli, L., (1992). Bayesian Methods for Mapping Disease Risk; In: *Geographical and Environmental Epidemiology: Methods for Small-Area Studies* (Eds., Elliott, P., Cuzick, J., English, D., and Stern R.), Oxford Press, Oxford.

- Clayton, D. G. and Kaldor, J., (1987). Empirical Bayes Estimates of Age-Standardized Relative Risks for Use in Disease Mapping, *Biometrics*,, **43**: 671-691.
- Cressie, N., (1993). *Spatial Statistics for Spatial Data*, New York, John Wiley.
- De Boor C., (1978) *A Practical Guide to Splines*. New York: Springer-Verlag New York, Inc.
- Glattre, E., (1989). *Atlas of Cancer Incidence in Norway 1970-1979*. Recent Results Cancer Res, **114**: 216-226.
- Lawson, A. B., (2001). *Statistical Methods in Spatial Epidemiology*, John Wiley.
- Lawson, A. B., Cressie, N., (2000). *Spatial Statistical Methods for Environmental Epidemiology*, Handbook of Statistics, **18**: 357-396.
- Lawson, A. B. and Williams, L. R., (2001). *An Introductory Guide to Disease Mapping*, John Wiley.
- Linde, V. D. A., Witzko, K. H. and Jaockel, K., H., (1995). Spatial-Temporal Analysis of Mortality Using Splines, *Biometrics*,, **51**: 1352-1360.
- Marshall, R. A., (1991). Review of Methods for the Statistical Analysis of Spatial Patterns of Disease, *Journal of Royal Statistical Society, Series A*,, **154**, 421-441.
- Matheron, G., (1963). Principles of Geostatistics, *Economic Geology*,, **58**: 1246-1266.
- Matheron, G., Le Krigeag, Universal, (1969). *Cahiers du Center de Morphologie Mathematique*,, No. 1 Fountainbleau, France.
- Matheron, G., (1973). The Intrinsic Random Functions and their Applications, *Adv. Appl. Prob.*, No. 5, 439-468.

- Mugglin, A. S., Cressie, N., and Gemmell, I., Hierarchical, (2000). Statistical Modeling of Influenza-Epidemic Dynamic in Space and Time, (Tentative title), in preparation.
- Pukkala, E., (1989). *Cancer Maps of Finland: An Example of Small-Area Based Mapping*. Recent Results Cancer Res, **114**: 208-215.
- Ribeiro, Jr. P. J. and Diggle, P. J., **geoR/geoS**: Functions for Geostatistical Analysis Using R or S-PLUS Technical Report: ST-99-09, Dep. Of Maths. And Stats. Lancaster University, Lancaster, UK.
<http://www.Maths.Lancs.ac.uk/~ribeiro/geoR.html>.
- Ripley, B., D., (1981). *Spatial Statistics*, New York: John Wiley & Son, 44-77.
- Snow, J. (1954). *On the Mode of Communication of Cholera*. Churchill Livingstone, London, 2nd edition.
- Stone, M. (1974). Cross-Validatory choice and Assement of Statistical Predictions. *Journal of the Royal Statistical Society, B.*, **36**: 111-133.
- Waston, G. S., (1972). Trend Surface Analysis and Spatial Correlation. *Geological Society of America*, (Special Paper), **146**: 39-46.
- Watson, D., F. and Philip, G., M., (1985). A refinement of Inverse Distance Weighted Interpolation, *Geoprocessing*, **2**: 315-327.
- Webster, R. and Oliver M., (2001). *Geostatistics for Environmental Scientists*, John Wiley.
- محمدزاده محسن (۱۳۷۸)، تابع کوواریانس تعمیم‌یافته برای کریگیند عالم اندیشه آماری، شماره دوم، ۲۳-۱۸.
- محمدزاده محسن، کاظم‌نژاد انوشیروان، فقیه‌زاده سقراط، واقعی‌یدا...، (۱۳۸۰)، تعداد لگ مناسب در مدل‌سازی تغییرنگار، ارائه شده به مجله علوم دانشگاه تربیت معلم.

شناسایی داده‌های دورافتادهٔ فضایی

محسن محمدزاده، علی محمدیان مصمم

P ۱۳۰۲۶

گروه آمار، دانشگاه تربیت مدرس

چکیده: یکی از مسائل مهم در تجزیه و تحلیل داده‌های فضایی شناسایی داده‌های دورافتادهٔ فضایی است. وقتی که بیش از یک داده دورافتاده در مجموعه داده‌ها باشد، اغلب روش‌هایی که تاکنون ارائه شده چهار مشکلاتی تحت عنوان درون آوری یا برون بری می‌شوند که منجر به شناسایی غلط یک داده معمولی عنوان داده دورافتاده یا بالعکس می‌شوند. در این مقاله یک روش تحلیل اکتشافی شناسایی داده‌های دورافتادهٔ فضایی ارائه و به کمک یک الگوریتم جستجوی پیشرو، مشاهدات بر اساس سازگاری‌شان با یک مدل مشخص مرتب و داده‌های دورافتادهٔ فضایی شناسایی می‌شوند.

واژه‌های کلیدی: داده‌های فضایی، کریگینگ و داده‌های دورافتادهٔ فضایی.

۱ مقدمه

در حالتی که مشاهدات دارای نوعی واپستگی ناشی از موقعیت آنها در فضای مورد مطالعه باشند، داده‌های فضایی^۱ نامیده می‌شوند. هدف این مقاله ارائه روشی برای شناسایی داده‌های دورافتاده^۲ در مجموعه داده‌های فضایی است. معمولاً داده‌هایی که نزدیک هم قرار دارند، شبیه یکدیگرند و هر چقدر موقعیت جغرافیایی آنها از یکدیگر دور می‌شود از همبستگی و شباهت آنها کاسته می‌شود. در آمار فضایی، مشاهداتی که نسبت به مقادیر همسایگی خود، فرین^۳ باشند، داده دورافتادهٔ فضایی نامیده می‌شود.

معمولًاً برای تجزیه و تحلیل آماری داده‌ها ابتدا تحلیل اکتشافی^۴ برای شناخت خصوصیات آنها انجام می‌گیرد. هاسلت (۱۹۹۱)، کرسی (۱۹۹۳) و پاناتیر (۱۹۹۶) برای کشف داده‌های دورافتاده از ابزارهای گرافیکی ساده نظری نمودارهای جعبه‌ای^۵ و ابرهای واریوگرام استفاده نموده‌اند. همچنین کریستنسن، جانسون و پیرسون (۱۹۹۲) و هاسلت و هایز (۱۹۹۸) نیز

^۱ Spatial Data

^۲ Outliers

^۳ Extreme

^۴ Exploratory Analysis

^۵ Pocket Plots

روش تشخیص نقاط دورافتاده را پس از برازش مدل به داده‌ها ارائه نمودند. در این مقاله داده‌های دورافتاده فضایی براساس الگوریتم جستجوی پیشو^۶، که مشاهدات را براساس سازگاری‌شان با یک مدل مشخص مرتب می‌کند، شناسایی می‌شوند. برتری این روش نسبت به روش‌های سنتی این است که براساس مدل پیشنهاد شده یک مرتب سازی برای داده‌ها فراهم می‌سازد و داده‌های دورافتاده چندگانه را بدون تاثیر مسائل درون‌آوری^۷ و برونوبری^۸، کشف می‌نماید. چون شناسایی داده‌های دورافتاده بر اساس باقیمانده پیش‌بینی‌ها انجام می‌شود، روش پیش‌بینی کریگینگ در بخش ۲ شرح داده خواهد شد. در بخش ۳ داده‌های دورافتاده چندگانه معرفی و روش‌های شناسایی و مشکلات درون‌آوری و برونوبری آنها در بخش ۴ مورد بررسی قرار می‌گیرد و نهایتاً بحث و نتیجه گیری در بخش ۵ ارائه خواهد شد.

۲ کریگینگ

معمولًاً در آمار فضایی، میدان تصادفی^۱

$$\{Z(s) : s \in D \subset \mathbb{R}^d\} \quad (1)$$

برای مدل سازی داده‌های فضایی در نظر گرفته می‌شود، که در آن D مجموعه اندیس‌گذار، فضای اقلیدسی \mathbb{R}^d بعدی و برای هر $s \in D$ یک متغیر تصادفی است. برای تخمین مقادیر میدان تصادفی (۱) روش کریگینگ به عنوان بهترین پیش‌بینی کننده خطی ناریب^۲ مورد استفاده قرار می‌گیرد. کریگینگ معمولی^۳ روشی برای پیش‌بینی کمیت‌های غیرقابل مشاهده از میدان تصادفی (۱) بفرم

$$Z(s) = \mu + \delta(s) \quad s \in D \quad \mu \in \mathbb{R} \quad (2)$$

است، که در آن μ ثابت نامعلوم و $\{\delta(s) : s \in D\}$ یک میدان تصادفی ایستای ذاتی با واریوگرام معلوم

$$2\gamma(s-t) = Var(Z(s) - Z(s+h)); \quad s \in D$$

^۶ Forward Search

^۷ Masking

^۸ Swamping

^۱ Random Field

^۲ Best Linear Unbiased Estimator

^۳ Ordinary Kriging

و میانگین صفر باشد. برای مشاهدات داده شده در مکانهای فضایی $\{s_1, \dots, s_n\}$ بهترین پیش‌بینی کننده خطی نا اریب در موقعیت $s \in D$ ، بصورت

$$\hat{Z}(s_0) = \underline{\lambda}' Z$$

محاسبه می‌شود، که در آن

$$\underline{\lambda}' = (\underline{\gamma} + \underline{\lambda} \frac{\underline{1} - \underline{\lambda}' \Gamma^{-1} \underline{\gamma}}{\underline{\lambda}' \Gamma^{-1} \underline{1}})' \Gamma^{-1} \quad (3)$$

بردار ضرایب، $Z' = [Z(s_1), \dots, Z(s_n)]$ بردار مشاهدات، $\underline{\lambda}$ بردار ستونی واحد، $\underline{\gamma}'$ یک ماتریس $n \times n$ بصورت زیر می‌باشد.

$$\Gamma = \begin{pmatrix} \gamma(s_1 - s_1) & \gamma(s_1 - s_2) & \dots & \gamma(s_1 - s_n) \\ \gamma(s_2 - s_1) & \gamma(s_2 - s_2) & \dots & \gamma(s_2 - s_n) \\ \vdots & \vdots & \ddots & \vdots \\ \gamma(s_n - s_1) & \gamma(s_n - s_2) & \dots & \gamma(s_n - s_n) \end{pmatrix}$$

مقدار $\hat{Z}(s_0)$ از مینیمم کردن میانگین مجدول خطای پیش‌بینی به روش لاگرانژ در میان تمام توابع خطی $\lambda' Z$ بدست می‌آید، بطوریکه برای بردار n بعدی λ ، $E[\lambda' Z] = E[Z(s_0)]$ باشد. در این صورت واریانس کریگینگ در موقعیت s_0 نیز بصورت زیر خواهد بود.

$$\sigma^2(s_0) \equiv E[Z(s_0) - \hat{Z}(s_0)]^2 = \underline{\gamma}' \Gamma^{-1} \underline{\gamma} - \frac{(\underline{\lambda}' \Gamma^{-1} \underline{\gamma} - \underline{1})^2}{\underline{\lambda}' \Gamma^{-1} \underline{1}}$$

۳ داده دورافتاده چندگانه

در مقالات روش‌های متفاوتی برای شناسایی داده‌های دورافتاده فضایی ارائه گردیده‌اند. کرسی (۱۹۹۳) با قرار دادن $\{s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_n\}$ بعنوان $\hat{Z}_i(S_{-i})$ پیش‌بینی کننده کریگینگ برای $Z(s_i)$ براساس مشاهدات در موقعیت‌های S_{-i} و واریانس کریگینگ $(S_{-i})^2 \sigma_i^2$ ، باقیمانده‌های استاندارد پیش‌بینی را بصورت

$$\tilde{e}_i(S_{-i}) = \frac{Z(s_i) - \hat{Z}_i(S_{-i})}{\sigma_i(S_{-i})} \quad i = 1, \dots, n. \quad (4)$$

تعریف نمود و استفاده از نمودارهای تشریحی باقیمانده‌های (۴) را برای کشف داده‌های دورافتاده فضایی پیشنهاد کرد. کریستنسن، جانسون و پیرسون (۱۹۹۲) نیز با تجزیه میدان تصادفی بفرم $Z(s) = \mu + \delta(s) + \varepsilon(s)$ و $\delta(s)$ مستقل هستند،

روشهای تشخیص را بر اساس باقیماندهای (۴) مورد بررسی قرار دادند. روشهای فوق عمدتاً برای مواردی قابل استفاده هستند که تنها یک داده دور افتاده در داده‌ها وجود داشته باشد. در غیر این صورت ممکن است یک داده دور افتاده موجب پنهان شدن داده‌های دور افتاده دیگر شود، که به آن مشکل درون آوری گفته می‌شود. یا امکان دارد مشاهده معمولی به عنوان داده دور افتاده شناسایی شود، که مشکل برونو برع نامیده می‌شود. برای نمایش این مشکلات یک میدان تصادفی به فرم $Z(s) = \mu + \delta(s) + \varepsilon(s)$ $s \in D$, با فرض گوسی بودن $(\cdot)^{\delta}$ و $(\cdot)^{\varepsilon}$ و با استفاده از شبهواریوگرام کروی^۱

$$\gamma(h) = \begin{cases} 0 & ||h|| = 0, \\ \theta_0 + \theta_1 [1.5||h||/\theta_2 - 0.5(||h||/\theta_2)^3] & 0 < ||h|| \leq \theta_2, \\ \theta_0 + \theta_1 & ||h|| \geq \theta_2, \end{cases} \quad (5)$$

برای پارامترهای $\theta_0 = 2$, $\theta_1 = 4$, $\theta_2 = 8$, میانگین $\mu = 10$ و واریانس $\text{Var}(\varepsilon(s_i)) = 1$ داده $n = 81$ در یک شبکه منظم 9×9 شبیه‌سازی و در جدول ۱ نشان داده شده است. در موقعیتها $s_1 = (1, 1)$, $s_2 = (1, 2)$, $s_3 = (1, 3)$ بترتیب مقادیر ۶, ۴

جدول ۱: موقعیتها فضایی در شبکه منظم 9×9

	۱	۲	۳	۴	۵	۶	۷	۸	۹
۱	۱۱/۴	۱۲/۴	۱۲/۸	۱۳/۲	۹/۸	۱۱/۴	۱۲/۰	۱۰/۸	۱۰/۵
۲	۱۱/۸	۱۲/۷	۱۰/۸	۱۳/۸	۱۴/۴	۱۱/۳	۸/۷	۷/۸	۷/۲
۳	۱۱/۲	۱۲/۳	۱۵/۶	۱۲/۳	۱۲/۳	۱۰/۲	۸/۹	۷/۷	۱۰/۱
۴	۱۲/۸	۹/۹	۱۲/۸	۹/۷۰	۱۰/۴	۸/۴	۵/۵	۷/۲	۱۰/۰
۵	۱۱/۶	۱۱/۹	۱۳/۸	۹/۹	۸/۹	۸/۵	۷/۵	۱۰/۶	۹/۱
۶	۱۱/۰	۱۱/۴	۱۴/۵	۱۲/۶	۱۱/۶	۱۱/۶	۸/۳	۸/۷	۱۰/۴
۷	۱۱/۴	۱۳/۱	۱۱/۳	۱۲/۴	۱۱/۳	۷/۴	۱۱/۰	۱۱/۴	۹/۵
۸	۱۰/۷	۱۲/۳	۱۲/۲	۱۴/۰	۱۱/۶	۱۱/۴	۱۰/۸	۱۱/۵	۱۰/۲
۹	۱۲/۲	۱۲/۶	۱۳/۸	۱۵/۰	۱۲/۸	۱۱/۶	۱۲/۲	۱۰/۹	۱۰/۸

و ۵ اضافه شده‌اند، بطوریکه داده‌های حاصل در این نقاط نسبت به عمدۀ داده‌ها، دور افتاده محسوب می‌شوند.

کوکل صاف و موده شیر و $\tilde{e}_i(S_{-i})$ اهدنامیقابی ابعاجی اهراد و مذ

$$C_i(S_{-i}) = -(\hat{\delta}(S_{-i}) - \hat{\delta})' \hat{T}^{-1} (\hat{\delta}(S_{-i}) - \hat{\delta}) \quad i = 1, \dots, n \quad (6)$$

^۱ Spherical Semivariogram

شکل ۱: نمودار جعبه‌ای الف: باقیمانده‌های پیش‌بینی شده ب: ریشه دوم فواصل کوک

که در آن، \hat{S} بردار مقادیر پیش‌بینی بر اساس S و $(S_{-i})\hat{\theta}$ بردار مقادیر پیش‌بینی بر اساس S_{-i} و \hat{T} برآورد ماتریس T است، در شکل ۱ نمایش داده شده است. همانطور که ملاحظه می‌شود، هیچ یک از مقادیر ناخالص موجود در موقعیتهای s_1, s_2 و s_3 بعنوان داده دور افتاده تشخیص داده نمی‌شود و مشکل درون آوری بروز می‌کند. تنها داده دور افتاده‌ای که در هر دو نمودار الف و ب شکل ۱ نشان داده می‌شود، مشاهده موقعیت s_{12} است که دارای مقدار نسبتاً کوچک $(S_{12})\hat{\theta}$ و مقدار بزرگ $C_{12}(S_{12})$ می‌باشد. چنین نتایج نامطلوب را می‌توان در تکنیک‌های تشریحی که توسط کرسی (۱۹۹۳) برای پیدا کردن بی‌نظمی‌های فضایی طراحی شده‌اند نیز ملاحظه نمود. در این تکنیک مقدار قدر مطلق تفاضلهای استاندارد شده در هر ردیف و ستون شبکه بصورت

$$u \equiv \sqrt{m}(\bar{Z} - \tilde{Z}) / (0.56\hat{\Psi}), \quad (7)$$

محاسبه می‌شوند، که در آن m تعداد موقعیتهای دارای مشاهده در یک سطر یا ستون مشخص می‌باشد و \bar{Z} و \tilde{Z} بترتیب میانگین، میانه و دامنه میان چارکی همان سطر یا ستون هستند. اگر مقدار $|u|$ نزدیک یا بیشتر از ۳ باشد، آنگاه آن ردیف یا ستون می‌تواند دارای داده دور افتاده باشد. در اینجا سطرهای ۷ و ۸ بخاراط داشتن مقادیر $4/3$ و $3/8$ و ستون ۱ بخاراط وجود مقدار $1/5$ مشکوک داشتن داده دور افتاده تشخیص داده می‌شوند، در حالیکه این سطرهای نقاط دور افتاده نیستند و تنها ستون ۱ دارای نقطه دور افتاده است. با این روش ستونهای ۲ و ۳ و سطر ۱ که شامل نقاط دور افتاده هستند، انتخاب نشده‌اند. لذا تکنیکی که این گونه داده‌های دور افتاده را به خوبی شناسایی کند ضروری می‌باشد. بعلاوه در این روش اطلاعات زیادی از طریق مرتب کردن مشاهدات از بیشترین سازگاری‌شان با مدل فضایی مشخص تا آنهایی که سازگاری کمتری با آن دارند ارائه

جدول ۲: قدرمطلق تفاضلهای استاندارد شده u (میانگین منهای میانه)

شماره سطرها	۱	۲	۳	۴	۵	۶	۷	۸	۹
$ u $	۱/۱	۰/۶	۱	۰/۸	۰/۵	۱/۴	۴/۳	۳/۴	۰/۵
شماره ستونها	۱	۲	۳	۴	۵	۶	۷	۸	۹
$ u $	۵/۱	۲/۳	۰/۴	۰/۲	۰/۴	۲	۱	۱/۸	۲/۱

می‌شود.

۴ روش جستجوی پیشرو

برای شناسایی نقاط دورافتاده چندگانه در مدل‌های رگرسیونی با خطاهای مستقل و آنالیز چند متغیره روش‌هایی در مقالات هادی (۱۹۹۲)، هادی و سیمونوف (۱۹۹۳)، اتکینسن (۱۹۹۴) و ریانی و اتکینسن (۱۹۹۹) ارائه گردیده است. برای شناسایی داده‌های دورافتاده چندگانه در مدل‌های پیش‌بینی فضایی از روش جستجوی پیشرو استفاده خواهد شد، که الگوریتم‌های آن با تعديل روش‌های فوق برای داده‌های فضایی طراحی می‌شوند. در این روش سریولی و ریانی (۱۹۹۹) یک زیرمجموعه مقدماتی از داده‌ها را که قادر نقاط دور افتاده هستند اختیار نموده و بر اساس مراحل چهارگانه الگوریتم زیر آنرا به مجموعه‌ای بزرگتر توسعه داده‌اند، بگونه‌ای که نقاط دور افتاده خارج از این مجموعه قرار خواهند گرفت.

مرحله ۰ (برآورد ماتریس Γ): در اغلب مسائل کاربردی تابع واریوگرام، که تعیین کننده ساختار وابستگی داده‌های فضایی است، نامعلوم می‌باشد. لذا بایستی این تابع بر اساس مشاهدات برآورد شود. کرسی و هاوکینز (۱۹۸۰) یک برآوردگر مقاوم برای واریوگرام بصورت

$$2\tilde{\gamma}(h) = \frac{(|N(h)|^{-1} \sum_{N(h)} \sqrt{|Z(s_i) - Z(s_j)|})^4}{.452 + .494/|N(h)|} \quad h \in \Re^2$$

معرفی نمودند، که در آن $\{(s_i, s_j); s_i - s_j = h \quad i, j = 1, \dots, n\}$ و $N(h) = \{(s_i, s_j); s_i - s_j = h\}$ تعداد زوجها با تأخیر h و $\sum_{N(h)}$ مجموع روی چنین زوجهایی است. هاوکینز و کرسی (۱۹۸۴) همچنین $\frac{[\tilde{Z}(h)]^4}{.457} = 2\tilde{\gamma}^*(h)$ را معرفی کردند که در آن $\tilde{Z}(h)$ میانه ریشه دوم تفاضلهای $\sqrt{|Z(s_i) - Z(s_j)|}$ روی زوجهای با تأخیر h و $.457$ ضریب تصحیح

اریبی است. کرسی (۱۹۹۳) نشان داد (۲۷) یکتابع معین منفی^۱ است. برای اینکه برآوردهای (۲۷) در شرط معین منفی صدق کنند، مدل‌های پارامتری برای برازش (۲۷) استفاده می‌شوند و ماتریس Γ بصورت $[\hat{\gamma}(s_i - s_j)]$ برآورد می‌گردد.

مرحله ۱ (انتخاب زیرمجموعه مقدماتی): وقتی مدل مورد استفاده برای پیش‌بینی شامل p پارامتر باشد، این الگوریتم با انتخاب یک زیرمجموعه شامل p موقعیت فضایی که عاری از نقاط دورافتاده باشند، شروع می‌شود. برای شناسایی این موقعیت‌ها از نمودارهای جعبه‌ای دو متغیره استفاده می‌شود. برای $n \leq p$ می‌توان p تایی‌های مجزای

$$S_{i_1, \dots, i_p}^{(p)} \equiv \{s_{i_1}, \dots, s_{i_p}\}; 1 \leq i_1, \dots, i_p \leq n \text{ و } i_j \neq i_{j'}$$

را که تعداد آنها $\binom{n}{p}$ می‌باشد، انتخاب کرد. فرض کنید $[i_1, \dots, i_p] = \hat{Z}_i(S_\ell^{(p)})$ یک پیش‌بینی کننده کریگینگ در محل s_i براساس مشاهدات در $S_\ell^{(p)}$ باشد، که در اینصورت باقیماندهای استاندارد $\tilde{e}_i(S_\ell^{(p)})$ خواهد بود. اکنون زیرمجموعه مقدماتی p تایی $S_*^{(p)}$ بگونه‌ای انتخاب می‌شود که در رابطه

$$\tilde{e}_{[med]}^2(S_*^{(p)}) = \min_\ell [\tilde{e}_{[med]}^2(S_\ell^{(p)})] \quad (8)$$

صدق کند، که در آن $\tilde{e}_{[\ell]}^2(S_\ell^{(p)})$ نشان دهنده مجذور ℓ امین باقیمانده مرتب شده در میان باقیماندهای n, \dots, n است و $med = p + \frac{n-p}{2}$ می‌باشد، بطوری که $\frac{n-p}{2}$ نشان دهنده جُز صحیح است.

معیار (۸) که به منظور برازش یک پیش‌بینی کننده نیرومند به داده‌ها انتخاب شده است، شبیه روش حداقل مجذور میانه است که برای مدل‌های رگرسیونی با خطاهای مستقل استفاده می‌شود، جز اینکه در اینجا خطاهای براساس $(S_\ell^{(p)}, \sigma_i)$ استاندارد می‌شوند. اگر n یا p بزرگ باشند، این فرایند ممکن است محاسبات زیادی داشته باشد. اتکینسن (۱۹۹۴) نشان داد که جستجوی پیشرو را می‌توان با انتخاب تصادفی زیرمجموعه مقدماتی انجام داد.

مرحله ۲ (پیشروی در جستجوی پیشرو): اگر بعد زیرمجموعه $S_*^{(m)}$ باشد، جستجوی پیشرو برای بعد $m+1$ با انتخاب ۱ مکان فضایی بطوریکه مجذور باقیماندها حداقل باشد با مرتب کردن مجذور باقیماندهای $\tilde{e}_i^2(S_*^{(m)})$ ، $i = 1, \dots, n$ ، فقط یک موقعیت جدید به زیرمجموعه انتخاب می‌شوند. در اکثر انتقالات از m به $m+1$ ، ممکن است دو یا چند موقعیت به $S_*^{(m)}$ ملحق می‌شود. اما برای مدل (۲) ممکن است دو یا چند موقعیت به $S_*^{(m)}$ ملحق یا چند موقعیت از آن جدا شوند. هر چند طبق تجربه چنین رخدادهایی وقتهای که الگوریتم از $S_*^{(p)}$

^۱ Negative Definite

شروع می‌شود، غیرطبیعی هستند و تنها وقتی اتفاق خواهد افتاد که جستجو شامل یک موقعیت متعلق به مجموعه داده دورافتاده فضایی باشد. می‌توان جستجو را بوسیله مرتب کردن مجدد باقیمانده‌ها که بصورت زیر تعریف می‌شوند،

$$r_i^*(S_*^{(m)}) = \begin{cases} e_i^*(S_*^{(m)}) & \text{اگر } s_i \in S_*^{(m)} \\ \tilde{e}_i^*(S_*^{(m)}) & \text{اگر } s_i \notin S_*^{(m)} \end{cases}$$

ادامه داد، که در آن $e_i(S_i) - \hat{Z}_i(S_*^{(m)}) \equiv Z(s_i)$ است. در این صورت با استفاده از $r_i^*(S_*^{(m)})$ بجای $\tilde{e}_i^*(S_*^{(m)})$ احتمال اینکه دو یا چند موقعیت به $S_*^{(m)}$ در اولین مراحل جستجوی پیشرو ملحق شوند تا حدی افزایش می‌یابد. بنابراین توصیه می‌شود این روش وقتی که زیرمجموعه مقدماتی بطور تصادفی انتخاب می‌شود، به منظور افزایش شанс دفع داده دورافتاده از زیرمجموعه مقدماتی اتخاذ شود.

مرحله ۳ (مرتب سازی داده‌های فضایی) : تا زمانیکه تمام مکانها وارد زیرمجموعه مقدماتی شوند، مرحله ۲ تکرار می‌شود. اگر در هر انتقال فقط یک مکان به $S_*^{(m)}$ وارد شود، الگوریتم یک مرتب سازی از داده‌ها را فراهم می‌کند. نقاط دورافتاده فضایی و همچنین مشاهدات دیگری که بی‌نظمی‌هایی دارند را می‌توان توسط تشریح گرافیکی از تغییرات آماره‌های موجود در جستجوی پیشرو را کشف نمود. برای این منظور نمودار باقیمانده‌های $\tilde{e}_i^*(S_*^{(m)})$ در مقابل m رسم می‌شوند.

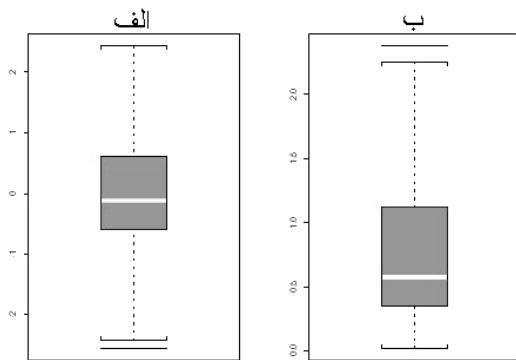
$$\begin{aligned} \tilde{e}_m^{(m+1)} &= \sqrt{\tilde{e}_{[m+1]}^*(S_*^{(m)}),} & \sigma_m^{(m+1)} &= \sigma_{[m+1]}^*(S_*^{(m)}) & m &= p+1, \dots, n-1 \\ \tilde{e}_m^{(max)} &= \sqrt{\tilde{e}_{[n]}^*(S_*^{(m)}),} & \sigma_m^{(max)} &= \sigma_{[n]}^*(S_*^{(m)}) & m &= p+1, \dots, n-1 \end{aligned}$$

اگریک یا چند مشاهده بی‌قاعده یا داده دورافتاده وجود داشته باشد، نمودارهای $\tilde{e}_m^{(m+1)}$ و $\sigma_m^{(m+1)}$ باقیستی یک قله در مرحله قبلی نسبت به اولین داده دورافتاده نشان دهند.

در رابطه با یک طبقه از داده دورافتاده، وقتی اولین داده دورافتاده به $S_*^{(m)}$ ملحق می‌شود منحنی‌های $\tilde{e}_m^{(max)}$ و $\sigma_m^{(max)}$ باقیستی یک شیب نزولی برای سهم اثر درون‌آوری داشته باشند. شکل ۲ نمودار باقیمانده‌های آماره‌های فوق را نشان می‌دهد. همانطور که ملاحظه می‌شود با نمایش توأم این نمودارها می‌توان اثر هر مشاهده را بر حسب ویژگی‌های داده‌های اصلی تعبیر کرد.

۵ نتیجه‌گیری

برای تشخیص داده‌های دورافتاده وقتی داده‌ها تنها شامل یک داده دورافتاده باشند روش‌های کلاسیک قابل استفاده هستند. اما چنانچه بیش از یک داده دورافتاده در مجموعه داده‌ها



شکل ۲: منحنی‌های الف: $\sigma_m^{(max)}$ ، ب: $\tilde{\sigma}_m^{(m+1)}$ ، ج: $\tilde{e}_m^{(max)}$

باشد، توان این روشها شدیداً کاهش می‌یابد. بعلاوه همانطور که نشان داده شد، این روشها در تشخیص داده‌های دور افتاده همواره موفق نیستند، چون با مشکلات درون آوری و برون بری مواجه می‌شوند. روش ارائه شده بر اساس الگوریتم جستجوی پیشرو داده‌های دور افتاده چندگانه را بدون اینکه تحت تأثیر اینگونه مسائل قرار گیرند شناسایی می‌کند. بعلاوه می‌تواند برای کشف غیرایستایی و دیگر بی‌نظمی‌های موضعی نیز بکار رود. چون روش ارائه شده یک مرتب سازی برای داده‌ها فراهم می‌سازد، حتی اگر داده‌ها شامل نقاط دور افتاده نباشند، این روش می‌تواند اطلاعات اکتشافی مفیدی را فراهم نماید.

مراجع

Atkinson, A. C. (1986), *Masking Unmasked* ", Biometrika,**73**, 533-541.

Atkinson, A. C. (1994), *Fast Very Robust Method for the Detection of Multivariate Outliers*, Journal of the American Statistical Association,**89**, 1329-1339.

Christensen, R., Johnson, W., and Pearson, L. M. (1992), *Prediction Diagnostics for Spatial Linear Models*, Biometrika,**79**, 583-591.

Cerioli, A. (1999), *The Ordering of Spatial Data and Detection of Multiple Outliers*, Journal of Computational and Graphical Statistics,**8**, 239-258.

- Cressie, N. (1993), *Statistics for Spatial Data*, New York: Wiley.
- Cressie, N. and Hawkins, D. M. (1980), *Roust Estimation of the Variogram:I*, Mathematical Geology,**12**, 115-125.
- Hadi, A. S. (1992), *Identifying Multiple Outliers in Multivariate Data*, Journal of the Royal Statistical Society, Ser. B,**54**, 761-771.
- Hadi, A. S., and Simonoff, J. S. (1993), *Procedures for the Identification of Multiple Outliers in Linear Models*, Journal of the American Statistical Association,**88**, 1264-1272.
- Haslett, J., Bradley, R., Craig, P., Unwin, A. and Wills, G. (1991), *Dynamic Graphics for Exploring Spatial Data With Application to Locating Global and Local Anomalies*, The American Statistician,**45**, 234-242.
- Haslett, J. and Hayes, K. (1998), *Residuals for the Linear Model whit Covariance Structure*, Journal of the Royal Statistical Society, Ser. B,**60**, 201-215.
- Hawkins, D. M. (1993), *The Accuracy of Elemental Sets Approximations for Regression*, Journal of the American Statistical Association,**88**, 580-589.
- Riani and Atkinson, A. C. (1999), *A Unified Approach to Multivariate Transformations and Multiple Outliers*, Submitted for Publication.

سرعت همگرایی الگوریتم EM و تسریع آن

سید محسن میرحسینی^۱، دکتر عین‌الله پاشا^۲

P11223

^۱ دانشکده ریاضی، دانشگاه یزد

^۲ دانشگاه تربیت معلم تهران

چکیده: یکی از روش‌های برآورد پارامتر بعد از مشاهده، روش ماکزیمم درستنمایی است؛ اگر بردار مشاهدات غیرکامل باشند که فقدان قسمتی از داده‌ها ممکن است از مکانیزم سانسور کردن یا از مقادیر گمشده و... باشد، می‌توان از الگوریتم EM استفاده نمود. در این مقاله، ابتدا به بررسی اجمالی الگوریتم EM و بیان برخی ویژگی‌ها و انتقادهای وارد بر آن می‌پردازیم. سپس مطالبی درباره همگرایی الگوریتم EM بیان می‌کنیم. با توجه به اینکه همگرایی الگوریتم EM به آرامی صورت می‌گیرد، با روش ایتنکن، همگرایی آن را تسریع می‌بخشیم و در ادامه به ارائه چند مثال در این زمینه می‌پردازیم.

واژه‌های کلیدی:تابع درستنمایی، الگوریتم EM ، همگرایی الگوریتم EM ^۱، تسریع همگرایی^۲ ایتنکن

۱ بخش اول: فرمول‌بندی الگوریتم EM

فرض کنید Y یک بردار تصادفی P بعدی با تابع چگالی احتمال $(p.d.f) g(y; \Theta)$ باشد که $\Theta = (\theta_1, \theta_2, \dots, \theta_d)^T$ بردار پارامترهای نامعلوم در تابع چگالی احتمال Y است. T -نشان دهنده ترانسپوزه بردار یا ماتریس است. فضای پارامتر را با Ω نمایش میدهیم. بردار Θ را می‌توان به روش ماکزیمم درستنمایی برآورد کرد. تابع درستنمایی برای Θ که از مقادیر مشاهده شده y تشکیل شده است برابر است با:

$$L(\Theta) = g(y; \Theta)$$

در بعضی مسائل، بردار شامل داده‌های غیرکامل است که فقدان قسمتی از داده‌ها ممکن است از مقادیر گمشده و یا از مکانیزم سانسور کردن و... ناشی شده باشد، بدست

^۱ convergence of EM algorithm

^۲ Aitken's acceleration method

آوردن برآورد ماکریم درستنایی (MLE)، کاری بس دشوار است. الگوریتم EM ، یکی از ابزاری است که در این زمینه کارآیی وسیعی دارد. $(\Theta; g_c(x))$ تابع چگالی احتمال بردار تصادفی X که نشان دهنده بردار داده‌های کامل x است. بنابراین لگاریتم درستنایی داده‌های کامل را با

$$\log L_c(\Theta) = \log g_c(x; \Theta)$$

نمایش می‌دهیم.

الگوریتم EM یک روش باتکرار برای برآورد پارامتر در داده‌های غیرکامل با استفاده از روش ماکریم درستنایی است. با توجه به اینکه $\log L(\Theta)$ شامل داده‌های غیرکاملی است، از تابع درستنایی $\log L_c(\Theta)$ استفاده می‌کیم. چون بردار داده‌های کامل x شامل تمام مشاهدات نمی‌باشد، بنابراین به جای $\log L_c(\Theta)$ ، امید ریاضی شرطی آن به شرط y را جایگزین می‌کنیم. به عبارت دیگر فرض کنید $(\Theta^{(k)}, \text{مقداری از } \Theta \text{ بعداز } k \text{ امین تکرار باشد. در گام } (1 + k)\text{ام گامهای } E \text{ و } M \text{ به صورت زیرند:}$

۱) گام E (گام امید ریاضی) : محاسبه $Q(\Theta; \Theta^{(k)})$ که

$$Q(\Theta; \Theta^{(k)}) = E_{\Theta^{(k)}} \{\log L_c(\Theta) | y\}$$

۲) گام M (گام ماکریم سازی) : انتخاب $\Theta^{(k+1)}$ برای هر مقداری از $\Omega \in \Theta$ که $Q(\Theta; \Theta^{(k)})$ را نسبت به Θ ماکریم کند. به عبارت دیگر:

$$Q(\Theta^{(k+1)}; \Theta^{(k)}) \geq Q(\Theta; \Theta^{(k)}) \quad \forall \Theta \in \Omega.$$

گامهای E و M مکرراً تکرار می‌شوند تا نهایه $L(\Theta^{(k)})$ به یک مقداری همگرا شود.

۲ بخش دوم: همگرایی یک دنباله EM به یک مقدار مانا

در سال ۱۹۷۷ Laird, Rubin, Dempster و نشان دادند که تابع درستنایی داده‌های غیرکامل $L(\Theta)$ بعد از یک تکرار، کاهشی نیست [۶، ۱۱]، بنابراین دنباله از بالا کراندار مقادیر درستنایی $\{L(\Theta^{(k)})\}$ ، به طور یکتاخت به تعدادی مقدار L^* همگرا می‌شود که در اغلب موارد L^* یک مقدار مانا است که θ^* در رابطه $0 = \partial L(\Theta)/\partial \theta$ یا به طور

معادل $\partial \log L(\Theta) / \partial \Theta = 0$ صدق می‌کند. در برخی دیگر از مواقع ممکن است L^* یک ماکریم نسبی باشد. در موقع نادری ممکن است Θ^* یک ماکریم کننده مطلق یا نسبی نباشد بلکه به عنوان مثال، یک نقطه زینی باشد که به نقطه شروع (Θ^0) بستگی دارد.

۳ بخش سوم: نگاشت EM

الگوریتم EM رامی‌توان با تعریف نگاشتی از فضای پارامتر $\Theta(\Omega)$ به خودش تعریف نمود:

$$\Theta^{(k+1)} = M(\Theta^{(k)}) \quad (k = 0, 1, 2, \dots)$$

اگر $\Theta^{(k)}$ به برخی مقادیر Θ^* همگرا شود و $M(\Theta)$ پیوسته باشد، آنگاه Θ^* باید در $\Theta^{(k+1)} = M(\Theta^*)$ صدق کند. بنابراین Θ^* نقطه ثابت شده نگاشت M است. با بسط تیلور $M(\Theta^{(k)})$ در همسایگی Θ^* داریم

$$\Theta^{(k+1)} - \Theta^* \approx J(\Theta^*)(\Theta^{(k)} - \Theta^*) \quad (1)$$

که $J(\Theta)$ ماتریس ژاکوبی $d \times d$ برای

$$M(\Theta) = (M_1(\Theta), M_2(\Theta), \dots, M_d(\Theta))^T$$

است که (i, j) امین درایه آن برابر است با $J_{ij}(\Theta) = \frac{\partial M_i(\Theta)}{\partial \Theta_j}$. بنابراین در یک همسایگی Θ^* ، الگوریتم EM اساساً یک تکرار خطی با ماتریس نرخ همگرایی $J(\Theta^*)$ است، بنابراین $J(\Theta^*)$ نوعاً ناصرف است. از این رو فرض می‌کنیم که $\Theta^{(k)}$ یک دنباله است به طوری که برای $\Theta^{(k+1)} = \Theta^*$ داریم

$$\frac{\partial Q(\Theta, \Theta^{(k)})}{\partial \Theta} = 0$$

در سال ۱۹۷۷ و همکارانش نشان دادند که اگر $\Theta^{(k)}$ همگرا به نقطه Θ^* باشد آنگاه

$$J(\Theta^*) = I - \mathcal{I}_c^{-1}(\Theta^*; y)I(\Theta^*; y) \quad (2)$$

$$I(\Theta; y) = \frac{-\partial \log L(\Theta)}{\partial \Theta \partial \Theta^T}$$

و

$$\mathcal{I}_c(\Theta; y) = E_\Theta(I_c(\Theta; X)|y)$$

و

$$I_c(\Theta; x) = \frac{-\partial^2 \log L_c(\Theta)}{\partial \Theta \partial \Theta^T}.$$

همچنین I یک ماتریس $d \times d$ است.

همگرایی الگوریتم EM به یک مقدار ماکزیمم نسبی برای $L(\Theta)$ وقتی رخ می‌دهد که $J(\Theta^*)$ دارای مقادیر ویژه بزرگتر از یک باشد. یک مقدار ویژه یک برای $J(\Theta^*)$ در همسایگی نقطه مرزی $L(\Theta^*)$ را نشان می‌دهد. در کل نشان داده می‌شود [۳]، $\partial^2 \log L(\Theta^*) / \partial \Theta \partial \Theta^T$ نیمه معین منفی است اگر معین منفی نباشد که در این حالت مقادیر ویژه $J(\Theta^*)$ ، به ترتیب همگی در فواصل $(1, 0]$ یا $[0, 1]$ می‌افتد.

۴ بخش چهارم: برخی ویژگیها و معایب الگوریتم EM

تعدادی از ویژگیها و محسان الگوریتم EM به صورت زیر است:

- ۱) با هر تکرار افزایشی درستنمایی EM ، به طور عددی پایدار است.
- ۲) تحت شرایط کلی مناسبی، حتماً به یک مقداری همگراست.
- ۳) به آسانی قابل استفاده است، زیرا روی محاسبات داده‌های کامل تکیه می‌کند.
- ۴) عموماً برای برنامه نویسی آسان است؛ چراکه هیچ نوع ارزیابی تابع درستنمایی و مشتقانش در آن دخیل نیست.
- ۵) به فضای ذخیره سازی کمی نیازدارد و عموماً می‌توان بایک کامپیوتر کوچک هم انجام داد.
- ۶) هزینه تکرار عموماً پایین است به طوری که می‌توان با تعداد تکرارهای بیشتر مورد نیاز، الگوریتم EM را در مقایسه با دیگر روش‌های مشابه توازن بخشید.
- ۷) کار آنالیزی مورد نیاز در این الگوریتم، خیلی ساده تر از روش‌های مشابه است
- ۸) با مشاهده افزایش یکنوا در درستنمایی روی تکرارها، به سادگی می‌توان همگرایی و خطاهای برنامه‌ریزی را برآنداز نمود.
- ۹) الگوریتم EM می‌تواند برای داده‌های گمشده، سانسور شده و... به کار رود.

برخی معایب و انتقادهای وارد بر این الگوریتم از این قرارند:

- ۱) الگوریتم EM ممکن است حتی در مسائلی که به ظاهر بی‌دردرسند؛ به آرامی همگرا شود.
- ۲) برخلاف اکثر روش‌ها، الگوریتم EM یک روش کلاسیک برای بدست آوردن یک تخمین از ماتریس برآورد پارامترها ارائه نداده است.

۳) الگوریتم EM نیز مانن سایر روش‌های دیگر همگرایی به یک مقدار ماکزیمم مطلق را تضمین نمی‌کند و ممکن است به یک مقدار ماکزیمم نسبی یا نقطه زینی نیز همگرا شود.

۴) در بعضی مسائل، گام E ممکن است به طور آنالیزی به سادگی قابل کنترل نباشد، هرچند در چنین وضعیت‌هایی تأثیر پذیری از یک روش مونت کارلو وجود دارد.

۵ بخش پنجم: روش تسریع ایتنکن چند متغیره

یکی از انتقادهایی که بر الگوریتم EM وارد است این است که همگرایی آن کاملاً کند صورت می‌گیرد. برخی روش‌ها جهت تسریع آن وجود دارد که یکی از آنها که اغلب از آن استفاده می‌شود؛ روش ایتنکن است که در اینجا بیان می‌کنیم.
فرض کنید $\Theta^* \rightarrow \Theta^{(k)}$, وقتی که $k \rightarrow \infty$. بنابراین Θ^* را می‌توان به صورت زیر نوشت

$$\Theta^* = \Theta^{(k)} + \sum_{h=1}^{\infty} (\Theta^{(h+k)} - \Theta^{(h+k-1)}). \quad (3)$$

همچنین داریم:

$$\Theta^{(k+h)} - \Theta^{(h+k-1)} = M(\Theta^{k+h-1}) - M(\Theta^{(h+k-2)}). \quad (4)$$

$$\approx J(\Theta^{(h+k-2)}) (\Theta^{(h+k-1)} - \Theta^{(h+k-2)}) \quad (5)$$

$$\approx J(\Theta^*) (\Theta^{(h+k-1)} - \Theta^{(h+k-2)}) \quad (6)$$

بنابراین برای k ‌های به اندازه کافی بزرگ، $J(\Theta^{(h+k)}) = J(\Theta^*)$. تقریب (۴) و (۵) با بسط سری تیلور خطی $M(\Theta^{(h+k-1)})$ حول نقطه $\Theta^{(h+k-2)}$ بدست آمده است. با تکرار عمل (۶) در (۳) داریم:

$$\begin{aligned} \Theta^* &= \Theta^{(k)} + \sum_{h=0}^{\infty} \{J(\Theta^*)\}^h (\Theta^{(k+1)} - \Theta^{(k)}) \\ &= \Theta^{(k)} + \{I - J(\Theta^*)\}^{-1} (\Theta^{(k+1)} - \Theta^{(k)}) \end{aligned} \quad (7)$$

بنابراین سری توانی $\{I - J(\Theta^*)\}^h$ همگرا می‌شود اگر همه مقادیر $J(\Theta^*)$ بیان صفر و یک باشد.

حال می‌خواهیم با استفاده از روش ایتنکن، سرعت همگرایی دنباله $\{\Theta_A^{(k)}\}$ را تسریع بخشیم که $\Theta_A^{(k)}$ به صورت زیر تعریف می‌شود

$$\Theta_A^{(k+1)} = \Theta_A^{(k)} + \{I - J(\Theta_A^{(k)})\}^{-1} (\Theta_{EMA}^{(k+1)} - \Theta_A^{(k)}) \quad (8)$$

که در تکرار $(k+1)$ با استفاده از $\Theta_A^{(k)}$ در تکرار EM بدست آمده است. این روش را تا زمانی که در تکرار $(k+1)$ ام، به تولید اولین مقدار $\Theta_{EMA}^{(k+1)}$ در فرآیند EM با استفاده از $\Theta^{(k)}$ بررسیم، ادامه می‌دهیم.

این روش در سال ۱۹۸۲ توسط لوئیس (*Louis*) برای تسريع بخشیدن به همگرایی الگوریتم EM پیشنهاد شد [۲].

لوئیس در سال ۱۹۸۲، برای برآورد $\Theta_A^{(k)}$ موفق به یافتن ارتباط بین (۲) و (۸) شد که داریم

$$\Theta_A^{(k+1)} = \Theta_A^{(k)} + I^{-1}(\Theta_A^{(k)}; y)\mathcal{I}_c(\Theta_A^{(k)}; y)\left(\Theta_{EMA}^{(k+1)} - \Theta_A^{(k)}\right). \quad (9)$$

ارتباط (۲) برای مقادیر Θ به غیر از Θ^* به عنوان تقریب نسبتاً مفیدی برای $MLE\Theta^*$ است و همچنین نباید برخی تکرارهای EM به طور عادت انجام شود. استفاده از (۹) به عنوان تقریبی معادل استفاده از الگوریتم نیوتون رافسون برای پیدا کردن صفر آماره نمره (داده‌های غیر کامل) $S(y; \Theta) = \partial \log L(\Theta) / \partial \Theta$ در سال ۱۹۸۹ *meilijson* توسط مورد توجه قرار گرفت. می‌توان نشان داد [۳] که

$$S(y; \Theta_A^{(k)}) \approx \mathcal{I}_c(\Theta_A^{(k)}; y)\left(\Theta_{EMA}^{(k+1)} - \Theta_A^{(k)}\right) \quad (10)$$

با جایگذاری این نتایج در (۹) داریم:

$$\Theta_A^{(k+1)} \approx \Theta_A^{(k)} + I^{-1}(\Theta_A^{(k)}; y)s(y; \Theta_A^{(k)}). \quad (11)$$

طرف راست (۱۱)، فرآیند بدست آمده در (۱) امین تکرار روش نیوتون رافسون بکار رفته برای یافتن صفر $S(y; \Theta)$ است.

بنابراین استفاده از روش ایتكن به فرم استاندارد شده به طور اساسی معادل با روش نیوتون رافسون بکار رفته برای $S(y; \Theta)$ است که توسط لوئیس (۱۹۸۲) بیان شد.

در سال ۱۹۸۹ *Meilijson* نشان داد که در حالتی که مشاهدات w_1, w_2, \dots, w_n مستقل و همتوزیع باشند، در روش نیوتون رافسون (۱۱)، می‌توان $I(\Theta; y)$ را با ماتریس اطلاع تجربی^۱ زیر تقریب کرد:

$$I_e(\Theta; y) = \sum_{j=1}^n s(w_j; \Theta)s(w_j; \Theta)^T - n^{-1}S(y; \Theta)S^T(y; \Theta) \quad (12)$$

که $s(w_j; \Theta)$ آماره نمره بر اساس مشاهده تکی w_j است.

^۱ empirical information matrix

۶ بخش ششم: مثالهایی از تسریع الگوریتم EM

۱.۶ مثال ۱ - یک آزمایش چند جمله‌ای

یک مثال جهت نشان دادن چگونگی استفاده از روش تسریع اینکن که توسط لوئیس در سال ۱۹۸۲ بیان شد؛ مثال معروف داده‌های یک توزیع چندجمله‌ای است که در مقاله [۱] *DLR* [بیان شده است که بردار مشاهده شده $y = (y_1, y_2, y_3, y_4)^T$ از یک توزیع چندجمله‌ای به ترتیب با احتمالهای

$$\frac{1}{2} + \frac{1}{4}\theta, \frac{1}{4}(1-\theta), \frac{1}{4}(1-\theta), \frac{1}{4}\theta$$

گرفته شده است. پارامتر θ براساس بردار y برآورد می‌شود؛ در صورتی که بردار $y = (125, 18, 20, 34)^T$ باشد آنگاه $0/6268215 = \hat{\theta}$ بدست می‌آید. حال اگر اولین خانه از خانه‌های چهارگانه اصلی را به دو خانه با احتمالهای $\frac{1}{4}$ و $\frac{1}{4}\theta$ تقسیم کنیم و y_{11} و y_{12} نشاندهنده تقسیم y_1 باشد به طوری که $y_1 = y_{11} + y_{12}$. کاملاً واضح است که برآورد θ به سادگی قابل محاسبه نیست. پس از به کاربردن الگوریتم EM [۶]، درگامهای E و M داریم:

$$\theta^{(k+1)} = \frac{y_{12}^{(k)} + y_4}{n - (y_{11}^{(k)})}$$

که

$$y_{11}^{(k)} = \frac{\frac{1}{4}y_1}{(\frac{1}{2} + \frac{1}{4}\theta^{(k)})}$$

و

$$y_{12}^{(k)} = y_1 - y_{11}^{(k)}.$$

با انتخاب $\theta^{(0)} = 0/626821484$ شرط y_1 در دنباله $\{\theta^{(k)}\}$ بعد از هشت مرحله تکرار به مقدار $\theta^{(8)} = 0/626821484$ دست می‌یابیم. جدول نتایج الگوریتم EM برای این مثال در زیر آورده شده است: برای تسریع همگرایی آن، با بکاربردن فرمول (۶) بعد از دو مرحله تکرار EM داریم:

$$\theta_A^{(2)} = \theta^{(2)} = 0/626328$$

$$I^{-1}(\theta_A^{(2)}; y) I_c(\theta_A^{(2)}; y) = \frac{434/79}{376/95}$$

$$= 1/153442.$$

k	تکرار	$\theta^{(k)}$	$\log L(\theta^{(k)})$
۰		۰/۵۰۰۰۰۰۰۰۰	۶۴/۶۲۹۷۴
۱		۰/۶۰۸۲۴۷۴۲۲	۶۷/۳۲۰۱۷
۲		۰/۶۲۴۳۲۱۰۵۱	۶۷/۳۸۲۹۲
۳		۰/۶۲۶۴۸۸۸۷۹	۶۷/۳۸۴۰۸
۴		۰/۶۲۶۷۲۷۷۳۲۳	۶۷/۳۸۴۱۰
۵		۰/۶۲۶۸۱۵۶۲۲	۶۷/۳۸۴۱۰
۶		۰/۶۲۶۸۲۰۷۱۹	۶۷/۳۸۴۱۰
۷		۰/۶۲۶۸۲۱۳۹۵	۶۷/۳۸۴۱۰
۸		۰/۶۲۶۸۲۱۴۸۴	۶۷/۳۸۴۱۰

جدول ۱: نتایج الگوریتم EM برای داده‌های مدل چند جمله‌ای

$\theta_A^{(۳)} = ۰/۶۲۶۳۳۸ + ۱/۱۵۳۴۴۲(\theta_{EMA}^{(۳)} - ۰/۶۲۶۳۳۸) = ۰/۶۲۶۸۲۱۶$
 و $\theta^{(۴)} = ۰/۶۲۶۸۱۲$. می‌توان ملاحظه نمود که $\theta_A^{(۳)}$ از $\theta^{(۴)}$ نزدیکتر است.
 که چهارمین تکرار با مقادیر اصلی (بدون تسریع) الگوریتم EM است به $\hat{\theta} = ۰/۶۲۶۸۲۱۵$ نزدیکتر است.

۲.۶ مثال ۲. آمیزه هندسی

در سال ۱۹۸۹، *Meilijson* مثالی دیگر در مورد آمیزه هندسی برای برآورد ML بردار پارامتر

$$\Theta = (\pi_1, p_1, p_2)^T$$

بیان داشت. در آمیزه هندسی دو پارامتری

$$f(w; \Theta) = \sum_{i=1}^2 \pi_i f(w_i; p_i)$$

که به ازای $i = 1, 2$ داریم:

$$f(w; p_i) = p_i(1 - p_i)^{w-1} \quad w = 1, 2, 3, \dots (0 \leq p_i \leq 1).$$

در حالت کلی می‌توان بردار داده کامل x را به دو قسمت داده‌های مشاهده شده (y) و داده‌های گمشده (z) تقسیم بندی نمود یعنی

$$x = (y^T, z^T)^T$$

که

$$z = (z_1^T, z_2^T, \dots, z_n^T)^T.$$

$$y = (w_1, w_2, \dots, w_n)^T$$

که وقتی $z_{ij} = (z_j)_i$ است که این مقدار مشاهده شده از i امین مؤلفه معلوم باشد مقداریک و وقتی گمشده باشد مقدار صفر را می‌گیرد. ($i = 1, 2; j = 1, 2, \dots, n$)

لگاریتم درستنماهی داده های کامل به صورت زیر است:

$$\log L_c(\Theta) = \sum_{i=1}^r \sum_{j=1}^n z_{ij} \{ \log \pi_i + \log p_i + (w_j - 1) \log(1 - p_i) \} \quad (13)$$

$$= \sum_{i=1}^r \{ n_i (\log \pi_i + \log p_i) + \left(\sum_{j=1}^n z_{ij} w_j - n_i \right) \log(1 - p_i) \}, \quad (14)$$

که

$$n_i = \sum_{j=1}^n z_{ij}$$

مطابق گام E در تکرار $(k+1)$ م الگوریتم EM ، تابع Q به صورت زیر است:

$$Q(\Theta; \Theta^{(k)}) = \sum_{i=1}^r \{ n_i^{(k)} (\log \pi_i + \log p_i) + \left(\sum_{j=1}^n z_{ij}^{(k)} w_j - n_i^{(k)} \right) \log(1 - p_i) \} \quad (15)$$

که

$$n_i^{(k)} = \sum_{j=1}^n z_{ij}^{(k)}$$

و

$$z_{ij}^{(k)} = \frac{\pi_i^{(k)} p_i^{(k)} (1 - p_i^{(k)})^{w_j - 1}}{\sum_{h=1}^r \pi_h^{(k)} p_h^{(k)} (1 - p_h^{(k)})^{w_j - 1}}$$

امید ریاضی شرطی Z_{ij} به شرط y باشد.

$$\pi_i^{(k+1)} = \frac{n_i^{(k)}}{n}$$

و

$$p_i(k+1) = \left(\sum_{j=1}^n z_{ij}^{(k)} w_j / n_i^{(k)} \right)^{-1} \quad (i = 1, 2)$$

تکرار کام	$\pi_1^{(k)}$	$p_1^{(k)}$	$p_2^{(k)}$
۰	۰/۲	۰/۱	۰/۲
۵	۰/۱۴۶	۰/۲۴۸	۰/۷۰۴
۱۰	۰/۲۰۱	۰/۲۶۴	۰/۷۶۷
۱۰۰	۰/۴۰۴۶۵	۰/۳۵۳۹۵	۰/۹۰۲۲۰

جدول ۲: نتایج الگوریتم EM برای مدل آمیزه هندسی

برای محاسبه ماتریس اطلاع تجربی $I_e(\Theta; y)$ با استفاده از تعریف (۱۲) باید (w_j) را محاسبه کنیم که $s(w_j)$ آماره نمره داده غیرکامل بر اساس مشاهده تکی j ام ($= j$) است. داریم:

$$\begin{aligned} s(w_j; \Theta) &= \partial \log L_j(\Theta) / \partial \Theta \\ &= E_{\Theta^{(k)}} \{ \partial \log L_{cj}(\Theta) / \partial \Theta | y \} \\ &= \partial Q_j(\Theta, \Theta^{(k)}) / \partial \Theta \end{aligned}$$

که

$$Q_j(\Theta, \Theta^{(k)}) = E_{\Theta^{(k)}} \{ \log L_{cj}(\Theta) | y \}.$$

بعد از محاسبات، ماتریس اطلاع تجربی $I_e(\Theta; y)$ به صورت زیر است:

$$s_1(w_j; \Theta) = (z_{\setminus j}^{(k)} - \pi_1) / (\pi_1 \pi_2)$$

$$s_{i+1}(w_j; \Theta) = z_{ij}^{(k)} (1 - w_j p_i) / \{p_i (1 - p_i)\} \quad (i = 1, 2).$$

روش‌های مختلفی را برای برآش مدل آمیزه هندسی به وسیله تابع $MLE\theta$ بدست آمده به صورت زیر شد: *Meilijson* درستنمایی با 500 مشاهده انجام داد که

$$\hat{\theta} = (0/40647, 0/35466, 0/90334)^T.$$

مقدار شروع θ برای الگوریتم EM در تکرارهای پنجم، دهم و سیم در جدول ۲ آورده شده است. همچنین مقادیر ترسیح الگوریتم EM نیز در جدول ۳ داده شده است.

تکرار ام $ k$	$\pi_1^{(k)}$	$p_1^{(k)}$	$p_2^{(k)}$
۰	۰/۱۴۶	۰/۲۴۸	۰/۷۰۴
۱	۰/۴۳	۰/۵۲	۰/۸۷
۲	۰/۶۰	۰/۴۹	۰/۹۸۶
۳	۰/۴۵	۰/۴۲	۰/۹۱
۴	۰/۴۲	۰/۳۷	۰/۹۱
۵	۰/۴۰۶	۰/۳۵۶	۰/۹۰۳
.	.	.	.
.	.	.	.
۸	۰/۴۰۶۴۷	۰/۳۵۴۶۶	۰/۹۰۳۲۴

جدول ۳: نتایج تسریع الگوریتم EM برای مدل آمیزه هندسی

مراجع

- Dempster, A.P., Laird, N.M., and Rubin, D.B. (1977). *Maximum Likelihood from Incomplete Data Via the EM Algorithm (with Discussion)*. Journal of the Royal Statistical Society Series B 39, 1-38.
- Louis, T.A. (1982). *Finding the obseved information matrix when using the EM algorithm* ., Journal of the Royal Statistical Society Series B 44, 226-233.
- McLachlan, G.J. , and Krishnan, T. (1995). *The EM Algorithm and Extensions.*, New York: Wiley. To be published.
- McLachlan, G.J. *On Aitken's Method and Other Approaches for Accelerating Convergence of the EM Algorithm*.
- Meilijson, I. (1989). *A Fast Improvement to the EM Algorithm on Its Own Terms*. Journal of the Royal Statistical Society Series B 51, 127-138.
- میرحسینی، سید محسن، الگوریتم EM و برآورد پارامترها در داده‌های غیرکامل، پایان نامه کارشناسی ارشد، ۱۳۷۸.

برآوردهای مینیماکس مجاز پارامتر مقیاس خانواده کای - دو تبدیل یافته در فضای پارامتری کراندار

نادر نعمت‌الهی^۱، محمد جعفری جوزانی^۲

P11086

^۱دانشگاه علامه طباطبائی

^۲دانشگاه شهید بهشتی

چکیده: روش‌های معمول برآوردهایی از قبیل روش حداکثر درستنمایی در فضاهای پارامتری کراندار تحت تابع زیان مربع خطأ، عموماً منجر به برآوردهای غیر مجاز می‌گردند. برآوردهایی به روش مینیماکس روش متداول دیگری است که در چنین فضاهایی مورد استفاده قرار می‌گیرد و عموماً برآوردهای مجاز را به دست می‌دهد. در این مقاله نتایج به دست آمده توسط ایدن (۱۹۹۵) در رابطه با برآوردهای پارامتر مقیاس توزیع گاما در فضای پارامتری کراندار تحت تابع زیان مربع خطأ پایای مقیاس را به خانواده کای - دو تبدیل یافته تعمیم می‌دهیم. در این خانواده برآوردهای مینیماکس مجاز پارامتر مقیاس را در فضاهای پارامتری کراندار تحت تابع زیان مربع خطأ پایای مقیاس به دست آورده و چند مثال ارایه می‌دهیم.

واژه‌های کلیدی: تابع زیان مربع خطأ پایای مقیاس، برآوردهای مینیماکس مجاز، فضای پارامتری کراندار، خانواده کای - دو تبدیل یافته.

۱ مقدمه

یکی از مسائلی که بیشتر موقع در عمل با آن رویرو هستیم مسئله استنباط کردن و به ویژه برآوردهایی در فضای پارامتری کراندار است. روش‌های معمول برآوردهایی از قبیل روش حداکثر درستنمایی در فضاهای پارامتری کراندار تحت تابع زیان مربع خطأ عموماً منجر به برآوردهای غیر مجاز می‌گردند. برآوردهایی به روش مینیماکس روش متداول دیگری است که در چنین فضاهایی مورد استفاده قرار می‌گیرد. مسئله برآوردهایی در فضاهای پارامتری کراندار در ابتدا توسط برانک (۱۹۵۵) و ون ایدن (۱۹۵۷) مورد بررسی قرار گرفت. این نویسندها توجه خود را روی وجود برآوردهای

حداکثر درستنمایی معطوف کردند. یکی از اولین کسانی که به مطالعه درباره برآورده مینیماکس مجاز در فضاهای پارامتری کراندار پرداخت کاتز(۱۹۶۱) بود. کاتز توانست یک برآورده مینیماکس مجاز برای میانگین توزیع نرمال در فضای پارامتری از پایین کراندار ، تحت تابع زیان مربع خطابه دست آورد. بعدها بیکل ، کسلا و استرادمن (۱۹۸۱) مطالعاتی درباره برآوردهای به روش مینیماکس در فضای پارامتری کراندار انجام دادند و بری (۱۹۹۳) برآورده مینیماکس مجاز برای توزیع نمایی با تکیه گاه $\theta \in [0, \infty)$ یا $a \in [a, b]$ یا $b \in \theta$ را به دست آورد. تابع زیان مورد استفاده دربیشتر مطالعات انجام شده، تابع زیان مربع خطابه است و برآوردهای توزیع زیان مربع خطای وزنی و در حالت خاص ، مربع خطای پایای مقیاس مورد توجه کمتری قرار گرفته است.

تحت تابع زیان مربع خطای پایای مقیاس، ون ایدن (۱۹۹۵ و ۲۰۰۰) برآورده مینیماکس مجاز برای پارامتر مقیاس توزیع گاما و برآورده مینیماکس برای پارامتر مقیاس توزیع F در فضای پارامتری کراندار را به دست آورد.

در این مقاله نتایج به دست آمده توسط ون ایدن (۱۹۹۵) را به خانواده کای - دو تبدیل یافته که بوسیله رحمان^۱ و گوپتا^۲ (۱۹۹۳) معرفی شده است ، تعمیم می دهیم و در این خانواده تحت فضای پارامتری کراندار برآورده مینیماکس مجاز پارامتر مقیاس را تحت تابع زیان مربع خطای مقیاس به دست می آوریم . بدین منظور دربخش ۲ نتایج به دست آمده توسط ون ایدن (۱۹۹۵) را ارایه می دهیم. دربخش ۳ خانواده کای - دو تبدیل یافته را معرفی کرده و چند توزیع متعلق به این خانواده را مشخص می کنیم . در بخش ۴ برآورده مینیماکس مجاز پارامتر مقیاس خانواده کای - دو تبدیل یافته در فضای پارامتری کراندار را به دست آورده و چند مثال ارایه می دهیم .

۲ برآورده مینیماکس مجاز پارامتر مقیاس توزیع گاما در فضای پارامتری از پایین کراندار

در این بخش نتایج به دست آمده توسط ون ایدن (۱۹۹۵) در رابطه با برآورده مینیماکس مجاز پارامتر مقیاس توزیع گاما در فضای پارامتری از پایین کراندار ، تحت تابع زیان مربع خطای پایای مقیاس را می آوریم. برای اثبات این نتایج می توان به ون ایدن (۱۹۹۵) مراجعه نمود.

^۱ Rahman

^۲ Gupta

فرض کنید X دارای توزیع گاما با تابع چگالی احتمال زیر باشد:

$$f_\theta(x) = \frac{1}{\Gamma(\alpha)} \cdot \frac{x^{\alpha-1}}{\theta^\alpha} \cdot e^{-\frac{x}{\theta}}, x > 0 \quad (1)$$

که در آن $\alpha > 0$ معلوم، θ پارامتر مقیاس نامعلوم، $a \in [a, \infty)$ و $a > 0$ معلوم است. تابع زیان مورد بررسی، تابع زیان مریع خطای پایای مقیاس است. یعنی

$$L(\theta, \delta) = (\frac{\delta}{\theta} - 1)^2, \quad \delta \geq a, \theta \geq a \quad (2)$$

همچنین فرض کنید $\hat{\theta}$ برآورده‌گری برای θ باشد که در شرط زیر صدق کند

$$P_\theta(\hat{\theta} \geq a) = 1, \quad \forall \theta \geq a \quad (3)$$

برآورده‌گر $\hat{\theta}$ که در این بخش معرفی می‌شود به صورت زیر می‌باشد

$$\delta(x) = \lim_{n \rightarrow \infty} \delta_n(x) \quad (4)$$

که در آن $(\delta_n(x), \text{برای هر } n = 1, 2, \dots, 1)$ یک برآورده‌گر بیز نسبت به توزیع پیشین زیر است

$$\begin{aligned} \pi_n(\theta) &= \frac{1}{n\theta} \left(\frac{a}{\theta} \right)^{\frac{1}{n}} \\ &= \frac{a^{\frac{1}{n}}}{n\theta^{1+\frac{1}{n}}}, \quad \forall \theta \geq a \end{aligned} \quad (5)$$

با به کار بردن (1) و (5) توزیع پیشین θ به شرط $x = x$ مناسب است با

$$\pi(\theta|x) \propto \frac{1}{\theta^{\alpha_n+1}} e^{-\frac{x}{\theta}}, \quad x > 0, \quad \theta \geq a \quad (6)$$

که در آن $\alpha_n = \alpha + \frac{1}{n}$. ون ایدن (۱۹۹۵) توسط قضایای زیر برآورد مینیماکس مجاز پارامتر θ را تحت تابع زیان (2) به دست آورد.

قضیه ۱ برای تابع چگالی احتمال (1) برآورده‌گر θ ، برای هر $x > 0$ ، نسبت به پیشین (5)، تحت تابع زیان (2) عبارت است از

$$\delta_n(x) = \frac{x}{\alpha+1+\frac{1}{n}} \left(1 + \frac{\left(\frac{x}{\alpha}\right)^{\alpha+1+\frac{1}{n}} e^{-\frac{x}{\alpha}}}{\int_0^{\frac{x}{\alpha}} t^{\alpha+1+\frac{1}{n}} e^{-t} dt} \right) \quad (7)$$

بعلاوه

$$r_n(\delta_n) < \infty , \quad \forall n = 1, 2, \dots$$

$$\lim_{n \rightarrow \infty} r_n(\delta_n) = \frac{1}{\alpha+1}$$

که در آن، r_n مخاطره بیز متناظر با δ_n است.

قضیه ۲ برای تابع چگالی احتمال (۱) با پارامتر مقیاس نامعلوم θ و پارامتر شکل معلوم a که در آن $a \geq \theta$ و $a > \theta$ معلوم است، برآورده شود.

$$\delta(x) = \frac{x}{\alpha+1} \left\{ 1 + \frac{\left(\frac{x}{a}\right)^{\alpha+1} e^{-\frac{x}{a}}}{\int_0^{\frac{x}{a}} t^{\alpha+1} e^{-t} dt} \right\} \quad (8)$$

تحت تابع زیان (۲)، یک برآورد مینیماکس مجاز برای θ است. بعلاوه مقدار مینیماکس برابر $\frac{1}{\alpha+1}$ می‌باشد که در نقطه $a = \theta$ و یا هنگامی که θ به بینهایت میل کند به دست می‌آید.

۳ خانواده کای - دو تبدیل یافته و چند حالت خاص آن

در این بخش خانواده کای - دو تبدیل یافته که توسط رحمان و گوپتا (۱۹۹۳) معرفی شده است را مورد بررسی قرار می‌دهیم. فرض کنید $\tilde{X} = (X_1, X_2, \dots, X_m)$ برداری تصادفی متعلق به خانواده نمایی یک پارامتری از توزیعها با تابع چگالی احتمال توأم زیر باشد

$$f(\tilde{x}, \eta) = \exp\{a(\tilde{x})b(\eta) + c(\eta) + h(\tilde{x})\} \quad (9)$$

برای مثال توزیعهای پواسن، دوجمله ای، نرمال و گاما با اختیار $a(\tilde{x})$ و $c(\eta)$ مناسب به خانواده فوق تعلق دارند. در قضیه زیر نشان می‌دهیم تحت وجود بعضی شرایط $h(\tilde{x}) = 2a(\tilde{x})b(\eta) - 2a(\tilde{x})$ دارای توزیع گاما است.

قضیه ۳: در خانواده نمایی یک پارامتری (۹) عبارت $2a(\tilde{X})b(\eta) - 2a(\tilde{X})$ دارای توزیع گاما با پارامترهای $\frac{k}{2}$ و 2 است اگر و تنها اگر

$$\frac{2c'(\eta)b(\eta)}{b'(\eta)} = k \quad (10)$$

که در آن k مثبت و مستقل از η می‌باشد.

برهان : فرض کنید $k = \frac{2c'(\eta)b(\eta)}{b'(\eta)}$ از آنجا که (۹) یک تابع چگالی احتمال است ، بنابراین

$$\int e^{a(x)b(\eta)+c(\eta)+h(x)} dx = 1$$

و یا به طور معادل

$$\int e^{a(x)b(\eta)+h(x)} dx = e^{-c(\eta)} \quad (11)$$

از طرفی با استفاده از رابطه (۱۰) داریم

$$c'(\eta) = \frac{k}{2} \frac{b'(\eta)}{b(\eta)}$$

و با انتگرال گیری از طرفین رابطه فوق نسبت به η داریم

$$c(\eta) = \frac{k}{2} Ln|b(\eta)| + k_1$$

که در آن k_1 ثابت انتگرال گیری است . بنابراین می توان رابطه (۱۱) را به صورت زیر بازنویسی کرد

$$\int e^{a(x)b(\eta)+h(x)} dx = e^{-\frac{k}{2}Ln|b(\eta)| - k_1}$$

قرار می دهیم $Y = -2a(X)b(\eta)$ عبارت است از

$$\begin{aligned} \varphi_Y(t) &= E \left[e^{-2ita(X)b(\eta)} \right] \\ &= e^{c(\eta)} \int e^{a(x)b(\eta)(1-2it)+h(x)} dx \\ &= e^{\frac{k}{2}ln|b(\eta)| - k_1} e^{\frac{k}{2}ln|b(\eta)|(1-2it) + k_1} \\ &= (1-2it)^{-\frac{k}{2}} \end{aligned}$$

که همان تابع مشخصه توزیع گاما با پارامترهای $\frac{k}{2}$ و ۲ است . از آنجا که تابع مشخصه به طور یکتا توزیع متغیر تصادفی را مشخص می کند لذا $-2a(X)b(\eta)$ دارای توزیع گاما با پارامترهای مذکور است .

برعکس فرض کنید $Y = -2a(X)b(\eta)$ دارای توزیع گاما با پارامترهای $\frac{k}{2}$ و ۲ باشد . تابع چگالی احتمال Y عبارت است از

$$f_\eta(y) = \frac{e^{-\frac{y}{2}} y^{\frac{k}{2}-1}}{\Gamma(\frac{k}{2}) 2^{\frac{k}{2}}}$$

بنابراین تابع چگالی احتمال $a(x)$ به صورت زیر است

$$\begin{aligned} f_\eta(a(x)) &= \frac{e^{a(x)b(\eta)} \{-b(\eta)\}^{\frac{k}{\gamma}} \{a(x)\}^{\frac{k}{\gamma}-1}}{\Gamma(\frac{k}{\gamma})} \\ &= e^{\{a(x)b(\eta)+\frac{k}{\gamma}Ln|b(\eta)|+(\frac{k}{\gamma}-1)Ln\{a(x)\}-Ln\{\Gamma(\frac{k}{\gamma})\}\}} \end{aligned}$$

که متعلق به خانواده نمایی یک پارامتری (۹) است و در آن

$$c(\eta) = \frac{k}{\gamma} Ln|b(\eta)| + k_1 \quad (12)$$

لازم به ذکر این نکته است که اگر $\eta > 0$ آنگاه می بایست $b(\eta)$ منفی باشد . در رابطه بالا k_1 ثابت انتگرالگیری می باشد و با مشتق گیری از طرفین رابطه (۱۲) داریم

$$\frac{2c'(\eta)b(\eta)}{b'(\eta)} = k$$

تعريف ۱: خانواده توزیعهای نمایی یک پارامتری با تابع چگالی احتمال (۹) که در شرط (۱۰) صدق می کند هنگامی که یک عدد صحیح مثبت باشد را خانواده توزیعهای کای - دو تبدیل یافته گویند.

برای مثال توزیع توانم یک نمونه تصادفی m تایی از توزیع لاغ نرمال $LN(0, \sigma^2)$ عبارت است از

$$f_{\sigma^2}(x) = (\frac{1}{2\pi})^{-\frac{m}{2}} (\sigma^2)^{-\frac{m}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^m (lnx_i)^2}$$

که متعلق به خانواده کای - دو تبدیل یافته با

$$k = m \quad , \quad a(x) = \frac{1}{2} \sum_{i=1}^m (lnx_i)^2$$

$$b(\sigma) = -\frac{1}{\sigma^2}, c(\sigma) = -mln\sigma, -2a(x)b(\sigma) = \frac{1}{\sigma^2} \sum_{i=1}^m (lnx_i)^2$$

می باشد. همچنین توزیع توانم یک نمونه تصادفی m تایی از توزیع وایبل $W(\alpha, \beta)$ با α معلوم عبارت است از

$$f_\beta(x) = \beta^m \alpha^m e^{-\beta \sum_{i=1}^m x_i^\alpha} (\prod_{i=1}^m x_i^{\alpha-1})$$

که متعلق به خانواده کای - دو تبدیل یافته با

$$k = 2m, a(\tilde{x}) = \sum_{i=1}^m x_i^\alpha, b(\beta) = -\beta, c(\beta) = m \ln \beta, -2a(\tilde{x})b(\beta) = 2\beta \sum_{i=1}^m x_i^\alpha$$

می باشد. در جدول (۱) بعضی از توزیعهای متعلق به خانواده کای - دو تبدیل یافته نشان داده شده است.

۴ برآورده مینیماکس مجاز پارامتر مقیاس در خانواده کای - دو تبدیل یافته

در این بخش برآورده مینیماکس مجاز پارامتر مقیاس θ که در آن $a \geq \theta > 0$ معلوم است در خانواده کای - دو تبدیل یافته تحت تابع زیان مربع خطای پایای مقیاس را به دست می آوریم.

در رابطه (۱۲) قرار می دهیم $\theta = -\frac{1}{b(\eta)}$ در این صورت

$$\begin{aligned} e^{c(\eta)} &= \left[-\frac{1}{b(\eta)} \right]^{-\frac{k}{\alpha}} e^{k_1} \\ &= \theta^{-\frac{k}{\alpha}} e^{k_1} \end{aligned}$$

بنابر این خانواده توزیعهای نمایی (۹) (کای - دو تبدیل یافته) را می توان به صورت زیر نوشت

$$\begin{aligned} e^{-\frac{a(\tilde{x})}{[-1/b(\eta)]} + c(\eta) + h(\tilde{x})} &= e^{h(\tilde{x}) + k_1} e^{c(\eta)} e^{-\frac{a(\tilde{x})}{[-1/b(\eta)]}} \\ &= c(\tilde{x}, m) \theta^{-\frac{k}{\alpha}} e^{-\frac{a(\tilde{x})}{\theta}} \end{aligned}$$

که در آن $c(\tilde{x}, m) = e^{h(\tilde{x}) + k_1}$ است. همچنین با توجه به اینکه

$$-2a(\tilde{X})b(\eta) = \frac{2a(X)}{\theta} \sim \Gamma\left(\frac{k}{\alpha}, 2\right)$$

بنابر این $a(X) \sim \Gamma(\frac{k}{\gamma}, \theta)$. در نتیجه خانواده کای - دو تبدیل یافته به صورت زیر تبدیل می شود

$$f_\theta(x) = c(x, m)\theta^{-\frac{k}{\gamma}} e^{-\frac{a(x)}{\theta}}, \quad k > 0, \theta > 0, a(X) \sim \Gamma(\frac{k}{\gamma}, \theta) \quad (13)$$

در این بخش با درنظر گرفتن توزیع پیشین (۵) و تابع زیان به فرم (۲) و شرایط (۳) و (۴) روی برآوردها ، می خواهیم با استفاده از نتایج به دست آمده توسط ون ایدن (۱۹۹۵) ، برآوردگر مینیماکس مجاز پارامتر θ را در خانواده کای - دو تبدیل یافته به دست آوریم.

قضیه ۴ برای خانواده توزیعهای کای - دو تبدیل یافته به فرم (۱۳) با $\theta \geq a > 0$ معلوم، برآورد بیز $\hat{\theta} = \frac{1}{b(\eta)}$ برای هر $x > 0$ ، نسبت به توزیع پیشین (۵)، تحت تابع زیان (۲) عبارت است از

$$\delta_n(x) = \frac{a(x)}{\frac{k}{\gamma} + 1 + \frac{1}{n}} \left(1 + \frac{\left(\frac{a(x)}{\theta} \right)^{\frac{k}{\gamma} + 1 + \frac{1}{n}} e^{-\frac{a(x)}{\theta}}}{\int_0^{\frac{a(x)}{\theta}} t^{\frac{k}{\gamma} + 1 + \frac{1}{n}} e^{-t} dt} \right) \quad (14)$$

به علاوه

$$r_n(\delta_n) < \infty, \quad \forall n = 1, 2, \dots$$

$$\lim_{n \rightarrow \infty} r_n(\delta_n) = \frac{2}{k+2} \quad (15)$$

که در آن، r_n مخاطره بیز متناظر با δ_n است.

اثبات : با استفاده از (۱۳) و (۵)، تابع چگالی پسین θ به شرط x متناسب است با

$$\pi(\theta|x) \propto \frac{1}{\theta^{\alpha_n + 1}} e^{-\frac{a(x)}{\theta}}, \quad x > 0, \theta \geq a$$

که در آن $\frac{1}{n} + \frac{k}{\alpha_n} = \alpha$. این تابع چگالی پسین همانند چگالی پسین (۶) می‌باشد که در آن x و α به ترتیب به $a(\tilde{x})$ و $\frac{k}{\alpha}$ تبدیل شده‌اند . بنابراین برآورده بیز (۱۴) و خواص آن را می‌توان به موازات اثبات قضیه (۱) ارایه شده توسط ون ایدن (۱۹۹۵) به دست آورد.

از قضیه (۴) برای هر $x > 0$ به دست می‌آوریم که

$$\delta(\tilde{x}) = \lim_{n \rightarrow \infty} \delta_n(\tilde{x}) = \frac{a(\tilde{x})}{\frac{k}{\alpha} + 1} \left\{ 1 + \frac{\left(\frac{a(\tilde{x})}{\alpha}\right)^{\frac{k}{\alpha}+1} e^{-\frac{a(\tilde{x})}{\alpha}}}{\int_0^{\frac{a(\tilde{x})}{\alpha}} t^{\frac{k}{\alpha}+1} e^{-t} dt} \right\} \quad (16)$$

مینیماکس و مجاز بودن برآورد (۱۶) در قضیه زیر آورده شده است . اثبات این قضیه شبیه اثبات قضیه (۲) می‌باشد که درون ایدن (۱۹۹۵) آمده است.

قضیه ۵ : برای خانواده توزیعهای تبدیل یافته به فرم (۱۳) با $a > 0$ و $\theta \geq 0$ معلوم، برآورد مینیماکس مجاز $\frac{1}{b(\eta)} - \theta = 0$ برای هر $\tilde{x} > a(\tilde{x})$ ، تحت تابع زیان به فرم (۲) به صورت (۱۶) می‌باشد. به علاوه مقدار مینیماکس برابر $\frac{2}{k+2}$ است. این مقدار در $\theta = a$ یا هنگامی که θ یه سمت بینهایت میل کند به دست می‌آید.

مثال ۱ : فرض کنید X یک متغیر تصادفی با توزیع $\Gamma(\alpha, \beta)$ باشد که در آن $\alpha > 0$ مقداری معلوم است در این صورت

$$f_\beta(x) = \frac{1}{\Gamma(\alpha)\beta^\alpha} x^{\alpha-1} e^{-\frac{x}{\beta}} , \quad x > 0$$

که متعلق به خانواده کای - دو تبدیل یافته با

$$k = 2\alpha, a(\tilde{x}) = x, b(\beta) = -\frac{1}{\beta}, c(\beta) = -\alpha \ln \beta, -2a(\tilde{x})b(\beta) = \frac{2x}{\beta}, \theta = \beta$$

است . همچنین $2\alpha = \frac{2c'(\beta)b(\beta)}{b'(\beta)}$ بنابر این طق قضیه (۵) برآورده مینیماکس مجاز پارامتر $\theta = \beta$ با شرط $\beta \geq a$ عبارت است از

$$\delta(x) = \frac{x}{\alpha+1} \left\{ 1 + \frac{\left(\frac{x}{\alpha}\right)^{\alpha+1} e^{-\frac{x}{\alpha}}}{\int_0^{\frac{x}{\alpha}} t^{\alpha+1} e^{-t} dt} \right\}$$

که همان برآورد به دست آمده در قضیه ۲ می باشد که توسط ون ایدن (۱۹۹۵) ارایه گردیده است.

مثال ۲ : فرض کنید X_1, X_2, \dots, X_m یک نمونه تصادفی از $(\sigma^2, N(\mu, \sigma^2))$ باشند. در این صورت

$$f_{\sigma^2}(x) = (\frac{1}{2\pi})^{-\frac{m}{2}} (\sigma^2)^{-\frac{m}{2}} e^{-\frac{\sum_{i=1}^m x_i^2}{2\sigma^2}}$$

که متعلق به خانواده کای - دو تبدیل یافته با

$$k = m, a(x) = \frac{1}{2} \sum_{i=1}^m x_i^2, b(\sigma) = -\frac{1}{\sigma^2},$$

$$c(\sigma) = -mln\sigma, -2a(x)b(\sigma) = \frac{\sum_{i=1}^m x_i^2}{\sigma^2}, \theta = \sigma^2$$

است. همچنین $\frac{\lambda c'(\sigma)b(\sigma)}{b'(\sigma)} = m$ بنابراین طبق قضیه (۵) برآورد مینیماکس مجاز پارامتر $\theta = \sigma^2$ با شرط $\sigma^2 \geq a$ عبارت است از

$$\delta(x) = \frac{\sum_{i=1}^m x_i^2}{m+2} \left\{ 1 + \frac{(\frac{\sum_{i=1}^m x_i^2}{\lambda})^{\frac{m}{2}} + 1}{\int_0^{\frac{1}{\lambda}} t^{\frac{m}{2}} + 1 e^{-t} dt} e^{-\frac{1}{\lambda} \sum_{i=1}^m x_i^2} \right\}$$

مثال ۳ : فرض کنید X_1, X_2, \dots, X_m یک نمونه تصادفی از $(\lambda, IG(\infty, \lambda))$ باشند. در این صورت

$$f_\lambda(x) = (\prod_{i=1}^m 2\pi x_i^2)^{-\frac{1}{2}} \lambda^{\frac{m}{2}} e^{-\frac{\lambda}{2} \sum_{i=1}^m \frac{1}{x_i}}$$

که متعلق به خانواده کای - دو تبدیل یافته با

$$k = m, a(x) = \frac{1}{2} \sum_{i=1}^m \frac{1}{x_i}, b(\lambda) = -\lambda,$$

$$c(\lambda) = \frac{m}{2} \ln \lambda, -2a(x)b(\lambda) = \lambda \sum_{i=1}^m \frac{1}{x_i}, \theta = \frac{1}{\lambda}$$

است. همچنین $m = \frac{2c'(\lambda)b(\lambda)}{b'(\lambda)}$ بنا بر این طبق قضیه (۵) برآورد مینیماکس مجاز پارامتر $\theta = \frac{1}{\lambda}$ با شرط $a \geq \frac{1}{\lambda}$ عبارت است از

$$\delta(x) = \frac{\sum_{i=1}^m \frac{1}{x_i}}{m+1} \left\{ 1 + \frac{\left(\sum_{i=1}^m \frac{1}{x_i} \right)^{\frac{m}{m+1}} e^{-\frac{1}{\lambda}} \sum_{i=1}^m \frac{1}{x_i}}{\int_0^{-\frac{1}{\lambda}} \sum_{i=1}^m \frac{1}{x_i} t^{\frac{m}{m+1}-1} e^{-t} dt} \right\}$$

در جدول (۱) بعضی از توزیعهای متعلق به خانواده توزیعهای کای - دو تبدیل یافته به همراه $a(x)$ و $k = \frac{1}{b(\eta)}$ آورده شده است. در این توزیعها برآورده (۱۶) برای θ مینیماکس و مجاز است.

مراجع

- Berry J. C. (1993). *Minimax Estimation of a Restricted Exponential location Parameter*. Statist. Decisions, 11, 307-316.
- Brunk, H. D. (1955). *Maximum Likelihood Estimates of Monotone Parameters*. Ann. Math. Statist., 26, 607-616.
- Katz, M.W. (1961). *Admissible and Minimax Estimates of Parameter in Truncated Spaces*. Ann. Math. Statist., 32, 136-142.
- Rahman, M.S and Gupta, R.P. (1993) *Family of Transformed Chi-square Distributions*. Commun. Statist. Theory Meth., 22 (1), 135-146.
- van Eeden, C. (1957). *Maximum Likelihood Estimation of Partially or Completely Ordered Parameters*. Proc. Kon. Ned. Akad. V.Wet. 60A, 128-136 and 201-211.
- van Eeden, C. (1995). *Minimax Estimation of a Lower Bounded Scale-parameter of a Gamma Distribution for Scale-invariant Squared-error Loss*. The Canadian Journal of Statistics, 23 (3), 245-256.
- van Eeden, C. (2000). *Minimax Estimation of a Lower-bounded Scale-parameter of an F-distribution*. Statistics & Probability Letters ,46 , 283-286.

Name of distribution with p.d.f	$a(X)$	$b(\cdot)$	$c(\cdot)$	$-\nabla a(X)b(\cdot)$	k	θ
<i>Normal</i> $\frac{e^{-\frac{x}{\sigma^2}}}{\sigma \sqrt{2\pi}}$	$\frac{1}{\sqrt{m}} \sum_{i=1}^m X_i^2$	$-\frac{1}{\sigma^2}$	$-m \ln \sigma$	$\frac{\sum_{i=1}^m X_i^2}{\sigma^2}$	m	σ^2
<i>Lognormal</i> $\frac{e^{-\frac{(Ln x)^2}{\sigma^2}}}{\sigma x \sqrt{2\pi}}$	$\frac{1}{\sqrt{m}} \sum_{i=1}^m (Ln X_i)^2$	$-\frac{1}{\sigma^2}$	$-m \ln \sigma$	$\frac{\sum_{i=1}^m (Ln X_i)^2}{\sigma^2}$	m	σ^2
<i>Exponential</i> $\frac{1}{\lambda} e^{-\frac{x}{\lambda}}$	$\sum_{i=1}^m X_i$	$-\frac{1}{\lambda}$	$-m \ln \lambda$	$\frac{\sum_{i=1}^m X_i}{\lambda}$	γm	λ
<i>Gamma</i> $x^{\alpha-1} \frac{e^{-\frac{x}{\beta}}}{\beta^\alpha \Gamma(\alpha)}$ with known α	$\sum_{i=1}^m X_i$	$-\frac{1}{\beta}$	$-m\alpha \ln \beta$	$\frac{\sum_{i=1}^m X_i}{\beta}$	γm	β
<i>Rayleigh</i> $x \frac{e^{-\frac{x^2}{\beta^2}}}{\beta^2}$	$\frac{1}{\sqrt{m}} \sum_{i=1}^m X_i^2$	$-\frac{1}{\beta^2}$	$-\gamma m \ln \beta$	$\frac{\sum_{i=1}^m X_i^2}{\beta^2}$	γm	β^2
<i>Pareto</i> $\alpha x^{-(\alpha+1)}$	$\sum_{i=1}^m \ln X_i$	$-\alpha$	$m \ln \alpha$	$\gamma \alpha \sum_{i=1}^m \ln X_i$	γm	$\frac{1}{\alpha}$
<i>Weibull</i> $\beta \alpha x^{\alpha-1} e^{-\beta x^\alpha}$ with known α	$\sum_{i=1}^m X_i^\alpha$	$-\beta$	$m \ln \beta$	$\gamma \beta \sum_{i=1}^m X_i^\alpha$	γm	$\frac{1}{\beta}$
<i>Maxwell</i> $\sqrt{\frac{2}{\pi}} \beta^{\frac{1}{2}} x^{\frac{1}{2}} e^{-\frac{\beta}{2} x^2}$	$\frac{1}{\sqrt{m}} \sum_{i=1}^m X_i^2$	$-\beta$	$\frac{\gamma m}{\sqrt{m}} \ln \beta$	$\beta \sum_{i=1}^m X_i^2$	γm	$\frac{1}{\beta}$
<i>Inverse Gaussian</i> $\sqrt{\frac{\lambda}{2\pi x^3}} e^{-\frac{\lambda(x-\mu)^2}{2\mu^2 x}}$ with known μ	$\frac{1}{\sqrt{m}} \sum_{i=1}^m \frac{(X_i - \mu)^2}{\mu^2 X_i}$	$-\lambda$	$m \frac{\ln \lambda}{\sqrt{m}}$	$\frac{1}{\sqrt{m}} \sum_{i=1}^m \frac{\lambda(X_i - \mu)^2}{\mu^2 X_i}$	m	$\frac{1}{\lambda}$

بعضی از توزیعهای متعلق به خانواده کای دو تبدیل یافته: جدول

مقدمه‌ای بر افشاری داده‌ها

دکتر حمیدرضا نواب پور، محمد بردبار

P ۱۴۲۷۲

دانشگاه علامه طباطبائی

چکیده: امروزه سازمانهای آماری در کشورهای مختلف اطلاعات بسیار زیادی در مورد جامعه خود جمع‌آوری می‌نمایند. این اطلاعات مواردی همچون جرائم، وضعیت مالی افراد، میزان درآمد شرکتها، موضوعاتی مثل وضعیت استخدامی و وضعیت بهداشت جامعه و ... را شامل می‌شوند. داده‌های جمع‌آوری شده جمع‌بندی و تحلیل می‌گردند و نتایج آن به دو صورت جدول یا داده‌های خرد^۱ منتشر می‌شوند تا کسانیکه برای تحقیق و برنامه‌ریزی به این اطلاعات نیاز دارند، به آنها دسترسی پیدا کنند. در تمامی قوانین عمومی آمار، اطلاعات فردی محظمانه تلقی می‌گردد زیرا که ممکن است در صورت انتشار به دلایل مختلف مورد سوءاستفاده قرار گیرند. این عمل نتایج زیان‌باری را برای پاسخ‌دهندگان و در نهایت برای جامعه در پی خواهد داشت. به همین دلیل سازمانهای آماری سعی می‌کنند که از راههای گوناگون جنبهٔ محظمانه بودن اطلاعات دریافتی از پاسخ‌دهندگان خود را حفظ نمایند و تا حد ممکن خطر افشاری اطلاعات فردی را کاهش دهند. در این مقاله سعی بر آن است تا مفهوم کلی افشاری داده‌ها، روش‌های گوناگون ارزیابی خطر آن و راههای کاهش این خطر که امروزه به موضوع بحث‌انگیز و مهمی تبدیل شده است مورد بحث قرار گیرد.

واژه‌های کلیدی: افشاء^۲، داده‌های جدولی^۳، داده‌های خرد.

۱ مقدمه

محظمانه بودن اطلاعات شخصی و محدودسازی افشاری آنها موضوعاتی ذاتاً آماری می‌باشد که تا این اواخر توجه تعداد محدودی از محققان آماری را به خود جلب کرده‌بودند. این وضعیت در دههٔ اخیر تغییر یافته است. قبل این مباحثت به بخش غیر آماری تنزل داده می‌شدن و در نتیجه شیوه‌هایی که برای پرداختن به آنها مورد استفاده قرار می‌گرفت قدیمی و غیرآماری بودند. بحث روی این موضوعات با یک مقاله توسط دالنیوس در سال ۱۹۷۷

^۱ Microdata

^۲ Disclosure

^۳ Tabular Data

و یک گزارش مشروح از زیر کمیته روش شناختی محدودسازی افشاء اطلاعات به سال ۱۹۷۸ در آمریکا بصورت جدی شروع شد و آمارشناسان به آهستگی به تحقیق در این موضوع پرداختند. با اینکه تا به حال تحقیقات زیادی روی این موضوعات صورت گرفته است اما هنوز هم باید برای رسیدن به سطح مطلوب تلاش کرد. ابتدا برخی مفاهیم مربوط به ادبیات افشاء داده‌ها مطرح می‌شود و در ادامه چند روش محدودسازی افشاء را مورد بررسی قرار می‌دهیم.

۲ تعاریف

۱.۲ حفظ محramانه بودن

بسیاری از داده‌هایی که توسط سازمانهای آماری جمع‌آوری می‌گردند یا مستقیماً از پاسخ‌دهندگان بدست می‌آیند که می‌تواند به صورت مصاحبهٔ حضوری، مصاحبهٔ تلفنی، مصاحبهٔ پستی، مصاحبهٔ اینترنتی و... باشد و یا اینکه برای کاهش زحمت پاسخگویان و مینمم کردن هزینه‌های جمع‌آوری داده‌ها از رکوردهای اداری (ثبتی) استفاده شود. در هر صورت، داده‌های جمع‌آوری شده به جنبه‌هایی از زندگی شخصی یا تجاری پاسخ‌دهندگان که به نظرشان این جنبه‌ها محramانه تلقی می‌شود، مربوط می‌شوند. حفاظت از جنبهٔ محramانه بودن داده‌ها یک مسئولیت مشخص و انکار ناپذیر سازمانهای آماری می‌باشد. این مسئولیت به ۳ دلیل زیراست:

الف- حفاظت از جنبهٔ محramانه بودن داده‌ها یک رفتار اخلاقی برای آمارشناسان در حرفهٔ آمار می‌باشد که به تصویب بین‌المللی رسیده است.^۱

ب- حفاظت از جنبهٔ محramانه بودن داده‌ها اغلب مورد توجه قوانین دولتی می‌باشد.

ج- حقوقدانان معتقدند که اغلب پاسخگویان، بدون تضمین اینکه پاسخ‌هایشان برای یک نهاد سوم فاش نمی‌شوند، حاضر به پاسخگویی نیستند و یا بدرستی پاسخ نمی‌دهند. اگر برای حفظ محramانه بودن داده‌های جمع‌آوری شده تضمینی ارائه نشود، اعتماد عمومی پاسخ‌دهندگان به سازمانهای آماری کاهش خواهد یافت و این امر باعث کاهش میزان پاسخگویی در آمارگیریها می‌شود که نتیجه آن افت کیفی سرشماریها و تحقیقات بعدی خواهد بود.

۲.۲ افشاء

بحث افشاء مربوط به برداشت نابجا (نامناسب) از اطلاعات موجود در داده‌های مربوط به یک شخص یا سازمان می‌باشد. دانکن در سال ۱۹۹۳ تعریفی از افشاء ارائه کرد که

^۱ International Statistical Institute 1986 and American Statistical Association 1989

نوع از افشاء را دربر دارد.

افشاء هنگامی اتفاق می‌افتد که یک موضوع داده‌ای از یک فایل منتشر شده داده‌ها مشخص شود (افشاء هویت^۱)، یا اینکه اطلاعات حساس در مورد یک موضوع داده‌ای توسط یک فایل داده‌های منتشرشده آشکار گردند (افشاء صفت^۲)، یا اینکه داده‌های منتشرشده این امکان را بوجود آورد که مقدار بعضی مشخصه‌های یک فرد یا واحد آماری با دقیقی بیشتر نسبت به دیگر حالات ممکن، تعیین شود (افشاء استباطی^۳). افشاء استباطی ممکن است شامل افشاء هویت یا افشاء صفت نیز باشد. انتشار داده‌های آماری قطعاً بعضی اطلاعات را در مورد موضوعات شخصی آشکار می‌کند. افشاء هنگامی اتفاق می‌افتد که اطلاعاتی که جنبهٔ محترمانه دارند آشکار شوند. گاهی اوقات افشاء تنها بر پایهٔ داده‌های منتشرشده اتفاق می‌افتد. گاهی اوقات افشاء، تبیجهٔ ترکیب داده‌های منتشرشده با اطلاعات عمومی در دسترس می‌باشد و گاهی اوقات افشاء تبیجهٔ ترکیب داده‌های منتشرشده با منابع دادهٔ بیرونی است که می‌تواند در دسترس عموم باشند.

۳.۲ دسترسی محدود و داده‌های محدود

برای برخورد با بحث افشاء داده‌ها دو فلسفه وجود دارد: یکی محدود کردن امکان دسترسی به داده‌ها و دیگری محدود کردن مقدار و شکل داده‌های منتشرشده است. کمبود اطلاعات در یک قلمرو معین می‌تواند فاجعه‌آفرین باشد. با ورود به عصر اطلاعات تمایل برای بدست آوردن هر چه بیشتر اطلاعات بطور سراسم آوری افزایش یافته است. با پیشرفت فوق العاده وسائل ارتباط جمعی در آینده نیز این تمایل باشد بیشتری ادامه خواهد یافت. به سبب اینکه داده‌های آماری که با هزینهٔ دولت جمع آوری می‌شوند، مایلک عمومی محسوب می‌گردند، عدم انتشار اطلاعات با موازین یک جامعه آزاد مغایرت دارد و در کل منطق قابل قبولی برای پیشگیری از افشاء اطلاعات محسوب نمی‌شود. محدود کردن دسترسی به اطلاعات موجود در داده‌ها تنها در صورتی باید موجه باشد که از راههای دیگر نتوان جنبهٔ محترمانه بودن اطلاعات را حفظ کرد. معمولاً سازمانهای آماری از گزارش مطالب کلیدی که موجب شناخت دقیق افراد می‌شود (شناسه‌های مستقیم) اجتناب می‌کنند. بعلاوه بسیاری از سازمانهای آماری به منظور کاهش خطر افشاء اطلاعات تنها نمونه‌هایی از اطلاعات سرشماری و زیرنمونه‌هایی از نمونه‌گیریهای گستردگی را گزارش می‌دهند. با اتخاذ چنین سیاستهایی و درکنار یک سری اقدامات قانونی و نیز اقداماتی در جهت افزایش امنیت سازمانی و کامپیوتري، شمار روزافزونی از سازمانهای آماری به خط مشی های

^۱ Identity Disclosure

^۲ Attribute Disclosure

^۳ Inferential Disclosure

روی می آورند که نیازمند تعدیل داده های آماری است. این اقدامات موضوع محدودسازی افشاء اطلاعات آماری می باشد. پس از یک سو باید حداقل اطلاعات در اختیار جامعه قرار گیرد و از سوی دیگر باید این تضمین داده شود که حریم شخصی افراد و سایر واحد های آماری در حد معقولی حفظ می شود. پس باید یک توازن بین این دو مسئله برقرار شود. زیر مجموعه ای از آمار که در این زمینه فعالیت دارد تحت عنوان مختلفی همچون کنترل افشاء آمار، جلوگیری از افشاء آمار، حفظ اطلاعات آماری و محرومانه بودن آمار عنوان می گردد.

۴.۲ جدولها و داده های خرد

انتخاب روش های محدودسازی افشاء آماری به طبیعت محصولات داده ای که باید جنبه محرومانه بودن آنها حفظ شود بستگی دارد. بیشتر داده های آماری به دو شکل جدولها یا فایلهای داده خرد منتشر می شوند. جدولها شامل دو گروه هستند: جدولهای فراوانی و جدولهای بزرگی^۱. برای هر گروه داده ها می توانند به شکل تعداد، نسبت یا درصد بیان شوند.

یک فایل داده های خرد شامل رکوردهای شخصی است که هر مقدار آن شامل متغیرهای یک شخص ، یک مؤسسه تجاری یا یک واحد دیگر می باشد. بعضی از فایلهای داده ها شامل مشخص کننده های صریحی مثل نام، آدرس یا شماره کد ملی می باشد. حذف کردن هر کدام از این مشخص کننده ها از فایل ، یک اقدام اولیه بدیهی برای مهیا کردن زمینه انتشار یک فایل داده ها است که بایستی جنبه محرومانه بودن اطلاعات شخصی آن حفظ شود.

۵.۲ پوشاندن ماتریسی

یک ردی کلی از روش های محدودسازی افشاء با نام پوشاندن ماتریسی^۲ شناخته می شود که خلاصه آن در زیر می آید.

فرض کنید ماتریس Z یک ماتریس $n \times p$ باشد که داده های مربوط به n شخص یا حالت را برای تعداد p متغیر یا صفت ارائه کرده است. پوشاندن ماتریسی بدین صورت عمل می کند: فرض کنید بتوان ماتریس Z را بصورت $Z = A Z B + C$ به ماتریس M تبدیل کرد که در آن A ماتریسی است که حالات را تبدیل می کند و B ماتریس تبدیل صفات یا متغیرها می باشد و C ماتریسی است که اعضای Z را مبهم می سازد. پوشاندن ماتریسی

^۱ Tables of Magnitude

^۲ Matrix Masking

شامل راهبردهای استاندارد گوناگون گسترهای می‌باشد که چند راهبرد مشهور آن در ادامه آمده است :

الف – انتشار زیرمجموعه‌ای از مشاهدات (حذف سطرهایی از Z).

ب – شبیه‌سازی داده‌ها (افزودن سطرهایی به Z).

پ – افزودن پرشیدگیهای تصادفی^۱ به Z .

ت – انتشار زیرمجموعه‌ای از متغیرها (حذف ستونهایی از Z).

ج – انتشار ماتریس واریانس – کوواریانس.

چ – مبادله مقادیر ستونی انتخاب شده‌ای با جفت‌هایی از سطرها.

ح – پنهان‌سازی خانه‌ای^۲ برای طبقه‌بندی‌های متقاطع از داده‌ها.

اگر چه در اینجا حالت دو بعدی داده‌ها را در نظر گرفته‌ایم و آن را بصورت یک ماتریس $n \times p$ نمایش دادیم ولی حقیقت این است که پوشاندن ماتریسی قابل تعمیم به ابعاد بالاتر نیز است.

۳ روش‌های محدودسازی افشاء

۱.۳ جدولهای فراوانی

۱.۱.۳ نمونه‌گیری به عنوان یک روش محدودسازی افشاء آماری

خطر افشاء در مورد داده‌های حاصل از یک بررسی مبتنی بر نمونه‌گیری تقریباً کمتر از خطر افشاء در مورد داده‌های حاصل از یک سرشماری می‌باشد. هر چه کسر نمونه‌گیری کوچکتر باشد احتمال کمتری وجود دارد که شخصی بتواند با استفاده از اطلاعات موجود در داده‌های نمونه به اطلاعات محروم‌های افراد یا سازمانهایی در جامعه دست پیدا کند. البته این موضوع ارتباط نزدیکی با طرح نمونه‌گیری و توزیع مشخصه مورد نظر دارد.

هرگاه جدولهای فراوانی را مستقیماً با استفاده از داده‌های سرشماری بسازیم، بایستی روش‌های محدودسازی افشاء را به کار ببریم. دو رده از قوانین محدودسازی افشاء برای جدولهای فراوانی وجود دارد. اولین رده شامل قواعد خاصی برای جدولهای معینی می‌باشد. دومین رده کلی تراست: اگر تعداد پاسخگویان در یک خانه جدول از تعداد مشخص شده‌ای

^۱ Random Perturbations

^۲ Cell Suppression

کمتر باشد، آن خانه، حساس تعريف میشود (قاعده آستانه‌ای)^۲.

۲.۱.۳ قواعد خاص

قواعد خاص روی سطح جزئیاتی که می‌تواند در یک جدول فراهم شود، محدودیتهای وضع می‌کند. به عنوان مثال، انتشار جدولهایی که مقدار یک خانه داخلی آنها با یک جمع حاشیه‌ای برابر است را ممنوع می‌کنند. این قواعد کلاً برای فراهم کردن داده‌هایی که از طرف سازمان آماری بصورت خاصی، حساس تلقی می‌شود مطرح است و از سازمانی به سازمان دیگر متفاوت می‌باشد.

۳.۱.۳ قاعده آستانه‌ای

با استفاده از قاعده آستانه‌ای اگر تعداد پاسخگویان در خانه جدولی از اعداد مشخص شده‌ای کمتر باشد، آن خانه را حساس می‌نامیم. برخی سازمانها خواستار وجود حداقل ۵ پاسخگو در یک خانه می‌باشند ویرخی ۳ پاسخگو و یک سازمان می‌تواند جدولهایی بسازد و گروهها را مانند حالت قبل ترکیب کند یا اینکه از فنون پنهان‌سازی خانه‌ای، گردکردن تصادفی، گردکردن کترل شده و ویرایش محروم‌بودن استفاده کند. جدول ۱ یک جدول ۱ پاسخگو با افشاء می‌باشد که مربوط به نوجوانان بزرگوار است. خانه‌های دارای کمتر از ۵ پاسخگو را حساس تعريف کرده‌ایم. به عنوان مثال با مراجعت به این جدول مشخص می‌شود که در ناحیه A تنها یک نوجوان بزرگوار وجود دارد که میزان تحصیلات سرپرستش خیلی بالا است. با اطلاعات عمومی که ممکن است در مورد تحصیلات سرپرستان خانوارها در ناحیه A در دسترس باشد، این نوجوان به دقت مشخص می‌شود.

الف – پنهان‌سازی^۱

یکی از عمومی‌ترین راه‌های مورد استفاده برای حفاظت خانه‌های حساس روش پنهان‌سازی است. واضح است که در یک سطر یا ستون با یک خانه حساس پنهان شده، حداقل یک خانه اضافی باید پنهان گردد، در غیر این صورت مقدار خانه حساس می‌تواند از طریق جمع حاشیه‌ای متناظر محاسبه گردد. بدین دلیل خانه‌های معین دیگری باید پنهان شوند. این عمل را پنهان‌سازی مکمل گویند. تا هنگامیکه انتخاب خانه‌ها برای پنهان‌سازی مکمل بصورت دستی انجام گیرد، تضمین اینکه نتیجه حاصله، حفاظت کافی را انجام دهد مشکل است. جدول ۲ مثالی از یک سیستم خانه‌های پنهان شده را برای داده‌های جدول

^۲ Threshold Rule

^۱ Suppression

۱ ارائه می‌دهد که در هر سطر و ستون حداقل ۲ خانه پنهان شده دارد. به نظر می‌رسد که این جدول خانه‌های حساس را حفظ می‌کند ولی عملاً این‌گونه نیست. به عنوان مثال با استفاده از رابطه:

$$(15 + D1 + D2 + D3) + (20 + D4 + D5 + 15) - (D1 + D4 + 10 + 14) \\ - (D2 + D5 + 10 + 17) = (20 + 55 - 35 - 30)$$

نتیجه می‌شود: $D3 = 1$. یعنی عملاً می‌توان مقدار $D3$ را با استفاده از جمعهای حاشیه‌ای تعیین کرد. پس هنوز مسأله افشاء وجود دارد.

این مثال نشان می‌دهد که انتخاب خانه‌ها برای پنهان‌سازی مکمل پیچیده‌تر از آن چیزی است که در ابتدا تصور می‌شد. روش‌های ریاضی برنامه‌ریزی خطی برای انتخاب خودکار خانه‌ها به منظور پنهان‌سازی مکمل و همچنین الگوی پنهان‌سازی مفروض به کار می‌روند. جدول ۳ نشان‌دهنده یک سیستم خانه‌های پنهان است که حفاظت کافی را برای خانه‌های حساس فراهم می‌کند. به هر حال جدول ۳ یکی از مشکلات پنهان‌سازی را نشان می‌دهد. از ۱۶ خانه اولیه تنها ۷ خانه منتشر شده‌اند و ۹ خانه دیگر پنهان شده‌اند.

ب - گردکردن تصادفی^۱

به منظور کاهش تعداد داده‌هایی که در صورت استفاده از روش پنهان‌سازی از دست می‌دهیم روش‌های دیگری برای حفاظت خانه‌های حساس در جدولهای فراوانی به کار می‌روند که به روش‌های پرشیدگی^۲ موسومند. روش گردکردن تصادفی و گردکردن کنترل شده مثالهایی از روش‌های پرشیدگی می‌باشند. در روش گردکردن تصادفی، مقادیر خانه‌ای گرد می‌شوند اما به جای استفاده از قراردادهای استاندارد گردکردن، یک تصمیم تصادفی برای گردکردن به بالا یا پایین اتخاذ می‌گردد. عمل گردکردن برای هر خانه در یک جدول بصورت جداگانه انجام می‌شود ولی جمعهای حاشیه‌ای بدون تغییر باقی می‌مانند. یک جدول که عمل گردکردن تصادفی در آن انجام گرفته است میتواند منجر به از دست دادن اطمینان عمومی به اعداد گردد.

ج - گردکردن کنترل شده^۳

^۱ Random Rounding

^۲ Perturbation Methods

^۳ Controlled Rounding

برای حل مشکل روش گردکردن تصادفی، روش گردکردن کنترل شده توسعه یافته است. این روش همانند روش گردکردن تصادفی است اما تاکید بر این است که مجموع مقادیر منتشرشده در هرسطر و ستون برابر جمعهای حاشیه‌ای مقناظر باشد. این روش در سالهای ۱۹۷۰ تا ۱۹۸۰ مورد تحقیق قرار گرفت و از آن در سرشماری سال ۱۹۹۰ آمریکا استفاده شد. یک مشکل این روش این است که به یک برنامه کامپیوتری مخصوص نیاز دارد که در حال حاضر این برنامه‌ها به گستردگی در دسترس نیست. مشکل دیگر آن این است که پاسخهای گردکردن کنترل شده برای جدولهای پیچیده ممکن است همواره وجود نداشته باشد.

د- ویرایش محرمانه بودن^۱

در این روش دورهای متفاوت وجود دارد. یکی برای داده‌های سرشماری استفاده می‌شود و دیگری برای داده‌های نمونه‌ای به کار می‌رود. هر دو فن، روشهای محدودسازی افشاء را برای فایلهای داده‌های خرد، قبل از اینکه این فایلهای برای تنظیم جدولها به کار روند استفاده می‌کنند و خود این فایلهای منتشر نمی‌شوند.

روش اول: برای فایل داده‌های خرد مربوط به سرشماریها ویرایش محرمانه بودن شامل تبادل داده‌ها یا معاوضه می‌باشد. مراحل ویرایش محرمانه بودن در زیر آمده است:

۱- گرفتن یک نمونه از رکوردها از فایل داده‌های خرد.

۲- پیدا کردن یک همانند برای این رکوردها در بعضی نواحی جغرافیایی دیگر (همانندیابی روی یک مجموعه مشخص از صفات مهم).

۳- تعویض تمامی صفات روی رکوردهای همانند شده.

اداره سرشماری برای حفاظت بیشتر از بلوکهای کوچک، کسر نمونه‌گیری را افزایش می‌دهد. پس از آماده شدن فایل، از آن برای آماده کردن جدولها استفاده می‌شود.

روش دوم: فایل داده‌های نمونه‌ای شامل داده‌های یک نمونه از جامعه می‌باشد و همانطور که قبل ذکر کردیم، انجام نمونه‌گیری تقریباً حفاظت افشاء را فراهم می‌کند. مطالعات نشان می‌دهد که این حفاظت به جز در نواحی کوچک جغرافیایی، کافی است. برای توضیح ویرایش محرمانه بودن به کار رفته در فایل داده‌های خرد سرشماری، از رکوردهای ساختگی برای ۲۰ شخص در ناحیه A که در جدولهای قبلی به کار رفته‌اند استفاده

^۱ Confidentiality Edit

می‌کنیم. جدول ۴، پنج متغیر را برای این اشخاص نشان می‌دهد.

برای به کار بردن روش ویرایش محترمانه بودن با استی مراحل زیر را طی کنیم :

- ۱— گرفتن یک نمونه از رکوردها از فایل داده‌های خرد (مثلاً یک نمونه ۱۰ درصدی).
- ۲— وقتی به جدولهایی با تفکیک ناحیه و سطح تحصیلات نیاز داریم، همانندی در ناحیه‌ای دیگر روی متغیرهای دیگری مثل رنگ پوست، جنسیت و درآمد پیدا می‌کنیم. (در نتیجه جمعهای حاشیه‌ای برای این متغیرها با عمل تبادل بدون تغییر خواهد ماند). یک همانند برای رکورد شماره ۴ (پیت) در ناحیه B پیدا شد. شخص همانند، آنفسو می‌باشد که سرپرست تحصیلات بالایی دارد. رکورد شماره ۱۷ (مايك) با جرج از ناحیه D همانند شد که سرپرست دارای تحصیلات متوسط است. همچنین بخشی از نمونه ۱۰ درصدی انتخاب شده از نواحی دیگر دارای رکوردهای نظری در ناحیه A می‌باشند. یک رکورد از ناحیه D جون با سرپرست دارای تحصیلات بالا با رکورد ویرجینیا همانند شد (شماره ۱۲). یک رکورد از ناحیه C (هیتر با سرپرست دارای تحصیلات پایین (با رکورد نانسی همانند شد (شماره ۲۰).
- ۳— پس از اینکه تمام همانندها ساخته شدند، صفات آنها روی رکوردهای همانند شده مبادله می‌شوند. فایل داده‌خرد مرتب شده در جدول ۵ آمده است.
- ۴— استفاده مستقیم از فایل داده‌های مبادله شده برای تولید جدولها (جدول ۶).

اگر چه برخی از خانه‌های جدول دارای تعداد پاسخگوی کمتر از ۵ می‌باشد اما چون مبادله صورت گرفته است، نمی‌توان رکورد یک شخص را در جامعه بطور دقیق تعیین کرد.

روش ویرایش محترمانه بودن این مزیت را دارد که می‌تواند به آسانی برای جدولهای چندبعدی آماده شود و همیشه می‌تواند از افشاری محترمانه بودن جلوگیری کند.

۲.۳ جدولهای بزرگی

جدولهای داده‌های بزرگی یک مجموعه منحصر به فرد از مسائل افشاء را دارند. داده‌های بزرگی عموماً کمیتهايی نامنفی می‌باشند که در تحقیقات یا سرشماریها از مؤسسه‌های تجاری، مزارع یا سازمانها گزارش می‌شوند. این داده‌ها ممکن است بر درآمد اشخاص، میزان فروش یا میزان درآمد مؤسسات دلالت کنند. توزیع این مقادیر معمولاً چوله می‌باشد که در آن تعداد کمی از اعضاء دارای مقادیر خیلی بزرگی هستند. محدودسازی افشاء در این جدولها برای اطمینان دادن اینکه داده‌های منتشرشده برای برآورد مقادیر گزارش شده

توسط بزرگترین واحد استفاده نخواهد شد، متوجه شده است. در این حالت امکان اینکه نمونه‌گیری به تنها یی بتواند حفاظت افساء را فراهم کند کمتر است، زیرا واحدهایی که به علت اندازه‌شان نمایانی بالایی دارند بوسیله نمونه‌گیری نمی‌توانند حفاظت شوند. قواعد پنهان‌سازی مقدماتی یا اندازه‌های حساس خطی برای تعیین اینکه آیا یک خانه جدولی داده شده می‌تواند اطلاعات پاسخ‌دهنده شخصی را نمایان کند، توسعه داده شده‌اند. چنین خانه‌ای، خانه حساس نامیده می‌شود و نمی‌تواند منتشر شود. قواعد پنهان‌سازی مقدماتی که عموماً برای تشخیص خانه‌های حساس استفاده می‌شوند، قاعده (n, k) ، قاعده p درصد و قاعده pq می‌باشد. همه این قواعد بر این پایه استوار هستند که امکان برآوردن دقیق مقدار گزارش شده توسط یک پاسخگو را برای شخصی دیگر مشکل نمایند. بزرگترین مقدار گزارش شده دارای بیشترین احتمال برای برآوردن دقیق می‌باشد. قواعد پنهان‌سازی مقدماتی برای داده‌های فراوانی نیز قابل استفاده می‌باشند. هرگاه خانه‌های حساس مشخص شوند تنها دو گزینه وجود خواهد داشت: بازسازی جدول و کاهش خانه‌ها تا اینکه خانه حساسی در جدول باقی نماند یا پنهان‌سازی خانه‌ای.

هرگاه خانه‌های حساس مشخص شوند آنها را از انتشار کنار می‌گذاریم. این خانه‌ها پنهان‌سازی‌های مقدماتی نامیده می‌شوند. خانه‌های دیگر که پنهان‌سازی‌های مکمل نامیده می‌شوند، انتخاب می‌گردند و پنهان می‌شوند تا مقدار خانه‌های حساس با توجه به جمعهای حاشیه‌ای منتشر شده قابل تعیین نباشند. مشکلات مربوط به پنهان‌سازی خانه‌ای در جدولهای بزرگی همانند جدولهای فراوانی می‌باشد. یک راه اجتناب از مشکلات مربوط به بکارگیری پنهان‌سازی خانه‌ای تعدادی از سازمانهای آماری بکار گرفته شده است که عبارت است از گرفتن اجازه کتبی از پاسخگویان برای انتشار یک خانه حساس.

۳.۳ داده‌های خرد

اطلاعات جمع‌آوری شده درباره مؤسسات در ابتدا بصورت داده‌های بزرگی هستند. این داده‌ها به شدت چوله می‌باشند و پاسخ‌گویانی که چشم آمد هستند، بواسطه اطلاعات دیگری که در دسترس عمومی می‌باشند، می‌توانند به آسانی مشخص شوند. یک نتیجه این می‌تواند باشد که فایلهای داده‌های خرد برای استفاده عموم منتشر نشوند. حفاظت یک فایل داده‌های خرد به علت وجود دیگر منابع داده‌ای بسیار مشکل می‌باشد. بعلاوه اندازه قابل قبولی از خطر افشاء برای یک فایل داده‌های خرد وجود ندارد. شیوه‌های حفاظت از فایلهای داده‌های خرد در ادامه توصیف شده‌اند. این شیوه‌ها توسط تمام سازمانهایی که فایلهای داده‌های خرد را برای استفاده عموم فراهم می‌کنند به کار می‌روند. به منظور کاهش خطر افشاء، تمام فایلهای داده‌های خرد مورد استفاده عمومی باید شرایط زیر را دارا باشند:

۱- باید تنها شامل داده‌های یک نمونه از جامعه باشند.

۲- شامل متغیرهای چشم آمد^۱ نباشند.

۳- جزئیات جغرافیایی محدود شده باشند.

۴- تعداد متغیرهای موجود در فایل محدود شده باشند.

روشهای اضافه‌ای که برای متغیرهای چشم آمد استفاده می‌شوند عبارتند از:

۱- کدبندی به بالا یا کدبندی به پایین.

۲- بازکدبندی^۲ به بازه‌ها یا گردکردن.

۳- افزودن یا ضرب کردن اعداد تصادفی (نوفه^۳).

۴- مبادله یا مبادله رتبه‌ای.

۵- انتخاب رکوردهایی به تصادف، حذف متغیرهای انتخاب شده و جانه‌ی برای آنها.

۶- تجمع^۴ در طول گروههای کوچک پاسخ‌گویان و جایگذاری یک مقدار گزارش شده شخصی با مقدار متوسط پاسخهای گروههای انتخاب شده.

این روشها را بایک مثال ساختنگی که در قسمتهای قبل استفاده کردیم توضیح می‌دهیم

^۱ High Visibility Variables

^۲ Recoding

^۳ Noise

^۴ Aggregation

۱.۳.۳ نمونه‌گیری، حذف مشخص‌کننده‌ها و محدودسازی جزئیات جغرافیایی

الف – تنها داده‌های یک نمونه از جامعه را در نظر می‌گیریم. برای مثال از یک نمونه ۱۰ درصدی از جامعه نوجوانان بزهکار استفاده نموده‌ایم. بخشی از جامعه آماری (ناحیه) A در جدول ۴ نشان داده شده است.

ب – حذف مشخص‌کننده‌های صریح. در این مثال مشخص‌کننده نام نوجوان است.

ج – محدود کردن جزئیات جغرافیایی. در این مثال تصمیم گرفتیم داده‌های ناحیه‌ای شخصی مربوط به یک ناحیه دارای کمتر از ۳۰ نوجوان بزهکار در جامعه را منتشر نکنیم. بنابراین نواحی A و C از جدول ۱ در جکیب شده و بصورت AC نشان داده شده است. در این مثال تنها ۵ متغیر برای هر نوجوان در نظر گرفته‌ایم. شاید به نظر برسد که این ۵ متغیر از مجموعه داده کاملتری شامل متغیرهای نام مادر، نام خواهر و برادر، تعداد خواهر و برادر، سن نوجوان، سن خواهر و برادر، آدرس مدرسه و غیره انتخاب شده‌اند. هر چه تعداد متغیرها در فایل داده‌های خرد مربوط به هر نوجوان بیشتر باشد، ترکیبات یکتای متغیرها احتمال تشخیص نوجوانی را با استفاده از دانش شخصی افزایش می‌دهد. محدود کردن تعداد متغیرها به ۵ متغیر احتمال چنین تشخیصی را کاهش خواهد داد.

۲.۳.۳ متغیرهای چشم آمد

ممکن است در جامعه اطلاعاتی موجود باشد که در دسترس دیگران است. مثلاً می‌توان با استفاده از داده‌های درآمد که در جدول ۷ آمده است، خانواده یک نوجوان مجرم را بصورت منحصر به فرد تشخیص داد. به عنوان مثال، کارفرمای سرپرست خانوار معمولاً حقوق دقیق او را می‌داند. چنین متغیرهایی، متغیرهای چشم آمد نامیده می‌شوند و احتیاج به محافظت اضافی دارند.

الف – کدبندی به بالا، کدبندی به پایین و بازکدبندی به بازه‌ها

مقادیر زیاد درآمد به بالا کدبندی می‌شوند، به این صورت که فقط نشان می‌دهیم که میزان درآمد بیشتر از ۱۰۰ هزار دلار در سال است. مقادیر کوچک درآمد به پایین کدبندی می‌شوند، بدین صورت که فقط مشخص می‌کنیم که درآمد کمتر از ۴۰ هزار دلار در سال می‌باشد و درآمدهای دیگر به بازه‌های ۱۰ هزار دلاری بازکدبندی می‌شوند. نتیجه این اقدامات فایلی می‌باشد که در جدول ۸ برای استفاده عموم آمده است. روش‌های کدبندی به بالا، کدبندی به پایین و بازکدبندی به بازه‌ها عمومی‌ترین روش‌های مورد استفاده برای فراهم کردن حفاظت متغیرهای چشم آمد در فایلهای داده‌های خرد می‌باشد.

ب – افزودن پرشیدگی تصادفی

یک روش دیگر برای مخفی کردن متغیرهای چشم آمد از قبیل درآمد، افزودن یک عدد تصادفی به مقادیر این متغیرها یا ضرب یک عدد تصادفی در مقادیر این متغیرها می‌باشد. به عنوان نمونه، در مثال قبل فرض کنید می‌خواهیم یک متغیر تصادفی از توزیع نرمال با میانگین صفر و انحراف استاندارد ۵ به درآمد اضافه کنیم. ممکن است این عمل در حین نمونه‌گیری، حذف مشخص کننده‌ها و محدود کردن جزئیات جغرافیایی انجام شود. نتیجه در جدول ۹ آمده است. برای تولید این جدول، ۱۴ عدد تصادفی از توزیع نرمال مشخص شده انتخاب گردید و به داده‌های درآمد در جدول ۷ اضافه شد.

ج - مبادله یا مبادله رتبه‌ای

مبادله کردن شامل انتخاب یک نمونه از رکوردها، یافتن یک همانند در بانک داده‌ها روی یک مجموعه از متغیرهای از پیش تعیین شده و مبادله کردن تمام متغیرها می‌باشد. مبادله کردن در بخش ویرایش محترمانه بودن برای جدولهای داده‌های اندازه مورد بحث قرار گرفت. در آن مثال رکوردها از نواحی مختلف مشخص بودند و روی متغیرهای رنگ پوست، جنسیت و درآمد همانندیابی شدند، آنگاه نام کوچک نوجوان و میزان تحصیلات سرپرست خانوار مبادله شدند. به منظور فراهم کردن محافظت اضافی برای متغیر درآمد در یک فایل داده‌های خرد، می‌توان روی متغیرهای رنگ پوست و تحصیلات سرپرست همانندیابی کنیم و آنگاه متغیر درآمد را مبادله نماییم. مبادله کردن رتبه‌ای یک راه برای استفاده از متغیرهای پیوسته برای تعریف جفت رکوردهای قابل مبادله ارائه می‌کند. به جای اینکه دقیقاً همانند متغیرها را بیابیم، آنها را بطور تقریبی و براساس یک لیست مرتب شده تحت متغیر پیوسته تعریف می‌کنیم. رکوردهایی که بر اساس این متغیر مرتب شده‌اند و دارای رتبه نزدیکی هستند را به عنوان جفت‌های قابل مبادله در نظر می‌گیریم. در این روش، متغیری که برای مرتب کردن در نظر می‌گیریم اغلب یکی از متغیرهایی است که مبادله خواهد شد.

د - حذف و جانهی برای رکوردهای انتخاب شده تصادفی

حذف و جانهی شامل انتخاب تعدادی از رکوردهای فایل داده‌های خرد، حذف متغیرهای انتخاب شده و جایگذاری آنها با مقادیر انتسابی می‌باشد. این فن با استفاده از داده‌های جدول ۷ توضیح داده می‌شود. ابتدا از هر ناحیه قابل انتشار یعنی AC، B و D یک رکورد بصورت تصادفی انتخاب می‌شود. در رکورد انتخاب شده مقدار درآمد با یک مقدار انتسابی جایگذاری می‌شود. اگر رکورد منتخب از ناحیه AC رکورد شماره ۲ و در ناحیه B رکورد شماره ۶ و در ناحیه D رکورد شماره ۱۳ باشد، مقدار انتسابی به این رکوردها ممکن است به ترتیب اعداد ۶۳، ۵۲ و ۴۹ باشد. این اعداد ساختگی هستند اما می‌توان دید که مقادیر انتسابی با متوسط میزان درآمد خانوارها در یک ناحیه با رنگ پوست و

سطح تحصیلات یکسان، برابر می‌باشد.

هـ – نامشخص کردن

در روش نامشخص کردن یک مقدار گزارش شده با یک مقدار متوسط جایگذاری می‌شود. راههای زیادی برای نامشخص کردن وجود دارد. گروههایی از رکوردها که برای متوسطگیری انتخاب می‌کنیم شاید با همانندیابی روی متغیرهای دیگر تعیین شوند و یا از طریق مرتب کردن داده‌ها روی متغیر مورد علاقه‌مان حاصل شوند. تعداد رکوردهای یک گروه که داده‌های آنها متوسطگیری می‌شود ممکن است ثابت یا تصادفی باشد. مقدار متوسط مربوط به یک گروه خاص ممکن است به همه اعضای گروه تخصیص داده شود و یا اینکه به عضو میانی گروه تخصیص یابد. به عنوان مثال این فن را برای نامشخص کردن داده‌های درآمد در مثال قبل تشریح می‌کنیم. در فایل داده‌های خرد کامل ممکن است روی متغیرهای مهمی همچون ناحیه، رنگ پوست و دو گروه تحصیلی (خیلی بالا و بالا) و (متوسط و پایین) همانندیابی را انجام دهیم. آنگاه نامشخص کردن می‌تواند متوسطگیری خانوارها در هر گروه باشد. با توجه به جدول ۴ شاید این گونه برداشت شود که باستثنی درآمد خانوار برای گروه شامل جان و سو با مقدار متوسط درآمدشان جایگذاری شود (۱۳۹)، درآمد خانوار برای گروه شامل جیم و پیت باید با مقدار متوسط درآمدشان جایگذاری شود (۸۲) و به همین صورت ادامه می‌دهیم. پس از نامشخص کردن، فایل داده‌ها می‌تواند به عنوان موضوعی برای نمونه‌گیری، حذف کردن مشخص‌کننده‌ها و محدودسازی جزئیات جغرافیایی به کار رود.

۴ خلاصه

در این مقاله سعی شد مفاهیم مربوط به ادبیات افشاءی داده‌ها مورد بررسی قرار گیرد و روش‌های استاندارد محدودسازی افشاء که توسط بسیاری از سازمانهای آماری در مورد جدولها و فایلهای داده‌های خرد استفاده می‌شود، با مثالهایی ساده معرفی گردند. هر کدام از روش‌های مطرح شده دارای ابعاد کاملتر و پیچیده‌تری است که نیازمند تحقیق و کار بسیار بیشتری می‌باشد.

ناحیه	بایین	متوسط	بالا	خیلی بالا	جمع
A	۱۵	۱	۳	۱	۲۰
B	۲۰	۱۰	۱۰	۱۵	۵۵
C	۳	۱۰	۱۰	۲	۲۵
D	۱۲	۱۴	۷	۲	۳۵
جمع	۵۰	۳۵	۳۰	۲۰	۱۳۵

جدول ۱: تعداد نوجوانان بزهکار به تفکیک ناحیه و سطح تحصیلات سرپرست (با افشاء)

ناحیه	بایین	متوسط	بالا	خیلی بالا	جمع
A	۱۵	D۱	D۲	D۳	۲۰
B	۲۰	D۴	D۵	۱۵	۵۵
C	D۶	۱۰	۱۰	D۷	۲۵
D	D۸	۱۴	۷	D۹	۳۵
جمع	۵۰	۳۵	۳۰	۲۰	۱۳۵

جدول ۲: تعداد نوجوانان بزهکار به تفکیک ناحیه و سطح تحصیلات سرپرست (با افشاء)

ناحیه	بایین	متوسط	بالا	خیلی بالا	جمع
A	۱۵	F	F	F	۲۰
B	۲۰	F	F	۱۵	۵۵
C	F	۱۰	۱۰	F	۲۵
D	F	۱۴	۷	F	۳۵
جمع	۵۰	۳۵	۳۰	۲۰	۱۳۵

جدول ۳: تعداد نوجوانان بزهکار به تفکیک ناحیه و سطح تحصیلات سرپرست (بدون افشاء)،
خانه‌هایی که با F نشان داده شده‌اند پنهان شده‌اند

رتبه پوست	درآمد سربرست	سطح تحصیلات سربرست	ناحیه	نوجوان	شماره
سیاه	۲۰۱	خیلی بالا	A	John	۱
سفید	۱۰۳	بالا	A	Jim	۲
سیاه	۷۷	بالا	A	Sue	۳
سفید	۶۱	بالا	A	Pete	۴
سفید	۷۲	متوسط	A	Ramesh	۵
سفید	۱۰۳	پایین	A	Danis	۶
سیاه	۹۱	پایین	A	Virgil	۷
سفید	۸۴	پایین	A	Wanda	۸
سفید	۷۵	پایین	A	Stan	۹
سیاه	۶۲	پایین	A	Irmie	۱۰
سفید	۵۸	پایین	A	Renee	۱۱
سیاه	۵۶	پایین	A	Virginia	۱۲
سیاه	۵۴	پایین	A	Mary	۱۳
سفید	۵۲	پایین	A	Kim	۱۴
سیاه	۵۵	پایین	A	Tom	۱۵
سفید	۴۸	پایین	A	Ken	۱۶
سفید	۴۸	پایین	A	Mike	۱۷
سیاه	۴۱	پایین	A	Joe	۱۸
سیاه	۴۴	پایین	A	Jeff	۱۹
سفید	۳۷	پایین	A	Nancy	۲۰

جدول ۴: رکوردهای نوجوانان بزهکار در ناحیه A (درآمد بر حسب هزار دلار است)

رتبه پوست	درآمد سربرست	سطح تحصیلات سربرست	ناحیه	نوجوان	شماره
سیاه	۲۰۱	خیلی بالا	A	John	۱
سفید	۱۰۳	بالا	A	Jim	۲
سیاه	۷۵	بالا	A	Sue	۳
سفید	۶۱	خیلی بالا	A	Alfonso	۴
سفید	۷۲	متوسط	A	Ramesh	۵
سفید	۱۰۳	پایین	A	Danis	۶
سیاه	۹۱	پایین	A	Virgil	۷
سفید	۸۴	پایین	A	Wanda	۸
سفید	۷۵	پایین	A	Stan	۹
سیاه	۶۲	پایین	A	Irmie	۱۰
سفید	۵۸	پایین	A	Renee	۱۱
سیاه	۵۶	بالا	A	June	۱۲
سیاه	۵۴	پایین	A	Mary	۱۳
سفید	۵۲	پایین	A	Kim	۱۴
سیاه	۵۵	پایین	A	Tom	۱۵
سفید	۴۸	پایین	A	Ken	۱۶
سفید	۴۸	متوسط	A	George	۱۷
سیاه	۴۱	پایین	A	Joe	۱۸
سیاه	۴۴	پایین	A	Jeff	۱۹
سفید	۳۷	پایین	A	Heather	۲۰

جدول ۵: رکوردهای نوجوانان بزهکار در ناحیه A پس از مبادله رکوردها (درآمد بر حسب هزار دلار است)

ناحیه	بایین	متوسط	بالا	خیلی بالا	جمع
A	۱۳	۲	۳	۲	۲۰
B	۱۸	۱۲	۸	۱۷	۵۵
C	۵	۹	۱۱	۰	۲۵
D	۱۴	۱۲	۸	۱	۳۵
جمع	۵۰	۳۵	۳۰	۲۰	۱۳۵

جدول ۶: تعداد نوجوانان بزرگار به تفکیک ناحیه و سطح تحصیلات سرپرست (بدون افساء)

ردیگ پوست	درآمد سرپرست	سطح تحصیلات	ناحیه	شماره
۱	۶۱	بالا	AC	
۲	۴۸	پایین	AC	
۳	۳۰	متوسط	AC	
۴	۵۲	متوسط	AC	
۵	۱۱۷	خیلی بالا	AC	
۶	۱۳۸	خیلی بالا	B	
۷	۱۰۳	خیلی بالا	B	
۸	۴۵	پایین	B	
۹	۶۲	متوسط	B	
۱۰	۸۵	بالا	B	
۱۱	۳۳	پایین	D	
۱۲	۵۱	متوسط	D	
۱۳	۵۹	متوسط	D	
۱۴	۷۲	بالا	D	

جدول ۷: رکوردهای نوجوانان بزرگار پس از حذف مشخص کننده‌ها (درآمد بر حسب هزار دلار است)

ردیگ پوست	درآمد سرپرست	سطح تحصیلات	ناحیه	شماره
۱	۶۰-۶۹	بالا	AC	
۲	۴۰-۴۹	پایین	AC	
۳	<۴۰	متوسط	AC	
۴	۵۰-۵۹	متوسط	AC	
۵	>۱۰۰	خیلی بالا	AC	
۶	>۱۰۰	خیلی بالا	B	
۷	>۱۰۰	خیلی بالا	B	
۸	۴۰-۴۹	پایین	B	
۹	۶۰-۶۹	متوسط	B	
۱۰	۸۰-۸۹	بالا	B	
۱۱	<۴۰	پایین	D	
۱۲	۵۰-۵۹	متوسط	D	
۱۳	۵۰-۵۹	متوسط	D	
۱۴	۷۰-۷۹	بالا	D	

جدول ۸: رکوردهای نوجوانان بزرگار پس از حذف مشخص کننده‌ها، محدود کردن جزئیات جغرافیایی و کدبندی (درآمد بر حسب هزار دلار است)

رتبه	نام	جنس	محل زندگی	ساختار خانوار	درآمد سالانه	وضعیت اقتصادی
۱	ABC	بازار	بازار	۶۱	سفید	
۲	ABC	پایین	پایین	۴۲	سفید	
۳	ABC	متوسط	متوسط	۳۲	سیاه	
۴	ABC	متوسط	متوسط	۵۲	سفید	
۵	ABC	خیلی بالا	خیلی بالا	۱۲۳	سفید	
۶	B	خیلی بالا	خیلی بالا	۱۲۸	سیاه	
۷	B	خیلی بالا	خیلی بالا	۹۴	سفید	
۸	B	پایین	پایین	۴۶	سفید	
۹	B	متوسط	متوسط	۶۱	سفید	
۱۰	B	بالا	بالا	۸۲	سفید	
۱۱	D	پایین	پایین	۳۱	سیاه	
۱۲	D	متوسط	متوسط	۵۲	سیاه	
۱۳	D	متوسط	متوسط	۵۵	سفید	
۱۴	D	بالا	بالا	۶۱	سیاه	

جدول ۹: رکوردهای نوجوانان بزهکار پس از افزودن پرشیدگی تصادفی (درآمد بر حسب هزار دلار است)

مراجع

- 1- Bader , E. *Basic right information-Basic right to data protection: A contradiction*, Bulletin of the international statistical institute 49th session, 1993, p.63.
- 2- Cox, L.H. and Zayatz, L.V. *An Agenda for Research in Statistical Disclosure Limitation*, Journal of Official Statistics, Vol.11, No.2 , 1995, pp. 205-220.
- 3- Dalenius T. Towards a methodology for statistical disclosure control. *Statistisk Tidskrift* 1977.
- 4- Duncan G.T., Jabin T.B., de Wolf V.A. (eds). *Private Lives and Public Policies: Confidentiality and Accessibility of Goverment Statistics*. Panel on Confidentiality and Data Access, Committee on National Statistics. National Academy Press: Washington, D.C., 1993.
- 5- Fienberg, S.E. *Conflicts Between the Needs for Access to Statistical Information and Demands for Confidentiality*, Journal of official Statistics Vol. 10, No. 2, 1994, pp.115-132.

- 6- Fienberg, S.E. and Makov, U.E. *Confidentiality , Uniqueness and Disclosure Limitation for Categorical Data*, Journal of official Statistics, Vol.14 No. 1998, pp.385-397.
- 7- Fienberg, S.E. *Confidentiality and data access in public health*, Statist. Med. 2001 ; 20: 1347-1356.
- 8- Subcommitee on Disclosure-Avoidance Techniques. Statistical Policy Working Paper No. 22: Report on statistical disclosure limitation methodology. Federal Committee on Statistical Methodology, Statistical Policy Office, Office of Information and Regulatory Affairs, Office of Management and Budget, Washington, D.C., 1994.
- 9- Subcommitee on Disclosure-Avoidance Techniques. Statistical Policy Working Paper No. 2: Report on statistical disclosure and disclosure avoidance techniques. Federal Committee on Statistical Methodology, Office of Federal Policy and Standards, U.S. Department of Commerce, Washington, D.C., 1978.

توسعه کنترل فرآیند آماری با استفاده از رویکرد بیز تجربی

رسول نورالسناء^۱، محمد رضا لطفی^۲

P ۱۲۰۴۸

^۱ دانشکده صنایع، دانشگاه علم و صنعت ایران

^۲ دانشگاه آزاد آزاد واحد علوم تحقیقات

چکیده: در این تحقیق، کنترل فرآیند آماری جهت یک فرآیند با توزیع پواسن دارای پارامتر λ که این پارامتر در طی زمان تغییر می‌کند، توسعه داده می‌شود. برای این پارامتر متغیر، توزیع گاما بعنوان توزیع پیشین منظور شده و جهت بروز نمایی آن از مشاهدات اخیر، یک برآورد پسین ارائه می‌گردد. ساختار فوق شرایط استفاده از تئوری بیز تجربی را فراهم می‌سازد. این موضوع مدلی بهینه کنترل فرآیند ایجاد می‌کند و با استفاده از روابط بازگشتی سرعت پردازش اطلاعات افزایش می‌یابد و حجم داده‌ها و ذخیره‌سازی اطلاعات را در شبکه‌های رایانه‌ای در فرآیندهای بزرگ و پیچیده کاهش می‌دهد.

واژه‌های کلیدی: کنترل فرآیند، بیز تجربی، توزیع پواسن، توزیع گاما، برآورد میانگین متحرک وزنی، برآوردهای بازگشتی.

۱ مقدمه

در کنترل فرآیند تولید با استفاده از روش‌های آماری، برای مشخصه‌های کیفی محصول توزیع احتمال مناسب منظور می‌گردد و سعی می‌شود با کنترل نمودن پارامترهای آن، عملکرد فرآیند تولید را کنترل نمود و از این طریق میزان ضایعات و محصولات معیوب را کاهش داد. از دهه ۱۹۲۰ که آقای والترووهارت کنترل فرآیند آماری را معرفی کرد معمولاً فرض اصلی، ثابت بودن پارامترهای توزیع مفروض بوده است. به دلیل پیچیده شدن فرآیندهای تولید ضرورت می‌یابد که این فرض مورد تجدیدنظر قرار گیرد. در سالهای اخیر تحقیقات مهمی در زمینه کاربرد تئوری بیز تجربی در کنترل فرآیند بعمل آمده است. بویژه یوسری و دیگران (۱۹۹۱) از رویکرد بیزی برای یک توزیع دوجمله‌ای با پارامتر متغیر p_t استفاده کرده‌اند. فلتز و شیا (۲۰۰۱) این رویکرد را برای توزیع نرمال چند متغیره با پارامتر متغیر μ بکار برdenد. ویلو و رودی (۲۰۰۱) از تئوری بیز در آزمون تایید قابلیت استفاده می‌کنند. همچنین مارتز و والر (۱۹۸۹) در کتاب تجزیه و تحلیل قابلیت بیزی کاربرد تئوری بیزی را در تجزیه و تحلیل قابلیت اعتماد و دوام معرفی کردند.

جین و گرینما (۱۹۹۳) کنترل کیفیت چند متغیره با رویکرد بیزی را ارائه دادند. چون مشخصه‌های کیفی در بسیاری از فرآیندها از توزیع پواسن پیروی می‌کند ما در این مقاله، کاربرد تئوری بیز در کنترل فرآیند آماری را با در نظر گرفتن توزیع پواسن با پارامتر متغیر λ توسعه می‌دهیم.

۲ مدل

اگر X_t تعداد نقص یا قطعات معیوب در زمان t و دارای توزیع پواسن با پارامتر λ_t باشد در اینصورتتابع احتمال X_t بشرح زیر است:

$$f(x_t|\lambda_t) = \frac{e^{-\lambda_t} \lambda_t^{x_t}}{x_t!} \quad x_t = 0, 1, 2, \dots \quad (1)$$

در این تحقیق فرض می‌شود پارامتر λ_t ثابت نمی‌باشد و در طی زمان تغییر می‌کند و یک متغیر تصادفی است و دارای توزیع گاما با پارامترهای α و β است. بنابراین :

$$f(\lambda_t) = \frac{1}{\beta^\alpha \Gamma(\alpha)} \lambda_t^{\alpha-1} e^{-\lambda_t/\beta} \quad \lambda_t > 0, \quad \alpha, \beta > 0. \quad (2)$$

در اینصورت به این تابع احتمال، توزیع احتمال پیشین پارامتر λ_t گفته می‌شود. میانگین و واریانس این توزیع پیشین که به آن میانگین و واریانس فرآیند نیز می‌گویند عبارت است از :

$$\mu = \alpha\beta \quad (3)$$

$$\gamma^2 = \alpha\beta \quad (4)$$

برای برآورد λ_t در زمان t و بروز نمایی آن از نمونه گیری و مشاهدات اخیر استفاده می‌شود. اگر X_t تعداد نقص در نمونه دوره t باشد، براساس تئوری بیز تجزیی تابع چگالی احتمال λ_t به شرط X_t بشرح زیر است و به آن توزیع احتمال پسین پارامتر λ_t گفته می‌شود.

$$f(\lambda_t|x_t) = \frac{f(\lambda_t)f(x_t|\lambda_t)}{\int f(\lambda_t)f(x_t|\lambda_t)d\lambda_t} \quad (5)$$

می‌توان نشان داد که این توزیع گاما با پارامترهای $\alpha + x_t$ و $\beta + 1$ است، درنتیجه میانگین و واریانس پسین λ_t می‌شود:

$$\mu_{\lambda_t} = \frac{\beta(x_t + \alpha)}{\beta + 1}, \quad \sigma_{\lambda_t}^2 = \frac{\beta^2(x_t + \alpha)}{(\beta + 1)^2} \quad (6)$$

اگر واریانس کل X_t با σ_T^2 نشان داده شود، این تغییر پذیری کل از دو منع تغییر پذیری شامل تغییر پذیری یا واریانس فرآیند (γ^2) و تغییر پذیری حاصل از نمونه‌گیری که به آن واریانس نمونه‌ای می‌گویند و در اینجا با σ_S^2 نشان می‌دهیم، تشکیل شده است. بنابراین :

$$\begin{aligned}\sigma_T^2 &= \sigma_S^2 + \gamma^2 = E_{\lambda_t} [V(X_t|\lambda_t)] + V_{\lambda_t} [E(X_t|\lambda)] \\ &= E_{\lambda_t} (\lambda_t) + V_{\lambda_t} (\lambda_t) \\ &= \alpha\beta + \alpha\beta^2\end{aligned}\quad (7)$$

بطور مشابه، منطقی است که بخشی از میانگین پسین λ_t ، تحت تاثیر نمونه‌گیری و بخشی دیگر ناشی از وضعیت فرآیند باشد، لذا می‌توان آنرا یک میانگین وزنی بصورت زیر در نظر گرفت:

$$\mu_{\lambda_t} = E(\lambda_t|x_t) = wx_t + (1-w)\mu \quad (8)$$

جهت تعیین مقدار مناسب برای w به این نکته توجه می‌شود که هر چقدر پراکندگی فرآیند یا نمونه‌گیری کمتر باشد باید اهمیت وزن آن بیشتر در نظر گرفته شود، لذا نسبت زیر مناسب خواهد بود.

$$w = \frac{\gamma^2}{\sigma_T^2} = \frac{\alpha\beta}{1+\beta} \quad (9)$$

با توجه به اینکه در همه موارد فوق پارامترهای α و β نقش اساسی دارند، برآورد آنها مهم خواهد بود.

۳ برآوردها

با توجه به روابط (۳) و (۷) جهت برآورد α و β می‌توان از میانگین توزیع پیشین (μ) و واریانس کل نمونه‌ای (σ_T^2) استفاده نمود، لذا ابتدا برآوردهای آنها را می‌گردد.

۱-۳ برآوردها

چون X_t یک برآورد ناریب برای μ است زیرا:

$$E(X_t) = E_{\lambda_t} [E(X_t|\lambda_t)] = E_{\lambda_t} (\lambda_t) = \alpha\beta = \mu$$

در نتیجه یک برآورد کننده مناسب و ناریب با واریانس کمتر برای μ را می‌توان \bar{x}_t را در نظر گرفت:

$$\hat{\mu}_T = \bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t \quad \text{Var}(\hat{\mu}_T) = \frac{\sigma_T^2}{T} \quad (10)$$

σ_T^2 برآورده است (۲-۳)

برآورده کننده زیر که می‌توان نشان داد یک برآورده کننده ناارباست، σ_T^2 را برآورد می‌کند.

$$\hat{\sigma}_T^2 = \frac{1}{T-1} \left[\sum_{t=1}^T (x_t - \hat{\mu}_T)^2 \right] = \frac{1}{T-1} \left[\sum_{t=1}^T x_t^2 - T\hat{\mu}_T^2 \right] \quad (11)$$

برآوردهای α و β (۳-۳)

براساس روابط (۳) و (۷) داریم:

$$\begin{aligned} \mu &= \alpha\beta \\ \sigma_T^2 &= \alpha\beta + \alpha\beta^2 \end{aligned}$$

در نتیجه :

$$\begin{aligned} \alpha &= \frac{\mu^2}{\sigma_T^2 - \mu} \\ \beta &= \frac{\sigma_T^2 - \mu}{\mu} \end{aligned}$$

و با قرار دادن برآوردهای μ و σ_T^2 برآوردهای α و β بشرح زیر تعیین می‌گردد:

$$\begin{aligned} \hat{\alpha} &= \frac{\hat{\mu}^2}{\hat{\sigma}_T^2 - \hat{\mu}} \\ \hat{\beta} &= \frac{\hat{\sigma}_T^2 - \hat{\mu}}{\hat{\mu}} \end{aligned} \quad (12)$$

برآوردهای α و β (۴-۳)

$$\gamma^2 = \alpha\beta^2 \rightarrow \hat{\gamma}^2 = \hat{\sigma}_T^2 - \hat{\mu}^2 \quad (13)$$

با توجه به اینکه $\gamma^2 > \alpha > \beta > 0$ است، بدیهی است برآوردهای فوق وقتی معتبر است که:

$$\hat{\sigma}_T^2 - \hat{\mu}^2 > 0 \rightarrow \hat{\sigma}_T^2 > \hat{\mu}^2$$

برآوردهای میانگین متحرک وزنی نمایی (EWMA) و رابطه بازگشتی

در برآوردهای قبلی کلیه T مشاهده x_1, \dots, x_T دارای تاثیرگذاری و اهمیت یکسان بوده است ، جهت واقعی تر نمودن و بهینه‌سازی بروز نمایی اطلاعات ، مناسب تراست که مشاهدات جدیدتر، اهمیت و وزن بیشتری داشته باشند، بدین منظور از برآورد EWMA استفاده می‌گردد

$$\hat{\mu}_T = \frac{\sum_{t=0}^{T-1} r^t x_{T-t}}{\sum_{t=0}^{T-1} r^t} = \frac{(1-r) \sum_{t=0}^{T-1} r^t x_{T-t}}{1-r^T} \quad (14)$$

$$\hat{\sigma}_T^2 = \frac{\sum_{t=0}^{T-1} r^t (x_{T-t} - \hat{\mu}_T)^2}{\sum_{t=0}^{T-1} r^t} = \frac{(1-r) \sum_{t=0}^{T-1} r^t x_{T-t} - \hat{\mu}_T^2}{1-r^T} \quad (15)$$

در محاسبات کامپیوتری بویژه در نمونه‌گیری با مشاهدات زیاد ، استفاده از روابط (14) و (15) زمان و حافظه و هزینه زیادی را شامل می‌گردد. جهت کاهش این موارد می‌توان روابط فوق را بصورت روابط بازگشتی زیر نوشت و از آن استفاده نمود.

$$\hat{\mu}_T = \frac{1-r}{1-r^T} x_T + \frac{r-r^T}{1-r^T} \hat{\mu}_{T-1} \quad (16)$$

$$\hat{\sigma}_T^2 = \frac{1-r}{1-r^T} x_T^2 + \frac{r-r^T}{1-r^T} (\hat{\sigma}_{T-1}^2 + \hat{\mu}_{T-1}^2) - \hat{\mu}_T^2 \quad (17)$$

بنابراین برآورد μ در زمان T یک ترکیب خطی از مشاهده فعلی (x_T) و برآورد μ در یک دوره زمانی قبل و بطور مشابه برآورد σ_T^2 یک ترکیب خطی از x_T^2 و برآوردهای میانگین واریانس در یک دوره قبل و میانگین دوره فعلی می‌باشد. در عمل برای اینکه مشاهدات جدید وزن بیشتری داشته باشند، $1 < r < 9/10$. در نظر گرفته می‌شود. در روابط فوق چون در هر دوره یک مشاهده به حجم نمونه اضافه می‌شود، می‌توان طول دوره را ثابت و بطول n در نظر گرفت در این صورت روابط بالا بصورت زیر در می‌آید.

$$\hat{\mu}_T = \frac{\sum_{t=0}^{n-1} r^t x_{T-t}}{\sum_{t=0}^{n-1} r^t} = \frac{1-r}{1-r^n} x_T + r \hat{\mu}_{T-1} + \frac{r^{n+1}-r^n}{1-r^n} x_{T-n} \quad (18)$$

$$\begin{aligned} \hat{\sigma}_T^2 = \frac{\sum_{t=0}^{n-1} r^t (x_{T-t} - \hat{\mu}_T)^2}{\sum_{t=0}^{n-1} r^t} &= \frac{1-r}{1-r^n} x_T^2 + r (\hat{\sigma}_{T-1}^2 + \hat{\mu}_{T-1}^2) \\ &+ \frac{r^{n+1}-r^n}{1-r^n} x_{T-n}^2 - \hat{\mu}_T^2 \end{aligned} \quad (19)$$

۵ برآوردهای حدی

وقتی تعداد مشاهدات (T یا n) زیاد باشد، به عبارت دیگر به سمت بینهایت میل کند در اینصورت روابط (۱۶) و (۱۸) به رابطه (۲۰) و روابط (۱۷) و (۱۹) به رابطه (۲۱) گرایش می‌یابد.

$$\hat{\mu}_T = (1 - r)x_T + r\hat{\mu}_{T-1} \quad (20)$$

$$\hat{\sigma}_T^2 = (1 - r)x_T^2 + r(\hat{\sigma}_{T-1}^2 + \hat{\mu}_{T-1}^2) - \hat{\mu}_T^2 \quad (21)$$

بنابراین برآورد کننده‌های حدی یک ترکیب محدب از مشاهدات یا برآوردهای فعلی و یک دوره قبل می‌باشد.

۶ برآوردهای کوتاه مدت و بلند مدت

در برآوردهای EWMA و بازگشتی، اگر r نزدیک به ۱ باشد، ضرایب r^t دیرتر به صفر میل می‌کند و از مشاهدات دورتر و گذشته بیشتر استفاده می‌گردد، لذا به این حالت، برآورد کننده بلند مدت گفته می‌شود. و اگر r نزدیک به $0/9$ باشد ضرایب r^t زودتر به صفر گرایش دارد و مشاهدات گذشته کمتر منظور می‌گردد و به این حالت برآورد کننده کوتاه مدت گفته می‌شود. مثال زیر یک نمونه از کاربردهای نتایج این تحقیق در بروز نمایی پارامتر λ جهت تجزیه و تحلیل آماری و تشکیل نمودارهای کنترل در فرایندهای تولیدی می‌باشد.

مثال: اگر در یک فرآیند تولید طی ۲۰ روز گذشته تعداد قطعات معیوب به ترتیب $52, 60, 65, 48, 70, 55, 63, 82, 62, 45, 50, 64, 57, 68, 65, 59, 61, 78, 72, 80, 79$ قطعه در روز باشد. براساس این مشاهدات ۲۰ روزه و مفروضات و روابط حاصله درین تحقیق، برآورد پارامترها و مشخصه‌های کیفی فرآیند تولید برای روزهای بعدی براساس برآوردهای بلند مدت ($r=0/99$) و کوتاه مدت ($r=0/90$) و دوره ثابت بطول $n=20$ (روابط (۱۸) و (۱۹)) بشرح زیراست.

روز T	مشاهده x_i	برآورد بلند مدت ($r = ۰/۹۹$)					
		$\hat{\mu}_T$	$\hat{\sigma}_T^2$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\mu}_{\lambda_t}$	$\hat{\sigma}_{\lambda_t}^2$
۲۰	۷۹	۶۳/۳	۱۰۹	۸۷/۱	۰/۷۲	۶۵/۲	۲۷/۴
۲۱	۷۰	۶۳/۶	۱۰۶/۲	۹۴/۲	۰/۶۷	۶۱/۹	۲۴/۸
۲۲	۷۳	۶۴/۳	۱۰۶/۲	۹۸/۷	۰/۶۵	۶۷/۶	۲۶/۶
۲۳	۸۸	۶۵/۳	۱۳۹/۶	۵۷/۴	۱/۱۴	۷۷/۵	۴۱/۳
۲۴	۵۳	۶۵/۴	۱۳۳/۳	۶۳	۱/۰۴	۵۹/۱	۳۰/۱
۲۵	۹۷	۶۷/۶	۱۷۸	۴۱/۴	۱/۶۳	۸۵/۸	۵۲/۲

روز T	مشاهده x_i	برآورد کوتاه مدت ($r = ۰/۹۰$)					
		$\hat{\mu}_T$	$\hat{\sigma}_T^2$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\mu}_{\lambda_t}$	$\hat{\sigma}_{\lambda_t}^2$
۲۰	۷۹	۶۳/۳	۱۰۹	۸۷/۱	۰/۷۲	۶۵/۲	۲۷/۴
۲۱	۷۰	۶۲/۹	۱۲۱	۶۸/۱	۰/۹۲	۶۱/۴	۲۹/۴
۲۲	۷۳	۶۴/۱	۱۱۸/۷	۷۵/۳	۰/۸۵	۶۸/۱	۳۱/۳
۲۳	۸۸	۶۶/۸	۱۵۷/۷	۴۹/۱	۱/۳۶	۷۹	۴۵/۵
۲۴	۵۳	۶۵/۵	۱۰۶/۱	۴۷/۴	۱/۳۸	۵۸/۲	۳۲/۸
۲۵	۹۷	۶۹/۲	۲۴۴	۲۷/۴	۲/۵۲	۸۹/۲	۶۲/۹

مراجع

Yousry, Strum, Feltz and Noorossana (1991), Process Monitoring in Real Time : Empirical Bayes approach-Discrete Case. Qual. Reliab. Engng. Int., 7, 123-132.

Feltz and shiau (2001) statistical process monitoring using an Empirical Bayes Multivaraiate process contorol chart, Qual. Reliab. Engng. Int., 17, 119-124.

Weilu and Rudy (2001), Reliability Demonstration Test for a Finite population; Qual. Reliab. Engng. Int., 17, 33-38.

Martz and Waller (1989) Bayesian Reliability Analysis, Krieger, New York, NY.

Jain, Alt and Grimsnaw (1993), Multivariate quality contorl-A Bayesian approach, ASQC Quality congress Transactions, Boston : 645-641.

تصمیم آماری بیزی بوسیله برآوردهای MLM در فضاهای غیر حقیقی

حسین نوری امامزاده، مهدی دوست پرست

P ۱۳۰ ۵۵

گروه آمار، دانشگاه فردوسی مشهد

چکیده: زمانی که از برآوردهای پارامتری مجھول سخن گفته می‌شود، آنچه در اذهان تجلی می‌یابد حقیقی بودن پارامتر مجھول است. اما باید متذکر شویم که در بسیاری از موارد پارامتر واقعی درفضای اقلیدسی نیست. در این گونه مسائل معمولاً از دو برآوردهای ماکزیمم چگالی پسین (MAP) که بالاترین چگالی احتمال پسین را انتخاب می‌کند و مینیمم مربعات خطای ($MMSE$) که متوسط مربع فاصله از مقادیر واقعی پارامتر را کمینه می‌کند، استفاده می‌شود.

در این مقاله برآوردهای جدیدی را با نام برآوردهای ماکزیمم جرم موضعی (MLM) که روی چگالی احتمال موضعی انتگرال می‌گیرد، معرفی کرده و ضمن مقایسه آنها با یکدیگر، نشان می‌دهیم که در اکثر موارد، این برآوردهای کارآئی از برآوردهای MAP و $MMSE$ است و حتی از الگوریتم‌های استاندارد مورد استفاده، مناسب تراست. مثالهایی که از فیزیک نور ارائه می‌شود صحبت مطالب فوق را تأیید می‌کند.

واژه‌های کلیدی: برآوردهای ماکزیمم پسین، برآوردهای حداقل مربعات، برآوردهای ماکزیمم جرم موضعی، چگالی پسین، چگالی پیشین، الگوریتم.

۱ مقدمه

اغلب الگوریتم‌های شهودی در چارچوب آمار بیزی بررسی می‌شوند. بطور کلی دو برآوردهای مورد استفاده است که عبارتند از: برآوردهای ماکزیمم (MAP) و برآوردهای مینیمم میانگین مربعات خطای ($MMSE$) ما نشان می‌دهیم که هیچکدام از این دو برآوردهای مناسب نیست چرا که برآوردهای MAP از ساختار توزیع احتمال پسین استفاده نمی‌کند و برآوردهای $MMSE$ نوع خطای را منعکس نمی‌کند. در این مقاله برآوردهای جدیدی تحت عنوان ماکزیمم جرم موضعی (MLM) را معرفی می‌کنیم که با انتگرال گیری نسبت به چگالی احتمال موضعی بدست می‌آید. برآوردهای MLM نسبت به ساختار موضعی چگالی پسین حساس است در حالیکه MAP این گونه نیست. MLM محتمل ترین جواب به

طور تقریبی صحیح را می‌یابد و در حالتی که مشاهدات دارای اغتشاش کم هستند نیز تقریب کارایی فراهم می‌کند. این برآورده‌گر جدید را در یک مسئله فیزیک نور (color constancy) به کار می‌گیریم. فرض کنید یک پرتوی تابش مجھول به یک سطح مجھول تابیده می‌شود و خصوصیات پرتوی بازنگاری تابش سنسور ثبت می‌شود. می‌توانیم براساس اطلاعات بدست آمده از سنسور خصوصیات پرتوی تابش و سطح بازنگاری مشخص کنیم که این کار را با الگوریتم MLM انجام داده و نشان می‌دهیم مناسب تر از دو برآورده‌گر قبلی است.

در این مقاله، در بخش ۲ مسئله اساسی را مطرح کرده و معایب MAP و $MMSE$ را بیان می‌کنیم. بخش ۳ شامل ایده کلی جهت بدست آوردن برآورده‌گر مناسب است و همچنین با یک مثال ساده توزیعهای احتمالی پیشین و پسین وتابع زیان را معرفی می‌کنیم. در بخش ۴ برآورد MLM را برای یک مفهوم عمیق تر شهودی یعنی *color constancy* بیان و به کار می‌گیریم. در بخش ۵ یک مقایسه کلی بین سایر الگوریتمهای رایج و الگوریتم MLM صورت گرفته است. در بخش ۶ شامل نتایج بدست آمده است.

۲ مسئله کلی

هدف استنباط در مورد خصوصیات جهان طبیعی از قبیل: اشکال، رنگها و... است که با توجه سنسورهای خاصی اندازه‌گیری می‌شوند. روش کلی در برخورد با این مسئله آنالیز بیزی است به طوریکه داده‌های شهودی با توزیع احتمالات پیشین برای یافتن توزیع احتمال پسین خواص طبیعت ترکیب می‌شود که آن را پارامترهای منظری (*sceneparameters*) می‌نامیم. نوعاً، هدف انتخاب بهترین برآورده‌گر پارامتر منظری است. برای این منظور یک معیار بهینگی انتخاب شده و از توزیع احتمال پیشین برای یافتن برآورده‌گر بهینه استفاده می‌کنیم. دو قاعده تصمیم که بطور عمومی مورد استفاده‌اند MAP و $MMSE$ هستند. قاعده MAP مقداری از پارامترهای منظری را انتخاب می‌کند که چگالی احتمال پسین را ماکزیمم می‌کند. این برآورده‌گر رابطه تنگاتنگی با روش درستنمایی ماکزیمم دارد و در بسیاری از الگوریتمها از آن استفاده می‌شود.

قاعده $MMSE$ مقداری از پارامترهای منظری را انتخاب می‌کند که میانگین مربعات فاصله از مقدار واقعی پارامتر منظری را مینیمم می‌کند. به سادگی می‌توان نشان داد که میانگین توزیع پسین، $MMSE$ است. همچنین این برآورده‌گر بطور وسیعی مورد استفاده قرار گرفته است.

هنگامی که چگالی احتمال پسین فضای پارامتر منظری را بخوبی مشخص کند قاعده $MMSE$ و MAP تقریباً جواب یکسانی می‌دهند و در صورتیکه این چگالی کمی ساده

تر باشد این دو برآورده رضایت بخش نیستند. برآورده MAP به جزئیات ساختار چگالی پسین حساس نیست زیرا در آن، نقطه ماکریم چگالی احتمال مد نظر است. برآورده $MMSE$ ممکن است جوابی خارج از فضای پارامتر ارائه دهد. بعلاوه در محاسبات این برآورده ممکن است نیاز به حل یک انتگرال روی فضای پارامتر پیدا کنیم.

۳ ایده

جهت مقایسه MAP و $MMSE$ و بدست آوردن برآورده جدید، تابع زیان در این فضا معرفی می‌کنیم. برای یک پارامتر منظری θ تابع زیان را با $L(\theta, \hat{\theta})$ نشان داده و آن را زیان حاصل از برآورد $\hat{\theta}$ وقتی که پارامتر واقعی θ باشد، می‌نامیم. تابع زیان منجر به روشی برای برآورد پارامترهای منظری از روی توزیع پسین می‌شود بگونه‌ای که پارامترهایی که متوسط زیان را مینیمیم می‌کنند را انتخاب می‌کند. قاعده MAP تابع زیان را بصورت $L(\theta, \hat{\theta}) = -\delta(\theta - \hat{\theta})$ تعریف می‌کند که ما آن را δ -می‌نامند. قاعده $MMSE$ تابع زیان را بصورت $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ در نظر می‌گیرد.

معتقدیم، هیچکدام از این توابع زیان، تقریب مناسبی برای زیان حاصل از برآورد نادرست نمی‌دهد. تابع زیان δ -اشاره دارد که خطاهای کوچک هم به اندازه خطاهای بزرگ زیان حاصل می‌کنند. تابع زیان مربع خطای بیان می‌کند که با افزایش اندازه خطای برآورد میزان زیان حاصل بشدت افزایش می‌یابد لذا به دنبال تابع زیانی هستیم که معایب توابع زیان فوق را نداشته باشد یعنی خطاهای کوچک به اندازه خطاهای بزرگ نباشد و افزایش خطای زیان حاصل را به شدت و غیر معقول افزایش ندهد زیرا در طبیعت نیز با مسائلی شویم اما با خطای زیاد ممکن است حتی جهت توب را تشخیص ندهیم.

ما یک تابع زیان معرفی می‌کنیم که برای خیلی از مسائل شهودی نیز مناسب باشد این تابع زیان با انتگرال گیری روی چگالی احتمال برای یافتن جرم احتمال موضعی در یک ناحیه مشخص می‌شود که ما آنرا تابع زیان جرم موضعی نامیده و برآورده حاصل را برآورده ماکریم جرم موضعی (MLM) می‌نامیم.

۱.۳ مثال

فرض کنیدیک عدد (y) مشاهده شده که حاصل ضرب دو عدد دیگر است یعنی $y = ab$. هدف تعیین مقادیر a و b است. در ارتباط با داده‌های بصری (شهودی) y می‌تواند به

عنوان یک تصویر حاصل از دو عامل دیگر (یا پارامتر منظری) در نظر گرفته شود. از نظر ریاضی این یک مسئله غیر خطی و با جواب نامتناهی است.

می توان فرض کرد $y = ab$ و $a, b > 0$ برقرار است. اگر بخواهیم به تنها یی از معلومات هندسی استفاده کنیم فقط می توان گفت که جواب در امتداد یک جزی از هذلولی $-ab = 1$ قرار می گیرد. یک آنالیز احتمالی جواب دقیق تری ارایه می دهد. می خواهیم توزیع احتمالی پسین $p(\theta|y)$ را بایابیم که $\theta = (a \ b)^T$. با استفاده از قضیه بیز می توان نوشت:

$$P(\theta|y) = \frac{P(y|\theta)\Pi(\theta)}{P_y(y)} \quad (1)$$

ثابت نرمال ساز $P_Y(y)$ مستقل از پارامتر θ است که به دنبال برآورد آن هستیم. $\Pi(\theta)$ توزیع پیشین پارامتر منظری θ است و $p(y|\theta)$ تابع درستنمایی است. اجازه می دهیم مشاهدات همراه با اغتشاش باشد، لذا مدل احتمالی است. تابع تحويلی $f(\theta)$ مشاهدات اختصاصی پارامترهای منظری θ بدون اغتشاش را بیان می کند. در این مثال تابع تحويلی عبارت است از

$$f(\theta) = a.b = \theta_1\theta_2$$

فرض می کنیم اغتشاش مشاهدات دارای توزیع نرمال با میانگین 0 و واریانس σ^2 باشد. پس تابع درستنمایی به صورت

$$P(y|\theta) = \frac{1}{(\sqrt{2\pi\sigma^2})^N} e^{-\frac{\|y-f(\theta)\|^2}{2\sigma^2}} \quad (2)$$

چون $0 < \theta_1, \theta_2 < 4$ توزیع پیشین را بصورت

$$P_\theta(\theta) = \frac{1}{16} \quad (3)$$

روی ناحیه $[0, 4] \times [0, 4]$ و برابر در سایر نقاط، تعریف می کنیم. با جایگذاری (2)، (3) در (1) داریم

$$P(\theta|y=1) = \begin{cases} C e^{\frac{-(1-ab)^2}{2\sigma^2}} & 0 < a, b < 4 \\ 0 & \text{در غیر این صورت} \end{cases}$$

که در آن C یک ثابت است. در شکل (1) این توزیع پسین نمایش داده شده است. نقاط در طول هذلولی $1 = ab$ تشکیل یک مرز بابا لاترین احتمال می دهند چون مشاهدات دارای اغتشاش هستند جفت پارامتر دارای احتمال ناصفراند. با توجه به شکل (b) ۱) توجه می کنیم که هر چند همه نقاط مرز دارای ارتفاع یکسانی هستند اما نزدیک نقطه (1, 1)

عريض تراز ساير نقاط درامتداد هذلولي است.

برآورد بهينه $\hat{\theta}$ بامي نيمم كردن متوسط زيان $R(\hat{\theta}|y)$ به دست مي آيدكه آن را يسک بيزى مى ناميم. داريم

$$R(\hat{\theta}|y) = \int L(\theta, \hat{\theta})P(\theta|y)d\theta \quad (4)$$

درحالت خاصى كه $L(\theta, \hat{\theta}) = L(\theta - \hat{\theta})$ ريسک بيزى بصورت

$$R(\hat{\theta}|y) = \int L(\theta - \hat{\theta})P(\theta|y)d\theta = L(-\hat{\theta})$$

مي باشد.

شكل ۲ زيان حاصل شده باقاعده MAP و ريسک بيزى آن را برای مسئله $y = ab$ با توزيع پيشين يكداخت را بدست مي دهد. برای برآوردگر MAP ريسک بيزى قرينه تابع احتمال پسین است. جهت واضح نشان دادن و مقاييسه باشكال (۱) آن را به طرف بالا رسم كرده‌aim. هر نقطه واقع بر مرز هذلولي دارای چگالي احتمال يكسان است. پس قاعده MAP برآورد منحصر به فردی نمى دهد. همچوين اين قاعده از همه تغييرات درعرض چشم پوشى مى كندكه مى تواند منبع معنى داري از اطلاعات باشد.

شكل ۳ تابع زيان برآوردگر $MMSE$ و ريسک بيزى آن رانشان مى دهد. توجه داريم كه برآورد $MMSE$ حاصل برابر $a = b = 1.3$ است كه نتيجه عجيبى است. چرا كه $a \times b \neq 1$

تابع زيانى را كه برای قاعده MLM تعريف مى كنيم نقص دو برآوردگر MAP و $MMSE$ را بطرف مى كند.

تابع زيان جرم موضوعى را با علامت منفى به فرم نرمال با همبستگى كم بصورت

$$L(\theta, \hat{\theta}) = -\exp\{-|K_L^{-\frac{1}{2}}(\theta - \hat{\theta})|^2\}$$

تعريف مى كنيم كه درآن :

$$|K^{-\frac{1}{2}}\theta|^2 = \theta^T K^{-1} \theta$$

چون مقادير ويژه ماترييس K_L به قدرکافى كوچك هستند. اين تابع زيان تقربياً جواب صحيح مى دهد. (شكل ۴). همچوين شكل ۴ متوسط زيان به روش MLM برای مسئله $1 = ab$ را نشان مى دهد كه برآورد MLM حاصل به صورت (۱، ۱) است. چنانچه مشاهده مى شود اين برآورد عيب هاي دو برآوردگر قبلى راندارد.

۲.۳ نتیجه

این مثال ساده چندین مطلب را نشان می دهد.

—تابع زیان دارای اهمیت زیادی است که می تواند برآورده بینه را تغییر دهد.

—برآورده MAP از عرض های مرز ماکریم درستنمایی چشم پوشی می کند.

—برآورده $MMSE$ که از نظر میانگین مربعات خطأ، بینه است جواب نادرست می دهد.

$$(1) \quad 1.3 \neq 1.3 \times 1.3$$

—برآورده MLM ساختار جرم احتمال موضعی رابه کارمی گیرد و برآورده بینه را رائمه می دهد. و همانطور که مشاهده شد از دو برآورده قبلی بهتر بود.

۴ کاربرد در *color constancy*

در برخورد پرتو تابش به برخی سطوح پرتوی بازتابش تلفیقی از رنگها است اما در برخی سطوح دیگر فقط شکست نور حاصل می شود. در اینجا هدف برآورده مشخصه های نور تابیده شده از روی اطلاعاتی است که سنسور از پرتو بازتابش می دهد، می باشد. به این معنی که نه پرتو تابش و نه سطح بازتابش معلوم هستند.

فرض کنیم مجموعه ای شامل N_L سطح داریم. قابلیت بازتاب r امین سطح را با یک بردار ستونی S_j نشان می دهیم. مؤلفه های این بردار ستونی نشان دهنده شکست نور مشاهده شده، منعکس شده در N_λ نوار طیف نوری است. به طور مشابه قدرت توزیع پرتو تابش $S.P.D$ را با بردار ستونی e نمایش می دهیم. $S.P.D$ نور رسیده از دستگاه c_j است که

$$c_j = diag(e) \times s_j$$

در سیستم شهودی از پرتو بازتابش c_j با N_r کلاس از فوتون سنسورهای خطی نمونه گیری می شود. حساسیت های فوتون سنسور را بamatrise R مشخص می کنیم. pq امین درایه R نشان دهنده حساسیت p امین سنسور به نور در q امین طول موج نوار است. بردار r_j نشان دهنده پاسخ هریک از کلاس های سنسور به پرتو بازتابیده c_j است. صرف نظر از خطای سنسور معادله تحويلی r_j به r امین سطح رنگی به عنوان تابعی از پارامترهای منظری، e و s_j بصورت

$$r_j = f(e_j, s_j) = R c_j = R diag(e)s_j \quad (6)$$

بالین نمادها، هدف یافتن پرتو تابش e و سطح بازتابش s_j از پاسخ r_j برای هر سطح است.

۱.۴ تعاریف و نماد

تعريف ۱

گوییم یک پرتو تابش دریک مدل خطی B_e بعدی صدق می‌کند اگر توان نوشت: $e = B_e W_e$ که در آن B_e یک ماتریس با N_s سطر و N_e ستون و W_e یک بردار ستونی N_e بعدی است. مؤلفه‌های W_e را وزن‌های مدل خطی برای پرتو e گوییم.

تعريف ۲

گوییم یک سطح بازتابش دریک مدل خطی B_s صدق می‌کند اگر توان نوشت: $s_j = B_s W_{sj}$ که در آن B_s یک ماتریس با $N - \lambda$ سطر و N_s ستون و W_{sj} یک بردار ستونی N_s بعدی است. CIE یک مدل خطی ۳ بعدی را برای نور خورشید استاندارد کرده است. $L.T.Maloney$ نشان داده است که در مدل‌های خطی با بعدهای کم (۳ بعدی) درصد زیادی از خطای اندازه‌گیری سطح بازتابش کاهش می‌یابد. در این نمایش معادله تحويلی (۶) به صورت:

$$r_j = f(W_e, W_{sj}) = R \operatorname{diag}(B_e W_e) B_s W_{sj} = R \operatorname{diag}(B_s W_{sj}) B_e W_e \quad (6)$$

۲.۴ الگوریتم‌ها

پیشرفته‌ترین الگوریتم‌ها در فیزیک نور تحت شرط خطی بودن مدل هستند. *Maloney* و *Wandell* نشان دادند که وقتی N_r کلاس از اشعه بازتابش داریم ممکن است معادله تحويلی معکوس گردد اگر اشعه تابش دریک مدل خطی N_r بعدی و سطح بازتابش دریک مدل خطی N_{r-1} بعدی قرار گیرند. از دیدگاه شهودی $r = 3$ است. پس نتایج حاصل شده از آنالیز *Maloney* و *Wandell* معتبرند اگر سطح بازتابش دریک مدل دو بعدی قرار گیرند که از نظر طبیعی چنین حالتی رخ نمی‌دهد. در الگوریتم *Buchsbaum* اشعه تابش و سطح بازتابش دریک مدل خطی N_r بعدی قرار می‌گیرند. در این الگوریتم فرض می‌شود که میانگین فاصله‌ای سطح بازتابش بازی اهمه مقادیر ثابت است. این فرض را *gray world* گوییم. که نتیجه می‌دهد که *color constancy* برای تصاویری که میانگین بازتابش آنها از میانگین فرض شده متفاوت باشد بایستی تفکیک شود. در اینجا *color constancy* را به عنوان یک مسئله برآورد آماری مطالعه می‌کیم. برای این منظور *Vhrel* و *Trussell* از روش ماکریم درستنماهی که ارتباط نزدیکی با برآورد *MAP* دارد استفاده کردند. *Freeman* و *Brainnad* *MMSE* محک را برآورد اشعه تابش به خدمت گرفته‌اند. در اینجا روش *MLM* را به کار گرفته و مزیت آن را به روشهای

فوق بیان می کنیم. برای تعیین الگوریتم باقیمانده احتمال پیشین را انتخاب کرده و سپس امیدریاضی تابع زیان را می نیمم کنیم.

۳.۴ تعیین تابع احتمال پیشین

یک مدل خطی ۳ بعدی برای سطح بازتابش و پرتوتابش به دست می آوریم. برای به دست آوردن چنین مدلی برروی داده های $K.L.Kelly$ تکنیک تجزیه و تحلیل مؤلفه های اصلی را به کار می گیریم. سپس بهترین مدل بازتابش وزنی رابرای هرسطح می یابیم. فرض می کنیم مدل خطی ۳ بعدی برای پرتوی تابش همان مدل CIE برای نور خورشید باشد. برای تولید توزیع پیشین پرتوهای خورشید به رنگهای به هم مرتبط که دمایشان مطابق با توزیع نرمال بامیانگین $K = 6500^{\circ}$ و انحراف استاندارد 4000° است تابیده می شود. در هر تابع محاسبه کرده و نتیجه را با توزیع نرمال مقایسه می کنیم. توجه داریم که برای صحیح بودن نتایج نیازی به نرمال بودن توزیع پیشین نیست.

۴.۴ تقریبی برای جرم احتمال موضعی

می خواهیم برآوردگر MLM را برای مسئله *color constancy* بکارگیریم. رسک بیزی که در جستجوی آن برای مینیمم کردن هستیم با انتگرال گیری روی همه فضای پارامتر منظری معادله (۴) بدست می آید. با ترکیب روابط (۱)، (۲) و (۴) داریم

$$R(\hat{\theta}|y) = -C \int \exp\left\{-\frac{\tau}{\gamma} |K_n^{-\frac{1}{\gamma}}(r-f(x))|^2\right\} P_\theta(\theta) \exp\left\{-\frac{\mu}{\gamma} |(K_L^{-\frac{1}{\gamma}}(\theta-\hat{\theta}))|^2\right\} d\theta \quad (8)$$

که در این رابطه $\frac{1}{\gamma} K_n$ کواریانس اغتشاش مشاهدات است که K_n ثابت فرض شده و τ زیاد می شود و بطور مشابه $\frac{1}{\gamma} K_L$ کواریانس تابع زیان موضعی است که K_L به شکل تابع زیان بستگی دارد و μ بزرگ می شود. در اینجا با این فرض مطالعه انجام می شود که $\frac{\mu}{\tau}$ کوچک باشد.

قضیه

انتگرال

$$I(\tau) = \int \exp\{-\tau\phi(\theta)\} g(\theta) d\theta \quad (9)$$

بطور تقریبی و با افزایش τ برابر است با

$$I(\tau) \approx \frac{e^{-\tau\phi(\theta_*)}}{\sqrt{|det(\phi_{\theta_i\theta_j}(\theta_*))|}} = \left(\frac{2\pi}{\tau}\right)^{\frac{n}{2}} g(\theta_*) \quad (10)$$

که در آن θ مینیمم $\phi(\theta)$ است.
با جایگذاری $\theta = P(\theta)$ و $g(\theta) =$

$$\phi(\theta) = \frac{1}{2} |K_n^{-\frac{1}{\tau}}(y - f(\theta))|^2 + \frac{\mu}{2\tau} |K_L^{-\frac{1}{\tau}}(\theta - \hat{\theta})|^2 \quad (11)$$

در قضیه فوق داریم

$$R(\hat{\theta}|y) \approx -Ce^{\{-\tau(\frac{1}{\tau}|K_n^{-\frac{1}{\tau}}(y-f(\theta_0))|^2 + \frac{\mu}{2\tau}|K_L^{-\frac{1}{\tau}}(\theta_0-\hat{\theta})|^2)\}} \times$$

$$\frac{P_\theta(\theta_0)}{\sqrt{|det(\phi_{\theta_i \theta_j}(\theta_0))|}} \quad (12)$$

با دوبار مشتق گیری از طرفین رابطه (11) داریم

$$\phi_{\theta_i \theta_j}(\theta_0) = f_i'^T K_n^{-1} f_j' - (y - f(\theta_0))^T K_n^{-1} f_{ij}'' + \frac{\mu}{\tau} \{K_L^{-1}\}_{ij} \quad (13)$$

که در آن

$$f_i' = \frac{\partial f(\theta)}{\partial \theta_i} |_{\theta=\theta_0}$$

$$f_{ij}'' = \frac{\partial^2 f(\theta)}{\partial \theta_i \partial \theta_j} |_{\theta=\theta_0} \quad (14)$$

برای بکارگیری رابطه (12) نیازمند به یافتن یک نقطه θ هستیم که $\phi(\theta)$ معادله (11) را می‌نیمم کند. اگر تنها توجه خود را به محاسبه ریسک در نقطه ماکزیمم موضعی معطوف کنیم در این صورت هردو جمله $\phi(\theta)$ با انتخاب $\theta_0 = \hat{\theta}$ بطور موضعی مینیمم شده‌اند. بنابراین با جایگذاری $\theta_0 = \hat{\theta}$ در معادله (12) برای محاسبه ریسک در نقطه ماکزیمم درستنمایی برآورد تقریباً خوبی خواهیم داشت.
در $\theta_0 = \hat{\theta}$ تفاوت بین ریسک $R(\hat{\theta}|y)$ در معادله (12) و توزیع پسین معادله (1) عبارت

$$\frac{1}{\sqrt{|det(\phi_{\theta_i \theta_j}(\hat{\theta}))|}}$$

است و این نشان می‌دهد که ریسک MLM به همان اندازه که به ارتفاع احتمال بستگی دارد به همان اندازه به عرض هم بستگی دارد.

۵ مقایسه الگوریتمها

چهار نمودار شکل ۵ عملکرد ۴ الگوریتم بازای یک پرتوی تابش را نشان می‌دهد. در هر نمودار پرتوی تابش واقعی با خط تیره و برآورد حاصل با خطهای نقطه چین بطور جداگانه رسم شده‌اند. برای رسم این برآوردها ابتدا مجموعه سطوحی از توزیع پیشین رسم شده و بازای آن پاسخ سنسور محاسبه گردیده است سپس به اندازه یک درصد، اغتشاش به پاسخ سنسور افزوده شده و هر الگوریتم به کار گرفته شده است. همانگونه که در شکل (۵ الف) دیده می‌شود برآوردهای حاصل از الگوریتم *MLM* بسیار نزدیک مقدار واقعی (خط تیره) قرار دارند. برآوردهای حاصل از الگوریتم *gray world* نسبتاً صحیح است اما دارای انحرافات زیادی است (شکل ۵ ب). برآورد *MAP* همچنان که انتظار می‌رفت از اطلاعات مرتبط با توزیع پیشین صرف نظر کرده است (شکل ۵ ج). یکی دیگر از الگوریتم‌هایی که مورد استفاده قرار می‌گیرد، الگوریتم *subspace* است. همانطور که در شکل (۵ د) می‌بینیم هیچ تضمینی برای صحیح بودن جواب حاصل از این الگوریتم وجود ندارد چون برآورد بسیار مغوش است.

به علاوه عملکرد الگوریتمها تحت برخی خطاها فرض شده در توزیع پیشین نیز بررسی شده‌اند. برای این منظور پرتوی تابش را در حالیکه سطح باز تابش، پرتوی باز تابش را منعکس می‌کند تغییر می‌دهیم در اینجا از پرتو تابش، نور خورشید K^{4000} و K^{10000} همراه با نور خورشید K^{6500} استفاده کردایم. نتایج با محاسبه ریشه میانگین مربع خطأ (*RMSE*) بین پرتوی تابش و برآورد حاصل به دست آمده‌اند. شکل (۶ الف) این خطا را برای چهار الگوریتم فوق و سه پرتوی تابش نشان می‌دهد. برای هر سه نوع پرتو الگوریتم *MLM* از سایر الگوریتم‌ها بهتر عمل می‌کنند و برآورد حاصل از آن مستقل از نوع پرتواست.

در قسمت بعد از نور خورشید K^{6500} استفاده شده ولی در انتخاب سطوح اربیی وجود دارد، این اربیی ممکن استفاده به دو صورت باشد: یکی، میانگین سطوح راطوری اربی کردیم که میانگین پاسخ‌های سنسور تحت اشعه K^{6500} همان پاسخ‌های سنسور تحت اشعه K^{4000} و بدون اربی سطح باشد. دیگری، میانگین سطوح را طوری اربی کردیم که میانگین پاسخ‌های سنسور تحت اشعه K^{6500} همان پاسخ‌های سنسور تحت اشعه K^{10000} و بدون اربی سطح باشد. با این فرآیند قصد داریم تا پایداری الگوریتم‌ها را تحت تخطی از شرایط *gray world* آزمون کنیم. شکل (۶ ب) *RMSE* را برای حالاتی که سطوح نااربیند با دوگونه اربی نشان می‌دهد. اربی سطوح در عملکرد الگوریتم‌ها تأثیر دارد. چنان‌که از شکل پیداست الگوریتم *MLM* به مرتب نسبت به سایر الگوریتم‌ها نسبت به تخطی از فرضیات *gray world* پایدارتر است.

۶ نتایج

ما این مسئله را از دیدگاه تئوری تصمیم بیزی بررسی کردیم. ثابت کردیم هیچ کدام از دو برآوردهای MAP و $MMSE$ که معمولاً مورد استفاده قرار می‌گیرند مناسب نیست. برای از بین بردن نقایص این دو برآوردهای یک برآوردهای جدید به نام MLM معرفی کردیم. این برآوردهای یک انتگرال را روی چگالی احتمال پسین موضعی بیشینه می‌کند.تابع زیان این برآوردهای متناسب با نوع مسئله‌ای که با آن روبرو هستیم در نظر گرفته می‌شود و محتمل‌ترین برآورد تقریباً درست را ارائه می‌دهد. با یک مثال نشان دادیم که برآوردهای MLM یک برآورد مناسب (بهینه) به دست می‌آورد در حالیکه برآوردهای MAP و $MMSE$ این گونه نبودند. همچنین الگوریتم MLM را با برآورد MAP مقایسه کرده و نشان دادیم که MLM به طور معنی داری حتی تحت تخطی از شرایط مفروض، بهتر از سایرین عمل می‌کند. از این‌رو معتقدیم که برآوردهای معرفی شده باعث بهبودی عملکرد الگوریتم‌های رایج می‌گردد.

مراجع

- [1] D.H.Brainard and W.T.Freeman (1994), *Bayesian method for recovering surface and illuminant properties from photosensor responses*, In proc. SPIE, volume 2179.
- [2] D.H.Brainard and B.A.Wandell and E.J.Chichilnisky (1993), *color constancy*, Current dir.in psychl. science.
- [3] R.A.Fisher. (1959), *statistical methods and scientific inference*, Hafner.
- [4] K.L.Kelly,K.S.Gibson, and D.Nickerson. (1943), *Tristimulus specification of the munsel of color from spectrophotometric measurements*, J.Op.Soc.Am.
- [5] L.T.Maloney (1986), *Evaluation of linear models of surface spectral reflectance with small numbers of parameters*, J.Op.Soc.Am.

تعیین مدل خوشبندی احتمالاتی براساس معیار اطلاع بیزی

محمد قاسم وحیدی‌اصل^۱، محسن محمدزاده^۲، محمد قربانی^۳

P11192

۱ گروه آمار، دانشگاه شهید بهشتی

۲ گروه آمار، دانشگاه تربیت مدرس

۳ گروه آمار، دانشگاه تبریز

چکیده: یکی از مسائل مهم در تحلیل داده‌های چند متغیره، پیدا کردن ارتباط بین آنهاست. ساده‌ترین روش کشف رابطه موجود بین داده‌ها، رسم نمودار پراکنش است. در بسیاری از مفاهیم پژوهشی، ژنتیکی و غیره، یکی از مسائل مهم، خوشبندی داده‌ها به گروههای همگن است. روش‌های ابتکاری مختلفی از جمله، روش‌های سلسه مراتبی با ماکسیمم کردن تشابهات درون‌گروهی، داده‌ها را به خوش‌هایی افزای می‌کنند. بدلیل وابسته بودن این روش‌ها به تعریف فاصله بین دو خوش و همچنین انتخاب یک مقدار آستانه برای تعیین تعداد خوش‌ها، محققان در انتخاب بهترین معیاری که تشابهات درون‌گروهی را ماکسیمم نماید با مشکل مواجه می‌باشند. در این مقاله روش «انتخاب مدل برای خوشبندی احتمالاتی با استفاده از معیار اطلاع بیزی^۱ (BIC)» مورد بررسی قرار می‌گیرد که در آن، فرضهای مختلف برای معیار تشابه نقشی ندارند و با تجزیه طیفی ماتریس کوواریانس می‌توان معیارهای ساده‌ای برای مشخص کردن حجم، شکل و جهت خوش‌ها به دست آورد. بعلاوه با استفاده از معیار BIC می‌توان بهترین مدل خوشبندی را انتخاب نمود.

کلید واژه: خوشبندی کردن، مدل‌های آمیخته، الگوریتم EM، تجزیه طیفی، معیار BIC.

۱ مقدمه

یکی از مسائل مهم در بسیاری از مطالعات ژنتیکی و پژوهشی، خوشبندی داده‌ها است، که در آن N مشاهده با M ویژگی به g گروه همگن افزای می‌شوند بطوریکه تشابهات درون گروهها ماکسیمم گردند. روش‌های مختلفی از جمله، روش‌های سلسه مراتبی برای خوشبندی کردن داده‌ها به کار می‌روند که در تعریف «فاصله بین دو خوش» باهم تفاوت دارند. در روش «خوشبندی با اتصال منفرد»، فاصله بین نزدیکترین عضوهای دو گروه

^۱ Bayesian Information Criterion

بعنوان فاصله بین دو خوشة تعریف می‌شود، در روش «خوشه بندی با اتصال کامل»، فاصله بین دورترین زوجها و در روش اتصال میانگین، متوسط فاصله بین تمام زوجهای موجود در داخل دو خوشه به عنوان معیار خوشه‌بندی مورد استفاده واقع می‌شود (هارتیگان، ۱۹۷۵). چون این روشها مدل خاصی نمی‌باشند، استنباط از نمونه به جامعه امکان‌پذیر نیست و تعداد خوشه‌ها نیز به صورت ابتکاری با تعریف آستانه‌ای دلخواه تعیین می‌گردد. بنابراین لازم است روش خوشه‌بندی حتی الامکان مبتنی بر سلیقه محقق نبوده و براساس مدل یا توزیع احتمالی باشد تا بتوان در مورد آن استنباط آماری انجام داد. در این صورت شرایطی فراهم خواهد شد که بتوان با تجزیه طیفی ماتریس کوواریانس عناصر داخل خوشه‌ها، شکل، حجم و جهت خوشه‌ها را مشخص نمود. عمدهاً مجموعه مشاهدات تحت بررسی، همگی از یک جامعه خاص نمی‌باشند، برای تشخیص این که هر مشاهده از کدام جامعه آمده است، منطقی است فرض شود که هر مشاهده براساس ویژگی‌ها و خصوصیات دارای توزیع احتمال خاصی است. بنابراین جامعه‌ای مرکب از چند زیرجامعه، دارای توزیع احتمالی آمیخته از توزیع‌های احتمال زیر جامعه‌ها خواهد بود، که در حالت کلی بشکل

$$f(x|\psi) = \lambda_1 f_1(x|\theta_1) + \dots + \lambda_g f_g(x|\theta_g)$$

است، که در آن برای $g = 1, \dots, G$ $f_j(\cdot)$ تابع چگالی مولفه‌ها، $\lambda_j \leq 1 < 0$ و $\sum_{j=1}^g \lambda_j = 1$ ، $\lambda = (\lambda_1, \dots, \lambda_g)$ ، $\theta = (\theta_1, \dots, \theta_g)$ و $\psi = (\lambda, \theta)$. از توابع چگالی آمیخته معروف که خیلی کاربرد دارد می‌توان به

$$\begin{aligned} f(x|\psi) &= \lambda\phi(x|\mu_1, \sigma_1^2) + (1-\lambda)\phi(x|\mu_2, \sigma_2^2) \\ \phi(x|\mu_j, \sigma_j^2) &= N(\mu_j, \sigma_j^2) \quad j = 1, 2 \end{aligned}$$

اشارة نمود. بنابراین خوشه‌بندی را می‌توان معادل با تفکیک توزیع آمیخته به مولفه‌های ساده دانست. قبل از بیان خوشه‌بندی براساس مدل می‌بایست یک سری فرضیات اساسی روی توزیع مولفه‌ها داشته باشیم. به عنوان مثال، یکی از این فرضها این است که توزیع هر مولفه تک مدلی است. بنابراین هدف از خوشه‌بندی، تجزیه مولفه‌های چند بعدی مبهم و آمیخته به مولفه‌های ساده تک مدلی است. به همین منظور در این مقاله برای خوشه‌بندی مشاهدات از مدل‌های آمیخته استفاده می‌شود تا بتوان تمام خصوصیات مشاهدات را در خوشه‌بندی کردن بکار گرفت.

۱.۱ تعیین معیارهای خوشه‌بندی براساس مدل

فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی از جامعه Π با زیر جامعه‌های $\Pi_1, \Pi_2, \dots, \Pi_g$ باشند و Σ_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 $\lambda_j > 0$ و $\sum_{j=1}^g \lambda_j = 1$ ، $\psi = (\lambda_1, \dots, \lambda_g, \theta_1, \dots, \theta_g)$ ، $\theta_j = (\mu_j, \Sigma_j)$ است و تابع درستنایی 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 خوشبندی حاصل از ماقسیم کردن تابع درستنایی

$$L(x | C) = \prod_{j=1}^g \lambda_j^{n_j} \prod_{X_i \in C_j} f(x_i | \theta_j) \quad (3)$$

معادل با معیار حاصل از ماقسیم کردن (2) خواهد بود (مکلن، ۱۹۸۲ و فرالی و رافتری، ۱۹۹۸). معمولاً مولفه اصلی X_i ها معلوم نیستند و برای مشخص کردن مولفه اصلی X_i ، متغیرهای گروه‌بندی Z_{ij} به صورت زیر تعریف می‌شود

$$Z_{ij} = \begin{cases} 1 & X_i \in \Pi_j \\ 0 & X_i \notin \Pi_j \end{cases}$$

براساس مشخصه‌های گروه‌بندی لگاریتم تابع درستنایی 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، در این صورت عناصر خوشه‌زام به صورت $\{j; z_{ij} > z_{ij'} \quad j \neq j'\}$ 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 z_{ij} تلقی نمود. چون

$$E(Z_{ij}) = P(X_i \in \Pi_j) = \frac{\lambda_j f(x_i | \theta_j)}{\sum_{j=1}^g \lambda_j f(x_i | \theta_j)}$$

لازم است برای هر j پارامترهای نامعلوم λ_j و θ_j برآورد شوند. برآورد این پارامترها بروش ماقسیم درستنایی مستلزم استفاده از روش‌های عددی است.

۱.۱.۱ برآورد پارامترها با استفاده از الگوریتم EM

الگوریتم امید ریاضی و ماسکیم‌سازی^۱ (EM) روشی کاربردی در محاسبات تکراری، برای به دست آوردن برآورد ماسکیم درستنایی پارامترها در توزیع‌های آمیخته است که در مسائل گوناگونی مانند داده‌های ناکامل، داده‌های گم شده وغیره کاربرد دارد. این الگوریتم در شرایطی به ویژه در هنگام مواجهه شدن با داده‌های ناکامل بهتر از الگوریتم نیوتن-رافسون کار می‌کند (مکلن و کریشنان، ۱۹۹۷). فرض کنید (ψ) مقدار اولیه ψ باشد، در این صورت مراحل الگوریتم EM به صورت زیر خواهد بود:

مرحله E : محاسبه امید ریاضی لگاریتم تابع درستنایی در نقطه (ψ) 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, X_2, \dots, X_n داده‌های ناکامل و $(X_i, Z_i) = Y_i$ داده‌های کامل در الگوریتم EM باشند. در این صورت برآورد پارامترهای (۴) با استفاده از این الگوریتم به صورت زیر خواهد بود.

$$\begin{aligned} \lambda_j^{(k+1)} &= \frac{1}{n} \sum_{i=1}^n z_{ij}^{(k)} \\ \mu_i^{(k+1)} &= \sum_{i=1}^n \frac{z_{ij}^{(k)} X_i}{\sum_{i=1}^n z_{ij}^{(k)}} \\ (\Sigma)^{(k+1)} &= \frac{1}{n} \sum_{j=1}^g \sum_{i=1}^n z_{ij}^{(k)} (X_i - \mu_j^{(k+1)}) (X_i - \mu_j^{(k+1)})' \\ z_{ij}^{(k)} &= \frac{\lambda_j^{(k)} f(x_i | \theta_j^{(k)})}{\sum_{j=1}^g \lambda_j^{(k)} f(x_i | \theta_j^{(k)})} \end{aligned}$$

با z_{ij} برآورد شده به وسیله الگوریتم EM می‌توان تحلیل خوش‌ای را مناسب با بیشترین مقدار $z_{ij}^{(k)}$ انجام داد (مکلن و کریشنان، ۱۹۹۷). اما همان طوری که بدان اشاره شد در علومی چون پزشکی و فتوگرافی یکی از اهداف خوش‌بندی تعیین حجم، شکل و جهت

^۱ Expectation-Maximization Algorithm

خوشه‌ها است. تعیین حجم، شکل و جهت خوشه‌ها با استفاده از تجزیه طیفی ماتریس کوواریانس نخستین بار توسط بانفیلد و رافتری (۱۹۹۳) بیان شد و سرانجام گیلز‌سیلوکس و گوورت (۱۹۹۵) آنرا تعمیم دادند. معیارهای تحلیل خوشه‌ای توسط فرالی (۱۹۹۹) در نرم‌افزار *MCLUST* گنجانده شد.

۱. تعیین حجم، شکل و جهت خوشه‌ها

تابع درستنماهی (۳) را می‌توان به صورت

$$L(\psi|Z, X) = \prod_{j=1}^g \lambda_j^{n_j} |\Sigma_j|^{-\frac{n_j}{2}} \exp \left\{ -\frac{1}{2} \sum_{X_i \in C_j} (X_i - \mu_j)' \Sigma_j^{-1} (X_i - \mu_j) \right\} \quad (5)$$

نوشت، که در آن C_j تعداد مشاهدات X_{ij} است که توسط Z_i به زیرجامعه Π_j تخصیص یافته است و n_j تعداد مشاهدات در C_j است. فرض کنید Σ_j , Σ_j , $j = 1, \dots, g$, ماتریس کوواریانس مولفه زام باشد. در این صورت تجزیه طیفی ماتریس متقابله Σ_j بصورت

$$\Sigma_j = D_j B_j D_j' = \sum_{k=1}^j \xi_j e_j e_j'$$

است، که در آن $D_j = (e_1, \dots, e_j)$ ماتریس متعامد از بردارهای ویژه یکه و B_j ماتریس قطری از مقادیر ویژه می‌باشد. تجزیه طیفی فوق را می‌توان به صورت $\Sigma_j = \xi_j D_j A_j D_j'$ نوشت که در آن $|A_j|^{\frac{1}{d}}$ و A_j ماتریس قطری از مقادیر ویژه نرمال شده می‌باشد. پارامتر ξ_j ، مشخص کننده حجم مولفه زام، D_j جهت آن و A_j شکل آن می‌باشد. فرض کنید برخی از این کمیت‌ها (حجم، شکل و جهت) بین خوشه‌ها متفاوت باشند. در این صورت می‌توان معیارهای ساده‌ای برای خوشه‌بندی به دست آورد که در اکثر علوم کاربرد دارند (بن اسماعیل و سیلوکس، ۱۹۹۶). در مدل‌های زیر علامت اختصاری E ، $(Equal)$ ، F ، $(Fixed)$ و V ، $(Vary)$ نشان دهنده ثابت بودن کمیت و متغیر بودن کمیت است.

در مدل EEE که با DAD' نیز نشان داده می‌شود حجم، شکل و جهت کلیه خوشه‌ها یکسان است و به راحتی ثابت می‌شود تخصیص بهینه \hat{Z} ، معادل با مینیمم کردن $|W|$ است.

در مدل VEE شکل و جهت خوشه‌ها ثابت ولی حجم خوشه‌ها متغیر است. فرض کنید $C = DAD'$ با جایگذاری $\Sigma_j = \xi_j C$ در (۵) داریم،

$$\log L(X|\psi) \propto -\frac{d}{2} \sum_{j=1}^g n_j \log(\xi_j) - \frac{1}{2} \sum_{j=1}^g \frac{1}{\xi_j} \text{tr}(W_j C^{-1})$$

که در آن $\hat{\lambda}_j = \frac{1}{n} \sum_{j=1}^g W_j$ و $\hat{\mu}_j = \frac{n_j}{n} \sum_{j=1}^g W_j$ ، به ترتیب برآوردهای ماکسیمم درستنمایی λ_j و μ_j و Σ هستند و $W_j = \sum_i (X_{ij} - \bar{X}_j)(X_{ij} - \bar{X}_j)'$ است. ماکسیمم کردنتابع درستنمایی (۵) با درنظرگرفتن مدل "VEE" هم ارز مینیمم کردن

$$F(\xi_1, \dots, \xi_g, C) = \sum_{j=1}^g \frac{1}{\xi_j} (tr(W_j C^{-1})) + d \sum_{j=1}^g n_j \log(\xi_j)$$

است و نشان داده می‌شود که مینیمم کردن آن باید به روش تکراری (الگوریتم EM) صورت گیرد. با فرض ثابت بودن ماتریس C مقداری از ξ_j که $F(\xi_1, \dots, \xi_g, C)$ را مینیمم می‌کند عبارت است از

$$\xi_j = \frac{tr(W_j C^{-1})}{dn_j} \quad (6)$$

فرض کنید ξ_j ها ثابت باشند در این صورت مینیمم کردن $F(\xi_1, \dots, \xi_g, C)$ ، معادل مینیمم کردن $tr(\sum_{j=1}^g (\frac{1}{\xi_j} W_j) C^{-1})$ است. سیلوکس و گورت (۱۹۹۵) نشان دادند ماتریس متقارن $M = \frac{Q}{|Q|^{\frac{1}{d}}}$ ، که در آن Q ماتریس متقارن و معین مثبت است و $tr(QM^{-1})$ را مینیمم می‌کند بنابراین ماتریسی که $(tr(\sum_{j=1}^g (\frac{1}{\xi_j} W_j) C^{-1}))$ را مینیمم کند بصورت

$$C = \frac{\sum_{j=1}^g (\frac{1}{\xi_j} W_j)}{|\sum_{j=1}^g (\frac{1}{\xi_j} W_j)|^{\frac{1}{d}}} \quad (7)$$

خواهد بود با توجه به معادلات (۶) و (۷) نمی‌توان فرم بسته‌ای از یک معیار خوشه‌بندی ارائه نمود و باستی از الگوریتم EM استفاده شود. مدل EVV (مدل EVV : با قرار دادن $C_j = D_j A_j D'_j$) مینیمم کردن (۵) معادل با مینیمم کردن

$$F(\xi, C_1, \dots, C_g) = \sum_{j=1}^g \frac{1}{\xi} (tr(W_j C_j^{-1})) + nd \log(\xi)$$

است. بنابراین $C_j = \frac{W_j}{|W_j|^{\frac{1}{d}}}$ و از طرفی ξ که تابع F را مینیمم می‌کند عبارت است از

$$\xi = \frac{\sum_{j=1}^g tr(W_j C_j^{-1})}{nd}$$

با جایگذاری C_j در معادله فوق داریم $\xi = \frac{\sum_{j=1}^g (|W_j|^{\frac{1}{d}})}{nd}$. بنابراین برای مدل "EVV" افزایی که $(|W_j|^{\frac{1}{d}})$ را مینیمم می‌کند، معادل با ماکسیمم درستنمایی است.

در مدل‌های قبلی شکل خوش‌ها چه ثابت و چه متغیر، بیضوی هستند. اکنون دو مدل که شکل خوش‌های آنها کروی است، معرفی می‌شوند. در مدل $\xi I = \sum_j \xi_j$ شکل خوش‌ها کروی و حجم آنها با هم مساوی است. با فرض این مدل، ماکسیمم کردن (۵) همارز مینیمم کردن

$$F(\xi) = \frac{1}{\xi} \operatorname{tr}(W) + nd \log \xi$$

می‌باشد و $F(\xi)$ تابع (ξ) را مینیمم می‌کند. با جایگذاری ξ در F داریم، $F(\xi) = \frac{\operatorname{tr}(W)}{nd} + nd \log(\frac{\operatorname{tr}(W)}{nd})$. پس به جای مینیمم کردن F می‌توان $nd \log(\operatorname{tr}(W))$ را مینیمم کرد، که کاربردی‌ترین و قدیمی‌ترین معیار خوش‌بندی است.

در مدل $\xi_j I = \sum_j \xi_j$ شکل خوش‌ها کروی ولی حجم آنها با هم متفاوت است. در این صورت ماکسیمم کردن (۵) همارز مینیمم کردن

$$F(\xi_1, \dots, \xi_g) = \sum_{j=1}^g \frac{1}{\xi_j} \operatorname{tr}(W_j) + d \sum_{j=1}^g n_j \log \xi_j$$

است. با مشتق گیری از F نسبت به ξ_j و مساوی صفر قرار دادن آن داریم

$$\xi_j = \frac{\operatorname{tr}(W_j)}{dn_j}$$

با جایگذاری ξ_j در F خواهیم داشت،

$$F(\xi_1, \dots, \xi_g) = \sum_{j=1}^g dn_j + d \sum_{j=1}^g n_j \log(\frac{\operatorname{tr}(W_j)}{dn_j}).$$

بنابراین مینیمم کردن F همارز مینیمم کردن $\sum_{j=1}^g n_j \log(\frac{\operatorname{tr}(W_j)}{n_j})$ می‌باشد. معیارهای خوش‌بندی براساس مدل را می‌توان به طور خلاصه در جدول ۱ ملاحظه نمود.

۲ انتخاب بهترین مدل با استفاده از معیار BIC

فرض کنید X_1, \dots, X_n یک نمونه تصادفی از توزیع آمیخته با تابع چگالی احتمال (۱) باشد. در استنباط بیزی برای مدل‌های آمیخته گاوی، یک روش ساده برای انتخاب بهترین مدل و تعیین تعداد مؤلفه‌ها محاسبه عامل بیزی^۱ است. عامل بیزی B_{10} برای مدل M_1 در مقابل مدل M_0 به صورت

$$B_{10} = \frac{P(\underline{X}|M_1)}{P(\underline{X}|M_0)}$$

^۱ Bayese Factor

معیار خوشبندی	جهت	شكل	حجم	شكل خوشها	مدل	علامت اختصاری
$tr(W)$	موجود نیست مساوی مساوی	کروی			ξI	$E I$
$\sum_{j=1}^g n_j \log(\frac{tr(W_j)}{dn_j})$	موجود نیست مساوی متغیر	کروی			$\xi_j I$	$V I$
$ W $	مساوی مساوی مساوی	بیضوی			$\xi D A D'$	$E E E$
$\sum_{j=1}^g n_j \log(\frac{ W_j }{n_j})$	متغیر متغیر ثابت	بیضوی			$\xi_j D_j A_j D'_j$	$V V V$
از الگوریتم	مساوی متغیر	بیضوی			$\xi_j D A D'$	$V E E$
از الگوریتم	متغیر ثابت	بیضوی			$\xi_j D_j A D'_j$	$V F V$
$\sum_{j=1}^g W ^\frac{1}{d}$	متغیر مساوی	بیضوی			$\xi D_j A_j D'_j$	$E V V$

جدول ۱ : معیارهای خوشبندی براساس مدل‌های آمیخته با مولفه‌های نرمال چند متغیره

تعریف می‌شود، که در آن

$$P(\underline{X}|M_j) = \int P(\underline{X}|\psi_j, M_j) P(\psi_j|M_j) d\psi_j. \quad (8)$$

و ψ بردار پارامتری (λ_j, θ_j) در مدل M_j و $P(\psi_j|M_j)$ تابع چگالی پیشین است. بنابراین عامل بیزی شانس پسین یک مدل دربرابر مدل دیگر می‌باشد (کاس و رافتري، ۱۹۹۵). هنگامی که از الگوریتم EM برای برآورد ماکسیمم درستنمایی پارامترهای مدل آمیخته استفاده شود یک روش کلاسیک برای تقریب (۸) استفاده از معیار اطلاع بیزی (BIC) است. این تقریب عبارت است از

$$2 \log P(\underline{X}|M_j) = 2 L_M(\underline{X}|\hat{\psi}) - m_M \ln(n) = BIC(M, g).$$

که در آن $L_M(\underline{X}|\hat{\psi})$ لگاریتم تابع درستنمایی و m_M تعداد پارامترهایی است که با در نظر گرفتن مدل M می‌بایست برآورد شوند. براین اساس مدلی که دارای بیشترین مقدار BIC باشد، بهترین خواهد بود. در حالیکه براساس عامل بیزی اگر $1 > B_1$ باشد آنگاه مدل M_1 بهتر از M_0 است

س اساردی دنبه شوخی اهشورحیرشته ارد: قی بنز مل گی اهدهداد ندرکی دنبه شوخ بل اشه (۹۷۹۱) (ایدرام رد اهدهداد نیا . دوشی م دلفتسا ۶۳۹۱) (رشیف قی بنز مل گی اهدهداد زا لدم هس زا ق بنز مل گ ۵۱ ° ارد اهگ بریلاگ و اهگ رسماک ضرعه و لوط مل ماشه ک مذاهده داد دنشابی می‌یابنیجرو و گ نرماگنر، رادراخ عوز

۱.۲ کاربرد الگوریتم EM در خوشه‌بندی براساس مدل

همانطوری که بیان شد در مرحله E الگوریتم EM ماتریس مشخصه گروه‌بندی، $\{Z_{ij}\}$ ، محاسبه می‌شود که در آن Z_{ij} ، برآورد احتمال شرطی تعلق مشاهده i ام به گروه k با انتخاب مقادیر اولیه برای پارامترها یا برآورد آنها است. سپس در «مرحله M » برآورد ماقسیم درستنمایی پارامترها مشخص می‌شود.

برای این منظور از نرم‌افزار $Mclust$ استفاده شده است. براساس یک سری محاسبات با $MCLUST$ داریم $z_{1,1} = 5.8e^{-0.22}$ ، $z_{1,2} = 3.17e^{-0.42}$ و $z_{1,3} = 1.96e^{-0.09}$. چون $z_{1,1}$ از $z_{1,2}$ و $z_{1,3}$ بزرگتر است مشاهده ۱۱۹ ام در خوشه اول قرار می‌گیرد، در حالی که مشاهده ۱۱۹ ام به خوشه سوم تعلق می‌باشد. به همین ترتیب براساس z برآورد شده، در مورد خوشه اصلی سایر مشاهدات نیز می‌توان تصمیم گرفت. براین اساس مشاهدات به گروهی تعلق می‌گیرند که دارای z بیشتری باشند. بعد از خوشه‌بندی کردن داده‌ها، می‌توان برآورد ماقسیم درستنمایی پارامترها را به دست آورد که به صورت زیر می‌باشد.

$$\hat{\mu} = \begin{pmatrix} 5.006 & 5.942231 & 6.574623 \\ 2.428 & 2.760757 & 2.980792 \\ 1.462 & 4.258718 & 5.539016 \\ 0.264 & 1.319203 & 2.024933 \end{pmatrix}$$

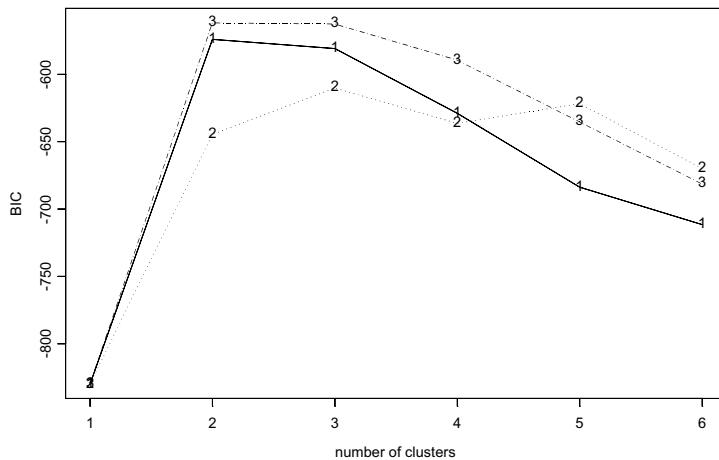
$$\hat{\Sigma} = \begin{pmatrix} 0.26393480 & 0.08984982 & 0.16965778 & 0.03933808 \\ 0.08984982 & 0.11194776 & 0.05112076 & 0.02997824 \\ 0.16965778 & 0.05112076 & 0.18653476 & 0.04197325 \\ 0.03933808 & 0.02997834 & 0.04197235 & 0.03971195 \end{pmatrix}$$

$$\hat{\lambda} = (0.32332323 \quad 0.3296191 \quad 0.3370475)$$

پس از خوشه‌بندی کردن داده‌ها توسط الگوریتم EM و برآورد پارامترها، به انتخاب بهترین مدل با استفاده از معیار BIC می‌پردازیم.

با توجه به شکل ۱ و جدول ۱ مقدار BIC برای مدل خوشه‌بندی $"VEV"$ با تعداد دو خوشه بیشتر از سایر حالت‌هاست. لذا این مدل با فرض دو خوشه بهترین مدل براساس معیار BIC در بین سه مدل فوق با کلیه حالات درنظر گرفته شده برای تعداد خوشه‌ها می‌باشد.

مدل	۱	۲	۳	۴	۵	۶
VVV	-۸۲۹.۹۷۸۲	-۵۷۴.۰۱۷۸	-۵۸۰.۸۳۸۹	-۶۲۸.۹۵۶۴	-۶۸۳.۸۱۱۴	-۷۱۱.۵۶۵۷
E EV	-۸۲۹.۹۷۸۲	-۶۴۴.۵۹۹۷	-۶۱۰.۰۸۳۶	-۶۴۵.۹۹۵۰	-۶۲۱.۶۹۰۱	-۶۶۹.۷۰۶۹
VEV	-۸۹۲.۹۷۸۲	-۵۶۱.۲۲۸۵	-۵۶۲.۵۵۰۷	-۵۸۹.۳۵۱۰	-۶۳۵.۲۰۵۱	-۶۸۱.۲۹۷۶



شکل ۱ : نمودار برای سه مدل (۱) ، (۲) و (۳) :

۳ بحث و تیجه‌گیری

در این مقاله روش خوشه‌بندی براساس مدل با مولفه‌های نرمال چند متغیره مورد بررسی قرار گرفت و نشان داده شد که این روش معادل با برآورد پارامترهای مولفه‌ها می‌باشد. سپس از الگوریتم EM برای برآورد پارامترها استفاده شد. همچنین بیان شد که در این روش فرضیات دلخواه اشخاص در مورد معیار تشابه نقشی ندارد و با تجزیه طیفی ماتریس کوواریانس می‌توان حجم، شکل و جهت خوشه‌ها را به دست آورد. بعلاوه می‌توان با استفاده از معیار BIC بهترین مدل را تعیین نمود که هم‌ارز بهترین مقدار ممکن برای تعداد خوشه‌ها نیز می‌باشد.

مراجع

- Banfield, J. D. and Raftery, A. E. (1993), Model-Based Gaussian and Non-Gaussian Clustering. *Biometrics*, 49, 803-821.
- Bensmail, H. and Celeux, G. (1996), Regularized Gaussian Discriminant Analysis Through Eigenvvalue Decomposition. *J. Amer. Stat. Assoc.*, 91, 1743-1748.
- Celeux, G. and Govert, G. (1995), Gaussian Parsimonious Clustering Models. *Pattern Recognition*, 28, 781-793.
- Fraly, C. and Raftery, A. E. (1998), How Many Clusters? Which Clustering Method? Answer via Model-Based Cluster Analysis. *Techincal Report*, No, 329. Seattle: Department of Statistics, University of Washington.
- Fraley, C. and Raftery, A. E. (1999), MCLUST: Software for Model-Based Cluster Analysis, *J. Classification*, 16, 297-306.
- Hartigan, J. A. (1975), *Clustering Algorithms*. Wiley, New York.
- Jeff C. F., (1983), On the Convergence of the EM Algorithm, *Annals of Statistcs*, 11, 95-103.
- Kass, R. E. and Raftery, A. E. (1995), Bayes Factors. *J. Amer. Stat. Assoc.*, 90, 773-775.
- Mardia, K. V. and Kent, J. M. (1979), *Multivariate Analysis*, Bibby, academic Press.
- 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, Amesterdam.
- McLachlan, G. J. and Krishnan, T. (1997), *The EM Algorithm and Extensions*. Wiley, New York.

کاربرد گرافهای تصادفی بازه‌ای در تحلیل خوشه‌ای

محمد قاسم وحیدی‌اصل^۱، محسن محمدزاده^۲، فرانک گودرزی^۲

P ۱۳۰۳۹

^۱ گروه آمار، دانشگاه شهید بهشتی

^۲ گروه آمار، دانشگاه تربیت مدرس

چکیده: در این مقاله، مدلی برای گرافهای تصادفی بازه‌ای معرفی شده و خواصی در ارتباط با این مدل عنوان می‌شود. سپس نتایج دقیقی برای توزیع متغیرهای تصادفی مربوط به همبندی این گراف تصادفی بازه‌ای مورد بحث قرار می‌گیرد.

در ادامه خواص مجانبی تعدادی آزمون ترکیبیاتی بالقوه سودمند برای تحلیل خوشه‌ای بر اساس نظریه گرافهای تصادفی بازه‌ای برای داده‌هایی در فضای k بعدی اقلیدسی توصیف شده، سپس مقایسه‌هایی از اثر مجانبی یک رده از این آزمونها ارائه می‌شوند. به عنوان یک مورد از کاربردهای ممکن، آشکارسازی آمیزه‌های توزیعهای احتمال مورد بحث قرار گرفته و مثالی عددی ارائه می‌شود.

واژه‌های کلیدی: گرافهای تصادفی بازه‌ای، تحلیل خوشه‌ای، مدل‌های آمیخته، آشکارسازی.

۱ بخش اول

مدلهای نظری گراف وقتی مفید هستند که قصد ما اکتشاف ساختارهای مجموعه‌های داده‌ها باشد. در تحلیل خوشه‌ای، ماتریس تشابه یک ساختار گراف را روی n شی که باید خوشبندی شوند (که به عنوان رأسهای $1, \dots, n$ تعبیر می‌شود) به وجود می‌آورد؛ دو رأس به وسیله یک یال به هم متصل می‌شوند اگر و تنها اگر شیء‌های متناظر به قدر کافی شباهت داشته باشند. با به کاربردن مفهوم گرافهای تصادفی، می‌توان آماره‌های آزمونی برای آزمون تصادفی بودن خوشه‌های حاصل به دست آورد.

موضوع گرافهای بازه‌ای در تقاطع مسیرهای ترکیبیات و احتمال قرار می‌گیرد. شاینرمن این دو حوزه را به امید آشکارسازی قضیه‌های جدید و کاربردها در هم آمیخته است. نه تنها این حوزه یک شاخه سحرآمیز و زیبای ریاضیات است، بلکه همچنین، کاربردهای زیادی از جمله در مطالعه پلیمریزاسیون، محاسبه کردن خطای تحمل و تحلیل احتمالاتی الگوریتم‌ها دارد.

گرافهای بازه‌ای، رده خاصی از گرافها با چندین خاصیت نظری جالب توجه را تشکیل

می‌دهند. به علاوه، آنها به عنوان مدل‌های مفیدی در خیلی از کاربردها به انضمام تخصیص فرکانس رادیو، زمانبندی پردازشگر، توالی زمانی باستان شناسی و زمان بندی چراغهای راهنمایی به کار می‌روند.

در این مقاله مدلی برای گرافهای تصادفی بازه‌ای معرفی می‌کنیم، این مدل "ایستا" است: برای هر عدد صحیح مثبت n ، تنها یک فضای احتمال برای گرافهای بازه‌ای روی n را برچسب گذاری شده معرفی می‌کنیم.

به عنوان بازه، یک بازهٔ بستهٔ حقیقی را در نظر می‌گیریم. گراف (садه، متناهی) G یک گراف بازه‌ای نامیده می‌شود اگر بتوان یک بازه به هر یک از رأسهای آن منتسب کرد، به طوری که دو رأس مجاور باشند اگر و تنها اگر بازه‌های متناظر با آنها اشتراک داشته باشند. به طور رسمی‌تر، فرض کنید \mathcal{I} مجموعهٔ همهٔ بازه‌های حقیقی را نشان دهد. نمایشی برای گراف بازه‌ای G تابع $f : V(G) \rightarrow \mathcal{I}$ است، به طوری که برای هر زوج از رأسهای متمایز v و w

$$v \sim w \Leftrightarrow f(v) \cap f(w) \neq \emptyset.$$

نماد $w \sim v$ به معنی آن است که v با w مجاور است. گرافهای بازه‌ای، گرافهایی هستند که نمایش بازه‌ای را بپذیرند (گیلمور و هافمن (۱۹۶۴)، گلومبیک (۱۹۸۰) را ببینید).

۲ مدل

مدل یکنواخت $\mathcal{IG}_{n,d}$. فرض کنید X_1, X_2, \dots, X_n متغیر تصادفی مستقل باشند که به طور یکنواخت روی بازهٔ واحد توزیع شده‌اند. یک گراف تصادفی $\mathcal{IG}_{n,d}$ به وسیلهٔ مجموعهٔ رئوس $\mathcal{V} = \{1, \dots, n\}$ ، متناظر با این مجموعه از n متغیر تصادفی، و مجموعهٔ یالهای $\mathcal{E} = \{(i, j) : |X_i - X_j| \leq d\}$ تعریف می‌شود، که $1 < d < 0$. برای x_1, \dots, x_n از متغیرهای تصادفی، یک گراف بازه‌ای را به عنوان تحققی از گراف تصادفی $\mathcal{IG}_{n,d}$ به دست می‌آوریم.

در اینجا مذکور می‌شویم که سطح فاصلهٔ d در $\mathcal{IG}_{n,d}$ معمولاً به عنوان تابعی مشخص از n در نظر گرفته می‌شود.

می‌توان به آسانی بررسی کرد که برای $\mathcal{IG}_{n,d}$

$$P((i, j) \in \mathcal{E}) = P(|X_i - X_j| \leq d) = 2d - d^2.$$

این تعریف مستلزم آن است که علاوه بر خواص نظری گراف و مشخصه‌ها، بتوان خواص مربوط به هندسه رأسها را مطالعه کرد، یعنی، تحقیقهای x_1, x_2, \dots, x_n روی بازه واحد. برای مثال، طول یک زیرگراف همبند، یا "اندازه" (تعداد رأسها) اولین مؤلفه بدیهی است که، پالهای این گرافهای تصادفی مستقل نیستند. برای d کوچک (مرحله تنک از تکامل $\mathcal{IG}_{n,d}$)، این وابستگی قوی نیست. بنابراین وقتی که احتمالهای یال $\sim P(n)$ از $2d(n)-d^2(n)$ به اندازه کافی کوچک هستند، انتظار می‌رود مشابهت‌هایی با $\mathcal{G}_{n,p}$ (گیلبرت ۱۹۵۹) را ببینید و وجود داشته باشد. اما، در حالت کلی، هر دو مدل، حتی وقتی d خیلی کوچک است، در بسیاری جهات متفاوت هستند.

۱.۲ خواص گرافهای تصادفی بازه‌ای

از مهمترین خواص گرافهای تصادفی دو قضیه زیر از شاینرمون (۱۹۸۸) هستند که آنها را بدون برهان در زیر نقل می‌کنیم.

قضیه ۱: تقریباً همه گرافها در G_n ، $P(d(v) \leq xn) = o(n^2) + \frac{n}{3}$ یال دارند.

قضیه ۲: فرض کنید $G \in G_n$ و $v \in V(G)$. برای مقدار ثابت x با $[0, 1]$ ، داریم

$$\lim_{n \rightarrow \infty} P(d(v) \leq xn) = \begin{cases} 1 - (1-x)^{\frac{x}{2}} & ; x \geq \frac{1}{2} \\ 1 - (1-x) \left\{ \frac{x}{2} - 2 \cos^{-1} \left[\frac{1}{\sqrt{2-x}} \right] \right\} - \sqrt{1-2x} & ; x < \frac{1}{2} \end{cases}$$

Δ و Δ مینیمم و ماکسیمم درجه‌های یک گراف را نشان می‌دهند. نتیجه فوق ایجاب می‌کند:

نتیجه ۱: برای هر $\varepsilon > 0$ ، تقریباً همه گرافهای بازه‌ای در $n < \delta < (1-\varepsilon)n$ صدق می‌کنند.

۳ نتایج دقیق برای گراف تصادفی بازه‌ای $\mathcal{IG}_{n,d}$

این بخش را با دنباله‌ای از لم‌ها آغاز می‌کنیم که بر اساس آنها گودهارت و یاورووسکی (۱۹۹۶) فرمولهای دقیقی را برای احتمالهای بعضی خواص اصلی همبندی $\mathcal{IG}_{n,d}$ به دست آورند. در اینجا منذکر می‌شویم که لم کلیدی، یعنی لم $??$ ، در واقع هم‌ارز با نتیجه مربوط به مسائل پوشش است. احتمالی که می‌خواهیم به دست آوریم مساوی با این احتمال است

که ۱ - کمان انتخاب شده به تصادف روی محیط یک دایره با پیرامون y به طور کامل دایره را پوشاند.

لم ۲ : فرض کنید $[x, x + y]$ یک زیربازه به طول y از $[1, 0]$ باشد. فرض کنید دو تا از k رأس مفروض در مرزهای این زیربازه قرار داده شوند و $\mathcal{A}_{k,y,d}$ عبارت از این پیشامد باشد که ۲ - k نقطه، متناظر با باقیمانده رأسها که به طور تصادفی از $[1, 0]$ استخراج می‌شوند، درون $[x, x + y]$ باشند و مرزها را به هم "متصل" کنند، یعنی، k رأس یک زیرگراف همبند به طول y تشکیل می‌دهند و فرض کنید $P(k, y, d) = P(\mathcal{A}_{k,y,d})$. آنگاه رابطه بازگشتی زیر برقرار است:

$$P(k, y, d) = yP(k-1, y, d) + ((k-1)d - y)P(k-1, y-d, d)$$

با شرایط اولیه

$$P(i, u, d) = \begin{cases} u^{i-1} & i \geq 2 \text{ و } u \leq d \\ 0 & i = 2 \text{ و } u > d \end{cases}$$

یک فرمول بسته برای این احتمال عبارت است از

$$P(k, y, d) = \sum_{j=0}^{\min\{k-1, [y/d]\}} \binom{k-1}{j} (-1)^j (y-jd)^{k-2}.$$

برهان. توجه می‌کنیم که $P(k, y, d)$ به x, y, d و انتخاب k رأس در فضای احتمال بستگی ندارد. بدیهی است که،

$$P(k, y, d) = 0 \quad y > (k-1)d \quad \text{و} \quad k = 2, 3, \dots$$

زیرا در این صورت زیربازه بزرگتر از آن است که مرزها را با به کار بردن $2-k$ رأس به هم متصل کند. از این پس فرض کنید $d \leq y \leq (k-1)d$ و \mathcal{B} عبارت از این پیشامد باشد که پس از انتخاب کردن $3-k$ نقطه در بازه به طول y ، از میان $2-k$ فاصله‌گذاری به وجود آمده بین این $3-k$ نقطه و نقاط مرزی، $3-k$ فاصله‌گذاری با طول کمتری مساوی d و دقیقاً یک فاصله‌گذاری با طول بزرگتر از d اما کوچکتر یا مساوی $2d$ وجود دارد. برای بیان عبارت زیر، فرض می‌شود که $\mathcal{A}_{k-1,y,d}$ به وسیله برخی از $3-k$ رأس مفروض از میان $2-k$ رأس دیگر تولید شود. آنگاه

$$P(\mathcal{A}_{k,y,d}) = P(\mathcal{A}_{k,y,d} | \mathcal{A}_{k-1,y,d})P(\mathcal{A}_{k-1,y,d}) + P(\mathcal{A}_{k,y,d} \cap \mathcal{A}_{k-1,y,d}^c \cap \mathcal{B})$$

و به وضوح

$$P(\mathcal{A}_{k,y,d} | \mathcal{A}_{k-1,y,d}) = y \quad \text{و} \quad \mathcal{A}_{k-1,y,d}^c \cap \mathcal{B} = \mathcal{B}$$

حال پیشامد $\mathcal{A}_{k,y,d} \cap \mathcal{B}$ را در نظر می‌گیریم. ملاحظه کنید که پس از حذف یک بازه به طول d از فاصله‌گذاری بزرگ یکتا، باید $k-2$ - نقطه مرزهای بازه به طول $y-d$ را به یکدیگر متصل کنند. بنابراین می‌توانیم ابتدا $k-3$ - نقطه مفروض را در بازه به طول $y-d$ به گونه‌ای انتخاب کنیم که پیشامد $\mathcal{A}_{k-1,y-d,d}$ برقرار باشد. سپس برای به دست آوردن \mathcal{B} می‌بایست یکی از $k-2$ - فاصله‌گذاری را با اضافه کردن بازهای به طول d بسط دهیم و سرانجام k -امین رأس را در فاصله‌گذاری "بزرگ" برای به دست آوردن $\mathcal{A}_{k,y,d}$ انتخاب کنیم. $k-3$ - نقطه در نظر گرفته شده بازه به طول $y-d$ را به $k-2$ - بازه از هم جدا به طولهای y_1, y_2, \dots, y_{k-2} به گونه‌ای تقسیم می‌کند که برای $i=1, 2, \dots, k-2$ داریم $y_i + y_{i+1} + \dots + y_{k-2} = y - d$ و $y_i \leq d$. از این رو ما $k-2$ -امکان مجزا برای ایجاد $\mathcal{A}_{k,y,d} \cap \mathcal{B}$ داریم. حال فرض کنید یکی از آنها را در نظر بگیریم: برای مثال، فرض کنید که i -امین بازه به طول y_i را به یک بازه به طول d بسط دهیم و آنگاه k -امین رأس را برای به دست آوردن $\mathcal{A}_{k,y,d}$ ارائه دهیم. این با احتمال y_i/d - انجام می‌شود. این بدان معنی است که برای یک انتخاب مفروض y_1, y_2, \dots, y_{k-2} به گونه‌ای که برای $i=1, 2, \dots, k-2$ داریم $y_1 + y_2 + \dots + y_{k-2} = y - d$ و $y_i \leq d$ را اضافه می‌کنیم و k -امین رأس را به گونه‌ای انتخاب می‌کنیم که با احتمال زیر روی دهد.

$$\sum_{i=1}^{k-2} (d - y_i) = (k-1)d - y,$$

که در واقع به انتخاب y_1, y_2, \dots, y_{k-2} بستگی ندارد. بنابراین،

$$P(\mathcal{A}_{k,y,d} \cap \mathcal{B}) = ((k-1)d - y)P(k-1, y-d, d)$$

و چون شرایط اولیه بدیهی هستند، اثبات فرمول بازگشتی برای $P(k-1, y-d, d)$ کامل می‌شود.

به آسانی می‌توان تحقیق کرد که فرمول برای $P(k, y, d)$ جواب فرمول بازگشتی ما با شرایط اولیه داده شده می‌باشد. به علاوه،

$$\sum_{j=0}^{k-1} \binom{k-1}{j} (-1)^j (y-jd)^{k-2} = d^{k-2} A_{k-1}(-\frac{y}{d}, \frac{y}{d} - k + 1; 0, -1),$$

که

$$A_N(a, b; r, s) = \sum_{i=0}^N \binom{N}{i} (i+a)^{i+r} (N-i+b)^{N-i+s}$$

نماد ری یوردان برای مجموع آبل^۱ است (ری یوردان (۱۹۶۸) را ببینید). با به کار بردن فرمول آبل

$$A_N(a, b; \circ, -1) = b^{-1} (a + b + N)^N,$$

به دست می آوریم

$$\sum_{j=0}^{k-1} \binom{k-1}{j} (-1)^j (y - jd)^{k-1} = 0 \quad y > (k-1)d \quad k = 2, 3, \dots \quad \text{اگر}$$

□ همچنین برای $y > (k-1)d$ نیز معتبر است.

حال این پیشامد را در نظر بگیرید که k رأس مفروض، یک زیرگراف همبند $\mathcal{IG}_{n,d}$ را با این خاصیت اضافی تشکیل دهنده که بین نقاط متوالی روی خط (تحقیق راسها)، راسهای دیگری از گراف بازه‌ای موجود نباشد. یک زیرگراف همبند با این خاصیت را یک زیرگراف همبند کامل می‌نامیم.

لم ۳: فرض کنید $\mathcal{CS}_k(u)$ عبارت از این پیشامد باشد که k رأس مفروض یک زیرگراف همبند کامل با طول کمتر از $u > 0$ را تشکیل دهنده. آنگاه

$$P(\mathcal{CS}_k(u)) = \frac{n-k+1}{\binom{n}{k}} \sum_{j=0}^{\min\{k-1, \lfloor u/d \rfloor\}} \binom{k-1}{j} (-1)^j \sum_{t=k-1}^n \binom{n}{t} (1-u)^{n-t} \times (u-jd)^t.$$

به ویژه، فرض کنید \mathcal{CS}_k پیشامد آن باشد که k رأس مفروض یک زیرگراف همبند کامل را تشکیل دهنده، یعنی، پیشامد $\mathcal{CS}_k(u)$ که $u \geq (k-1)d$. آنگاه

$$P(\mathcal{CS}_k) = \frac{n-k+1}{\binom{n}{k}} \sum_{j=0}^{\min\{k-1, \lfloor 1/d \rfloor\}} \binom{k-1}{j} (-1)^j (1-jd)^n.$$

برهان. با استفاده از تعریف، $P(\mathcal{CS}_1(u)) = P(\mathcal{CS}_1) = 1$ و به آسانی می‌توان تحقیق کرد که فرمولهای مفروض برای $k = 1$ نیز مساوی ۱ است. از این رو، می‌توان فرض

^۱ Abel

کرد که $2 \geq k$. لم فوق را می‌توان توسط انتگرال‌گیری با به کار بردن لم λ ثابت کرد.
داریم

$$\begin{aligned} P(\mathcal{CS}_k(u)) &= k(k-1) \int_0^u \int_0^{1-y} dx P(k, y, d)(1-y)^{n-k} dy \\ &= k(k-1) \int_0^u P(k, y, d)(1-y)^{n-k+1} dy. \end{aligned}$$

ابتدا قسمت دوم لم را اثبات می‌کنیم، زیرا با به کار بردن کران بالای ۱ در انتگرال فوق روابط به سادگی محاسبه می‌شوند. (پایان اثبات لم λ را ببینید). با تقسیم کردن بازه $[1, u]$ به زیربازه‌ها، به دست می‌آوریم،

$$\begin{aligned} P(\mathcal{CS}_k) &= k(k-1) \sum_{t=0}^{\lfloor 1/d \rfloor - 1} \int_{td}^{(t+1)d} \sum_{j=0}^M \binom{k-1}{j} (-1)^j (y-jd)^{k-2} (1-y)^{n-k+1} dy \\ &\quad + k(k-1) \int_{\lfloor 1/d \rfloor d}^1 \sum_{j=0}^M \binom{k-1}{j} (-1)^j (y-jd)^{k-2} (1-y)^{n-k+1} dy \end{aligned}$$

با $1 = M = \min\{t, k-1\}$. حال می‌توان مجموع پیرون انتگرال را به داخل برد. آنگاه پس از تغییر دادن ترتیب مجموعیابی به دست می‌آوریم:

$$\begin{aligned} P(\mathcal{CS}_k) &= k(k-1) \sum_{j=0}^{M^*} \binom{k-1}{j} (-1)^j \int_{jd}^1 (y-jd)^{k-2} (1-y)^{n-k+1} dy \\ &= k(k-1) \sum_{j=0}^{M^*} \binom{k-1}{j} (-1)^j \int_0^1 z^{k-2} (1-z)^{n-k+1} (1-jd)^n dz \\ &= k(k-1) \sum_{j=0}^{M^*} \binom{k-1}{j} (-1)^j (1-jd)^n \frac{\Gamma(k-1)\Gamma(n-k+2)}{\Gamma(n+1)}, \end{aligned}$$

که در آن $M^* = \min\{k-1, \lfloor 1/d \rfloor\}$. چون برای هر عدد صحیح نامنفی m $\Gamma(m+1) = m!$ ، اثبات برای $u > (k-1)d$ کامل می‌شود.
به همین شیوه برای حالت کلی به دست می‌آوریم:

$$\begin{aligned} P(\mathcal{CS}_k(u)) &= k(k-1) \sum_{j=0}^{\min\{k-1, \lfloor u/d \rfloor\}} \binom{k-1}{j} (-1)^j \\ &\quad \times \sum_{i=0}^{n-k+1} \binom{n-k+1}{i} (1-u)^{n-k+1-i} (u-jd)^{k+i-1} \frac{\Gamma(k-1)\Gamma(i+1)}{\Gamma(k+i)}, \end{aligned}$$

که برهان را کامل می‌کند. \square

حال فرض کنید که k رأس مفروض، یک زیرگراف همبند کامل تشکیل دهنده، به گونه‌ای که چپ ترین رأس درون فاصله d از \circ قرار گرفته و رأسهای دیگری قبل از آنها روی خط قرار نمی‌گیرند. چنین زیرگرافی، یک زیرگراف مرزی همبند کامل نامیده می‌شود.

لم ۴ : فرض کنید \mathcal{CBS}_k عبارت از این پیشامد باشد که k رأس مفروض یک زیرگراف مرزی همبند کامل را تشکیل دهنده. آنگاه

$$P(\mathcal{CBS}_k) = \frac{1}{\binom{n}{k}} \sum_{j=0}^{\min\{k, \lfloor 1/d \rfloor\}} \binom{k}{j} (-1)^j (1-jd)^n.$$

برهان. فرض کنید یک رأس اضافی را در نقطه صفر قرار دهیم. آنگاه پیشامد \mathcal{CBS}_k بر آن اشاره دارد که این رأس همراه با k رأس مفروض، یک زیرگراف همبند را تشکیل دهنده. بنابراین مجدداً با به کار بردن لم ۲؟، شبیه اثبات لم ۲؟، داریم

$$\begin{aligned} P(\mathcal{CBS}_k) &= k \int_0^1 P(k+1, y, d) (1-y)^{n-k} dy \\ &= k \sum_{j=0}^{\min\{k, \lfloor 1/d \rfloor\}} \binom{k}{j} (-1)^j \int_0^{1-jd} y^{k-1} (1-jd-y)^{n-k} dy \\ &= k \sum_{j=0}^{\min\{k, \lfloor 1/d \rfloor\}} \binom{k}{j} (-1)^j (1-jd)^n \frac{\Gamma(k)\Gamma(n-k+1)}{\Gamma(n+1)} \square \end{aligned}$$

لم ۵ : فرض کنید \mathcal{CC}_k پیشامد آن باشد که k رأس مفروض، یک مؤلفه همبند را تشکیل دهنده. آنگاه، برای $n < k$ ، احتمال این پیشامد به صورت زیر داده می‌شود

$$\begin{aligned} P(\mathcal{CC}_k) &= \frac{1}{\binom{n}{k}} \sum_{j=0}^{\min\{k, \lfloor 1/d \rfloor - 1\}} \binom{k}{j} (-1)^j (1-(j+1)d)^n \\ &\quad + \frac{n-k+1}{\binom{n}{k}} \sum_{j=0}^{\min\{k-1, \lfloor 1/d \rfloor - 1\}} \binom{k-1}{j} (-1)^j (1-(j+2)d)^n \end{aligned}$$

و

$$P(\mathcal{CC}_n) = P(\mathcal{CS}_n).$$

برهان. با استفاده از لم \mathcal{P}_2 ، احتمال آنکه تمام n رأس از یک زیربازه به طول $1 - 2d$ انتخاب شوند و در مجموع k رأس مفروض از آنها یک زیرگراف همبند کامل تشکیل دهنده، به صورت زیر است

$$\frac{n-k+1}{\binom{n}{k}} \sum_{j=0}^{\min\{k-1, \lfloor 1/d \rfloor - 1\}} \binom{k-1}{j} (-1)^j (1 - 2d - jd)^n$$

زیرا تحت فرض فوق، عملاً با یک زیرگراف بازه‌ای \mathcal{IG}_{n,d^*} سر و کار داریم، که $d^* = d/(1 - 2d)$. حال فرض کنید این بازه به طول $1 - 2d$ را با افزودن دو بازه به طول d به بازه اولیه $[1, 0]$ بسط دهیم، یکی قبل از نقطه منتهی‌الیه سمت چپ زیرگراف و دیگری بعد از نقطه منتهی‌الیه سمت راست. آنگاه زیرگراف همبند، یک مؤلفه $\mathcal{IG}_{n,d}$ خواهد بود. بنابراین فرمول فوق احتمال مطلوب را برای حالت خاصی نشان می‌دهد که هر دو نقاط انتهایی مؤلفه، درون فاصله d از مرزهای بازه واحد واقع نیستند. سه حالت دیگری که اینجا در نظر گرفته می‌شوند عبارتند از: دو، هم ارزی بالفعل به واسطه مقارن، که درست یکی از نقاط انتهایی درون فاصله d از حاشیه بازه واحد است؛ و آخرین حالت موردعی است که هر دو نقاط انتهایی دارای این خاصیت می‌باشند. اما، آخرین حالت تنها برای $k = n$ امکان پذیر است، که بدان معنی است که، درست یک مؤلفه داریم؛ احتمال برای $k = n$ با استفاده از لم \mathcal{P}_2 داده می‌شود. برای دو حالت مقارن، لم \mathcal{P}_2 به شیوه‌ای همانند با حالت نخست به کار برده می‌شود؛ یعنی هم اکنون n نقطه را روی یک زیربازه به طول d قرار می‌دهیم. با استفاده از لم \mathcal{P}_2 ، احتمال اینکه تمام n رأس از یک زیربازه به طول $1 - d$ انتخاب شوند و k رأس مفروض آنها یک زیرگراف مرزی همبند کامل تشکیل دهنده عبارت است از

$$\frac{1}{\binom{n}{k}} \sum_{j=0}^{\min\{k, \lfloor 1/d \rfloor - 1\}} \binom{k}{j} (-1)^j (1 - d - jd)^n.$$

حال، زیربازه به طول $d - 1$ را به بازه $[1, 0]$ بسط می‌دهیم، با افزودن یک بازه به طول d ، یا قبل از نقطه منتهی‌الیه سمت چپ زیرگراف یا بعد از نقطه منتهی‌الیه سمت راست که واثبات قضیه فوراً نتیجه می‌شود. \square

لم ۶: فرض کنید \mathcal{P}_2 عبارت از این پیشامد باشد که دو رأس مفروض به وسیله یک راه به هم متصل شوند. آنگاه

$$P(\mathcal{P}_2) = 1 - \frac{1}{\binom{n}{2}} \sum_{j=1}^{\min\{n-1, \lfloor 1/d \rfloor\}} \binom{n+1}{j+2} (-1)^{j+1} (1 - jd)^n.$$

برهان. بدیهی است که، اگر دو رأس توسط یک راه به هم متصل شوند، آنگاه آنها در مؤلفه یکسانی می‌باشند و بالعکس. آنگاه

$$P(\mathcal{P}_2) = 2 \sum_{k=0}^{n-1} \binom{n-2}{k} \frac{1}{(k+1)(k+2)} P(\mathcal{CS}_{k+2}).$$

از این رو با استفاده از لم ??، بعد از تغییر دادن ترتیب مجموعیابی، به دست می‌آوریم

$$\begin{aligned} P(\mathcal{P}_2) &= \frac{1}{\binom{n}{2}} \sum_{j=1}^{\min\{n-1, \lfloor 1/d \rfloor\}} (-1)^j (1-jd)^n \sum_{k=j-1}^{n-2} \binom{k+1}{j} (n-k-1) \\ &\quad + \frac{1}{\binom{n}{2}} \sum_{k=0}^{n-2} (n-k-1). \end{aligned}$$

آخرین مجموع $\binom{n}{2}$ است و چون

$$\sum_{k=j-1}^{n-2} \binom{k+1}{j} (n-k-1) = \sum_{k=0}^{n-j-1} \binom{n-k}{j+1} = \binom{n+1}{j+2},$$

□

د (برای اتحادها برای مثال گولد (۱۹۷۲) را ببینید).

حال گودهارت و یاواروسکی (۱۹۹۶) نتایج اصلی اشان را در رابطه با توزیعهای متغیرهای تصادفی بیان می‌کنند: تعداد C_n مؤلفه، تعداد C_n^k مؤلفه با اندازه k ، و تعداد C_n^1 رأس منزوی.

قضیه ۳ : فرض کنید C_n تعداد مؤلفه‌ها در یک گراف تصادفی بازه‌ای $\mathcal{IG}_{n,d}$ باشد، توزیع احتمال گستتهٔ متغیر تصادفی C_n با رابطه زیر داده می‌شود.

$$P(C_n = r) = \sum_{j=r-1}^{\min\{n-1, \lfloor 1/d \rfloor\}} \binom{n-1}{j} \binom{j}{r-1} (-1)^{j+r-1} (1-jd)^n$$

برای $1 \leq r \leq \min\{n-1, \lfloor 1/d \rfloor\} + 1, 2, \dots, n$. این گشتاور فاکتوریل برای تعداد مؤلفه‌های کاهش یافته به اندازهٔ ۱ عبارت است از

$$E_t(C_n - 1) = (n-1)_t (1-td)^n$$

هرگاه $0 < td < 1$ ، و در غیر این صورت صفر می‌شود. برهان. این قضیه با استفاده از روش فاصله‌گذاری اثبات می‌شود. نقطه، بازه را به $n-1$

۱ زیربازه متوالی بین کوچکترین نقطه و بزرگترین نقطه افزای می‌کنند که طولهای آن زیربازه‌ها فاصله گذاری ها نامیده می‌شوند (باربور، هولست و جانسون ۱۹۹۲) را ببینید).
 ۱ $n - 1$ متغیرهای I_1, I_2, \dots, I_{n-1} را به گونه‌ای تعریف کنید که برای $i = 1, 2, \dots, n - 1$

$$I_i = \begin{cases} 1 & \text{اگر } i \text{ این فاصله گذاری از } d \text{ تجاوز کند} \\ 0 & \text{در سایر جاهای} \end{cases}$$

به آسانی می‌توان دید که

$$C_n = 1 = I_1 + I_2 + \dots + I_{n-1}.$$

فرض کنید $[1/d] < j$. پیشامدی را در نظر بگیرید که j فاصله گذاری مفروض، مثلاً فاصله گذاری i_1, i_2, \dots, i_j از d بزرگتر باشند. بدیهی است که، در این مدل احتمال این پیشامد به انتخاب i_1, i_2, \dots, i_j بستگی ندارد. می‌توان n نقطه به طور یکنواخت و مستقل از زیربازه به طول $jd - 1$ با احتمال $(1 - jd)^n$ استخراج کرد. سپس فاصله گذاریهای i_1, i_2, \dots, i_j را، هر کدام به اندازه d بسط می‌دهیم که تمامی توزیعهای ممکن n نقطه در بازه واحد را که برای آن فاصله گذاریهای i_1, i_2, \dots, i_j بزرگتر از d هستند را خواهد داد. بنابراین، احتمال پیشامدی که در نظر گرفته می‌شود، عبارت است از

$$P(I_{i_1} = 1, I_{i_2} = 1, \dots, I_{i_j} = 1) = (1 - jd)^n.$$

حال برهان قضیه ؟؟ با استفاده از کاربرد متناول فرمول رد و قبول کامل می‌شود (بلباش (۱۹۸۵) را ببینید).

نتیجه ۲: احتمال آنکه یک گراف تصادفی بازه‌ای $\mathcal{IG}_{n,d}$ همبند باشد عبارت است از

$$P(\mathcal{IG}_{n,d} \text{ همبند باشد}) = \sum_{j=0}^{\min\{n-1, [1/d]\}} \binom{n-1}{j} (-1)^j (1 - jd)^n,$$

که برای $1/(n-1) < d$ به صورت زیر کاهاش می‌پابد

$$P(\mathcal{IG}_{n,d} \text{ همبند باشد}) = n! d^{n-1} \left(1 - \frac{1}{2}(n-1)d\right).$$

برهان. چون

$$P(\mathcal{IG}_{n,d} \text{ همبند باشد}) = P(C_n = 1),$$

اولین قسمت نتیجه یک پیامد مستقیمی از قضیه ?? است و نیز می‌توان آن را با استفاده از لم ?? نیز اثبات کرد ($P(\mathcal{CS}_n)$).

برای $1/d < n - 1/(n-1)$ کاوش می‌یابد. در نتیجه

$$\begin{aligned} P(C_n = 1) &= d^n \sum_{j=0}^{n-1} \binom{n-1}{j} (-1)^j \left(\frac{1}{d} - j\right)^n \\ &= -d^n \sum_{j=0}^{n-1} \binom{n-1}{j} \left(j - \frac{1}{d}\right)^{j+1} \left(\frac{1}{d} - n + 1 + n - 1 - j\right)^{n-1-j} \end{aligned}$$

با استفاده مجدد از نماد ری یورдан برای جمعهای آبل (اثبات لم ?? را ببینید) اتحادهای دیگر آبل، می‌توان احتمال همبندی را به صورت زیر دوباره نویسی کرد.

$$\begin{aligned} P(C_n = 1) &= -d^n A_{n-1}(-1/d, 1/d - n + 1; 1, 0) = -d^n (\alpha + \beta(-1/d))^{n-1} \\ &= -d^n \sum_{i=0}^{n-1} \binom{n-1}{i} i! (n-1-i)! (n-1-i-1/d) \\ &= -d^n \sum_{i=0}^{n-1} (n-1)! (n-1-i-1/d) = n! d^{n-1} \left(1 - \frac{1}{d}(n-1)d\right), \end{aligned}$$

که α و $\beta(x)$ متغیرهای ظاهری هستند به طوری که

$$(\beta(x))^i = \beta_i(x) = i!(x+i) \quad \text{و} \quad \alpha^i = \alpha_i = i!$$

□

(ری یوردان ۱۹۶۸) را ببینید.

یکی از نتایج مستقیم لم ?? فرمولی برای مقدار مورد انتظار تعداد C_n^k مؤلفه به اندازه k است. با شیوه‌ای نظری لم ?? و با به کار بردن همان اتحادهای دو جمله‌ای قضیه بعدی به دست می‌آید.

قضیه ۴ : فرض کنید C_n^k تعداد مؤلفه‌های به اندازه k در گراف تصادفی بازه‌ای باشد. آنگاه، تعداد مورد انتظار مؤلفه‌های به اندازه بزرگتر از m به صورت زیر داده می‌شود

$$\begin{aligned} \sum_{k=m+1}^n E(C_n^k) &= \sum_{j=0}^{\min\{m+1, \lfloor 1/d \rfloor\}} \binom{m+1}{j} (-1)^j (1-jd)^n \\ &\quad + (n-m) \sum_{j=0}^{\min\{m, \lfloor 1/d \rfloor - 1\}} \binom{m}{j} (-1)^j (1-(j+1)d)^n. \end{aligned}$$

برای به دست آوردن توزیع احتمال دقیق تعداد مؤلفه‌هایی با یک اندازه مفروض ابتدا حالت خاص رأسهای منزوی را در نظر می‌گیریم.

قضیه ۵: توزیع تعداد C_n^1 رأسهای منزوی در یک گراف تصادفی بازه‌ای $\mathcal{IG}_{n,d}$ به صورت زیر داده می‌شود

$$\begin{aligned} P(C_n^1 = r) &= \sum_{k=r}^{\min\{n-1, \lfloor 1/d \rfloor\}} \binom{k}{r} (-1)^{k-r} \sum_{j=j_0}^k \binom{n-k-1}{k-j} \binom{k+1}{j} \\ &\quad \times (1 - (2k-j)d)^n + \binom{n}{r} (-1)^{n-r} (1 - (n-1)d)_+^n \end{aligned}$$

با $\{1\}$ ، که برای $a_0 = \max\{0, 2k - \lfloor 1/d \rfloor, 2k - n + 1\}$ مثبت، $a_+ = a$ و در غیر این صورت $a_0 = 0$. به علاوه، t امین گشتاور فاکتوریل برای تعداد رأسهای منزوی به صورت زیر داده می‌شود

$$E_t(C_n^1) = t! \sum_{j=j_0}^t \binom{n-t-1}{t-j} \binom{t+1}{j} (1 - (2t-j)d)^n$$

با $\{1\}$ و در غیر این صورت صفر می‌شود.

قضیه ۶: امین گشتاور فاکتوریل برای تعداد مؤلفه‌های همبند به اندازه k به صورت زیر داده می‌شود

$$\begin{aligned} E_t(C_n^k) &= t! \sum_{j=j_0}^t \binom{n-kt-1}{t-j} \binom{t+1}{j} \sum_{i=0}^{(k-1)t} \binom{(k-1)t}{i} (-1)^i \\ &\quad \times (1 - (2t-j+i)d)_+^n \end{aligned}$$

با $\{1\}$ و برای $a_0 = a$ مثبت، $a_+ = a$ و در غیر این صورت $a_0 = 0$.

۴ مفاهیم نظریه گراف در فضای k بعدی

فرض کنید $F_{\underline{X}}$ یکتابع توزیع تجمعی روی فضای اقلیدسی k بعدی، E_k ، و نسبت به اندازه لبگ k بعدی مطلقاً پیوسته باشد. تابع چگالی احتمال متناظر را با $f_{\underline{X}}$ نشان می‌دهیم. فرض کنید یک نمونه تصادفی به اندازه n از $F_{\underline{X}}$ به دست آمده است و یافته‌ها را با $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ نشان دهد.

حال گراف (V, E) چنین تعریف می‌شود: $G_n = (V, E)$ یک مجموعه با $|V| = n$ و یک مجموعه از جفت‌های (بی‌ترتیب) عناصر V است. عناصر V رأس‌های گراف G_n و جفت‌ها در E به عنوان یالهای G_n نامیده می‌شوند. بدون از دست دادن کلیت، می‌توانیم فرض کنیم $\{1, 2, \dots, n\} = V$. برای اهداف مد نظر، یک فاصله ρ روی E_k و یک آستانه $d > 0$ انتخاب می‌کنیم. بنابراین برای $j \neq i$, $(i, j) \in E$ قرار دهید اگر و تنها اگر $\rho(\underline{x}_i, \underline{x}_j) \leq d$. چون $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ تحقیقاتی متغیرهای تصادفی هستند، مجموعه E یک مجموعه تصادفی و گراف G_n یک گراف تصادفی است. به ویژه، این گرافها تعمیمهایی از گرافهای بازاری هستند. مشخصاً اگر I_1, I_2, \dots, I_n بازه‌هایی روی خط حقیقی باشند، آنگاه گراف بازاری $G(I_n)$ به وسیله $\{(i, j) \in V \mid 1 \leq i < j \leq n \text{ و } I_i \cap I_j \neq \emptyset\}$ تعریف می‌شوند. بنابراین برای مدل بررسی، اگر $i = k$ ، بازه‌های I_i ، بازه‌های $[x_i - \frac{d}{2}, x_i + \frac{d}{2}]$ هستند.

تعریف ۱: فرض کنید $K_{m,d}$ یک زیرگراف کامل از $V_m \subset V$ با $|V_m| = m$ است، اگر همه $\binom{m}{2}$ جفت عناصر V_m در E باشند. اگر $1 \leq m = K_{1,d}$ یک رأس است، اگر $2 \leq m = K_{2,d}$ یک یال است و اگر $3 \leq m = K_{3,d}$ یک مثلث نامیده می‌شود.

تعریف ۲: یک زیرگراف کامل از مرتبه m یک زیرگراف کامل بیشینه نامیده می‌شود، و با $K_{m,d}^*$ نشان داده می‌شود، هرگاه رأسی در $V \setminus V_m$ موجود نباشد به طوری که از متصل کردن آن رأس به V_m یک زیرگراف کامل از مرتبه $m+1$ منتج شود.

درجه رأس v , $v = 1, 2, \dots, n-1$, در G_n تعداد یالهای G_n است که v بر آنها واقع است. اگر $v = n$, آنگاه آن رأس منزوی نامیده می‌شود.

۵ توزیعهای احتمال برای مشخصه‌های حقیقی گرافهای بازاری

حال احتمال آن را که یک مجموعه معین از m رأس، یک $K_{m,d}$ تشکیل دهنده توصیف می‌کنیم. بدون از دست دادن کلیت، فرض می‌کنیم که این رأسها با $1, 2, \dots, m$ برچسب

گذاری می‌شوند.

احتمال آنکه یک مجموعه معین از m رأس، یک $K_{m,d}$ تشکیل دهنده صورت زیر است.

$$P\left\{\max_{1 \leq i \leq m} X_i - \min_{1 \leq i \leq m} X_i \leq d\right\} = m \int_{-\infty}^{\infty} \{F(x+d) - F(x)\}^{m-1} f(x) dx. \quad (1)$$

احتمال آنکه یک رأس معین دارای درجه $v = 0, 1, \dots, n-1$ باشد عبارت است از

$$\begin{aligned} P\left\{\text{دارد } v \text{ رأس } 1 \text{ درجه}\right\} &= \binom{n-1}{v} \int_{-\infty}^{\infty} \{F(x+d) - F(x-d)\}^v \\ &\times \{1 - F(x+d) + F(x-d)\}^{n-v-1} f(x) dx. \end{aligned} \quad (2)$$

نتیجه ۳: احتمال آنکه یک رأس معین یک $K_{1,d}^*$ تشکیل دهد (یعنی منزوی باشد) عبارت است از

$$P\left\{1 \text{ منزوی}\right\} = \int_{-\infty}^{\infty} \{1 - F(x+d) + F(x-d)\}^{n-1} f(x) dx. \quad (3)$$

با توجه به تعریف رأس منزوی اگر در فرمول (۲)، $v = 0$ را قرار دهیم رابطه (۳) به دست خواهد آمد.

۱.۵ خواص مجانبی گرافهای تصادفی بازه‌ای

برای به دست آوردن تقریب‌های مجانبی برای توزیعهای بالا، به فرضهایی درباره رفتار تابع چگالی احتمال $f_X(x)$ نیاز داریم. از این رو فرض می‌کنیم که تابع چگالی احتمال روی هر زیرمجموعه فشرده از مجموعه حامل X به طور یکنواخت پیوسته است و فرض می‌کنیم که $f'_X(x)$ موجود و روی مجموعه حامل X به طور یکنواخت کراندار است. حال رفتار مجانبی توزیعهای احتمال معرفی شده در بخش قبلی را، تحت شرایط $n \rightarrow \infty$ و $d \rightarrow 0$ (معمولًا، d تابعی مفروض از n است، $d = d(n)$) آزمون می‌کنیم. همچنین شرایط نظم را برای $f_X(x)$ داده شده در بالا در نظر می‌گیریم. فرض کنید $(A_t)_{t \in T}$ خانواده‌ای از پیشامدها باشد که به پaramتر t وابسته است. می‌گوییم احتمال مجانبی A_t ، برای $t \rightarrow t'$ ، است هرگاه $\lim_{t \rightarrow t'} P\{A_t\}/p = 1$.

احتمال مجانبی $(t \rightarrow d \text{ هرگاه } \infty \rightarrow n) \text{ آنکه رأسهای } 1, 2, \dots, m \text{ یک } K_{m,d}$ تشکیل دهنند،

$$md^{m-1} \int_{-\infty}^{\infty} f^m(x) dx \quad (4)$$

است.

میانگینها و واریانس‌های مجانبی تعداد زیرگرافهای کامل از مرتبه m به صورت زیر داده می‌شوند.

$$E\{|K_{m,d}|\} \sim \binom{n}{m} P\{K_{m,d}\} = \binom{n}{m} m d^{m-1} \int_{-\infty}^{\infty} f^m(x) dx \quad (5)$$

و

$$Var\{|K_{m,d}|\} \sim \binom{n}{2m-1} (2d)^{2m-2} \int_{-\infty}^{\infty} f^{2m-1}(x) dx. \quad (6)$$

توجه کنید که هرگاه $\rightarrow d$ ، احتمال آنکه رأسهای معین یک $K_{m,d}$ تشکیل دهند به صفر میل می‌کند. فرمول (4) همچنین به صفر میل می‌کند، اما نسبت آنها به سمت واحد میل می‌کند.

احتمال مجانبی $(\rightarrow d)$ هرگاه $\rightarrow \infty$ ، به طوری که $\rightarrow (nd)$ آنکه رأس 1 دارای درجه v باشد

$$\binom{n-1}{v} (2d)^v \int_{-\infty}^{\infty} f^{v+1}(x) \{1 - 2ndf(x)\} dx \quad (7)$$

است.

۲.۵ آشکارسازی آمیزه‌های توزیعهای احتمال

یک مسئله که دقیقاً به آشکارسازی خوشها مربوط می‌شود، آشکارسازی آمیزه‌های توزیعهای احتمال است. یک آمیزه از k تابع چگالی احتمال یک تابع چگالی احتمال به شکل

$$f_X(x) = \sum_{i=1}^k \alpha_i f_i(x),$$

است که در آن $\sum_{i=1}^k \alpha_i = 1$ و $\alpha_i > 0$. به منظور اجتناب از بعضی پیچیدگیهای ریاضی، آن توزیعهای احتمالی را به کار می‌بریم که تابع چگالی احتمالشان ($f_i(x)$) بوده و دو به دو مطلقاً پیوسته و متمایز باشند. (یعنی یک مجموعه $A_{(i,j)}$ از اندازه لبگ مثبت وجود دارد که برای آن $\int_{A_{(i,j)}} f_i(x) dx \neq \int_{A_{(i,j)}} f_j(x) dx$ برای $i < j$ و $i \neq j$). فرض کنید X_1, X_2, \dots, X_n یک نمونه تصادفی از $f_X(x)$ است و فرض کنید $d > k$ یک آستانه معین باشد. تقریبات مجانبی که می‌باشد اینجا به کار روند تحت فرض آنکه $d \rightarrow \infty$ معتبر هستند، به طوری که $\rightarrow \infty$ $n^m d^{m-1}$. بنابراین،

برای مثال، احتمال مجانبی آنکه m تحقق به طور تصادفی انتخاب شده به یک زیرگراف کامل از مرتبه m منجر شود به صورت

$$P\{K_{m,d}\} \sim md^{m-1} \sum_{j_1, \dots, j_k} \binom{m}{j_1, \dots, j_k} \prod_{i=1}^k \left\{ \alpha_i^{j_i} \int_{-\infty}^{\infty} f_i^{j_i}(x) dx \right\},$$

است، که در آن $\sum_{l=1}^k j_l = m$ و $i = 1, \dots, k$ ، $j_i \geq 0$ می‌باشد.

هریس و گودهارت (۱۹۹۸) روشی را به شرح زیر پیشنهاد می‌کنند:

ابتدا فرض می‌شود که آمیزه‌ای وجود ندارد، یعنی، $\alpha_1 = \dots = \alpha_k = 1$. مقدار مورد انتظار واریانس $K_{m,d}$ تحت این فرض برآورد شده و تعداد مشاهده شده زیرگرافهای کامل از مرتبه m با این میانگین مقایسه می‌شود. اگر این عدد اساساً (نسبت به واریانس) از مقدار مورد انتظار برآورد شده متفاوت باشد، آنگاه نتیجه خواهد شد که یک آمیزه نابدیهی وجود دارد. می‌توان با تکرار این عمل فرض سازگاری داده‌ها را با آمیزه‌ای از دو توزیع آزمون نمود.

آزمونهای دنباله‌ای توصیف شده به طور تصادفی مستقل نیستند، هرچند که بیشتر روش‌های عملی تحلیل داده‌های آماری در حقیقت به همین روش انجام می‌شوند.

۳.۵ آمیزه‌هایی از توزیعهای نمایی

کاربرد زیر از این روش‌ها احتمالاً مفهوم ضعیفی در تحلیل خوش‌های دارد، اما در نظریه قابلیت اعتماد و تحلیل ریسک به طور طبیعی بروز می‌کند. ما در اینجا، با به کار بردن آزمونهایی براساس نظریه زیرگرافهای کامل از مرتبه m ، آشکار کردن وجود آمیزه‌هایی از توزیعهای نمایی را بررسی می‌کنیم. یک روش طبیعی که می‌توان مطرح کرد استفاده کردن از آزمون نسبت درستنمایی و نظریه مجانبی متناظر است. متأسفانه در مورد آمیزه‌ها، شرایط نظم مورد نیاز برای نظریه مجانبی نسبت درستنمایی برقرار نیست. بنابراین روش کنونی ممکن است راه حل مناسبی باشد.

در نتیجه فرض کنید

$$f_i(x) = \lambda_i e^{-\lambda_i x}, \lambda_i > 0, x > 0, i = 1, 2, \dots, k.$$

از این رو به دست می‌آوریم

$$P\{K_{m,d}\} \sim md^{m-1} \sum_{j_1, \dots, j_k} \binom{m}{j_1, \dots, j_k} \prod_{i=1}^k \left\{ \alpha_i^{j_i} \lambda_i^{j_i} \middle/ \sum_{l=1}^k j_l \lambda_l \right\}.$$

یک مورد خاص در نظر گرفته می‌شود. به منظور سادگی محاسبات در این مثال فرض می‌کنیم $k = 2$ و $m = 2$ ؛ یعنی یک آمیزه از دو توزیع نمایی در نظر می‌گیریم و تعداد

یالهای تحقیق یافته را می‌شماریم.

مثال ۱: فرض کنید X طول عمر تصادفی یک دستگاه و دارای توزیع نمایی باشد باشد. پس از مدتی گمان می‌شود که منبع تولید دومی با نرخ شدت متفاوتی معرفی شده است. بنابراین ممکن است برخی از داده‌ها از منبع تولید اول و برخی دیگر از داده‌ها از منبع تولید دوم آمده باشند، و بنابراین داده‌های مشاهده شده ممکن است از یک آمیزه از دو توزیع نمایی باشند. بدون از دست دادن کلیت مسئله، می‌توانیم فرض کیم $\lambda_1 = \lambda$ باشد. برای ارائه یک مثال عددی، 300 مشاهده شبیه‌سازی شدند، 183 مشاهده از یک توزیع نمایی باشد $\lambda_1 = 117$ و 117 مشاهده از توزیع نمایی با $\lambda_2 = 2$. از این رو $\alpha_1 = 0.61$ و $\alpha_2 = 0.39$. داده‌های طول عمر شبیه‌سازی شده یک میانگین نمونه‌ای $\bar{x} = 57.8$ و یک واریانس نمونه‌ای $s^2 = 72.9$ دارند. آستانه d در سطح 0.05 اختیار شده است (البته این مقدار عمدتاً بستگی به نظرپژوهشگر دارد) که به تعداد 563 یال مشاهده شده منتج می‌شود. در مجموع، گشتاور دوم نمونه‌ای طول عمرها $m_2 = 1/335$ و گشتاور سوم نمونه‌ای $m_3 = 2/628$ است. برای این آمیزه خاص، میانگین نظری طول عمر $\mu = 58.5$ و واریانس نظری $\sigma^2 = 76.7$ است. از این رو داده شبیه‌سازی شده توافق معقولی با نظریه دارد.

برای مثال ۱، فرض کردیم که $\lambda_1 = 1$ معلوم و λ_2 نامعلوم است. فرض صفر طبیعی $H_0: \lambda_1 = \lambda$ است، یعنی، که هیچ آمیزشی موجود نیست. اگر این درست باشد، آنگاه مقدار مجانبی برای تعداد مورد انتظار یالها $448/5$ است، به همین نحو، مقدار مجانبی برای واریانس تعداد یالها تحت فرض صفر $1/594$ است که یک انحراف استاندارد $24/4$ را نشان می‌دهد. بنابراین می‌توان تشخیص داد که تعداد مشاهده شده یالها به اندازه کافی بزرگ است تا فرض صفر رد شود. برای آمیزه خاص به کار رفته، تعداد مورد انتظار یالها $587/8$ است. بنابراین واضح است که برای مثال ۱، آشکارسازی وجود یک آمیزه، را می‌توان با به کار بردن تعداد یالها انجام داد. شخص می‌باشد α و λ_2 را از روی داده‌ها برآورد کند. چندین روش وجود دارد که شخص ممکن است به کار برد. نوعاً شخص داده طول عمر را مورد استفاده قرار می‌دهد که یکی از فنهای برآورد آماری استاندارد نظری روش حداقل درستنمایی یا روش گشتاورها را به کار می‌گیرد. به دلیل سادگی محاسباتی، روش گشتاورها مورد استفاده قرار گرفته، که به برآوردهای $0/55$ و $\alpha = 1/96$ منجر می‌شود.

مراجع

- Barbour, A. D., Holst, L. and Janson, S., (1992). Poisson Approximation, Clarendon Press, Oxford.
- Bollobas, B., (1985). Random Graphs, Academic Press, New York.
- Gilbert, E. N., (1959). Random Graphs, *Ann. Math. Stat.*, 30, 1141-1144.
- Gilmore, P. and Hoffman, A., (1964). A Characterization of Comparability Graphs and of Interval Graphs, *Canad. J. Math.* 16, 539-548.
- Godehardt, E., (1990). Graphs as Structural Models, *Vieweg, Braunschweig*.
- Godehardt, E. and Jaworski, J., (1996). On the Connectivity of a Random Interval Graph, *Random Structures and Algorithms*, Vol. 9, Nos. 1 and 2, 137-161.
- Golumbic, M. C., (1980). Algorithmic Graph Theory and Perfect Graphs , (Academic Press).
- Gould, H. W., (1972). Combinatorial Identities, *Morgantown Printing and Binding, Morgantown, WV*.
- Harris, B. and Godehardt, E., (1998). Probability Models and Limit Theorems for Random Interval Graphs with Applications to Cluster Analysis, *Classification, Data Analysis, and Data Highways*. Springer, Berlin-Heidelberg-New York, 54-61.
- Harris, B. and Godehardt, E., (1999). The Comparative Efficacy of Some Combinatorial Tests for Detection of Clusters and Mixture of Probability Distributions, *Classification in the Information Age*. Springer, 295-301.
- Riordan, J., (1968). Combinatorial Identities, Wiley, New York.
- Scheinerman, E. R., (1990). An Evolution of Interval Graphs, *Discrete Math.*, 82, 287-302.
- Scheinerman, E. R., (1988). Random Interval Graphs, *Combinatorica* 8, 357-371.

نمایه کلید واژه

- رکوردهای پائین (بالا)، ۱۴
رگرسیون پواسن، ۱۲۸
روشهای احتمالاتی، ۶۶
سرو زمانی، ۸۰
سل ریوی، ۱۵۰
شبکه‌های عصبی، ۸۰
- طرح سکه اریب Efron، ۵۸
طرح سکه اریب تطبیقی، ۵۸
طرح کاسه‌ای، ۵۸
کریگینگ عام، ۱۵۰
- ماتریس اطلاع فیشر، ۲۶
متغیر کمکی، ۴۱
مجموعه‌های فازی، ۸۰
مدل‌های آمیخته، ۱۸۶
مدلهای تعیین یافته خطی، ۱۲۸
معیار *BIC*, ۱۸۶
- نایستایی، ۱۵۰
نابرازا، ۹۲
نقشه آماری، ۱۵۰
- نمونه‌گیری با احتمالات نابرابر، ۴۱
واژه کلیدی، ۱۱۱
- رگرسیون دوجمله‌ای منفی، ۱۲۸
اریبی انتخاب، ۵۸
استراتژی احتمال متناسب، ۵۸
استراتژی همگرا، ۵۸
اعداد استرلینگ نوع اول، ۱۴
الگوریتم، ۱۷۵
الگوریتم *EM*, ۱۸۶
- برآوردگر حداقل مربعات، ۱۷۵
برآوردگر درستنامایی ماکزیمم، ۲۶، ۱۴، ۱۴
برآوردگر رگرسیونی، ۴۱
برآوردگر رگرسیونی عام، ۴۱
برآوردگر گشتاوری، ۱۴
برآوردگر ماکزیمم پسین، ۱۷۵
برآوردگر ماکزیمم جرم موضعی، ۱۷۵
برآوردگر نااریب، ۱۴
برآوردگر هارویتز-تامپسون، ۴۱
- تابع درستنامایی، ۲۶
تجزیه طیفی، ۱۸۶
تحلیل داده‌های گستته، ۱
تقعر (تحدّب)، ۱۴
- چگالی پسین، ۱۷۵
چگالی پیشین، ۱۷۵
- حاملگی ناخواسته، ۱۲۸
- خطای طبقه‌بندی نادرست، ۱
خوشبندی کردن، ۱۸۶
- داده‌های زبانی، ۸۰