

آشنایی با روش دلتا برای تعیین توزیع‌های مجانبی و کاربردهایی از آن

منوچهر خردمندیا^۱

چکیده

فرض کنید که یک آماره دارای توزیع مجانبی نرمال است. آگرستی (۱۹۹۰) نشان داد که توابع بسیاری از آن آماره نیز به طور مجانبی نرمال هستند. وی یک روش ساده برای به دست آوردن این توزیع‌ها ارایه نمود. در مقاله حاضر کاربردهایی را از روش مذکور، در الگوسازی لگاریتم خطی و در الگوسازی لجیت معرفی می‌کنیم.

۱ مقدمه

با این روش به ایران‌پناه و پاشا (۱۳۷۶) مراجعه کنید. جک نایف نیز روش دیگری برای برآورد یک برآوردگر است. برای آشنایی با این روش به نیر و مند (۱۳۷۶) مراجعه کنید. روش‌های بوت استرپ و جک نایف، به خصوص در مواردی که شرایط لازم برای به کارگیری روش دلتا موجود نیست مفیدند. برای تعیین توزیع مجانبی تابعی از یک آماره می‌توان تحت شرایط خاصی از روش دلتا استفاده کرد. در این مقاله ضمن معرفی این روش، شرایط خاص مذکور را روشن می‌کنیم. این روش کاربردهای فراوان دارد. در مقاله حاضر بعضی از کاربردهای این روش را که در تحلیل‌های چندمتغیره گستته به کار می‌آیند معرفی می‌کنیم. از مراجع مهم مربوط به این روش می‌توان از رائو (۱۹۷۲) و آگرستی (۱۹۹۰) نام برد. بوت استرپ روش دیگری برای برآورد واریانس و توزیع نمونه‌ای آماره‌هاست. برای آشنایی

۱ دکتر منوچهر خردمندیا، گروه آمار، دانشگاه اصفهان

۲ روش دلتا برای تعیین توزیع مجانبی

است. وقتی که $(\cdot)g'$ و $\sigma = \sigma$ در نقطه θ پیوسته باشند، $\sigma(T_n)g'(T_n)$ یک برآورده سازگار برای $\sigma(\theta)g'(\theta)$ است، بنابراین در این حالت

$$g(T_n) \sim AN[g(\theta), \frac{\sigma^2(t_n)}{n}(g'(t_n))^2] \quad (3)$$

لذا فاصله زیر، یک فاصله اطمینان ۹۵ درصد برای θ است

$$g(t_n) \pm 1/96 \sigma(t_n) | g'(t_n) | / \sqrt{n} \quad (4)$$

که در آن t_n یافته‌ای از T_n است. برای توضیحات نظری بیشتر به فصل ۱۲ کتاب آگرستی (۱۹۹۰) مراجعه کنید. در این مرجع، نتیجه (۲)، روش دلتا برای تعیین توزیع مجانبی نامیده شده است. در همین مرجع، تاریخچه‌ای طولانی از فرایند تکامل این روش ارایه شده است.

برای یافتن وجه تسمیه اصطلاح روش دلتا، با بررسی متون آماری در دسترس، به دایرة المعارف اصطلاحات آماری می‌رسیم: در این مرجع، اصطلاح روش دلتا مستقیماً تعریف و تشریح نشده است. ولی در صفحه ۶۴۶ از جلد هشتم این دایرة المعارف، با ارجاع به فصل ۱۴ کتاب بیشاب و همکاران (۱۹۷۵) از روش دلتا به عنوان نامی جانشین برای روش دیفرانسیل آماری نام برده است. در این کتاب، روش دیفرانسیل آماری به عنوان روشی برای دستیابی به تقریبی برای مقدار مورد انتظار تابعی از یک متغیر تصادفی تعریف شده و آمده است که انحراف متغیرهای تصادفی را از مقادیر مورد انتظارشان $(\mu_j - r_j)$ ، که در روش دیفرانسیل آماری به کار می‌آید، اغلب با نماد δX نشان می‌دهند و لذا وجه تسمیه روش دلتا حاصل می‌شود. روشی که آگرستی (۱۹۹۰) تحت عنوان روش دلتا برای تعیین توزیع مجانبی معرفی نموده، شکل تعديل و تکامل یافته‌ای از روش دیفرانسیل آماری (روش دلتا) است. در بخش بعدی، به منظور رعایت اختصار به جای اصطلاح دقیق ولی طولانی روش دلتا برای تعیین توزیع مجانبی، اصطلاح روش دلتا را به کار می‌بریم.

یک دنباله X_n از متغیرهای تصادفی را به طور مجانبی نرمال (AN) با میانگین μ_n و انحراف معیار $\sigma_n > 0$ گویند، هرگاه $(X_n - \mu_n)/\sigma_n \xrightarrow{d} N(0, 1)$ ، یعنی هرگاه تابع چگالی احتمال $(X_n - \mu_n)/\sigma_n$ به توزیع $N(0, 1)$ همگرا باشد. در این صورت می‌نویسیم $X_n \sim AN(\mu_n, \sigma_n^2)$.

فرض کنید T_n یک آماره است که بستگی به n (اندازه نمونه) دارد و برای n های بزرگ تقریباً دارای توزیع نرمال با میانگین θ و انحراف معیار \sqrt{n}/σ است. به عبارت دقیق‌تر، وقتی $n \rightarrow \infty$ ، تابع چگالی احتمال $\sqrt{n}(T_n - \theta)/\sigma$ به توزیع $N(0, 1)$ همگراست. یعنی

$$T_n \sim AN(\theta, \sigma^2/n). \quad (1)$$

تحت شرایط فوق اگر $(\cdot)g$ تابعی از T_n باشد به طوری که $g(\theta)$ حداقل دوبار بر حسب θ مشتق پذیر باشد آنگاه

$$g(T_n) \sim AN[g(\theta), \frac{\sigma^2}{n}(g'(\theta))^2] \quad (2)$$

که در آن $dg(\theta)/d\theta = g'(\theta) = dg(t)/dt$. برای اثبات نتیجه (۲) ملاحظه کنید که بسط تیلور تابع g حول $\theta = t$ عبارت است از

$$g(t) = g(\theta) + (t - \theta)g'(\theta) + (t - \theta)^2 g''(\theta^*)/2$$

که در آن $t \leq \theta^* \leq \theta$ و $g''(\theta)$ مشتق دوم $g(\theta)$ در نقطه $\theta = \theta^*$ است. در صورتی که در رابطه اخیر متغیر تصادفی T_n را به جای t قرار دهیم، داریم:

$$\begin{aligned} \sqrt{n}[g(T_n) - g(\theta)] &= \sqrt{n}(T_n - \theta)g'(\theta) \\ &+ \sqrt{n}(T_n - \theta)^2 g''(\theta^*)/2. \end{aligned}$$

اما وقتی $n \rightarrow \infty$ ، جمله دوم یعنی $\sqrt{n}(T_n - \theta)^2 g''(\theta^*)/2$ در احتمال به صفر می‌گراید. بنابراین دو متغیر تصادفی $\sqrt{n}(T_n - \theta)g'(\theta)$ و $\sqrt{n}[g(T_n) - g(\theta)]$ دارای توزیع حدی یکسانی هستند و لذا نتیجه (۲) حاصل می‌شود.

با توجه به اینکه $\sigma^2(\theta) = \sigma^2$ و $g'(\theta)$ به طور کلی به پارامتر نامعلوم θ بستگی دارند، واریانس مجانبی $(T_n - \theta)$ نیز نامعلوم

عددی، استنباط آماری انجام داد. از آنجا که تکیه‌گاه γ شامل صفر است، واریانس دقیق $\text{ln}(Y)$ وجود ندارد. ولی با افزایش n احتمال $\gamma = 0$ به سرعت به صفر می‌گراید.

(ب) برای تعیین توزیع مجانبی Y ابتدا توزیع مجانبی Y/n را تعیین می‌کنیم. در اینجا داریم $g(\theta) = 1/\theta = 1/\mu$ و $T_n = Y/n$. لذا بر اساس روش دلتا می‌توان نوشت

$$\text{n}/Y \sim AN[1/\mu, 1/n\mu^2].$$

از رابطه‌ی اخیر به سهولت رابطه‌ی زیر نتیجه می‌شود:

$$1/Y \sim AN(1/m, 1/m^2). \quad (6)$$

بنابراین فاصله $1/y \pm 1/96(1/y)^{1/2}$ یک فاصله اطمینان ۹۵ درصد برای $1/m$ می‌باشد. که در آن y یافته‌ای از Y و برآورد ML پارامتر $m = n\mu$ است.

اگر برای داده‌های یک جدول توافقی، یک الگوی اشباع شده‌ی خطی را با پاسخ پواسون و تابع پیوند وارون در نظر گیریم، با استفاده از (۶)، بدون توصل به روش‌های عددی، می‌توان نسبت به پارامترهای الگو استنباط آماری انجام داد.

(ج) برای تعیین توزیع مجانبی \sqrt{Y} ابتدا توزیع مجانبی $\sqrt{Y/n}$ را تعیین می‌کنیم. در اینجا داریم $g(\theta) = \sqrt{\theta} = \sqrt{\mu}$ و $T_n = Y/n$. لذا بر اساس روش دلتا می‌توان نوشت

$$\sqrt{Y/n} \sim AN[\sqrt{\mu}, 1/4n]$$

در نتیجه

$$\sqrt{Y} \sim AN[\sqrt{m}, 1/4]. \quad (7)$$

۳ کاربردهایی از روش دلتا برای تعیین توزیع مجانبی

(i) فرض کنید متغیر تصادفی Y مجموع n پواسون مستقل هر یک با میانگین $\mu > 0$ است. بنابراین $Y \sim \text{Poisson}(n\mu)$ و می‌خواهیم توزیع مجانبی هر یک از آماره‌های زیر را به دست آوریم.

$$\sqrt{Y} \quad (ج) \quad 1/Y \quad (ب) \quad \text{ln}(Y) \quad (الف)$$

براساس قضیه حد مرکزی می‌دانیم که

$$Y/n \sim AN[\mu, \mu/n]$$

(الف) برای تعیین توزیع مجانبی $\text{ln}(Y)$ ابتدا توزیع مجانبی $\text{ln}(Y/n)$ را تعیین می‌کنیم. در اینجا بر اساس نمادهای اختصاری بخش قبل داریم

$$T_n = Y/n, \quad \theta = \mu$$

$$\sigma^2(\theta) = \sigma^2(\mu) = \mu$$

$$\frac{\sigma^2(\mu)}{n} (g'(\mu))^2 = \frac{\mu}{n} \left(\frac{1}{\mu}\right)^2 = \frac{1}{n\mu}.$$

بنابراین بر اساس روش دلتا می‌توان نوشت

$$\text{ln}(Y/n) \sim AN[\ln(\mu), 1/n\mu].$$

از رابطه‌ی اخیر بلافاصله رابطه زیر نتیجه می‌شود:

$$\text{ln}(Y) \sim AN[\ln(m), 1/m] \quad (5)$$

که در آن $m = E(Y) = n\mu$. بنابراین فاصله ۹۵ درصد برای $\text{ln}(y) \pm 1/96(1/y)^{1/2}$ یک فاصله اطمینان ۹۵ درصد برای $\text{ln}(m)$ است که در آن y یافته‌ای از Y و برآورد ML پارامتر m می‌باشد.

با استفاده از (۵) می‌توان نسبت به پارامترهای الگوی خطی اشباع شده، بدون توصل به روش‌های

بنابراین بر اساس روش دلتا داریم

$$\ln\left(\frac{Y}{n-Y}\right) \sim AN\left[\ln\left(\frac{\pi}{1-\pi}\right), \frac{1}{n\pi(1-\pi)}\right] \quad (8)$$

لذا فاصله زیریک فاصله اطمینان ۹۵ درصد برای

$$\text{لذا } \text{Logit}(\pi) = \ln(\pi/(1-\pi))$$

$$\ln\left(\frac{y}{n-y}\right) \pm 1/96 \frac{n}{y(n-y)}$$

که در آن y یافته‌ای از Y و برآورد ML پارامتر $n\pi$ است.

با استفاده از (۸)، بدون توصل به روش‌های عددی می‌توان راجع به پارامترهای الگوی اشباع شده‌ی لجیت با متغیرهای توضیحی رسته‌ای استنباط آماری انجام داد. از آنجا که احتمال $Y = n$ یا $Y = 0$ مثبت است، لجیت می‌تواند برابر $-\infty$ یا $+\infty$ شود. ولی با افزایش n احتمال $Y = n$ یا $Y = 0$ به سرعت به صفر می‌گراید.

(ب) برای تعیین توزیع مجانبی $\sin^{-1}(\sqrt{Y/n})$ داریم

$$g(\theta) = \sin^{-1}(\sqrt{\theta}) = \sin^{-1}(\sqrt{\pi})$$

بنابراین بر اساس روش دلتا می‌توان نوشت

$$\sin^{-1}(\sqrt{Y/n}) \sim AN[\sin^{-1}(\sqrt{\pi}), 1/4n] \quad (9)$$

ملاحظه می‌شود که برای پاسخ‌های دو جمله‌ای تبدیل $\sin^{-1}(\sqrt{Y/n})$ طریقه‌ای برای تثبیت واریانس است.

ملاحظه می‌شود که واریانس مجانبی \sqrt{Y} عددی ثابت است. بنابراین برای پاسخ‌های پواسون، تبدیل ریشه دوم طریقه‌ای برای تثبیت واریانس است.

(ii) فرض کنید متغیر تصادفی Y مجموع n برنولی مستقل هر یک با احتمال موفقیت $\pi > 0$ است. بنابراین $Y \sim Bin(n, \pi)$ و می‌خواهیم توزیع مجانبی هر یک از دو آماره زیر را به دست آوریم.

$$\text{(الف) } \ln \frac{Y}{n-Y} \quad \text{(ب) } \sin^{-1}(\sqrt{Y/n})$$

براساس قضیه حد مرکزی می‌دانیم که برای $1 < \pi < 0$ داریم

$$\frac{Y}{n} \sim AN\left[\pi, \frac{\pi(1-\pi)}{n}\right]$$

(الف) برای تعیین توزیع مجانبی $\ln(Y/(n-Y))$ بر اساس نمادهای اختصاری بخش ۲ داریم

$$Tn = Y/n, \quad \theta = \pi$$

$$\sigma^2(\theta) = \sigma^2(\pi) = \pi(1-\pi)$$

$$g(\theta) = g(\pi) = \pi/(1-\pi)$$

مراجع

[1] Agresti, A., *Categorical Data Analysis*, John Wiley & Sons, (1990).

[2] Bishop, Y. M. M., Fienberg, S. E. and Holland, P. W., *Discrete Multivariate Analysis*, MIT Press, Cambridge, MA, (1975).

[3] *Encyclopedia of Statistical Sciences*, Vol. 8, John Wiley & Sons, (1981).

[4] Rao, C. R., *Linear Statistical Interface and its Applications*, 2nd Ed., John Wiley & Sons, (1973).

[5] ایران‌پناه، ن. و پاشا، ع. (۱۳۷۶)، آشنایی با الگوریتم بوت استرپ، اندیشه آماری، سال دوم، شماره اول.

[6] نیرومند، ح. (۱۳۷۶)، آشنایی با جک نایف، اندیشه آماری، سال دوم، شماره دوم.

یک نشریه علمی انجمن پزشکان قلب آمریکا [۱]، آماری را پیرامون استفاده غیرصحیح از علم آمار ارائه داده است که جای تأمل دارد! در این نشریه آمده است که:

«تقریباً نیمی از مقالاتی که در نشریات علمی پزشکی چاپ شده و در آنها تجزیه و تحلیل آماری انجام گرفته است، متاسفانه از روش‌های نادرست استفاده نموده‌اند.

این مقاله باعث شد که هیأت تحریریه نشریه در داوری مقالات بیشتر دقت کنند، اما این سؤال باقی مانده است که آیا سایر نشریات هم چنین خواهند کرد؟

بسیار جای تأسف است که محققین علم آمار با استفاده از پیشرفته‌ترین مباحث ریاضی از جمله؛ نظریه اندازه، توپولوژی، آنالیز تابعی و فرایندهای تصادفی در صدد پیشبرد علم آمار و تقویت مبانی این علم هستند و از طرف دیگر مقالات علمی هنوز در استفاده از روش‌های ساده‌ای مانند آزمون t دچار اشتباه می‌شوند. تأسف‌بارتر این است که نه تنها دانشمندان سایر علوم آزمونهای آماری را صحیح به کار نمی‌برند، بلکه حتی آمارشناسان نیز چنین خطاهایی را مرتكب می‌شوند!»

[1] Glantz (1980), *How to Detect, Correct and Prevent Errors in the Medical Literature*, Biostatistics , Vol. 61, 1-7.