

شناسایی نقاط دورافتاده چندمتغیره با استفاده از فواصل استوار

اباذر خلجی، زهرا شاه‌محمدی
گروه آمار، دانشگاه اصفهان

چکیده

برای شناسایی نقاط دورافتاده یک‌متغیره روش‌هایی وجود دارد که در حالت چندمتغیره کارآمد نیستند. علاوه بر این زمانی که چندین مشاهده دورافتاده در مجموعه داده‌های چندمتغیره وجود دارد، شناسایی نقاط دورافتاده چندمتغیره امری دشوار و نیازمند روش‌های محاسباتی می‌باشد. روش‌های شناسایی کلاسیک اغلب قادر به شناسایی این نقاط نیستند. زیرا بر پایه بردار میانگین و ماتریس کوواریانس آلوده به شناسایی نقاط دورافتاده می‌پردازند. به همین دلیل با مساله درون‌آوری مواجه می‌شوند. برای مقابله با درون‌آوری فواصل باید بر اساس برآوردهای استوار محاسبه شوند. در این مقاله از برآوردهای کمینه درمینان ماتریس کوواریانس و بیضی‌گون کمینه حجم که دارای بیشینه فروریختگی می‌باشند در محاسبه فواصل بهره‌جسته‌ایم.

واژه‌های کلیدی: نقاط دورافتاده، فاصله ما‌هالانویس، فاصله استوار.

۱ مقدمه

مشاهدات دارای خطاهای اندازه‌گیری و جمع‌آوری هستند (رنچر، [۷]). اگر آلودگی مشاهدات ناشی از خطاهای اندازه‌گیری و ثبتي بوده و امکان دسترسی مجدد برای اصلاح آن وجود نداشته باشد، باید آن‌ها را حذف کرد، اما در غیر این صورت چون حاوی اطلاعات در مورد نمونه هستند حذف آن‌ها جایز نیست.

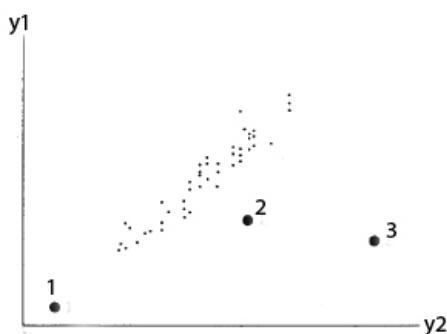
وجود نقاط دورافتاده در نمونه‌های یک یا چندمتغیره باعث تخریب برآوردها (بردار) مکانی و (ماتریس) پراکندگی می‌شوند، توزیع داده‌ها را تغییر داده و منجر به توزیع‌هایی با دم‌های سنگین می‌شوند، میزان تغییرات را افزایش داده و همبستگی بین متغیرها را در هم می‌ریزند. بنابراین در مدل سازی آماری و تحلیل داده‌ها بررسی

نقاط دورافتاده^۱ مقادیری در مجموعه داده‌ها هستند که غالباً از الگوی اکثر مشاهدات پیروی نمی‌کنند. هاوکینز [۴] نقاط دورافتاده را این‌گونه تعریف می‌کند:

«نقاط دورافتاده مشاهداتی هستند که به قدری دورتر از سایر مشاهدات قرار می‌گیرند که گویی با مکانیسم دیگری تولید شده‌اند.»

گاهی نقاط در اثر خطاهای اندازه‌گیری، ثبت نادرست یا برخی عوامل دیگر آلوده می‌شوند. بعضی از محققین عقیده دارند که در یک پژوهش معمولاً بیش از ۱۰ درصد

^۱Outliers



شکل ۱: نمایش سه نوع آلودگی در نمونه دومتغیره

تعریف ۱ یک بیضی‌گون در فضای p بعدی به صورت

$$(\underline{x} - \underline{b})^T A^{-1} (\underline{x} - \underline{b}) = d^2,$$

تعریف می‌شود، که در آن p نشان دهنده محور مختصات، A یک ماتریس $p \times p$ معین مثبت (و در کاربردهای آماری متقارن)، \underline{b} مرکز بیضی‌گون و d^2 یک عدد حقیقی مثبت می‌باشد.

تعریف ۲ اگر \bar{x} و S به ترتیب بردار میانگین و ماتریس کواریانس یک نمونه n تایی باشد آنگاه فاصله

$$MD = \sqrt{(\underline{x} - \bar{x})^T S^{-1} (\underline{x} - \bar{x})},$$

یک فاصله ماهالانوبیس^۲ برای \underline{x} است.

در واقع یک بیضی‌گون مشخص کننده مجموعه تمامی نقاطی از فضای p بعدی است که مربع فاصله ماهالانوبیسی آن نقاط تا مرکز \bar{x} برابر مقدار ثابت d می‌باشد.

وجود نقاط دورافتاده در مجموعه داده‌ها از اهمیت بالایی برخوردار است. در حالت یک‌بعدی (تک‌متغیره) امکان شناسایی نقاط دورافتاده از روی نمودار پراکنش داده‌ها وجود دارد و در حالت دو و سه‌متغیره نیز اگر نمودارهای پراکنش را در فضاهای دو یا سه‌بعدی ترسیم نماییم، احتمال شناسایی نقاط دورافتاده وجود دارد ولی لزوماً همیشه این کار شدنی نیست. بنابراین در داده‌های چندمتغیره نیازمند روش‌های دیگری می‌باشیم.

در داده‌های چندمتغیره مسئله شناسایی نقاط دورافتاده از اهمیت و دشواری‌های بیشتری در مقایسه با حالت یک‌متغیره برخوردار است، زیرا

۱. برای $p > 2$ نمودار پراکنش داده‌ها قادر به شناسایی نقاط دورافتاده نمی‌باشد.

۲. هر مشاهده به صورت یک بردار بوده و ممکن است یک خطای (اندازه‌گیری، ثبیتی و ...) بزرگ در یکی از اعضای بردار یا خطاهای کوچک در چندین عضو بردار رخ دهد.

۳. داده‌های چندمتغیره را نمی‌توان مانند داده‌های تک‌متغیره به صورت صعودی مرتب کرد و مقادیر کوچک و بزرگی را که در دو انتها قرار می‌گیرند مشخص نمود.

۴. نقاط دورافتاده چندمتغیره ممکن است به واسطه تغییر در میانگین، واریانس یا همبستگی ایجاد شده باشند. برای مثال در شکل ۱ که از کتاب رنچر [۷] اقتباس شده، مشاهده ۱ با تغییر اندک میانگین و واریانس هر دو متغیر y_1 و y_2 به وجود آمده است، مشاهده ۲ به واسطه تغییر در همبستگی بین متغیرها ایجاد شده و مشاهده ۳ با تغییر میانگین، واریانس و همبستگی بین متغیرها به وجود آمده است.

^۲Mahalanobis Distance

۲ شناسایی به وسیله مربع فاصله ماهالانوبیس

نتایج قابل اعتماد نیازمند مقدار برشی^۵ هستیم تا مقادیر MSD_i را با آن بسنجیم.

اگر بردار میانگین و ماتریس کوواریانس جامعه یعنی μ و Σ معلوم باشند آن گاه MSD_i دارای توزیع χ_p^2 (مربع کای با p درجه آزادی) می باشد و اگر از برآوردگرهای کلاسیک \bar{x} و S در محاسبه فاصله ماهالانوبیس استفاده کنیم، فاصله دارای توزیع مجانبی χ_p^2 خواهد بود، [۶]. بنابراین اگر مربع فاصله ماهالانوبیس مشاهده i ام از مقدار برشی $\chi_{\alpha, p}^2$ بزرگتر باشد آن گاه مشاهده i ام به عنوان مشاهده دورافتاده شناسایی می شود.

درواقع MSD_i فاصله مشاهده \bar{x}_i را تا مرکز همه داده ها نه مرکز انبوهه اصلی اندازه گیری می کند. به دلیل اینکه نقاط دور افتاده در محاسبه \bar{x} و S نیز حضور دارند دیگر با فاصله های بزرگ مواجه نخواهیم بود و لذا استفاده از معیار MSD درون آوری^۶ را در پی خواهد داشت. به عبارت دیگر زمانی که با چندین داده دورافتاده مواجه هستیم این مشاهدات آلوده بر میانگین تاثیر می گذارند و باعث می شوند که میانگین به مشاهدات دورافتاده نزدیک شود. نزدیکی میانگین به مشاهدات دورافتاده و اثر این مشاهدات بر ماتریس کوواریانس باعث کاهش فاصله ماهالانوبیس شده و در نهایت از شناسایی نقاط دورافتاده جلوگیری می کند. با این ویژگی ها چون برآوردگرهای کلاسیک تحت تاثیر نقاط دورافتاده هستند قادر به شناسایی نقاط پرت نخواهیم بود، به این ترتیب مقصود از درون آوری پنهان شدن مشاهدات دورافتاده به واسطه برآوردگرهای مکانی و پراگندگی آلوده می باشد. برای کنکاش بیشتر در مورد درون آوری به خلجی و خردمندیا [۱] رجوع کنید.

می دانیم مشاهداتی که از انبوهه^۳ اصلی مشاهدات دور هستند دورافتاده تلقی می شوند. اگر توزیع مولد داده های نمونه نرمال چند متغیره باشد، در یک نمودار پراکنش چند بعدی انبوهه داده ها به شکل بیضی گون است. بیضی گونی که در مرکز آن فشردگی داده ها زیاد و هر چه از مرکز دور شویم فشردگی داده ها کمتر می شود. بنابراین اگر تعداد متغیرها ۱، ۲ و یا ۳ باشد با رسم نمودار پراکنش دوری بعضی مشاهدات از انبوهه اصلی به صورت بصری قابل مشاهده است. در حالت سه بعدی با چرخاندن نمودار، انبوهه بیضی گون و نقاط دور از آن انبوهه آشکار می شوند. در حالت $p > 3$ این دوری قابل رویت نیست و لذا باید به طریق قابل اعتمادی اندازه گیری شود. روش های شهودی در شناسایی نقاط دورافتاده زمانی که تعداد متغیرها بیشتر از ۳ باشد غیر قابل اعتماد یا امکان ناپذیر هستند. از این رو به کمک روش های دیگر که غالباً محاسباتی هستند به شناسایی نقاط دورافتاده می پردازند. اولین روش که به ذهن خطور می کند، استفاده از مربع فاصله ماهالانوبیس^۴ می باشد که برای مشاهده i ام به صورت

$$MSD_i = (\bar{x}_i - \bar{x})^T S^{-1} (\bar{x}_i - \bar{x}), \quad (1)$$

تعریف می شود که \bar{x} میانگین حسابی و S ماتریس کوواریانس نمونه ای مجموعه داده های X می باشد. طبیعی است که مشاهداتی که دارای فواصل بزرگ می باشند به عنوان نقاط دورافتاده شناسایی شوند. اما برای رسیدن به

^۵Cutoff Value

^۶Masking

^۳Cloud

^۴Mahalanobis Squared Distance

۳ شناسایی به وسیله مربع فاصله استوار

می‌دهیم. به این ترتیب مربع فاصله استوار^۹ مشاهده i ام عبارت است از

$$RSD_i = (\underline{x}_i - \tilde{\underline{\mu}}_{(MCD)})^T \tilde{\Sigma}_{(MCD)}^{-1} (\underline{x}_i - \tilde{\underline{\mu}}_{(MCD)}).$$

روسو و زومرن [۹] پیشنهاد کردند از مقدار برشی $\chi_{\alpha,p}^2$ استفاده شود و مقدار $\alpha = 0.025$ را به عنوان یک مقدار مرسوم معرفی کردند. بنابراین مشاهده \underline{x}_i دورافتاده تلقی می‌شود اگر مربع فاصله استوار آن از مقدار برشی $\chi_{0.025,p}^2$ بیشتر باشد.

با احتمال زیاد می‌توانیم امیدوار باشیم که برآوردهای MCD ، پارامتر مکانی و پراکندگی را اندازه‌گیری می‌کند که حدود نیمی از داده‌های سالم را در بر دارد. بنابراین استفاده از برآوردهای استوار منجر به تولید فواصل استواری شده که تحت تاثیر نقاط دورافتاده نبوده و لذا استفاده از آنها درون‌آوری را از بین می‌برد. برای روشن شدن مطلب به مثال‌های زیر توجه کنید. لازم به ذکر است که در این مثال‌ها از فواصل (و نه از مربع آنها) و مقدار برشی $\sqrt{\chi_{\alpha,p}^2}$ برای شناسایی نقاط دورافتاد استفاده شده است.

مثال ۱: داده‌های وزن مغز و بدن

این مجموعه داده‌ها شامل وزن مغز و وزن بدن ۲۸ گونه جانوری است، که توسط دانشگاه ایلینوی جمع‌آوری شده است. وزن بدن جانوران بر حسب کیلوگرم و وزن مغز آنها بر حسب گرم محاسبه شده است. داده‌های خام تحت عنوان *Animals* در نرم‌افزار R موجود می‌باشند. در این مثال با استفاده از برنامه‌ای که در پیوست ارائه

\bar{x} مرکز همه داده‌ها است ولی مرکز انبوهه اصلی داده‌ها نیست. به طریق مشابه S نیز اطلاعات پراکندگی همه داده‌ها را در بر دارد و فاقد اطلاعات انبوهه اصلی می‌باشد. بنابراین با توجه به مباحث بخش قبل معقول به نظر می‌رسد که به جای برآوردهای کلاسیک \bar{x} و S از برآوردهای استوار^۷ استفاده شود که در برابر نقاط دورافتاده مقاوم بوده و حاوی اطلاعات انبوهه اصلی داده‌ها می‌باشند. در این مقاله برآوردهای استوار بردار میانگین و ماتریس کوواریانس را به ترتیب با $\tilde{\underline{\mu}}$ و $\tilde{\Sigma}$ نشان می‌دهیم.

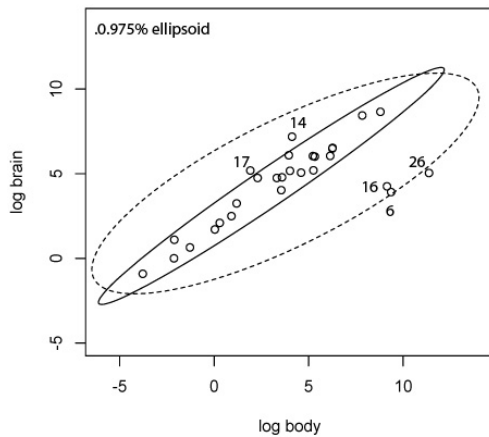
کمپل [۲] پیشنهاد کرد که برآوردهای M را جایگزین \bar{x} و S کنیم که بهبود حائز اهمیتی حاصل شد. اما نقطه فروریختگی این برآوردها حداکثر برابر $\frac{1}{p+1}$ می‌باشد (روسو [۸]) که با افزایش بعد داده‌ها به شدت دچار فروریختگی می‌شود.

در پژوهش‌های بعد برآوردهای بردار مکانی و ماتریس پراکندگی که دارای بیشینه فروریختگی بودند جایگزین \bar{x} و S شدند. برای اولین بار استاهل [۱۰] و دونوو و هابر [۳] چنین برآوردهایی را توصیه کردند. در این بخش ما از برآوردهای کمینه دترمینان ماتریس کوواریانس $\wedge(MCD)$ که دارای نقطه فروریختگی 50° درصدی می‌باشند، استفاده می‌کنیم. ایده اولیه این روش متعلق به روسو و زومرن [۹] می‌باشد. برآورد استوار مبتنی بر روش MCD را با $\tilde{\underline{\mu}}_{(MCD)}$ و $\tilde{\Sigma}_{(MCD)}$ نشان

^۷Robust

^۸Minimum Covariance Determinant

^۹Robust Squared Distance



شکل ۲: بیضی‌گون ۹۷/۵ درصد برای لگاریتم داده‌های وزن مغز و بدن

شده فواصل ماهالانوبیس (MD) و فاصله استوار (RD) را محاسبه کرده و نتایج را در جدول ۱ قرار داده‌ایم. با استفاده از مقدار برشی $\sqrt{\chi_{0.05,2}^2} = 2.72$ و فاصله ماهالانوبیس مشاهدات تنها مشاهده ۲۶ که مربوط به نوعی دایناسور (*Brachiosaurus*) می‌باشد به‌عنوان نقطه دورافتاده شناسایی می‌شود. اما با استفاده از فاصله استوار مشاهده‌های ۱۶، ۶ و ۲۶ که سه دایناسور با مغز کوچک و بدن سنگین هستند و مشاهدات ۱۴ و ۱۷ که مربوط به بشر و نوعی میمون است به‌عنوان داده دورافتاده شناسایی می‌شوند.

جدول ۱: فاصله ماهالانوبیس و فاصله استوار برای داده‌های وزن مغز و بدن

مشاهده	MD	RD	مشاهده	MD	RD
۱	۱,۰۱	۰,۸۶	۱۵	۱,۷۶	۱,۷۲
۲	۰,۷	۱,۴۸	۱۶	۲,۳۷	۹,۵۸
۳	۰,۳	۰,۲۴	۱۷	۱,۲۲	۳,۵۷
۴	۰,۳۸	۰,۵۲	۱۸	۰,۲	۱,۳۲
۵	۱,۱۵	۱,۱۱	۱۹	۱,۸۶	۱,۷۴
۶	۲,۶۴	۱,۰۶	۲۰	۲,۲۷	۲,۰۳
۷	۱,۷۱	۱,۸۱	۲۱	۰,۸۳	۰,۷۴
۸	۰,۷۱	۰,۶۹	۲۲	۰,۴۲	۰,۴۲
۹	۰,۸۶	۱,۰۵	۲۳	۰,۲۶	۰,۹۴
۱۰	۰,۸۰	۲,۰۸	۲۴	۱,۰۵	۲,۲۹
۱۱	۰,۶۹	۰,۸۹	۲۵	۱,۵۹	۱,۵۲
۱۲	۰,۸۷	۱,۰۴	۲۶	۲,۹۱	۱۱,۴۷
۱۳	۰,۶۸	۰,۶۸	۲۷	۱,۵۸	۲,۰۵
۱۴	۱,۷۲	۴,۲۶	۲۸	۰,۳۹	۱,۶۸

نکته ۱: شناسایی نقاط دورافتاده زمانی که $\frac{n}{p}$ کوچک باشد امری دشوار است. زیرا در این حالت نقاط کمتری در محاسبه MCD شرکت دارند. با کمک شیبه سازی و به‌عنوان یک قاعده سرانگشتی توصیه می‌شود که اندازه نمونه حداقل باید ۵ برابر تعداد متغیرها باشد.

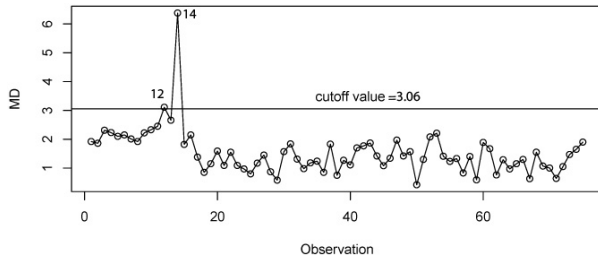
مثال ۲: داده‌های هاوکینز و همکاران

در این مثال با استفاده از داده‌هایی که توسط هاوکینز و همکاران [۵] تولید شده است، عملکرد فاصله استوار در شناسایی نقاط دورافتاده و مقابله با درون‌آوری را زمانی که تعداد مشاهدات زیاد است را مورد بررسی قرار می‌دهیم. سه متغیر اول این مجموعه تشکیل مجموعه داده‌ای با $n = 75$ و $p = 3$ را می‌دهد که نتایج بررسی در شکل-های ۳ و ۴ آورده شده است.

با توجه به شکل ۳ که نشان دهنده فاصله‌های ماهالانوبیس است. با در نظر گرفتن مقدار برشی $3/06$ تنها قادر به شناسایی دو مشاهده ۱۲ و ۱۴ شدیم. زیرا در این روش به علت استفاده از برآوردهای کلاسیک در فرمول فاصله ماهالانوبیس با فواصل بزرگ مواجه

در شکل ۲ بیضی ۹۷/۵ درصد برای لگاریتم داده‌های وزن مغز و بدن ارائه شده است. بیضی‌گون ماهالانوبیس با خط چین و بیضی‌گون استوار با خط ممتد مشخص شده که داده‌های خارج از بیضی‌گون‌ها پرت تلقی می‌شوند. با توجه به مثال ۱ می‌توان نتیجه گرفت که فاصله ماهالانوبیس به علت درون‌آوری قادر به شناسایی نقاط دور افتاده نبوده که با استفاده از فاصله استوار این مشکل مرتفع می‌شود.

نخواهیم شد و بسیاری از مشاهدات دورافتاده به خاطر درون آوری شناسایی نمی شوند. اما در شکل ۴ چون برآورد را از بین می برد.



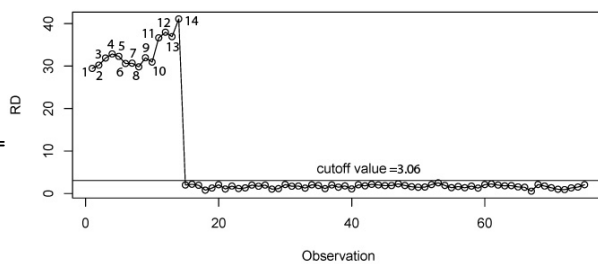
شکل ۳: نمودار فاصله ماهالانوبیس برای داده‌های هاوکینز

پیوست

این برنامه برای محاسبه و مقایسه فاصله ماهالانوبیس و فاصله استوار طراحی شده است. ابتدا ماتریس مشاهدات X را به عنوان ورودی دریافت کرده و با استفاده از بسته $MASS$ برآوردگرهای MVE و MCD را محاسبه می کند. در انتها فاصله ماهالانوبیس و فواصل استوار مبتنی بر برآوردگرهای استوار MVE و MCD را همراه با نمودار این فاصله‌ها ارائه می دهد. در این نمودارها فاصله‌ها با مقدار برشی $\sqrt{\chi_{\alpha, 25, p}^2}$ سنجیده شده و نقاطی که بالای خط قرار می گیرند، دورافتاده تلقی می شوند.

بردار میانگین و ماتریس کوواریانس را با روش استوار MCD محاسبه کرده و این برآوردگرها را در فرمول فاصله ماهالانوبیس جایگزین برآوردگرهای کلاسیک نموده ایم درون آوری از بین رفته و مشاهدات ۱ تا ۱۴ به عنوان نقاط دورافتاده شناسایی می شوند. همان طور که مشاهده می کنیم فاصله این مشاهدات به طور معنی داری از سایر مشاهدات بیشتر است و مشاهداتی که به عنوان مشاهده معمولی شناخته می شوند تقریباً فاصله یکسانی با انبوهه داده‌ها دارند.

```
RD=function(x){
library(MASS)
n=dim(x)[1];p=dim(x)[2];h=floor((n+p+1)/2)
#####
# compute of MVE &MCD
x.mve=cov.rob(x,h,method="mve")
x.mcd=cov.rob(x,h,method="mcd")
md=numeric();rd.mve=numeric()
rd.mcd=numeric()
#####
# compute of mahalanobis & robust distance
for(i in 1:n){
md[i]=round(sqrt(t(x[i,]-(colMeans(x)))*%
solve(cov(x))%*(x[i,]-(colMeans(x))))),2)
rd.mve[i]=round(sqrt(t(x[i,]-(x.mve$center))*%
solve(x.mve$cov)*%*(x[i,]-(x.mve$center))),2)
rd.mcd[i]=round(sqrt(t(x[i,]-(x.mcd$center))*%
solve(x.mcd$cov)*%*(x[i,]-(x.mcd$center))),2)
}
df=data.frame(MD=md,RD_MVE=rd.mve,RD_MCD=rd.mcd)
#####
```



شکل ۴: نمودار فاصله استوار برای داده‌های هاوکینز

نکته ۲ برآوردگرهای بیضی گون کمینه حجم (MVE)^۱ نیز خواصی مشابه برآوردگرهای MCD دارند و غالباً

^۱ Minimum Volume Ellipsoid

[6] Mardia, k., Kent, J. and Bibby, J. (1979). *Multivariate Analysis*, New York: Academic Press.

[7] Rencher, A. C. (2002). *Methods of Multivariate Analysis (2ndED.)*. New York: Wiley - Interscience.

[8] 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.

[9] Rousseeuw, P. J. and van Zomeren, B. C. (1990). *Unmasking Multivariate Outliers and Leverage Points*. *Journal of the American Statistical Association*, 85, 633-639.

[10] Stahel, W. A. (1981). *Robuste Schätzungen: Infinitesimale Optimalität und Schätzungen von Kovarianzmatrizen*. Unpublished Ph.D.thesis, ETH Zurich.

```
# plot of mahalnobis & robust distance
par(mfrow=c(2,1))
plot(md,type='o')
abline((sqrt(qchisq(.975,p))),0)
plot(rd.mve,type='o',lty=2)
abline((sqrt(qchisq(.975,p))),0)
df
}
```

مراجع

[۱] خلجی، ا. و خردمندنی، م. (۱۳۹۴). اثر ماسک در شناسایی نقاط دورافتاده چندمتغیره. دهمین سمینار احتمال و فرایندهای تصادفی، ۳۶۵-۳۶۹.

[2] Campbell, N. A. (1989). *Bushfire Mapping Using NOAA AVHRR Data*. Technical Report, CSIRO.

[3] Donoho, D. L. and Huber, P. J. (1983). *The Notion of Breakdown Point*. In *A Festschrift for Erich L. Lehmann*, eds. Bickel, P. J. , Doksum, K. A. and Hodges, J. L. , 157-184.

[4] Hawkins, D. M. (1980). *Identification of Outliers*, London: Chapman and Hall.

[5] 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.