

## تبدیل مشخصه‌ای توزیع نرمال

نلی کریمانی

گروه آمار، دانشگاه پیام نور، ایران.  
neli\_karimani@yahoo.com

حسین جباری خامنه‌ای

گروه آمار، دانشگاه تبریز  
h\_jabbari@tabrizu.ac.ir

لیدر نوایی

گروه آمار، دانشگاه پیام نور، ایران.  
leadernavaei@yahoo.com

### چکیده

در این مقاله ما یک تبدیل جدید از متغیرهای تصادفی  $RVs$  که توزیع نرمال را مشخص می‌کند، پیشنهاد می‌کنیم. آن تبدیل به ما اجازه می‌دهد که  $n$  تا متغیر تصادفی نرمال هم‌توزیع و مستقل  $iid$  که دارای میانگین و واریانس نامعلوم هستند را به  $(n-2)$  تا متغیر نرمال جدید هم‌توزیع و مستقل با میانگین صفر زمانی که همان واریانس نامعلوم را در خود نگه دارد، تبدیل کنیم. این متعلق به کلاسی از تبدیلات است که برای کاهش تعدادی از پارامترهای نامعلوم یا از بین بردن تمامی آنها طراحی شده است. چند اظهار نظر قدیمی درباره روش‌هایی برای از بین بردن پارامترها در توزیع نرمال داده شده‌اند و دو کاربرد ممکن از تبدیل جدید شرح داده شده است.

واژگان کلیدی: آزمون نرمال بودن؛ پارامترهای مزاحم؛ تبدیلات متغیر تصادفی؛ مشخص سازی نرمال بودن.

### ۱. پیش‌گفتار

اولین مسئله مورد توجه، آزمون نرمال بودن می‌باشد. بیشترین مسائل نیکویی برازش که به شیوه واقعی به دست می‌آید شامل پارامترهای مزاحم می‌باشد. همانطور که فرض صفر مورد آزمون، فرض‌های مرکب هستند. آزمون‌های کلوموگروف، اندرسون-دارلینگ (۱۹۵۲) و سایر آزمون‌هایی که بر اساس انحراف  $EDF$  (تابع چگالی تجربی) هستند آزمون‌هایی خوب و دقیق می‌باشند که فقط نیاز به معلومات پارامتر دارند. تعدیل‌های بیشتری به این آزمون‌ها برای رفع این عیب (نیاز به معلومات) پیشنهاد داده شده‌اند و نیز برای بیشتر آزمون‌های دیگری که بر اساس انحراف  $EDF$  نیستند، فرمول‌بندی شده‌اند. به هر حال، تاکنون اکثریت آزمون‌های به کار رفته برای نرمال بودن (برای مثال، آزمون های جاکو-بر [۷]؛ لی لی فورس [۸]؛ شاپیرو-ویلک [۹]؛ اندرسون-دارلینگ تعدیل یافته [۱]؛ استیفن [۱۰]؛ گریجو و همکاران [۳] و غیره) راه حل‌های پیشنهادی را که درست باشند را تضمین نمی‌کنند. بنابراین، یک راه طبیعی برای رسیدن به این نیاز، کاهش فرض مرکب که معادل یک فرض ساده است می‌باشد (بوندسون [۲]؛ سورگو-سیشادری [۴]). اولین مرحله برای بدست آوردن این نیاز، به کاربرد تبدیل پیشنهاد شده اینجا بر روی داده‌های اصلی است. بدیهی است بعد از این تبدیل، یک تیمار اضافی از داده‌ها بایستی به منظور از بین بردن پارامترهای مقیاسی همچنان انجام داده شود. دومین استفاده ممکن از تبدیل مان (در اولین مرحله از تحلیل) به مسئله معروف بئرنس-فیشر مربوط می‌شود. تصور اصلی از تبدیل پیشنهاد شده آن است که به ما اجازه می‌دهد که یک مجموعه داده‌های نرمال تبدیل شده برای هر نمونه اصلی به دست آوریم که با میانگین صفر و همان واریانس‌های اصلی است، حتی اگر میانگین‌های دو نمونه اصلی به شدت متفاوت باشند. اولین کاهش به همان میانگین‌ها در مدتی که واریانس‌های اصلی به دست آمده است را حفظ می‌کند، ممکن است که روش‌های خودگردان و جایگشتی برای به دست آوردن توزیع دقیق صفر از هر آماره را به کار برد. همان طور که این تکنیک هنوز تحت بررسی است، ما این وضع را در این مقاله فعلی بیشتر بحث نمی‌کنیم.

### ۲. نتایج اصلی

تا آنجایی که می‌دانیم اولین مثال از یک تکنیک برای از بین بردن پارامتر مکانی مزاحم در هلمرت [۵، ۶] بود، که "تبدیل هلمرت" نامیده می‌شود. هنگام روبرو شدن با نمونه تصادفی به اندازه  $n$  از هر توزیع نرمال، آن تبدیل به ما اجازه می‌دهد که  $(n-1)$  مقادیر جدید هم‌توزیع و مستقل به طور نرمال توزیع شده را به دست آوریم که با میانگین صفر و همانند داده‌های اصلی دارای همان واریانس نامعلوم هستند. اگر  $X_1, \dots, X_n$ ،  $n$  تا متغیرهای تصادفی هم‌توزیع و مستقل از توزیع نرمال با میانگین  $\mu$  واریانس نامعلوم

$\sigma^2$  باشد  $X \sim N(\mu, \sigma^2)$  متغیرهای تبدیل یافته هلمرت  $Z_1, Z_2, \dots, Z_n$  همانند زیر تعریف شده‌اند:

$$Z_1 = (X_1 - X_2) / \sqrt{2}$$

⋮

$$Z_i = (X_1 + \dots + X_2 + iX_{i+1}) / \sqrt{i(i+1)}$$

$$Z_{n-1} = (X_1 + \dots + X_2 + (n-1)X_n) / \sqrt{n(n-1)}$$

که به اندازه  $(n-1)$  هم توزیع و مستقل هستند که  $Z \sim N(0, \sigma^2)$ . این تبدیل در نوع خودش مهم است، زیرا منجر می‌شود که بدانیم زمانی که  $Z_i^2 = \sum_{i=1}^n (X_i - \bar{X})^2$ ، واریانس نمونه  $s^2$  برابر است با:  $s^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n-1)$  که از توزیع کای دو  $\chi_{n-1}^2$  پیروی می‌کند و مستقل از میانگین نمونه  $\bar{X}$  می‌باشد.

**قضیه ۱.۲.** فرض کنید  $X_1, \dots, X_n$ ، برای  $n > 3$  متغیرهای تصادفی مستقل و هم توزیع باشند. تبدیل زیر را در نظر بگیرید:

$$Y_i = (X_i - \bar{X}) + (X_c - X_k) / \sqrt{2(n-2)}, \quad i = 1, 2, \dots, n, \quad i \neq c, k \quad (1)$$

جایی که  $i \leq c \neq k \leq n$  به صورت دو عدد صحیح دلخواه انتخاب شده‌اند و  $\bar{X} = \sum_{i \neq (c,k)}^n X_i / (n-2)$  که  $(n-2)$  تا متغیرهای تصادفی تبدیل شده به صورت *i.i.d* هستند:  $Y_i \sim N(0, \sigma^2)$ ،  $i = 1, 2, \dots, n$ ،  $i \neq c, k$  اگر و تنها اگر  $X_i \sim N(\mu, \sigma^2)$ . اثبات. برای قسمت شرط لازم، کافی است معادله (۱) به روش زیر دوباره مرتب‌سازی شود:

$$Y_i = \frac{(n-3)X_i - \sum_{j \neq (i,c,k)}^n X_j}{n-2} + \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \quad (2)$$

از این فرمول‌سازی مجدد واضح است که معادله (۱) یک ترکیب خطی از متغیرهای تصادفی مستقل که مجموع آنها به صورت نرمال توزیع شده است زیرا مؤلفه‌های آن به صورت نرمال مستقل هستند. به هر حال، اگر  $Y_i$  عمومی به صورت نرمال توزیع شده باشد، بر اساس قضیه‌ی لوکاکس (۱۹۵۶) که به نو به خود نتیجه‌ی گذشته به دست آمده توسط کرامر (۱۹۳۶) را تعمیم می‌دهد، که مؤلفه‌های آن بایستی همچنان به صورت نرمال توزیع شده باشد. برای قسمت شرط کافی، بایستی نشان داده شود که

الف) میانگین:  $E(Y_i) = 0$ ؛ ب) واریانس:  $Var(Y_i) = \sigma^2$ ؛ ج)  $Y_i$  ها به طور تصادفی مستقل هستند.

در الف) بوسیله ساختاری که داریم:  $E(X_i - \bar{X}) = 0$  و  $\sum_{i \neq (c,k)}^n (X_i - \bar{X}) = 0$  و چون  $E(X_j - \bar{X}) = E(X_i - \bar{X})$ ، نتیجه می‌شود که

$$E \left[ \sum_{i \neq (c,k)}^n (X_i - \bar{X}) \right] = (n-2)E(X_i - \bar{X}) = 0$$

$$E \left[ \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \right] = \frac{E(X_c) - E(X_k)}{\sqrt{2(n-2)}} = 0 \text{ علاوه بر این } E(X_i - \bar{X}) = 0$$

در ب) متغیرهای تصادفی  $\sum_{i \neq (c,k)}^n (X_i - \bar{X})^2 \sim \sigma^2 \chi_{n-3}^2$  بنابراین

$$Var(Y_i) = Var(X_i - \bar{X}) + Var \left[ \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \right] = \sigma^2$$

در ج) ابتدا بایستی ثابت شود که هر دو متغیر تصادفی عمومی  $Y_i$ ،  $Y_j$ ،  $i \neq j$  ناهمبسته هستند. با امتحان کردن رفتار  $(n-2)$  تا متغیرهای تصادفی  $(X_i - \bar{X})$  برای بررسی اینکه آنها به هم وابسته هستند، راحت است. چون  $\sum_{i \neq (c,k)}^n (X_i - \bar{X}) = 0$  بنابراین ضرب دومی در  $(X_i - \bar{X})$  و با به کار بردن عامل میانگین، به دست می‌آوریم:

$$E \left[ (X_1 - \bar{X})(X_i - \bar{X}) + \dots + (X_i - \bar{X})^2 + (X_n - \bar{X})(X_i - \bar{X}) \right] = 0$$

$$E \left[ (X_1 - \bar{X})(X_i - \bar{X}) \right] = E \left[ (X_j - \bar{X})(X_i - \bar{X}) \right], \quad \forall i \neq j$$

$$(n-3)E \left[ (X_j - \bar{X})(X_i - \bar{X}) \right] + E \left[ (X_i - \bar{X})^2 \right] = 0$$

و همچنین

که نتیجه می‌دهد:

بنابراین،  $E \left[ (X_j - \bar{X})(X_i - \bar{X}) \right] = -\sigma^2 / (n-2)$ ؛ آنگاه:

$$\begin{aligned} Cov(Y_i, Y_j) &= E \left\{ \left[ (X_i - \bar{X}) + \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \right] \left[ (X_j - \bar{X}) + \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \right] \right\} \\ &= -\frac{\sigma^2}{(n-2)} + \frac{\sigma^2}{(n-2)} = 0 \end{aligned}$$

دوماً، برای اثبات استقلال تصادفی هر زوج  $(Y_i, Y_j)$ ,  $i \neq j$  باید نشان داده شود که متغیرهای تصادفی به صورت نرمال دو متغیره توزیع شده‌اند. در واقع اگر توزیع‌های حاشیه ای از نرمال دو متغیره ناهمبسته باشند، آنگاه آنها از هم مستقل هستند. برای نشان دادن این استقلال کافی است، نشان دهیم که متغیر تصادفی عمومی  $\zeta = aY_i + bY_j$  برای هر زوج  $(a, b) \in [R^2 - (0, 0)]$  به صورت نرمال توزیع شده‌اند. با استفاده از رابطه‌ی (۲) داریم:

$$\begin{aligned} \zeta &= a \left[ \frac{(n-3)X_i - \sum_{m \neq (i,c,k)}^n X_m}{n-2} + \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \right] \\ &+ b \left[ \frac{(n-3)X_j - \sum_{m \neq (i,c,k)}^n X_m}{n-2} + \frac{(X_c - X_k)}{\sqrt{2(n-2)}} \right] \\ &= \frac{[(n-3)a-b]X_i + [(n-3)b-a]X_j - (a+b)\sum_{m \neq (i,c,k)}^n X_m}{n-2} + \frac{(a+b)(X_c - X_k)}{\sqrt{2(n-2)}} \end{aligned}$$

حال ما بایستی بررسی کنیم که آیا در ادامه سه معادله‌ی وابسته به ضرایب در  $(a, b)$ ، برای حداقل یک جفت از اعداد حقیقی بواسطه  $(0, 0)$  به طور همزمان برقرار هستند:  $a+b=0$ ;  $(n-3)b-a=0$ ;  $(n-3)a-b=0$ ؛ چون  $n > 3$ ، بررسی اینکه تنها زوج  $(0, 0)$  به طور مشترک در سه معادله صدق می‌کند، آسان است. این برهان قضیه را  $\square$  به اتمام می‌رساند.

توجه ۱.۲. میانگین نمونه تبدیل شده  $Y_i$  ها عبارتند از:  $\bar{Y} = \frac{X_c - X_k}{\sqrt{2(n-2)}}$  این نتیجه‌ی ساده می‌تواند هنگام روبرو شدن با نمونه‌گیری مجدد و روش‌های جایگشتی مفید باشد. ارزش مشاهده آن را دارد که عمل ثابت بودن دوعدد صحیح دلخواه  $c$  و  $k$  معادل است با معرفی یک عنصر خارجی از تصادف بودن.

### ۳. مقایسه توان بین تبدیلات $Z$ و $Y$

اندازه‌های دو نمونه  $n = 20, 50$  در نظر گرفته شده‌اند و خطای نوع اول برای داشتن یک مقایسه شفاف‌تر در  $(\alpha = 0/20)$  ثابت در نظر گرفته شده است. برای هر توزیع و اندازه نمونه، نمونه‌های تصادفی  $500$  تایی استخراج شده‌اند و دو تبدیل برای هر نمونه بکاربرده شده و آماره آزمون کولموگروف استفاده شده است. نتیجه‌ها نایبستی به عنوان عملکردی از یک آزمون واقعی برای نرمال بودن تفسیر شوند، بلکه تنها به عنوان یک روش برای مقایسه تبدیل  $Z$  در مقابل تبدیل  $Y$  می‌باشند. در واقع، پارامتر مقیاسی معمولاً مجهول می‌باشد و برای وقتی که آماره آزمون در جایی که پارامترها در  $H_0$  باید به صورت کاملاً معلوم استفاده شده باشند، یک تیمار اضافی برای آزمون نرمال بودن لازم است. نتیجه‌های شبیه‌سازی توصیف شده در جدول ۱ ارائه شده‌اند. برای  $n = 20$  می‌تواند مشاهده شود که  $Y$  به طور واضح و یکنواخت قوی‌تر از  $Z$  است؛ با  $n = 50$  هنوز نسبت به  $Z$  قوی‌تر است اما اختلاف بین دو تبدیل تمایل به کاهش دارد. بنابراین نتیجه‌ها در موافقت با فرضیات قرار داده شده بالا ظاهر می‌شوند. بعضی از نویسندگان (به طور مثال، سورگو و همکاران [۴]) در نظر گرفته‌اند که تبدیل  $Z$  اولین قدم برای آزمون نرمال بودن می‌باشد. چند نتیجه بهتر می‌تواند احتمالاً بوسیله جایگزینی  $Z$  با تبدیل  $Y$  ها به دست آید.

جدول ۱: مقایسه توان بین تبدیلات  $Z$  و  $Y$  تحت غیر نرمال بودن با استفاده از آزمون کولموگروف ( $\alpha = 0/20$ )

|                        | $n=20$ |       | $n=50$ |       |
|------------------------|--------|-------|--------|-------|
|                        | $Y_s$  | $Z_s$ | $Y_s$  | $Z_s$ |
| یکنواخت                | 28/7   | 23/2  | 33/0   | 29/1  |
| نمایی                  | 58/7   | 45/1  | 90/3   | 81/7  |
| لاپلاس                 | 42/0   | 30/6  | 45/5   | 42/0  |
| استیودنت $t(\epsilon)$ | 34/6   | 26/6  | 40/3   | 37/1  |
| کاما(۲)                | 46/0   | 32/4  | 78/0   | 72/6  |
| بتا (۲,۱)              | 32/8   | 27/2  | 39/3   | 34/3  |

## مراجع

- [1] T. W. Anderson and D. A. Darling, *Asymptotic theory of certain "goodness-of-fit" criteria based on stochastic processes.*, Ann. Mathemat. Statist. **23** (1952), 193–212.
- [2] L. Bondesson, *To reduce a composite hypothesis to a simple one by sampling from the structural or conditional distribution*, Scand. J. Statist. **9(3)** (1982), 129–138.
- [3] S. Grego, I. Pegoraro and F. Pesarin, *Problemi di verifica d'ipotesi di normalità: approfondimenti ed innovazioni*, Unpublished dissertation, Padua University: Padua, Italy (1997).

- [4] M. Csorgo and V. Seshadri, *On the problem of replacing composite hypotheses by equivalent simple ones*, Rev. Int. Statist. Institut. **38** (1970), 351–68.
- [5] F. R. Helmert, *Über die Berechnung des wahrscheinlichen Fehlers aus einenendlichen Anzahl wahrer Beobachtungsfehler*, Zeitschrift fuer angewandte Mathematik und Physik **20** (1875), 300–303.
- [6] F. R. Helmert, *Über die Wahrscheinlichkeit der Potenzsummen der Beobachtungsfehler und über einige damit im Zusammenhang stehende Fragen*, Zeitschrift fuer angewandte Mathematik und Physik **21** (1876), 192–218.
- [7] C. M. Jarque and A. K. Bera, *A test for normality of observations and regression residuals*, Rev. Int. Statist. Institut. **55** (1987), 163–172.
- [8] H. Lilliefors, *On the Kolmogorov–Smirnov test for normality with mean and variance unknown*, J. Amer. Statist. Assoc. **62** (1967), 399–402.
- [9] S. S. Shapiro and M. B. Wilk, *An analysis of variance test for normality (complete samples)*, Biometrika **52** (1965), 591–611.
- [10] M. A. Stephens, *EDF statistics for goodness of fit and some comparisons*, J. Amer. Statist. Assoc. **69** (1974), 730–737.