







دانشگاه پیام نور مرکز شیراز  
دانشکده علوم گروه آمار

عنوان پایان نامه

## ماتریس کوواریانس مجانبی میانه اوجا

برای دریافت درجه کارشناسی ارشد

رشته آمار ریاضی

استاد راهنما:

دکتر نرگس عباسی

استاد مشاور:

دکتر علی رضا نعمت‌اللهی

مؤلف:

صفیه آتش‌روز

تیر 1389

## چکیده

در بسیاری از تحلیل‌های آماری، برخی از داده‌ها نقش بسزایی در تحلیل‌ها ایجاد می‌کنند و شاید باعث اختلال در نتیجه‌گیری‌ها شوند از جمله‌ی این داده‌ها، نقاط پرت و داده‌های آلوده می‌باشند که در کران نفوذپذیری تأثیر گذارند و همچنین برآوردهای پایدار از پارامترهای مکان چندمتغیره و ترکیبی از آنها در مسائل آماری مورد تردید قرار می‌گیرند. در این پایان‌نامه، میانه‌ی پایای تعمیم‌یافته چندمتغیره خطی را که توسط اوجا پیشنهاد شده معرفی و سپس به مقایسه با میانه‌ی فاصله‌ای و میانه در حالت کروی، و بیضوی پرداخته و تأثیر این نقاط پرت و آلوده را روی آنها بررسی می‌کنیم. سپس به معرفی برآوردهای ماتریس کوواریانس میانه پرداخته و با ارائه دو روش اندازه‌گیری و احتمال پوشش تجربی (CP) به مقایسه عددی برآوردها می‌پردازیم. اوجا در واقع میانه در حالت چند متغیره را به میانه‌ی فاصله‌ای در حالت یک‌متغیره ساده می‌کند و نشان می‌دهد که به طور مجانبی به کارایی میانه‌ی فاصله‌ای است و دارای خصوصیات کارا و پایدار خوش‌رفتاری است. از این رو در این پایان‌نامه به کارایی مجانبی، تابع نفوذپذیری و اثرات نقاط فروریزش در حالت کلی و برای میانه اوجا، اشاره خواهیم کرد. سپس سه مفهوم اساسی، علامت و رتبه چندمتغیره، ماتریس کوواریانس رتبه و علامت و برآوردهای آن، همچنین پایداری و کارایی این برآوردها، ارائه خواهد شد.

نهایتاً با توجه به مدل‌های چندمتغیره متقارن و راه‌های گوناگون برآورد ماتریس کوواریانس علامت مربوط به این مدل به مقایسه عددی این برآوردها می‌پردازیم و همچنین وابستگی این برآوردها را به نقاط فروریزش در نمونه‌های متناهی و همچنین اربیبی، آنها را نشان می‌دهیم سپس پیشنهادی برای از بین بردن میزان اربیبی ارائه خواهیم کرد.

## فهرست مطالب

صفحه	عنوان
1	مقدمه
	<b>فصل اول</b>
3	مفاهیم
4	توزیع بیضی‌گون و کروی
4	1-1 معرفی توزیع بیضی‌گون
6	1-1-1 معرفی برخی از توزیع‌های کروی
7	2-1 داده‌های پرت
7	1-2-1 روش تعیین داده‌های پرت
9	3-1 داده‌های آلوده
10	4-1 نقاط فروریزش
12	5-1 کوواریانس رتبه و علامت
12	1-5-1 رتبه و علامتهای حاشیه‌ای
21	2-5-1 رتبه و علامت فاصله‌ای
28	3-5-1 رتبه و علامت‌های اوجا
32	6-1 آماره بوت استرپ
33	1-6-1 یک مثال در حالت دوبعدی
34	7-1 آشنایی با آماره جک‌نایف و ارائه‌ی چند مثال مربوط به آن
37	8-1 تجزیه چولسکی
38	9-1 توزیع بادم‌های سنگین
39	10-1 آزمون نوع‌والد
40	11-1 آماره یو
43	12-1 پایایی

صفحه	عنوان
45	13-1 بحث و نتیجه‌گیری
<b>فصل دوم</b>	
46	میان‌اوجا
47	1-2 مقدمه
47	2-2 میان‌اوجا
50	3-2 تابع هدف یا آماره آزمون (آزمون علامت) در حالت سه‌متغیره
54	4-2 تابع هدف یا آماره آزمون علامت در حالت $k$ متغیره تعمیم یافته
57	5-2 خصوصیات مجانبی میان‌اوجا ی تعمیم یافته
58	6-2 کارایی مجانبی در حالتی از نرمال چند متغیره
63	7-2 نقاط فروریزش میان‌اوجا ی اوجا
65	8-2 ارائه یک مثال کاربردی
68	9-2 بحث و نتیجه‌گیری
<b>فصل سوم</b>	
69	برآوردهای ماتریس کوواریانس علامت
70	1-3 مقدمه
70	2-3 توزیع مجانبی
72	3-3 برآوردهایی از کوواریانس مجانبی $\theta$
76	3-3 استواری در برآورد کوواریانس
78	4-3 ماتریس کوواریانس و برآورد آن برای مدل بیضی گون
83	2-4-3 اثبات همگرایی در احتمال ماتریس $A$ به ماتریس $A$
84	5-3 مطالعه ی کارایی و استواری و مقایسه این دو کمیت در میانگین های برآورد شده
85	تعریف فاصله تحمل

صفحه	عنوان
93	6-3 بحث و نتیجه‌گیری
	<b>فصل چهارم</b>
94	مقایسه برآوردگرهای ماتریس کوواریانس
95	1-4 مقدمه
95	2-4 برآوردگرهای ماتریس کوواریانس
96	3-4 تاثیر نقاط پرت در برآورد ماتریس کوواریانس در حالت دوبعدی
	4-4 ارزیابی نمونه متناهی و بررسی تاثیرات نقاط پرت و آلوده در برآوردهای ماتریس
98	کوواریانس در حالت سه بعدی
100	5-4 ارزیابی برآوردگرها به دو روش اندازه‌گیری و احتمال پوشش تجربی $CP^2$
102	6-4 بحث و نتیجه‌گیری
<b>103</b>	<b>پیشنهاد:</b>
104	منابع
111	پیوست
116	واژگان
120	چکیده انگلیسی

## فهرست نمودارها

- نمودار 1-5-1 نمودار  $R(x)$  مربوط به مثال 1-1-5-1 ..... 13
- نمودار 2-5-1 نمودار  $D(x)$  مربوط به مثال 2-1-5-1 ..... 14
- نمودار 3-5-1 نمودار  $D2(x)$  مربوط به مثال 2-2-5-1 ..... 25
- نمودار 4-5-1 نمودار  $H3x$  مربوط به مثال 1-3-5-1 ..... 29
- نمودار 5-5-1 نمودار  $D3x$  مربوط به مثال 1-3-5-1 ..... 29
- نمودار 1-7-2 نمودار مقایسه میانه فاصله‌ای و میانه اوجا و بردار میانگین تحت تأثیر نقاط فروریزش با  $n$  داده مطلوب و  $m$  داده نامطلوب ..... 64
- نمودار 1-5-3 نمودار تحمل بیضی‌گون ماتریس کوواریانس نمونه‌ای و ماتریس کوواریانس فاصله‌ای  $\tau$ -کندال و ماتریس کوواریانس رتبه‌ای اوجا با حضور داده‌های پرت. .... 88
- نمودار 2-5-3 نمودار حساسیت ماتریس‌های کوواریانس علامت نمونه‌ای و ماتریس کوواریانس فاصله‌ای  $\tau$ -کندال و ماتریس کوواریانس رتبه‌ای اوجا تحت نفوذ داده‌های پرت. .... 90



## فهرست جداول

- جدول 1-1-1 زیر گروههایی از توزیع های کروی  $p$  بعدی ..... 6
- جدول 1-6-2 مشاهدات برای 14 درخت در 4 جهت متفاوت ..... 63
- جدول 1-7-2 معرفی تعداد داده های مطلوب و نسبت تعداد مطلوب به کل، برای چهار جامعه متفاوت ..... 65
- جدول 1.8.2 مشاهدات برای 14 درخت در چهار جهت متفاوت ..... 66
- جدول 2-8-2 مقادیر ملاحظه شده برای  $x_1$  و  $x_2$  و  $x_3$  و  $q_1$  و  $q_2$  و  $q_3$  های متناظر با آنها ..... 67
- جدول 1-5-3، کارایی میانه اوجا و فاصله ای نسبت به اولین بردار ویژه و معیارهای شرطی ..... 93
- جدول 1-3-4 مقایسه مدلی از ماتریس کوواریانس نمونه متناهی و متوسط مقدار ماتریس کوواریانس  $V$  و  $V_S$  در حالت  $N_2(0, I_2)$  با حضور صفر و سه داده پرت ..... 98
- جدول 1-4-4 مقایسه مدل ماتریس کوواریانس نمونه متناهی و متوسط ماتریس کوواریانس  $V$  و  $V_S$  در حالت  $N_3(0, I_3)$  با حضور صفر داده پرت و 10 درصد آلودگی ..... 100

## برنامه‌های کامپیوتری

- 113..... برنامه کامپیوتری با استفاده از نرم‌افزار میپل مربوط به نمودار 1-5-1.....
- 113..... برنامه کامپیوتری با استفاده از نرم‌افزار میپل مربوط به نمودار 2-5-1.....
- 114..... برنامه کامپیوتری با استفاده از نرم‌افزار میپل مربوط به نمودار 3-5-1.....
- 115..... برنامه کامپیوتری با استفاده از نرم‌افزار میپل مربوط به نمودار 4-5-1.....
- 116..... برنامه کامپیوتری با استفاده از نرم‌افزار میپل مربوط به نمودار 5-5-1.....

## مقدمه

در بعضی از تحلیل‌های آماری ممکن است با داده‌هایی روبه‌رو شویم که در نتیجه‌گیری‌هایمان تأثیرات بسزایی دارند و باعث اختلال در تحلیل می‌گردند و در این صورت ما را در تصمیم‌گیری دچار تردید و شک می‌کنند، از جمله داده‌های پرت و داده‌های آلوده، بنابراین بایستی به دنبال راه حل و پیشنهادی برای حل این موضوع بود. هدف این پایان‌نامه معرفی میانه‌ی پیشنهاد شده توسط اوجا (اوجا 1999) و خصوصیات مربوط به آن، ماتریس کوواریانس علامت (هتمنسپرگر 1994) می‌باشد، سپس برآوردهایی برای ماتریس کوواریانس و مقایسه عددی آنها (اوجا، هتمنسپرگر و نادار 2003)، تأثیر نقاط پرت و آلوده روی میانه اوجا و برآورد ماتریس کوواریانس میانه اوجا (اوجا، کوئینن و ویسوری 2000)، ارائه خواهد شد. در این میان به مسائلی همچون کارایی مجانبی (نینیما و اوجا 1994)، نقاط فروریزش (نینیما و اوجا 1994)، کران نفوذپذیری (نینیما و اوجا 1994) تعریف توزیع‌های بیضوی (اوجا و کوئینن و ویسوری 2000)، معرفی برخی از توزیع‌های کروی، مقایسه میانه چندمتغیره در حالت بیضوی با میانه تعمیم‌یافته معرفی شده توسط اوجا (نینیما و اوجا 1984) و همچنین رتبه و علامت میانه فاصله‌ای و میانه اوجا (اوجا و کوئینن و ویسوری 2000)، مطالعه پایداری و کارایی برآوردهای ماتریس کوواریانس (اوجا و کوئینن و ویسوری 2000) خواهیم پرداخت.

در فصل اول به مفاهیم اولیه که در متن این پایان‌نامه به دفعات مورد استفاده قرار گرفته است، اشاره می‌کنیم و در فصل دوم به معرفی میانه اوجا، تابع اصلی یا تابع هدف در حالت سه‌متغیره و  $k$  متغیره و ارائه یک مثال می‌پردازیم. سپس خصوصیات مجانبی میانه تعمیم‌یافته و کارایی مجانبی درحالتی از نرمال چندمتغیره و تابع نفوذپذیری میانه‌ی اوجا و سپس نتایج مربوط به آن می‌پردازیم.

در فصل سوم ماتریس کوواریانس علامت و رتبه‌ای (هتمنسپرگر 1999) و مدل‌های بیضی‌گون (اوجا و کوئینن و ویسوری 2000) را شرح داده و با ارائه‌ی برآوردهایی برای ماتریس کوواریانس (اوجا و هتمنسپرگر 2003) به مطالعه کارایی و پایداری آنها می‌پردازیم در این میان چند قضیه و اثبات مربوط به آن ارائه می‌گردد.

در فصل چهارم با تعریف داده‌های آلوده و تأثیرات آنها به بررسی اریبی این برآوردها برای نمونه‌های متناهی می‌پردازیم و یافته‌هایمان را از برآورد ماتریس کوواریانس در حالت دو و سه‌متغیره به مقایسه عددی گذاشته و با بیان نتایج، پیشنهادهایی ارائه خواهیم کرد.

# فصل اول

مفاهيم

در این فصل برخی از مفاهیم به کار گرفته در این رساله آورده شده است و برای درک و تصور بهتر مطالب در برخی از بخش‌ها مثال‌های عددی به همراه نمودار ارائه گردیده است.

## توزیع بیضی‌گون و کروی

توزیع‌های کروی<sup>1</sup> و بیضی‌گون<sup>2</sup> به ترتیب تعمیم توزیع نرمال استاندارد چندمتغیره‌ی  $N_d(0, I_d)$  و توزیع نرمال چندمتغیره‌ی  $N_d(\mu, \Sigma)$  می‌باشند (رجوع کنید به اندرسن<sup>3</sup> (1993)، لانگ<sup>4</sup> (1979)، آلکین<sup>5</sup> (1992)، یان و بنتلر<sup>6</sup> (2004)، زلنر<sup>7</sup> (1976) ) می‌توان بسیاری از روش‌های قدیمی تحلیل‌های آماری را مستقیماً یا با اندکی تغییر در مورد جوامع با توزیع بیضی‌گون به کار برد.

### 1-1 معرفی توزیع بیضی‌گون

یک بردار تصادفی با توزیع بیضی‌گون، مانند  $\mathbf{X}$ ، به وسیله یک سه‌تایی  $(\mu, \mathbf{D}, g)$  مشخص می‌شود و به صورت  $X \sim EC_p(\mu, \mathbf{D}, g)$  نمایش داده می‌شود که در آن  $E(\mathbf{X}) = \mu$  و  $Var(\mathbf{X}) = \Sigma = a\mathbf{D}$  که  $a$  یک مقدار ثابت مثبت و  $g$  تابعی است که یک خانواده خاص از توزیع‌های بیضی‌گون را مشخص می‌نماید، تابع چگالی بردار تصادفی بیضی‌گون در صورت وجود به صورت زیر نوشته می‌شود:

$$f_{\mathbf{X}}(\mathbf{x}) = |\mathbf{D}|^{-\frac{1}{2}} g[(\mathbf{x} - \mu)' \mathbf{D}^{-1} (\mathbf{x} - \mu)]$$

به عنوان مثال خانواده نرمال چندمتغیره نامنفرد (دارای ماتریس کوواریانس وارون‌پذیر) به این کلاس تعلق دارد، که در آن  $a = 1$  و  $g(t) = \exp\left\{-\frac{t}{2}\right\}$  در نظر گرفته می‌شود.

توزیع‌های کروی حالت خاصی از توزیع‌های بیضی‌گون هستند، بدین صورت که اگر در یک توزیع بیضی‌گون  $\mu = 0$  و  $\mathbf{D} = \mathbf{I}_p$  در نظر گرفته شود، توزیع حاصل توزیع کروی خواهد بود. در نتیجه اگر  $Y$  دارای توزیع کروی باشد، آنگاه با انتخاب یک تبدیل مناسب، توزیع  $\mathbf{X} = \mu + \mathbf{B}Y$  بیضی‌گون خواهد بود که در آن  $\mathbf{D} = \mathbf{B}\mathbf{B}'$  و  $Var(\mathbf{X}) \propto \mathbf{B}\mathbf{B}'$  می‌توان بردار تصادفی بیضی‌گون  $\mathbf{X}_{p \times 1}$  را به صورت  $\mathbf{X} = \mathbf{R}\mathbf{B}\mathbf{U} + \mu$  نیز نمایش داد. که در آن  $\mathbf{U}$  یک بردار

<sup>1</sup>Spherical Distribution

<sup>2</sup> Empirical Distribution

<sup>3</sup> -Anderson

<sup>4</sup> -lang

<sup>5</sup> -Olkin

<sup>6</sup> -Bentler and yuan

<sup>7</sup> -Zellner

تصادفی است که به طور یکنواخت روی سطح کره واحد در  $R^p$  توزیع شده است و  $R$  یک متغیر تصادفی مثبت و مستقل از  $U$  می باشد که با استفاده از آن می توان به تابع  $g$  دست پیدا کرد. به عنوان مثال اگر  $R^2$  دارای توزیع کای دو با  $p$  درجه آزادی باشد، آنگاه  $\mathbf{X}$  دارای توزیع نرمال  $p$  متغیره است.

فانگ<sup>8</sup> (1990) با فرض  $X \sim EC_p(\boldsymbol{\mu}, \mathbf{D}, g)$  و  $f(\cdot)$  تابع چگالی متغیر تصادفی  $R$ ، رابطه میان  $f(\cdot)$  و  $g(\cdot)$  را به صورت زیر بیان کرد.

$$f(r) = \frac{2 \pi^{\frac{n}{2}}}{\Gamma\left(\frac{n}{2}\right)} r^{n-1} g(r^2)$$

---

<sup>8</sup> Fang

## 1-1-1 معرفی برخی از توزیع‌های کروی

جدول زیر به اختصار برخی از توزیع‌های کروی را معرفی می‌نماید:

نوع	تابع چگالی $f(x)$ یا تابع مشخصه $\Psi(t)$
نوع کوتر	$f(x) = c(x'x)^{N-1} \exp(-r(x'x)^s), r, s > 0, 2N + p > 2$
نرمال چند متغیره	$f(x) = c \exp(-1/2 x'x)$
پیرسن نوع $V$	$f(x) = c(1 + \frac{x'x}{s})^{-N}, s > 0, N > \frac{p}{2}$
$t$ چند متغیره	$f(x) = c(1 + \frac{x'x}{s})^{-\frac{(p+m)}{2}}, m > 0, m \in \mathbb{Z}$
کوشی چند متغیره	$f(x) = c(1 + \frac{x'x}{s})^{-\frac{(p+1)}{2}}, s > 0$
پیرسن نوع $2$	$f(x) = c(1 - x'x)^m, m > 0$
لجستیک	$f(x) = c \frac{\exp(-x'x)}{[1 + \exp(-x'x)]^2}$
بسل چند متغیره ( $k_\alpha(\cdot)$ که در آن مبین تابع بسل نوع سوم است)	$f(x) = c \left(\frac{\ x\ }{\beta}\right)^\alpha k_\alpha\left(\frac{\ x\ }{\beta}\right), \alpha > -p/2, \beta > 0$
ترکیب مقیاس که در آن $G(t)$ یک تابع توزیع پیوسته است.	$f(x) = c \int_0^\infty t^{-\frac{p}{2}} \exp\left(-\frac{x'x}{2t}\right) dG(t)$
قوانین پایا	$\Psi(t) = \exp\left\{r(t't)^{\frac{\alpha}{2}}\right\}, 0 < \alpha \leq 2, r < 0$
یکنواخت چند متغیره که در آن $\mathbf{0}_1^F(\cdot, \cdot)$ یک تابع فوق هندسی تعمیم یافته است	$\Psi(t) = \mathbf{0}_1^F\left(\frac{p}{2}; -\frac{1}{4} \ t\ ^2\right)$

جدول (1-1-1): زیر گروههایی از توزیع‌های کروی  $p$  بعدی

همانگونه که عنوان شد بردار تصادفی  $\mathbf{X}$ ، با توزیع بیضی‌گون، می‌تواند به صورت  $\mathbf{X} = \mathbf{RBU} + \boldsymbol{\mu}$  نمایش داده شود. که در آن  $\mathbf{B}$  ماتریس تبدیلات،  $\mathbf{U}$  یک بردار تصادفی است که به طور یکنواخت روی سطح کره واحد در  $R^p$  توزیع شده و  $R$  یک متغیر تصادفی مستقل از  $\mathbf{U}$  می‌باشد پس برای مشخص شدن یک توزیع بیضی‌گون باید  $E(\mathbf{X})$  و  $Var(\mathbf{X})$  و توزیع  $R$  معلوم باشد. میانگین و واریانس را می‌توان با استفاده از میانگین و واریانس نمونه‌ای داده‌ها یا میانگین و واریانس نمونه‌ای داده‌های پیراسته<sup>9</sup>، که با حذف درصدی از نقاط پرت و دور افتاده به دست می‌آید برآورد نمود.

کوک<sup>10</sup> (1994 و 1998) بیان داشت که انتخاب توزیع  $R$  اختیاری است، بدین معنا که با انتخاب‌های مختلف برای توزیع  $R$  می‌توان توزیع‌های مختلفی به دست آورد، مثلاً برای دستیابی به توزیع نرمال، متغیر  $R^2$  باید دارای توزیع کای دو باشد. در عمل برای دستیابی به یک انتخاب

<sup>9</sup> Trimmed Data

<sup>10</sup> - Cook

مناسب برای توزیع  $R$  می‌توان از توزیع تجربی شعاع‌های نمونه<sup>11</sup>، که از رابطه زیر به دست می‌آیند، استفاده نمود. (در این رساله به دلیل عدم نیاز به اجرای این مراحل اشاره نخواهد شد).

$$r_i = \left[ (x_i - \bar{x})' \widehat{\Sigma}^{-1} (x_i - \bar{x}) \right]^{\frac{1}{2}}$$

میانگین پیراسته: تمام مشاهدات کوچکتر از چارک اول و بزرگتر از چارک سوم را حذف، و میانگین داده‌های بین  $Q_1$  و  $Q_3$  را محاسبه می‌کند.

## 2-1 داده‌های پرت

بارنت<sup>12</sup> (1994) داده‌ای که بیشترین فاصله از بقیه مشاهدات دارد را به عنوان داده پرت معرفی کرد. و گرابس<sup>13</sup> (1965) داده پرت را مشاهده‌ی دورافتاده‌ای که انحراف آن از بقیه داده‌ها به طور محسوس و برجسته قابل مشاهده است تعریف کرد.

### 1-2-1 روش تعیین داده‌های پرت

با توجه به بررسی‌های وسیع در زمینه‌های داده‌های پرت و روش‌های مورد نظر و بحث‌های مرتبط در این زمینه، چندین روش عمومی برای تعیین داده‌های پرت ارائه می‌شود: استفاده از روش استانداردسازی ( $Z$ )، نمودار جعبه‌ای، نمودار پراکنش، آزمون گراب، نمودار هیستوگرام، و همچنین استفاده از چارک‌ها در تعیین داده‌های پرت کم و زیاد.

#### روش استانداردسازی:

مطابق قضیه چبی شف همه مشاهدات (تقریباً 99/8 درصد) در مجموعه داده‌ها دارای مقادیر ( $Z$  نمرات) کمتر از 3 هستند، که در فاصله  $(\bar{X} \pm 3S)$  قرار می‌گیرند که در آن  $\bar{X}$  میانگین و  $S$  انحراف معیار نمونه می‌باشد. بنابراین مشاهدات دارای مقدار قدرمطلق  $Z$  بزرگتر از 3، پرت خواهند بود.

<sup>11</sup> Empirical Distribution of the sample Radii

<sup>12</sup> - Barnet

<sup>13</sup> - Grubbs



### روش رسم نمودار جعبه‌ای:

روش‌های آماری زیادی برای رسم نمودار جعبه‌ای وجود دارند. بعد از رسم نمودار جعبه‌ای برای هر مجموعه داده، مشاهداتی که ما بین دیواره‌های داخلی و خارجی قرار می‌گیرند، مشکوک به پرت هستند. مشاهداتی که بیرون از این دیواره‌های خارجی قرار می‌گیرند، جزو داده‌های پرت هستند. برای داده‌های بزرگ این نمودار با استفاده از برنامه‌های کامپیوتری رسم می‌گردد.

### روش رسم نمودار پراکنش:

نمودار پراکنش از شیوه‌های ساده دیگر برای تعیین داده‌های پرت است، بدین صورت که بعد از رسم نمودار پراکنش بین دو متغیر وابسته  $Y$  و مستقل  $X$ ، اگر داده‌ای پرت باشد روی خط برازش شده برای داده قرار نگرفته و در فاصله‌ی دورتری از خط قرار می‌گیرند. لذا خط رگرسیونی رسم شده برای داده‌ها با احتساب مقدار داده‌پرت دارای برازش خوب نبوده و بدون در نظر گرفتن داده‌پرت، برازش مدل دقیق‌تری خواهد شد.

### روش آزمون گراب:

این آزمون توسط گراب در سال 1969 بیان گردید. این آزمون دارای مقدار آماره‌ای است که بر اساس اطلاعات موجود در نرم‌افزارهای مربوطه قابل محاسبه می‌باشد، لذا بعد از محاسبه مقدار آماره آزمون گراب، اگر این مقدار، از مقدار استاندارد تک تک داده‌ها بیشتر باشد آن داده از دید آزمون گراب پرت به حساب می‌آید. انجام آزمون گراب با استفاده از برنامه‌های کوییک‌کالک<sup>14</sup> و دیتاپلات<sup>15</sup>، صورت می‌گیرد که آن‌لین قابل استفاده است.

### روش نمودار هیستوگرام:

در داده‌های دارای توزیع نرمال اگر نمودار هیستوگرام رسم گردد، وجود یک داده پرت می‌تواند الگوی توزیع داده‌ها را تغییر داده و داده‌ها شکل یک توزیع چوله‌دار به خود بگیرند. همچنین وجود یک داده‌پرت باعث می‌گردد که آن داده در فاصله‌ی دورتری از بقیه داده‌ها قرار گرفته و بعد از رسم نمودار هیستوگرام به عنوان یک نقطه دور در طرف چپ و یا راست توزیع قرار بگیرد، و نوع شکل رسم شده را تحت تأثیر قرار دهد.

---

14- Quick calc

15- Data plot

## روش استفاده از چارک‌ها:

اگر فرض کنیم که داده‌ها دارای چارک اول  $Q_1$  و چارک سوم  $Q_3$  با دامنه‌ی میان‌چارکی  $IQR = Q_3 - Q_1$  باشند در این صورت مشاهداتی را که در نامساوی‌های  $X_i < Q_1 - 1/5IQR$  یا  $Q_1 + 1/5IQR < X_i$  قرار دارند جزو داده‌های پرت خفیف و مشاهداتی را که در نامساوی‌های  $X_i < Q_1 - 3IQR$  یا  $Q_1 + 3IQR < X_i$  قرار می‌گیرند جزو داده‌های پرت شدید هستند.

## 3-1 داده‌های آلوده

اغلب کسری از داده‌ها با کیفیت پایین، فاسد یا معیوب جمع‌آوری می‌شوند که به این نوع داده‌ها، داده‌های آلوده<sup>16</sup> گویند. گاهی این داده‌ها تأثیر چندانی در استنباط‌هایمان ندارند، اما گاهی نه‌تنها نقش بسزایی دارند بلکه تحلیل‌ها و استنباط‌هایمان را دچار اختلال می‌کنند بنابراین باید روش و ابزاری برای محدود کردن این تأثیرات در نظر گرفت. (برای درک تفاوت داده‌آلوده و داده‌پرت فرض کنید یک نمونه 20 تایی از داده‌های دو بعدی نرمال داشته باشید، چنانچه 2 یا 3 داده به جای توزیع نرمال دارای توزیع دیگری باشند و در این داده‌ها قرار بگیرند حکم داده‌آلوده دارند. اما داده‌پرت داده‌ای است که خیلی دورتر و یا به عبارتی خیلی پرت‌تر از داده‌های داده شده قرار می‌گیرد).

در این رساله تأثیر داده‌های پرت و داده‌های آلوده روی میانه‌های فاصله‌ای، بیضی‌گون و اوجا را مورد بررسی قرار داده و ابزارها و پیشنهادهایی برای حل این نابسامانی‌ها ارائه خواهیم کرد.

## 4-1 نقاط فروریزش

برای نمونه‌ی مرتب شده‌ی  $X_{(1)} < \dots < X_{(n)}$  با اندازه‌ی  $n$  فرض کنید  $T_n$  آماره‌ی مربوط به این نمونه باشد.  $T_n$  به ازای  $0 \leq b \leq 1$ ، دارای نقطه‌ی فروریزش<sup>17</sup>  $b$ ، است اگر برای هر  $\varepsilon > 0$ ، داشته باشیم:

$$\lim_{x_{((1-b)n)} \rightarrow \infty} T_n = \infty, \quad \lim_{x_{((1-b-\varepsilon)n)} \rightarrow \infty} T_n < \infty$$

توجه داشته باشید که این شرایط در مورد نقاط نمونه‌ای  $X_{(1)}, \dots, X_{(n)}$  است، نه متغیرهای  $X_{(1)} < \dots < X_{(n)}$ .

<sup>16</sup> -Contamination datas

<sup>17</sup> -Breakdown points

الف) اگر  $T_n = \bar{x}_n$  میانگین نمونه‌ای باشد، می‌توان نشان داد که  $b = 0$  است. برای این منظور فقط شرایط اول را در نظر می‌گیریم. بنابراین برای  $b = 0$ ، شرط اول به  $x(n)$  وابسته است و داریم:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{\sum_{i=1}^n x(i)}{n} = \frac{\sum_{i=1}^{n-1} x(i)}{n} + \frac{x(n)}{n}$$

هنگامی که  $x(n) \rightarrow \infty$ ، اولین خط کسری (نسبت بیان شده در بالا) ثابت باقی می‌ماند و دومین نسبت به سمت  $\infty$  میل می‌کند، بنابراین برای  $b = 0$ ،  $\lim T_n = \infty$  است. در نتیجه مقدار نقطه فروریزش برای میانگین نمونه‌ای صفر است. ب) اگر  $T_n = M_n$  میانه‌ی نمونه‌ای باشد، آنگاه  $b = \frac{1}{2}$  است. بنابراین برای  $n=2m$  (های زوج) داریم:

$$M_n = \frac{1}{2}(x_{(m)} + x_{(m+1)})$$

برای  $b = \frac{1}{2}$  اولین شرط عبارت است از:

$$x\left(\frac{n}{2}\right) = x(m)$$

$$, x(m) \rightarrow \infty$$

$$\Rightarrow x(m+1) \rightarrow \infty$$

بنابراین

$$M_n \rightarrow \infty$$

$$\Rightarrow \lim T_n = \infty$$

حال اگر  $n$  فرد باشد. داریم

$$M_n = x\left(\frac{(n+1)}{2}\right)$$

برای  $b = \frac{1}{2}$ ، از آنجا که  $\frac{(n+1)}{2} > \frac{(n)}{2}$  است، بنابراین

$$M_n = x\left(\frac{(n+1)}{2}\right) \rightarrow \infty$$

برای بررسی کردن شرط دوم  $\varepsilon > 0$  را ثابت قرار داده، در نتیجه

$$(1 - (b - \varepsilon))n = \left(\frac{n}{2}\right) + \varepsilon n$$

حال برای  $n > \frac{1}{\varepsilon}$  داریم:

$$\left(\frac{n}{2}\right) + \varepsilon n > \frac{n}{2} + 1 \quad \text{و} \quad \left(\frac{n}{2}\right) + \varepsilon n > \frac{n+1}{2}$$

بنابراین چه مقدار  $n$  زوج باشد و چه مقدار آن فرد باشد، آماره‌های ترتیبی‌ای که به  $\infty$  میل

می‌کنند مقدار حد  $M_n$  است. بنابراین نقطه فروریزش برای میانه  $\frac{1}{2}$  است.

توجه داشته باشید که نقطه فروریزش نمی‌تواند از 50% تجاوز کند، در صورت رخ دادن چنین

حالتی می‌گوییم داده‌ها، آلوده شده‌اند. بنابراین ماکسیم نقطه فروریزش 0/5 است. برای مثال،

میانه دارای یک نقطه فروریزش 0/5 است و میانگین پیراسته  $q$  ام، دارای یک نقطه فروریزش  $q$

است. آماره‌هایی با نقاط فروریزش زیاد آماره‌های استوار نامیده می‌شوند. (کسلا و برگر 1900)