

آزمون همزمان استقلال برای زیربردارهای چند بردار با بُعد نسبتاً بالای نرمال چندمتغیره

داریوش نجارزاده

گروه آمار، دانشکده علوم ریاضی، دانشگاه تبریز

چکیده: آزمون فرض استقلال میان زیربردارهای یک بردار p -متغیره، به عنوان پیش‌نیاز بسیاری از آزمون‌های آماری، همواره مورد توجه بوده است. وقتی اندازه نمونه n در مقایسه با بُعد p خیلی بزرگ است، آزمون نسبت درستی با توزیع تقریبی χ^2 دو، عملکرد قابل قبولی دارد. برای "داده‌های با بُعد نسبتاً بالا" که در آنها n در قیاس با p چندان بزرگ نیست، تقریب χ^2 دو برای توزیع آماره آزمون نسبت درستی کارایی لازم را ندارد. به عنوان یک حالت جامع‌تر، در این مقاله، آزمونی همزمان در k جامعه p -متغیره نرمال با بُعد نسبتاً بالا که در هر جامعه آزمون استقلال میان زیربردارهای دلخواه آزموده می‌شود، مد نظر قرار گرفته است. به منظور آزمون این فرض، یک تقریب نرمال برای توزیع آماره آزمون نسبت درستی تحت فرض صفر بدست آمده است. علاوه بر این، به منظور تصدیق عملکرد بهتر تقریب نرمال پیشنهادی بر تقریب χ^2 دوی کلاسیک، مطالعه شبیه‌سازی انجام شده است. در پایان، کاربردی از روش پیشنهادی بر مجموعه داده سرطان پرستات ارائه شده است.

واژه‌های کلیدی: توزیع نرمال چندمتغیره، آزمون نسبت درستی، داده‌های با بُعد نسبتاً بالا، آزمون استقلال، تابع گامای چندمتغیره.

۱ مقدمه

در تحلیل‌های چند متغیره در حالت یک جامعه ($k = 1$)، آزمون استقلال زیربردارهای یک بردار p -متغیره نرمال با بعد بالا همواره پیش‌نیاز بسیاری از آزمون‌ها است. در بسیاری از مواقع لازم است همین آزمون به طور همزمان روی چندین ($k > 1$) جامعه p -متغیره نرمال اجرا شود. کاربردهایی از این نوع آزمون استقلال چه در حالت تک جامعه و چه در حالت چندین جامعه بر داده‌های ریزآرایه^۱، داده‌های مالی^۲، داده‌های مصرف^۳، داده‌های ساخت پیشرفته^۴ و داده‌های چندرسانه‌ای^۵ و غیره مشهود است (مائو، ۲۰۱۸؛ چن و همکاران، ۲۰۱۸؛ لئونگ و درتون، ۲۰۱۸). به عنوان مثال، در تحلیل داده‌های ریزآرایه بررسی اینکه بین ژن‌های مختلف استقلال وجود دارد یا نه، همواره از اهمیت خاصی برای تحلیل‌های آتی برخوردار است.

برای آزمون استقلال زیربردارها در حالت $k > 1$ ، روش آزمون نسبت درست‌نمایی^۶ (LRT) با توزیع مجانبی خی‌دو پیشنهاد شده است. ویلکس (۱۹۳۸) نشان داد برای اندازه‌های نمونه‌ای n_1, \dots, n_k بزرگتر از بُعد p با نسبت‌های $\frac{p}{n_i}$ نزدیک به صفر، توزیع تحت فرض صفر آماره LRT، به توزیع مجانبی خی‌دو (χ^2) همگراست. در حالت “داده‌های با بُعد نسبتاً بالا”^۷، با این تعریف که در آن n_i ها از بُعد p بزرگ‌اند و نسبت‌های $\frac{p}{n_i}$ اعدادی نزدیک به یک هستند و همچنین داده‌های با بعد بالا، با این تعریف که در آن n_i ها از بُعد p کوچک‌اند، تقریب خی‌دو به آزمونی با اندازه^۸ بسیار بزرگتر از سطح معنی‌داری اسمی α یا به طور معادل آزمونی با اندازه‌ای متورم منجر می‌شود. این ایراد در مورد تقریب خی‌دو آماره LRT در آزمون فرض‌های دیگر نیز مشاهده شده است، که از این جمله می‌توان به آزمون فرض برابری ماتریس کوواریانس با ماتریس همانی یا همان آزمون گرویت^۹ برای داده‌های با بُعد نسبتاً بالا در کار بای و همکاران (۲۰۰۹) اشاره کرد که در آن برای رفع مشکل متورم بودن اندازه آزمون

¹Microarray data

²Financial data

³Consumer data

⁴Modern manufacturing data

⁵Multimedia data

⁶Likelihood Ratio Test

⁷Moderately high dimensional data

⁸Test size

⁹Sphericity test

LRT به کمک نظریه ماتریس‌های تصادفی و آماره‌های طیفی خطی^{۱۰} روشی برای تصحیح این آماره ارائه شده است. موارد مشابه دیگری را می‌توان در پژوهش‌های اسکات (۲۰۰۵، ۲۰۰۷)، لدویت و ولف (۲۰۰۲)، چن و همکاران (۲۰۱۰)، ژیانگ و همکاران (۲۰۱۲)، ژیانگ و یانگ (۲۰۱۳) یافت. روش کلاسیک LRT برای آزمون استقلال میان زیربردارهای $k = 1$ بردار p -متغیره نرمال با بُعد نسبتاً بالا، قبلاً توسط ژیانگ و یانگ (۲۰۱۳) مورد مطالعه قرار گرفته است. آنها نشان دادند که برای حالتی که بُعد p متناسب با اندازه نمونه n به صورت $y \in (0, 1] \rightarrow \frac{p}{n}$ رشد می‌کند، تقریب خرد و عملاً قابل استفاده نبوده و در این حالت توزیع تحت فرض صفر آماره LRT به جای توزیع خرد به یک توزیع نرمال همگرا است.

به عنوان توسییی از یافته‌های ژیانگ و یانگ (۲۰۱۳) به حالت بیش از یک جامعه، در این مقاله، آزمون همزمان استقلال میان زیربردارهای $k > 1$ بردار p -متغیره نرمال با بُعد نسبتاً بالا مورد مطالعه قرار گرفته است. به بیان دیگر، با این فرض که به ازای هر $k = 1, \dots, k$ بردار تصادفی X_i دارای توزیع نرمال p -متغیره با بردار میانگین μ_i و کوواریانس Σ_i ، یا به اختصار توزیع $N_p(\mu_i, \Sigma_i)$ است و $X_i' = (X_1^{(i)'}, \dots, X_{k_i}^{(i)'})'$ افزای از X_i به $1 \leq k_i \leq p$ زیربردار باشد، هدف این مقاله آزمون استقلال زیربردارهای $X_1^{(i)}, \dots, X_{k_i}^{(i)}$ به طور همزمان در همه k جامعه بر اساس نمونه‌های تصادفی با اندازه‌های نمونه‌ای $n_i > p + 1$ ، $i = 1, \dots, k$ با نسبت‌های $\frac{p}{n_i}$ ها نزدیک به یک است. در این حالت، ثابت می‌شود که توزیع آماره LRT به یک توزیع نرمال همگرا خواهد شد. توجه شود که حالت داده‌های با بعد بالا در این مقاله مورد مطالعه نیست.

در بخش ۲، ضمن استخراج آماره LRT، همگرایی توزیع این آماره به یک توزیع نرمال، با میانگین و واریانس خوش‌تعریف، اثبات شده است. در بخش ۳، یک مطالعه شبیه‌سازی به منظور بررسی اندازه و توان آزمون‌های مورد مطالعه انجام شده است. نتایج شبیه‌سازی حاکی از این است که آزمون معرفی شده در این مقاله بر آزمون کلاسیک خرد برتری دارد. مثالی کاربردی از روش پیشنهادی روی مجموعه داده سرطان پرستات در بخش ۴ بررسی می‌شود. در بخش ۵ به بحث و نتیجه‌گیری پرداخته می‌شود.

¹⁰Linear spectral statistics

۲ بخش نظری

فرض کنید بردار X_i دارای توزیع $N_p(\mu_i, \Sigma_i)$ است. اگر X_i به $1 \leq k_i \leq p$ زیربردار به صورت $X_i' = (X_1^{(i)'}, \dots, X_{k_i}^{(i)'})'$ افراز شود، آنگاه $X_r^{(i)} \in \mathbb{R}^{p_r^{(i)}}$ و $(p_1^{(i)}, \dots, p_{k_i}^{(i)})$ افزای از p با $p = \sum_{r=1}^{k_i} p_r^{(i)}$ است. متناظر با افراز X_i بردار میانگین μ_i و ماتریس کوواریانس Σ_i به ترتیب به صورت $\mu_i' = (\mu_1^{(i)'}, \dots, \mu_{k_i}^{(i)'})'$ و $\Sigma_i := (\Sigma_{\ell \times m}^{(i)})_{k_i \times k_i}$ افراز خواهند شد، که در آن $\Sigma_{\ell \times m}^{(i)} = \text{Cov}(X_\ell^{(i)}, X_m^{(i)})$ درایه ماتریسی روی سطر ℓ و ستون m ماتریس بلوکی حاصل از افراز Σ_i است. با این قراردادهای، فرض استقلال زیربردارهای $X_1^{(i)}, \dots, X_{k_i}^{(i)}$ یا همان H_0 فرضی است که در این مقاله آزمون می‌شود. این فرض را به شکل معادل می‌توان به صورت قابلیت نوشتن چگالی X_i به شکل حاصل ضرب چگالی‌های $X_1^{(i)}, \dots, X_{k_i}^{(i)}$ نیز بیان کرد. از این رو می‌توان فرض H_0 را به صورت

$$H_0 : f_{X_i}(\mathbf{x}_i; \mu_i, \Sigma_i) = \prod_{r=1}^{k_i} f_{X_r^{(i)}}(x_r^{(i)}; \mu_r^{(i)}, \Sigma_{rr}^{(i)}), \quad i = 1, \dots, k, \quad (1)$$

نیز نوشت. حال، فرض کنید به ازای هر $i = 1, \dots, k$ بردارهای x_{i1}, \dots, x_{in_i} نمونه‌ای تصادفی با اندازه n_i از X_i هستند. بنابراین تحت فرض صفر (۱) تابع درست‌نمایی این مشاهدات به صورت

$$L(\mu_1, \dots, \mu_k; \Sigma_1, \dots, \Sigma_k) = \prod_{i=1}^k \prod_{j=1}^{n_i} f_{X_{ij}}(x_{ij}; \mu_i, \Sigma_i),$$

است. می‌توان نشان داد (اندرسون، ۲۰۰۳) که

$$\sup_{\mu_i, \Sigma_i, i=1, \dots, k} L(\mu_1, \dots, \mu_k; \Sigma_1, \dots, \Sigma_k) = \prod_{i=1}^k (\sqrt{\pi} e n_i^{-1})^{-\frac{n_i p}{4}} |\mathbf{W}_i|^{-\frac{n_i}{4}},$$

به گونه‌ای که $\mathbf{W}_i = \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)' := (\mathbf{W}_{\ell \times m}^{(i)})_{k_i \times k_i}$ به سادگی می‌توان دید که تحت فرض H_0 در (۱)، ماتریس کوواریانس Σ_i برابر ماتریس قطری بلوکی $\Sigma_i^{H_0}$ به صورت

$$\Sigma_i^{H_0} = \begin{bmatrix} \Sigma_{11}^{(i)} & \circ & \dots & \circ \\ \circ & \Sigma_{22}^{(i)} & \dots & \circ \\ \vdots & \vdots & \ddots & \vdots \\ \circ & \circ & \dots & \Sigma_{k_i k_i}^{(i)} \end{bmatrix}.$$

خواهد بود. بنابراین، تحت فرض H_0 در (۱)،

$$\begin{aligned} & \sup_{\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{H_0}, i=1, \dots, k} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k; \boldsymbol{\Sigma}_1^{H_0}, \dots, \boldsymbol{\Sigma}_k^{H_0}) \\ &= \sup_{\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{H_0}, i=1, \dots, k} \prod_{i=1}^k \prod_{j=1}^{n_i} f_{\mathbf{X}_{ij}}(\mathbf{x}_{ij}, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{H_0}) \\ &= \prod_{i=1}^k \prod_{r=1}^{k_i} \sup_{\boldsymbol{\mu}_r^{(i)}, \boldsymbol{\Sigma}_{rr}^{(i)}, i=1, \dots, k} \prod_{j=1}^{n_i} f_{\mathbf{X}_{rj}^{(i)}}(\mathbf{x}_{rj}^{(i)}; \boldsymbol{\mu}_r^{(i)}, \boldsymbol{\Sigma}_{rr}^{(i)}) \\ &= \prod_{i=1}^k \prod_{r=1}^{k_i} \left(\sqrt{\pi} e n_i^{-1} \right)^{-\frac{n_i p_r^{(i)}}{\nu}} \left| \mathbf{W}_{rr}^{(i)} \right|^{-\frac{n_i}{\nu}}. \end{aligned}$$

در نتیجه، آماره LRT عبارت خواهد بود از:

$$\begin{aligned} W_n &= \frac{\sup_{\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{H_0}, i=1, \dots, k} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k; \boldsymbol{\Sigma}_1^{H_0}, \dots, \boldsymbol{\Sigma}_k^{H_0})}{\sup_{\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, i=1, \dots, k} L(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k; \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k)} \\ &= \frac{\prod_{i=1}^k \prod_{r=1}^{k_i} \left(\frac{\sqrt{\pi} e}{n_i} \right)^{-\frac{n_i p_r^{(i)}}{\nu}} \left| \mathbf{W}_{rr}^{(i)} \right|^{-\frac{n_i}{\nu}}}{\prod_{i=1}^k \left(\frac{\sqrt{\pi} e}{n_i} \right)^{-\frac{n_i p_i}{\nu}} \left| \mathbf{W}_i \right|^{-\frac{n_i}{\nu}}} \\ &= \prod_{i=1}^k \left(\frac{\left| \mathbf{W}_i \right|}{\prod_{r=1}^{k_i} \left| \mathbf{W}_{rr}^{(i)} \right|} \right)^{\frac{n_i}{\nu}}, \end{aligned} \quad (2)$$

که در آن $\mathbf{n} = (n_1, \dots, n_k)$. توجه شود که آماره W_n تنها در حالت $\min_{1 \leq i \leq k} n_i > p$ تعریف می‌شود. نظریه عام آزمون‌های نسبت درست‌نمایی (رائو، ۲۰۰۹) بیان می‌کند که وقتی اندازه‌های نمونه‌ای n_i ، $i = 1, \dots, k$ از بُعد p بزرگ‌اند و نسبت‌های $\frac{p}{n_i}$ اعدادی نزدیک به صفر $(\frac{p}{n_i} \rightarrow 0)$ هستند، توزیع تحت فرض صفر آماره $-2 \log W_n$ به توزیع χ^2 با $df = \frac{1}{\nu} (kp^2 - \sum_{i=1}^k \sum_{r=1}^{k_i} p_r^{(i)})$ درجه آزادی همگراست.

به هر حال، برای داده‌های با بُعد (نسبتاً) بالا، تقریب χ^2 عملاً غیر قابل استفاده است. به منظور غلبه بر این مشکل، در قالب قضیه زیر تقریب بهتری از تقریب χ^2 برای توزیع تحت فرض صفر آماره LRT ارائه می‌شود. در ادامه، به منظور تسهیل نمایش و کار با روابط حدی فرض شده است که n_i به ازای هر $i = 1, \dots, k$ به صورت $n_i = n_i(p)$ وابسته است.

قضیه ۱. فرض کنید به ازای هر $i = 1, \dots, k$ ، رابطه

$$n_i = n_i(p) > 1 + p = 1 + \sum_{j=1}^{k_i} p_j^{(i)}$$

برقرار باشد، که در آن $(p_1^{(i)}, \dots, p_{k_i}^{(i)})$ افزایی از p است و به ازای هر $r = 1, \dots, k_i$

$$\lim_{p \rightarrow \infty} \frac{p_r^{(i)}}{n_i} = y_r^{(i)} \in (0, 1).$$

آنگاه تحت فرض H_0 در (۱)، وقتی $p \rightarrow \infty$ آماره $\frac{\log W_n - \mu_n}{n\sigma_n}$ در توزیع به $N(0, 1)$ همگراست، به گونه‌ای که روابط

$$\begin{aligned} \mu_n &= \frac{1}{\sqrt{p}} \sum_{i=1}^k \left[\left(\sum_{r=1}^{k_i} r_{n_i-1, p_r^{(i)}}^{\chi} (p_r^{(i)} - n_i + 1/5) - r_{n_i-1, p}^{\chi} (p - n_i + 1/5) \right) n_i \right] \\ \sigma_n^{\chi} &= \frac{1}{\sqrt{p}} \sum_{i=1}^k \left[\left(r_{n_i-1, p}^{\chi} - \sum_{r=1}^{k_i} r_{n_i-1, p_r^{(i)}}^{\chi} \left(\frac{n_i}{n} \right)^{\chi} \right) \right] \end{aligned}$$

برقرارند، که در آنها $n = \sum_{i=1}^k n_i$ و $r_{x,y} = \sqrt{-\log(1 - \frac{y}{x})}$

برهان. هاردی و همکاران (۱۹۸۸) برای اعداد حقیقی a_1, \dots, a_q بزرگتر از -1 ، که همگی یا

$$\text{مثبت‌اند یا منفی، ثابت کردند } 1 + \sum_{i=1}^q a_i > \prod_{i=1}^q (1 + a_i) \text{ یا}$$

$$\log \left(\prod_{i=1}^q (1 + a_i) \right) - \log \left(1 + \sum_{i=1}^q a_i \right) > 0.$$

با تثبیت $i \in \{1, \dots, k\}$ و تعریف $a_i = -\frac{p_r^{(i)}}{n_i - 1}$ و $q = k_i$ ، مشاهده می‌شود

$$r_{n_i-1, p}^{\chi} - \sum_{r=1}^{k_i} r_{n_i-1, p_r^{(i)}}^{\chi} = \log \left(\prod_{r=1}^{k_i} \left(1 - \frac{p_r^{(i)}}{n_i - 1} \right) \right) - \log \left(1 - \sum_{r=1}^{k_i} \frac{p_r^{(i)}}{n_i - 1} \right) > 0.$$

در نتیجه، $\sigma_n^{\chi} > 0$ از این وقتی $p \rightarrow \infty$:

$$\frac{n}{p} = \sum_{i=1}^k \frac{n_i}{p} = \sum_{i=1}^k \left(\sum_{r=1}^{k_i} \frac{p_r^{(i)}}{n_i} \right)^{-1} \rightarrow \sum_{i=1}^k \left(\sum_{r=1}^{k_i} y_r^{(i)} \right)^{-1} = \sum_{i=1}^k \frac{1}{y_i} := \frac{1}{y},$$

که در آن $\frac{n_i}{n} = \frac{(\frac{n_i}{p})}{(\frac{n}{p})} \rightarrow \frac{y}{y_i} \in (0, 1]$ در نتیجه، $y \in (0, \frac{1}{k}]$ و $y_i = \sum_{r=1}^{k_i} y_r^{(i)} \in (0, 1]$ وقتی

$p \rightarrow \infty$ بنا براین، $\sigma_n^{\chi} = \lim_{p \rightarrow \infty} \sigma_n^{\chi}$ به ازای $\max_{1 \leq i \leq k} y_i < 1$ برابر

$$\frac{1}{\sqrt{p}} \sum_{i=1}^k \left[\left(\sum_{r=1}^{k_i} \log(1 - y_r^{(i)}) - \log(1 - y_i) \right) \left(\frac{y}{y_i} \right)^{\chi} \right]$$

و به ازای $\max_{1 \leq i \leq k} y_i = 1$ برابر $+\infty$ خواهد بود. در حقیقت برای حالت دوم، از اینکه $y_r^{(i)} \in (0, 1)$ و $\lim_{x \rightarrow 1^-} \log(1-x) = -\infty$ واضح است که حد برابر $+\infty$ خواهد شد. حال با تثبیت s با شرط $|s| < \frac{\sigma}{\sqrt{y}}$ ، t به شکل $t = t_n = \frac{s}{n\sigma_n}$ تعریف می‌شود. واضح است که $-\frac{\sigma}{\sqrt{y}} < -\frac{n\sigma_n}{\sqrt{p+1}} \rightarrow -\frac{\sigma}{\sqrt{y}} < s$ وقتی $p \rightarrow \infty$. این موضوع نتیجه می‌دهد که برای p به اندازه کافی بزرگ $s < -\frac{n\sigma_n}{\sqrt{p+1}}$ یا $-\frac{1}{\sqrt{p+1}} < \frac{s}{n\sigma_n}$ این نابرابری در کنار این واقعیت که

$$\max_{1 \leq i \leq k} \left\{ \frac{1}{n_i} \right\} < \frac{1}{p+1} \quad \text{یا} \quad \max_{1 \leq i \leq k} \left\{ \frac{p}{n_i} - 1 \right\} < -\frac{1}{p+1} < -\frac{1}{\sqrt{p+1}},$$

نتیجه می‌دهد که $t = t_n = \frac{s}{n\sigma_n} > \max_{1 \leq i \leq k} \left\{ \frac{p}{n_i} - 1 \right\}$ حال برای یک مقدار ثابت i و حالتی که $y_i < 1$ ، $\frac{r_{n_i-1,p}^{\sqrt{y_i}}}{\sigma_n^{\sqrt{y_i}}}$ برای $\max_{1 \leq i \leq k} y_i < 1$ به مقدار $\frac{-\log(1-y_i)}{\sigma^{\sqrt{y_i}}}$ و به ازای $\max_{1 \leq i \leq k} y_i = 1$ به مقدار 0 همگرا است. این نتیجه می‌دهد که $\frac{r_{n_i-1,p}^{\sqrt{y_i}}}{\sigma_n^{\sqrt{y_i}}}$ به ازای $\max_{1 \leq i \leq k} y_i < 1$ و به ازای $\max_{1 \leq i \leq k} y_i = 1$ همگرا خواهد بود. علاوه بر این، برای حالت $y_i = 1$ ،

$$\sqrt{2}\sigma_n^{\sqrt{y_i}} \geq (r_{n_i-1,p}^{\sqrt{y_i}} - \sum_{r=1}^{k_i} r_{n_i-1,p_r^{(i)}}^{\sqrt{y_i}}) \left(\frac{n_i}{n} \right)^{\sqrt{y_i}},$$

یا به طور معادل

$$\frac{r_{n_i-1,p}^{\sqrt{y_i}}}{\sigma_n^{\sqrt{y_i}}} \leq \sqrt{2 \left(\frac{n_i}{n} \right)^{\sqrt{y_i}} + \sigma_n^{-2} \sum_{r=1}^{k_i} r_{n_i-1,p_r^{(i)}}^{\sqrt{y_i}}} \rightarrow \frac{\sqrt{2}}{y} < \infty,$$

وقتی $p \rightarrow \infty$ بنابراین،

$$\limsup_{p \rightarrow \infty} \frac{r_{n_i-1,p}^{\sqrt{y_i}}}{\sigma_n^{\sqrt{y_i}}} < \infty \Rightarrow \frac{1}{\sigma_n} = O\left(\frac{1}{r_{n_i-1,p}^{\sqrt{y_i}}}\right),$$

یا

$$\frac{tn_i}{\sqrt{2}} = \frac{s}{\sqrt{2}} \frac{n_i}{n} \frac{1}{\sigma_n} = O(1) \frac{1}{\sigma_n} = O\left(\frac{1}{r_{n_i-1,p}^{\sqrt{y_i}}}\right).$$

به طور مشابه، از آنجا که $-\log(1-x)$ برای $x < 1$ یک تابع صعودی است،

$$p_r^{(i)} < p \Leftrightarrow r_{n_i-1,p_r^{(i)}}^{\sqrt{y_i}} < r_{n_i-1,p}^{\sqrt{y_i}}.$$

در نتیجه،

$$\limsup_{p \rightarrow \infty} \frac{r_{n_i-1,p_r^{(i)}}^{\sqrt{y_i}}}{\sigma_n^{\sqrt{y_i}}} < \limsup_{p \rightarrow \infty} \frac{r_{n_i-1,p}^{\sqrt{y_i}}}{\sigma_n^{\sqrt{y_i}}} < \infty,$$

یا

$$\frac{1}{\sigma_n} = O\left(\frac{1}{r_{n_i-1, p_r^{(i)}}}\right) \Rightarrow \frac{tn_i}{\gamma} = O\left(\frac{1}{r_{n_i-1, p}}\right).$$

حال بنا بر قضیه ۱۱.۲.۳ در مؤیرهد (۱۹۸۲)، وقتی H_0 صحیح است، t امین گشتاور W_n برابر

$$\begin{aligned} E[W_n^t] &= \prod_{i=1}^k E\left[\left(|W_i| \prod_{r=1}^{k_i} |W_{rr}^{(i)}|^{-1}\right)^{\frac{tn_i}{\gamma}}\right] \\ &= \prod_{i=1}^k \left(\frac{\Gamma_{p_i}\left(\frac{n_i-1}{\gamma} + \frac{tn_i}{\gamma}\right)}{\Gamma_{p_i}\left(\frac{n_i-1}{\gamma}\right)} \prod_{r=1}^{k_i} \frac{\Gamma_{p_r^{(i)}}\left(\frac{n_i-1}{\gamma}\right)}{\Gamma_{p_r^{(i)}}\left(\frac{n_i-1}{\gamma} + \frac{tn_i}{\gamma}\right)}\right), \end{aligned} \quad (۳)$$

خواهد بود، که در آن برای عدد مختلط z با ویژگی $Re(z) > \frac{1}{\gamma}(p-1)$ ، تابع گامای چند متغیره^{۱۱} $\Gamma_p(z) = \pi^{\frac{p(p-1)}{\gamma}} \prod_{j=1}^p \Gamma\left(z - \frac{1}{\gamma}(j-1)\right)$ (مؤیرهد، ۱۹۸۲، صفحه ۶۲) به صورت $\Gamma_p(z)$ تعریف می‌شود. بدیهی است که امید ریاضی (۳) تنها زمانی وجود دارد که $\frac{n_i-1}{\gamma} + \frac{tn_i}{\gamma} > \frac{p-1}{\gamma}$ یا $i \in \{1, 2, \dots, k\}$ با $t > \max_{1 \leq i \leq k} \left\{ \frac{p}{n_i} - 1 \right\}$ با توجه به اینکه

$$t = t_n = \frac{s}{n\sigma_n} > \max_{1 \leq i \leq k} \left\{ \frac{p}{n_i} - 1 \right\}, \quad \frac{tn_i}{\gamma} = O\left(\frac{1}{r_{n_i-1, p}}\right),$$

و $\frac{tn_i}{\gamma} = O\left(\frac{1}{r_{n_i-1, p_r^{(i)}}}\right)$ وقتی که $p \rightarrow \infty$ ، با استفاده از لم ۵.۴ ژیانگ و یانگ (۲۰۱۳)، به برابری‌های

$$\begin{aligned} \log \frac{\Gamma_p\left(\frac{n_i-1}{\gamma} + \frac{tn_i}{\gamma}\right)}{\Gamma_p\left(\frac{n_i-1}{\gamma}\right)} &= \frac{tpn_i}{\gamma} \log\left(\frac{n_i-1}{\gamma e}\right) + r_{n_i-1, p}^{\gamma} \frac{n_i^{\gamma} t^{\gamma}}{\gamma} \\ &\quad - r_{n_i-1, p}^{\gamma} (p - n_i + 1) \frac{tn_i}{\gamma} + o(1), \end{aligned}$$

و

$$\begin{aligned} \log \frac{\Gamma_{p_r^{(i)}}\left(\frac{n_i-1}{\gamma} + \frac{tn_i}{\gamma}\right)}{\Gamma_{p_r^{(i)}}\left(\frac{n_i-1}{\gamma}\right)} &= \frac{tp_r^{(i)} n_i}{\gamma} \log\left(\frac{n_i-1}{\gamma e}\right) + r_{n_i-1, p_r^{(i)}}^{\gamma} \frac{n_i^{\gamma} t^{\gamma}}{\gamma} \\ &\quad - r_{n_i-1, p_r^{(i)}}^{\gamma} (p_r^{(i)} - n_i + 1) \frac{tn_i}{\gamma} + o(1), \end{aligned}$$

¹¹Multivariate gamma function

حاصل می‌شود. بنابراین،

$$\begin{aligned} \log E[W_n^t] &= \sum_{i=1}^k \left[\log \frac{\Gamma_p\left(\frac{n_i-1}{\nu} + \frac{tn_i}{\nu}\right)}{\Gamma_p\left(\frac{n_i-1}{\nu}\right)} - \sum_{r=1}^{k_i} \log \frac{\Gamma_{p_r^{(i)}}\left(\frac{n_i-1}{\nu} + \frac{tn_i}{\nu}\right)}{\Gamma_{p_r^{(i)}}\left(\frac{n_i-1}{\nu}\right)} \right] \\ &= \sum_{i=1}^k \left[\frac{tpn_i}{\nu} \log\left(\frac{n_i-1}{\nu e}\right) + r_{n_i-1,p}^\nu \frac{n_i t^\nu}{\nu} \right] \\ &\quad - \sum_{i=1}^k \left[r_{n_i-1,p}^\nu (p - n_i + 1/\omega) \frac{tn_i}{\nu} + o(1) \right] \\ &\quad - \sum_{i=1}^k \sum_{r=1}^{k_i} \left[\frac{tp_r^{(i)} n_i}{\nu} \log\left(\frac{n_i-1}{\nu e}\right) + r_{n_i-1,p_r^{(i)}}^\nu \frac{n_i t^\nu}{\nu} \right] \\ &\quad + \sum_{i=1}^k \sum_{r=1}^{k_i} \left[r_{n_i-1,p_r^{(i)}}^\nu (p_r^{(i)} - n_i + 1/\omega) \frac{tn_i}{\nu} + o(1) \right] \\ &= \frac{1}{\nu} n^\nu \sigma_n^\nu t^\nu + \mu_n t + o(1), \end{aligned}$$

وقتی $p \rightarrow \infty$ از آنجا که $t = t_n = \frac{s}{n\sigma_n}$

$$\log E\left[e^{\frac{\log W_n}{n\sigma_n} s}\right] = \log e^{\frac{1}{\nu} s^\nu + \frac{\mu_n}{n\sigma_n} s + o(1)},$$

یا به طور معادل

$$E\left[e^{\frac{\log W_n - \mu_n}{n\sigma_n} s}\right] = e^{\frac{1}{\nu} s^\nu + o(1)},$$

وقتی $p \rightarrow \infty$ یعنی، وقتی $p \rightarrow \infty$ در توزیع به $N(0, 1)$ همگراست.

به طور خلاصه، قضیه فوق بیان می‌کند که در حالت داده‌های با بُعد نسبتاً بالا، توزیع تقریبی مناسب برای لگاریتم آماره نسبت درست‌نمایی $-2 \log W_n$ توزیع نرمال با میانگین $-2\mu_n$ و واریانس $4n^2\sigma_n^2$ است. در بخش بعد، به مقایسه این تقریب با تقریب کلاسیک خی‌دو برای $-2 \log W_n$ پرداخته می‌شود.

۳ مطالعه شبیه‌سازی

در این بخش به منظور مقایسه اندازه و توان آزمون‌های LRT برای فرض (۱) بر اساس توزیع‌های تقریبی خی‌دو (χ^2) و نرمال (N) یک مطالعه شبیه‌سازی انجام شده است. در این شبیه‌سازی‌ها،

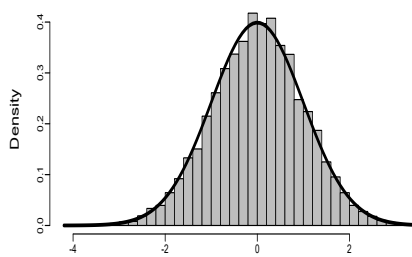
ماتریس‌های $p \times p$ J_p و I_p به ترتیب به عنوان ماتریسی با تمامی درایه‌های برابر با عدد یک و ماتریس همانی تعریف شده است. همچنین فرض می‌شود که ماتریس‌های کوواریانس Σ_i ها به ازای ثابت $0 \leq \rho < 1$ ، همگی برابر $\rho J_p + (1 - \rho)I_p$ هستند. در حقیقت در اینجا توان آزمون به عنوان تابعی از ρ بررسی شده است که در آن به ازای $\rho = 0$ توزیع تحت فرض صفر و با دور شدن ρ از عدد ۰ توزیع‌های تحت فرض مقابل بدست می‌آیند. علاوه بر این، بردارهای میانگین μ_i ها همگی برابر بردار صفر، سطح معنی‌داری اسمی α برابر 0.05 و $k = 3$ است. همچنین، در هر یک از این k جامعه نرمال، به منظور آزمون همزمان استقلال زیربردارها به تفکیک هر یک از این جوامع، افزایش‌های p یا همان ابعاد زیربردارها در تمامی k گروه یکسان در نظر گرفته شده است. تمامی برنامه‌های کامپیوتری در نرم‌افزار آماری R نوشته شده و در صورت درخواست، از طرف نویسنده در اختیار خواننده قرار می‌گیرد.

حال برای هر ترکیب از پارامترهای شبیه‌سازی ρ ، n_1 ، n_2 ، n_3 و p که در جدول ۱ آمده‌اند، با استفاده از ۱۰۰۰۰ نمونه از توزیع $(N_p(0, \rho J_p + (1 - \rho)I_p))$ ، مقادیر شبیه‌سازی شده اندازه آزمون $\hat{\varphi}_M(0)$ و $\hat{\varphi}_M(\rho)$ و توان آزمون $\hat{\varphi}_M(\rho)$ ، برای $\rho > 0$ تقریب $\{\chi^2, N\}$ $M \in \{\chi^2, N\}$ محاسبه و در جدول ۱ آورده شده است. برای نمونه‌های شبیه‌سازی شده تحت فرض صفر ($\rho = 0$) و نمونه‌های شبیه‌سازی شده تحت فرض مقابل ($\rho > 0$) به ترتیب مقادیر $\hat{\varphi}_M(0)$ و $\hat{\varphi}_M(\rho)$ برای $M \in \{\chi^2, N\}$ برابر نسبت تعداد دفعاتی که در این ۱۰۰۰۰ بار فرض صفر رد شده، محاسبه شده است. همانطور که در جدول ۱، ملاحظه می‌شود تقریب نرمال معرفی شده و تقریب کلاسیک χ^2 دو برای اندازه‌های نمونه‌ای n_i بزرگ و بُعد p کوچک رفتاری شبیه یکدیگر دارند. چنین وضعیتی را می‌توان در جدول ۱ به ازای $(150, 150, 150) = n = 6v$ ، $p = 5$ و $\rho = 0, 0.05, 0.6$ ملاحظه کرد. تحت فرض H_0 با افزایش p و نزدیک شدن آن به مقادیر اندازه‌های نمونه‌ای n_i (حالت داده‌های با بُعد نسبتاً بالا)، اندازه آزمون با تقریب χ^2 ؛ یعنی، $\hat{\varphi}_{\chi^2}(0)$ به جای اینکه به مقدار اسمی $\alpha = 0.05$ نزدیک باشد، مقادیری نزدیک به عدد یک قبول می‌کند (!) در حالی که اندازه آزمون روش معرفی شده؛ یعنی، $\hat{\varphi}_N(0)$ ، در حدود مقدار اسمی 0.05 است (جدول ۱ را به ازای $\rho = 0$ و مقادیر بزرگ p مشاهده کنید). بنابراین تقریب χ^2 دو در مقایسه با تقریب نرمال برای بُعد نسبتاً بالا عملاً غیر قابل استفاده است. نگاهی بر جدول ۱ نشان می‌دهد که به ازای انحرافات کوچک از فرض صفر به مانند $\rho = 0.05$ ، توان آزمون

جدول ۱: اندازه آزمون $\hat{\varphi}_M(\circ)$ و توان آزمون $\hat{\varphi}_M(\rho)$ برای تقریب $M \in \{\chi^2, N\}$

$n = 4\nu$				$n = \nu$				ρ
$\hat{\varphi}_N(\rho)$	$\hat{\varphi}_{\chi^2}(\rho)$	افراز p	p	$\hat{\varphi}_N(\rho)$	$\hat{\varphi}_{\chi^2}(\rho)$	افراز p	p	
۰/۰۴۶	۰/۰۶۶	۳, ۲	۵	۰/۰۴۸	۰/۱۵۴	۲, ۱, ۲	۵	
۰/۰۵	۱	۱۵, ۱۷, ۱۸	۵۰	۰/۰۵۲	۰/۶۶۱	۴, ۳, ۳	۱۰	۰
۰/۰۵۵	۱	۱۵, ۳۰, ۳۰, ۲۰	۹۵	۰/۰۶۲	۱	۴, ۵, ۶, ۵	۲۰	
۰/۱۳۸	۰/۱۹۳	۳, ۲	۵	۰/۰۶۸	۰/۱۹۶	۲, ۱, ۲	۵	
۰/۴۷۴	۱	۱۵, ۱۷, ۱۸	۵۰	۰/۰۷۰	۰/۷۳۹	۴, ۳, ۳	۱۰	۰/۰۵
۰/۲۲۳	۱	۱۵, ۳۰, ۳۰, ۲۰	۹۵	۰/۰۶۸	۱	۴, ۵, ۶, ۵	۲۰	
۱	۱	۳, ۲	۵	۱	۱	۲, ۱, ۲	۵	
۱	۱	۱۷, ۱۸, ۱۵	۵۰	۱	۱	۴, ۳, ۳	۱۰	۰/۶
۱	۱	۱۵, ۳۰, ۳۰, ۲۰	۹۵	۱	۱	۴, ۵, ۶, ۵	۲۰	
$n = 6\nu$				$n = 2\nu$				
$\hat{\varphi}_N(\rho)$	$\hat{\varphi}_{\chi^2}(\rho)$	افراز p	p	$\hat{\varphi}_N(\rho)$	$\hat{\varphi}_{\chi^2}(\rho)$	افراز p	p	
۰/۰۴۲	۰/۰۵۷	۱, ۴	۵	۰/۰۴۷	۰/۰۸۹	۲, ۳	۵	
۰/۰۵۴	۱	۱۲, ۲۴, ۱۳, ۲۱, ۱۰	۸۰	۰/۰۵۱	۰/۹۶۳	۵, ۷, ۸	۲۰	۰
۰/۰۵۴	۱	۵۰, ۷۰, ۲۵	۱۴۵	۰/۰۶۱	۱	۱۷, ۱۵, ۱۳	۴۵	
۰/۱۷۴	۰/۲۱۴	۱, ۴	۵	۰/۰۷۹	۰/۱۴۹	۲, ۳	۵	
۰/۹۴۷	۱	۱۲, ۲۴, ۱۳, ۲۱, ۱۰	۸۰	۰/۱۲۷	۰/۹۹۴	۵, ۷, ۸	۲۰	۰/۰۵
۰/۲۲۴	۱	۵۰, ۷۰, ۲۵	۱۴۵	۰/۰۹۲	۱	۱۷, ۱۵, ۱۳	۴۵	
۱	۱	۱, ۴	۵	۱	۱	۲, ۳	۵	
۱	۱	۱۲, ۲۴, ۱۳, ۲۱, ۱۰	۸۰	۱	۱	۵, ۷, ۸	۲۰	۰/۶
۱	۱	۵۰, ۷۰, ۲۵	۱۴۵	۱	۱	۱۷, ۱۵, ۱۳	۴۵	

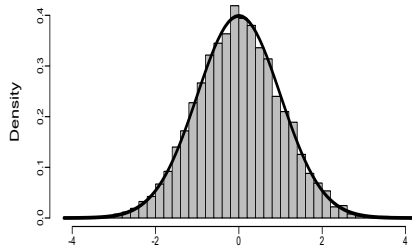
با تقریب خردی کلاسیک به شکل کاذب بسیار بالاتر از توان آزمون پیشنهادی قرار می‌گیرد، که این به دلیل متورم یا بالا بودن خطای نوع اول آزمون با تقریب خردی دو است. به هر حال با افزایش ρ به مقدار $\rho = 0.6$ و در نتیجه دور شدن بیشتر از فرض صفر استقلال، توان دو آزمون رفته رفته بر یکدیگر منطبق و به عدد یک همگرا می‌شوند. نتایج حاصل از جدول ۱ برای نمونه‌های تحت فرض H_0 را می‌توان در شکل‌های ۱ و ۲ نیز مشاهده کرد.



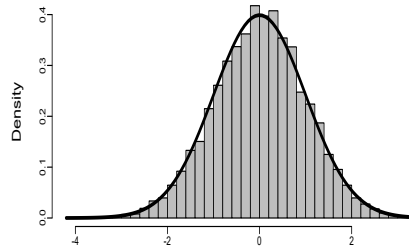
(ب)



(الف)



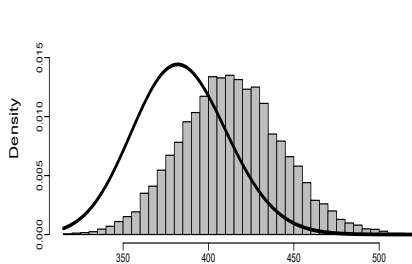
(د)



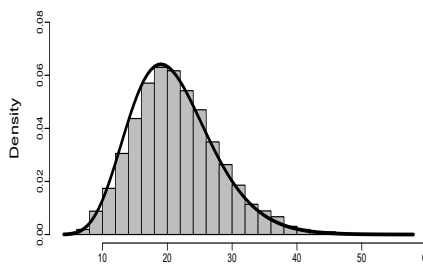
(ج)

شکل ۱: بافت‌نگار مقادیر شبیه‌سازی شده و چگالی $N(0, 1)$ به ازای $k = 3$ ، الف- بُعد ۵ و افراز (۱, ۳)؛ ب- بُعد ۲ و افراز (۸, ۸, ۴)؛ ج- بُعد ۴ و افراز (۱۰, ۱۰, ۲۰)؛ د- بُعد ۱۰۰ و افراز (۳۰, ۳۰, ۴۰).

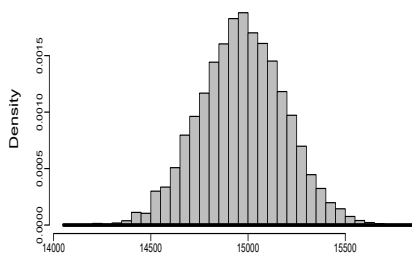
در این شکل‌ها، به ازای مقادیر مختلف از پارامترهای شبیه‌سازی، نمودارهای بافت‌نگار بدست آمده از ۱۰۰۰۰ مقدار شبیه‌سازی شده از $\frac{\log W_n - \mu_n}{n\sigma_n}$ و $-2 \log W_n$ به ترتیب به همراه منحنی‌های متناظر با چگالی‌های $N(0, 1)$ رسم شده است. شکل ۱ نشان می‌دهد که بافت‌نگار شبیه‌سازی شده



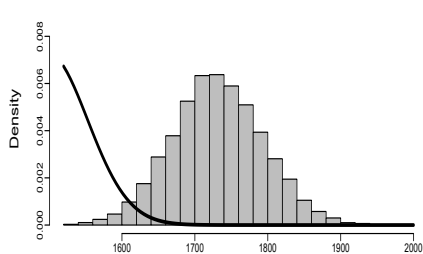
(ب)



(الف)



(د)



(ج)

شکل ۲: بافت‌نگار مقادیر شبیه‌سازی‌شده و چگالی χ^2_{df} به ازای $k = 3$ ، الف- بُعد ۵ و افراز $(1, 1, 3)$ ؛ ب- بُعد ۲۰ و افراز $(8, 8, 4)$ ؛ ج- بُعد ۴۰ و افراز $(10, 10, 20)$ ؛ د- بُعد ۱۰۰ و افراز $(30, 30, 40)$.

$\frac{\log W_n - \mu_n}{n\sigma_n}$ و تابع چگالی $N(0, 1)$ با افزایش p بر یکدیگر منطبق می‌شوند؛ این در حالی است که در شکل ۲ بافت‌نگار شبیه‌سازی‌شده $-2 \log W_n$ با افزایش p به مرور از چگالی χ^2_{df} دور می‌شود. به عبارت دیگر، در حالتی که p متناسب با n_i ‌ها افزایش می‌یابد یا نسبت‌های $\frac{p}{n_i}$ ‌ها به یک نزدیک می‌شوند، تقریب χ^2 در مقایسه با تقریب نرمال بدتر می‌شود. جالب این است که حتی به ازای بُعد کوچک نیز تقریب نرمال عملکرد خوبی دارد.

۴ تحلیل داده‌های سرطان پروستات

در این بخش کاربردی از روش پیشنهادی بر مجموعه داده‌های سرطان پروستات یا “سطح بیان ژن تومور پروستات”^{۱۲} (سینگ و همکاران، ۲۰۰۲؛ دتلینگ و بوهمن، ۲۰۰۲) ارائه شده است. این مجموعه از داده‌ها که در بسته spls^{۱۳} از نرم‌افزار R نیز قابل دسترس است، شامل $n_1 = 52$ نمونه از بافت‌های تومور پروستات و $n_2 = 50$ نمونه از بافت‌های سالم است. بر روی این دو جامعه از بافت‌ها، برای هر نمونه سطح بیان 6033 ژن به کمک تکنولوژی آفی‌متریکس^{۱۴} اندازه‌گیری و ثبت شده است. سینگ و همکاران (۲۰۰۲) بر اساس همبستگی‌های بین مقادیر سطوح بیان و همچنین همبستگی بین این سطوح با متغیرهای بالینی و آسیب‌شناختی (به مانند سن بیمار، نتیجه آزمون‌های تشخیصی PSA^{۱۵} و GS^{۱۶}، غیره) به تشریح رفتار بالینی سرطان پروستات پرداختند.

به دلیل اهمیت همبستگی‌های درونی موجود در هر یک از این دو جامعه، در سطح آزمون $\alpha = 0.05$ استقلال همزمان زیربردارهایی از $p = 48$ سطح بیان ژن نخست از بین 6033 سطح بیان مطالعه می‌شود. در این مطالعه، به عنوان یک حالت خاص از فرض (۱)، استقلال تمامی 48 سطح بیان ژن یا همان آزمون گرویت ژیانگ و یانگ (۲۰۱۳) به طور همزمان در هر دو جامعه بافت‌های سالم و بافت‌های سرطانی بررسی شده است. توجه شود که این حالت (به بخش ۲ مراجعه شود) افزاز مشترکی از $p = 48$ سطح بیان در دو جامعه به صورت $p = (p_1, \dots, p_{48}) = (1, \dots, 1)$ را نتیجه می‌دهد. همانگونه که مشاهده می‌شود در این حالت، نسبت‌های $\frac{p}{n_1} = \frac{48}{52} = 0.923$ و $\frac{p}{n_2} = \frac{48}{50} = 0.96$ حاصل می‌شوند که اعدادی نزدیک به یک‌اند (یعنی داده‌های با بُعد نسبتاً بالا) و در نتیجه روش LRT با تقریب خوبی دو به منظور آزمون همزمان استقلال زیربردارهای در هر دو جامعه بافت‌های سالم و بافت‌های سرطانی عملاً غیر قابل استفاده است. نتایج حاصل از آزمون پیشنهادی و روش کلاسیک LRT به شرح زیر است. مقادیر آماره‌های آزمون LRT با تقریب خوبی و تقریب پیشنهادی نرمال به ترتیب برابر $11354.41 = -2 \log(W_n)$ و $-44.692 = \frac{\log W_n - \mu_n}{n\sigma_n}$ ($\sigma_n = 0.593$ و $\mu_n = -2165.82$) بدست می‌آیند که به ترتیب در مقایسه با مقادیر بحرانی

¹²Prostate tumor gene expression level dataset

¹³<https://CRAN.R-project.org/package=spls>

¹⁴Affymetrix technology

¹⁵Prostate-specific antigen test

¹⁶Gleason score

در جامعه بافت‌های سالم و هم در جامعه بافت‌های سرطانی را نتیجه می‌دهد. همانطور که ملاحظه می‌شود، روش کلاسیک LRT در مقایسه با روش پیشنهادی این مقاله با شدت بالاتری فرض صفر را رد می‌کند که این به دلیل متورم بودن خطای نوع اول این آزمون برای داده‌هایی نسبتاً بالا است؛ که می‌توانست حتی در صورت مستقل بودن مؤلفه‌ها نیز رأی به عدم استقلال مؤلفه‌ها دهد! به عبارتی دیگر، این مثال بیان می‌کند که در حالت داده‌های با بُعد نسبتاً بالا (حالت‌هایی که نسبت بُعد به اندازه نمونه کمتر از یک و نزدیک یک‌اند) باید در استفاده از روش LRT با تقریب کلاسیک خی‌دو تجدید نظر نمود.

۵ بحث و نتیجه‌گیری

در این مقاله، آزمون همزمان استقلال زیربردارهای k بردار نرمال p -متغیره با بُعد نسبتاً بالا مورد مطالعه قرار گرفت. در حالتی که n_i ها از بُعد p بزرگ‌اند و نسبت $\frac{p}{n_i}$ عددی کوچک و نزدیک به صفر است، این آزمون را می‌توان به کمک روش نسبت‌درست‌نمایی با تقریب قابل قبول خی‌دو انجام داد. اما برای داده‌های با بُعد نسبتاً بالا؛ یعنی، حالتی که در آن n_i ها از بُعد p بزرگ‌اند و نسبت‌های $\frac{p}{n_i}$ اعدادی نزدیک به یک هستند، روش آزمون کلاسیک LRT با تقریب خی‌دو به جای اینکه به مقدار اسمی $\alpha = 0.05$ نزدیک باشد، مقادیری نزدیک به عدد یک قبول می‌کند که این خود به معنای غیر قابل استفاده بودن این آزمون برای داده‌های با بُعد نسبتاً بالا است. برای این حالت، در این مقاله نشان داده شد که تقریب مناسب‌تر برای توزیع تحت فرض صفر آماره LRT به جای تقریب خی‌دو تقریب نرمال است. نتایج شبیه‌سازی حاکی از این بود که برای داده‌های با بُعد نسبتاً بالا تقریب خی‌دو کارایی نداشته و تقریب نرمال جایگزین مناسبتری برای آن است. از این رو، به منظور بررسی همزمان استقلال زیربردارهای چندین بردار نرمال p -متغیره، که استقلال همزمان مؤلفه‌های این بردارها حالت خاصی از آن است، برای داده‌های با بُعد نسبتاً بالا، استفاده از روش پیشنهادی این مقاله توصیه می‌شود. توجه شود که برای داده‌های با بُعد بالا که در آن اندازه‌های نمونه‌ای از بُعد کوچکتراند، روش این مقاله کارساز نیست و نویسنده در حال تحقیق بر روی این حالت است.

تقدیر و تشکر

نویسنده مقاله از دو داور محترم، سردبیر و ویراستار مجله به دلیل تلاش ایشان در راستای ارتقای کیفی مقاله و برطرف نمودن ایرادات احتمالی و ارائه بهتر مقاله کمال تشکر و قدردانی را دارد.

مراجع

- Anderson, T. W. (2003), *An Introduction to Multivariate Statistical Analysis (3rd ed.)*. John Wiley & Sons, New York.
- Bai, Z., D., Jiang, J., Yao, F. and Zheng. S. (2009), Corrections to LRT on Large-Dimensional Covariance Matrix by RMT, *The Annals of Statistics*, **37**, 3822-3840.
- Chen, X. and Liu, W. (2018), Testing Independence with High-Dimensional Correlated Samples, *The Annals of Statistics*, **46**, 866-894.
- Chen, S. X., Zhang, L. X. and Zhong, P. S. (2010), Tests for High-Dimensional Covariance Matrices, *Journal of the American Statistical Association*, **105**, 810-819.
- Detting, M. and Bühlmann, P. (2002), Supervised Clustering of Genes, *Genome biology*, **3**, Research0069-1.
- Hardy, G., Littlewood, J., and Pólya, G. (1988), *Inequalities*, Reprint of the 1952 Edition, Cambridge Mathematical Library.
- 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 Yang, F. (2013), Central Limit Theorems for Classical Likelihood Ratio

- Tests for High-Dimensional Normal Distributions, *The Annals of Statistics*, **41**, 2029-2074.
- Ledoit, O. and Wolf, M. (2002), Some Hypothesis Tests for the Covariance Matrix when the Dimension is Large Compared to the Sample Size, *The Annals of Statistics*, **30**, 1081-1102.
- Leung, D., and Drton, M. (2018), Testing Independence in High-Dimensions with Sums of Rank Correlations, *The Annals of Statistics*, **46**, 280-307.
- Mao, G. (2018), Testing Independence in High-Dimensions using Kendall's tau, *Computational Statistics & Data Analysis*, **117**, 128-137.
- Muirhead, R. J. (1982), *Aspects of Multivariate Statistical Theory*. Wiley Series in Probability and Mathematical Statistics, John Wiley & Sons, Inc., New York.
- Rao, C. (2009), *Linear Statistical Inference and its Applications*, Wiley Series in Probability and Statistics. Wiley.
- Schott, J. R. (2005), Testing for Complete Independence in High-Dimensions, *Biometrika*, **92**, 951-956.
- Schott, J. R. (2007), A Test for the Equality of Covariance Matrices when the Dimension is Large Relative to the Sample Sizes, *Computational Statistics & Data Analysis*, **51**, 6535-6542.
- 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.
- Wilks, S. (1938), The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses, *The Annals of Mathematical Statistics*, **9**, 60-62.

Simultaneous Test for Independence Among Subvectors of Several Moderately High Dimensional Multivariate Normal Distributions

Najarzadeh, D.

Department of Statistics, University of Tabriz, Tabriz, Iran.

Abstract: Testing the Hypothesis of independence of a p -variate vector subvectors, as a pretest for many others related tests, is always as a matter of interest. When the sample size n is much larger than the dimension p , the likelihood ratio test (LRT) with chisquare approximation, has an acceptable performance. However, for moderately high-dimensional data by which n is not much larger than p , the chisquare approximation for null distribution of the LRT statistic is no more usable. As a general case, here, a simultaneous subvectors independence testing procedure in all k p -variate normal distributions is considered. To test this hypothesis, a normal approximation for the null distribution of the LRT statistic was proposed. A simulation study was performed to show that the proposed normal approximation outperforms the chisquare approximation. Finally, the proposed testing procedure was applied on prostate cancer data.

Keywords: Multivariate normal distribution, likelihood ratio test, high-dimensional data, testing independence, multivariate Gamma function.

Mathematics Subject Classification (2010): 62H15, 62H10, 60F05.