



یک آزمون جدید نیکویی برازش با استفاده از تابع مشخصه‌ی تجربی

مینا توحیدی* و مهدی سلمان‌پور

دانشگاه شیراز

چکیده. تابع مشخصه در تعیین تابع توزیع احتمال نقش مهم و کلیدی دارد و به‌طور منحصربه‌فردی تابع احتمال یک متغیر تصادفی به‌وسیله‌ی آن تعیین می‌شود. اگر $c(t)$ و $c_0(t)$ به ترتیب تابع مشخصه‌ی توابع توزیع $F(\cdot)$ و $F_0(\cdot)$ باشند، آنگاه فرض صفر $F(x) = F_0(x), \forall x \in \mathbb{R}$ را می‌توان به فرض $c_n(t) = \frac{1}{n} \sum_{j=1}^n \exp\{itX_j\}$ تبدیل کرد و از تابع مشخصه‌ی تجربی، به همین دلیل، بین سال‌های ۱۹۷۲ تا ۱۹۹۳ بسیاری در آزمون نیکویی برازش توزیع استفاده نمود. از پژوهشگران از تابع مشخصه‌ی تجربی برای آزمون فرض‌های مختلف آماری استفاده کردند. در بیشتر آزمون‌های معرفی شده، مقایسه‌هایی بین $c_n(t)$ و $c_0(t)$ ، برای تعداد کمی از مقادیر t ، صورت گرفته است و این امر موجب سازگار نبودن آزمون‌ها شده است. ما در این مقاله سعی کرده‌ایم که مقایسه‌ی بین تابع مشخصه‌ی تجربی و تابع مشخصه‌ی توزیع جامعه تحت فرض صفر را برای تعداد زیادی از مقادیر t انجام دهیم و بدین ترتیب آزمونی را معرفی کرده‌ایم که بسیار پرتوان‌تر از آزمون‌های ناپارامتری قبلی بوده است.

واژگان کلیدی. آزمون نیکویی برازش؛ آزمون سازگار؛ بردارهای ویژه؛ تابع مشخصه؛ تابع مشخصه‌ی تجربی؛ روش مؤلفه‌های اصلی؛ قضیه‌ی حد مرکزی چندمتغیره؛ مقادیر ویژه.

* نویسنده‌ی عهده‌دار مکاتبات.

۱ مقدمه

فرض کنید که X متغیری تصادفی در فضای احتمال (Ω, A, P) با تابع توزیع $F(x)$ باشد و بخواهیم فرض X دارای توزیع شناخته شده $F_0(x)$ را آزمون کنیم. توزیع $F_0(x)$ را توزیع برانزده بر X (یا توزیع جور با یافته‌های X) می‌گویند. بر اساس n مشاهده مستقل و هم‌توزیع ناشناخته $F(x)$ ، آزمون‌هایی را که برای آزمون فرض $H_0: F(x) = F_0(x), \forall x \in \mathbb{R}$ در مقابل فرض $H_A: F(x) \neq F_0(x), \exists x \in \mathbb{R}$ به کار برده می‌شوند، آزمون‌های نیکویی برآزش توزیع می‌نامند.

تاکنون آزمون‌های ناپارامتری بسیاری برای آزمون فرض‌های بالا معرفی شده است که از جمله می‌توان به آزمون‌های کولموگروف-اسمیرنوف، کرامر-ون میسن، اندرسن-دارلینگ و واتسون اشاره کرد. با این وجود، آماردانان در مسائل مختلف آماری، همواره به دنبال یافتن آزمون‌های پرتوان‌تری بوده‌اند. همچنین در بعضی موارد، صورت بسته‌ای برای تابع توزیع $F_0(x)$ وجود ندارد و در نتیجه انجام آزمون‌های معروف یاد شده، امکان‌پذیر نیست. این مشکلات آماردانان را بر آن داشت که از تابع مشخصه‌ی تجربی در آزمون فرض‌های آماری استفاده کنند.

می‌دانیم که تابع مشخصه در تعیین تابع توزیع احتمال نقش مهم و کلیدی دارد و به وسیله‌ی آن می‌توان تابع احتمال یک متغیر تصادفی را تعیین کرد. به همین دلیل، بین سال‌های ۱۹۷۲ تا ۱۹۹۳ پژوهشگران در مقاله‌های گوناگون، با استفاده از تابع مشخصه‌ی تجربی، آزمون‌های نیکویی برآزش را با در نظر گرفتن توزیع‌های مشخص $F_0(x)$ ارائه کرده‌اند.

هدف ما در این مقاله، طراحی یک آزمون نیکویی برآزش کلی (با در نظر گرفتن هر توزیع مشخص $F_0(x)$) بر اساس تابع مشخصه‌ی تجربی است که پرتوان‌تر از آزمون‌های ارائه شده‌ی قبلی باشد. در بخش دوم، به معرفی تابع مشخصه‌ی تجربی و برخی از ویژگی‌های آن خواهیم پرداخت. آماره‌ی کلی آزمون و توزیع حدی آن را در بخش سوم ارائه خواهیم کرد و در بخش چهارم بر اساس تابع مشخصه‌ی تجربی و استفاده از روش مؤلفه‌های اصلی آماره‌ی جدیدی را معرفی خواهیم کرد. در پایان در بخش پنجم، با استفاده از شبیه‌سازی مونت کارلو در یک مثال، پرتوان‌تر بودن آزمون جدید را نسبت به آزمون‌های قبلی نشان خواهیم داد.

۲ تابع مشخصه‌ی تجربی و برخی ویژگی‌های آن

همان‌طور که می‌دانید تابع مشخصه‌ی یک متغیر تصادفی یک بعدی با تابع توزیع $F(x)$ به صورت زیر تعریف می‌شود:

$$c(t) = E[\exp\{itX\}] = \int \exp\{itx\}dF(x),$$

و به طور منحصر به فردی تابع توزیع توسط تابع مشخصه تعیین می‌شود. این تابع دارای دو خاصیت مهم زیر است:

آ) $c(t) = E\{\cos(tX)\} + iE\{\sin(tX)\}$ که در آن $E\{\cos(tX)\}$ را قسمت حقیقی $(\operatorname{Re} c(t))$ و $E\{\sin(tX)\}$ را قسمت موهومی $(\operatorname{Im} c(t))$ تابع مشخصه می‌نامیم.

ب) تابع مشخصه برای هر متغیر تصادفی X موجود است و برای هر $t \in \mathbb{R}$ ، $|c(t)| \leq 1$.

یک برآوردگر سازگار برای $c(t)$ ، با در دست داشتن یک نمونه‌ی تصادفی X_1, \dots, X_n از توزیع $F(x)$ ، تابع مشخصه‌ی تجربی است که به صورت زیر تعریف می‌شود:

$$c_n(t) = \frac{1}{n} \sum_{j=1}^n \exp\{itX_j\} = \frac{1}{n} \sum_{j=1}^n \cos(tX_j) + i \left\{ \frac{1}{n} \sum_{j=1}^n \sin(tX_j) \right\}.$$

یکی از مهم‌ترین خواص احتمالاتی تابع مشخصه‌ی تجربی در لم زیر آمده است:

لم ۱ اگر $c_n(t)$ تابع مشخصه‌ی تجربی نمونه‌ی تصادفی X_1, \dots, X_n از توزیع $F(x)$ با تابع مشخصه‌ی $c(t)$ باشد، آنگاه همگرایی‌های زیر با احتمال ۱ برقرار است:

$$c_n(t) \rightarrow c(t),$$

$$\operatorname{Re} c_n(t) \rightarrow \operatorname{Re} c(t),$$

$$\operatorname{Im} c_n(t) \rightarrow \operatorname{Im} c(t),$$

که در آن‌ها $\operatorname{Re} c_n(t) = \frac{1}{n} \sum_{j=1}^n \cos(tX_j)$ و $\operatorname{Im} c_n(t) = \frac{1}{n} \sum_{j=1}^n \sin(tX_j)$ است.

برهان. با استفاده از قانون قوی اعداد بزرگ برای متغیرهای تصادفی $\exp\{itX_1\}, \dots, \exp\{itX_n\}$ و با توجه به متناهی بودن $|c_n(t)|$ و $|c(t)|$ ، این لم به سادگی اثبات می‌شود.

لم اخیر بیان می‌کند که در صورت بزرگ بودن اندازه‌ی نمونه، مقدار $c_n(t)$ به $c(t)$ نزدیک خواهد بود. با توجه به این که فرض $H_0: F(x) = F_0(x), \forall x$ را می‌توان به صورت $H_1: c(t) = c_0(t), \forall t \in \mathbb{R}$ نوشت، پس می‌توان آماره‌ی آزمون نیکویی برازش توزیع را بر اساس تابع مشخصه‌ی تجربی $c_n(t)$ ، طرح ریزی کرد.

۳ آزمون‌های نیکویی برازش بر اساس تابع مشخصه‌ی تجربی

برای بررسی آزمون‌های نیکویی برازش بر اساس تابع مشخصه‌ی تجربی ابتدا لازم است مفهوم سازگاری آزمون را روشن سازیم.

تعریف ۱ یک آزمون فرض را سازگار گویند، هرگاه با افزایش اندازه‌ی نمونه، توان آزمون به سمت یک میل کند. به عبارت دیگر توان مجانبی آزمون برابر یک شود.

فیگین و هتکت (۱۹۷۶) آزمون فرض H_0 را با مقایسه‌ی $\text{Im } c_n(t)$ (یا $\text{Re } c_n(t)$) با $\text{Im } c_0(t)$ (یا $\text{Re } c_0(t)$) برای یک مقدار ثابت t انجام دادند. آن‌ها با انتخاب یک t مناسب، به آزمون پرتوان‌تر از آزمون کرامر-ون میسر دست یافتند. فیورورگر و موریکا (۱۹۷۷) یک آزمون را برای فرض تقارن ارائه کردند. آن‌ها آماره‌ی خود را بر اساس انتگرال توان دوم اختلاف بین $\text{Im } c_n(t)$ و صفر بیان کردند (تابع مشخصه، یک تابع حقیقی است اگر و فقط اگر توزیع مربوط متقارن باشد) سپس هال و ولش (۱۹۸۳) آزمون‌های را برای آزمون فرض نرمال بودن توزیع ارائه دادند. در بیش‌تر آزمون‌های معرفی شده در مقاله‌های پژوهشگران، مقایسه‌هایی بین $c_n(t)$ و $c_0(t)$ برای برخی مقادیر t انجام شده است، در نتیجه آزمون‌های ارائه شده سازگار نیستند. در سال ۱۹۸۸ بارینس و هنز آماره‌ای را بر اساس انتگرال توان دوم اختلاف بین $c_n(t)$ و $c_0(t)$ طراحی کردند و سازگاری آزمون مطرح شده را اثبات کردند، اما محاسبه‌ی مقدار آماره بسیار مشکل بود.

در سال ۱۹۹۷ فن برای رسیدن به یک آزمون سازگار، آماره‌ای را بر اساس توان دوم اختلاف بین $c_n(t)$ و $c_0(t)$ در نقاط t_1, \dots, t_m معرفی کرد و برای اثبات سازگاری آزمون، فرض کرد که $m_n \equiv m \rightarrow +\infty$ وقتی که $n \rightarrow +\infty$. برای توضیح چگونگی انجام این آزمون، فرض کنید بردار $\mathbf{t}_m = (t_1, \dots, t_m)$ حاوی m نقطه‌ی یاد شده، باشد. با تعریف بردارهای $\mathbf{Z}_n(\mathbf{t}_m)$ و $\mathbf{Z}(\mathbf{t}_m)$ به شکل زیر:

$$\mathbf{Z}_n(\mathbf{t}_m) = \begin{pmatrix} \text{Re } c_n(t_1) \\ \text{Re } c_n(t_2) \\ \vdots \\ \text{Re } c_n(t_m) \\ \text{Im } c_n(t_1) \\ \vdots \\ \text{Im } c_n(t_m) \end{pmatrix}, \quad \mathbf{Z}(\mathbf{t}_m) = \begin{pmatrix} \text{Re } c(t_1) \\ \vdots \\ \text{Re } c(t_m) \\ \text{Im } c(t_1) \\ \vdots \\ \text{Im } c(t_m) \end{pmatrix}.$$

به راحتی دیده می‌شود که $\mathbf{Z}_n(\mathbf{t}_m)$ برآوردگری ناریب و سازگار برای $\mathbf{Z}(\mathbf{t}_m)$ است، اگر تابع مشخصه‌ی توزیع جامعه $c(t)$ باشد. بنا بر این، به منظور بررسی فرض صفر، باید بردارهای $\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}(\mathbf{t}_m)$ را مورد ارزیابی قرار داد. تفاضل این دو بردار برابر است با:

$$\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}(\mathbf{t}_m) = \frac{1}{n} \sum_{j=1}^n \begin{pmatrix} \cos(t_1 X_j) - E\{\cos(t_1 X_j)\} \\ \vdots \\ \cos(t_m X_j) - E\{\cos(t_m X_j)\} \\ \sin(t_1 X_j) - E\{\sin(t_1 X_j)\} \\ \vdots \\ \sin(t_m X_j) - E\{\sin(t_m X_j)\} \end{pmatrix} = \frac{1}{n} \sum_{j=1}^n \mathbf{Y}_j(\mathbf{t}_m),$$

که در آن

$$(\mathbf{Y}_j(\mathbf{t}_m))' = (\cos(t_1 X_j) - E\{\cos(t_1 X_j)\}, \dots, \sin(t_m X_j) - E\{\sin(t_m X_j)\}).$$

تحت فرض صفر، امید ریاضی بردار تصادفی $\mathbf{Y}_j(\mathbf{t}_m)$ برابر با بردار صفر $2m$ بعدی است. اگر ماتریس واریانس-کوواریانس $\mathbf{Y}_j(\mathbf{t}_m)$ را با Ω_0 نشان دهیم، این ماتریس را می‌توان به صورت زیر افزایش کرد:

$$\Omega_0 = \begin{pmatrix} \Omega_{0,11} & \Omega_{0,12} \\ \Omega_{0,21} & \Omega_{0,22} \end{pmatrix},$$

که درایه‌های این زیر ماتریس‌ها برابر است با:

$$\begin{aligned} (\Omega_{0,11})_{ij} &= \text{cov}\{\cos(t_i X_1), \cos(t_j X_1)\} \\ &= E\{\cos(t_i X_1) \cos(t_j X_1)\} - E\{\cos(t_i X_1)\} E\{\cos(t_j X_1)\} \\ &= \frac{1}{2} [\text{Re } c.(t_i + t_j) + \text{Re } c.(t_i - t_j) \\ &\quad - 2\text{Re } c.(t_i) \text{Re } c.(t_j)], \quad i, j = 1, \dots, m, \quad i \neq j \end{aligned}$$

و به صورت مشابه:

$$\begin{aligned}
 (\Omega_{\cdot 11})_{ii} &= \text{var} \{ \cos(t_i X_1) \} \\
 &= \frac{1}{\nu} \left\{ 1 + \text{Re } c_{\cdot}(\nu t_i) - \nu (\text{Re } c_{\cdot}(t_i))^2 \right\}, \quad i = 1, \dots, m \\
 (\Omega_{\cdot 12})_{ij} &= \text{cov} \{ \cos(t_i X_1), \sin(t_j X_1) \} \\
 &= \frac{1}{\nu} \left\{ \text{Im } c_{\cdot}(t_i + t_j) - \text{Im } c_{\cdot}(t_i - t_j) \right. \\
 &\quad \left. - \nu \text{Re } c_{\cdot}(t_i) \text{Im } c_{\cdot}(t_j) \right\}, \quad i, j = 1, \dots, m \\
 (\Omega_{\cdot 22})_{ij} &= \text{cov} \{ \sin(t_i X_1), \sin(t_j X_1) \} \\
 &= \frac{1}{\nu} \left\{ \text{Re } c_{\cdot}(t_i - t_j) - \text{Re } c_{\cdot}(t_i + t_j) \right. \\
 &\quad \left. - \nu \text{Im } c_{\cdot}(t_i) \text{Im } c_{\cdot}(t_j) \right\}, \quad i, j = 1, \dots, m \\
 (\Omega_{\cdot 22})_{ii} &= \text{var} \{ \sin(t_i X_1) \} \\
 &= \frac{1}{\nu} \left\{ 1 - \text{Re } c_{\cdot}(\nu t_i) - \nu (\text{Im } c_{\cdot}(t_i))^2 \right\}, \quad i = 1, \dots, m \\
 \Omega_{\cdot 21} &= \Omega'_{\cdot 12}.
 \end{aligned}$$

اگر $\mathbf{Z}_{\cdot}(\mathbf{t}_m)$ مقدار $\mathbf{Z}(\mathbf{t}_m)$ تحت فرض صفر باشد، آنگاه آماره‌ی آزمون که در حقیقت همان آماره‌ی والد برای بردار تصادفی $\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_{\cdot}(\mathbf{t}_m)$ می‌باشد، به صورت زیر معرفی می‌شود:

$$(1) \quad T_n^{\circ} = (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_{\cdot}(\mathbf{t}_m))' \Omega_{\cdot}^{-1} (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_{\cdot}(\mathbf{t}_m)).$$

فن (۱۹۹۷) با تعریف ماتریس وزن قطری $W(\mathbf{t}_m, \theta)$ وقتی که اعضای روی قطر $(w_j, j = 1, \dots, 2m)$ غیر منفی باشند و W وابسته به m یا θ (نه لزوماً هر دو) بوده و $\theta = \theta_n \rightarrow 0$ وقتی که $n \rightarrow +\infty$ آماره‌ی جدید T_n را معرفی کرد:

$$(2) \quad T_n = (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_{\cdot}(\mathbf{t}_m))' \Omega_{\cdot}^{-\frac{1}{\nu}} W(\mathbf{t}_m, \theta) \Omega_{\cdot}^{-\frac{1}{\nu}} (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_{\cdot}(\mathbf{t}_m)),$$

و در قضیه‌ی توزیع حدی T_n را ارائه داد:

قضیه‌ی ۱ (فن (۱۹۹۷)) فرض کنید W ماتریسی قطری باشد که به θ وابسته نباشد و m عدد صحیح مثبت و متناهی باشد، در این صورت تحت فرض صفر nT_n دارای توزیع حدی $\sum_{j=1}^{2m} w_j \chi_{(1),j}^2$ خواهد

بود که در آن $\chi^2_{(1),j}$ متغیرهای تصادفی مستقل χ^2 با یک درجه‌ی آزادی است یعنی:

$$nT_n \xrightarrow{d} \sum_{j=1}^{2m} w_j \chi^2_{(1),j},$$

و همچنین اگر W ماتریسی خودتوان با رتبه‌ی k باشد ($k \leq 2m$) آن‌گاه:

$$nT_n \xrightarrow{d} \chi^2_{(k)}.$$

برهان. با استفاده از قضیه‌ی حد مرکزی چندمتغیره، بردار تصادفی $\sqrt{n}\{\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_0(\mathbf{t}_m)\}$ دارای توزیع مجانبی نرمال با میانگین صفر و واریانس Ω است و در نتیجه قضیه به‌سادگی اثبات می‌شود.

در انجام این آزمون، یافتن مقدار m مناسب و همچنین انتخاب نقاط t_1, \dots, t_m دارای اهمیت ویژه‌ای است. این نقاط باید به‌گونه‌ای انتخاب شوند که بررسی تابع مشخصه‌ی تجربی در نقاط t_1, \dots, t_m بتواند تعیین‌کننده‌ی تابع توزیع مناسب باشد. هر چند نویسندگان بسیاری آزمون نیکویی برازش بر اساس تابع مشخصه‌ی تجربی را مورد بررسی قرار داده‌اند، با این وجود روش کاربردی روشنی برای انتخاب m و نقاط t_i ارائه نشده است و در بیش‌تر مقاله‌ها m را عدد یک یا دو انتخاب کرده‌اند! اوبانک و لارسیا (۱۹۹۲) m را عددی در نظر گرفتند که $(R(m))$ ، حاصل میانگین توان دوم خطای برآوردگر تابع چگالی را مینیمم کند.

$$R(m) = -(n+1) \sum_{j=1}^m \frac{\tilde{a}_{j,n}^2}{n} + \frac{m}{n-1} + \sum_{j=1}^m \frac{\tilde{a}_{2j,n}^2}{n-1},$$

که در آن $(j\pi F_0(x_k))$ $\sqrt{2} \cos(j\pi F_0(x_k))$ و $\tilde{a}_{j,n} = \frac{1}{n} \sum_{k=1}^m \sqrt{2} \cos(j\pi F_0(x_k))$ تابع توزیع تحت فرض صفر می‌باشد. با وجود آن که محاسبات لازم برای یافتن m از طریق پیشنهاد بالا، به‌وسیله‌ی نرم‌افزارهای ریاضی چندان دشوار نیست، ولی m ای که از این راه به‌دست می‌آید، غالباً بسیار بزرگ (تقریباً برابر با اندازه‌ی نمونه) و یا بسیار کوچک (تقریباً برابر با یک) خواهد بود. کوچک بودن m ، دقت آزمون را کاهش می‌دهد و بزرگ بودن m ، محاسبات مربوط به معکوس کردن ماتریس واریانس-کوواریانس Ω را مشکل می‌کند. در بسیاری از موارد به حالت‌هایی برمی‌خوریم که به علت خطای محاسبات و خطای گرد کردن، ماتریس Ω معکوس‌پذیر نیست.

فن (۱۹۹۶) روش یافتن m مناسب را در آزمون نیکویی برازش توزیع نرمال استاندارد بیان کرد اما این روش قابل تعمیم به سایر توزیع‌ها نیست.

وجود این مشکلات، ما را بر آن داشت که روشی برای یافتن m ابداع کنیم که بتواند در آزمون نیکویی برازش هر توزیعی به کار آید. ما این کار را با استفاده از روش مؤلفه‌های اصلی انجام داده‌ایم و آماره‌ی

جدیدی برای آزمون فرض صفر ارائه داده‌ایم که در بخش‌های بعد در مورد آن صحبت خواهیم کرد.

۴ آماره‌ی مؤلفه‌های اصلی والد در آزمون نیکویی برازش

در محاسبه‌ی آماره‌ی والد که در رابطه‌ی (۱) معرفی شد، به معکوس ماتریس $2m \times 2m$ بعدی Ω نیاز داریم. اگر m بزرگ باشد، در بیش‌تر موارد عملی، به دلیل وجود مقادیر بسیار کوچک در بین مقادیر ویژه ماتریس Ω ، این ماتریس معکوس را حتا با نرم‌افزارهای جدید رایانه‌ای نمی‌توان به‌دست آورد. بنا بر این، ما توجه خود را به‌سوی متغیرهایی معطوف می‌کنیم که بیش‌ترین واریانس را دارند، زیرا متغیرهایی با واریانس پایین را می‌توان به‌عنوان متغیرهایی ثابت در نظر گرفت که در ماتریس واریانس-کوواریانس Ω تأثیری ندارند. بدین وسیله می‌توانیم به‌راحتی مسئله‌ی خود را در یک زیرفضا با بعد کم‌تر مورد مطالعه قرار دهیم.

این کار را به‌کمک روش مؤلفه‌های اصلی می‌توان انجام داد. یعنی ترکیب‌های خطی از $\mathbf{Y}_z(\mathbf{t}_m)$ را در نظر می‌گیریم که بیش‌ترین تأثیر را در ماتریس واریانس-کوواریانس Ω داشته باشند. اگر $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{2m}$ مقادیر ویژه ماتریس Ω باشند، واریانس مؤلفه‌های اصلی (ترکیب‌های خطی از $\mathbf{Y}(\mathbf{t}_m)$) برابر با مقادیر λ_i خواهد بود. k را به‌گونه‌ای می‌یابیم که نسبت $\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^{2m} \lambda_i}$ یعنی سهم k مؤلفه‌ی اصلی اول از واریانس کل جامعه، مقدار نسبتاً بالایی (حدود ۹۰ تا ۹۹ درصد) باشد. سپس در مطالعات خود فقط از k مؤلفه‌ی اصلی اول استفاده می‌کنیم. اگر $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ بزرگ‌ترین k مقدار ویژه ماتریس Ω بوده و $\beta_1, \dots, \beta_k, k$ بردار ویژه‌ی مربوط به آن‌ها باشند، آنگاه آماره‌ی «مؤلفه‌های اصلی والد» را به‌صورت زیر تعریف می‌کنیم:

$$(۳) \quad T_n^* = (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_0(\mathbf{t}_m))' B_k \Lambda_k^{-1} B_k' (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_0(\mathbf{t}_m)),$$

که در آن Λ_k یک ماتریس قطری $k \times k$ است که اعضای روی قطر آن مقادیر ویژه $\lambda_1 \geq \dots \geq \lambda_k$ هستند. همچنین B_k ماتریسی $2m \times k$ بعدی است که ستون‌های آن k بردار ویژه‌ی مربوط به $\lambda_1, \dots, \lambda_k$ هستند.

قضیه‌ی زیر توزیع مجانبی nT_n^* را مشخص می‌کند که می‌توان در مسائل عملی از آن استفاده کرد.

قضیه‌ی ۲ اگر بردار تصادفی \mathbf{v}_n برابر با $B_k' (\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_0(\mathbf{t}_m))$ باشد آنگاه آماره‌ی مؤلفه‌های اصلی والد برابر است با:

$$T_n^* = \mathbf{v}_n' \Lambda_k^{-1} \mathbf{v}_n,$$

و تحت فرض صفر توزیع حدی nT_n^* توزیع χ^2 با k درجه‌ی آزادی است یعنی:

$$nT_n^* \xrightarrow{d} \chi_{(k)}^2.$$

برهان. با توجه به خواص مقادیر ویژه و بردارهای ویژه‌ی یک ماتریس، می‌دانیم که برای $i = 1, \dots, k$ ، $\beta_i' \Omega_0 \beta_i = \lambda_i$ و برای $i, j = 1, \dots, k, i \neq j$ ، $\beta_i' \Omega_0 \beta_j = 0$ با استفاده از قضیه‌ی حد مرکزی چندمتغیره، تحت فرض صفر، توزیع مجانبی $\sqrt{n} \mathbf{v}_n$ نرمال با میانگین صفر و ماتریس واریانس Λ_k می‌باشد یعنی:

$$\sqrt{n} \mathbf{v}_n \xrightarrow{d} N(\mathbf{0}, \Lambda_k),$$

یا

$$\sqrt{n} \Lambda_k^{-\frac{1}{2}} \mathbf{V}_n \xrightarrow{d} N(\mathbf{0}, I),$$

و در نتیجه:

$$nT_n^* = n \mathbf{v}_n' \Lambda_k^{-1} \mathbf{v}_n \xrightarrow{d} \sum_{j=1}^k \chi_{(1),j}^2 = \chi_{(k)}^2.$$

۵ شبیه‌سازی و مقایسه‌ی آزمون جدید با آزمون‌های پیشین

به‌طور خلاصه می‌توان الگوریتم به‌دست آوردن آماره‌ی بالا و انجام آزمون نیکویی برازش توزیع را با در دست داشتن n مشاهده از توزیع جامعه، به‌شکل زیر نوشت:

(۱) با در نظر گرفتن $(F_0(\cdot))$ به‌عنوان توزیع جامعه تحت فرض صفر مقادیر t_1, \dots, t_m را

به‌وسیله‌ی رابطه‌های $t_i = F_0^{-1}\left(\frac{i}{m+1}\right)$ ، $i = 1, \dots, m$ به‌دست می‌آوریم. m را مساوی n در نظر می‌گیریم.

(۲) بردارهای $\mathbf{Z}_n(t_m)$ و $\mathbf{Z}_0(t_m)$ را تشکیل می‌دهیم.

(۳) بردار $2m$ بعدی $\mathbf{Z}_n(t_m) - \mathbf{Z}_0(t_m)$ را مشخص می‌کنیم.

(۴) ماتریس Ω_0 را می‌سازیم و مقادیر ویژه‌ی مربوط به آن‌ها را به‌دست می‌آوریم.

(۵) مقدارهای ویژه را به‌صورت نزولی مرتب می‌کنیم و مقدار k مناسب را به‌گونه‌ای می‌یابیم که

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^{2m} \lambda_i} \geq 0.99.$$

(البته انتخاب کردن ۰/۹۹ اختیاری است و جهت افزایش دقت می توان از کران های بالاتر نیز استفاده کرد.)

(۶) ماتریس های B_k و Λ_k را تشکیل داده و با محاسبه ی بردار v_n مقدار آماره ی T_n^* را به دست می آوریم.

(۷) p -مقدار آزمون را با استفاده از فرمول $p = P(\chi_{(k)}^2 > nT_n^*)$ مقدار برای رد یا قبول فرض صفر محاسبه می کنیم.

اگرچه به دست آوردن p -مقدار و مقایسه ی آن با میزان با معنایی α راهی برای پذیرش یا رد فرض صفر می باشد، اما برای بررسی برتری آزمون ارائه شده نسبت به آزمون های ناپارامتری دیگر نیاز به محاسبه ی توان آزمون داریم. توان آزمون با مشخص شدن فرض مقابل، قابل دسترسی است. در نتیجه باید فرض صفر $H_0: F(x) = F_0(x)$ را در مقابل فرض مشخص $H_1: F(x) = F_1(x)$ آزمون کنیم که در آن $F_0(\cdot)$ و $F_1(\cdot)$ توابع توزیع کاملاً مشخصی هستند.

در جدول ۱ نتایج مقادیر توان آزمون جدید در آزمودن فرض صفر نرمال در مقابل توزیع های متفاوتی، بر اساس ۱۰۰۰ نمونه ی ۵۰ تایی آورده شده است. مقادیر توان برای آزمون های ناپارامتری کولموگروف-اسمیرنوف (KS)، کرامر-ون میسز (CM)، اندرسن-دارلینگ (AD) و واتسون (WA) محاسبه شده است.

همه ی آماره ی آزمون های بالا بر اساس تابع توزیع تجربی بوده و به صورت زیر تعریف شده اند. آزمون کولموگروف-اسمیرنوف:

$$KS = \max \{ D^+, D^- \},$$

$$D^+ = \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - \hat{F}(X_{(j)}) \right\},$$

$$D^- = \max_{1 \leq j \leq n} \left\{ \hat{F}(X_{(j)}) - \frac{j-1}{n} \right\}.$$

آزمون کرامر-ون میسز:

$$CM = \frac{1}{\sqrt{2n}} + \sum_{j=1}^n \left(\hat{F}(X_{(j)}) - \frac{j-1}{2n} \right)^2.$$

آزمون اندرسن-دارلینگ:

$$AD = -n - \frac{1}{n} \sum_{j=1}^n \left[(2j-1) \log \hat{F}(X_{(j)}) + \{2(n-j)+1\} \log \{1 - \hat{F}(X_{(j)})\} \right].$$

آزمون واتسون:

$$WA = \frac{1}{12n} + \sum_{j=1}^n \left\{ \hat{F}(X_{(j)}) - \frac{j-1}{2n} \right\}^2 - n \left\{ \sum_{j=1}^n \frac{\hat{F}(X_{(j)})}{n} - \frac{1}{2} \right\}^2,$$

که در تمام آن‌ها $\hat{F}(\cdot)$ تابع توزیع X و $X_{(j)}$ ، j امین آماره‌ی مرتب‌شده است. شبیه‌سازی با روش مونت کارلو در سطح ۵٪ و با استفاده از نرم‌افزار S-plus و Maple(Waterloo) انجام گرفته است.

جدول ۱. درصد توان ۵ آزمون ناپارامتری در سطح ۵٪ برای فرض نرمال در برابر فرض‌های مختلف آماری

	KS	CM	WA	AD	NEW
Alternative					
Cauchy	۵۶	۵۹	۷۹	۴۶	۹۹
Exp(۱)	*	*	*	*	*
t(۳)	۱۱	۱۲	۱۷	۳۴	۷۳
Logistic	۵۳	۵۷	۹۰	۷۳	۹۹
U(-۲, ۲)	۲۳	۲۸	۶۸	۴۴	۷۶
Lognormal	*	*	*	*	*
Laplace	۶	۶	۱۰	۳۶	۷۶
Gumbel	۶۶	۷۸	۲۲	۹۲	۹۰
Triangular	*	*	*	*	*
GH(۰٫۲, ۰٫۵)	۱۴	۱۴	۱۶	۵۲	۶۷

همچنین برای فرض مقابل به ترتیب توزیع‌های زیر در نظر گرفته شده است:

- توزیع کوشی استاندارد؛
- توزیع نمایی با پارامتر ۱؛
- توزیع t با سه درجه‌ی آزادی؛
- توزیع لوژستیک؛
- توزیع لگ نرمال؛
- توزیع نمایی دوگانه (لاپلاس) استاندارد؛
- توزیع گامبل (مقدار کرانگین)؛

- توزیع یکنواخت در بازه $(-2, 2)$ ؛
- توزیع مثلثی در بازه $(-1, 1)$ ؛
- توزیع توکی $g-h$ با نماد $GH(g,h)$ و به صورت $X = \exp\left\{\frac{hZ^r}{r}\right\} \frac{\exp\{gZ\}-1}{g}$ که در آن $Z \sim N(0, 1)$.

در جدول ۱، NEW نشان‌دهنده‌ی آزمون جدید و نماد * نشان‌دهنده‌ی توان 10% است. نتایج اصلی برگرفته از جدول ۱ به شرح زیر است:

(۱) در بین آزمون‌هایی که بر پایه‌ی توزیع تجربی بنا شده‌اند، آزمون واتسون و اندرسن-دارلینگ بتوان تر از کولموگروف-اسمیرنوف هستند. این موضوع به خصوص در فرض‌های مقابل با دم‌های سنگین همچون کوشی، لوژستیک و گامبل مشهودتر است.

(۲) در فرض مقابل گامبل آزمون جدید ارائه شده به خوبی آزمون اندرسن-دارلینگ عمل می‌کند و در فرض‌های مقابل دیگر، آزمون جدید بتوان تر از آزمون‌های دیگری است که بر پایه‌ی تابع توزیع تجربی استوارند. این موضوع ادعای ما را مبنی بر بهینه بودن روش جدید نشان می‌دهد.

مرجع‌ها

- Baringhaus, L.; Henze, N. (1988). A consistent test for multivariate normality based on the empirical characteristic function, *Metrika* **35**, 339-348.
- Eubank, R.L.; LaRiccia, V.N. (1992). Asymptotic comparison of Cramer-Von Mises and nonparametric function estimation techniques for testing goodness-of-fit, *Ann. Statist.* **20**, 2071-2086.
- Fan, J. (1996). Test of significance based on wavelet thresholding and Neyman's truncation, *J. Amer. Statist. Assoc.* **91**, 674-688.
- Fan, Y. (1997). Goodness-of-fit tests for a multivariate distribution by the empirical characteristic function, *J. Multivariate Statist.* **62**, 36-63.
- Feigin, P.D.; Heathcote, C.R. (1976). The empirical characteristic function and the Cramer-Von Mises statistics, *Sankhyā* **38**, 309-325.
- Feuerverger, A.; Mureika, R.A. (1977). The empirical characteristic function and its applications, *Ann. Statist.* **5**, 88-97.
- Hall, P.; Welsh, A.H. (1983). A test for normality based on the empirical characteristic function, *Biometrika* **70**, 723-726.

Lukacs, E. (1970). Characteristic function, 2nd ed. Charies Griffin, London.

دریافت: ۲۱ تیر ۱۳۸۴
آخرین اصلاح: ۹ اردیبهشت ۱۳۸۵

مهدی سلمان پور	مینا توحیدی
گروه آمار، دانشکده‌ی علوم،	گروه آمار، دانشکده‌ی علوم،
دانشگاه شیراز،	دانشگاه شیراز،
چهارراه ادبیات،	چهارراه ادبیات،
شیراز، ایران.	شیراز، ایران.
پيام‌نگار: <i>mhi-salman@hotmail.com</i>	پيام‌نگار: <i>mtowhidi@susc.ac.ir</i>