

مجله علوم آماری، پاییز و زمستان ۱۴۰۰

جلد ۱۵، شماره ۲، ص ۴۴۳ - ۴۶۲

DOI: 10.29252/jss.15.2.443

مقاله پژوهشی

شناسائی نقاط دورافتاده چندمتغیره بر اساس تابع دورافتادگی ژرفا-مبنا

سکینه دهقان، محمدرضا فریدروحانی

گروه آمار، دانشکده علوم ریاضی، دانشگاه شهید بهشتی

تاریخ دریافت: ۱۳۹۹/۱۲/۰۸ تاریخ پذیرش و انتشار: ۱۴۰۰/۰۲/۰۴

چکیده: تابع ژرفا با در نظر گرفتن ویژگی‌های هندسی مجموعه داده‌های چندمتغیره و رتبه‌بندی مشاهدات ابزار مناسبی را در آمار ناپارامتری چندمتغیره فراهم آورده است. به عبارت دیگر، این تابع منجر به مرتب‌سازی از مرکز به بیرون نقاط چندمتغیره می‌شود. از آن‌جا که دورافتادگی نقاط به طور اجتناب‌ناپذیری وابسته به ترتیب داده‌ها است، این مرتب‌سازی می‌تواند راهی برای شناسایی نقاط دورافتاده فراهم کند. در این مقاله، بر اساس مفهوم تابع ژرفا، یک روش ناوردای آفین برای شناسائی نقاط دورافتاده چندمتغیره بیان می‌شود. ویژگی مطلوب ناوردای آفین تضمین می‌کند که نقطه دورافتاده تحت هرگونه تبدیل از محورهای مختصات کماکان به‌عنوان دورافتاده شناسایی شود. پیاده‌سازی این روش نسبت به بیشتر روش‌های چندمتغیره که دارای پیچیدگی محاسباتی هستند، ساده‌تر است. بر اساس مطالعات شبیه‌سازی عملکرد روش پیشنهادی بر اساس توابع ژرفای مختلف مورد بررسی قرار گرفته است. سرانجام، روش بیان شده برای داده‌های مسکن شهرهای منتخب ایران در سال ۱۳۹۷، بکار برده می‌شود.

واژه‌های کلیدی: تابع ژرفا، تابع دورافتادگی، دورافتاده، ناوردایی آفین.

۱ مقدمه

در بسیاری از مسائل استنباطی، شناسایی نقاط دورافتاده^۱ نمونه، به دو دلیل از اهمیت فراوانی برخوردار است. اولاً حضور نقاط دورافتاده در یک مجموعه داده اثر نامطلوبی بر استنباط از آن مجموعه داده گزارده و منجر به برآوردگرها، آماره‌های آزمون، بازه‌های اطمینان و ... نامناسبی می‌شود. ثانیاً برخی موارد نقاط دورافتاده اطلاعات مفیدی از مدل داده‌ها به تحلیل‌گر می‌دهند و از این روش‌شناسی آن‌ها برای تحلیل‌گر دارای اهمیت است. شناسایی نقاط دورافتاده در پزشکی، کشف تقلب، یافتن خطای اندازه‌گیری، رکورگیری‌های ورزشی و ... کاربرد دارد. برای شناسایی نقاط دورافتاده در حالت یک‌متغیره، روش‌های متعددی وجود دارد (هاوکینز، ۱۹۸۰؛ بارت و لويس، ۱۹۹۴). دورافتاده در نظر گرفتن نقاط خارج از بازه‌ای با مرکز میانگین مشاهدات و کران‌هایی که با اضافه و کم کردن سه برابر انحراف معیار مشاهدات از میانگین حاصل می‌شوند، یکی از متداولترین روش‌ها برای شناسایی نقاط دورافتاده یک‌متغیره است. این روش برگرفته از ایده بازه اطمینان با ضریب اطمینان ۹۹٫۸۷ برای میانگین توزیع نرمال است. البته برخی از نویسندگان با توجه به نوع مساله از ۲٫۵ یا ۲ برابر انحراف استاندارد حول میانگین استفاده می‌کنند (میلر، ۱۹۹۱). میانگین و انحراف معیار مشاهدات برآوردگرهای استواری از پارامترهای مکان و مقیاس نیستند و عملکرد این روش که معمولاً تحت برقراری فرض نرمال از آن استفاده می‌شود، تحت تاثیر نقاط دورافتاده است. یک روش جایگزین مناسب، استفاده از بازه‌ای بر اساس برآوردگر استوار میانه مشاهدات برای پارامتر مکان و برآوردگر استوار میانه قدرمطلق انحراف نقاط از میانه برای پارامتر مقیاس است. بسط روش‌های یک‌متغیره به چندمتغیره با مشکلات زیادی همراه است. با به کارگیری روش‌های یک‌متغیره و بررسی دورافتادگی حاشیه‌ای نقاط چندمتغیره، می‌توان به طور مجزا در هر بعدی نسبت به شناسایی نقاط دورافتاده اقدام کرد. روش بررسی دورافتادگی حاشیه‌ای نمی‌تواند بسیاری از نقاط دورافتاده را شناسایی کند و همچنین ویژگی‌های هندسی مجموعه داده‌ها را در نظر نمی‌گیرد. تجسم مجموعه داده‌ها برای ابعاد بیشتر از سه بسیار سخت است، بنابراین روش تصویری‌سازی عملاً محدود به بعد سه و کمتر است و برای ابعاد بالاتر کاربردی ندارد. همچنین می‌توان با روش‌های خوشه‌بندی چندمتغیره، داده‌ها را به خوشه‌هایی تقسیم کرد و خوشه‌هایی که فقط شامل یک یا دو نقطه هستند، به صورت دورافتاده در نظر گرفت. به دلیل سختی خوشه‌بندی، زمان‌بر بودن محاسبات رایانه‌ای، وابستگی به بعد داده‌ها، اندازه نمونه و ... این روش خیلی به‌کار نمی‌رود. شناسایی نقاط دورافتاده چندمتغیره بر اساس فاصله ماله‌الانویس توسط راک و وودراف (۱۹۹۶) معرفی شد. با اینکه این روش به طور گسترده‌ای به کار می‌رود، ولی چنانچه

¹Outlier

فرض تقارن بیضوی برای مشاهدات برقرار نباشد، دارای عملکرد مناسبی نیست. به منظور افزایش کارایی این روش در حالتی که فرض تقارن بیضوی برقرار نیست، **سرفلینگ و موزندر (۲۰۱۳)** با به کارگیری روش جستجوی تصویر، شناساگر دورافتادگی معرفی شده توسط **راک و وودراف (۱۹۹۶)** را تعمیم دادند. **کان و همکاران (۲۰۱۵)** با به کارگیری چندین فاصله، ابزارهای تشخیصی متعددی را برای شناسایی نقاط دورافتاده چندمتغیره معرفی کردند.

با معرفی تابع ژرفای نیم‌فضا^۱ توسط **توکی (۱۹۷۵)** در بسیاری از مباحث آمار ناپارامتری از این تابع استفاده می‌شود. تابع ژرفا^۲ با در نظر گرفتن ویژگی‌های هندسی مجموعه داده‌های چندمتغیره و رتبه‌بندی مشاهدات چندمتغیره ابزار مناسبی را فراهم آورده است. به عبارت دیگر، در حالت چندمتغیره، ژرفای داده‌ها می‌تواند برای اندازه‌گیری ژرفا یا دورافتادگی یک نقطه چندمتغیره نسبت به توزیع آن به‌کار رود که قابلیت مزبور منجر به مرتب‌سازی مرکز به بیرون نقاط نمونه می‌شود. واضح است، این مرتب‌سازی می‌تواند راهی برای شناسایی نقاط دورافتاده فراهم کند و نقطه دورافتاده تحت هرگونه تبدیل از محورهای مختصات کماکان به‌عنوان دورافتاده شناسایی شود. با توجه به ویژگی‌های مطلوب تابع ژرفا، پژوهش‌هایی برای شناسایی نقاط دورافتاده انجام گرفته است. **چن و همکاران (۲۰۰۸)** تابع ژرفای فضائی هسته‌ای^۳ را معرفی کردند و بر اساس آن یک شناساگر برای تعیین نقاط دورافتاده تعیین کردند. آن‌ها نشان دادند که این تابع ویژگی‌های هندسی توزیع را کاملاً در نظر گرفته و دارای عملکرد مناسبی در تشخیص نقاط دورافتاده است. البته روش معرفی شده توسط **چن و همکاران (۲۰۰۸)** دارای پیچیدگی‌های محاسباتی در تعیین آستانه مورد نیاز برای شناسایی نقاط دورافتاده است. **دنگ و سرفلینگ (۲۰۱۰)** با به کارگیری توابع ژرفا، شناساگرهایی با نقطه فروریزش^۴ بالا معرفی کردند. نقطه فروریزش که عبارت است از نسبتی از نمونه که می‌تواند آلوده شود بدون این‌که بعضی نقاط دورافتاده مخفی شوند (به صورت غیر دورافتاده شناسایی شوند)، معیاری برای بررسی استواری^۵ شناساگرهای دورافتادگی محسوب می‌شود (**روزبه و امینی، ۱۳۹۸؛ روزبه و معنوی، ۱۳۹۹**). در ادامه، **دوودو و چاکرابوتی (۲۰۱۳)** شناساگرهای دنگ و سرفلینگ را برای حالت خاص توزیع چوله نرمال به‌کار بردند. **فن (۲۰۱۶)** با محاسبه ژرفای نقاط اصلی نمونه نسبت به مجموعه نقاط اصلی نمونه و نقاطی که نسبت به ژرفترین نقطه متقارن شده‌اند، چنین استدلال کرد که می‌توان عملکرد روش شناسایی نقاط دورافتاده بر اساس تابع ژرفا را بهبود بخشید. روش وی با در نظر گرفتن نمونه متقارن شده، دارای پیچیدگی

^۱Half space depth

^۲Depth function

^۳Kernelized spatial depth function

^۴Breakdown point

^۵Robustness

محاسباتی است و نسبت به شناساگرهای **دنگ و سرفلینگ** (۲۰۱۰) باعث بهبود عملکرد چشم‌گیری در زمینه شناسائی نقاط دورافتاده نشده است.

در بخش ۲ مفهوم ژرفای داده و انواع تابع ژرفا مرور می‌شوند، سپس تابع دورافتادگی ژرفا-مبنا معرفی می‌شود. در بخش ۳ شناساگرهای دورافتادگی^۱ بر اساس توابع ژرفا بر اساس یک مدل آلودگی معرفی شده و به طور دقیق مقدار آستانه مورد نیاز برای شناسائی نقاط دورافتاده مشخص می‌شود. مطالعات شبیه‌سازی برای بررسی و مقایسه عملکرد توابع ژرفای مختلف در بخش ۴ ارائه می‌شود. در بخش ۵ بر اساس داده‌های قیمت مسکن شهرهای منتخب ایران که توسط مرکز آمار ایران ثبت شده‌اند، شهرهایی که از نظر ارزش مالی مسکن نسبت به سایر شهرها دورافتاده محسوب می‌شوند، شناسائی و نتایج تحلیل می‌شوند.

۲ توابع ژرفا و دورافتادگی

فرض کنید اندازه احتمال P بر فضای \mathcal{R}^p داده شده و F تابع توزیع متناظر با آن باشد. ژرفای داده، اندازه‌ای از میزان ژرف بودن یا دورافتاده بودن یک نقطه نسبت به توده داده‌ها یا توزیع F است. هر تابع $D(x, F)$ که ترتیبی از مرکز به بیرون از نقاط x متعلق به \mathcal{R}^p فراهم کند، تابع ژرفای متناظر با F نامیده می‌شود. بنابراین بر اساس تابع ژرفای مفروض، نقاط p متغیره به نحوی رتبه‌بندی می‌شوند که ژرفای بزرگتر نشان‌دهنده مرکزی‌تر بودن نقاط و ژرفای کوچکتر نشان‌دهنده دورافتاده‌تر بودن آن‌ها است. تحت تابع ژرفای مفروض، نقطه یا میانگین نقاط با بیشترین مقدار ژرفا، مرکز را تشکیل می‌دهند.

تعریف ۱. (ژو و سرفلینگ، ۲۰۰۰) اگر X یک بردار تصادفی p بعدی بر فضای احتمال (Ω, \mathcal{F}, P)

باشد، تابع $D(\cdot, F) : \mathcal{R}^p \rightarrow \mathcal{R}$ ، ژرفای آماری است، اگر در ویژگی‌های زیر صدق کند:

الف- ناوردایی آفین: $D(x, F)$ مستقل از دستگاه مختصات باشد. یعنی برای هر بردار تصادفی x در \mathcal{R}^p ، هر ماتریس A نامنفرد $p \times p$ و هر بردار p بعدی b ، $D(Ax + b, F_{Ax+b}) = D(x, F_X)$ ،

که در آن F_X و F_{Ax+b} به ترتیب توابع توزیع بردارهای تصادفی X و $AX + b$ هستند.

ب- ماکسیم شدن در مرکز: اگر F نسبت به θ متقارن باشد، آنگاه $D(x, F)$ در θ ماکسیم است.

ج- به صفر رسیدن در بینهایت: اگر $\|x\| \rightarrow \infty$ آنگاه $D(x, F) \rightarrow 0$.

د- یکنوایی نسبت به ژرفترین نقطه: اگر برای هر بردار تصادفی x در \mathcal{R}^p ، $D(\theta, F) \geq D(x, F)$ ،

آنگاه برای هر $\alpha \in [0, 1]$ ، $D(\theta + \alpha(x - \theta), F) \geq D(x, F)$.

¹Outlier identifiers

به منظور تصویر و مرتب‌سازی داده‌های چندمتغیره توکی (۱۹۷۵) تابع ژرفای نیم‌فضا، (HD)، را معرفی کرد. به دنبال آن تابع ژرفای سادگی^۱ (LD)؛ (لیو، ۱۹۸۸)،، تابع ژرفای ماهالانوبیس^۳ (MD)؛ (لیو و همکاران، ۱۹۹۳)، تابع ژرفای تصویر-پایه^۱ (PD)؛ (ژو، ۲۰۰۳)، و تابع ژرفای فضائی^۲ (SD)؛ (وردی و ژانگ، ۲۰۰۰)، پیشنهاد شد، که در ادامه معرفی می‌شوند.

تعریف ۲. (توکی، ۱۹۷۵) تابع ژرفای نیم‌فضا برای $\mathbf{x} \in \mathcal{R}^p$ نسبت به توزیع F به صورت

$$HD(\mathbf{x}, F) = \inf \{P(H) : H \text{ نیم‌فضایی بسته شامل } \mathbf{x} \text{ در فضای } p \text{ بعدی است}\}$$

تعریف می‌شود، که در آن $P(H)$ احتمال نیم‌فضای H توسط توزیع F است. ژرفای نیم‌فضای نقطه^۲ \mathbf{x} نسبت به نمونه^۲ $\mathbf{X}_1, \dots, \mathbf{X}_n$ در \mathcal{R}^p عبارت است از کمترین کسر نمونه در بین تمامی نیم‌فضاهای بسته شامل نقطه^۲ \mathbf{x} . به عبارت دیگر

$$HD(\mathbf{x}, F_n) = \frac{1}{n} \min_{\|\mathbf{u}\|=1} \#\{i : \mathbf{u}^T \mathbf{X}_i \leq \mathbf{u}^T \mathbf{x}, i = 1, \dots, n\}$$

که در آن F_n تابع توزیع تجربی نمونه^۲ $\mathbf{X}_1, \dots, \mathbf{X}_n$ است.

تعریف ۳. بردارهای k_1, \dots, k_t در \mathcal{R}^p به طور آفین^۳ مستقل هستند، اگر و تنها اگر $t - 1$ بردار $k_t - k_1, \dots, k_t - k_{t-1}$ به طور خطی مستقل باشند. در فضای اقلیدسی p بعدی \mathcal{R}^p ، مجموعه‌ای از نقاط که ترکیبی محدب از $p + 1$ نقطه به طور آفین مستقل باشند را سادک می‌نامند. به عبارت دیگر اگر نقاط به طور آفین مستقل را با k_1, \dots, k_{p+1} نشان داده شوند، سادک تولید شده توسط این نقاط به صورت

$$S(k_1, k_2, \dots, k_{p+1}) = \left\{k : k = a_1 k_1 + \dots + a_{p+1} k_{p+1}, a_i \geq 0, \sum_{i=1}^{p+1} a_i = 1\right\}.$$

تعریف می‌شود، که در آن k_1, \dots, k_{p+1} رئوس سادک نام دارند.

^۲Simplicial depth

^۳Mahalanobis depth

^۱Projection-based depth

^۲Spatial depth

^۳Affine invariance

تعریف ۴. تابع ژرفای سادگی (لیو، ۱۹۸۸) برای $\mathbf{x} \in \mathcal{R}^p$ نسبت به توزیع F به صورت

$$LD(\mathbf{x}, F) = P\{\mathbf{x} \in S[\mathbf{X}_1, \dots, \mathbf{X}_{p+1}]\},$$

تعریف می‌شود، که در آن $\mathbf{X}_1, \dots, \mathbf{X}_{p+1}$ مشاهدات مستقل از توزیع F هستند و $S[\mathbf{x}_1, \dots, \mathbf{x}_{p+1}]$ نشان‌دهندهٔ سادک بسته در \mathcal{R}^p با رئوس $\mathbf{x}_1, \dots, \mathbf{x}_{p+1}$ است. ژرفای سادگی نقطهٔ \mathbf{x} نسبت به نمونهٔ $\mathbf{X}_1, \dots, \mathbf{X}_n$ در \mathcal{R}^p برابر با

$$LD(\mathbf{x}, F_n) = \binom{n}{p+1}^{-1} \sum_* I(\mathbf{x} \in S[\mathbf{X}_{i_1}, \dots, \mathbf{X}_{i_{p+1}}])$$

است، که در آن $*$ شامل همهٔ زیرمجموعه‌های $p+1$ عضوی از $\mathbf{X}_1, \dots, \mathbf{X}_n$ و $I(\cdot)$ تابع نشان‌گر است.

تعریف ۵. (لیو و همکاران، ۱۹۹۳) تابع ژرفای ماهالانوبیس برای $\mathbf{x} \in \mathcal{R}^p$ نسبت به توزیع F

$$MD(\mathbf{x}, F) = \left[1 + (\mathbf{x} - \boldsymbol{\mu}_F)^T \sum_F^{-1} (\mathbf{x} - \boldsymbol{\mu}_F) \right]^{-1}$$

است، که در آن $\boldsymbol{\mu}_F$ و $\boldsymbol{\Sigma}_F$ به ترتیب بردار میانگین و ماتریس کوواریانس توزیع F هستند. برای به دست آوردن نسخهٔ نمونه‌ای ژرفای ماهالانوبیس نقطهٔ $\mathbf{x} \in \mathcal{R}^p$ کافی است به جای $\boldsymbol{\mu}_F$ و $\boldsymbol{\Sigma}_F$ برآوردهای نمونه‌ای آن‌ها جایگزین شوند، که یک انتخاب معمول به ترتیب بردار میانگین نمونه $\bar{\mathbf{X}}$ و ماتریس کوواریانس نمونه‌ای \mathbf{S} است.

تعریف ۶. فرض کنید \mathbf{X} دارای توزیع F باشد و

$$\tilde{O}_P(\mathbf{x}, F) = \sup_{\|\mathbf{u}\|=1} \left| \frac{\mathbf{u}^T \mathbf{x} - \mu(F_{\mathbf{u}^T \mathbf{X}})}{\sigma(F_{\mathbf{u}^T \mathbf{X}})} \right| \quad (1)$$

که در آن $\mu(\cdot)$ و $\sigma(\cdot)$ به ترتیب اندازهٔ یک متغیرهٔ مکان و مقیاس توزیع F هستند. هرگاه هر دو مقدار صورت و مخرج کسر (۱) صفر باشند، یعنی $\mathbf{u}^T \mathbf{x} - \mu(F_{\mathbf{u}^T \mathbf{X}}) = \sigma(F_{\mathbf{u}^T \mathbf{X}}) = 0$ ، آن‌گاه صفر در نظر گرفته می‌شود. ژرفای تصویر-مینا (ژو، ۲۰۰۳) برای $\mathbf{x} \in \mathcal{R}^p$ نسبت به

توزیع F به صورت $PD(\mathbf{x}, F) = [1 + \tilde{O}_P(\mathbf{x}, F)]^{-1}$ تعریف می‌شود.

برای به دست آوردن نسخه نمونه‌ای ژرفای تصویر-مبنا به جای توزیع F همتای نمونه‌ای آن، F_n جایگزین می‌شود. هر زوج خاصی که برای برآورد (μ, σ) به کار رود، تابع ژرفای نمونه‌ای تصویر-مبنا خاصی را نتیجه می‌دهد. اما به نظر انتخابی مناسب می‌تواند زوج $(Med, MAD) = (\mu, \sigma)$ باشد، که در آن Med میانه و MAD میانه قدرمطلق انحراف از میانه هستند. برای نمونه یک‌متغیره $\mathbf{Y}_N = \{Y_1, \dots, Y_N\}$ در R ، به صورت زیر تعریف می‌شوند.

$$Med(\mathbf{Y}_N) = \text{median} \{Y_i, 1 \leq i \leq N\}, \quad (2)$$

$$MAD(\mathbf{Y}_N) = \text{median} \{|Y_i - Med(\mathbf{Y}_N)|, 1 \leq i \leq n\}. \quad (3)$$

تعریف ۷. تابع علامت برداری در R^p به صورت

$$\mathbf{S}(\mathbf{x}) = \begin{cases} \frac{\mathbf{x}}{\|\mathbf{x}\|} & \mathbf{x} \neq \mathbf{0} \\ \mathbf{0} & \mathbf{x} = \mathbf{0} \end{cases}$$

تعریف می‌شود، که در آن $\|\mathbf{x}\|$ نرم اقلیدسی \mathbf{x} به صورت $\sqrt{\mathbf{x}^T \mathbf{x}}$ است. اگر \mathbf{X} دارای توزیع F باشد، آن‌گاه ژرفای فضایی (وردی و ژانگ، ۲۰۰۰) نقطه \mathbf{x} نسبت به توزیع F به صورت

$$SD(\mathbf{x}, F) = 1 - \|E(\mathbf{S}(\mathbf{x} - \mathbf{X}))\|$$

تعریف می‌شود. مقدار تابع فضایی نمونه‌ای نقطه \mathbf{x} ، نسبت به نمونه $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ نیز عبارتست از

$$SD(\mathbf{x}, F_n) = 1 - \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{S}(\mathbf{x} - \mathbf{X}_i) \right\|.$$

با توجه به تعریف توابع ژرفا واضح است، نقطه‌ای با مقدار ژرفای کمتر به عنوان دورافتاده تلقی می‌شود. بر این اساس، **دنگ و سرفلینگ** (۲۰۱۰)، برای به دست آوردن تابع دورافتادگی ابتدا تابع ژرفا را بر بازه $\mathbf{0}$ تا 1 استاندارد کرده، سپس از رابطه $O(\mathbf{x}, F) = 1 - D(\mathbf{x}, F)$ تابع دورافتادگی متناظر را تعریف کردند. به طور مشابه، **دهقان و فریدروحانی** (۲۰۱۹) برای تابع ژرفای استاندارد شده بر بازه $\mathbf{0}$

تا ۱، تابع دورافتادگی را به صورت $O(\mathbf{x}, F) = \frac{1}{D(\mathbf{x}, F)} - 1$ تعریف کردند، که مقادیر خود را در بازه $[0, \infty)$ اختیار می‌کند. در فضای \mathcal{R}^p ، منحنی‌های تراز متناظر تابع $g(\mathbf{x})$ عبارتند از رده‌هایی هم‌ارزی از نقاط \mathbf{x} که برای آن‌ها مقدار $g(\mathbf{x})$ یکسان است. به طور مثال اگر $g(\cdot)$ تابع چگالی احتمال باشد، منحنی‌های تراز، مقادیر چگالی احتمال روی نقاط واقع در رده هم‌ارزی را مشخص می‌کنند. این رده‌های هم‌ارزی لزوماً با رده‌های دورافتادگی یکسان نیستند. چنان‌چه $g(\cdot)$ را تابع ژرفا یا تابع دورافتادگی در نظر بگیریم، رده‌های هم‌ارزی متناظر، شامل نقاطی با مقدار ژرفای برابر خواهند بود. منحنی‌های تراز در تحلیل چندمتغیره نوعاً به دو دلیل توصیف چگالی احتمال و شناسایی نقاط دورافتاده به‌کار می‌روند. مجموعه

$$\{\mathbf{x} \in \mathcal{R}^p : O(\mathbf{x}, F) = t\} \quad (۴)$$

مجموعه سطح یا منحنی تراز دورافتادگی t نامیده می‌شود.

۳ کشف نقاط دورافتاده چندمتغیره

برای تعیین این‌که نقطه $\mathbf{x} \in \mathcal{R}^p$ نسبت به تابع توزیع F دورافتاده است یا خیر، ابتدا تابع دورافتادگی $O(\cdot, F)$ و مقدار آستانه‌ای λ مشخص می‌شود. سپس نقاطی که مقدار دورافتادگی آن‌ها از این مقدار آستانه‌ای تجاوز کند، نقاط دورافتاده نسبت به توزیع F نامیده می‌شوند. به عبارت دیگر، هر نقطه \mathbf{x} متعلق به $\text{out}(\lambda, F) = \{\mathbf{x} : O(\mathbf{x}, F) > \lambda\}$ ناحیه λ -دورافتاده از توزیع F نامیده می‌شود. حال باید مقدار آستانه λ تعیین شود. از طرف دیگر، در عمل، توزیع F معمولاً نامعلوم است و باید از تابع دورافتادگی نمونه‌ای استفاده کرد. ابتدا مدل آلودگی $F = (1 - \varepsilon)G + \varepsilon H$ را برای توزیع F در نظر می‌گیریم، که G توزیع ایده‌آل معلوم مدل، H منبع نامعلوم آلودگی و ε نسبت و یا احتمال دورافتادگی است. مشاهدات کرانگین نسبت به توزیع G به اشتباه به عنوان دورافتاده و نسبت به توزیع H به درستی به عنوان دورافتاده تلقی می‌شوند. از این‌رو مشاهدات کرانگین نسبت به توزیع G مثبت نادرست و نسبت به توزیع H مثبت درست، نامیده می‌شوند. نرخ مثبت نادرست^۱ را حاصل ضرب احتمال این‌که \mathbf{X} از توزیع G تولید شده باشد در احتمال این‌که تحت توزیع G به عنوان λ -دورافتاده شناسایی شود، تعریف می‌کنیم. بنابراین نرخ مثبت نادرست به صورت $(1 - \varepsilon) P_G(O(\mathbf{X}, G) > \lambda)$ ، تعریف می‌شود و به

^۱False positive rate

طور مشابه نرخ مثبت درست^۲ $\varepsilon P_H(O(\mathbf{X}, H) > \lambda)$ است. آستانه λ باید به گونه‌ای انتخاب شود که برای λ به اندازه کافی بزرگ، $P_G(O(\mathbf{X}, G) > \lambda)$ مقداری کوچک و برای λ به اندازه کافی کوچک، $P_H(O(\mathbf{X}, H) > \lambda)$ مقداری بزرگ باشد. بنابراین می‌توان به طور تقریبی روابط

$$(1 - \varepsilon) P_G(O(\mathbf{X}, G) > \lambda) \cong P_G(O(\mathbf{X}, G) > \lambda),$$

$$\varepsilon P_H(O(\mathbf{X}, H) > \lambda) \cong \varepsilon,$$

را در نظر گرفت. اگر $\alpha = P_G(O(\mathbf{X}, G) > \lambda)$ ، در این صورت به طور تقریبی α و ε به ترتیب نشان‌دهنده نرخ مثبت نادرست و نرخ مثبت درست هستند. اگر α از قبل تعیین شده باشد، آنگاه

$$\lambda = F_{O(\mathbf{X}, G)}^{-1}(\alpha). \quad (5)$$

بنابراین آستانه λ چند $\alpha - 1$ م توزیع $O(\mathbf{X}, G)$ تحت توزیع ایده‌آل G است. مطلوب آن است که نرخ مثبت نادرست نسبت به نرخ مثبت درست مقداری کوچک باشد، بنابراین اگر $\delta = \alpha/\varepsilon$ ، آنگاه δ باید مقداری کوچک باشد. در عمل δ و ε از قبل تعیین شده هستند، و لذا رابطه (۵) به صورت

$$\lambda = F_{O(\mathbf{X}, G)}^{-1}(1 - \delta\varepsilon) \quad (6)$$

است. برای مثال اگر $\delta = 0.1$ و نسبت دورافتادگی $\varepsilon = 0.25$ باشد، آنگاه $F_{O(\mathbf{X}, G)}^{-1}(0.975)$ است. منطقی است، نسبت آلوده شده از نمونه با افزایش اندازه نمونه کاهش یابد. در این صورت، می‌توان ε را برابر $\frac{c}{\sqrt{n}}$ که c مقداری ثابت است، تعریف کرد. بنابراین از رابطه (۶) داریم

$$\lambda_n = F_{O(\mathbf{X}, G)}^{-1}(1 - c\delta/\sqrt{n}), \quad (7)$$

برای $n = 100, 500, 1000$ و $c = 1.5$ مقدار λ_n به ترتیب برابر با $F_{O(\mathbf{X}, G)}^{-1}(0.985)$ ، $F_{O(\mathbf{X}, G)}^{-1}(0.993)$ و $F_{O(\mathbf{X}, G)}^{-1}(0.995)$ است. در عمل به دلیل نامعلوم بودن توزیع F و تابع توزیع تابع دورافتادگی، حتی پس از تعیین مقادیر δ و c نمی‌توان از رابطه (۷) برای تعیین آستانه λ_n استفاده کرد. بدین منظور $\text{out}(\lambda_n, F)$ با شناساگر دورافتادگی بر پایه $\mathbf{X}_1, \dots, \mathbf{X}_n$ (یا ناحیه دورافتادگی نمونه) به صورت

²True positive rate

$\mathbf{x} \in \mathcal{R}^p$ در این صورت دورافتادگی $OR(\lambda_n, F_n) = \{\mathbf{x} : O(\mathbf{x}, F_n) > \lambda_n\}$ برآورد می‌شود. نسبت به نمونه و مقدار آستانه تعیین شده λ_n که بر اساس نمونه حاصل می‌شود، تعیین می‌شود.

الگوریتم ۱. الگوریتم تعیین آستانه λ_n :

- فرض کنید $\mathbf{X}_1, \dots, \mathbf{X}_n$ نمونه تحت بررسی باشد و \hat{G}_n نشان دهنده تابع توزیع تجربی آن باشد.
- گام ۱- مقادیر دورافتادگی نمونه‌ای $O(\mathbf{X}_1, \hat{G}_n), \dots, O(\mathbf{X}_n, \hat{G}_n)$ محاسبه شود.
- گام ۲- بر اساس δ و c که از قبل تعیین شده اند، مقدار $1 - \delta \frac{c}{\sqrt{n}}$ محاسبه شود.
- گام ۳- بر اساس تابع توزیع تجربی $O(\mathbf{X}_1, \hat{G}_n), \dots, O(\mathbf{X}_n, \hat{G}_n)$ چندک $1 - \delta \frac{c}{\sqrt{n}}$ محاسبه شود که این مقدار برابر با λ_n است.

روش شناسائی نقاط دورافتاده که بیان شد، ناوردای آفین است. به عبارت دیگر، بر اساس این روش، نقطه دورافتاده تحت هرگونه تبدیل از محورهای مختصات کماکان به‌عنوان دورافتاده شناسایی می‌شود.

گزاره ۱. اگر تابع ژرفای به کار رفته در تعریف تابع دورافتادگی، در ویژگی ناوردایی آفین صدق کند، آنگاه نقطه‌ای که تحت الگوریتم بالا دورافتاده شناسایی شود، تحت هرگونه تبدیل از محورهای مختصات کماکان به‌عنوان دورافتاده شناسایی می‌شود.

برهان: فرض کنید $\mathbf{X}_1, \dots, \mathbf{X}_n$ نمونه تحت بررسی باشد، و نقطه \mathbf{X}_j تحت الگوریتم ۱ به عنوان دورافتاده شناسایی شود. برای $i = 1, \dots, n$ تبدیل $\mathbf{Y}_i = \mathbf{A}\mathbf{X}_i + \mathbf{b}$ که در آن \mathbf{A} ماتریس $p \times p$ نامنفرد و \mathbf{b} بردار p بعدی دلخواه است، در نظر گرفته می‌شود. چنانچه نشان دهیم \mathbf{Y}_j نیز بر اساس نمونه $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ و تحت الگوریتم ۱ به عنوان دورافتاده شناسایی می‌شود، گزاره اثبات شده است. از ویژگی ناوردایی آفین تابع ژرفا، رابطه $O(\mathbf{X}_i, \hat{G}_n) = O(\mathbf{Y}_i, \hat{U}_n)$ ، به ازای $i = 1, \dots, n$ نتیجه خواهد شد، که در آن \hat{U}_n تابع توزیع تجربی $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ است. بر اساس الگوریتم ۱ مقدار آستانه بر اساس \mathbf{Y}_i ها نیز برابر با λ_n است. در نتیجه \mathbf{Y}_j به عنوان دورافتاده شناسایی می‌شود.

۴ مطالعه شبیه‌سازی

در این بخش عملکرد روش بیان شده در بخش قبل بر اساس توابع ژرفای مختلف مورد ارزیابی و مقایسه قرار می‌گیرد. شناسایی نقاط دورافتاده بر اساس توابع ژرفای ماهالانوبیس بر اساس برآوردگر بردار میانگین

نمونه‌ای و ماتریس واریانس کوواریانس نمونه‌ای، ژرفای ماهالانوبیس بر اساس برآوردگرهای استوار میانه و میانه قدرمطلق انحراف نقاط از میانه، فضائی، تصویرمبنا، نیم‌فضا و سادکی محاسبه شده و به ترتیب با MO, RMO, SO, PO, HO و LO نشان داده می‌شوند. عملکرد این شناساگرها با دو شناساگر رقیب معرفی شده توسط راک و وودراف (۱۹۹۶) و سرفلینگ و موزوندر (۲۰۱۳) که به ترتیب با MD و $RMSP$ نشان داده می‌شوند، مقایسه می‌شود. همان‌گونه که در بخش مقدمه بیان شد، MD بر اساس فاصله ماهالانوبیس تعریف شده و منحنی تراز بیضوی برای مشاهدات در نظر می‌گیرد. از طرف دیگر، $RMSP$ با تعمیم MD و به کارگیری روش جستجوی تصویر تعریف شده است.

از توزیع نرمال استاندارد دو متغیره به عنوان یک توزیع متقارن بیضوی، توزیع یکنواخت دو متغیره بر مربع $[0, 1]^2$ به عنوان یک توزیع متقارن مرکزی که تقارن بیضوی ندارد، و توزیع چوله نرمال دو متغیره استاندارد با پارامتر چولگی 0.25 در هر دو مؤلفه، نمونه‌هایی با اندازه 100 تولید می‌شود. به منظور درک شهودی بهتر این توزیع‌ها، نمودار منحنی تراز و تابع چگالی آن‌ها در شکل ۱ ارائه شده است. یک مدل آلودگی با $\delta = 0.1$ و $c = 1.5$ در نظر گرفته می‌شود. به عبارت دیگر با انتخاب این مقادیر، نرخ مثبت درست تقریبی تحت مدل آلوده و نرخ مثبت غلط تقریبی تحت مدل غیرآلوده به ترتیب برابر با $\varepsilon_{100} = 0.15$ و $\alpha_{100} = 0.15$ حاصل می‌شوند. بنابراین مقدار آستانه تجربی به صورت چندک 0.885 توزیع تجربی تابع دورافتادگی نمونه‌های اولیه تولید شده است. دو الگو برای آلوده کردن 15 درصد از مشاهدات در نظر گرفته می‌شود.

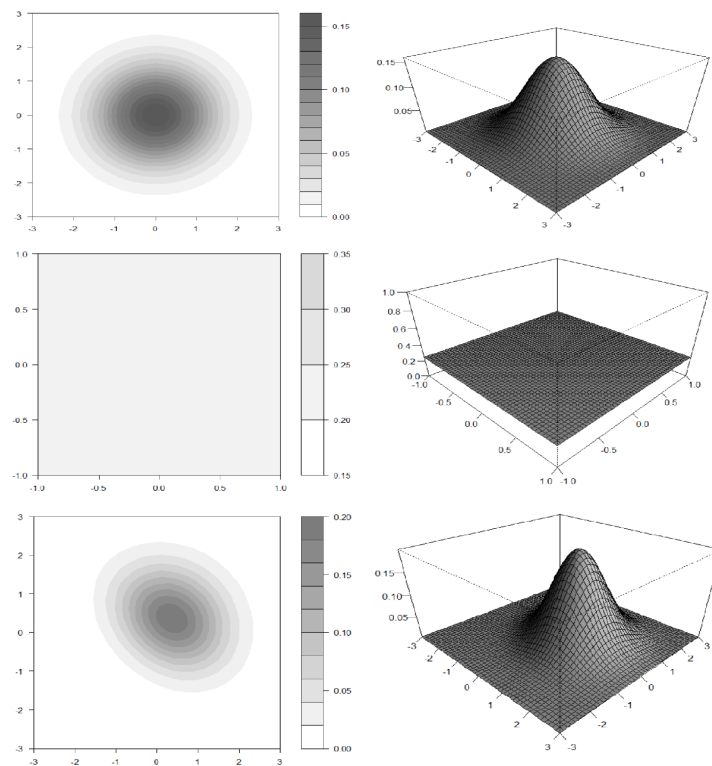
الگوی A: مشاهدات $\mathbf{X}_1, \dots, \mathbf{X}_n$ با $K\mathbf{X}_1, \dots, K\mathbf{X}_n$ که K یک عدد ثابت است، جایگزین می‌شوند. در مطالعات شبیه سازی این بخش مقدار K برابر با 5 در نظر گرفته می‌شود. این نوع آلوده کردن بر اساس تغییر در مقیاس مشاهدات اعمال می‌شود.

الگوی B: مشاهدات $\mathbf{X}_1, \dots, \mathbf{X}_n$ با $\mathbf{X}_1 + \mathbf{K}, \dots, \mathbf{X}_n + \mathbf{K}$ که \mathbf{K} یک بردار دوبعدی ثابت است، جایگزین می‌شوند. در مطالعات شبیه سازی این بخش K برداری با مؤلفه‌های برابر با 4 در نظر گرفته می‌شود. این نوع آلوده کردن بر اساس تغییر در مکان مشاهدات اعمال می‌شود.

دو معیار حساسیت^۱ و دقت^۲ که توسط فن (۲۰۱۶) برای ارزیابی عملکرد روش‌های شناسائی نقاط دورافتاده به کار رفته‌اند، برای شناساگرهای به کار رفته در این مقاله محاسبه می‌شوند. بر اساس جدول ۱، حساسیت و دقت به ترتیب برابر با $\frac{a}{n_1}$ و $\frac{d}{n_2}$ هستند. به عبارت دیگر، حساسیت و دقت به ترتیب نشان‌دهنده توانائی شناسایی صحیح یک نقطه دورافتاده به عنوان دورافتاده و توانائی شناسائی صحیح

¹Sensitivity

²Specificity



شکل ۱. چپ- نمودار تراز و راست- از بالا به پایین نمودارهای توابع چگالی توزیع‌های نرمال دومتغیره استاندارد، یکنواخت دو متغیره بر مربع $[0, 1]^2$ و چوله نرمال دومتغیره استاندارد با پارامتر چولگی ۰.۲۵.

جدول ۱. محاسبه معیارهای حساسیت و دقت

وضعیت	دورافتاده واقعی	غیردورافتاده واقعی
شناسائی شده به عنوان دورافتاده	a	b
شناسائی شده به عنوان غیردورافتاده	c	d
	$n_1 = a + c$	$n_2 = b + d$

یک نقطه غیردورافتاده به عنوان غیردورافتاده، هستند. حال بر اساس موارد بیان شده روند شبیه‌سازی به صورت زیر انجام می‌شود:

۱- از توزیع مورد نظر نمونه‌ای به اندازه ۱۰۰ استخراج می‌شود.

۲- اگر \hat{G}_n نشان دهنده تابع توزیع تجربی نمونه باشد. آنگاه $O(\mathbf{X}_1, \hat{G}_n), \dots, O(\mathbf{X}_n, \hat{G}_n)$ بر اساس تابع دورافتادگی مورد نظر محاسبه می‌شوند.

۳- چندک ۰/۹۸۵ نمونه‌ای $O(\mathbf{X}_1, \hat{G}_n), \dots, O(\mathbf{X}_n, \hat{G}_n)$ محاسبه می‌شود که این مقدار برابر با λ_n است.

۴- الگوی آلودگی مورد نظر (الگوی A یا B) برای جایگذاری ۱۵ مشاهده آخر از مجموعه نمونه 100 تایی تولید شده در نظر گرفته می‌شود.

۵- تابع دورافتادگی برای هر عضو نمونه (کل اعضا شامل آلوده شده و غیرآلوده) محاسبه شده و با مقایسه مقدار دورافتادگی هر عضو با آستانه موردنظر، معیارهای حساسیت و دقت حاصل می‌شوند.

۶- مراحل ۱ تا ۵، ۱۰۰۰ بار تکرار می‌شود و میانگین موارد حاصل شده در مرحله ۴ به عنوان برآورد نهایی معیارهای حساسیت و دقت گزارش می‌شوند.

نتایج در جدول ۲ برای توزیع نرمال استاندارد دو متغیره، توزیع یکنواخت دو متغیره بر مربع $[0, 1]^2$ و توزیع چوله نرمال دو متغیره استاندارد با پارامتر چولگی 0.25 در هر دو مؤلفه، آمده است. $\mathbb{X}_n^A, \mathbb{X}_n^B$ و \mathbb{X}_n^B به ترتیب نمونه اصلی، نمونه آلوده شده تحت الگوی A و نمونه آلوده شده تحت الگوی B هستند. برای ستون‌های \mathbb{X}_n انتظار می‌رود، مقدار حساسیت حدود 0.15 و مقدار دقت حدود 0.985 باشد، که بررسی نتایج، عملکرد قابل قبول همه شناساگرها بجز شناساگرهای مبتنی توابع دورافتادگی نیم‌فضا و سادگی را نشان می‌دهد. برای \mathbb{X}_n^A و \mathbb{X}_n^B واضح است، هرچه مقادیر حساسیت و دقت به یک نزدیک‌تر باشند، نشان دهنده عملکرد بهتر روش موردنظر است. همانطور که ملاحظه می‌شود، عملکرد شناساگرهای مبتنی بر توابع دورافتادگی، تحت تاثیر تابع ژرفای به کار رفته در آن‌ها است. واضح است تغییر و آلودگی در داده‌ها، بیشترین تاثیر را بر روی روش‌هایی که مبتنی بر برآوردگرهای استوار نیستند، مانند MO ، SO و MD که بر اساس برآوردگرهای ناستوار پارامتر مکان و پراکندگی تعریف می‌شوند، دارد. نتایج نشان می‌دهند، همان‌گونه که انتظار می‌رود، شناساگر رقیب MD که بر اساس فاصله ماهالانوبیس تعریف می‌شود، عملکرد مشابه با MO دارد. از طرف دیگر، برای توزیع‌های غیربیضوی یکنواخت و چوله نرمال، $RMSP$ عملکرد بهتری نسبت به MD دارد. از طرف دیگر ژرفای سادگی و نیم‌فضا بر اساس موقعیت نسبی یک نقطه نسبت به تابع توزیع تعریف می‌شوند. بنابراین شناساگرهای تعریف شده بر اساس آن‌ها ساختار هندسی نمونه را در نظر می‌گیرند، اما عملکرد مناسبی، در شناسایی صحیح دورافتادگی و یا عدم دورافتادگی نقاط ندارند.

شناساگرهای معرفی شده بر اساس تابع ژرفای تصویرمبنا، PO ، و تابع ژرفای ماهالانوبیس با برآوردگرهای استوار، RMO ، تحت هر دو الگوی آلودگی A و B عملکرد بسیار بهتری نسبت به سایر شناساگرها دارند. همچنین لازم به ذکر است، برخلاف توابع ژرفای سادگی و نیم‌فضا، محاسبه توابع ژرفای

جدول ۰۲. حساسیت و دقت تجربی برای توزیع‌های دو متغیره مختلف و مدل‌های آلودگی A و B

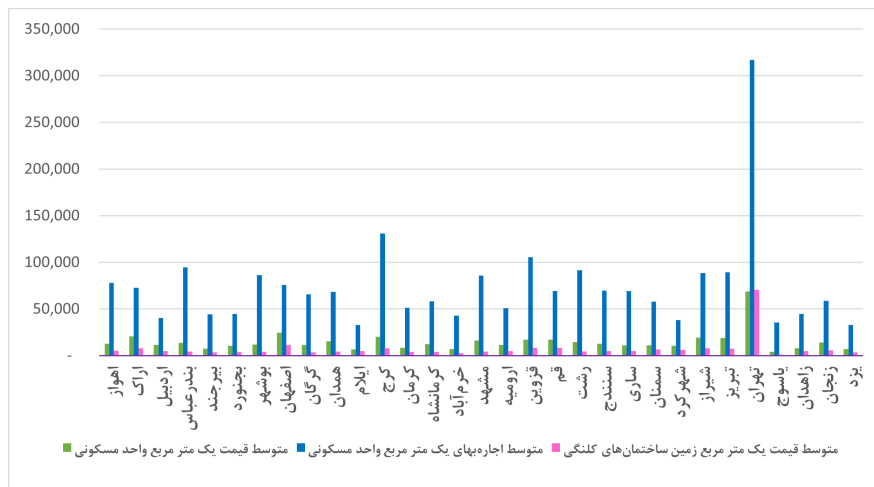
دقت			حساسیت			تابع دورافتادگی	توزیع دو متغیره
\bar{X}_n^B	\bar{X}_n^A	\bar{X}_n	\bar{X}_n^B	\bar{X}_n^A	\bar{X}_n		
۰.۹۹۴	۱	۰.۹۸۰	۰.۹۱۰	۰.۵۰۴	۰.۰۱۷	MO	نرمال
۰.۹۸۸	۰.۹۹۶	۰.۹۸۰	۰.۹۱۸	۰.۷۸۹	۰.۰۱۷	RMO	
۰.۹۸۴	۱	۰.۹۸۰	۰.۷۳۰	۰.۵۳۳	۰.۰۱۷	SO	
۰.۹۹۲	۰.۹۹۶	۰.۹۸۰	۰.۹۵۴	۰.۷۹۶	۰.۰۱۸	PO	
۰.۹۳۶	۰.۹۹۹	۰.۹۱۰	۰.۲۴۴	۰.۴۱۰	۰.۰۸۷	HO	
۰.۹۳۵	۰.۹۹۹	۰.۹۱۰	۰.۲۴۴	۰.۴۱۰	۰.۰۸۷	LO	
۰.۹۹۴	۱	۰.۹۸۰	۰.۹۱۱	۰.۵۰۹	۰.۰۱۷	MD	
۰.۹۹۴	۱	۰.۹۸۰	۰.۹۰۹	۰.۵۰۲	۰.۰۱۶	RMSP	
۰.۹۸۷	۱	۰.۹۸۰	۰.۷۴۹	۰.۷۶۶	۰.۰۲۲	MO	یکنواخت
۰.۸۳۵	۰.۹۹۹	۰.۹۸۰	۱	۰.۹۲۰	۰.۰۲۲	RMO	
۰.۹۵۰	۱	۰.۹۸۰	۰.۲۵۰	۰.۸۰۹	۰.۰۲۲	SO	
۰.۹۸۴	۰.۹۹۹	۰.۹۸۰	۱	۰.۹۱۶	۰.۰۲۱	PO	
۰.۹۲۰	۱	۰.۸۸۲	۰.۲۹۲	۰.۴۶۳	۰.۱۱۸	HO	
۰.۹۲۰	۱	۰.۸۸۱	۰.۲۹۳	۰.۴۶۴	۰.۱۱۸	LO	
۰.۹۸۱	۱	۰.۹۷۷	۰.۷۳۹	۰.۷۵۶	۰.۰۲۱	MD	
۰.۹۸۳	۰.۹۹۹	۰.۹۸۰	۰.۹۵۲	۰.۸۴۶	۰.۰۱۹	RMSP	
۰.۹۹۴	۱	۰.۹۸۰	۰.۲۰۰	۰.۵۱۱	۰.۰۲۰	MO	چوله‌نرمال
۰.۹۸۶	۰.۹۹۶	۰.۹۸۰	۰.۹۹۵	۰.۷۹۶	۰.۰۲۰	RMO	
۰.۹۸۴	۱	۰.۹۸۰	۰.۰۷۳	۰.۵۳۶	۰.۰۱۹	SO	
۰.۹۸۴	۰.۹۹۶	۰.۹۸۰	۰.۹۷۹	۰.۸۰۳	۰.۰۲۱	PO	
۰.۹۳۶	۰.۹۹۸	۰.۹۱۱	۰.۲۴۶	۰.۴۰۷	۰.۰۸۸	HO	
۰.۹۳۶	۰.۹۹۸	۰.۹۱۱	۰.۲۴۶	۰.۴۰۶	۰.۰۸۸	LO	
۰.۹۹۰	۱	۰.۹۸۰	۰.۷۹۳	۰.۶۱۱	۰.۰۲۱	MD	
۰.۹۸۳	۰.۹۹۹	۰.۹۷۹	۰.۹۲۱	۰.۷۶۳	۰.۰۲۰	RMSP	

ماهالانوبیس و تصویرمبنا دارای پیچیدگی محاسباتی نیست. بنابراین به طور کلی می‌توان روش ناوردای آفین پیشنهادی بر اساس توابع ژرفای تصویر-مبنا و ماهالانوبیس با برآوردگرهای استوار را از جهات توانائی شناسائی نقاط دورافتاده و مرتبه محاسباتی به عنوان مطلوب‌ترین شناساگرها در مقایسه با سایر شناساگرهای مبتنی بر تابع ژرفا و همچنین شناساگرهای رقیب در نظر گرفت.

۵ تحلیل داده‌های مسکن

مرکز آمار ایران، برآورد پارامترهایی نظیر متوسط قیمت یک مترمربع واحد مسکونی، متوسط اجاره‌بهای یک متر مربع واحد مسکونی و متوسط قیمت یک مترمربع زمین ساختمان‌های مسکونی کلنگی را از طریق طرح‌های نمونه‌گیری در طی چندین دوره زمانی ارائه داده است. ساختمان مسکونی کلنگی، منظور آن دسته از ساختمان‌های مسکونی است که به منظور تخریب کامل بنا و احداث بنای جدید بر روی زمین آن، مورد

معامله قرار گرفته است (این دسته از ساختمان‌ها می‌توانند پس از خرید مورد بازسازی و استفاده قرار گیرند). در این بخش به دنبال شناسایی شهرهایی از کشور ایران هستیم که به طور کلی در زمینه ارزش مالی واحد مسکونی متفاوت از سایر شهرها هستند. سه متغیر متوسط قیمت یک مترمربع واحد مسکونی، متوسط اجاره‌بهای یک متر مربع واحد مسکونی و متوسط قیمت زمین ساختمان‌های کلنگی در سال ۱۳۹۷ به تفکیک شهرهای منتخب اهواز، اراک، اردبیل، بندرعباس، بیرجند، بجنورد، بوشهر، اصفهان، گرگان، همدان، ایلام، کرج، کرمان، کرمانشاه، خرم‌آباد، مشهد، ارومیه، قزوین، قم، رشت، سنندج، ساری، سمنان، شهرکرد، شیراز، تبریز، تهران، یاسوج، زاهدان، زنجان و یزد، در نظر گرفته می‌شوند. در شکل ۲ نمودار داده‌های موردنظر مشاهده می‌شود. همان‌گونه که مشاهده می‌کنید در هر سه متغیر شهر تهران اختلاف



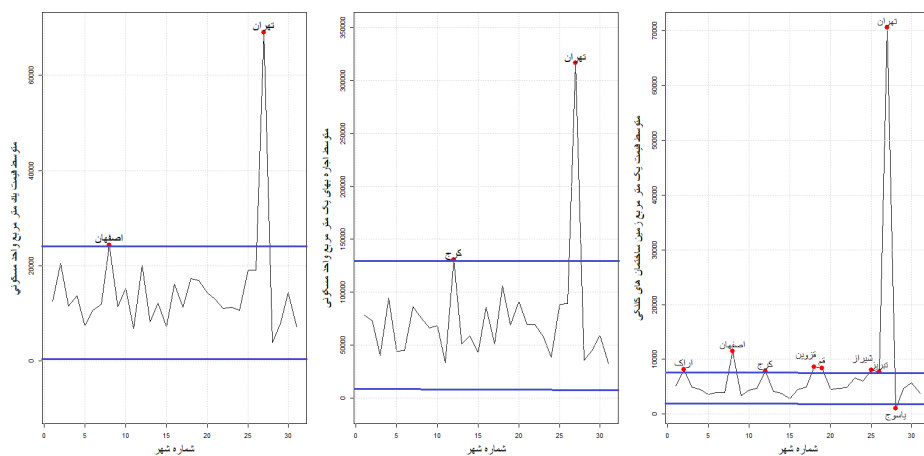
شکل ۲. متوسط قیمت یک مترمربع زمین ساختمان‌های کلنگی، متوسط قیمت اجاره مسکن و متوسط قیمت یک مترمربع واحد مسکونی در سال ۱۳۹۷ به تفکیک شهرهای منتخب

معنی‌داری با سایر شهرها دارد، به‌طوری‌که در مقایسه با این اختلاف، تفاوت میان سایر شهرها ناچیز به نظر می‌رسد. ابتدا، بر اساس روش‌ها و آزمون‌های معروف در حالت یک‌متغیره به شناسایی نقاط دورافتاده می‌پردازیم. ابتدا بر اساس آزمون کولموگروف-اسمیرنوف و شاپیرو-ویلک نرمال بودن متغیرها بررسی شده است که نرمال بودن هیچ‌یک از سه متغیر مورد نظر معنی‌دار نشد. بنابراین با استفاده از روش‌های استوار یک‌متغیره به دنبال شناسایی نقاط دورافتاده یک‌متغیره هستیم. یکی از روش‌های استوار، بازه‌ای با مرکز میانه مشاهدات و کران‌هایی بر اساس میانه انحراف نقاط از میانه است. به عبارت دیگر اگر Med نشان دهنده میانه مشاهدات یک‌متغیره و MAD میانه قدرمطلق انحراف نقاط از میانه باشد، که در بخش ۲ به

ترتیب در معادلات (۲) و (۳) تعریف شده‌اند، آنگاه نقاطی که در خارج از بازه

$$(Med - 3MAD, Med + 3MAD) \quad (۸)$$

قرار می‌گیرند، به عنوان دورافتاده در نظر گرفته می‌شوند. بر اساس ترتیبی که در بالا اسم شهرها آمده است، شماره ۱ تا ۳۱ به آن‌ها اختصاص داده شده است. در شکل ۳ برای هر یک از متغیرها کران پائین و بالای این بازه با خطوط افقی نمایش داده شده است. همانگونه که مشاهده می‌شود، سه شهر تهران، اصفهان و کرج برای حداقل دو متغیر در خارج از بازه قرار گرفته‌اند و به عنوان دورافتاده در نظر گرفته می‌شوند. از طرف دیگر شهر یاسوج برای هر سه متغیر دارای کمترین مقدار در بین شهرها است و برای متغیر متوسط قیمت یک مترمربع زمین ساختمان‌های کلنگی به عنوان نقطه دورافتاده شناسائی شده است. حال آزمون‌های



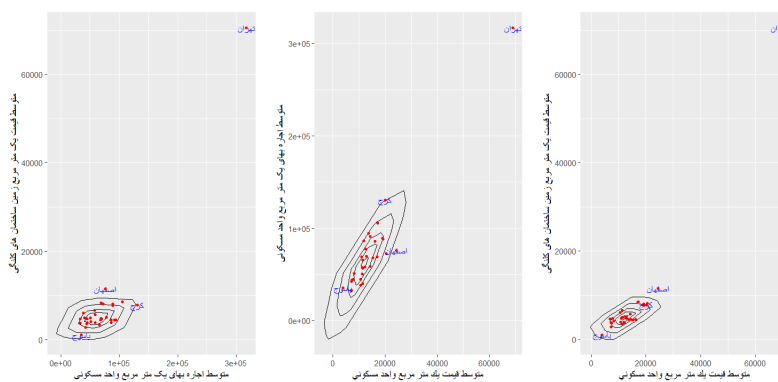
شکل ۳. نمودار پراکندگی به تفکیک متغیرها، خطوط افقی کران‌های پائین و بالای بازه تعریف شده در (۸) را برای هر متغیر نشان می‌دهند.

معروف ناپارامتری یک‌متغیره شناسائی نقاط دورافتاده، **گرایس** (۱۹۵۰) و **کای‌دو** (دیکسون، ۱۹۵۰)، نیز برای هر یک از متغیرها به کار می‌رود. همان‌گونه که در جدول ۳ ملاحظه می‌شود، شهر تهران بر اساس هر دو آزمون و برای هر سه متغیر به عنوان دورافتاده شناسائی شده است. بر اساس نتایج می‌توان استدلال کرد که شهرهای اصفهان و یاسوج نیز از نظر دورافتادگی نسبت به سایر شهرها در رتبه‌های بعدی قرار گرفته‌اند.

جدول ۳. شهرهای شناسائی شده به عنوان دورافتاده بر اساس آزمون گرایس و کای دو

متغیر	آزمون گرایس	آزمون کای دو
متوسط قیمت یک مترمربع واحد مسکونی	تهران	تهران، اصفهان، یاسوج
متوسط اجاره‌بهای یک متر مربع واحد مسکونی	تهران	تهران، کرج، قزوین
متوسط قیمت یک مترمربع زمین ساختمان‌های کلنگی	تهران، اصفهان	تهران، اصفهان، یاسوج

حال روش ناپارامتری چندمتغیره بیان شده در این مقاله بر روی داده‌ها به‌کار برده می‌شود. با توجه به نتایج شبیه‌سازی و عملکرد مناسب تابع ژرفای تصویر-مبنا، تابع دورافتادگی بر اساس آن محاسبه می‌شود. در شکل ۴، خطوط تراز که در معادله (۴) تعریف شده‌اند، بر اساس مقدار تابع دورافتادگی تصویر-مبنا برای داده‌های دومتغیره و مقادیر t برابر با $0.07, 0.06, 0.05$ و 0.3 برای هر یک از دو متغیر ممکن رسم شده است. همان‌گونه که مشاهده می‌شود، شهرهای تهران و اصفهان در خارج از منحنی تراز 0.07 و شهر کرج نیز بر روی این منحنی قرار دارد. همچنین مقدار تابع دورافتادگی تصویر-مبنا برای داده‌های سه‌متغیره محاسبه شده است و مقادیر آن در جدول ۴ آمده است. با در نظر گرفتن نرخ مثبت نادرست 0.05 ، مقدار آستانه تجربی به صورت چندک 0.85 توزیع تجربی مقادیر دورافتادگی به دست آمده است. چندک تجربی 0.85 مقادیر دورافتادگی برابر با 0.861 است. بنابراین بر اساس مقادیر دورافتادگی شهرهای تهران و اصفهان از نظر ارزش مالی مسکن نسبت به سایر شهرهای کشور، دورافتاده محسوب می‌شوند.



شکل ۴. منحنی تراز $0.07, 0.06, 0.05$ و 0.3 بر اساس تابع دورافتادگی تصویر-مبنا

جدول ۴. مقدار تابع دورافتادگی تصویر-مبنا

شهر	تهران	اصفهان	کرج	اراک	قزوین	یاسوج	قم	بندعباس	شیراز	بوشهر	تهریز
PO	۰.۹۸۴	۰.۸۷۹	۰.۸۴۶	۰.۸۰۱	۰.۸۱۰	۰.۷۹۵	۰.۷۸۷	۰.۷۶۸	۰.۷۵۷	۰.۷۵۴	۰.۷۴۶
شهر	رشت	همدان	مشهد	ایلام	خرم آباد	شهرکرد	اردبیل	یزد	زاهدان	سمنان	بیرجند
PO	۰.۷۴۵	۰.۷۴۲	۰.۷۳۴	۰.۷۰۹	۰.۶۹۲	۰.۶۸۷	۰.۶۸۵	۰.۶۶۸	۰.۶۶۷	۰.۶۶۳	۰.۶۵۷
شهر	اهواز	زنجان	کرمانشاه	بجنورد	گرگان	کرمان	ساری	سندج	ارومیه		
PO	۰.۶۴۷	۰.۶۴۳	۰.۶۳۹	۰.۶۲۹	۰.۶۱۷	۰.۶۰۳	۰.۵۹۴	۰.۵۰۰	۰.۴۷۹		

۶ بحث و نتیجه‌گیری

در این مقاله، یک روش برای شناسایی نقاط دورافتاده چندمتغیره بر اساس تابع ژرفا معرفی شد. روش پیشنهادی دارای ویژگی مطلوب ناوردای آفین است و پیاده‌سازی آن نسبت به بیشتر روش‌های چندمتغیره که دارای پیچیدگی محاسباتی هستند، ساده‌تر است. بر اساس مطالعات شبیه‌سازی، عملکرد روش پیشنهادی بر اساس توابع ژرفای مختلف مورد بررسی قرار گرفت و همچنین با شناساگرهای رقیب مقایسه شد. نتایج شبیه‌سازی نشان داد، روش پیشنهادی بر اساس توابع ژرفای تصویر-مبنا و ماهالانوبیس با برآوردگرهای استوار، دارای عملکرد مطلوب‌تری در مقایسه با سایر شناساگرها است. در ادامه، روش پیشنهادی بر اساس تابع ژرفای تصویر-مبنا برای داده‌های ارزش مالی واحد مسکونی ایران در سال ۱۳۹۷ به کار رفت. نتایج نشان داد، شهرهای تهران و اصفهان را می‌توان از نظر ارزش مالی واحد مسکونی متفاوت از سایر شهرهای ایران در نظر گرفت.

مراجع

- روزبه، م. و امینی، م. (۱۳۹۸)، برآوردگر استوار مرزبندی شده تعمیم‌یافته محتمل در مدل رگرسیون نیمه‌پارامتری، مجله علوم آماری، (۲)۱۳، ۴۶۰-۴۴۱.
- روزبه، م. و معنوی، م. (۱۳۹۹)، مدل‌سازی سن تقویمی به روش رگرسیون ستیغی کمترین توان‌های دوم پیراسته، مجله علوم آماری، (۲)۱۴، ۴۲۸-۴۰۹.
- Barnett, V. and Lewis, T. (1994), *Outliers in Statistical Data*, John Wiley & Sons.
- Chen, Y., Dang, X., Peng, H. and Bart, L. (2008), *Outlier Detection with*

the Kernelized Spatial Depth Function, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **31**(2), 288-305.

Dang, X. and Serfling, R. (2010), Nonparametric Depth-based Multivariate Outlier Identifiers, and Masking Robustness Properties, *Journal of Statistical Planning and Inference*, **140**(1), 198–213.

Dehghan, S. and Faridrohani, M. R. (2019), Affine Invariant Depth-Based Tests for the Multivariate One-Sample Location Problem, *Test*, **28**(3), 671-693.

Dixon, W. J. (1950), Analysis of Extreme Values. *The Annals of Mathematical Statistics*, **21**(4), 488-506.

Dovoedo, Y. H. and Chakraborti, S. (2013), Outlier Detection for Multivariate Skew-Normal Data: a Comparative Study, *Journal of Statistical Computation and Simulation*, **83**(4), 773-783.

Fan, Yi. (2016), *New Nonparametric Approaches for Multivariate and Functional Data Analysis in Outlier Detection, Construction of Tolerance Tubes, and Clustering*, Diss. Rutgers University-Graduate School-New Brunswick.

Grubbs, F. E. (1950), Sample Criteria for Testing Outlying Observations. *The Annals of Mathematical Statistics*, **21**(1), 27-58.

Liu, R. Y. (1988), On a Notion of Simplicial Depth, *Proceedings of the National Academy of Sciences*, **85**(6), 1732–1734.

Liu, R. Y. and Singh, K. (1993), A Quality Index Based on Data Depth and Multivariate Rank Tests, *Journal of the American Statistical Association*, **88**(421), 252–260.

- Hawkins, D. M. (1980), *Identification of Outliers*, Chapman and Hall.
- Kannan, K., Senthamarai, and Manoj, K. (2015), Outlier Detection in Multivariate Data, *Applied Mathematical Sciences*, **9**(47), 2317-2324.
- Miller, J. (1991), Reaction Time Analysis with Outlier Exclusion: Bias Varies with Sample Size. *The Quarterly Journal of Experimental Psychology*, **43**(4), 907–912.
- Serfling, R. and Mazumder, S. (2013), Computationally Easy Outlier Detection via Projection Pursuit with Finitely Many Directions, *Journal of Nonparametric Statistics*, **25**(2), 447-461.
- Rocke, D. M. and Woodruff, D. L. (1996), Identification of Outliers in Multivariate Data, *Journal of the American Statistical Association*, **91**(435), 1047-1061.
- Tukey, J. W. (1975), Mathematics and the Picturing of Data, *Proceedings of the international congress of mathematicians*, **2**, 523–531.
- Vardi, Y. and Zhang, C. H. (2000), The Multivariate L1-Median and Associated Data Depth, *Proceedings of the National Academy of Sciences*, **97**(4), 1423-1426.
- Zuo, Y. and Serfling, R. (2000), General Notions of Statistical Depth Function, *Annals of statistics*, **28**(2), 461–482.
- Zuo, Y. (2003), Projection-based Depth Functions and Associated Medians. *The Annals of Statistics*, **31**(5) 1460–1490.

Journal of Statistical Sciences, Autumn and Winter, 2021
Vol. 15, No. 2, pp 443-462
DOI: 10.29252/jss.15.2.443

Multivariate Outlier Detection Based on Depth-Based Outlyingness Function

Dehghan, S., Faridrohani, M. R.

Department of Statistics, Faculty of Mathematical Sciences, University of Shahid Beheshti, Tehran, Iran.

Abstract: The concept of data depth has provided a helpful tool for non-parametric multivariate statistical inference by taking into account the geometry of the multivariate data and ordering them. Indeed, depth functions provide a natural centre-outward order of multivariate points relative to a multivariate distribution or a given sample. Since the outlyingness of issues is inevitably related to data ranks, the centre-outward ordering could provide an algorithm for outlier detection. In this paper, based on the data depth concept, an affine invariant method is defined to identify outlier observations. The affine invariance property ensures that the identification of outlier points does not depend on the underlying coordinate system and measurement scales. This method is easier to implement than most other multivariate methods. Based on the simulation studies, the performance of the proposed method based on different depth functions has been studied. Finally, the described method is applied to the residential houses' financial values of some cities of Iran in ۲۰۱۸.

Keywords: Depth function, Outlyingness function, Outlier, Affine invariance.

Mathematics Subject Classification (2010): 62G99, 62H99.