

## استفاده از روش‌های تک‌متغیره و چندمتغیره به‌منظور شناسایی داده‌های پرت در منشأیابی رسوبات، مطالعه موردی: حوزه آبخیز تنگ بستانک

احمد نوحه‌گر<sup>۱</sup>، محمد کاظمی<sup>۲\*</sup>، سیدجواد احمدی<sup>۳</sup>، حمیدغلامی<sup>۴</sup> و رسول مهدوی<sup>۵</sup>  
<sup>۱</sup>استاد، دانشکده محیط زیست، دانشگاه تهران، <sup>۲</sup>دانشجوی دکتری آبخیزداری، دانشکده منابع طبیعی، دانشگاه هرمزگان، <sup>۳</sup>دانشیار،  
پژوهشکده چرخه سوخت سازمان انرژی اتمی و <sup>۴</sup>استادیار، دانشکده منابع طبیعی، دانشگاه هرمزگان

تاریخ پذیرش: ۱۳۹۵/۱۲/۰۹

تاریخ دریافت: ۱۳۹۵/۰۶/۲۷

### چکیده

کارایی روش منشأیابی با ردیاب‌ها یا انگشت‌نگاری به‌عنوان روشی موفق و مؤثر برای تعیین منابع رسوبات به اثبات رسیده است. اولین و مهم‌ترین مرحله روش منشأیابی رسوب، انتخاب ترکیب مناسبی از ردیاب است که قادر به جداسازی منابع رسوب باشند. وجود داده‌های پرت بر انتخاب ترکیب مناسبی از ردیاب‌ها اثر گذاشته، ممکن است مانع انتخاب متغیرهای مهم (ردیاب‌های مؤثر) شده و توان جداسازی یا درصد طبقه‌بندی صحیح را کاهش دهد. بنابراین داده‌های یادشده باید شناسایی و در صورت وجود شواهد کافی دال بر پرت بودن، نسبت به تصحیح یا حذف آن‌ها اقدام شود. در مطالعه حاضر هدف شناسایی داده‌های پرت در بین مجموعه ردیاب‌های اندازه‌گیری شده در حوزه آبخیز تنگ بستانک، برای تشخیص بهترین ترکیب ردیاب‌ها بود. بر این اساس از روش‌های تک‌متغیره شناسایی داده‌های پرت همچون آزمون گراب، آزمون گوس، آزمون دیکسون، نمودار جعبه‌ای، میانه به اضافه یا منهای میانه انحراف‌های تمام داده‌ها از میانه و میانگین به اضافه و منهای سه برابر انحراف از معیار داده‌ها و نیز از روش‌های چندمتغیره شناسایی داده‌های پرت همچون تحلیل مؤلفه‌های اصلی، فاصله ماهالانوبیس، مربع فاصله ماهالانوبیس، نمودار چندک مربع فاصله ماهالانوبیس به روی درجه آزادی در برابر توزیع مربع کای، نمودارهای جعبه‌ای مربع فاصله ماهالانوبیس استفاده شد. در مجموع داده‌ای پرت شناخته شد که دست‌کم نیمی از روش‌های مذکور به پرت بودن آن اذعان داشته باشند. نتایج نشان داد روش میانه به اضافه و منهای میانه انحراف‌های تمام داده‌ها از میانه تعداد بیشتری از داده‌ها را به‌عنوان داده پرت معرفی می‌کند و همچنین روش‌های چندمتغیره اشتراک کمتری در تشخیص داده‌های پرت با یکدیگر دارند. آزمون‌های تک‌متغیره اجماع بهتری نسبت به شناسایی و معرفی داده‌های پرت دارند. برای استفاده از روش‌های تک-متغیره برای شناسایی داده‌های پرت روش‌های میانه به اضافه یا منهای میانه انحراف‌های تمام داده‌ها از میانه، نمودار جعبه‌ای و آزمون دیکسون به‌ترتیب حساسیت آن‌ها پیشنهاد می‌شود. همچنین نتایج نشان داد، بیشینه اجماع روش‌های به‌کار رفته برای روش‌های تک‌متغیره چهار و برای روش‌های چندمتغیره دو مورد هست و در کل اجماع نیمی از روش‌ها برای پرت بودن داده‌ها مشاهده نشد.

واژه‌های کلیدی: آنالیز مؤلفه‌های اصلی، انگشت‌نگاری، ردیاب، فاصله ماهالانوبیس، گراب

**مقدمه**

گذاشته و ممکن است مانع انتخاب متغیرهای مهم شود (Wiegand و همکاران، ۲۰۰۹). همچنین قابل ذکر است که یکی از اصول مدل‌های مرکب در شناسایی مناطق مولد رسوب نرمال بودن متغیرها (ردیاب‌ها) می‌باشد.

در تحلیل‌های چندمتغیره علاوه بر این که هریک از متغیرها باید از توزیع نرمال تبعیت کنند، ترکیب متغیرها نیز باید از توزیع نرمال (نرمال چندمتغیره) پیروی کند از سوی دیگر بیشتر عناصر ژئوشیمیایی که از مهم‌ترین ردیاب‌ها در منشأیابی محسوب می‌شوند از توزیع نرمال تبعیت نمی‌کنند (Zhang و همکاران، ۲۰۰۸) و اگر این چولگی ناشی از عدم وجود داده‌های پرت باشد، مشکلی وجود ندارد یا به عبارتی اگر تبعیت نکردن مجموعه متغیرها از فرض نرمال ناشی از وجود داده‌های پرت نبوده و به داده‌های حد مربوط باشد، نتایج روش‌های چندمتغیره معتبر خواهد بود (Tabachnick و Fidell، ۱۹۹۶). بنابراین، ضروری به نظر می‌رسد که برای صحت و دقت طبقه‌بندی ابتدا داده‌های پرت را شناسایی و سپس وارد مراحل بعد شد (Hair و همکاران، ۱۹۹۸).

داده‌های پرت در عناصر ژئوشیمی و آلی خاک علاوه بر اشتباه یا خطا در اندازه‌گیری ممکن است بر اثر عواملی نظیر کانی‌سازی، آلتراسیون و یا فعالیت‌های انسانی باشد. تاکنون روش‌های زیادی برای تشخیص داده‌های پرت به کار گرفته شده است اما هیچ‌کدام مقبولیت جهانی نیافته‌اند (Reimann و همکاران، ۲۰۰۵). بنابراین برای اطمینان و بررسی بیشتر می‌توان نتایج اجماع چندین روش مختلف را به عنوان ملاک تصمیم‌گیری مد نظر قرار داد. هدف از تحقیق حاضر استفاده از روش‌های تک‌متغیره و چند-متغیره برای شناسایی داده‌های پرت در منشأیابی رسوب در حوزه آبخیز تنگ بستانک می‌باشد.

**مواد و روش‌ها**

**مشخصات منطقه مورد پژوهش:** حوضه مورد مطالعه در این بررسی، تحت عنوان حوزه آبخیز تنگ بستانک در حدود ۸۰ کیلومتری شمال غرب شهرستان شیراز و در موقعیت جغرافیایی  $33^{\circ} 16'$  تا  $33^{\circ} 36'$  و  $52^{\circ} 13'$  تا  $52^{\circ} 33'$  شرقی و  $30^{\circ} 16'$  تا  $30^{\circ} 25'$  شمالی واقع

منشأیابی رسوب از اصول اولیه کنترل و مبارزه با فرسایش خاک محسوب می‌شود، زیرا با شناسایی مناطق برداشت می‌توان به جای پرداختن به معلول‌ها، علت‌ها را شناسایی نمود و فعالیت‌های اجرایی مبارزه با فرسایش را در مناطق برداشت متمرکز کرد (Feng و همکاران، ۲۰۱۱). لازمه اجرای برنامه‌های حفاظت خاک و کنترل رسوب، کسب اطلاعات از اهمیت نسبی منابع رسوب و سهم آن‌ها در تولید رسوب و در نتیجه شناسایی مناطق بحرانی در داخل آبخیز است (Walling و همکاران، ۲۰۰۸). به دلیل مشکلات مرتبط با روش‌های سنتی برای تعیین منابع اولیه رسوب داخل یک آبخیز، روش‌های منشأیابی به‌عنوان یک روش غیرمستقیم جمع‌آوری اطلاعات مورد توجه قرار گرفته است.

تلفیق مدل‌های ترکیبی در ارتباط با ردیاب‌های مرکب در اواخر دهه ۱۹۸۰ و اوایل دهه ۱۹۹۰ به مطالعات ردیابی منبع ورود کرد و امکان به‌دست آوردن تخمین‌های سهم نسبی از منابع مختلف را فراهم نمودند (Walling، ۲۰۰۵). در مطالعات امروزی از روش منشأیابی مرکب استفاده می‌شود، به این شکل که در مرحله اول، ترکیبی مناسب از ردیاب‌ها که قادر به جداسازی منابع رسوب (نظیر واحدهای سنگ‌شناسی، کاربری‌های اراضی و انواع فرسایش) باشند انتخاب می‌شود. در مرحله دوم، ترکیب یاد شده برای تعیین سهم نسبی هر یک از منابع رسوب با استفاده از مدل‌های ترکیبی چندمتغیره به کار می‌رود. منشأیابی مرکب عبارت است از بررسی چندین منبع بالقوه رسوب از طریق چندین گروه متفاوت از خصوصیات که برای افزایش تشخیص بین منابع و اجتناب از منبع رسوب غیرواقعی است و شامل طیف وسیعی از خصوصیات مختلف شناختی می‌باشد (Collins و همکاران، ۲۰۱۰). در این مدل‌های ترکیبی شناسایی مناطق مولد رسوب اولین و مهم‌ترین قدم، تشخیص ترکیب بهینه از عناصر ژئوشیمیایی و آلی می‌باشد. در این راستا وجود داده‌های پرت مانع از شناسایی ترکیب مؤثر و بهینه از عناصر خاکی به‌عنوان ردیاب خواهد شد (Walling و همکاران، ۲۰۰۸؛ Collins و Walling، ۲۰۰۷) و یا به عبارتی وجود داده‌های پرت بر انتخاب ترکیب مناسبی از متغیرها اثر

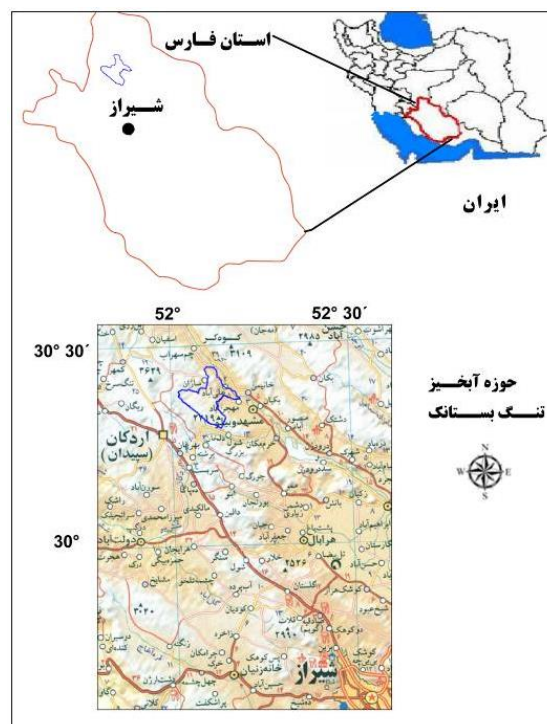
برداشت و نمونه‌ها طوری انتخاب شدند که معرف تغییرات در نوع کاربری‌ها باشند. کلاً تعداد ۴۳ نمونه برداشت شد و بعد از خشک کردن نمونه‌ها در هوای آزاد و دمای اتاق از الک‌های ۶۲/۵، ۷۵، ۱۵۰، ۳۰۰، ۶۰۰، ۱۱۸۰ و ۱۷۰۰ میکرون عبور داده شدند. سپس به‌روی مقدار خاکی که کمتر از ۶۳ میکرون بود عملیات آزمایشگاهی زیر صورت گرفت.

مرحله یک: ابتدا سه گرم از هر نمونه وزن شده و در محلول Aqua Regia (ترکیب اسید کلریدریک (HCl) و اسید نیتریک (HNO<sub>3</sub>) با نسبت ۱:۳ یعنی سه قسمت اسید کلریدریک و یک قسمت اسید نیتریک) هضم شد. در این مطالعه از ۲۱ میلی‌لیتر اسید کلریدریک و هفت میلی‌لیتر اسید نیتریک استفاده شد (Collins و همکاران، ۲۰۱۰ و ۲۰۱۲؛ Nosrati و همکاران، ۲۰۱۱). مرحله دو: سپس نمونه‌ها به مدت دو ساعت در دمای ۹۵ درجه سانتی‌گراد به‌روی حمام آبی قرار داده شد. مرحله سه: در این مرحله ابتدا نمونه‌ها از کاغذ صافی واتمن عبور داده و سپس به‌منظور اطمینان از شفافیت و عاری بودن نمونه‌ها از ذرات معلق دوباره از کاغذ صافی استات سلولز ۰/۲ میکرومتر عبور داده شد. مرحله چهار: اندازه‌گیری به‌وسیله ICP-Mass.

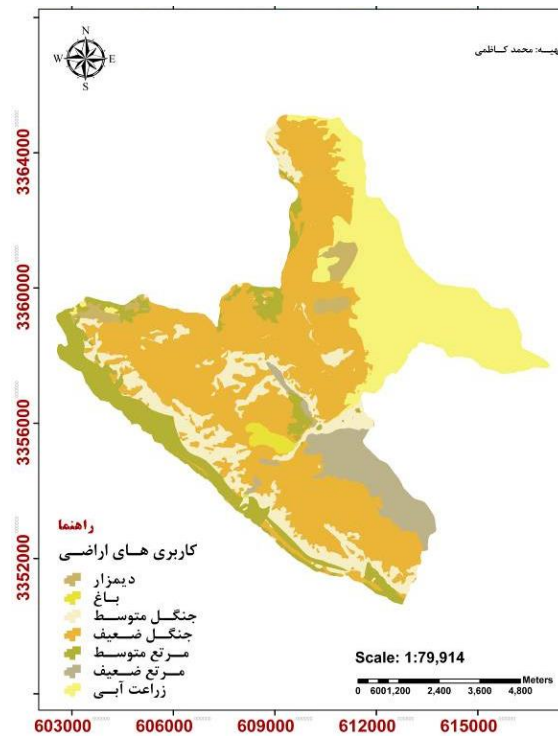
شده است. این حوضه از نظر تقسیمات حوزه‌های آبخیز کشوری، جزء زیرحوضه آبخیز نیریز و شیراز بوده که آب‌های آن پس از وارد شدن به رودخانه کر، وارد دریاچه بختگان می‌شود. میانگین بارش سالانه این حوضه ۶۰۹ میلی‌متر می‌باشد، این حوضه طبق روش اقلیم‌نمای دومارتن اصلاح شده دارای اقلیم مدیترانه‌ای سرد می‌باشد. شکل ۱، موقعیت منطقه و راه‌های دسترسی به آن را نشان می‌دهد.

در این مطالعه نقشه کاربری اراضی تحت کاربری‌های اراضی زراعی، جنگل با پوشش ضعیف، جنگل با پوشش متوسط، باغات، مرتع با پوشش گیاهی ضعیف، مرتع با پوشش گیاهی متوسط و دیم‌زارها به‌عنوان واحدهای مولد رسوب و مطابق با روش الگوریتم بیشینه تشابه (ML) با تصویر ماهواره لندست ۸ سنجنده OLI تهیه شد.

اندازه‌گیری عناصر ژئوشیمیایی: آماده‌سازی نمونه‌ها به‌منظور اندازه‌گیری عناصر ژئوشیمیایی در چهار مرحله به‌ترتیب زیر انجام شد. از هر یک از واحدهای کاربری اراضی نمونه خاک از عمق صفر تا پنج سانتی‌متری و حدوداً به اندازه دو کیلوگرم با یک بیلچه استیل برداشت شد (Walling و همکاران، ۱۹۹۹). از هر واحد کاری با توجه به وسعت دست‌کم پنج نمونه



شکل ۱- موقعیت حوزه آبخیز تنگ بستانک و راه‌های دسترسی به آن



شکل ۲- نقشه واحدهای کاربری اراضی به‌عنوان منابع مختلف تولید کننده رسوب

مشخص کردن داده‌های پرت هست. روش سنتی در این مورد میانگین به‌اضافه و منهای سه برابر انحراف از معیار<sup>۱</sup> می‌باشد (Chiang و همکاران، ۲۰۰۳). داده‌های بزرگ‌تر از میانگین به‌اضافه سه برابر انحراف معیار و کوچک‌تر از میانگین منهای سه برابر انحراف معیار پرت محسوب می‌شوند. در این روش ۹۹/۸ درصد از داده‌ها پرت محسوب نمی‌شوند. روش مذکور دارای دقت پایینی می‌باشد و چون این روش تحت تأثیر داده‌های پرت است (در محاسبه میانگین و انحراف معیار از تمام داده‌ها از جمله داده‌های پرت استفاده می‌شود)، از این‌رو، روش‌های دیگری از جمله میانه به‌اضافه یا منفی میانه انحراف‌های تمام داده‌ها از میانه<sup>۲</sup> و نمودار جعبه‌ای<sup>۳</sup> (Tukey، ۱۹۷۷؛ Reimann و همکاران، ۲۰۰۵) ارائه شده که تحت تأثیر داده‌های پرت قرار نمی‌گیرند.

میانه انحراف تمام داده‌ها از میانه (MAD) از رابطه (۱) محاسبه می‌شود. در منشأیابی رسوبات با استفاده از مدل‌های ترکیبی داده‌های نرمال و یا دارای چولگی (ناشی از داده‌های حد) پیش فرض ورود به این مدل‌ها

در این مطالعه تعداد ۱۷ عنصر شامل ۱۵ عنصر ژئوشیمیایی و دو عنصر آلی استفاده شد. عناصر ژئوشیمیایی شامل باریم، کادمیوم، کروم، مس، لیتیم، منگنز، نئودیوم، نیکل، فسفر، سیلیکون، استرانسیوم، تیتانیوم، وانادیوم، زینک و دو ماده آلی نیتروژن و کربن بودند که در آزمایشگاه سازمان انرژی اتمی کشور اندازه‌گیری شدند. قابل ذکر است که اندازه‌گیری میزان کربن آلی به روش والکر و بلاک و نیتروژن کل به روش کجدال انجام شده است (Alihyai و Behbahanizade، ۱۹۹۳).

**روش‌های تشخیص داده پرت:** روش‌های تشخیص داده پرت را می‌توان در مجموع به سه دسته تقسیم کرد. روش‌های تک‌متغیره، دومتغیره و چندمتغیره (Hair و همکاران، ۱۹۹۸). در تحقیق حاضر از روش-های تک‌متغیره و چندمتغیره استفاده شد.

**روش‌های تک‌متغیره شناسایی داده‌های پرت:** این روش‌ها را می‌توان به دو گروه روش‌های دامنه‌ای و آزمون‌های آماری تقسیم نمود. در روش‌های دامنه توزیع مشاهدات بررسی شده و داده‌های خارج از یک دامنه معین به‌عنوان داده پرت تلقی می‌شوند. مهم-ترین موضوع در این ارتباط تعیین دامنه یاد شده برای

<sup>۱</sup>  $X \pm 3S$

<sup>۲</sup>  $\text{Median} \pm 3\text{MAD}$

<sup>۳</sup> Box plot

این روش فرض صفر این است که مقدار حدی با بقیه نمونه‌های جمعیت تفاوت ندارد و فرض یک این است که مقدار حدی با بقیه نمونه‌های جمعیت تفاوت دارد. آماره آزمون دیکسون از رابطه (۲) محاسبه می‌شود (تعداد نمونه‌ها در هر ردیاب باید بیشینه ۱۰۰ و یا کمتر از ۱۰۰ باشد).

$$Q = \frac{x_n - x_{n-1}}{x_n - x_1} \quad (2)$$

که در آن،  $x_n$  میزان داده حد،  $x_{n-1}$  میزان داده قبل از مقادیر حدی و  $x_1$  کمترین مقدار طبقه‌بندی شده در نمودار شاخه و برگ مقادیر حد می‌باشد.

آزمون‌های دیگری نظیر آزمون گراب<sup>۲</sup> نیز برای تشخیص یک‌متغیره داده‌های پرت استفاده می‌شود (Lalor و Zhang، ۲۰۰۱). در آزمون گراب فرض بر این است که داده‌ها از توزیع نرمال پیروی می‌کنند و اگر توزیع داده‌ها غیر از این و خصوصاً توزیع متقارن باشد (مثلاً توزیع لوگ نرمال)، نتایج این آزمون اشتباه خواهد بود. در آزمون یاد شده، در هر مرحله یک داده پرت تشخیص داده می‌شود. در صورتی که داده پرتی شناسایی شود، داده یاد شده حذف می‌شود و آزمون برای بقیه داده‌ها دوباره انجام می‌شود. این کار آن قدر ادامه می‌یابد تا هیچ داده پرتی وجود نداشته باشد. فرض صفر این است که هیچ نوع داده پرتی وجود ندارد و فرض مخالف این است که دست‌کم یک داده پرت وجود دارد. آماره آزمون گراب (G) از رابطه (۳) محاسبه می‌شود (در این آزمون هم تعداد نمونه‌ها در هر ردیاب باید در نهایت ۱۰۰ باشد).

$$G = \frac{\max |X_i - \bar{X}|}{S} \quad (3)$$

که در آن،  $X_i$  کوچک‌ترین یا بزرگ‌ترین داده،  $\bar{X}$  میانگین داده‌ها و S انحراف از معیار داده‌ها می‌باشند. فرض صفر موقعی رد خواهد شد که (فرض مخالف یا کمینه وجود یک داده پرت) شرط زیر برقرار باشد (رابطه ۴).

$$G > \frac{(n-1)}{\sqrt{n}} \sqrt{\frac{t^2 \frac{\alpha}{(2n, n-2)}}{n-2+t^2 \frac{\alpha}{(2n, n-2)}}} \quad (4)$$

که در آن، n اندازه نمونه و  $t^2 \frac{\alpha}{(2n, n-2)}$  مقدار بحرانی آماره توزیع t استیودنت با درجه آزادی n-2 و سطح معناداری  $\alpha/2n$  می‌باشد. روش دیگری که برای

می‌باشند. اما اگر چولگی و عدم تبعیت از توزیع نرمال داده‌ها به‌خاطر وجود داده‌های پرت باشد، نتایج مدل-های ترکیبی دارای اشکال می‌باشد و لازم است برای انتخاب بهترین ترکیب از ردیاب‌ها (ترکیب بهینه ردیاب‌ها) داده‌های پرت حذف شوند (Tabachnick و Fidell، ۱۹۹۶). به‌عبارت دقیق‌تر، داده‌ها (مقادیر ردیاب‌ها) برای ورود به تابع تحلیل تشخیص (DFA) و طبقه‌بندی صحیح هر ردیاب در هر منبع باید عاری از داده‌های پرت باشند.

$$MAD = 1.482 \text{Median}(X_i - X_{\text{Median}}) \quad (1)$$

مقدار ثابت ۱/۴۸۲ برای تبدیل MAD به برآورد نارایی از انحراف معیار (امید ریاضی انحراف معیار نمونه برابر با انحراف معیار جامعه) داده‌های نرمال است (Chiang و همکاران، ۲۰۰۳). نمودار جعبه‌ای نیز از روش‌های دامنه محسوب می‌شود. این روش نموداری برای نشان دادن موقعیت، پراکنش و چولگی داده‌ها می‌باشد (Tukey، ۱۹۷۷) و به فراوانی برای تشخیص داده‌های پرت استفاده می‌شود (Reimann و همکاران، ۲۰۰۵). این نمودار با استفاده از یک مستطیل (باکس) و دو خط یا میله در دو طرف مستطیل و به‌وسیله میانه و چارک‌های اول ( $Q_1$ ) و سوم ( $Q_3$ ) و کمترین و بیشترین مقادیر رسم می‌شود. طول مستطیل برابر با تفاوت بین چارک سوم و چارک اول (IQR) است. در یک نوع از نمودار جعبه‌ای که از آن برای تشخیص داده‌های پرت استفاده می‌شود، داده‌هایی که کوچک‌تر از  $Q_1 - 1.5IQR$  و نیز بزرگ‌تر از  $Q_3 + 1.5IQR$  باشند جزء داده‌های پرت خفیف و داده‌هایی که کوچک‌تر از  $Q_1 - 3IQR$  و بزرگ‌تر از  $Q_3 + 3IQR$  باشند، پرت قوی محسوب می‌شود. از آزمون‌های آماری نظیر آزمون دیکسون<sup>۱</sup> برای تشخیص داده پرت استفاده می‌شود. در این آزمون با فرض تبعیت داده‌ها از توزیع نرمال به‌شرطی که مقدار محاسباتی Q بزرگ‌تر از Q بحرانی (Q جدول) باشد داده پرت محسوب می‌شود (Alfassi، ۲۰۰۵؛ Rorabacher، ۱۹۹۱). در روش فوق برای شناسایی داده پرت لازم است تا ابتدا مقادیر حد با توجه به نمودار شاخه و برگ شناخته شوند (EPA، ۲۰۰۶). در

<sup>2</sup> Grubbs' test

<sup>1</sup> Dixon' Q test

که در آن،  $X_i$  بردار متغیرها برای مشاهده نام،  $\bar{X}$  بردار میانگین متغیرها (مرکز ثقل مشاهدات) و  $C$  ماتریس کوواریانس نمونه است. خصوصیات مربع فاصله ماهالانوبیس ( $MD^2$ ) به گونه‌ای است که اجازه استفاده از آزمون‌های آماری از جمله  $t$  و آزمون مربع کای را برای بررسی داده‌های پرت می‌دهد. Hair و همکاران (۱۹۹۸) از مقایسه  $MD^2/df$  و توزیع  $t$  استفاده کردند. هرگاه مقدار یاد شده با توجه به درجه آزادی ( $df$ ) (تعداد ردیاب‌ها منهای یک) و سطح معناداری مورد نظر بیشتر از  $t$  جدول باشد، داده مربوطه پرت محسوب می‌شود. بعضی از محققان از جمله Rousseeuw و Van Driessen (۱۹۹۹) از مقایسه  $MD^2$  و توزیع مربع کای استفاده کردند. همچنین، قابل ذکر است که روش‌های گرافیکی (Garrett, ۱۹۸۹؛ Filzmoser و همکاران، ۲۰۰۵؛ Zhang و همکاران، ۲۰۰۸) که فاصله ماهالانوبیس را در مقابل توزیع مربع کای رسم می‌کنند نیز ارائه شده‌اند.

تجزیه به مؤلفه‌های اصلی، روشی برای کاهش تعداد ابعاد داده‌های مورد مطالعه است. در تحلیل یاد شده امکان تبدیل تعداد زیادی از متغیرهای اولیه به تعداد معدودی از متغیرهای جدید (ابعاد یا مؤلفه‌های اصلی) که بیشترین واریانس مشاهده شده در متغیرهای اولیه را بیان کرده باشد بررسی می‌شود. از روش یاد شده به‌منظور کشف روابط و همبستگی بین متغیرها نیز استفاده می‌شود (Hair و همکاران، ۱۹۹۸). در سال‌های اخیر از روش تجزیه به مؤلفه‌های اصلی برای تشخیص داده‌های پرت نیز استفاده کرده‌اند (Zhang و همکاران، ۱۹۹۹؛ Chiang و همکاران، ۲۰۰۳؛ Lalor و Zhang، ۲۰۰۱) در این پژوهش از این روش نیز به‌علت توان بالای آن در تشخیص داده‌های پرت، سادگی و قابل انجام بودن آن با نرم افزارهای آماری استفاده شد. تمام مؤلفه‌های اصلی که بیش از ۹۹ درصد واریانس متغیرها (ردیاب‌ها) را بیان کرده باشد انتخاب می‌شوند. از فاصله نمرات نمونه به‌عنوان معیاری برای بررسی داده‌های پرت که از رابطه زیر (رابطه ۷) برآورد می‌شود، استفاده می‌شود (Zhang و همکاران، ۱۹۹۹).

$$DSC_i = \sqrt{X_1^2 + X_2^2 + \dots + X_n^2} \quad (7)$$

شناسایی داده‌های پرت استفاده می‌شود، به روش J.Gauss مشهور است و مطابق با رابطه (۵) محاسبه می‌شود.

$$g = \frac{x_{extr} - \bar{x}}{s} \quad (5)$$

که در آن،  $x_{extr}$  مقادیر حدی،  $S$  انحراف از معیار و  $x$  میانگین داده‌ها می‌باشد. چنانچه مقدار  $g$  محاسباتی از رابطه (۵) از مقدار  $G$  جدول بزرگ‌تر باشد، آن‌گاه داده پرت محسوب می‌شود (Szalma، ۱۹۸۴).

### روش‌های چندمتغیره شناسایی داده‌های پرت:

تشخیص چندمتغیره داده‌های پرت شامل بررسی چندمتغیره هر یک از مشاهدات بر اساس ترکیبی از متغیرها است. چون بیشتر تحلیل‌های چندمتغیره دارای بیش از دو متغیر است، تعیین داده‌های پرت از نظر ترکیبی از متغیرها نیز ضروری است. به‌منظور تشخیص چندمتغیره داده‌های پرت باید از روش‌ها یا معیارهایی استفاده شود که موقعیت چند بعدی (فضایی) هر یک از مشاهدات را نسبت به یک نقطه مشترک نشان دهند (Hair و همکاران، ۱۹۹۸). روش‌های مختلفی نظیر تجزیه به مؤلفه‌های اصلی (Causinus، ۲۰۰۳)، روش‌های مبتنی بر فاصله ماهالانوبیس (Filzmoser و همکاران، ۲۰۰۵) بدین منظور ارائه شده‌اند. از میان این روش‌ها فاصله ماهالانوبیس معروف‌تر است (Hair و همکاران، ۱۹۹۸). در روش‌های یاد شده، فاصله ماهالانوبیس به‌عنوان معیاری از موقعیت چند بعدی هر یک از مشاهدات نسبت به مرکز ثقل کل مشاهدات عمل می‌کند. به عبارت دیگر فاصله یاد شده، معیاری از فاصله هر یک از مشاهدات در فضای چند بعدی از مرکز میانگین تمام مشاهدات است. برتری عمده فاصله ماهالانوبیس نسبت به سایر فاصله‌ها، در نظر گرفته شدن ماتریس کواریانس در آن است (Filzmoser و همکاران، ۲۰۰۵). چون شکل و اندازه داده‌های چندمتغیره به وسیله ماتریس کواریانس تعیین می‌شود، برای یک نمونه چندمتغیره  $P$  (تعداد متغیرها) بعدی، فاصله ماهالانوبیس برای مشاهده نام از رابطه (۶) به‌دست می‌آید.

$$MD_i = [(X_i - \bar{X})^T C^{-1} (X_i - \bar{X})]^{1/2} \quad (6)$$

اما مقدار G بیشینه در این نقاط برابر ۲/۷۷ هست که از مقدار جدول کوچکتر بوده، بنابراین، این نقاط پرت محسوب نمی‌شوند. برای عنصر Nd (نئودیمیوم) نقطه ۱۶ بحرانی محسوب می‌شود، مقدار G بیشینه ۳/۳۴ هست که از میزان جدول هم بزرگتر هست و این نقطه برای این عنصر پرت محسوب می‌شود.

برای عنصر Ni (نیکل) در نمودار شاخه و برگ مقدار بحرانی دیده نمی‌شود. میزان G بیشینه آن برابر ۲/۸۵ هست که از مقدار جدول کوچکتر می‌باشد و داده پرتی محسوب نمی‌شود. برای عنصر P (فسفر) نقطه ۲۰ مقداری بحرانی دارد و G بیشینه معادل سه است که از میزان جدول بزرگتر است و داده پرت محسوب می‌شود. برای عنصر Si (سیلیکون) نقطه ۱۶ دارای مقدار بیشینه می‌باشد و بحرانی تلقی می‌شود که میزان G آن معادل ۲/۴۴ هست و چون کوچکتر از مقدار G جدول است، داده پرتی محسوب نمی‌شود. برای عنصر Sr (استرانسیم) مقادیر بحرانی وجود ندارد و میزان G بیشینه برابر ۱/۱۴ است که نشان می‌دهد برای این عنصر در هیچ نقطه‌ای مقدار پرتی وجود ندارد. برای عنصر Ti (تیتانیوم) در نمودار شاخه و برگ مقدار بحرانی در نقطه ۲۰ مشاهده می‌شود که G بیشینه معادل ۳/۲۱ هست و داده برای این نقطه پرت محسوب می‌شود.

برای عنصر V (وانادیوم) مقادیر بحرانی وجود ندارد و مقادیر بیشینه G برابر ۱/۵ و کمینه آن معادل ۱/۱۴ است که کوچکتر از مقدار G جدول است و داده پرتی برای این عنصر دیده نمی‌شود. برای عنصر Zn (زینک) مقادیر بحرانی وجود ندارد. مقدار G بیشینه معادل ۲/۲۵ و مقدار G کمینه ۱/۸ است که حکایت از عدم وجود داده پرت دارد. برای عنصر N (نیتروژن) مقادیر بحرانی دیده نمی‌شود. مقدار G بیشینه معادل دو و مقدار G کمینه ۱/۵ است که حکایت از عدم وجود داده پرت دارد. برای عنصر C (کربن) مقادیر بیشینه وجود ندارد و مقدار G بیشینه معادل ۱/۶۶ و مقدار G کمینه ۱/۴۹ است که عدم وجود داده پرت را نشان می‌دهد. بنابراین مطابق با این روش برای عناصر Nd (نقطه ۱۶)، P (نقطه ۲۰) و Ti (نقطه ۲۰) داده‌های پرت مشاهده شد.

که در آن، DSC فاصله نمرات نمونه  $n$ ، تعداد مؤلفه‌های استخراجی و  $X_1$ ،  $X_2$  و  $X_n$  به ترتیب برابر با نمرات مؤلفه اول، دوم و  $n$ ام برای نمونه  $n$ ام می‌باشند. هرچه فاصله نمرات نمونه بیشتر باشد، امکان پرت بودن نمونه‌ها بیشتر است. همچنین، در تشخیص چندمتغیره داده‌های پرت از مربع فاصله ماهالانوبیس (مقایسه  $MD^2/df$  و توزیع  $t$ ) نمودار جعبه‌ای مربع فاصله ماهالانوبیس و نمودار چندک-چندک<sup>۱</sup> مربع فاصله ماهالانوبیس در مقابل مربع توزیع کای استفاده شد. بر اساس نمودار چندک-چندک نمونه‌هایی که از روند کلی داده‌ها تبعیت نکنند، پرت محسوب می‌شوند (درجه آزادی تعداد ردیاب‌ها منهای یک می‌باشد).

### نتایج و بحث

**نتایج آزمون گراب (G):** برای عنصر Ba (باریم) طبق نمودار شاخه و برگ (برای هر ردیاب) نقطه ۲۴ یک مقدار حدی (EXTEREME) است. G جدول معادل ۲/۹ می‌باشد و بیشینه مقدار G برابر ۲/۶۹ هست. چون مقدار G کمتر از G جدول هست لذا فرض اول قبول است و داده پرت نیست. کمترین میزان G مربوط به نقطه ۲۷ می‌باشد که مقدار آن معادل ۱/۴۵ است. برای این نقطه هم فرض اول قبول و نقطه پرت محسوب نمی‌شود. برای عنصر Cd (کادمیوم) مقدار حدی وجود ندارد. برای نقطه ۳۳ مقدار G کمترین مقدار هست (۱/۵) که کوچکتر از مقدار جدول می‌باشد و فرض اول قبول می‌شود و داده پرتی مشاهده نمی‌شود.

برای عنصر Co (کبالت) نمودار شاخه و برگ مقدار حدی نشان نمی‌دهد و داده پرتی وجود ندارد. برای عنصر Cu (مس) نمودار شاخه و برگ مقدار حدی نشان نمی‌دهد. مقدار G دست‌کم معادل ۱/۳۹ هست که از مقدار G جدول کمتر بوده، آزمون نشان می‌دهد داده پرتی نیست. برای عنصر Li (لیتیم) مقدار G بیشینه و کمینه به ترتیب معادل ۲/۱۹ و ۱/۳۹ هست که مقدار G جدول کوچکتر بوده و داده پرتی وجود ندارد. برای عنصر Mn (منگنز) در نمودار شاخه و برگ برای نقاط ۱۱، ۱۵ و ۲۴ مقادیر بحرانی دیده می‌شود.

<sup>1</sup> Quantile-quantile plot

محسوب می‌شود. برای عنصر Ti نقطه ۲۰ با مقدار ۰/۵۷ که بزرگ‌تر از مقدار جدول یعنی ۰/۲۳۱ است، داده پرتی محسوب می‌شود. همچنین، برای عناصر V، Zn، N و C داده پرتی مشاهده نشد. بنابراین مطابق با این روش عناصر Ba (نقاط ۲۴ و ۱۵)، Co (نقطه ۲۴)، Mn (نقاط ۱۵ و ۲۴)، Nd (نقطه ۱۶)، Si (نقطه ۱۱) و Ti (نقطه ۲۰) دارای داده‌های پرت می‌باشند.

نتایج آزمون گوس: مقادیر محاسبه شده برای مقادیر بیشینه در عناصر Ba (۲/۶۹)، Cd (۲/۰۳)، Co (۲/۱)، Cr (۲/۴۴)، Cu (۲/۱۲)، Li (۲/۲۹)، Mn (۲/۷۷)، Nd (۳/۳۴)، Ni (۲/۸۵)، P (۳)، Si (۲/۴۴)، Ti (۳/۲۱)، V (۲/۰۳)، Zn (۲/۲۴)، N (۱/۶) و C (۱/۶۶) محاسبه شد که کلیه مقادیر محاسبه شده کمتر از مقدار جدول برای آزمون گوس یعنی مقدار ۳/۵۳ می‌باشد. بنابراین مطابق با این روش هیچ‌کدام از ردیاب‌ها دارای داده پرت نیستند.

نتایج آزمون میانگین به اضافه و منهای سه برابر انحراف از معیار: مطابق با این روش برای عناصر Nd (نقطه ۳۵)، P (نقطه ۱) و Ti (نقطه ۱) داده پرت وجود دارد و برای بقیه عناصر در هیچ نقطه‌ای داده پرت مشاهده نشد.

نتایج آزمون میانه  $\pm 3$  برابر MAD: این آزمون نشان داد، برای همه عناصر اندازه‌گیری شده دارای مقادیر پرت هستند و تعداد بیشتری از عناصر را در نقاط مختلف پرت تلقی کرد. نتایج حاصل از این روش در جدول زیر (جدول ۱) آمده است. در این روش کمترین داده‌های پرت مربوط به عناصر نئودیمیوم و نیتروژن هست.

همان‌گونه که در جدول ۱ مشاهده می‌شود، برای عنصر Ti نقطه ۲۰ به‌وسیله چهار روش پرت معرفی شده است. برای عنصر Si نقطه ۱۱ به‌وسیله سه روش نقطه پرت معرفی شده است. برای عنصر P نقطه ۲۰ به‌وسیله سه روش نقطه پرت معرفی شده است. برای عنصر Nd نقطه ۱۶ به‌وسیله چهار روش نقطه پرت معرفی شده است و برای عنصر Ba نقاط ۱۵ و ۲۴ به‌وسیله سه روش به‌عنوان نقطه پرت و داده پرت معرفی شده است.

مقایسه مقادیر مربع فاصله مایلانویس بروی درجه آزادی ( $MD^2/df$ ) با مقدار t مشاهداتی از جدول

نتایج نمودار جعبه‌ای: مطابق با نتایج این نمودار برای عنصر Ba نقاط ۱۵ و ۲۴ پرت محسوب می‌شوند. برای Cd نمودار جعبه‌ای داده پرتی را نشان نمی‌دهد. برای عنصر Co نمودار جعبه‌ای داده پرتی را نشان نمی‌دهد. برای عناصر Cu، Li، نقاط پرتی نشان داده نشد. برای عنصر Mn نقاط ۱۱، ۱۵ و ۲۴ به‌عنوان داده‌های پرت شناخته شدند. برای عنصر Nd نمودار جعبه‌ای نقطه ۱۶ را داده پرت معرفی کرد. برای عنصر Ni مطابق با نمودار جعبه‌ای داده پرتی دیده نمی‌شود. برای عنصر P نقطه ۲۰ داده پرت محسوب شده، برای عنصر Si در نمودار جعبه‌ای نقاط ۱۱ و ۱۶ داده پرت محسوب می‌شوند. برای عنصر Sr داده پرتی وجود ندارد. برای عنصر Ti مطابق با نمودار جعبه‌ای نقطه ۲۰ داده پرت تلقی شده است. برای عناصر V، Zn، N و C مطابق با نمودار جعبه‌ای نقطه پرتی دیده نمی‌شود. بنابراین مطابق با این روش برای عناصر Ba (نقاط ۲۴ و ۱۵)، Mn (نقاط ۱۱، ۱۵ و ۲۴)، Nd (نقطه ۱۶)، P (نقطه ۲۰)، Si (نقاط ۱۱ و ۱۶) و Ti (نقطه ۲۰) داده‌های پرت وجود دارد.

نتایج آزمون دیکسون<sup>۱</sup>: برای عنصر Ba در نقطه ۲۴ نقطه بحرانی شناخته شده است، میزان محاسبه شده معادل ۰/۴۴ است و مقدار قرائت شده از جدول (با توجه به تعداد نقاط نمونه‌برداری شده برای هر ردیاب و سطح اطمینان یک درصد) برابر ۰/۲۳۱ است. بنابراین این نقطه داده پرت محسوب می‌شود. همچنین برای نقطه ۱۵، مقدار محاسبه شده ۰/۳۹ است که از مقدار جدول بزرگ‌تر است و داده پرت محسوب می‌شود. برای عنصر Cd مقادیر بحرانی وجود ندارد و مقدار محاسباتی معادل ۰/۱۸۳ است که کوچک‌تر از مقدار جدول است و داده پرت دیده نمی‌شود. برای عنصر Co نقطه ۲۴ با مقدار ۰/۳۱ داده پرت محسوب می‌شود. برای عناصر Cr، Cu و Li داده پرت وجود ندارد. برای عنصر Mn نقاط ۲۴ و ۱۵ با مقادیر ۰/۳۷ و ۰/۹۹ داده‌های پرت محسوب می‌شوند. برای عنصر Nd نقطه ۱۶ با مقدار ۰/۹۴ پرت محسوب می‌شود. برای عناصر Ni و P داده پرت دیده نمی‌شود. برای عنصر Si نقطه ۱۱ با مقدار ۰/۵۷ داده پرت

<sup>1</sup> Dixon



رسوب، تنها نقطه ۲۲ را به عنوان داده پرت نشان داد. همچنین روش فاصله مالهالانویس نیز این نقطه را به عنوان داده پرت معرفی نمود. بنابراین همان طور که مشاهده می شود، تنها دو روش یاد شده بر پرت بودن نقطه ۲۲ اجماع دارند و بقیه روش ها چنین استدلالی ندارند.

همان گونه که مشاهده شد، نقطه ۲۲ تنها به وسیله روش میان به اضافه یا منهای میان انحراف های تمام داده ها از میان پرت شناخته شد و بقیه روش های تک-متغیره آن را پرت تشخیص ندادند. با این حساب تنها سه روش از مجموع ۱۱ روش به کار رفته بر پرت بودن این نقطه تأکید دارند و برای بقیه نقاط هم تنها چهار روش از مجموع ۱۱ روش به کار رفته بر پرت بودن داده ها دلالت دارند که این حکایت از عدم توافق ۵۰ درصد بر پرت بودن داده هاست.

همان طور که مشاهده می شود، در کل هیچ داده پرتی وجود ندارد و همه اندازه گیری ها و مقادیر صحیح و قابل استناد برای ورود به مدل های ترکیبی و مراحل بعد هستند. با توجه به این که روش های آزمون گراب، گوس و میانگین به اضافه و منهای سه برابر انحراف معیار تحت تأثیر داده های پرت و غیرعادی (داده های بسیار کوچک و بسیار بزرگ) می باشد (چون میانگین و انحراف معیار با استفاده از تمام داده ها برآورد می شوند) (Reimann و همکاران، ۲۰۰۵) کارایی زیادی در تشخیص داده های پرت و غیرعادی ندارند (Alijanpour و Hakimkhani، ۲۰۱۰). جدول ۱ نیز این نتیجه را تأیید می کند، به طوری که آزمون گوس و میانگین به اضافه و منهای سه برابر انحراف از معیار، هیچ نمونه ای را در مقایسه با آزمون گراب پرت معرفی نکرده و حساسیت آزمون گراب در مقایسه با این دو روش به مراتب بیشتر است. روش های میان به اضافه یا منفی میان انحراف های تمام داده ها از میان و نمودار جعبه ای تحت تأثیر داده های پرت نبوده، می توانند تعداد بیشتری نمونه پرت را شناسایی کنند (Chiang و همکاران، ۲۰۰۳). آزمون دیکسون تعداد داده های پرت بیشتری را در مقایسه با آزمون های گوس، گراب و میانگین به اضافه و منهای سه برابر انحراف از معیار، نشان می دهد.

t استیوندت برای درجه آزادی ۱۶ در سطح اطمینان ۰/۰۱ (مقدار قرائت شده برابر ۲/۵۸ بود)، نشان داد همه داده ها کمتر از مقدار جدول هستند و هیچ داده ای مطابق با این روش پرت محسوب نمی شود (درجه آزادی تعداد ردیاب ها منهای یک و سطح اعتماد داده ها یک درصد می باشد). برای بررسی بیشتر نمودار جعبه ای مربع فاصله مالهالانویس برای واحد کاربری اراضی نقطه ۲۲ پرت تشخیص داده شده است (شکل ۳).

نمودار چندک چندک مربع فاصله مالهالانویس که در مقابل توزیع مربع کای مرتب شده در شکل ۴ آورده شده است. طبق شکل نشان داده شده هیچ نمونه یا داده ای پرت تشخیص داده نشده است. پراکنش تمام نقاط (نمونه ها) از یک روند خاص پیروی می کند و انحراف قابل توجهی مشاهده نمی شود.

نتایج تجزیه تحلیل مؤلفه های اصلی نشان داد که نمونه های شماره ۴۳، ۱ و ۳۵ پرت هستند و فاصله DSC آن ها از بقیه نمونه ها بیشتر است. در تحلیل مؤلفه های اصلی تعداد چهار مؤلفه اصلی بیش از ۹۹ درصد از واریانس متغیرها (ردیاب ها) را بیان کردند. هر چه فاصله نمرات نمونه بزرگ تر باشد، امکان پرت بودن داده بیشتر است. همچنین، نمودار جعبه ای فاصله نمرات نمونه ها نیز برای تشخیص داده پرت ترسیم شد. در این نمودار برای فواصل نمونه ها هیچ داده ای پرت تشخیص داده نشد (شکل ۶). بنابراین، مقادیر مشکوک به پرت بودن مطابق با پراکنش فواصل نمرات نمونه پرت نیستند (شکل ۵).

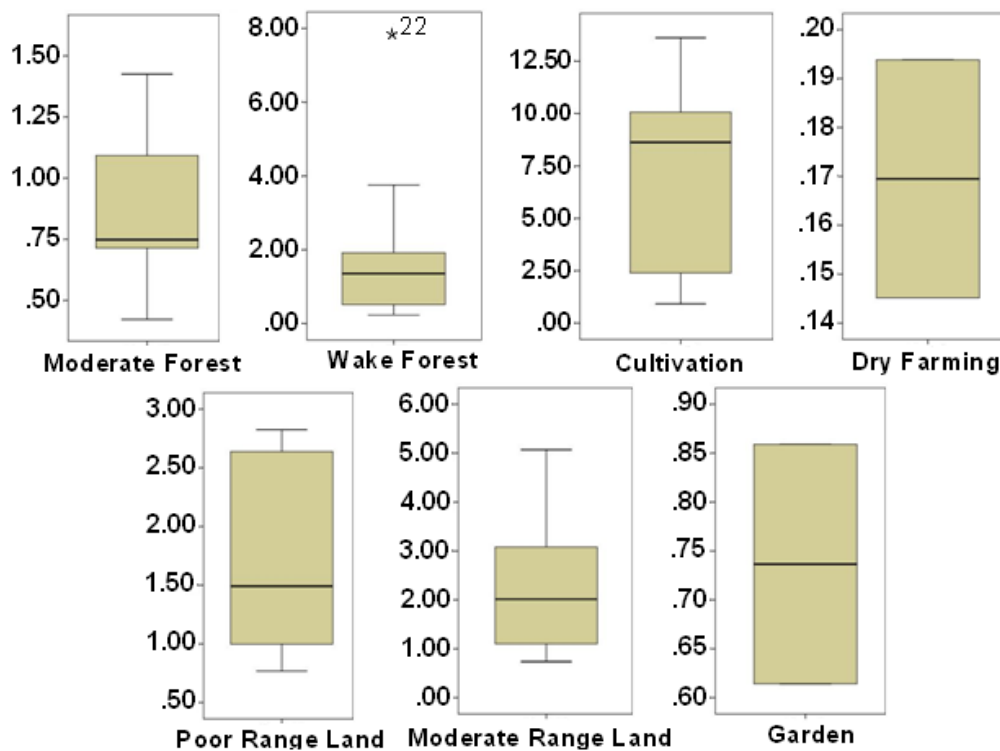
برای تشخیص داده پرت با استفاده از فاصله مالهالانویس از جدول خی دو (مربع کای) در سطح احتمال یک درصد و با درجه آزادی ۱۶ (مقدار قرائت شده ۵/۸۱) استفاده شد. نتایج مقایسه فاصله مالهالانویس و مقدار مربع کای در جدول ۳ ارائه شده است. همان طور که نتایج بالا نشان داد، روش چندمتغیره تحلیل مؤلفه های اصلی هیچ داده ای را پرت معرفی نکرد. همچنین روش مربع فاصله مالهالانویس نیز داده ای را پرت نشان نداد. نمودار چندک چندک مربع فاصله مالهالانویس در مقابل توزیع مربع کای نیز داده ای را پرت نشان نداد. نمودار جعبه ای مربع فاصله مالهالانویس برای هر منبع

جدول ۱- نتایج حاصل از بررسی داده های پرت با روش های تک متغیره

عنصر	میانگین	انحراف از معیار	میانگین	میانچه	میانچه انحراف های تمام داده ها از میانچه	آزمون گراب	نمودار جمع های مقادیر استاندارد	آزمون گوس	آزمون دیکسون	آزمون میانگین به اضافه یا منهای سه برابر انحراف از معیار	میانچه به اضافه یا منهای میانچه انحراف های تمام داده ها از میانچه
باریم	۱۶/۵۷	۸/۳	۱۶/۰۳	۶/۹۸	-	۲۴ و ۱۵	-	۲۴ و ۱۵	۲۴ و ۱۵	۲۸، ۲۶، ۲۴، ۱۵، ۲۰، ۲۲، ۱۱، ۱۰، ۷، ۰، ۲۷، ۳۴	
کادمیوم	۰/۴۴۴	۰/۲۷۹	۰/۳۸	۰/۳۴	-	-	-	-	-	۳۱، ۳۳، ۲۱، ۱۳، ۵، ۴	
کیالت	۳/۳۲۹	۱/۷۶۹	۳/۳۳	۲/۳۴	-	-	-	۲۴	۲۴	۳۲، ۳۴، ۲۵، ۲۴، ۱۵، ۱	
کروم	۱۵/۷۶۷	۷/۱۸۴	۱۵/۵۸	۷/۴۲	-	-	-	-	-	۲۷، ۲۵، ۲۳، ۱۹، ۳۴، ۱۷، ۱۵، ۶، ۵، ۰، ۳۳، ۳۲	
مس	۶/۳۵۳	۳/۰۹۶	۶/۱	۳/۴	-	-	-	-	-	۲۷، ۲۳، ۱۹، ۴۳، ۳۵، ۳۴، ۳۳، ۲۸، ۱۵، ۴۰	
لیتیم	۳/۴۶۴	۱/۸۱۷	۳/۴۴	۲/۲۶	-	-	-	-	-	۴۰، ۳۲، ۲۵، ۲۳، ۱۹، ۴۳، ۳۵، ۲۶، ۱۵	
منگنز	۱۵۶/۴۷	۹۳/۳۷	۱۲۰/۳	۷۶/۶۱	-	۲۴ و ۱۱، ۱۵	-	۲۴ و ۱۵	۲۴ و ۱۵	۳۷، ۲۸، ۲۶، ۲۴، ۲۲، ۲۰، ۱۵، ۱۱، ۶، ۰، ۳۴	
نئودیم	۳/۰۸۴	۱/۵۶	۳/۱۶	۱/۹۷	-	۱۶	۱۶	۱۶	۱۶	۲۰، ۱۰، ۱۶	
نیکل	۲۶/۱۲۷	۱۳/۹۳۵	۲۵/۶۸	۱۶/۱۹	-	-	-	-	-	۳۳، ۳۲، ۲۷، ۲۵، ۴۳، ۳۷، ۳۵، ۳۴، ۰	
فسفر	۱۲۹/۲۱۴	۸۳/۳۳۳	۱۰۰/۳۱۶	۷۹/۵۲	-	۲۰	۲۰	۱۱	۱	۲۵، ۳۴، ۲۸، ۲۶، ۲۴، ۲۰، ۱۸، ۱۵، ۱۱، ۲۵، ۴۳، ۳۶	
سیلیکون	۳/۴۹۲	۱/۹۷۳	۳/۱۳	۱/۷	-	۱۶ و ۱۱	-	-	-	۳۵، ۳۳، ۳۲، ۲۵، ۲۱، ۲۰، ۱۸، ۱۶، ۱۱، ۶، ۴۳، ۳۸	
استرانسیوم	۳۷/۷۲۶	۲۴/۲۱۷	۲۵/۶۱	۱۸/۳۴	-	-	-	-	-	۳۴، ۳۳، ۳۲، ۲۵، ۲۳، ۲۱، ۱۹، ۱۷، ۱۴، ۴۲، ۴۰، ۳۹، ۳۸، ۳۷	
تیتانیوم	۸/۷۰۵	۵/۰۹۲	۷/۵۸	۵/۱۴	-	۲۰	۲۰	۲۰	۱	۳۰، ۲۹، ۲۷، ۲۶، ۲۵، ۲۰، ۱۸، ۱۷، ۶، ۲، ۴۱، ۳۶، ۳۴	
وانادیوم	۱۰/۷۳۷	۴/۹۳	۹/۹۹	۵/۵۲	-	-	-	-	-	۲۷، ۲۳، ۴۳، ۳۵، ۳۴، ۳۴، ۱۵، ۱۱، ۶، ۰، ۳۳، ۳۲	
زینک	۱۴/۳۵۶	۷/۱۷۱	۱۲/۰۵	۶/۶۳	-	-	-	-	-	۳۴، ۲۸، ۲۶، ۲۴، ۲۲، ۲۰، ۱۸، ۱۱، ۱۵، ۲۵، ۴۳، ۳۵	
نیتروژن	۰/۰۸	۰/۰۲۴	۰/۰۸۱۹	۰/۰۲۹	-	-	-	-	-	۴۱، ۳۷، ۲۹، ۲۷، ۲۶، ۲۴، ۲۲، ۱۳، ۱۱، ۰، ۴	
کربن	۲۷/۰۷۶	۸/۰۸	۲۷	۱۱/۱۱	-	-	-	-	-	۱۷، ۱۳، ۱۱، ۴	

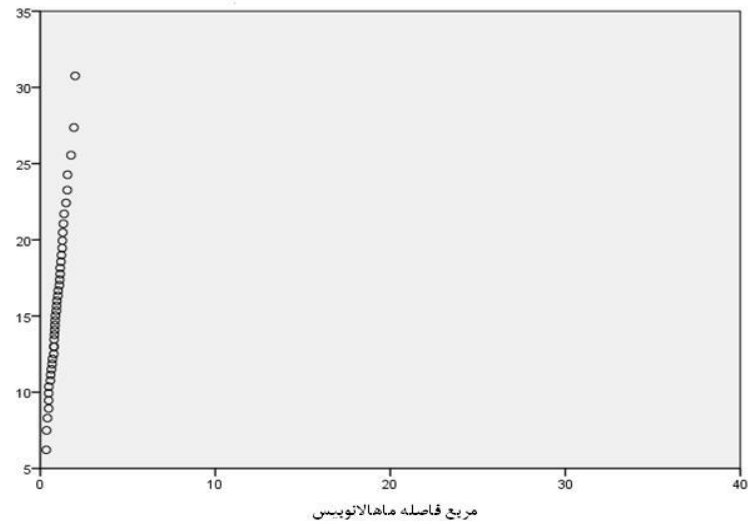
جدول ۲- مربع فاصله ماھالانوبیس و مربع فاصله ماھالانوبیس به روی درجه آزادی برای تشخیص داده‌های پرت

MD <sup>2</sup> /df	MD <sup>2</sup>	شماره نمونه	MD <sup>2</sup> /df	MD <sup>2</sup>	شماره نمونه
۰/۰۵۹۱	۱/۰۰۶	۲۳	۰/۰۹۳۲	۱/۵۸۴	۱
۰/۰۷۹۸	۱/۳۵۶	۲۴	۰/۰۱۳۶	۰/۲۳۲	۲
۰/۱۱۸۴	۲/۰۱	۲۵	۰/۰۳۱۸	۰/۵۴	۳
۰/۰۷۲۲	۱/۲۲۸	۲۶	۰/۱۱۲۷	۱/۹۱۶	۴
۰/۰۷	۱/۱۹۱	۲۷	۰/۰۳۶۷	۰/۶۲۴	۵
۰/۱۰۳۱	۱/۷۵۳	۲۸	۰/۰۵۰۵	۰/۸۵۸	۶
۰/۰۳۰۲	۰/۵۱۴	۲۹	۰/۰۲۴۷	۰/۴۲۱	۷
۰/۰۴۳۳	۰/۷۳۶	۳۰	۰/۰۴۵۱	۰/۷۶۷	۸
۰/۰۸۳۸	۱/۴۲۴	۳۱	۰/۰۵۸۶	۰/۹۹۷	۹
۰/۰۵۴	۰/۹۱۸	۳۲	۰/۰۳۶۱	۰/۶۱۴	۱۰
۰/۵۷۰۱	۹/۶۹۱	۳۳	۰/۰۹۹۹	۱/۶۹۹	۱۱
۰/۱۴۰۹	۲۱/۳۹۵	۳۴	۰/۰۵۸۳	۰/۹۹۲	۱۲
۰/۵۲۹۹	۹/۰۰۸	۳۵	۰/۰۴۳۹	۰/۷۴۷	۱۳
۰/۰۴۰۶	۰/۶۹۱	۳۶	۰/۱۵۵۱	۲/۶۳۷	۱۴
۰/۲۲	۳/۷۵۵	۳۷	۰/۱۵۵۱	۲/۶۳۷	۱۵
۰/۸۰۰۸	۱۳/۶۱	۳۸	۰/۰۱۷۳	۰/۲۹۴	۱۶
۰/۵۹۱۹	۱۰/۰۶	۳۹	۰/۰۱۱۴	۰/۱۹۳	۱۷
۰/۱۹۴۵	۳/۳۰۷	۴۰	۰/۰۰۸۵	۰/۱۴۵	۱۸
۰/۲۹۸۳	۵/۰۷۵	۴۱	۰/۰۴۵۷	۰/۷۷۸	۱۹
۰/۱۶۶۱	۲/۸۲۳	۴۲	۰/۱۶۷۸	۲/۸۵۳	۲۰
۰/۴۴۴۸	۷/۵۶۱	۴۳	۰/۰۲۵۱	۰/۴۲۶	۲۱
			۰/۴۶۴۷	۷/۹	۲۲

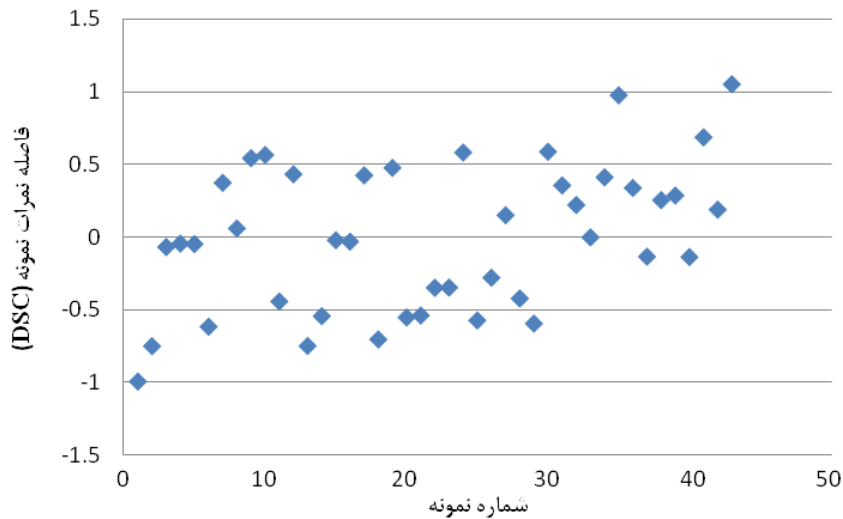


شکل ۳- نمودار جعبه‌ای مربع فاصله ماھالانوبیس برای منبع‌های زراعت آبی، جنگل ضعیف، جنگل متوسط، باغات، مرتع ضعیف، مرتع متوسط و دیهمزار

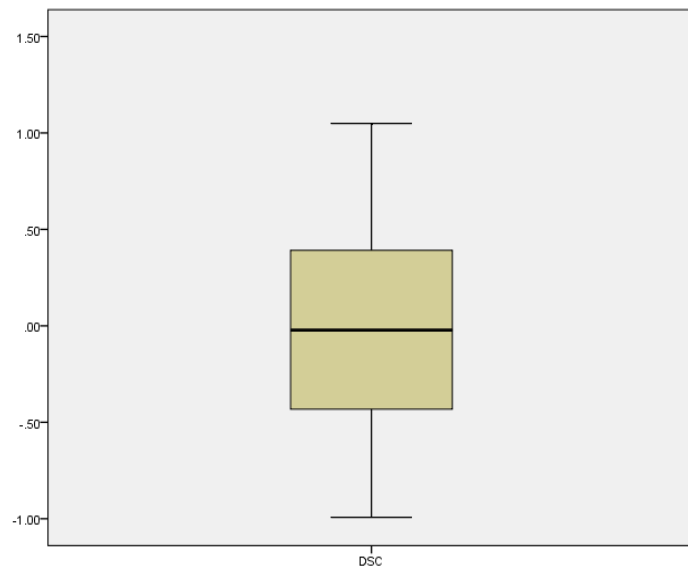
توزیع کای اسکور



شکل ۴- پراکنش مقادیر مربع فاصله ماهالانویس در مقابل مقادیر مورد انتظار توزیع مربع کای



شکل ۵- پراکنش فواصل نمرات نمونه (روش تجزیه مؤلفه‌های اصلی)



شکل ۶- نمودار جعبه‌ای فواصل نمرات نمونه‌ها

جدول ۳- داده‌های پرت مطابق با روش فاصله ماهالانوبیس

نام عنصر	نقاط پرت	میزان فاصله ماهالانوبیس داده پرت
Ba	۳۵	۷/۲۴۱
Cd	-	-
Co	-	-
Cr	۴۰	۵/۹۵
Cu	-	-
Li	۲۲	۵/۱۹۸
Mn	۴۳	۷/۶۷۷
Nd	۳۳	۱۱/۱۶۶
Ni	-	-
P	۳۹ و ۳۸	۷/۰۴ و ۸/۹۹۸
Si	۴۳ و ۳۳	۵/۸۵۷ و ۵/۹۳۱
Sr	-	-
Ti	۳۸	۱۰/۳
V	-	-
Zn	-	-
N	-	-
C	-	-

Hair و همکاران (۱۹۹۸) در مورد مربع فاصله ماهالانوبیس، Garrett (۱۹۸۹) و Filzmoser و همکاران (۲۰۰۵) در مورد نمودار چندک چندک به نتایج مشابهی دست یافتند. قابل ذکر است که آزمون‌های تک‌متغیره همچون گراب و گوس در تعداد نقاط برداشت شده برای هر ردیاب محدودیت دارند، در حالی که این محدودیت تعداد داده برای آزمون‌های چندمتغیره وجود نداشت. روشی همچون فاصله ماهالانوبیس اگرچه روشی مناسب برای تشخیص داده‌های پرت است، اما از آنجا که در به‌دست آوردن فاصله داده‌های پرت نیز مؤثرند این داده‌ها اثر منفی بر الگوریتم می‌گذارند. همچنین، قابل ذکر است که در کل با توجه به این که هیچ یک از نمونه‌ها از نظر اجماع روش‌های یک‌متغیره و از نظر تمام روش‌های بررسی چندمتغیره، پرت نیستند، بنابراین شواهد کافی مبنی بر پرت بودن و عضو جامعه نبودن هیچ نمونه‌ای وجود نداشته و نمی‌توان نسبت به حذف نمونه‌ای اقدام کرد.

داده‌های مورد استفاده در این پژوهش از نظر وجود خطا در ثبت، آماده‌سازی نمونه‌ها، روش‌های اندازه‌گیری عناصر و غیره کاملاً کنترل شده و تقریباً فاقد اشتباهات ناشی از عوامل یاد شده می‌باشند و از

قابل ذکر است که این روش بر اساس مقادیر حدی و مقدار فاصله آن‌ها از عمده داده‌ها کار می‌کند. همان‌طور که مشاهده شد دو روش باقی‌مانده از روش‌های تک‌متغیره یعنی نمودارهای جعبه‌ای و روش میانه به‌اضافه یا منهای میانگین انحراف‌های تمام داده‌ها از میانه تعداد قابل توجهی از نمونه‌ها را پرت نشان دادند. خصوصاً روش میانه به‌اضافه یا منهای میانگین انحراف‌های تمام داده‌ها از میانه. بنابراین این دو روش کارایی بیشتری نسبت به چهار روش قبل دارند.

برای استفاده از روش‌های تک‌متغیره به‌منظور شناسایی داده‌های پرت روش‌های میانه به‌اضافه یا منهای میانگین انحراف‌های تمام داده‌ها از میانه، نمودار جعبه‌ای و آزمون دیکسون به‌ترتیب پیشنهاد می‌شود. Reimann و همکاران (۲۰۰۵) نیز نتایج مشابهی در این زمینه دست یافتند. بررسی چندمتغیره داده‌های پرت به‌روش چندک چندک مربع فاصله ماهالانوبیس در برابر مربع کای، مربع فاصله ماهالانوبیس در مقایسه با توزیع  $t$  استیودنت نشان داد، هیچ داده‌ای پرت نیست. همچنین روش تحلیل مؤلفه‌های اصلی تنها سه نقطه را به‌عنوان داده پرت نشان داد و نمودار جعبه‌ای فاصله ماهالانوبیس در منابع تولید رسوب تنها یک نقطه را به‌عنوان داده پرت نشان داد.

کودپاشی مزارع) است. با توجه به نتایج حاصل از تحقیق حاضر پیشنهاد می‌شود در زمینه منشأیابی رسوب و استفاده از مدل‌های ترکیبی در مشخص کردن سهم منابع رسوب و در جهت عدم ورود داده‌های پرت به این روند، از اجماع چندین روش شناسایی داده‌های پرت استفاده شد و تنها یک شاخص یا یک آزمون ملاک عمل قرار نگیرد. همان‌گونه که قبلاً نیز اشاره شد، تمامی روش‌های شناسایی داده‌های پرت مقبولیت جهانی ندارند و نیاز است که حصول پرت بودن مشاهده‌ای از چندین روش مختلف آشکار شود. با توجه به جمیع مطالب بالا می‌توان گفت، یک داده پرت مشاهده‌ای است که به‌طور غیرعادی و یا اتفاقی از وضعیت کلی داده‌های تحت آزمایش و نسبت به قاعده‌ای که بر اساس آن تحلیل می‌شوند، تفاوت معنی‌داری داشته باشند و چندین روش بر انحراف آن از وضعیت عمومی داده‌ها ادعان داشته باشند.

سوی دیگر نمونه‌های خاک از محل‌های مناسب (معرف جامعه) و غیرآلوده برداشت شده‌اند. از این‌رو، امکان وجود داده‌های پرت بسیار کم است.

نمونه‌هایی که در بعضی از روش‌های مورد استفاده به‌عنوان پرت شناسایی شده‌اند، در واقع جزء داده‌های حد محسوب می‌شوند. فرق داده‌های حد با داده‌های پرت این است که داده‌های حد عضو توزیع اصلی (جامعه) داده‌های مورد مطالعه بوده و در فاصله دورتری از مرکز توزیع قرار دارند (خیلی بزرگ یا خیلی کوچک هستند)، ولی داده‌های پرت عضو توزیع اصلی نبوده و جز و یک یا چند توزیع متفاوت هستند (Reimann و همکاران، ۲۰۰۵؛ Filzmoser و همکاران، ۲۰۰۵). در واقع فرایند یا فرایندهای ایجاد کننده داده‌های پرت متفاوت از داده‌های اصلی است. هم‌آن‌طور که گفته شد داده‌های پرت در ژئوشیمی علاوه بر اشتباه یا خطا، اغلب ناشی از فرایندهای کانی‌سازی، آلتراسیون و فعالیت‌های انسانی (از جمله

#### منابع مورد استفاده

1. Aliehyai, M. and A.A. Behbahanizade. 1993. Methods of chemical analysis of soil. Institute of Soil and Water Research, Bulletin No 893, 26 pages (in Persian).
2. Alfassi, Z.B., Z. Boger and Y. Statistical. 2005. Treatment of analytical data: outliers (Chapter 6). CRC Press, 512 pages.
3. Caussin, H., M. Fekri, S. Hakam and A. Ruiz-Gazen. 2003. A monitoring display of multivariate outliers. Computational Statistics and Data Analysis, 44: 237-252.
4. Chiang, L.H., R.J. Pell and M.B. Seasholtz. 2003. Exploring process data with the use of robust outlier detection algorithms. Journal of Process Control, 13: 437-449.
5. Collins, A.L., D.E. Walling, L. Webb and P. King. 2010. Apportioning catchment scale sediment sources using a modified composite fingerprinting technique incorporating property weightings and prior information. Geoderma, 155: 249-261.
6. Collins, A.L. and D.E. Walling. 2007. Sources of fine sediment recovered from the channel bed of lowland groundwater-fed catchments in the UK. Geomorphology, 88: 120-138.
7. Collins, A.L., Y. Zhang, D. McChesney, D.E. Walling, S.M. Haley and P. Smith. 2012. Sediment source tracing in a lowland agricultural catchment in southern England using a modified procedure combining statistical analysis and numerical modelling. Science of the Total Environment, 414: 301-317.
8. EPA. 2006. Data quality assessment: statistical methods for practitioners EPA QA/G-9S, EPA/240/B-06/003, U.S. Environmental Protection Agency, Office of Environmental Information, Washington DC.
9. Feng, J., Z.G. Hu, J.T. Ju and L.P. Zhu. 2011. Variations in trace element (including rare earth element) concentrations with grain sizes in loess and their implications for tracing the provenance of eolian deposits. Quaternary International, 236: 116-126.
10. Filzmoser, P., R.G. Garrett and C. Reimann. 2005. Multivariate outlier detection in exploration geochemistry. Computers and Geosciences, 31: 579-587.
11. Garrett, R.G. 1989. The chi-square plot: a tool for multivariate outlier recognition. Journal of Geochemical Exploration, 32: 319-341.
12. Hakimkhani, Sh. and A. Alijanpour. 2010. Detection of outliers in the sediment fingerprinting method. Water and Soil Conservation, 17(1): 23-43 (in Persian).
13. Hair, J.F., R.E. Andersen, R.L. Tatham and W.C. Black. 1998. Multivariate data analysis. Rentic Hall, Upper Saddle River, New Jersey.

14. Szalma, J. 1984. Mérési eredmények kiértékelésének alapjai (Introduction to evaluation of measured data). Tankönyvkiadó, Budapest.
15. Reimann, C., P. Filzmoser and R.G. Garrett. 2005. Background and threshold: critical comparison of methods of determination. *Science of the Total Environment*, 346: 1-16.
16. Rousseeuw, P.J. and B.C. Van Zomeren. 1990. Unmasking multivariate outliers and leverage points. *Journal of the American Statistical Association*, 85: 633-651.
17. Rousseeuw, P.J. and K. Van Driessen. 1999. A fast algorithm for the minimum covariance determinant estimator. *Technometrics*, 41: 212-223.
18. Nosrati, K., G. Govers, H. Ahmadi, F. Sharifi, M.A. Amoozegar and R. Merckx. 2011. An exploratory study on the use of enzyme activities as sediment tracers: biochemical fingerprints? *International Journal of Sediment Research*, 26(2): 136-151.
19. Rorabacher, D.B. 1991. Statistical Treatment for rejection of deviant values: critical values of dixon Q parameter and related subrange ratios at the 95 percent confidence level. *Anal. Chem*, 83(2): 139-146.
20. Tabachnick, B.G. and L.S. Fidell. 1996. Using multivariate statistics. Harper Collins College Publishers, New York, 963 pages.
21. Tukey, J.W. 1977. Exploratory data analysis. Addison-Wesley Publication, 205-235.
22. Lalor, G.C. and C. Zhang. 2001. Multivariate outlier detection and remediation in geochemical databases. *The Science of the Total Environment*, 281: 99-109.
23. Walling, D.E., A.L. Collins and R.W. Stroud. 2008. Tracing suspended sediment and particulate phosphorus sources in catchments. *Journal of Hydrology*, 350: 274-289.
24. Walling, D.E. 2005. Tracing suspended sediment sources in catchments and river systems. *Science of the Total Environment*, 344: 159-184.
25. Walling, D.E., P.N. Owens and G.J.L. Leeks. 1999. Fingerprinting suspended sediment sources in the catchment of the River Ouse, Yorkshire, UK. *Hydrological Processes*, 13: 955-975.
26. Wiegand, P., R. Pell and E. Comas. 2009. Simultaneous variable selection and outlier detection using a robust genetic algorithm. *Chemo metrics and Intelligent Laboratory Systems*, 98(2): 108-114.
27. Zhang, C., D. Fay, D. McGrath, E. Grennan and O.T. Carton. 2008. Statistical analyses of geochemical variables in soils of Ireland. *Geoderma*, 146: 378-390.
28. Zhang, C.S., P.M. Wong and O. Selinus. 1999. A comparison of outlier detection methods: exemplified with an environmental geochemical dataset. P 183-187, In: *Proceeding of the 6th International Conference on Neural Information Processing*, Perth, Australia.

## Using univariate and multivariate methods to detect outliers in sediment fingerprinting method, case study: Tange Bostanak Watershed

Ahmad Nohegar<sup>1</sup>, Mohamad Kazemi<sup>\*2</sup>, Seyed Javad Ahmadi<sup>3</sup>, Hamid Gholami<sup>4</sup>, Rasool Mahdavi<sup>5</sup>  
<sup>1</sup> Professor, Faculty of Environment Sciences, Tehran University, Iran, <sup>2</sup> PhD Student, Faculty of Natural Resources, Hormozgan University, Iran, <sup>3</sup> Associate Professor, Fuel Cycle Research Institute of Atomic Energy Organization, Iran and <sup>4</sup> and <sup>5</sup> Assistant Professor, Faculty of Natural Resources, Hormozgan University, Iran

Received: 17 September 2016

Accepted: 27 February 2017

### Abstract

Efficiency of sediment fingerprinting by using tracers as a successful method to determine the sources of sediment has been proved. Selection of the suite subset of tracers, capable of discriminating sediment sources, is the first and the most important step in the sediment fingerprinting method. The presence of outliers affects the selection of the suite subset and possibly prevents picking the important tracers and reducing the accuracy of classification. Therefore, the outliers must be detected in order to be corrected or omitted, if enough evidences were present. The present study aims to detect outliers in the subset of tracers, to identify the best combination. For detecting outliers, We used univariate methods such as Grubbs test, Gauss test, Dioxin test, box plot, the Median  $\pm$  3MAD, the mean  $\pm$  3standard deviation and also multivariate methods such as squared Mahalanobis distance, separate box plots of squared Mahalanobis distance for each of sediment sources, principal component analysis and plot of the squared Mahalanobis distances against the quantiles of the chi-square distribution. we consider an observation as the outlier that at least half of these methods have detected it as an outlier. The results showed that Median  $\pm$  3MAD method introduced a larger number of data as outliers. Methods of multivariate outlier detection has low agreement with each other. Univariate methods to identify outliers show higher agreement with each other. To use univariate analysis techniques to detect outliers namely Median  $\pm$  3MAD, box plot, and Dioxin one can recommended to test their sensitivity. The results also showed that the maximum consensus for univariate analysis techniques is four samples (observations) and for multivariate methods is two samples (observations). In general, there is no observation that is identified as an outlier by half of the used methods.

**Key words:** Grubbs' Test, Mahalanobis Distance, PCA, Sediment Fingerprinting, Tracers

---

\* Corresponding author: mohamad.kazemi86@gmail.com