

آزمونی برای استقلال در داده‌های نرمال با بُعد بالا

داریوش نجارزاده

گروه آمار، دانشکده علوم ریاضی، دانشگاه تبریز

تاریخ دریافت: ۱۳۹۷/۱۰/۲۱ تاریخ پذیرش: ۱۳۹۸/۱۰/۲۰

چکیده: فرضیه استقلال کامل داده‌ها پیش‌نیاز بسیاری از استنباط‌های آماری است. روش‌های آزمون کلاسیک پاسخ‌گوی بررسی چنین فرضی در داده‌های با ابعاد بالا نیست. در این مقاله آماره آزمونی ساده برای وجود استقلال کامل در داده‌های نرمال با بُعد بالا معرفی و با استفاده از نظریه مارتینگل‌ها مجاناً نرمال بودن توزیع این آماره ثابت شده است. به منظور ارزیابی عملکرد آزمون پیشنهادی و مقایسه آن با روش‌های موجود، مطالعه‌ای شبیه‌سازی انجام شده است. نتایج شبیه‌سازی نشان می‌دهد که آزمون پیشنهادی دارای خطای نوع اول تجربی با میانگین خطای نسبی کوچکتر نسبت به آزمون‌های موجود است. کاربردی از روش پیشنهادی بر مجموعه داده سطح بیان ژن تومور پروستات معرفی شده است. **واژه‌های کلیدی:** آزمون استقلال کامل، توزیع نرمال چندمتغیره، داده‌های با بُعد بالا، نظریه مارتینگل.

۱ مقدمه

استنباط‌های آماری چند متغیره کلاسیک برای داده‌های با ابعاد بالا قابل استفاده نیستند (ساراناداسا و بای، ۱۹۹۶؛ دمپستر، ۱۹۵۸، ۱۹۶۰). چرا که بیشتر این آزمون‌ها بر این اساس طراحی شده‌اند که بُعد p ثابت و در مقایسه با اندازه نمونه n کوچک است. برای داده‌های با بُعد بالا این شرط دیگر برقرار نیست، چرا که در چنین داده‌هایی بُعد در مقایسه با اندازه نمونه به مراتب عدد بزرگتری است. نمونه‌هایی از چنین داده‌ها را می‌توان در داده‌های ریزآرایه^۱، زیستی، مالی، و چندرسانه‌ای یافت (ژیراد، ۲۰۱۴). در

آدرس الکترونیکی نویسنده مسئول مقاله: داریوش نجارزاده، d_najarzadeh@tabrizu.ac.ir

کد موضوع‌بندی ریاضی (۲۰۱۰): 60F05, 62H10, 62H15

¹Microarray

ادبیات چندمتغیره با بُعد بالا، همواره آزمون استقلال بردارهای تصادفی دارای توزیع نرمال چندمتغیره از اهمیت خاصی برخوردار بوده است. چن و مودهولکار (۱۹۹۰) توزیع مجموع توان‌های دوم تبدیل‌های Z فیشر برای آزمون فرضیه استقلال کامل را بررسی کردند. سیرواستاوا (۲۰۰۵، ۲۰۰۶) آزمون‌هایی در مورد ماتریس کوواریانس و از جمله آزمون استقلال در داده‌های نرمال با بُعد بالا را معرفی کرد. رید و سیرواستاوا (۲۰۱۲) آزمون‌هایی برای استقلال زیربردارهای یک بردار نرمال با بُعد بالا را مورد مطالعه قرار دادند. آزمون همزمان استقلال برای زیربردارهای چند بردار نرمال با بُعد نسبتاً بالا توسط نجارزاده (۱۳۹۸) معرفی شد. آزمون‌های نسبت‌درست‌نمایی در رابطه با ماتریس کوواریانس داده‌های نرمال با بُعد بالا توسط جیانگ و همکاران (۲۰۱۲) و جیانگ و کی (۲۰۱۵) معرفی شد. جیانگ و یانگ (۲۰۱۳) قضیه‌های حدی مرکزی برای آزمون‌های نسبت‌درست‌نمایی کلاسیک، از جمله آزمون استقلال، در جوامع نرمال با بُعد بالا را معرفی کردند.

این آزمون به صورت زیر معرفی می‌شود. فرض کنید X_1, \dots, X_n نمونه‌ای تصادفی از توزیع نرمال p -متغیره با ماتریس کوواریانس Σ است. در این حالت، فرضیه استقلال کامل به صورت

$$H_0: \rho_{ij} = 0, \quad 1 \leq i < j \leq p, \quad (1)$$

تعریف می‌شود، که در آن ρ_{ij} درایه (i, j) ام ماتریس همبستگی

$$\rho = (\text{diag} \Sigma^{-\frac{1}{2}}) \Sigma (\text{diag} \Sigma^{-\frac{1}{2}}) = [\rho_{ij}]_{p \times p}$$

است، که در آن $\text{diag}(\Sigma)$ ماتریس قطری با درایه‌های روی قطر اصلی ماتریس Σ است. فرض کنید S و $R = (\text{diag}(S)^{-\frac{1}{2}}) S (\text{diag}(S)^{-\frac{1}{2}}) = [r_{ij}]_{p \times p}$ به ترتیب ماتریس‌های کوواریانس و همبستگی نمونه‌ای حاصل از X_1, \dots, X_n باشند. در این صورت، بر اساس آزمون نسبت درست‌نمایی فرضیه H_0 رد می‌شود اگر مقدار آماره

$$-\left(n - \frac{2p + 1}{\epsilon}\right) \log(\det(R))$$

از صدک $(1 - \alpha)100$ درصد توزیع $\chi_{p(p-1)/2}^2$ بزرگ‌تر باشد (مؤرهد، ۲۰۰۵، قضیه ۵.۲۰.۱۱) که در آن $\det(R)$ دترمینان ماتریس R است. این آزمون برای داده‌های با ابعاد بالا که در آنها $p > n$ ، غیر قابل تعریف است، زیرا ماتریس همبستگی نمونه‌ای R به ازای $p > n$ ، ماتریسی تکین با $\det(R) = 0$ است. به منظور ارائه آزمون مناسب در داده‌های با ابعاد بالا برای فرضیه (۱) پژوهش‌های انجام شده

است. در این راستا، سیرواستاوا (۲۰۰۵) ثابت کرد که تحت فرضیه (۱) وقتی p و n هر دو به بینهایت میل می‌کنند، آماره آزمون به صورت

$$T^* = \frac{(n-1)t_{np} - \frac{p(p-1)}{2}}{\sqrt{p(p-1)}} \xrightarrow{d} N(0, 1),$$

است، که در آن $t_{np} = \sum_{i=2}^p \sum_{j=1}^{i-1} r_{ij}^2$ و \xrightarrow{d} نماد همگرایی در توزیع است. سیرواستاوا (۲۰۰۵) نشان داد همگرایی توزیع T^* کند است، برای غلبه بر این مشکل، آماره آزمون

$$T = \frac{(n-1)(\hat{\gamma}_3 - 1)}{2\sqrt{1 - \frac{\hat{a}_{30}}{p\hat{a}_{10}}}}$$

را پیشنهاد داد که در آن

$$\hat{\gamma}_3 = \frac{(n-1)[tr(S^3) - tr(S)^3/(n-1)]}{(n-2)tr(\text{diag}(S)^2)}, \quad \hat{a}_{30} = \frac{n-1}{p(n+1)}tr(\text{diag}(S)^3),$$

و $\hat{a}_{10} = \frac{1}{p}tr(\text{diag}(S)^2)$. با فرض اینکه به ازای هر $i = 1, \dots, 8$ حد $\frac{tr(\Sigma^i)}{p}$ وقتی $p \rightarrow \infty$ موجود و مثبت است و همچنین $n = O(p^\delta)$ برای $0 < \delta \leq 1$ ، سیرواستاوا (۲۰۰۵) نشان داد که $T \xrightarrow{d} N(0, 1)$ وقتی $p \rightarrow \infty$ و $n \rightarrow \infty$. از آنجا که T بر اساس ماتریس کوواریانس نمونه‌ای S ساخته می‌شود، تحت تبدیلات مقیاس از متغیرها ناوردا^۲ نخواهد بود. اسکات (۲۰۰۵) با در نظر گرفتن آماره t_{np} نشان داد که تحت فرضیه H_0 و این فرض که $\gamma \in (0, \infty)$ $\lim_{n,p \rightarrow \infty} (\frac{p}{n}) \rightarrow \gamma$ همواره

$$t_{np}^* = \frac{t_{np} - \frac{p(p-1)}{2(n-1)}}{\sigma_{np}^2} \xrightarrow{d} N(0, 1),$$

²Invariance

که در آن $\sigma_{np}^2 = \frac{p(p-1)(n-2)}{(n-1)^2(n+1)}$. با الهام از آزمون اسکات (۲۰۰۵)، مائو (۲۰۱۴) آماره

$$T_{np} = \sum_{i=2}^p \sum_{j=1}^{i-1} \frac{r_{ij}^2}{1 - r_{ij}^2}$$

را برای آزمون فرضیه (۱) معرفی کرده و نشان داد وقتی $\lim_{n,p \rightarrow \infty} (\frac{p}{n}) \rightarrow \gamma \in (0, \infty)$ ، آنگاه

$$T_{np}^* = \frac{T_{np} - \frac{p(p-1)}{2(n-2)}}{\delta_{np}^2} \xrightarrow{d} N(0, 1),$$

که در آن $\delta_{np}^2 = \frac{p(p-1)(n-3)}{(n-2)^2(n-6)}$. اخیراً، چانگ و کویی (۲۰۱۸) فرض $\lim_{n,p \rightarrow \infty} (\frac{p}{n}) \rightarrow \gamma \in (0, \infty)$ بیان شده در آزمون‌های اسکات (۲۰۰۵) و مائو (۲۰۱۴) را حذف و ثابت کردند بدون در نظر گرفتن این قید و تنها با وجود فرض‌های $p \rightarrow \infty$ و $n \rightarrow \infty$ نیز توزیع مجانبی T_{np}^* و t_{np}^* تحت فرضیه صفر همچنان نرمال خواهد بود. بعلاوه نشان دادند که تحت فرضیه (۱)، توزیع مجانبی آماره‌های

$$t_{np}^c = \sqrt{p(p-1)} t_{np}^* + \frac{p(p-1)}{2},$$

$$T_{np}^c = \sqrt{p(p-1)} T_{np}^* + \frac{p(p-1)}{2},$$

کای دو با $p(p-1)/2$ درجه آزادی خواهد بود.

در بخش ۲، آماره‌ای جدید برای آزمون فرضیه (۱) معرفی و ثابت می‌شود تحت فرضیه صفر وقتی $p = p_n \rightarrow \infty$ و $n \rightarrow \infty$ ، دارای توزیع نرمال است. این آماره که بر اساس ماتریس همبستگی طراحی شده ناوردای مکان-مقیاس^۳ است. در بخش ۳، مطالعه شبیه‌سازی به منظور بررسی و مقایسه خطای نوع اول و همچنین توان آزمون پیشنهادی با آزمون‌های بر مبنای آماره‌های T_{np}^* ، t_{np}^* و T_{np}^c و t_{np}^c ارائه می‌شود. کاربردی از آزمون پیشنهادی بر مجموعه داده‌های حاصل از یک مطالعه واقعی در بخش ۴ معرفی می‌شود. در بخش ۵ بحث و نتیجه‌گیری ارائه می‌شود.

³Location-scale invariance

۲ آماره آزمون و توزیع مجانبی آن

اسکات (۲۰۰۵) آماره $t_{np} = \sum_{i=2}^p \sum_{j=1}^{i-1} r_{ij}^{\frac{1}{\nu}}$ و مائو (۲۰۱۴) آماره $T_{np} = \sum_{i=2}^p \sum_{j=1}^{i-1} \frac{r_{ij}^{\frac{1}{\nu}}}{1-r_{ij}^{\frac{1}{\nu}}}$ را برای آزمون فرضیه (۱) پیشنهاد دادند. شهود به کار رفته در آماره‌های فوق و این نکته که $\sum_{i=2}^p \sum_{j=1}^{i-1} \sqrt{|\rho_{ij}|} = 0$ اگر و تنها اگر $\rho_{ij} = 0$ ، ما را به آماره‌ای دیگر برای آزمون فرضیه H_0 به صورت

$$D_{np} = \sum_{i=2}^p \sum_{j=1}^{i-1} \sqrt{|r_{ij}|}$$

رهنمون ساخت. توجه داریم که تحت فرضیه $\rho_{ij} = 0$ همواره

$$E|r_{ij}|^k = E\left[\left(r_{ij}^{\frac{1}{\nu}}\right)^{\frac{k}{\nu}}\right] = \frac{\Gamma\left(\frac{n-1}{\nu}\right)\Gamma\left(\frac{k+1}{\nu}\right)}{\sqrt{\pi}\Gamma\left(\frac{k+n-1}{\nu}\right)} = \Gamma_n(k), \quad (2)$$

(مؤرهد، ۲۰۰۵). بنابراین تحت فرضیه H_0

$$ED_{np} = \sum_{i=2}^p \sum_{j=1}^{i-1} E\sqrt{|r_{ij}|} = \frac{p(p-1)}{\nu} \Gamma_n\left(\frac{1}{\nu}\right) = \mu_{np}.$$

و می‌توان نوشت $r_{ij} = \mathbf{U}_i^T \mathbf{U}_j$ (اسکات، ۲۰۰۵)، که در آن بردارهای تصادفی مستقل $\mathbf{U}_1, \dots, \mathbf{U}_p$ با توزیع یکنواخت روی سطح گوی $(n-1)$ -بُعدی با شعاع یک، $S_{n-1}(1)$ ، و دارای تابع چگالی احتمال

$$f_{\mathbf{U}}(\mathbf{u}) = \frac{\Gamma\left(\frac{n-1}{\nu}\right)}{\sqrt{\pi} \frac{n-1}{\nu}} I(\mathbf{u}^T \mathbf{u} - 1), \quad \mathbf{u} \in R^{n-1},$$

هستند (فولند، ۲۰۱۳)، که در آن $I(x)$ تابع نشانگر روی نقطه $x = 0$ است. بنابراین با توجه به نحوه نمایش $r_{ij} = \mathbf{U}_i^T \mathbf{U}_j$ برای $i_1 \neq i_2$ و $j_1 \neq j_2$ متغیرهای تصادفی $r_{i_1 j_1}$ و $r_{i_2 j_2}$ مستقل‌اند. در حالی که تنها یکی از $i_1 \neq i_2$ یا $j_1 \neq j_2$ برقرار باشد، $r_{i_1 j_1}$ و $r_{i_2 j_2}$ حداقل به شکل مجانبی مستقل هستند (سیرواستاوا، ۲۰۰۵). چون نتایج به صورت مجانبی بدست آمده است، فرض می‌شود برای حداقل یکی از دو حالت $i_1 \neq i_2$ یا $j_1 \neq j_2$ متغیرهای تصادفی $r_{i_1 j_1}$ و $r_{i_2 j_2}$ همواره مستقل هستند.

لم ۱: فرض کنید $U = [U_1, \dots, U_m]^T$ یک بردار تصادفی با توزیع یکنواخت روی $S_m(1)$ باشد. در این صورت برای هر بردار ثابت $v \in S_m(1)$ داریم $E[|v^T U|] = \Gamma_{m+1}(1)$.

برهان: از آنجا که v و U هر دو بردارهایی روی $S_m(1)$ هستند، $v^T v = U^T U = 1$ در نتیجه

$$v^T U = (v^T v)(U^T U) \cos(\theta) = \cos(\theta),$$

که در آن θ زاویه بین دو بردار مذکور است. بنابراین، $v^T U$ تنها به زاویه بین v و U و نه به مکان آنها روی $S_m(1)$ وابسته است. از این رو می‌توان فرض کرد $v = [1, 0, \dots, 0]^T \in S_m(1)$ از سوی دیگر در مختصات قطبی می‌توان نوشت:

$$U_j = \begin{cases} \cos(\theta_j) \sin(\theta_1) \cdots \sin(\theta_{j-1}) & j = 1, \dots, m-1 \\ \sin(\theta_{m-1}) \sin(\theta_1) \cdots \sin(\theta_{m-2}) & j = m \end{cases}$$

که در آن $\theta_j \in [0, \pi]$ ، $j = 1, \dots, m-2$ و $\theta_{m-1} \in [0, 2\pi]$ (سزبوسکی، ۱۹۹۸). در نتیجه

$$\begin{aligned} E[|v^T U|] &= \int_{\mathbb{R}^m} |u_1| f_U(u) dU \\ &= \int_0^\pi \cdots \int_0^\pi \int_0^{2\pi} |\cos(\theta_1)| \prod_{j=1}^{m-2} \sin^{m-j-1}(\theta_j) \frac{\Gamma(\frac{m}{2})}{2\pi^{\frac{m}{2}}} d\theta_{m-1} \cdots d\theta_1 \\ &= \frac{\Gamma(\frac{m}{2})}{\pi^{\frac{m}{2}-1}} \left[\prod_{j=2}^{m-2} \int_0^\pi \sin^{m-j-1}(\theta_j) d\theta_j \right] \int_0^\pi |\cos(\theta_1)| \sin^{m-2}(\theta_1) d\theta_1. \quad (3) \end{aligned}$$

با تغییر متغیر $u = \sin^2(\theta_j)$ داریم:

$$\begin{aligned} \int_0^{\frac{\pi}{2}} \sin^{m-j-1}(\theta_j) d\theta_j &= \frac{1}{2} \int_0^1 u^{\frac{m-j}{2}-1} (1-u)^{\frac{1}{2}-1} du = \frac{1}{2} \frac{\Gamma(\frac{m-j}{2}) \sqrt{\pi}}{\Gamma(\frac{m-j+1}{2})} \\ \int_{\frac{\pi}{2}}^\pi \sin^{m-j-1}(\theta_j) d\theta_j &= \frac{1}{2} \frac{\Gamma(\frac{m-j}{2}) \sqrt{\pi}}{\Gamma(\frac{m-j+1}{2})}. \end{aligned}$$

در نتیجه

$$\int_0^\pi \sin^{m-j-1}(\theta_j) d\theta_j = \frac{\Gamma(\frac{m-j}{2})\sqrt{\pi}}{\Gamma(\frac{m-j+1}{2})}. \quad (4)$$

دوباره با تغییر متغیر $u = \sin^2(\theta_1)$ داریم:

$$\int_0^{\frac{\pi}{2}} \cos(\theta_1) \sin^{m-j-1}(\theta_1) d\theta_1 = \frac{1}{m-1} \text{ و } \int_{\frac{\pi}{2}}^\pi \cos(\theta_1) \sin^{m-j-1}(\theta_1) d\theta_1 = -\frac{1}{m-1}.$$

پس

$$\int_0^\pi |\cos(\theta_1)| \sin^{m-2}(\theta_1) d\theta_1 = \frac{2}{m-1}. \quad (5)$$

حال با جایگذاری (۴) و (۵) در معادله (۳) داریم:

$$E[|v^T U|] = \frac{\Gamma(\frac{m}{2})}{\sqrt{\pi}\Gamma(\frac{m+1}{2})} = \Gamma_{m+1}(1).$$

در نتیجه با استفاده از لم ۱ و اینکه $r_{i_1 j_1}$ و $r_{i_2 j_2}$ برای حداقل یکی از دو حالت $i_1 \neq i_2$ یا $j_1 \neq j_2$ همواره مستقل هستند، تحت فرضیه (۱) می‌توان نوشت:

$$\begin{aligned} \text{Var} D_{np} &= \sum_{i=2}^p \sum_{j=1}^{i-1} \text{Var} \sqrt{|r_{ij}|} \\ &= \sum_{i=2}^p \sum_{j=1}^{i-1} [E|r_{ij}| - \Gamma_n^2(\frac{1}{p})] \\ &= \frac{p(p-1)}{2} [\Gamma_n(1) - \Gamma_n^2(\frac{1}{p})] = \tau_{np}^2. \end{aligned} \quad (6)$$

قضیه ۱: تحت فرضیه (۱)، وقتی p و n به ∞ همگرا می‌شوند، آنگاه

$$D_{np}^* = \frac{D_{np} - \mu_{np}}{\tau_{np}} \xrightarrow{d} N(0, 1).$$

برهان: فرض کنید

$$d_{nl} = D_{nl} - \mu_{nl} = \sum_{i=2}^{\ell} \left(\sum_{j=1}^{i-1} \sqrt{|r_{ij}|} - (i-1)\Gamma_n\left(\frac{1}{\nu}\right) \right)$$

$$X_{nl} = d_{nl} - d_{n(\ell-1)} = \sum_{j=1}^{\ell-1} \sqrt{|r_{\ell j}|} - (\ell-1)\Gamma_n\left(\frac{1}{\nu}\right), \quad \ell = 2, \dots, p,$$

که در آن $d_{n1} = 0$. چون تحت فرضیه H_0 می‌توان نوشت $r_{ij} = \mathbf{U}_i^T \mathbf{U}_j$ ، با فرض $F_{n,\ell-1}$ و بنابه لم داریم ۱،

$$E[|r_{\ell j}| | F_{n,\ell-1}] = E[|\mathbf{U}_i^T \mathbf{U}_j| | F_{n,\ell-1}] = \Gamma_n(1), \quad j = 1, \dots, \ell-1$$

که در آن $E[\cdot | \cdot]$ امید ریاضی شرطی است. در نتیجه

$$E[X_{nl} | F_{n,\ell-1}] = E[d_{nl} | F_{n,\ell-1}] - d_{n(\ell-1)} = 0$$

یا $E[d_{nl} | F_{n,\ell-1}] = d_{n(\ell-1)}$ دنباله $\{d_{nl}, \ell = 2, \dots, k\}$ یک مارتینگل (هال و هی دی، ۱۹۸۰) با تفاضل‌های X_{n2}, \dots, X_{nk} است. علاوه بر این چون τ_{np} مقداری ثابت است، به طور مشابه برای $n \geq 2$ دنباله $\{d_{nl}, \ell = 2, \dots, k\}$ نیز یک مارتینگل با تفاضل‌های $Y_{n2} = \frac{X_{n2}}{\tau_{np}}, \dots, Y_{nk} = \frac{X_{nk}}{\tau_{np}}$ است. حال بنابر قضیه ۲.۳ در مک‌لش (۱۹۷۴) برای اینکه

$$D_{np}^* = \sum_{\ell=2}^p Y_{n\ell} \xrightarrow{d} N(0, 1)$$

کافیست ثابت شود وقتی $p, n \rightarrow \infty$ آنگاه

$$\sum_{\ell=2}^p E[Y_{n\ell}^2] \rightarrow 0 \quad \text{و} \quad E\left(\sum_{\ell=2}^p Y_{n\ell}^2 - 1\right)^2 \rightarrow 0.$$

حال تحت فرضیه (۱) با استفاده از (۲) داریم:

$$\begin{aligned}
 E[X_{n\ell}^*] &= E\left\{\left[\sum_{j=1}^{\ell-1} \sqrt{|r_{\ell j}| - (\ell-1)\Gamma_n\left(\frac{1}{\varphi}\right)}\right]^2\right\} \\
 &= \sum_{j=1}^{\ell-1} E|r_{\ell j}| + \sum_{j_1=1}^{\ell-1} \sum_{\substack{j_2=1 \\ j_2 \neq j_1}}^{\ell-1} E\sqrt{|r_{\ell j_1}|} E\sqrt{|r_{\ell j_2}|} - [(\ell-1)\Gamma_n\left(\frac{1}{\varphi}\right)]^2 \\
 &= (\ell-1)\Gamma_n(1) + \Gamma_n^*\left(\frac{1}{\varphi}\right)(\ell-1)(\ell-2) - ((\ell-1)\Gamma_n\left(\frac{1}{\varphi}\right))^2 \\
 &= (\ell-1)[\Gamma_n(1) - \Gamma_n^*\left(\frac{1}{\varphi}\right)]
 \end{aligned}$$

که در آن برابری سوم از (۲) و رابطه

$$\begin{aligned}
 \sum_{j_1=1}^{\ell-1} \sum_{\substack{j_2=1 \\ j_2 \neq j_1}}^{\ell-1} E\sqrt{|r_{\ell j_1}|} E\sqrt{|r_{\ell j_2}|} &= \sum_{j_1=1}^{\ell-1} \sum_{j_2=1}^{\ell-1} \underbrace{E\sqrt{|r_{\ell j_1}|} E\sqrt{|r_{\ell j_2}|} I(j_2 \neq j_1)}_{\Gamma_n^*\left(\frac{1}{\varphi}\right)} \\
 &= \Gamma_n^*\left(\frac{1}{\varphi}\right) \sum_{j_1=1}^{\ell-1} \sum_{j_2=1}^{\ell-1} (1 - I(j_2 = j_1)) \\
 &= \Gamma_n^*\left(\frac{1}{\varphi}\right) [(\ell-1)^2 - (\ell-1)] \\
 &= \Gamma_n^*\left(\frac{1}{\varphi}\right)(\ell-1)(\ell-2).
 \end{aligned}$$

بدست آمده است. پس $\sum_{\ell=2}^p EX_{n\ell}^* = \tau_{np}^*$. فرض کنید $U_{\ell j} = \sqrt{|r_{\ell j}| - \Gamma_n\left(\frac{1}{\varphi}\right)}$ ، به گونه‌ای که $EU_{\ell j} = 0$ و $j = 1, \dots, \ell-1$. با توجه به آنچه که در مورد استقلال $r_{\ell j_1}$ و $r_{\ell j_2}$ برای $j_1 \neq j_2$ تحت فرضیه H_0 بیان شد، $U_{\ell j_1}$ و $U_{\ell j_2}$ حداقل مجانباً مستقل خواهند بود. حال از آنجا که $X_{n\ell} = \sum_{j=1}^{\ell-1} U_{\ell j}$ می‌توان نشان داد

$$X_{n\ell}^* = \sum_{j=1}^{\ell-1} U_{\ell j}^* + \epsilon \sum_{j_1=2}^{\ell-1} \sum_{j_2=1}^{j_1-1} U_{\ell j_1}^* U_{\ell j_2}^* + U_0,$$

که در آن U ، جمله‌ای با $EU = 0$ تحت فرضیه صفر است. پس

$$EX_{n\ell}^* = \sum_{j=1}^{\ell-1} EU_{\ell j}^* + \epsilon \sum_{j_1=2}^{\ell-1} \sum_{j_2=1}^{j_1-1} E[U_{\ell j_1}^*] E[U_{\ell j_2}^*]. \quad (7)$$

اما بر اساس رابطه (۲) می‌توان نوشت

$$\begin{aligned} EU_{\ell j}^* &= E(\sqrt{|r_{\ell j}|} - \Gamma_n(\frac{1}{\nu}))^\nu \\ &= E|r_{\ell j}|^\nu - \nu \Gamma_n(\frac{1}{\nu}) E|r_{\ell j}|^{\nu-1} + \epsilon \Gamma_n^\nu(\frac{1}{\nu}) E|r_{\ell j}| - \nu \Gamma_n^\nu(\frac{1}{\nu}) E\sqrt{|r_{\ell j}|} + \Gamma_n^\nu(\frac{1}{\nu}) \\ &= \Gamma_n(\nu) - \nu \Gamma_n(\frac{1}{\nu}) \Gamma_n(\frac{\nu-1}{\nu}) + \epsilon \Gamma_n^\nu(\frac{1}{\nu}) \Gamma_n(1) - \nu \Gamma_n^\nu(\frac{1}{\nu}) = c_\epsilon \end{aligned} \quad (8)$$

$$\begin{aligned} EU_{\ell j}^* &= E(\sqrt{|r_{\ell j}|} - \Gamma_n(\frac{1}{\nu}))^\nu \\ &= E|r_{\ell j}| - \nu \Gamma_n(\frac{1}{\nu}) E\sqrt{|r_{\ell j}|} + \Gamma_n^\nu(\frac{1}{\nu}) \\ &= \Gamma_n(1) - \Gamma_n^\nu(\frac{1}{\nu}) = c_\nu. \end{aligned} \quad (9)$$

از جایگذاری (۸) و (۹) در رابطه (۷)

$$EX_{n\ell}^* = (\ell - 1)c_\epsilon + \nu c_\nu^{\nu-1} (\ell - 1)(\ell - \nu), \quad \ell = 2, \dots, k.$$

در نتیجه

$$\sum_{\ell=2}^p E[X_{n\ell}^*] = \frac{p(p-1)}{\nu} (c_\epsilon + \nu(p-\nu)c_\nu^{\nu-1}) = \frac{(c_\epsilon + \nu(p-\nu)c_\nu^{\nu-1})}{\frac{p(p-1)}{\nu} [\Gamma_n(1) - \Gamma_n^\nu(\frac{1}{\nu})]^\nu}.$$

حال بنابه لم ۱۰۵ در جیانگ و یانگ (۲۰۱۳)

$$\frac{\Gamma(\frac{n-1}{\nu})}{\Gamma(\frac{k+n-1}{\nu})} = (\frac{n-1}{\nu})^{-\frac{k}{\nu}} + o(1),$$

وقتی $n \rightarrow \infty$ در نتیجه،

$$\Gamma_n(k) = \frac{\Gamma(\frac{n-1}{\nu})\Gamma(\frac{k+1}{\nu})}{\sqrt{\pi}\Gamma(\frac{k+n-1}{\nu})} = \frac{\Gamma(\frac{k+1}{\nu})}{\sqrt{\pi}} \left(\frac{n-1}{\nu}\right)^{-\frac{k}{\nu}} + o(1), \quad (10)$$

وقتی $n \rightarrow \infty$ حال با به کارگیری (۱۰) در روابط (۸) و (۹)، وقتی $n \rightarrow \infty$ داریم:

$$c_{\nu} = k_0 \left(\frac{n-1}{\nu}\right)^{-1} + o(1) \quad \text{و} \quad c_{\nu} = k_1 \left(\frac{n-1}{\nu}\right)^{-\frac{1}{\nu}} + o(1),$$

که در آن k_0 و k_1 اعدادی ثابت و متناهی هستند. پس،

$$\sum_{\ell=\nu}^p E[Y_{n\ell}^{\nu}] = \frac{k_0 + \nu(p-\nu)k_1^{\nu}}{\frac{p(p-1)}{\nu} \left(\frac{1}{\sqrt{\pi}} - \frac{\Gamma^{\nu}(\frac{\nu}{\nu})}{\pi}\right)} + o(1) \rightarrow 0. \quad (11)$$

وقتی p و n به بینهایت همگرا می‌شود. همچنین، از استقلال $Y_{n\ell_1}$ و $Y_{n\ell_2}$ برای $\ell_1 \neq \ell_2$ نتیجه می‌شود

$$\begin{aligned} 0 &\leq E\left[\left(\sum_{\ell=\nu}^p Y_{n\ell}^{\nu} - 1\right)^{\nu}\right] = \sum_{\ell=\nu}^p E[Y_{n\ell}^{\nu}] + \sum_{\substack{\ell_1=\nu \\ \ell_2=\nu \\ \ell_2 \neq \ell_1}}^p \sum_{\ell_2=\nu}^p E[Y_{n\ell_1}^{\nu} Y_{n\ell_2}^{\nu}] - 1 \\ &= \sum_{\ell_1=\nu}^p E[Y_{n\ell_1}^{\nu}] + \sum_{\substack{\ell_1=\nu \\ \ell_2=\nu \\ \ell_2 \neq \ell_1}}^p \sum_{\ell_2=\nu}^p E[Y_{n\ell_1}^{\nu}] E[Y_{n\ell_2}^{\nu}] - 1 \\ &\leq \sum_{\ell=\nu}^p E[Y_{n\ell}^{\nu}] + \left(\sum_{\ell=\nu}^p E[Y_{n\ell}^{\nu}]\right)^{\nu} - 1 \\ &= \sum_{\ell=\nu}^p E[Y_{n\ell}^{\nu}]. \end{aligned} \quad (12)$$

بنابراین، از روابط (۱۱) و (۱۲) خواهیم داشت

$$E\left[\left(\sum_{\ell=\nu}^p Y_{n\ell}^{\nu} - 1\right)^{\nu}\right] \rightarrow 0. \quad (13)$$

۲۴۴ آزمون‌های استقلال در داده‌های نرمال با بُعد بالا

در نهایت، بنابر قضیه ۲.۳ (مک‌لش، ۱۹۷۴)، وقتی $p, n \rightarrow \infty$ ، از روابط (۱۱) و (۱۳) نتیجه می‌شود

$$D_{np}^* = \frac{D_{np} - \mu_{np}}{\tau_{np}} \xrightarrow{d} N(0, 1).$$

۳ مطالعه شبیه‌سازی

در این بخش به منظور بررسی عملکرد آزمون D_{np}^* و همچنین مقایسه آن با آزمون‌های t_{np}^c ، T_{np}^* ، t_{np}^* و T_{np}^c مطالعه شبیه‌سازی انجام شده است. برآورد اندازه و توان برای هر یک از آزمون‌ها بر اساس ۱۰۰۰۰ تکرار نمونه‌گیری مستقل و در سطح معنی‌داری $\alpha = 0.05$ محاسبه شده است. در اینجا بُعد p و اندازه نمونه n مقادیری به ترتیب از مجموعه‌های $\{5, 50, 500\}$ و $\{6, 60, 300\}$ در نظر گرفته شده است. برای بررسی عملکرد آزمون‌ها از دیدگاه اندازه آزمون، معیار میانگین خطای نسبی به صورت

$$\text{ARE} = \frac{100}{M\alpha} \sum_{j=1}^M |\hat{\alpha}_j - \alpha|, \quad (14)$$

استفاده شده است، که در آن M برابر کلیه ترکیبات از مقادیر p و n (در اینجا، $M = 9$) و $\hat{\alpha}_1, \dots, \hat{\alpha}_M$ مقادیر اندازه‌های تجربی یا برآوردشده متناظر با ترکیبات مد نظر است. مقادیر کوچک (۱۴) برای یک آزمون بیانگر عملکرد بهتر آن آزمون در حفظ سطح معنی‌داری α است. در این مطالعه به ازای هر مقدار از بُعد p ، بردار میانگین برابر بردار صفر و ماتریس کوواریانس به صورت $\Sigma^{(\phi)} = (1 - \phi)I_p + \phi J_p$ انتخاب شده است، که در آن $0 < \phi < 1$ مقداری ثابت، I_p ماتریس همانی $p \times p$ و $J_p = \mathbf{1}_p \mathbf{1}_p^T$ که $\mathbf{1}_p$ بردار ستونی $1 \times p$ از یک‌ها است. در حقیقت، مقدار پارامتر ϕ میزان گسیختگی از فرضیه H_0 را نشان می‌دهد. در این مطالعه، ϕ عضوی از مجموعه $\{0, 0.2, 0.5\}$ انتخاب شده است. برای هر ترکیب از مقادیر n ، p و ϕ ، فرایند نمونه‌گیری با اندازه نمونه n از توزیع نرمال p -متغیره $N_p(0, \Sigma^{(\phi)})$ به تعداد ۱۰۰۰۰ بار تکرار و مقادیر اندازه و توان تجربی آزمون‌های t_{np}^c ، T_{np}^* ، t_{np}^* و T_{np}^c محاسبه شده است. برای نمونه‌های تحت فرضیه (۱) ($\phi = 0$) اندازه‌های تجربی آزمون‌ها در جدول ۱ و برای نمونه‌های تحت فرضیه مقابل ($\phi = 0.2, 0.5$) توان‌های تجربی متناظر با این آزمون‌ها در جداول ۲ و ۳ آورده شده است. در حالتی که هر دو اندازه نمونه n و بُعد p اعداد بزرگی هستند، تمامی آزمون‌ها دارای اندازه‌های تجربی نزدیک به مقدار اسمی ۰.۰۵ هستند. برای داده‌های با بُعد بالا که در آن اندازه نمونه n

جدول ۱: اندازه‌های آزمون تجربی به ازای $\phi = 0$

D_{np}^*	T_{np}^c	t_{np}^c	T_{np}^*	t_{np}^*	p	n
۰/۰۵۱۱	۰/۰۴۶۲	۰/۰۴۶۵	۰/۰۴۲۸	۰/۰۴۴۶	۵	
۰/۰۴۸۹	۰/۱۰۷۴	۰/۰۶۴۳	۰/۱۳۴۳	۰/۰۴۴۹	۵۰	۶
۰/۰۴۷۳	۰/۱۵۸۶	۰/۰۶۵۶	۰/۲۶۳۰	۰/۰۴۱۴	۵۰۰	
۰/۰۵۰۰	۰/۰۵۳۱	۰/۰۵۱۴	۰/۰۴۵۶	۰/۰۴۳۶	۵	
۰/۰۵۱۲	۰/۰۵۱۵	۰/۰۵۱۱	۰/۰۵۲۳	۰/۰۵۰۳	۵۰	۶۰
۰/۰۴۸۶	۰/۰۵۲۰	۰/۰۵۰۵	۰/۰۵۱۴	۰/۰۴۹۰	۵۰۰	
۰/۰۴۹۸	۰/۰۵۰۶	۰/۰۴۹۶	۰/۰۴۲۹	۰/۰۴۲۲	۵	
۰/۰۴۹۷	۰/۰۵۴۷	۰/۰۵۴۶	۰/۰۵۰۹	۰/۰۵۰۹	۵۰	۳۰۰
۰/۰۵۱۰	۰/۰۴۷۹	۰/۰۴۷۰	۰/۰۴۹۰	۰/۰۴۹۱	۵۰۰	
۲/۰۰۰۰	۴۰/۸۴۴۴	۹/۸۶۶۷	۷۱/۴۶۶۷	۸/۰۸۸۹	ARE	

در مقایسه با بُعد p بسیار کوچک است (به عنوان مثال، $n = 6$ و $p = 500$)، اندازه تجربی آزمون‌های T_{np}^* و T_{np}^c به ترتیب با مقادیر $0/2630$ و $0/1586$ بسیار بالاتر از سطح آزمون اسمی $\alpha = 0/05$ است. این در حالی است که اندازه تجربی آزمون‌های t_{np}^* و به ویژه D_{np}^* نزدیک به مقدار اسمی $\alpha = 0/05$ است. در این میان، برای ابعاد بالاتر، آزمون t_{np}^c در رابطه با حفظ سطح معنی‌داری اسمی $0/05$ نسبت به آزمون‌های t_{np}^* و D_{np}^* عملکرد ضعیف‌تر و نسبت به آزمون‌های T_{np}^* و T_{np}^c عملکرد بهتری دارد. با افزایش اندازه نمونه، هر پنج آزمون سطح معنی‌داری اسمی $0/05$ به خوبی حفظ نموده و اندازه‌های تجربی یکسانی را ارائه داده‌اند. به هر حال برای داده‌های با ابعاد بالا که در آنها اختلاف بُعد با اندازه نمونه اعدادی بسیار بزرگ است استفاده از آزمون‌های T_{np}^* ، t_{np}^c و T_{np}^c توصیه نمی‌شود. همچنین با توجه به یافته‌های جدول ۱، روش آزمون معرفی‌شده در این مقاله مبنی بر آماره D_{np}^* در مقایسه با سایر آزمون‌ها از کمترین مقدار میانگین خطای نسبی با مقدار ARE برابر ۲ برخوردار است. به عبارت دیگر، اندازه آزمون تجربی روش معرفی‌شده در این مقاله در مقایسه با سایر روش‌های موجود به مقدار $\alpha = 0/05$ نزدیک‌تر است.

در حالت کلی، با ملاک مقدار میانگین خطای نسبی، به لحاظ حفظ سطح معنی‌داری اسمی به ترتیب آزمون‌های D_{np}^* ، t_{np}^* ، t_{np}^c و T_{np}^c از اولویت برخوردارند. لازم به ذکر است که بر اساس شبیه‌سازی‌هایی که به دلیل حجم زیاد در این مقاله ذکر نشده‌اند نیز همواره اولویت‌بندی فوق حفظ شده است. با توجه به جداول ۲ و ۳ تمامی آزمون‌ها نسبت به انحرافات بسیار کوچک بسیار حساس‌اند به طوری که مثلاً با انحرافی به اندازه $\phi = 0/05$ از فرضیه صفر برای اندازه نمونه متوسط،

جدول ۲: توان‌های آزمون تجربی به ازای $\phi = 0.02$

D_{np}^*	T_{np}^c	t_{np}^c	T_{np}^*	t_{np}^*	p	n
0.529	0.461	0.470	0.427	0.445	5	
0.516	0.1133	0.702	0.1401	0.498	50	6
0.998	0.2028	0.1405	0.2696	0.1057	500	
0.510	0.605	0.583	0.531	0.515	5	
0.928	0.1582	0.1557	0.1147	0.1124	50	60
0.9222	0.9735	0.9726	0.9624	0.9613	500	
0.674	0.911	0.901	0.800	0.788	5	
0.6361	0.8167	0.8165	0.7513	0.7510	50	300
1	1	1	1	1	500	

جدول ۳: توان‌های آزمون تجربی به ازای $\phi = 0.05$

D_{np}^*	T_{np}^c	t_{np}^c	T_{np}^*	t_{np}^*	p	n
0.521	0.463	0.513	0.436	0.465	5	
0.789	0.1387	0.1110	0.1476	0.793	50	6
0.3830	0.4189	0.4588	0.4310	0.4083	500	
0.743	0.1081	0.1041	0.977	0.951	5	
0.7024	0.8309	0.8302	0.7797	0.7783	50	60
1	1	1	1	1	500	
0.2826	0.3992	0.3987	0.3754	0.3749	5	
1	1	1	1	1	50	300
1	1	1	1	1	500	

$n = 60$ ، توان‌های تجربی نزدیک به یک برای ابعاد $p \geq 50$ حاصل شده است. در ابعاد بالاتر به دلیل افزایش تعداد انحرافات هر چند کوچک از فرضیه صفر و در نتیجه بزرگ تر شدن مجموع این انحرافات توان‌های بالاتر از طرف تمامی آزمون‌ها به دست آمده است. به هر حال برای همین انحرافات کوچک، آزمون‌های T_{np}^c و T_{np}^* دارای توان بالاتری نسبت به آزمون‌های t_{np}^c ، t_{np}^* و D_{np}^* هستند و در اکثر مواقع آماره آزمون معرفی شده دارای توان کمتری نسبت به سایر آزمون‌ها است. این امر با توجه به بالا بودن مقدار میانگین خطای نسبی سایر آزمون‌ها امری قابل پیش‌بینی است. در

حقیقت، آزمون‌هایی با اندازه‌های بزرگتر از سطح اسمی α همواره توان‌های بزرگ کاذب و فریبنده دارند (کریستنسن و رنچر، ۱۹۹۷؛ جیانگ و یانگ، ۲۰۱۳؛ ژانگ و همکاران، ۲۰۱۷). با توجه به جداول ۲ و ۳، همواره افزایش اندازه نمونه یا بُعد یا هر دو منجر به توان‌های بالاتر از تمامی آزمون‌ها می‌شود.

۴ کاربرد روش پیشنهادی

این بخش به ارائه کاربردی از روش آزمون پیشنهادی بر مجموعه داده سطح بیان ژن تومور پروستات اختصاص یافته است^۴ (سینگ و همکاران، ۲۰۰۲؛ دتلینگ و بوهمن، ۲۰۰۲). این داده‌ها را می‌توان در بسته `spls`^۵ از نرم‌افزار R یافت. این مجموعه شامل اندازه‌گیری سطح بیان ۶۰۳۳ ژن به کمک فناوری آفی‌متریکس^۶ روی دو نمونه با اندازه‌های $n_1 = 52$ و $n_2 = 50$ به ترتیب از بافت‌های تومور پروستات و بافت‌های سالم است. سینگ و همکاران (۲۰۰۲) بر اساس همبستگی‌های بین مقادیر سطوح بیان و همچنین همبستگی بین این سطوح با متغیرهای بالینی و آسیب‌شناختی (مانند سن بیمار، نتیجه آزمون‌های تشخیصی PSA و رتبه GS، غیره) به تشریح رفتار بالینی سرطان پروستات پرداختند.

در سطح $\alpha = 0.05$ استقلال همزمان $p = 1000$ سطح بیان ژن نخست از بین ۶۰۳۳ سطح بیان در هر دو جامعه بافت‌های سالم (I) و بافت‌های سرطانی (II) بررسی شده است. نتایج حاصل از آزمون‌های استقلال D_{np}^* ، T_{np}^c ، t_{np}^c ، T_{np}^* ، t_{np}^* روی دو جامعه در جدول ۴ آمده است. همان‌طور که ملاحظه

جدول ۴: آزمون استقلال در داده‌های تومور پروستات

مقدار آماره			
روش آزمون	بافت‌های سرطانی	بافت‌های سالم	مقدار بحرانی
$t_{n,p}^*$	۲۴۴۲٫۶۰۲	۳۰۱۲٫۱۰۴	۱٫۹۵۹۹
$T_{n,p}^*$	9.5367×10^{15}	۳۶۵۷۵٫۸۷	۱٫۹۵۹۹
$t_{n,p}^c$	۱۳۴۴۸٫۲۹	۱۶۲۹۲٫۹۵	۱۲۵۵۷۲٫۷
$T_{n,p}^c$	4.7636×10^{18}	۱۸۳۹۴۳٫۸۸	۱۲۵۵۷۲٫۷
$D_{n,p}^*$	۷۳۸٫۶۷۹۷	۸۳۲٫۲۲۱۱	۱٫۹۵۹۹

می‌شود فرضیه استقلال سطوح بیان این ۱۰۰۰ ژن در هر دو جامعه بافت‌های دارای تومور پروستات و بافت‌های سالم به شدت رد می‌شود و همبستگی خطی نسبتاً قوی بین سطوح بیان این ژن‌ها وجود دارد.

^۴Prostate Tumor Gene Expression dataset

^۵<https://CRAN.R-project.org/package=spls>

^۶Affymetrix technology

۵ بحث و نتیجه‌گیری

در این مقاله، آزمون استقلال در داده‌های نرمال با بُعد بالا بررسی و آزمون‌های جدید ارائه شد. با شبیه‌سازی، اندازه و توان آزمون پیشنهادی با آزمون‌های موجود محاسبه و نشان داده شد که آزمون معرفی شده به ویژه در حالت داده‌های با بُعد بالا، در مقایسه با دیگر آزمون‌ها دارای کمترین مقدار میانگین خطای نسبی است. در واقع خطای نوع اول واقعی سایر آزمون‌ها با مقدار اسمی α متفاوت است. چنین آزمون‌هایی به شکل کاذب دارای توان بالاتری هستند. آزمون پیشنهادی به دلیل اندازه آزمون واقعی نزدیک به سطح اسمی دارای کمترین مقدار توان آزمون در مقایسه با سایر آزمون‌ها بود. البته این اختلاف در توان با سایر آزمون‌ها، تنها برای انحرافات جزئی از فرضیه صفر بود و نه برای انحرافات میانگین و بالاتر که این حالت‌ها توان تمامی آزمون‌ها برابر می‌شد. برای داده‌های بُعد بالا، استفاده از آزمون پیشنهادی به دلیل مقدار میانگین خطای نسبی کم و به عبارت دیگر اندازه آزمون نزدیک به سطح اسمی α توصیه می‌شود.

مراجع

- نجاززاده، د. (۱۳۹۸)، آزمون همزمان استقلال برای زیربردارهای چند بردار با بُعد نسبتاً بالای نرمال چندمتغیره، مجله علوم آماری ایران، ۱۳، ۲۱۷-۲۳۳.
- Anderson, T. W. (2003), *An Introduction to Multivariate Statistical Analysis (3rd ed.)*, John Wiley & Sons, New York.
- Bai, Z., and Saranadasa, H. (1996), Effect of High Dimension: by an Example of a Two Sample Problem, *Statistica Sinica*, **6**, 311-329.
- Chang, S., and Qi, Y. (2018), On Schott's and Mao's Test Statistics for Independence of Normal Random Vectors, *Statistics & Probability Letters*, **140**, 132-141.
- Chen, S., and Mudholkar, G. S. (1990), Null Distribution of the Sum of Squared Z-Transforms in Testing Complete Independence, *Annals of the Institute of Statistical Mathematics*, **42**, 149-155.
- Christensen, W. F., and Rencher, A. C. (1997), A Comparison of Type I Error Rates and Power Levels for Seven Solutions to the Multivariate Behrens-Fisher Problem, *Communications in Statistics-Simulation and Computation*, **26**, 1251-1273.
- Dempster, A. P. (1958), A High Dimensional Two Sample Significance Test, *The Annals of Mathematical Statistics*, 995-1010.
- Dempster, A. P. (1960), A Significance Test for the Separation of Two Highly Multivariate Small Samples, *Biometrics*, **16**, 41-50.

- Dettling, M., and Bühlmann, P. (2002). Supervised Clustering of Genes. *Genome biology* **3**, research0069-1.
- Folland, G. B. (2013), *Real Analysis: Modern Techniques and Their Applications*, John Wiley & Sons, New York.
- Giraud, C. (2014), *Introduction to High-Dimensional Statistics*, Chapman and Hall/CRC, Hoboken, New Jersey.
- Gupta, A., and Song, D. (1997), Lp-Norm Spherical Distribution, *Journal of Statistical Planning and Inference*, **60**, 241-260.
- Jiang, D., Jiang, T. and Yang, F. (2012), Likelihood Ratio Tests for Covariance Matrices of High-Dimensional Normal Distributions, *Journal of Statistical Planning and Inference*, **142**, 2241-2256.
- Jiang, T., and Qi, Y. (2015), Likelihood Ratio Tests for High-Dimensional Normal Distributions, *Scandinavian Journal of Statistics*, **42**, 988-1009.
- Jiang, T., and Yang, F. (2013), Central Limit Theorems for Classical Likelihood Ratio Tests for High-Dimensional Normal Distributions, *The Annals of Statistics*, **41**, 2029-2074.
- Liang, J., Tang, M. L. and Chan, P. S. (2009), A Generalized Shapiro-Wilk W Statistic for Testing High-Dimensional Normality, *Computational Statistics & Data Analysis*, **53**, 3883-3891.
- Mao, G. (2014), A New Test of Independence for High-Dimensional Data, *Statistics & Probability Letters*, **93**, 14-18.
- McLeish, D. L. (1974), Dependent Central Limit Theorems and Invariance Principles, *The Annals of Probability*, **2**, 620-628.
- Muirhead, R. J. (2005), *Aspects of Multivariate Statistical Theory*, John Wiley & Sons, Inc, Hoboken, New Jersey.
- Hall, P., and Heyde, C. C. (2014), *Martingale Limit Theory and Its Application*, Academic press, New York.
- Schott, J. R. (2005), Testing for Complete Independence in High Dimensions, *Biometrika*, **92**, 951-956.
- Singh, D., Febbo, P. G., Ross, K., Jackson, D. G., Manola, J., Ladd, C., Tamayo, P., Renshaw, A. A., D'Amico, A. V., Richie, J. P. and Lander, E. S. (2007). Gene Expression Correlates of Clinical Prostate Cancer Behavior, *Cancer cell*, **1**, 203-209.
- Srivastava, M. S. (2005), Some Tests Concerning the Covariance Matrix in High Dimensional Data, *Journal of the Japan Statistical Society*, **35**, 251-272.
- Srivastava, M. S. (2006), Some Tests Criteria for the Covariance Matrix with Fewer Observations than the Dimension, *Acta et Commentations Universitatis Tartuensis de Mathematica*, **10**, 77-93.

- Srivastava, M. S., and Reid, N. (2012), Testing the Structure of the Covariance Matrix with Fewer Observations than the Dimension, *Journal of Multivariate Analysis*, **112**, 156-171.
- Szabowski, P. (1998), Uniform Distributions on Spheres Infinite Dimensional and Their Generalizations, *Journal of Multivariate Analysis*, **64**, 103-117.
- Tan, M., Fang, H. B., Tian, G. L. and Wei, G. (2005), Testing Multivariate Normality in Incomplete Data of Small Sample Size, *Journal of Multivariate Analysis*, **93**, 164-179.
- Zhang, J. T., Guo, J. and Zhou, B. (2017), Linear Hypothesis Testing in High-dimensional One-way MANOVA, *Journal of Multivariate Analysis*, **155**, 200-216.