

برآورد استوار کمینه دترمینان ماتریس کواریانس

اباذر خلجی، زهرا شاه محمدی
گروه آمار، دانشگاه اصفهان

چکیده

نقاط دورافتاده مشاهداتی هستند که از الگوی اکثر داده‌ها پیروی نمی‌کنند. به همین دلیل باعث تخریب برآوردهای بردار مکانی و ماتریس پراکندگی می‌شوند. در این مقاله روش کمینه دترمینان ماتریس کواریانس برای برآورد بردار مکانی و ماتریس پراکندگی معرفی می‌شود. این روش از بین تمام زیرنمونه‌های h نقطه‌ای از n نقطه‌ای، زیر نمونه‌ای را که دارای حداقل دترمینان ماتریس کواریانس است، اساس محاسبه برآورد بردار مکان و ماتریس پراکندگی قرار می‌دهد. برآوردهای بدست آمده از روش کمینه دترمینان ماتریس کواریانس در مقابل داده‌های دورافتاده مقاوم هستند و علاوه بر این دارای خاصیت هم‌پایایی آفین و بیشینه فروریزش نیز می‌باشند. مشکل عمده این روش زمان بر بودن و حجم بالای محاسبات آن می‌باشد، زیرا تعداد زیرنمونه‌های مورد نیاز گاهی اوقات بسیار زیاد می‌شود. برای برطرف نمودن این مشکل الگوریتم‌های سریع معرفی می‌شود که با حداقل باز نمونه‌گیری به نتیجه مطلوب می‌رسند.

واژه‌های کلیدی: برآورد استوار، بیشینه فروریزش، هم‌پایایی آفین.

۱ مقدمه

سه متغیره نیز اگر نمودارهای پراکنش را در فضاها دو یا سه بعدی ترسیم نماییم، احتمال شناسایی نقاط دورافتاده وجود دارد ولی لزوماً همیشه این کار شدنی نیست. در شناسایی نقاط دورافتاده زمانی که چند نقطه دورافتاده در مجموع داده‌ها وجود داشته باشد به علت درون‌آوری^۱ شناسایی دیگر امکان‌پذیر نیست. پس مناسب‌تر است از برآوردهای استوار^۲ بردار مکانی و ماتریس پراکندگی برای شناسایی نقاط دورافتاده بهره جست. در این مقاله علاوه بر معرفی روش استوار کمینه

وجود نقاط دورافتاده در نمونه‌های یک یا چندمتغیره باعث تخریب برآوردهای اعم از برآوردگر (بردار) مکان و (ماتریس) پراکندگی می‌شود. گاهی اوقات در نمونه‌های یک بعدی از میانه و گاهی نیز از میانگینی که با حذف درصدی از نقاط حدود بالایی و پایینی به دست آمده به عنوان برآورد مکان استفاده می‌شود و سپس با استفاده از همان نمونه تعدیل شده برآورد پراکندگی به دست می‌آید. در حالت تک متغیره امکان شناسایی نقاط دورافتاده از روی نمودار پراکنش داده‌ها نیز وجود دارد. در حالت دو و

^۱Masking

^۲Robust

بلکه بیش از اندازه زیاد و ناممکن نبودن است.

۲ روش کمینه دترمینان ماتریس کوواریانس

تعریف ۱.۲. اگر n مشاهده p متغیره داشته باشیم، در بین تمام زیرنمونه‌های h نقطه‌ای ($p < h < n$)، بردار میانگین و ماتریس کوواریانس زیرنمونه‌ای که دترمینان ماتریس کوواریانس آن کمینه باشد، برآورد MCD به ترتیب برای μ و Σ می‌باشد [۷].

کمیت h حجم زیرنمونه انتخابی برای محاسبه برآوردگرهای استوار $\bar{\mu}$ و $\bar{\Sigma}$ می‌باشد و نشان دهنده h مشاهده بالقوه خوب است که در محاسبه برآوردگرهای استوار شرکت دارند. اما $n-h$ مشاهده باقیمانده که باعث افزایش واریانس کل یا واریانس تعمیم یافته می‌شوند در محاسبه برآوردگرهای استوار شرکت داده نمی‌شوند و به این ترتیب برآوردگرهای استوار در برابر نقاط دورافتاده مقاوم می‌شوند. داوینز [۲] شرح داده است که با قرار دادن $h = \lfloor \frac{n+p+1}{4} \rfloor$ برآوردگرهای استوار بهبود می‌یابند. لذا ما نیز در محاسبات و مثال‌ها همین مقدار از h را در نظر گرفته‌ایم. در این جا مقصود از علامت [.] جزء صحیح است به طوری که h یک عدد صحیح مثبت است. در عمل h را می‌توان هر عدد صحیح در فاصله $\lfloor \frac{n+p+1}{4} \rfloor \leq h \leq n$ انتخاب نمود. یکی دیگر از انتخاب‌های مرسوم $h = 0.75n$ می‌باشد که منجر به افزایش کارایی برآوردگرها در مقابل کاهش نقطه فروریزش آن‌ها می‌شود. اگر $h = n$ آنگاه برآورد MCD مکانی همان بردار میانگین معمولی داده‌ها و برآورد MCD پراکندگی همان برآورد ماتریس کوواریانس نمونه است. در واقع با افزایش h برآوردگرهای

دترمینان ماتریس کوواریانس^۳ (MCD) برای برآورد μ و Σ الگوریتم‌های سریع برای محاسبه این روش نیز ارائه شده است. این روش به دنبال برآورد استواری است که تحت تاثیر نقاط دورافتاده نباشد. این روش نخستین بار توسط روسو [۷] معرفی شده است.

در پاراگراف قبل به استواری یک برآوردگر اشاره شد که در این جا از معیاری ساده و شهودی برای تعریف واژه استوار استفاده می‌کنیم. گوییم یک روش برای برآورد یا آزمون استوار است اگر با حذف یا اصلاح درصد کوچکی از مجموعه داده‌ها نتایج برآورد یا آزمون تغییر قابل توجهی نکنند. علاوه بر این‌ها برآوردگرهای MCD دارای خواص زیر نیز هستند:

۱. دارای نقطه فروریزش^۴ 50% درصدی می‌باشند. نقطه فروریزش یک برآوردگر حداقل نسبتی از داده‌ها است که آلوده بودن آنها می‌تواند به برآوردی نامعقول منتهی گردد [۹].

۲. هم‌پایای آفین^۵ هستند، به این معنی که تغییر مقیاس اندازه‌گیری بر خصوصیت برآوردگرها اثر نمی‌گذارد [۱۰].

۳. کارایی بالایی دارند، یعنی واریانس تعمیم یافته یا واریانس کل این برآوردگرها نسبت به سایر برآوردگرها کمتر است. معمولاً رابطه معکوسی بین فروریزش و کارایی برقرار است به این شکل که افزایش فروریختگی با کاهش کارایی همراه است.

۴. در زمانی معقولی قابل محاسبه هستند، در این جا مقصود از زمانی معقول، کوتاه بودن زمان نیست

^۳Minimum Covariance Determinant

^۴Breakdown Point

^۵Affine Equivariance

روش MVE صفر است.

۲. روسو و فن درایسن [۸] نشان دادند که فاصله‌های استواری که بر پایه روش MCD محاسبه می‌شوند بسیار دقیق‌تر از فاصله‌های استوار محاسبه شده بر پایه روش MVE می‌باشند و به همین دلیل بهتر می‌توانند نقاط پرت را شناسایی کنند.

مثال ساده زیر روشن‌گر بسیاری از ایده‌های مربوط به روش MCD می‌باشد.

مثال ۱.۲. ماتریس داده‌های زیر را در نظر بگیرید.

$$x = \begin{bmatrix} 4 & 15 & 6 & 12 & 5 \\ 13 & 25 & 12 & 15 & 17 \end{bmatrix}^T$$

در اینجا داریم، $n = 5$ و $p = 2$. بنابراین h برابر است با:

$$h = \left\lfloor \frac{n+p+1}{2} \right\rfloor = 4.$$

تعداد کل زیرنمونه‌های چهارتایی برابر ۵ = $\binom{5}{4}$ می‌باشد. این پنج زیر نمونه 2×4 در زیر آورده شده‌اند:

$$y_1 = \begin{bmatrix} 15 & 25 \\ 6 & 12 \\ 12 & 15 \\ 5 & 17 \end{bmatrix}, \quad y_2 = \begin{bmatrix} 4 & 13 \\ 6 & 12 \\ 12 & 15 \\ 5 & 17 \end{bmatrix},$$

استوار به برآوردهای معمولی نزدیک می‌شوند که منجر به کاهش نقطه فروریزش می‌شود. تعیین مقدار h بستگی به میزان آلودگی داده‌ها و هدف تحلیل‌گر دارد. بنابراین برای محاسبه برآوردهای کمینه دترمینان برای μ و Σ لازم است که تمام N زیرنمونه ممکن تعیین شود که در آن N برابر است با:

$$N = \binom{n}{h} = \frac{n!}{h!(n-h)!}.$$

سپس زیر نمونه‌ای که برای آن دترمینان ماتریس کوواریانس کمینه است، برای محاسبه $\bar{\mu}$ و $\bar{\Sigma}$ مورد استفاده قرار گیرد. واضح است که اگر n بزرگ باشد، تعداد زیرنمونه‌های h عضوی نیز افزایش یافته و به تبعیت از آن حجم محاسبات نیز بسیار زیاد می‌شود. اغلب تلاش‌های دو دهه اخیر مربوط به کم کردن حجم محاسبه‌ها از طریق روش‌های متنوع زیر نمونه گیری تصادفی بوده است. امروزه به دلایل مختلفی از روش MCD به جای روش بیضی‌گون کمینه حجم^۶ (MVE) استفاده می‌شود که بعضی از این دلایل عبارتند از:

۱. روش MCD دارای کارایی بهتری است چون برآوردهای این روش دارای توزیع جانبی نرمال هستند ولی در مقابل روش MVE دارای نرخ همگرایی پایینی در توزیع است [۷]. برای مثال کارایی جانبی ماتریس پراکندگی روش MCD با انتخاب $h = 0.75n$ و $p = 10$ برابر ۴۴ درصد می‌باشد و کارایی ماتریس کوواریانس باز وزن دهی شده با وزن‌هایی که از روش MCD به دست می‌آیند به ۸۳ درصد می‌رسد، در حالی که این مقدار برای

^۶Minimum Volume Ellipsoid

در این جا برآورد MCD عیناً برابر برآورد MVE می باشد. هر چند در اغلب کاربردها این دو برآورد برابر هم هستند ولی همیشه این دو برآورد یکسان نیستند. اگر الگوریتم محاسبه ای فوق را بجای $h = 4$ برای $h = 3$

اجرا کنیم زیر ماتریس $\begin{bmatrix} 4 & 6 & 12 \\ 13 & 12 & 15 \end{bmatrix}^T$ به عنوان زیر ماتریس h نقطه ای برای برآورد استوار حاصل می شود.

$$y_3 = \begin{bmatrix} 4 & 13 \\ 15 & 25 \\ 12 & 15 \\ 5 & 17 \end{bmatrix}, \quad y_4 = \begin{bmatrix} 4 & 13 \\ 15 & 25 \\ 6 & 12 \\ 5 & 17 \end{bmatrix},$$

$$y_5 = \begin{bmatrix} 4 & 13 \\ 15 & 25 \\ 6 & 12 \\ 12 & 15 \end{bmatrix},$$

ماتریس کواریانس داده های X تحت تبدیل $AX + b$ به صورت ASA^T است که در آن S ماتریس کواریانس داده ها می باشد. طبق روسو [۷] دترمینان ماتریس کواریانس داده های تبدیل شده برابر است با:

که دترمینان ماتریس کواریانس این زیرنمونه ها برابر است با

$$\det(ASA^T) = \det^2(A) \times \det(S),$$

$$\det(S_{y_1}) = 356/39, \quad \det(S_{y_2}) = 61/5,$$

چون $\det^2(A)$ ثابت است، می توان نتیجه گرفت برآوردگرهای MCD هم پایای آفین هستند. علاوه بر این ها روسو [۷] نشان داده است که برای هر اندازه نمونه n و هر بعد داده های p ، نقطه فروریزش برآوردگرهای MCD برابر است با:

$$\det(S_{y_3}) = 366, \quad \det(S_{y_4}) = 158/17,$$

$$\det(S_{y_5}) = 262/39,$$

که S_{y_i} ها نشان دهنده ماتریس کواریانس زیرنمونه i ام می باشد. با توجه به این که زیرنمونه y_2 دارای کمترین واریانس تعمیم یافته است، همین زیرنمونه را برای محاسبه برآورد MCD برای μ و Σ انتخاب می کنیم. بنابراین بردار میانگین و ماتریس کواریانس زیر نمونه دوم برآوردگر MCD به ترتیب برای μ و Σ می باشند. یعنی،

$$\frac{\lfloor \frac{n}{p} \rfloor - p + 1}{n}.$$

اگر n را به سمت بینهایت میل دهیم مقدار فروریزش به ۵۰ درصد همگرا می شود که برابر با فروریزش برآوردگرهای MVE می باشد.

$$\tilde{\mu} = \bar{y}_2 = \begin{bmatrix} 6/75 \\ 14/25 \end{bmatrix},$$

مثال ۲.۲. در این مثال به کمک شبیه سازی و با استفاده از نرم افزار R عملکرد روش MCD را در آلودگی های مختلف مورد بررسی قرار می دهیم. ابتدا ماتریس $X_{6 \times 3}$ را از توزیع $N_3(\mu, \Sigma)$ تولید می کنیم. که در آن μ و Σ برابر است با:

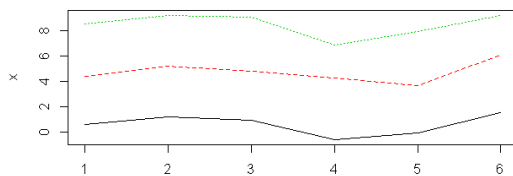
$$\tilde{\Sigma} = S_{y_2} = \begin{bmatrix} 12/9167 & 1/4167 \\ 1/4167 & 4/9167 \end{bmatrix}.$$

$$X_1 = \begin{bmatrix} 0/6 & 4/37 & 8/48 \\ 1/23 & 5/19 & 9/17 \\ 0/92 & 4/76 & 9/52 \\ -0/57 & 4/26 & 6/82 \\ -0/53 & 3/65 & 7/89 \\ 1/55 & 6/57 & 9/15 \end{bmatrix}.$$

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

یکی از مشاهدات (مشاهده چهارم) را با مشاهده‌ای که هر بار به روشی آلوده شده، جایگزین می‌کنیم. با توجه به ابعاد ماتریس مشاهدات $h = 5$ و تعداد کل زیر نمونه‌های 3×5 برابر شش می‌باشد. شش زیر نمونه y_1, y_2, y_3, y_4, y_5 و y_6 به ترتیب از حذف $x_1^T, x_2^T, x_3^T, x_4^T, x_5^T$ و x_6^T از ماتریس مشاهدات x بدست می‌آیند که نشان دهنده سطر i ام ماتریس x می‌باشد. لازم به ذکر است که خروجی برنامه شامل بردار d و برآوردگرهای استوار مکانی و پراکندگی MCD می‌باشد. بردار d بیانگر مشاهدات انتخاب شده در زیر نمونه‌ای می‌باشد که کمترین دترمینان ماتریس کوواریانس را داراست.

نمودار دنباله‌ای مشاهدات در شکل ۱ ملاحظه می‌شود. با اجرای برنامه آرایه‌های زیر حاصل



شکل ۱: نمودار پراکنش برای آلودگی میانگین

می‌شود:

$$\underline{\tilde{\mu}} = \begin{bmatrix} 0/85 \\ 4/81 \\ 8/74 \end{bmatrix}, \quad \tilde{\Sigma} = \begin{bmatrix} 0/37 & 0/54 & 0/32 \\ 0/54 & 0/82 & 0/44 \\ 0/32 & 0/44 & 0/31 \end{bmatrix},$$

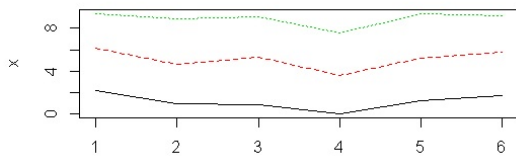
$$\underline{d} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 5 \\ 6 \end{bmatrix}.$$

با توجه به بردار d زیر نمونه‌ای که مبنای محاسبه برآورد $\underline{\mu}$ و Σ قرار گرفته شامل مشاهده چهارم که آلوده شده بود، نیست. می‌توان نتیجه گرفت روش MCD به درستی قادر به شناسایی نقاط پرت می‌باشد زیرا از نمونه‌ای که فاقد

۱. آلوده کردن میانگین: در این شبیه‌سازی بردار میانگین متناظر با چهارمین مشاهده را به صورت زیر آلوده می‌کنیم،

$$\underline{\mu}_{outlier} = \underline{\mu} + (0, 1, 0)^T.$$

اما ماتریس کوواریانس را آلوده نمی‌کنیم. به عبارت دیگر برای مشاهده چهارم، میانگین دومین متغیر را به اندازه یک انحراف معیار افزایش می‌دهیم. به این ترتیب ماتریس مشاهدات زیر حاصل می‌گردد.



شکل ۲: نمودار پراکنش برای آلودگی ماتریس کوواریانس

مشاهده آلوده است در برآورد μ و Σ استفاده می‌کند.

۲. آلوده کردن ماتریس کوواریانس: در شبیه‌سازی دیگری مشاهدات را از توزیع $N_3(\mu, \Sigma)$ که در بالا شرح داده شد تولید کردیم، اما مشاهده آلوده (مشاهده چهارم) را بر اساس ماتریس کوواریانس زیر تولید نمودیم،

$$\tilde{\Sigma}_{(MCD)} = \begin{bmatrix} 0.29 & 0.27 & 0.07 \\ 0.27 & 0.35 & 0.09 \\ 0.07 & 0.09 & 0.05 \end{bmatrix},$$

$$\Sigma_{out} = \begin{bmatrix} 1 & -0.9 & 0.9 \\ -0.9 & 2 & 0.9 \\ 0.9 & 0.9 & 1 \end{bmatrix},$$

ولی بردار میانگین را تغییر ندادیم. به این ترتیب ماتریس داده‌های زیر حاصل گردید،

$$\underline{d} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 5 \\ 6 \end{bmatrix}.$$

$$X_2 = \begin{bmatrix} 2.14 & 6.18 & 9.32 \\ 0.98 & 4.62 & 8.93 \\ 0.83 & 5.32 & 9.04 \\ 0.02 & 3.57 & 7.50 \\ 1.26 & 5.21 & 9.35 \\ 1.70 & 5.75 & 9.15 \end{bmatrix}.$$

با توجه به بردار \underline{d} زیرنمونه‌ای که مبنای محاسبه برآورد μ و Σ قرار گرفته شامل مشاهده چهارم که آلوده شده بود، نیست. می‌توان نتیجه گرفت روش MCD در این نوع آلودگی نیز به درستی قادر به شناسایی نقاط پرت می‌باشد زیرا از نمونه‌ای که فاقد مشاهده آلوده است در برآورد μ و Σ استفاده می‌کند. لازم به ذکر است که با ۱۰۰ بار تکرار شبیه سازی، روش MCD در ۹۱ درصد مواقع قادر به ارائه برآوردهای استوار بود. تعداد دفعات باقیمانده ناشی از انحرافات شبیه سازی و خطاهای روش می‌باشد.

درواقع واریانس متغیر دوم را دو برابر و همبستگی بین متغیر اول و دوم را معکوس نموده‌ایم. اگر داده‌ها را در دو بعد تصور کنیم نمودار پراکنش داده‌ها به شکل بیضی خواهد بود که این نوع آلودگی مشاهده پرت را بر روی خطی موازی محور اصلی بیضی اما با فاصله بیشتری قرار می‌دهد که باعث می‌شود مشاهده آلوده دورتر از سایر مشاهدات قرار بگیرد [۶]. نمودار دنباله‌ای ماتریس داده‌های فوق در شکل ۲ ملاحظه می‌شود.

۳ الگوریتم پرسرعت محاسبه

MCD

با اجرای برنامه داریم،

برآوردگر MCD از طریق زیر نمونه‌ای که دارای کمترین دترمینان کوواریانس است محاسبه می‌شود و تعداد زیر نمونه‌های ممکن با توجه به مقادیر حجم نمونه n و بعد

$$\tilde{\mu} = \begin{bmatrix} 1.38 \\ 5.42 \\ 9.14 \end{bmatrix},$$

$$\{d(i); i \in H_{new}\} = \{d_{(1)} \leq d_{(2)} \leq \dots \leq d_{(h)}\},$$

به طوری که $d_{(1)} \leq d_{(2)} \leq \dots \leq d_{(n)}$ فاصله‌های مرتب شده‌اند. اگر $\underline{\mu}_{new}$ و S_{new} را بر اساس H_{new} محاسبه کنیم، آنگاه داریم،

$$\det(S_{new}) \leq \det(S).$$

شرط تساوی برقرار است اگر و فقط اگر $\underline{\mu}_{new} = \underline{\mu}$ و $S_{new} = S$ باشد.

اثبات. برای اثبات به روسو و فن‌درایسن [۸] مراجعه کنید.

الگوریتم:

اکنون به شرح مراحل اجرای الگوریتم سریع می‌پردازیم:

۱. انتخاب یک زیر نمونه که این انتخاب به دو صورت امکان پذیر است.

- انتخاب یک زیر نمونه h تایی تحت عنوان H_{old} به صورت تصادفی تا جایی که $\det(S_{old}) > 0$.

- انتخاب یک زیر نمونه $p + 1$ تایی y به صورت تصادفی و محاسبه $\underline{\mu}_0 = \text{mean}(y)$ و $S_0 = \text{Cov}(y)$. اگر $\det(S_0) = 0$ آنگاه y را توسط اضافه نمودن یک مشاهده تصادفی دیگر بسط داده و آنقدر اضافه نمودن مشاهده را ادامه می‌دهیم تا زمانی که $\det(S_0) > 0$ شود. آنگاه فاصله‌های زیر را محاسبه می‌نماییم،

$$d_i^0 = (\underline{x}_i - \underline{\mu}_0)^T S_0^{-1} (\underline{x}_i - \underline{\mu}_0), \quad i = 1, \dots, n,$$

متغیر p تعیین می‌شود. از این رو در نمونه‌های با حجم بزرگ اغلب با کثرت زیر نمونه‌ها مواجه هستیم که باعث افزایش مدت زمان مورد نیاز برای محاسبه برآوردگرهای MCD می‌شود. بنابراین ناچار به استفاده از الگوریتم‌هایی هستیم که بتوانیم برآوردگرها را با حداکثر دقت و در حداقل زمان ممکن محاسبه کنیم. تمامی الگوریتم‌هایی که در رابطه با محاسبه برآوردگرهای MCD وجود دارند به نوعی در پی کوتاهترین زمان ممکن از طریق خلاصه سازی در زیر نمونه‌ها برای به دست آوردن کمترین دترمینان ماتریس کوواریانس می‌باشند. از این رو روسو و فن‌درایسن [۸] الگوریتم مؤثری را با استفاده از قضیه‌ای که در زیر به آن می‌پردازیم، برای محاسبه برآوردگرهای MCD پیشنهاد کردند و از آن به عنوان الگوریتم سریع نام بردند.

قضیه ۱.۳. فرض کنید $X_{n \times p} = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n]^T$ ، ماتریس مشاهدات n مشاهده p متغیره باشد و H یک زیر نمونه h عضوی تصادفی از بین کلیه زیر نمونه‌های ممکن باشد. $\underline{\mu}$ و S به ترتیب بردار میانگین و ماتریس کوواریانس این زیر نمونه باشند.

$$\underline{\mu} = \frac{1}{h} \sum_{i \in H} x_i, \quad S = \frac{1}{h} \sum_{i \in H} (x_i - \underline{\mu})(x_i - \underline{\mu})^T.$$

اگر $\det(S) \neq 0$ ، فاصله زیر را تعریف می‌کنیم،

$$d(i) = \sqrt{(\underline{x}_i - \underline{\mu})^T S^{-1} (\underline{x}_i - \underline{\mu})}, \quad i = 1, 2, \dots, n.$$

اکنون H_{new} زیر نمونه h عضوی جدید را طوری انتخاب می‌کنیم که داشته باشیم،

زیر نمونه گزارش می‌کنیم. در غیر این صورت به مرحله دوم بازگشته و با تعویض S_{new} و μ_{new} به ترتیب با S_{old} و μ_{old} الگوریتم را تکرار می‌کنیم.

۷. مرحله دوم تا ششم را با توجه به هزینه‌ها و زمان در دسترس آنقدر ادامه می‌دهیم تا به جواب همگرا برسیم.

همان طور که مشاهده می‌شود این الگوریتم تمامی زیرنمونه‌ها را مورد بررسی قرار نمی‌دهد. بلکه تا آن جایی که ممکن است از تعداد زیرنمونه‌ها می‌کاهد. با توجه به اینکه این الگوریتم فقط با زیرنمونه‌ای که مقدار $\sum_{i=1}^h d(i)$ برای آن کمترین مقدار باشد کار را ادامه می‌دهد، در مقایسه با سایر الگوریتم‌ها از خلاصه سازی و در نتیجه از سرعت بیشتری برخوردار است.

در رابطه با الگوریتم فوق نکات زیر قابل توجه می‌باشند:

- انتخاب عدد کوچک $p + 2$ به عنوان اندازه نمونه تصادفی با این انگیزه صورت می‌گیرد که احتمال قرار گرفتن یک نقطه دورافتاده در نمونه به حداقل خود برسد [۳].

- منطقی به نظر می‌رسد که تکرار اجرای الگوریتم‌های فوق تا رسیدن به زیرمجموعه‌ای که شامل اکثر مشاهدات بالقوه خوب است ادامه یابد، زیرا این زیرمجموعه دارای کمینه دترمینان ماتریس کوواریانس است و اساس محاسبه برآورد استوار بردار میانگین و ماتریس کوواریانس قرار می‌گیرد. بنابراین برای افزایش دقت الگوریتم‌ها و تضمین رسیدن آن‌ها به یک مینیمم مطلق باید این الگوریتم‌ها در تعداد تکرار بالا انجام شوند، در تکرارهای کم

و سپس آن‌ها را به صورت صعودی به ترتیب زیر مرتب کرده،

$$(d_{\circ})_{(1)} \leq (d_{\circ})_{(2)} \leq \dots \leq (d_{\circ})_{(n)},$$

و در نهایت مشاهده‌های متناظر با h فاصله اول را به صورت زیر در زیر نمونه قرار می‌دهیم،

$$H_{old} = \{\underline{x}_{(1)}, \underline{x}_{(2)}, \dots, \underline{x}_{(h)}\}.$$

۲. محاسبه S_{old} و μ_{old} به همراه کلیه فاصله‌ها، یعنی،

$$d_{old}(i) = \sqrt{(\underline{x}_i - \mu_{old})^T S_{old}^{-1} (\underline{x}_i - \mu_{old})}.$$

۳. فاصله‌ها را مرتب کرده و به صورت زیر نمایش می‌دهیم،

$$(d_{old})_{(1)} \leq (d_{old})_{(2)} \leq \dots \leq (d_{old})_{(n)}.$$

۴. مشاهده‌های متناظر با h فاصله اول را در زیر نمونه H_{new} قرار می‌دهیم،

$$H_{new} = \{\underline{x}_{(1)}, \underline{x}_{(2)}, \dots, \underline{x}_{(h)}\}.$$

۵. S_{new} و μ_{new} را به صورت زیر محاسبه می‌کنیم،

$$\mu_{new} = \text{mean}(H_{new}), \quad S_{new} = \text{Cov}(H_{new}).$$

۶. اگر $\det(S_{old}) = \det(S_{new})$ یا $\det(S_{new}) = 0$ آن‌گاه الگوریتم را متوقف و H_{new} را به عنوان بهترین

مقایسه شده‌اند. همان گونه ملاحظه می‌شود در تمام حالت‌های روش استوار برآورد واریانس هر یک از متغیرها از برآورد واریانس روش ML کوچکتر است. واریانس تعمیم یافته که بر اساس روش استوار بدست آمده نیز از واریانس تعمیم یافته بدست آمده از روش ML کمتر است. واریانس متغیرها در جدول ۱ در داخل پرانتز قرار گرفته‌اند.

۵ دستور برنامه R برای محاسبه برآوردگرهای MVE و MCD

نرم‌افزار R با استفاده از بسته ($MASS$) و دستور $cov.rob()$ ، قادر به محاسبه برآوردگرهای استوار MVE و MCD می‌باشد. دستور به صورت زیر می‌باشد،

```
cov.rob(x, cor = FALSE,
        quantile.used = floor((n + p + 1)/2),
        method = c("mve", "mcd")),
```

که در آن x ماتریس داده‌ها، عملگر منطقی cor با اختیار دو مقدار $TRUE$ و $FALSE$ در خروجی برنامه به ترتیب اجازه نمایش و عدم نمایش ماتریس همبستگی را می‌دهد. با استفاده از کد $quantile.used$ می‌توان حجم زیرنمونه مورد استفاده را مشخص کرد و با استفاده از کد $method$ یکی از دو روش MCD یا MVE انتخاب می‌شود. برای جزئیات بیشتر می‌توان به $help$ برنامه R مراجعه کرد.

۶ نتیجه گیری

با توجه به مباحثی که عنوان شد اگر مجموعه داده‌ها آلوده باشد، برآوردگرهای کلاسیک تحت تاثیر نقاط دور افتاده

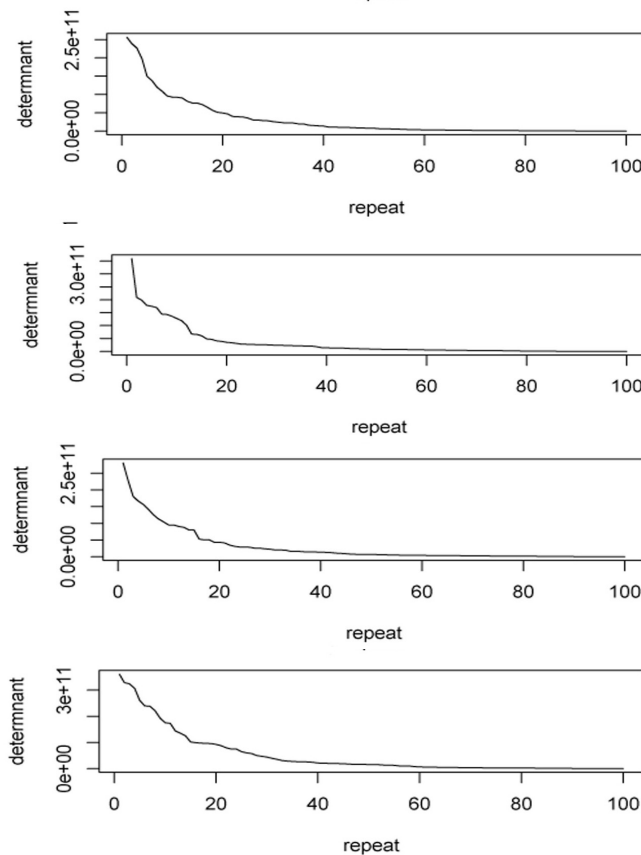
مینیم‌های موضعی به عنوان تقریبی مناسب از مینیمم مطلق استفاده می‌گردد.

۴ بررسی داده‌های واقعی

کمپبل [۱] مجموعه داده‌هایی متشکل از $n = 38$ مشاهده در $p = 5$ بعد را ارائه کرده است. در متون آماری نشان داده شده که ۱۳ مورد از این مشاهدات به شماره‌های ۷ تا ۱۱ و ۳۱ تا ۳۸ نقاط دورافتاده هستند. (در این مورد مراجع [۵] و [۴] را ملاحظه فرمایید.) لازم به ذکر است که روك و ودروف [۱۰] نیز به همین نتایج رسیدند با این تفاوت که مشاهدات ۷، ۱۰، ۱۱ و ۳۱ را مشکوک به دورافتاده بودن و سایر مشاهدات باقیمانده را نقاط دورافتاده شدید^۷ نامیدند. ما از این داده‌ها برای ارائه مثالی عملی از برآوردگر MCD بهره برده‌ایم. نکته بسیار جالب توجه این است که با تعداد تکرار نه چندان زیاد (مثلاً ۵۰ تکرار)، رویت حتی یک مورد از ۱۳ مورد نقطه دورافتاده در باز نمونه نهایی پیشامدی بسیار نادر است. عملاً با افزایش تعداد تکرارها می‌توان احتمال تحقق یافتن این پیشامد را به صفر نزدیک کرد. توجه کنید که از دیدگاه نظری برآوردگرهای MCD یکتا نیستند ولی احتمال یافتن دو زیر مجموعه h عضوی که دترمینان ماتریس کوواریانس آنها برابر کمینه دترمینان ماتریس کوواریانس باشد بسیار ناچیز است. در شکل ۳ در ۴ مورد با ۱۰۰ تکرار بر اساس الگوریتم سریع روسو و فن‌درایسن دروند کاهش دترمینان ماتریس کوواریانس و سرعت همگرایی الگوریتم در روش MCD ملاحظه می‌شود.

در جدول ۱ برآوردهای ML با برآوردهای استوار

^۷Extrem Outlier



شکل ۳: روند کاهش حجم بر اساس الگوریتم سریع روسو و فن درایسن در ۱۰۰ تکرار

قرار گرفته و کارایی لازم را ندارند. اگر آلودگی ناشی از محاسبه می‌کند. یعنی اگر مجموعه داده‌ها در فضای چند خطاهای اندازه‌گیری و ثبتي باشد می‌توان پس از حذف بعدی به صورت یک توده تصور شود در برآورد از نمونه‌ای این نقاط از برآوردگرهای کلاسیک بهره جست اما به دلیل اینکه شناسایی نقاط دور افتاده چندمتغیره بسیار سخت‌تر می‌باشد محاسبه می‌شود.

مراجع

- [1] Campbell, N. A. (1989). Bushfire mapping using NOAA AVHRR data. Technical Report, CSIRO.
- [2] Davies, P. L. (1987). Asymptotic bi-

از روشی که در این مقاله مورد بحث قرار گرفته حتی زمانی که مشکوک به وجود نقاط دورافتاده در مجموعه داده‌ها هستیم مفید بوده و برآوردگرها (بردار میانگین و ماتریس پراکندگی) را بر اساس داده‌هایی که به هم نزدیکتر هستند بنا بر این نیازمند روش‌هایی هستیم که در برابر نقاط دور افتاده برآوردگرهای معقولي را ارائه دهند. از این رو

جدول ۱: مقایسه برآوردهای ML و استوار

روش ML					برآورد
۱۰۳/۶	۱۲۹/۱	۲۸۸/۶	۲۲۷/۸	۲۸۶/۶	میانگین
(۴۰۶)	۵۶۵	-۲۰۹۱	-۶۳۹	-۵۱۶	ماتریس کواریانس
۵۶۵	(۱۲۲۵)	-۳۲۵۸	-۱۱۸۴	-۹۴۳	
-۲۰۹۱	-۳۲۵۸	(۳۱۴۰۵)	۱۱۰۶۰	۹۰۲۱	
-۶۳۹	-۱۱۴۸	۱۱۰۶۰	(۴۱۰۳)	۳۳۴۰	
-۵۱۶	-۹۴۳	۹۰۲۱	۳۳۴۰	(۲۷۲۲)	
۱۱۹۵۲۰۰۰۰۰۰۰۰					واریانس تعمیم یافته
استوار					برآورد
۱۰۹/۴	۱۴۹/۶	۲۶۲/۲	۲۱۵/۳	۲۷۷/۱	میانگین
(۳۸۶)	۴۸۱	۲۰۰۰	-۱۶۶	-۲۸۹	ماتریس کواریانس
۴۸۱	(۶۰۴)	۲۴۶۶	-۱۹۱	-۳۴۶	
۲۰۰۰	۲۴۶۶	(۹۷۶۷)	-۱۰۱۲	-۱۶۶۹	
-۱۶۶	-۱۹۱	-۱۰۱۲	(۲۳۶)	۲۸۰	
-۲۸۹	-۳۴۶	-۱۶۶۹	۲۸۰	(۳۷۷)	
۱۲۳۶۹/۶۷					واریانس تعمیم یافته

(1995). The Behavior of the stahel-donoho robust multivariate location and shape. Journal of the American Statistical Association, 4, 51-67.

havior of s-stimators of multivariate location parameteres and dispersion matrices. The Annals of Statistics, 15, 1269-1292.

[6] Rencher, A. C. (2002). Methods of multivariate analysis (2nd ED.). New York: Wiley - interscience.

[3] Hawkins, D. M. and Olive, D. J. (2002). Inconsistency of resampling algorithms for high-breakdown regression estimators and new algorithm (with discussion). Journal of the American Statistical Association, 97, 136-159.

[7] Rousseeuw, P. J. (1985). Multivariate estimation with high breakdown point. In Mathematical Statistics and Application, VOL. B. eds. W. Grossmann, G. Pflug, I. vineze and W. Werz, Dordrecht :Reidel.

[4] Kosinski, A. S. (1999). A procedure for the detection of multivariate outlier. Computational Statistic and Data Analysis, 29, 145-161.

[8] Rousseeuw, P. J. and Van Driessen, K. (1999). A fast algorithm for the min-

[5] Maronna, R. A. and Yohai, V. J.

- imum covariance determinant estimator. *Technometrics*, 41(3), 212-223.
- [9] Sakata, S. and With, H. (1998). Breakdown point. *Encyclopedia of Statistical Sciences*, 2, 84-89.
- [10] Woodruff, D. L, and Rocke, D. M. (1996). Identification of outlier in multivariate data. *Journal of the American Statistical Association*, 91(435), 1047-1061.