

یک آزمون جدید نیکویی برآزش 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 منحصر به فردی تابع احتمال یک متغیر تصادفی به وسیله‌ی 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 $c_n(t) = \frac{1}{n} \sum_{j=1}^n \exp\{itX_j\}$, $\forall t \in \mathbb{R}$ تبدیل کرد و از تابع مشخصه‌ی تجربی، $c(t) = c_0(t)$. در آزمون نیکویی برآزش توزیع استفاده نمود. به همین دلیل، بین سال‌های ۱۹۹۳ تا ۱۹۷۲ بسیاری از پژوهشگران از تابع مشخصه‌ی تجربی برای آزمون فرض‌های مختلف آماری استفاده کردند. در بیش‌تر آزمون‌های معرفی شده، مقایسه‌هایی بین (t) و $c_n(t)$ برای تعداد کمی از مقادیر t ، صورت گرفته است و این امر موجب سازگار نبودن آزمون‌ها شده است. ما در این مقاله سعی کردیم 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 کردیم که بسیار برتوان تراز آزمون‌های ناپaramتری قبلی بوده است.

وازگان کلیدی. آزمون نیکویی برآزش؛ 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 $(x) F$ باشد و بخواهیم فرض X دارای توزیع شناخته شده‌ی $(x) F$ را آزمون کنیم. توزیع $(x) F$ را توزیع برازنده بر X (یا توزیع جور با یافته‌های X) می‌گویند. بر اساس n مشاهده‌ی مستقل و هم توزیع ناشناخته‌ی $(x) F$, آزمون‌هایی راکه برای آزمودن فرض $H_1 : F(x) \neq F_{\circ}(x)$, $\exists x \in \mathbb{R} : F(x) = F_{\circ}(x)$, $\forall x \in \mathbb{R} H_0 : F(x) = F_{\circ}(x)$ بهکار بردۀ می‌شوند، آزمون‌های نیکوبی برازش توزیع می‌نامند.

تاکنون آزمون‌های ناپارامتری بسیاری برای آزمودن فرض‌های بالا معرفی شده است که از جمله می‌توان به آزمون‌های کولموگروف-اسمیرنوف، کرامر-ون میسز، 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 $(x) F$ 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 بهوسیله‌ی آن می‌توان تابع احتمال یک متغیر تصادفی را تعیین کرد. بهمین دلیل، بین سال‌های ۱۹۷۲ تا ۱۹۹۳ پژوهشگران در مقاله‌های گوناگون، با استفاده از تابع مشخصه‌ی تجربی، آزمون‌های نیکوبی برازش را با در نظر گرفتن توزیع‌های مشخص $(x) F$ ارائه کرده‌اند.

هدف ما در این مقاله، طراحی یک آزمون نیکوبی برازش کلی (با در نظر گرفتن هر توزیع مشخص $(x) F$) بر اساس تابع مشخصه‌ی تجربی است که پرتوان‌تر از آزمون‌های ارائه شده قبلي 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 $(x) F$ به صورت زیر تعریف می‌شود:

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

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

(آ) $E\{\cos(tX)\} + iE\{\sin(tX)\}$ را قسمت حقیقی

(Re $c(t)$) و $E\{\sin(tX)\}$ را قسمت موهومی (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 برآورده سازگار برای $c(t)$, با در دست داشتن یک نمونه 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:

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

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

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

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

برهان. با استفاده از قانون قوی اعداد بزرگ برای متغیرهای تصادفی $\{X_1, \dots, X_n\}$, 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 این‌که فرض $\forall x : F(x) = F_{\circ}(x)$, $\forall t \in \mathbb{R} : c(t) = c_{\circ}(t)$, نوشت، پس می‌توان آماره‌ی آزمون نیکویی برآش توزیع را براساس تابع مشخصه‌ی تجربی $c_n(t)$, طرح‌ریزی 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، قوان آزمون به‌سمت یک میل کند. به عبارت دیگر قوان مجذبی 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 $\text{Im } c_{n+1}(t)$ (یا $\text{Re } c_{n+1}(t)$) برای یک مقدار ثابت t انجام دادند. آن‌ها با انتخاب یک t مناسب، به آزمونی پرتوان‌تر از آزمون کرامر-ون می‌سیزدست یافته‌ند. فیورورگر و موریکا (۱۹۷۷) یک آزمون را برای فرض مقارن ارائه کردند. آن‌ها آماره‌ی خود را بر اساس انتگرال قوان دوم اختلاف بین $\text{Im } c_n(t)$ و صفر بیان کردند (تابع مشخصه، یک تابع حقیقی است اگر و فقط اگر توزیع مربوط مقارن باشد) سپس هال و ولش (۱۹۸۳) آزمونی را برای آزمودن فرض نرمال بودن توزیع ارائه دادند. در بیش‌تر آزمون‌های معرفی شده در مقاله‌های پژوهشگران، مقایسه‌هایی بین $c_n(t)$ و $c_{n+1}(t)$ برای برخی مقادیر t انجام شده است، در نتیجه آزمون‌های ارائه شده سازگار نیستند. در سال ۱۹۸۸ بارینس و هنر آماره‌ای را بر اساس انتگرال قوان دوم اختلاف بین $c_n(t)$ و $c_{n+1}(t)$ 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 قوان دوم اختلاف بین $c_n(t)$ و $c_{n+1}(t)$ در نقاط t_1, \dots, t_m معرفی کرد و برای اثبات سازگاری آزمون، فرض کرد که $m \equiv m_n \rightarrow +\infty$ 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 $\mathbf{Z}_n(\mathbf{t}_m)$ برآورده‌ی ناریب و سازگار برای $\mathbf{Z}(\mathbf{t}_m)$ است، اگر تابع مشخصه‌ی توزیع جامعه $c(t)$ باشد. بنا بر این، به‌منظور 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 $\mathbf{0}$ بعدی است. اگر ماتریس واریانس-کوواریانس $(\mathbf{Y}_j(\mathbf{t}_m))'$ را با Ω نشان دهیم، این ماتریس را می‌توان به صورت زیر افزایش کرد:

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

که درایه‌های این زیر ماتریس‌ها برابر است با:

$$\begin{aligned} (\Omega_{*11})_{ij} &= \text{cov}\{\cos(t_i X_1), \cos(t_j X_1)\} \\ &= E\{\cos(t_i X_1) \cos(t_j X_1)\} - E\{\cos(t_i X_1)\} E\{\cos(t_j X_1)\} \\ &= \frac{1}{4} [\text{Re } c_*(t_i + t_j) + \text{Re } c_*(t_i - t_j) \\ &\quad - 2\text{Re } c_*(t_i)\text{Re } c_*(t_j)], \quad i, j = 1, \dots, m, \quad i \neq j \end{aligned}$$

و به صورت مشابه:

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

اگر $\mathbf{Z}_{\circ}(\mathbf{t}_m)$ مقدار $\mathbf{Z}(\mathbf{t}_m)$ تحت فرض صفر باشد، آنگاه آماره‌ی آزمون که در حقیقت همان آماره‌ی $\Omega_{\circ 21}$ برای بردار تصادفی $\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_{\circ}(\mathbf{t}_m)$ می‌باشد، به صورت زیر معرفی می‌شود:

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

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 $\theta = \theta_n \rightarrow m$ یا $\theta = \theta_n \rightarrow +\infty$ و وقتی که $\theta = \theta_n \rightarrow 0$ بوده و آماره‌ی جدید T_n را معرفی کرد:

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

و در قضیه‌ای 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 مشیت 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 $\chi_{(1),j}^{\star}$ متغیرهای تصادفی مستقل χ^{\star} با یک درجه‌ی آزادی است یعنی:

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

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

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

برهان. با استفاده از قضیه‌ی حد مرکزی چندمتغیره، بردار تصادفی $\{\mathbf{Z}_n(\mathbf{t}_m) - \mathbf{Z}_\circ(\mathbf{t}_m)\}$ 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 نویسنده‌گان بسیاری آزمون نیکویی برآش براساس تابع مشخصه‌ی تجربی را مورد بررسی قرار داده‌اند، با این وجود روش کاربردی روشی 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 گرفته‌اند که $(R(m))$ حاصل میانگین توان دوم خطای برآورده‌گر تابع چگالی را مینیمم کند.

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

که در آن $(x_k) = \frac{1}{n} \sum_{k=1}^m \sqrt{2} \cos(j\pi F_\circ(x_k))$ و F_\circ 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 بدست می‌آید، غالباً بسیار بزرگ (قریباً برابر با اندازه‌ی نمونه) و یا بسیار کوچک (قریباً 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 برآش 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 برآش 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 داده ایم که در بخش های 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 حتی با نرم افزارهای جدید رایانه‌ای نمی‌توان بدست آورد. بنا بر این، ما توجه خود را به سوی متغیرهایی معطوف می‌کنیم که بیشترین واریانس را دارند، زیرا متغیرهایی با واریانس پایین را می‌توان به عنوان متغیرهایی ثابت در نظر گرفت که در ماتریس واریانس-کوواریانس Ω تأثیری ندارند. بدین وسیله می‌توانیم به راحتی مسئله‌ی خود را در یک زیرفضا با بعد کمتر مورد مطالعه قرار دهیم.

این کار را به کمک روش مؤلفه‌های اصلی می‌توان انجام داد. یعنی ترکیب‌های خطی از $\mathbf{Y}_j(t_m)$ 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 $\mathbf{Y}(t_m)$) برابر با مقادیر λ خواهد بود. k را به گونه‌ای می‌یابیم که نسبت $\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^{2m} \lambda_i}$ یعنی سهم مؤلفه‌ی اصلی اول از واریانس کل جامعه، مقدار نسبتاً بالای (حدود ۹۰ تا ۹۹ درصد) باشد. سپس در مطالعات خود فقط از k مؤلفه‌ی اصلی اول استفاده می‌کنیم. اگر $k \geq \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ باشد، بزرگترین k مقدار ویژه ماتریس Ω بوده و β_1, \dots, β_k بردار ویژه مربوط به آنها باشند، آنگاه آماره‌ی «مؤلفه‌های اصلی والد» را به صورت زیر تعریف می‌کنیم:

$$(۳) \quad T_n^* = (\mathbf{Z}_n(t_m) - \mathbf{Z}_*(t_m))' B_k \Lambda_k^{-1} B_k' (\mathbf{Z}_n(t_m) - \mathbf{Z}_*(t_m)),$$

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 $k \times 2m$ بعدی است که ستون‌های آن k بردار ویژه 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 v_n برابر با $(\mathbf{Z}_n(t_m) - \mathbf{Z}_*(t_m))' B_k'$ باشد آنگاه آماره‌ی مؤلفه‌های اصلی والد برابر است با:

$$T_n^* = v_n' \Lambda_k^{-1} v_n,$$

و تحت فرض صفر توزیع حدی nT_n^* ، توزیع $\chi_{(k)}^2$ 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 $i = 1, \dots, k$ $\beta_i' \Omega \beta_i = \lambda_i$ و برای $j = 1, \dots, k, i \neq j$ $\beta_i' \Omega \beta_j = 0$. 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

به طور خلاصه می‌توان الگوریتم به دست آوردن آماره‌ی بالا و انجام آزمون نیکویی برآش توزیع را با در دست داشتن n مشاهده از توزیع جامعه، به شکل زیر نوشت:

- (۱) با در نظر گرفتن (F_0) به عنوان توزیع جامعه تحت فرض صفر مقادیر t_1, t_2, \dots, t_m را به وسیله‌ی رابطه‌های $t_i = F_0^{-1}\left(\frac{i}{m+1}\right), i = 1, \dots, m$ به دست می‌آوریم. m را مساوی n در نظر می‌گیریم.
- (۲) بردارهای (t_m) و $Z_n(t_m)$ را تشکیل می‌دهیم.
- (۳) بردار $Z_n(t_m) - Z_n(t_m)$ را مشخص می‌کنیم.
- (۴) ماتریس Ω را می‌سازیم و مقادیر ویژه‌ی مربوط به آنها را به دست می‌آوریم.
- (۵) مقادارهای ویژه را به صورت نزولی مرتب می‌کنیم و مقدار k مناسب را به‌گونه‌ای می‌یابیم که

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

(البته انتخاب کردن $\alpha = 0.05$ اختیاری است و جهت افزایش دقت می‌توان از کران‌های بالاتر نیز استفاده کرد).

(۶) ماتریس‌های B_k و Λ_k را تشکیل داده و با محاسبه‌ی بردار v_n مقدار آماره‌ی T_n^* را به دست می‌آوریم.

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

اگرچه به دست آوردن p -مقدار و مقایسه‌ی آن با میزان با معنایی α راهی پذیرش یا رد فرض صفر می‌باشد، اما برای بررسی برتری آزمون ارائه شده نسبت به آزمون‌های ناپارامتری دیگر نیاز به محاسبه‌ی توان آزمون داریم. توان آزمون با مشخص شدن فرض مقابل، قابل دسترسی است. در نتیجه باید فرض صفر آزمون کنیم که در آن $F(x) = F_1(x)$: $F(x) = F_1(x)$ آزمون کنیم که در آن $F(x) = F_1(x)$ و $F_1(x) = F(x)$ توابع توزیع کاملاً مشخصی هستند.

در جدول ۱ نتایج مقادیر توان آزمون جدید در آزمودن فرض صفر نیمال در مقابل توزیع‌های متفاوتی، بر اساس 1000 نمونه 50 تایی آورده شده است. مقادیر توان برای آزمون‌های ناپارامتری کولموگروف-اسمیرنوف (KS)، کرامر-ون میسز (CM)، اندرسن-دارلینگ (AD) و واتسون (WA) محاسبه شده است.

همه‌ی آماره‌ی آزمون‌های بالا بر اساس تابع توزیع تجربی بوده و به صورت زیر تعریف شده‌اند. آزمون کولموگروف-اسمیرنوف:

$$\begin{aligned} KS &= \max \{D^+, D^-\}, \\ D^+ &= \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - \hat{F}(X_{(j)}) \right\}, \\ D^- &= \max_{1 \leq j \leq n} \left\{ \hat{F}(X_{(j)}) - \frac{j-1}{n} \right\}. \end{aligned}$$

آزمون کرامر-ون میسز:

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

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 $\hat{F}(.)$ تابع توزیع $X_{(j)}$ ، زامین آماره‌ی ترتیبی 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

Alternative	KS	CM	WA	AD	NEW
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 مقابله‌ی ترتیب توزیع‌های زیر در نظر گرفته شده است:

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

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

در جدول ۱، NEW نشان‌دهنده‌ی آزمون جدید و نماد * نشان‌دهنده‌ی توان 100% 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 پرتوان‌تر از کولموگروف-اسمیرنوف هستند. این موضوع به خصوص در فرض‌های مقابله با دم‌های سنگین همچون کوشی، لوژستیک و گامبل مشهودتر است.
- (۲) در فرض مقابله گامبل آزمون جدید ارائه شده به خوبی آزمون اندرسن-دارلینگ عمل می‌کند و در فرض‌های مقابله دیگر، آزمون جدید پرتوان‌تر 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. Charles Griffin, London.

دریافت:	۱۳۸۴ تیر ۲۱
آخرین اصلاح:	۱۳۸۵ اردیبهشت ۹

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