

آزمون پرت بودن گروهی از مشاهدات چند متغیره

اباذر خلجی^۱، منوچهر خردمندینا^۲

چکیده:

فرض کنید m نمونه تصادفی n تایی از توزیع $N_p(\mu, \Sigma)$ داریم که مستقلند و می‌خواهیم فرض پرت بودن i امین نمونه n تایی را بیازمائیم. چنین آزمونی در بررسی‌های فاز اول کنترل فرآیند آماری چند متغیره با مشاهدات گروهی می‌تواند مفید باشد. امروزه مشهور است یک آماره بتا برای چنین آزمونی وجود دارد و از طریق رابطه توزیع بتا با توزیع F از آماره F می‌توان استفاده نمود. در متون آماری اثبات روشن و دقیقی برای توزیع آماره آزمون یافت نمی‌شود و در مواردی نیز، اثبات نادرست یا ناقص است. در این مقاله اثبات دقیق و نسبتاً روشن ارائه می‌دهیم و از طریق شبیه سازی قابلیت‌ها و ضعف‌های آزمون را مورد بررسی قرار می‌دهیم.

واژه‌های کلیدی: نقاط پرت، کنترل فرآیند چند متغیره، فاز اول، T^2 یت هتلینگ

۱ مقدمه

در مقاله حاضر اثباتی دقیق و نسبتاً روشن ارائه می‌دهیم.

یکی از کاربردهای گسترده آماره مورد مطالعه این مقاله، در کنترل فرآیندهای آماری در فاز اول است که در آن کنار گذاشتن مشاهدات غیر معمول به منظور اطمینان از تحت کنترل بودن فرآیند، مشابه برخورد با نقاط پرت می‌باشد [۱].

در مدل سازی آماری و تحلیل داده‌ها اغلب این پرسش مطرح می‌شود که آیا در مجموعه داده‌ها نقاط پرت وجود دارد؟ نقاط پرت مشاهده‌هایی هستند که از الگوی اکثر مشاهدات پیروی نمی‌کنند. وجود نقاط پرت در نمونه‌های یک یا چند متغیره باعث تخریب برآوردگرها می‌شود و علاوه بر آن در یک فرآیند کنترل وجود مشاهدات غیر معمول، میزان تغییرات را افزایش داده و همبستگی بین متغیرها را در هم می‌ریزد.

۲ چند قضیه کلیدی

در این بخش ۳ قضیه کلیدی ارائه می‌شود که در بخش‌های بعد مورد استفاده قرار می‌گیرند.

قضیه ۱.۲. اگر \bar{X} و S به ترتیب بردار میانگین و ماتریس کوواریانس یک نمونه تصادفی n تایی از $N_p(\mu, \Sigma)$ باشد، آن‌گاه:

$$\bar{X} \sim N_p(\mu, \frac{1}{n}\Sigma) \quad \perp \quad (n-1)S \sim W_p(\Sigma, n-1)$$

که نماد \perp بیانگر استقلال می‌باشد.

قضیه ۲.۲. اگر $F \sim F_{d_1, d_2}$ آن‌گاه

$$B = \frac{d_1 F}{d_2 + d_1 F} \sim \text{Beta}\left(\frac{d_1}{2}, \frac{d_2}{2}\right)$$

قضیه ۳.۲. اگر $Y \sim N_p(0, \Sigma)$ مستقل از $W \sim W_p(\Sigma, p)$ باشد، آن‌گاه:

$$T^2 = kY^T W^{-1} Y \sim \frac{kp}{k-p+1} F_{p, k-p+1}$$

به نظر می‌رسد اولین کار مهم در زمینه شناسایی نقاط پرت را ویلکس (۱۹۶۳) انجام داده است [۹]. سریواستاوا و ونروسن (۱۹۹۸) در حالت مشاهدات تکی نرمال چند متغیره یک آماره بتا برای آزمون پرت بودن مشاهده i ام بدست آوردند و با استفاده از رابطه توزیع بتا با توزیع F یک آماره F معرفی کردند [۷]. تعمیم آزمون پرت بودن از حالت مشاهدات تکی به مشاهدات گروهی در متون کنترل فرآیند آماری چند متغیره مورد توجه بوده است. به نظر می‌رسد که اولین تلاش مهم در این راستا توسط آلت (۱۹۷۳) انجام شده ولی ایشان به جای توزیع بتا اشتباهاً به توزیع T^2 هتلینگ رسیده است [۳]. تریسی و همکاران (۱۹۹۲) به آماره درست بتا رسیده‌اند ولی اثباتی که ارائه کرده‌اند مختصر است [۸].

^۱دانشجو کارشناسی ارشد آمار دانشگاه اصفهان

^۲استادیار گروه آمار دانشگاه اصفهان

$$F = \frac{k-p+1}{p} \mathbf{Y}^T \mathbf{W}^{-1} \mathbf{Y} \sim F_{p, k-p+1}$$

برای اثبات قضایای فوق به کتاب خردمندیا و علیرضایی [۲]

مراجعه کنید.

$$(m-1)S = (m-2)S_{(i)} + \frac{m}{m-1}(\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T$$

بر اساس قضیه ۱.۲ $\bar{\mathbf{X}}_{(i)} \sim N_p(\boldsymbol{\mu}, \Sigma)$ و مستقل از $(m-2)S_{(i)} \sim W_p(\Sigma, m-2)$ است.

$$\left. \begin{array}{l} S_{(i)} \perp \bar{\mathbf{X}}_{(i)} \\ S_{(i)} \perp \mathbf{X}_i \end{array} \right\} \Rightarrow S_{(i)} \perp (\mathbf{X}_i - \bar{\mathbf{X}}_{(i)})$$

با توجه به اینکه $\sqrt{\frac{m-1}{m}}(\mathbf{X}_i - \bar{\mathbf{X}}) = \sqrt{\frac{m-1}{m}}(\mathbf{X}_i - \bar{\mathbf{X}}_{(i)})$ می توان نوشت:

$$(m-2)S_{(i)} \perp \sqrt{\frac{m}{m-1}}(\mathbf{X}_i - \bar{\mathbf{X}})$$

اکنون تمامی شرایط برای استفاده از قضیه ۴.۲ برقرار است.

در واقع در این جا با علائم بکار رفته در قضیه ۴.۲ داریم

$\mathbf{Y} = \sqrt{\frac{m}{m-1}}(\mathbf{X}_i - \bar{\mathbf{X}})$ و $W_1 = (m-2)S_{(i)}$ ، $W = (m-1)S$ که W_1 مستقل از \mathbf{Y} است. حال با توجه به این قضیه می توان نوشت:

$$B_i = \frac{m}{(m-1)^2} (\mathbf{X}_i - \bar{\mathbf{X}})^T S^{-1} (\mathbf{X}_i - \bar{\mathbf{X}})$$

$$B_i \sim \text{Beta}\left(\frac{p}{2}, \frac{m-p-1}{2}\right)$$

که $i = 1, \dots, m$ حال با توجه به قضیه ۲.۲ در نهایت داریم: □

$$F_i = \left(\frac{m-p-1}{p}\right) \frac{B_i}{1-B_i} \sim F_{p, m-p-1}$$

اگر مقدار آماره F_i از $F_{\alpha, p, m-p-1}$ بزرگتر باشد می توان نتیجه گرفت که فاصله \mathbf{X}_i به طور معنی داری از میانگین مشاهدات دور است و لذا فرض H_0 مبنی بر پرت نبودن مشاهده i ام رد می شود. لازم به ذکر است که آماره فوق توسط سریواستاوا و ونروسن (۱۹۹۸) به عنوان آماره ای برای شناسایی یک نقطه پرت در مجموعه داده های نرمال چند متغیره معرفی شده است [۷]. در واقع با توجه به مطالب فوق در این بخش قضیه زیر را اثبات کردیم.

قضیه ۱.۳. اگر $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m$ یک نمونه تصادفی از $N_p(\boldsymbol{\mu}, \Sigma)$ باشد، آن گاه:

$$B_i = \frac{m}{(m-1)^2} (\mathbf{X}_i - \bar{\mathbf{X}})^T S^{-1} (\mathbf{X}_i - \bar{\mathbf{X}})$$

قضیه ۴.۲. فرض کنید $B = \mathbf{Y}^T \mathbf{W}^{-1} \mathbf{Y}$ که در آن بردار \mathbf{Y} دارای توزیع $N_p(0, \Sigma)$ و ماتریس W دارای توزیع $W_p(\Sigma, k)$ می باشد. اگر بتوان نوشت $W = W_1 + \mathbf{Y}\mathbf{Y}^T$ به طوری که $W_1 \sim W_p(\Sigma, k-1)$ و مستقل از بردار \mathbf{Y} باشد، آنگاه [۲]:

$$B \sim \text{Beta}\left(\frac{p}{2}, \frac{k-p}{2}\right)$$

اثبات. با استفاده از فرمول وارون جمع دو ماتریس (سیرل [۶] ۱۹۸۲) می توان نوشت:

$$\begin{aligned} B &= \mathbf{Y}^T \mathbf{W}^{-1} \mathbf{Y} = \mathbf{Y}^T (\mathbf{W}_1 + \mathbf{Y}\mathbf{Y}^T)^{-1} \mathbf{Y} \\ &= \mathbf{Y}^T \left[\mathbf{W}_1^{-1} - \frac{\mathbf{W}_1^{-1} \mathbf{Y}\mathbf{Y}^T \mathbf{W}_1^{-1}}{1 + \mathbf{Y}^T \mathbf{W}_1^{-1} \mathbf{Y}} \right] \mathbf{Y} \\ &= \frac{\mathbf{Y}^T \mathbf{W}_1^{-1} \mathbf{Y}}{1 + \mathbf{Y}^T \mathbf{W}_1^{-1} \mathbf{Y}} \end{aligned}$$

با توجه قضیه ۳.۲ داریم:

$$F = \frac{k-p}{p} \mathbf{Y}^T \mathbf{W}_1^{-1} \mathbf{Y} \sim F_{p, k-p}$$

بنابراین با استفاده از قضیه ۲.۲ می توان نوشت:

$$B = \frac{pF}{(k-p) + pF} \sim \text{Beta}\left(\frac{p}{2}, \frac{k-p}{2}\right)$$

۳ آزمون پرت بودن یک مشاهده تکی

در این بخش ابتدا حالتی که مجموعه داده ها تنها دارای یک مشاهده پرت است، مورد مطالعه قرار می گیرد و علاقمند به بررسی پرت بودن مشاهده i ام هستیم لذا فرض های آماری زیر را در نظر گرفته آماره آزمون را همراه با اثبات دقیق ارائه می دهیم.

$$\begin{cases} H_0: \text{پرت نبودن مشاهده } i \text{ ام} \\ H_1: \text{پرت بودن مشاهده } i \text{ ام} \end{cases}$$

برای این منظور فرض می کنیم $\bar{\mathbf{X}}$ و S به ترتیب بردار میانگین و ماتریس کوواریانس یک نمونه تصادفی m تایی از $N_p(\boldsymbol{\mu}, \Sigma)$ باشد. فرض کنید $\bar{\mathbf{X}}_{(i)}$ و $S_{(i)}$ به ترتیب بردار میانگین و ماتریس کوواریانس پس از حذف مشاهده i ام باشد، می توان نوشت:

آن‌گاه:

$$F_i = \left(\frac{m-p-1}{p}\right) \frac{B_i}{1-B_i} \sim F_{p,m-p-1} \quad B_i \sim \text{Beta}\left(\frac{p}{\gamma}, \frac{m-p-1}{\gamma}\right)$$

حال با توجه به قضیه ۲.۲ در نهایت داریم:
اگر مقدار مشاهده شده F_i بزرگتر از $F_{\alpha,p,m-p-1}$ باشد، فرض پرت نبودن نمونه n تایی i ام در سطح α رد می‌شود.

مثال ۲.۴. داده‌های هاوکینز و همکاران (۱۹۸۴)

این مجموعه داده‌ها شامل ۷۵ مشاهده ۳ متغیره است. روسو و زومرن (۱۹۹۰) ادعا می‌کنند که ۱۵ مشاهده اول این مجموعه داده‌ها، پرت هستند [۵]. برای دسترسی به داده‌های خام می‌توان به هاوکینز و همکاران مراجعه کرد [۴]. برای بررسی صحت این ادعا ما در این مثال مجموعه داده‌ها را به $m = 5$ گروه $n = 15$ تایی تقسیم کرده و میانگین هر گروه را محاسبه کردیم که بردارهای زیر حاصل می‌شود.

$$\bar{X}_1 = \begin{bmatrix} 10/01 \\ 20/49 \\ 29/43 \end{bmatrix}, \bar{X}_2 = \begin{bmatrix} 1/67 \\ 1/99 \\ 1/56 \end{bmatrix}, \bar{X}_3 = \begin{bmatrix} 1/43 \\ 1/61 \\ 1/65 \end{bmatrix}$$

$$\bar{X}_4 = \begin{bmatrix} 1/90 \\ 1/75 \\ 1/81 \end{bmatrix}, \bar{X}_5 = \begin{bmatrix} 1/03 \\ 1/69 \\ 1/70 \end{bmatrix}$$

براساس این بردارها آرایه‌های زیر حاصل می‌گردد.

$$\bar{\bar{X}} = \begin{bmatrix} 3/12 \\ 5/6 \\ 7/23 \end{bmatrix}, S_{\bar{\bar{X}}} = \begin{bmatrix} 3/34 & 28/47 & 41/24 \\ 28/47 & 67/88 & 94/67 \\ 41/24 & 94/67 & 137/83 \end{bmatrix}$$

حال با استفاده از قضیه ۱.۴ آماره آزمون برای گروه اول برابر است با

$$B_1 = \frac{5}{(5-1)^2} (\bar{X}_1 - \bar{\bar{X}}) S_{\bar{\bar{X}}}^{-1} (\bar{X}_1 - \bar{\bar{X}})^T = 0.9999$$

$$F_1 = \left(\frac{5-3-1}{3}\right) \frac{0.9999}{0.0001} = 20.396/59$$

که از $F_{7.5,3,1} = 215/71$ بزرگتر می‌باشد و فرض صفر مبنی بر پرت نبودن گروه اول رد می‌شود. نتایج آزمون در سطح ۰/۰۵ در شکل ۱ آورده شده است که با توجه به شکل مقدار آماره آزمون برای گروه اول بالاتر از مقدار بحرانی (خط ممتد) قرار گرفته است و بیانگر این واقعیت است که آزمون صحت ادعای روسو و زومرن مبنی بر پرت بودن ۱۵ مشاهده اول را تایید می‌کند.

۴ آزمون پرت بودن گروهی از مشاهدات چند متغیره

فرض کنید m نمونه تصادفی n تایی از $N_p(\mu, \Sigma)$ داریم که مستقلند و می‌خواهیم فرض‌های زیر را بیازمائیم:

$$\begin{cases} H_0: \text{نمونه } n \text{ تایی } i \text{ ام پرت نیست} \\ H_1: \text{نمونه } n \text{ تایی } i \text{ ام پرت است} \end{cases}$$

فرض کنید \bar{X}_i بردار میانگین i امین نمونه n تایی است $(i = 1, \dots, m)$ و قرار دهید:

$$\bar{\bar{X}} = m^{-1} \sum_{i=1}^m \bar{X}_i$$

$$S_{\bar{\bar{X}}} = (m-1)^{-1} \sum_{i=1}^m (\bar{X}_i - \bar{\bar{X}})(\bar{X}_i - \bar{\bar{X}})^T$$

همچنین فرض کنید $\bar{\bar{X}}_{(i)}$ و $S_{(i)}$ به ترتیب بردار میانگین و ماتریس کوواریانس میانگین‌های نمونه‌ای پس از حذف \bar{X}_i باشند، می‌توان نوشت:

$$\bar{\bar{X}}_{(i)} = \frac{m}{m-1} \bar{\bar{X}} - \frac{1}{m-1} \bar{X}_i$$

$$(m-1)S_{\bar{\bar{X}}} = (m-2)S_{(i)} + \frac{m}{m-1} (\bar{X}_i - \bar{\bar{X}})(\bar{X}_i - \bar{\bar{X}})^T$$

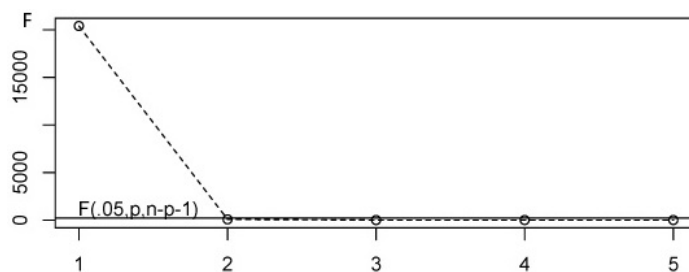
طبق قضیه ۱.۲ $\bar{\bar{X}}_{(i)} \sim N_p(\mu, \frac{1}{n(m-1)}\Sigma)$ و مستقل از $(m-2)S_{(i)} \sim W_p(\frac{1}{n}\Sigma, m-2)$ است. با کمی دقت ملاحظه می‌شود که شرایطی همانند قضیه ۱.۳ برقرار است. بنابراین بدون این که لازم باشد اثبات جدیدی ارائه شود، صرفاً براساس قضیه ۱.۳ می‌توان قضیه زیر را نتیجه گرفت.

قضیه ۱.۴. اگر $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_m$ یک نمونه تصادفی از $N_p(\mu, \frac{1}{n}\Sigma)$ باشد، آن‌گاه:

$$B_i = \frac{m}{(m-1)^2} (\bar{X}_i - \bar{\bar{X}}) S_{\bar{\bar{X}}}^{-1} (\bar{X}_i - \bar{\bar{X}})^T$$

آن‌گاه:

$$B_i \sim \text{Beta}\left(\frac{p}{\gamma}, \frac{m-p-1}{\gamma}\right)$$



شکل ۱. نتایج آزمون F برای داده‌های گروه بندی شده هاوکینر و همکاران که در آن نقاط بیانگرمقدار آماره آزمون برای گروه‌ها و خط ممتد مقدار بحرانی آزمون می‌باشد

شده است. در جدول ۱ درصد مواردی که آزمون نمونه n تایی پرت را بدرستی تایید کرده ارائه گردیده است.

جدول ۱: درصد تشخیص‌های صحیح در آلودگی میانگین

	m	n			
		۴	۱۰	۲۰	۳۰
a_1	۵	۸/۳	۱۱	۱۴/۲	۱۵/۸
	۱۰	۱۸/۷	۱۹/۲	۷۶/۵	۸۹/۹
	۳۰	۲۸/۱	۷۴/۱	۹۶	۹۹/۸
	۱۰۰	۳۶/۱	۷۷/۵	۹۸/۲	۹۹/۹
a_2	۵	۱۴/۹	۲۲/۱	۳۳/۷	۳۸/۷
	۱۰	۸۷/۱	۱۰۰	۱۰۰	۱۰۰
	۳۰	۹۸/۵	۱۰۰	۱۰۰	۱۰۰
	۱۰۰	۹۹/۶	۱۰۰	۱۰۰	۱۰۰
a_3	۵	۱۳/۴	۲۲/۶	۳۳/۳	۴۰/۴
	۱۰	۸۴/۵	۹۹/۵	۱۰۰	۱۰۰
	۳۰	۹۸/۵	۱۰۰	۱۰۰	۱۰۰
	۱۰۰	۹۹/۳	۱۰۰	۱۰۰	۱۰۰

با توجه به جدول ۱ با افزایش m و n عملکرد آزمون بهبود می‌یابد که افزایش m عملکرد آزمون را نسبت به n با نرخ سریع‌تری بهبود می‌بخشد. بردار a_3 نیز نسبت به دو بردار دیگر آلودگی بیشتری تولید می‌کند زیرا در این حالت با افزایش میانگین یک عضو همبستگی بین متغیرها نیز تحت

مثال ۳.۴. مطالعات شبیه سازی

در این مثال به کمک شبیه سازی و با استفاده از نرم افزار R عملکرد آزمون پرت بودن گروهی از مشاهدات را در آلودگی‌های مختلف مورد بررسی قرار می‌دهیم، به این صورت که گروه پرت را با میانگین، واریانس و همبستگی دیگری تولید می‌کنیم. در واقع در آلودگی واریانس و همبستگی ماتریس کوواریانس را آلوده می‌کنیم. در هر اجرا m نمونه n تایی (با مقادیر مختلف m و n) از توزیع $N_3(\mu, \Sigma)$ تولید می‌کنیم.

$$\mu = \begin{bmatrix} 1 \\ 5 \\ 9 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0/9 & 0/9 \\ 0/9 & 1 & 0/9 \\ 0/9 & 0/9 & 1 \end{bmatrix}$$

که نمونه n تایی m ام در همه موارد آلوده می‌باشد. در نهایت برای هر یک از سه حالت آلودگی میانگین، آلودگی واریانس و آلودگی همبستگی نتایج هزار تکرار آزمون را ارائه می‌دهیم.

۱. آلودگی میانگین

در این حالت $m-1$ نمونه n تایی اول را از توزیع $N_3(\mu, \Sigma)$ تولید کرده‌ایم ولی آخرین نمونه را با میانگین $\mu_{outlier} = \mu + a_i$ تولید کرده‌ایم. در جدول ۱ برای سه مقدار

$$a_1 = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}^T \quad a_2 = \begin{bmatrix} 1 & 1 & 0 \end{bmatrix}^T$$

$$a_3 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}^T$$

و برای مقادیر $m = 5, 10, 30, 100$ و $n = 4, 10, 20, 30$ نتایج هزار تکرار آزمون در سطح $0/05$ ارائه

تأثیر قرار می‌گیرد.

۲. آلودگی واریانس

در این حالت ۱- m نمونه n تایی اول را از توزیع $N_3(\mu, \Sigma)$ و آخرین نمونه را از توزیع $N_3(\mu, k\Sigma)$ برای مقادیر مختلف k تولید کرده‌ایم. در جدول ۲ برای مقادیر $k = 2, 5, 10$ ، $n = 4, 10, 20, 30$ و $m = 5, 10, 30, 100$ تکرار آزمون در سطح 0.05 و درصد مواردی که آزمون نمونه n تایی پرت را بدرستی تایید کرده ارائه گردیده است.

جدول ۲: درصد تشخیص‌های صحیح در آلودگی واریانس

m	n				
	۴	۱۰	۲۰	۳۰	
$k = 2$	۵	۷/۱	۷/۴	۷/۹	۷/۶
	۱۰	۱۵/۵	۱۳/۸	۱۴/۹	۱۵
	۳۰	۲۲/۴	۲۲/۹	۲۴/۹	۲۳/۱
	۱۰۰	۲۳/۷	۲۴/۵	۲۷/۳	۲۳/۵
$k = 5$	۵	۹/۶	۹/۷	۸/۶	۸/۶
	۱۰	۴۴/۹	۴۴/۹	۴۷/۲	۴۱/۹
	۳۰	۶۰/۲	۶۱/۶	۶۰/۲	۶۲/۹
	۱۰۰	۶۴/۱	۶۳	۶۴/۱	۶۵
$k = 10$	۵	۱۴/۷	۱۲/۹	۱۴/۳	۱۳/۳
	۱۰	۶۶/۴	۶۷/۹	۶۸/۲	۶۹/۳
	۳۰	۸۲	۸۲/۵	۸۰	۸۳
	۱۰۰	۸۳/۴	۸۶	۸۴/۲	۸۶/۵

۳. آلودگی همبستگی

در این حالت نیز ۱- m نمونه n تایی اول را از توزیع $N_3(\mu, \Sigma)$ تولید کرده‌ایم ولی آخرین نمونه را ماتریس کوواریانس دیگری از بین سه ماتریس کوواریانس زیر تولید کرده‌ایم.

$$\Sigma_1 = \begin{bmatrix} 1 & 0.5 & 0.5 \\ 0.5 & 1 & 0.5 \\ 0.5 & 0.5 & 1 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\Sigma_3 = \begin{bmatrix} 1 & -0.5 & -0.5 \\ -0.5 & 1 & -0.5 \\ -0.5 & -0.5 & 1 \end{bmatrix}$$

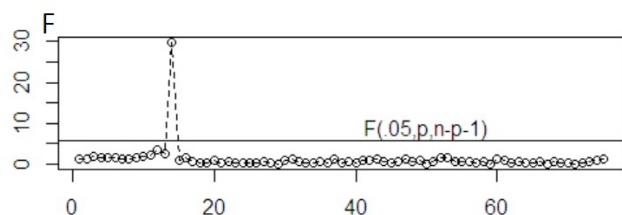
در جدول ۳ برای سه مقدار و برای مقادیر $n = 4, 10, 20, 30$ و $m = 5, 10, 30, 100$ تکرار آزمون در سطح 0.05 ارائه شده است. در جدول ۳ درصد مواردی که آزمون نمونه n تایی پرت را بدرستی تایید کرده ارائه گردیده است.

جدول ۳: درصد تشخیص‌های صحیح در آلودگی همبستگی

m	n				
	۴	۱۰	۲۰	۳۰	
Σ_3	۵	۱۲	۱۳/۴	۱۲/۸	۱۴/۳
	۱۰	۶۰/۲	۶۱/۱	۶۲	۶۳/۹
	۳۰	۷۴/۱	۷۲/۲	۷۶/۲	۷۳/۳
	۱۰۰	۷۲	۷۴/۲	۷۴/۷	۷۵/۷
Σ_2	۵	۱۰/۴	۱۳	۱۱/۳	۱۱/۶
	۱۰	۴۹/۸	۵۲	۵۳/۴	۴۹/۴
	۳۰	۶۳	۶۶/۱	۶۷	۶۴/۲
	۱۰۰	۶۳/۱	۶۶/۵	۶۷	۶۹/۶
Σ_1	۵	۸/۷	۸/۶	۹	۹/۲
	۱۰	۳۲/۳	۳۱/۱	۳۲/۹	۳۱/۲
	۳۰	۴۳/۹	۴۵/۲	۴۲/۲	۴۲/۳
	۱۰۰	۴۷/۵	۴۶/۱	۵۰	۴۹

باتوجه به جداول فوق افزایش مقدار m و n موجب بهبود عملکرد آزمون می‌شود. ولی به نظر می‌رسد که حداقل در شرایط مثال

آزمون در شکل ۲ آورده شده است که تنها مشاهده ۱۴ بالا تر از مقدار بحرانی (خط ممتد) قرار گرفته است به عبارت دیگر سایر مشاهدات پرت نیستند که خلاف واقعیت و نتایج مثال ۲.۴ می باشد.



شکل ۲. نتایج آزمون مشاهدات تکی برای داده‌های هاوکینر

بنابراین زمانی که با چندین داده پرت در مجموعه داده‌ها مواجه هستیم بهتر است از آزمون پرت بودن گروهی از مشاهدات چند متغیره استفاده کنیم زیرا با این آزمون داده‌های پرت در یک گروه قرار گرفته و باعث می‌شوند که فاصله آن‌ها با میانگین داده حفظ شود و به تبع آن آماره آزمون از مقدار بحرانی بزرگتر شده و به نتیجه مطلوب برسیم.

حاضر، این آزمون در تشخیص آلودگی ضریب همبستگی ضعیف است ولی در تشخیص آلودگی در بردار میانگین عملکرد خوبی دارد. این نکته حائز اهمیت است که آلودگی‌های بکارگرفته شده در مثال حاضر بسیار خفیف می‌باشد که در عمل و در شرایط طبیعی با آلودگی‌های شدیدتری روبرو خواهیم بود.

۵ نتیجه گیری

داده‌های مثال ۲.۴ را در نظر بگیرید، این مجموعه داده‌های مشهور به منظور بررسی روش‌های شناسایی نقاط پرت و آزمون آن‌ها در بسیاری از مقالات علمی و پژوهشی مورد استفاده قرار گرفته است که اکثر پژوهشگران پرت بودن ۱۵ مشاهده اول این مجموعه داده‌ها را تایید کرده‌اند. حال اگر از آزمون پرت بودن یک مشاهده تکی (بخش ۳) را برای این مجموعه داده‌ها استفاده کنیم مشاهده می‌شود که این آزمون تنها پرت بودن مشاهده ۱۴ را تایید می‌کند. نتیجه محاسبات آماره آزمون برای تک تک مشاهدات و مقدار بحرانی

مراجع

- [۱] توانگر، م، طالبی، ه، علامت ساز، م. (۱۳۹۰)، خواص توزیعی آماره T^2 هتلینگ در کنترل کیفیت چند متغیره. اندیشه آماری
- [۲] خردمندیا، م، علیرضایی، م، مقدمه‌ای بر آمار چند متغیره، انتشارات نگارخانه، اصفهان، ۱۳۹۲، (۱۶)، ۱۷-۲۷.
- [3] Alt, F. B. (1973). Aspects of Multivariate Control Charts. *Master of Science. Georgia Institute of Technology*, 6
- [4] Hawkins, D. M. , Bradu, D. and Kass, G. V. (1984). Location of Several Outlier in Multiple Regression Data Using Elemental Sets. *technometrics*, 26, 197-208.
- [5] Rousseeuw, P. J. and van Zomeren, B. C. (1990). Unmasking Multivariate Outliers and Leverage Points. *Journal of the American Statistical Association*. 85, 633-639.
- [6] Searle, S. (1928). *Matrix Algebra Useful for Statistics*. New York: John Wiley and Sons.
- [7] Srivastava, M. S. and von Rosen, D. (1998). Outlier in Multivariate Regression Models, *J. Multivariate Anal.*, 65, 195-208.
- [8] Tracy, N. D. ; Young, J. C. and Mason, R. L. (1992). Multivariate Control Chart for Individual Observation. *Journal of Quality Technology*, 24, 88-95.
- [9] Wilks, S. S. (1963). Multivariate statistical outliers, *Sankhya: The Indian Journal of Statistics*, 25(4), 407-426.