

کاربرد الگوریتم‌های داده‌کاوی در تفکیک منابع رسوبی حوضه آبخیز نوده گناباد

مهدی بشیری^{۱*}، مهسا آریاپور^۲، علی گل‌کاریان^۳

تاریخ دریافت: ۱۳۹۶/۱۲/۱۴

تاریخ پذیرش: ۱۳۹۷/۴/۱۱

چکیده

لازمه اجرای برنامه‌های کنترل رسوب، شناسایی اهمیت نسبی منابع رسوب، میزان مشارکت آن‌ها و در نتیجه شناسایی مناطق بحرانی آبخیزهاست. در این پژوهش از الگوریتم‌های داده‌کاوی برای تفکیک منابع رسوبی حوضه نوده گناباد در استان خراسان رضوی با کمک متغیرهای ژئوشیمیایی، دانه‌بندی و سنگ‌شناسی استفاده شد. یازده الگوریتم برای طبقه‌بندی در نرم‌افزار MATLAB برنامه‌نویسی و نتایج براساس ضریب تبیین و میانگین مربع خطا با یکدیگر مقایسه شد. بررسی غلظت عناصر ژئوشیمیایی در هفت واحد زمین‌شناسی حوضه نشان داد که عناصر Ca، Fe، Mg و Al دارای بیشترین و عناصر B و Co دارای کمترین غلظت در نمونه‌های خاک است. ارزیابی کلی الگوریتم‌های طبقه‌بندی در مرحله آموزش نشان داد که الگوریتم‌های تحلیل ممیزی، جنگل تصادفی، k نزدیک‌ترین همسایه و ماشین‌های بردار پشتیبان با توابع خطی، چندجمله‌ای، چندگانه و شعاع مبنا با حداکثر مقدار ضریب تبیین ($R^2=1$) و حداقل مقدار میانگین مربع خطا ($MSE=0$)، دقیق‌ترین الگوریتم‌ها در تفکیک منابع رسوبی هستند و روش درخت رگرسیونی ضعیف‌ترین عملکرد را دارد. در مرحله آزمون نیز ماشین‌های بردار پشتیبان با تابع شعاع مبنا، دقیق‌ترین الگوریتم و درخت طبقه‌بندی با بالاترین خطا، ناکارآمدترین الگوریتم بود. همچنین ورود متغیرهای ژئوشیمیایی منجر به بالاترین دقت در تفکیک منابع رسوبی شد و متغیرهای دانه‌بندی کمترین دقت تفکیک را باعث شد.

واژه‌های کلیدی: الگوریتم‌های طبقه‌بندی، منشأیابی، حوضه نوده، غلظت عناصر.

۱. استادیار گروه مرتع و آبخیزداری، دانشکده کشاورزی و منابع طبیعی، دانشگاه تربت حیدریه، خراسان رضوی / m.bashiri@torbath.ac.ir

۲. کارشناس ارشد آبخیزداری، دانشکده کشاورزی و منابع طبیعی، دانشگاه تربت حیدریه، خراسان رضوی

۳. استادیار گروه مرتع و آبخیزداری، دانشکده منابع طبیعی و محیط زیست، دانشگاه فردوسی مشهد، خراسان رضوی

مقدمه

فرسایش خاک پس از رشد جمعیت، دومین چالش مهم زیست‌محیطی جهان است که به دلیل داشتن اثرات چندجانبه آشکار و پنهان زیست‌محیطی و اجتماعی، یکی از فرایندهای پیچیده و خطرناک محیطی است. فرسایش خاک و حمل رسوب پدیده‌ای است که در نواحی مختلف و در مقیاس‌های متفاوتی اتفاق می‌افتد (عرب‌خدری، ۲۰۱۴). شناخت عوامل مؤثر در افزایش تخریب خاک به منظور درک صحیح‌تر فرایند فرسایش و تخریب خاک در تمامی اقالیم، به منظور جلوگیری و کاهش روند تخریب امری ضروری است (تورنیل^۱ و همکاران، ۲۰۰۸). در این میان شناخت منشأ رسوب از جمله سهم زیرحوضه‌ها، گامی اساسی برای سیاست‌گذاری مدیریت حوزه‌های آبخیز است. تعیین این سهم کاربرد مؤثری در مدل‌سازی، شناخت الگوهای پراکنش، تدوین روش‌های مهار رسوب و تخمین میزان فرسایش خاک دارد. بدین منظور یافتن روش‌هایی برای برآورد علمی و دقیق‌تر میزان فرسایش و ماهیت حرکت، رسوبات به دو شکل بار معلق و بستر انتقال می‌یابند و از نظر منبع پیدایش نیز در دو گروه بار رسوبی حوضه‌ای و آبراهه‌ای تقسیم می‌شوند. آگاهی از مقدار و تغییرات مقادیر رسوب در تنظیم تغییرات بستر، کیفیت آب، رفتار هیدرولیکی و انجام پروژه‌های حفاظت آب و خاک و مدیریت و برنامه‌ریزی آن کمک شایانی می‌کند (میلهوس^۲، ۱۹۸۸). امروزه حجم زیاد داده‌های ذخیره‌شده و نیز گستردگی ابعاد آن سبب شده است که در بسیاری از موارد، روش‌های آماری به‌تنهایی قادر به کشف خصوصیات داده‌ها نباشند. داده‌کاوی^۳ به‌عنوان یک راه‌حل برای این مسائل بوده و عبارت است از استخراج اطلاعات و دانش و کشف الگوهای پنهان از یک پایگاه داده بسیار بزرگ که این الگوها و دانش‌ها معمولاً مستتر در داده‌ها هستند (چان و لویز^۴، ۲۰۰۲). داده‌کاوی در اواخر دهه ۱۹۸۰ پدیدار شد و سپس در دهه ۱۹۹۰ گام‌های

بلندی در این شاخه از علم برداشته شد و انتظار می‌رود همچنان به رشد و پیشرفت خود ادامه دهد. تاکنون مطالعات مختلفی در زمینه مدل‌سازی فرسایش خاک و تولید رواناب و رسوب با استفاده از الگوریتم‌های داده‌کاوی صورت گرفته است. با توجه به اینکه هدف از مطالعه حاضر، تفکیک و طبقه‌بندی منابع رسوبی حوضه و ارزیابی عملکرد الگوریتم‌های طبقه‌بندی در شناسایی منابع رسوبی با استفاده از ردیاب‌های طبیعی حوضه است، در ادامه به تحقیقات مرتبط با پژوهش حاضر اشاره می‌شود.

حرما^۵ و همکاران (۲۰۱۵) مدل‌سازی رواناب و رسوب را در حوضه‌ای واقع در نپال به وسیله شبکه‌های عصبی مصنوعی انجام دادند. نتایج ایشان نشان داد شبکه عصبی مصنوعی^۶ کارایی بسیار خوبی در این زمینه دارد. کومارگویال^۷ (۲۰۱۴) طی پژوهشی عملکرد مدل درخت M5 و روش رگرسیونی^۸ را در مقایسه با شبکه عصبی مصنوعی برای پیش‌بینی آورد رسوبی حوزه آبخیز ناگوا در هند مورد بررسی قرار داد. نتایج نشان‌دهنده برتری عملکرد مدل M5 و رگرسیونی نسبت به روش دیگر بود. حدادچی^۹ و همکاران (۲۰۱۳) در مناطق آبرفتی، به اجرای تکنیک انگشت‌نگاری رسوب با استفاده از ردیاب‌ها پرداختند. در این مطالعه، منابع رسوبی با استفاده از هفت مدل آمیخته، ارزیابی و با کمک الگوریتم ژنتیک بهینه‌سازی گردید. نتایج نشان داد که دامنه خطا در تمامی مدل‌ها در محدوده قابل قبول قرار دارد. دمیرسی و بالتاسی^{۱۰} (۲۰۱۳) طی پژوهشی در ساکری منتوفی پورت^{۱۱} آمریکا، از روش‌های منطق فازی، رگرسیون و منحنی سنجه رسوب برای برآورد مقدار غلظت رسوب معلق و دمای آب که به‌طور پیوسته در مدت زمان پنج سال تهیه شده بود، استفاده کردند. نتایج حاکی از برتری دقت روش منطق فازی نسبت به سایر روش‌ها بود. میسرا^{۱۲} و همکاران (۲۰۰۹)، با هدف شبیه‌سازی رواناب و رسوب با استفاده از دو روش ماشین بردار پشتیبان و

5. Harma

6. Artificial Neural Network

7. Kumar goyal

8. Regression Method

9. Haddadchi

10. Demirci and Baltaci

11. Sacremento Feoport

12. Misra

1. Turnbull

2. Milhouse

3. Data mining

4. Chan and Lewis

ارائه کرده است. گرچه با توجه به مرور منابع، عمدتاً روش‌های داده‌کاوی در زمینه مدل‌سازی رسوب به کار رفته‌اند، تکنیک‌های طبقه‌بندی رسوبات با استفاده از الگوریتم‌های موجود، چندان توسعه پیدا نکرده‌اند و در زمینه طبقه‌بندی منابع رسوبی به کاربرد روش‌های سنتی مانند رگرسیون و نهایتاً آنالیز تشخیص اکتفا شده است. لذا این پژوهش با هدف تفکیک منابع رسوبی تجمع پیدا کرده در پشت سد خاکی حوزه آبخیز نوده گناباد (خروجی حوضه) با استفاده از علم نوین داده‌کاوی صورت گرفت تا بتواند در زمینه منشأیابی رسوب و تعیین بهترین الگوریتم در این زمینه مورد کاربرد قرار گیرد؛ چراکه این الگوریتم‌ها در زمینه‌های مختلف عملکرد موفق داشته‌اند. بنابراین انتظار می‌رود که در تفکیک منابع رسوبی نیز موفق عمل نمایند.

مواد و روش‌ها

حوزه آبخیز نوده با مساحتی در حدود ۲۲۹۸/۳۹ هکتار در حوزه آبخیز کویر مرکزی و از نظر کشوری در محدوده شهرستان گناباد در استان خراسان رضوی واقع شده است. موقعیت جغرافیایی حوزه آبخیز نوده بین طول جغرافیایی ۲۴° ۲۰' ۵۸" تا ۲۴° ۲۴' ۳۸" و عرض جغرافیایی ۱۳° ۳۴' ۵۰" تا ۳۴° ۲۰' ۵۰" است. این حوضه دارای هفت واحد زمین‌شناسی است که شامل نهشته‌های آبرفتی قدیمی (Qtl)، ریولیت، آلکالی ریولیت، آلکالی تراکیت، ریوداسیت، داسیت، لاتیت، توف اسید (Er)، دولومیت (Tr.sh)، آهک اوریتولین‌دار (KI)، داسیت، آندزی بازالت، ریولیت (Ed)، آندزیت، داسیت، ریوداسیت، لاتیت، بازالت (Ec) و گدازه‌های آتشفشانی (Ea) است. همچنین میزان بارندگی سالانه آن ۲۵۴ میلی‌متر است و در خروجی این حوضه، یک سد خاکی احداث شده است. روش کار این پژوهش شامل دو بخش آنالیزهای آزمایشگاهی و آماری است. از آنجایی که داده‌های ژئوشیمیایی رسوبات برای مطالعه آن‌ها از دید رسوبی و زیست‌محیطی بسیار حائز اهمیت است و ترکیب شیمیایی این رسوبات حاکی از پایایی عوامل مختلف زمین‌شناسی مانند جایگاه زمین‌ساختی، ترکیب سنگ منشأ، شدت هوازدگی، بلوغ بافتی و کانی‌شناسی طی حمل و رسوب‌گذاری است

شبکه عصبی به مطالعه حوضه‌ای در هند پرداختند. در این مطالعه، سری‌های زمانی داده به مجموعه‌های آموزش، واسنجی و اعتبارسنجی تقسیم شد و عملکرد مدل‌ها با استفاده از ضریب همبستگی بررسی گردید. نتایج نشان داد که ماشین بردار پشتیبان بهبود قابل توجهی در آموزش، کالیبراسیون و اعتبارسنجی در مقایسه با مدل شبکه عصبی دارد و می‌تواند به‌عنوان جایگزین مناسبی در برآورد رواناب و پیش‌بینی رسوب مورد استفاده قرار گیرد.

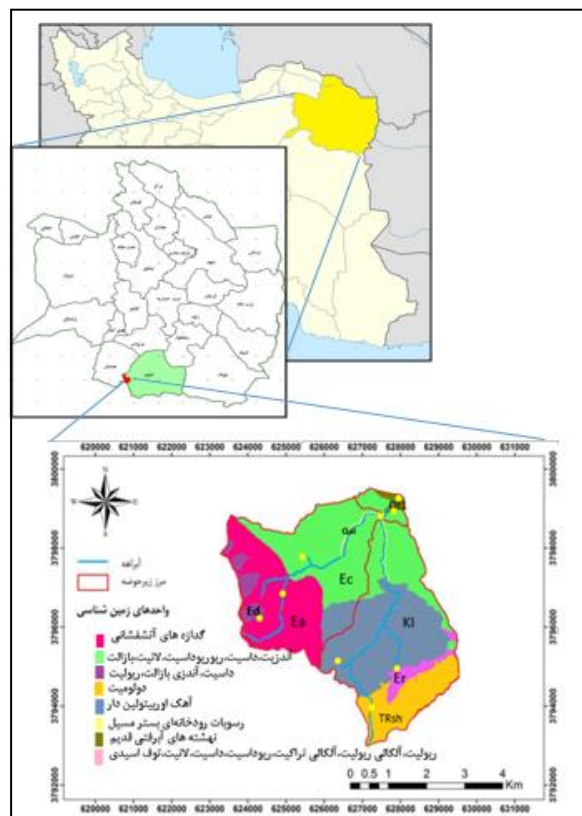
حیات‌زاده و همکاران (۲۰۱۵) طی پژوهشی به پیش‌بینی بار رسوبی حوضه باغ‌عباس با استفاده از دو روش شبکه عصبی و روش‌های رگرسیونی پرداختند. در این تحقیق، داده‌های ۱۳۶ واقعه دبی جریان و پارامترهای مورفولوژیکی مورد تحلیل قرار گرفت. نتایج نشان داد بهترین روش برای برآورد رسوب حوضه، روش شبکه عصبی به همراه داده‌های ژئومورفولوژیکی حوضه است. کاکائی لخدانی (۲۰۱۳) در پژوهشی، توانایی مدل‌های شبکه عصبی و ماشین بردار پشتیبان^۱ در برآورد رسوب معلق روزانه رودخانه دوبرج واقع در غرب ایران، بررسی و با روش‌های رگرسیونی مقایسه کردند. نتایج حاکی از برتری مدل‌های یادشده نسبت به مدل‌های رگرسیونی بود. جودی و ستاری (۱۳۹۵) در مقایسه برآورد رسوبات رودخانه‌ای نشان دادند روش‌های داده‌کاوی عملکرد بهتری را نسبت به روش‌های سنتی و رگرسیونی نشان دادند. در پژوهش آن‌ها روش‌های سنتی تفاوت چشمگیری نداشته و در برآورد بار رسوبی کل عملکرد مشابهی را داشتند.

بسالت‌پور و همکاران (۲۰۱۶) در بررسی اثرات ویژگی‌های خاک بر فرسایش‌پذیری حوضه کارون به این نتیجه رسیدند عملکرد مدل SVM نسبت به مدل رگرسیون خطی که حتی قادر به تشخیص روابط و اثرات فاکتورهای موثر بر فرسایش حوضه نبوده، بسیار مناسب معرفی شده است. با توجه به نتایج پژوهش‌های انجام‌شده با روش‌های مختلف در زمینه فرسایش و رسوبات حوضه می‌توان نتیجه گرفت الگوریتم‌های داده‌کاوی در اکثر موارد، دقت بالاتر و عملکرد بهتری را نسبت به سایر روش‌های آماری و سنتی

نمونه‌برداری از واحدهای زمین‌شناسی درون آبراهه‌هایی که فقط از همان واحد سرچشمه می‌گرفتند، انجام شد. در نهایت از هر مجموعه ۵ تایی نمونه واحدهای زمین‌شناسی، ۴ تکرار به صورت دوبه‌دو پس از حمل به آزمایشگاه با یکدیگر مخلوط شدند که از ۴۰ نمونه برداشت‌شده، ۲۲ نمونه (یک نمونه خروجی و ۲۱ نمونه منابع رسوب) باقی ماند. برای بررسی دانه‌بندی و بافت خاک از هر نمونه، دو زیرنمونه ۲۰۰ گرمی و ۵۰ گرمی تهیه شد. ابتدا زیرنمونه‌های ۲۰۰ گرمی از هفت الک با منافذ درشت‌تر از ۶۳ میکرون سری الک ادن - ونت‌ورث^۲ به روش الک‌تر عبور داده شد. همچنین زیرنمونه‌های ۵۰ گرمی که به روش الک خشک از رسوبات زیر ۲ میلی‌متر تهیه شده بود، با استفاده از هیدرومتر و روش ASTM^۳ (۲۰۰۸) برای اندازه‌گیری درصد رس، سیلت و شن استفاده شد و سپس پارامترهای دانه‌بندی ذرات رسوب، با استفاده از نرم‌افزار GRADISTAT نسخه ۸ تعیین گردید. برای مطالعات مورفوسکوپی و سنگ‌شناسی در نمونه‌های رسوب، ابتدا سنگ‌های شاخص واحدهای مختلف به روش هدف‌دار (جمع‌آوری خرده‌سنگ‌های شاخص واحدهای کاری) از داخل نمونه‌ها جمع‌آوری شد و بعد از دسته‌بندی (از لحاظ نوع سنگ و کانی‌های موجود در آنها) و شناسایی نوع سنگ، از هریک از نمونه رسوبات واحدها و خروجی حوزه یک زیرنمونه شامل ۱۰۰ عدد خرده‌سنگ به روش تصادفی روی الک با قطر منفذ ۶۳ میکرون تهیه و اقدام به تطابق آن‌ها با سنگ‌های شناسایی شده در زیر باینوکولر^۴ شد.

در مرحله بعد، برای بررسی غلظت عناصر موجود در نمونه‌ها ابتدا از هر نمونه، یک زیرنمونه یک گرمی از رسوبات زیر الک ۶۳ میکرون به روش الک خشک تهیه شد و بعد از هضم اسیدی، به حجم رسانیدن و عبور از کاغذ صافی (یاب^۵ و همکاران، ۲۰۰۲)، با دستگاه ICP OES^۶ غلظت عناصر اندازه‌گیری شد. برای این منظور ابتدا ۳۲ عنصر در نمونه رسوبات مخزن سد اندازه‌گیری و سپس عناصری که در

(سایگل^۱، ۲۰۰۲)، در بخش آنالیزهای آزمایشگاهی پژوهش حاضر به شناسایی و تعیین غلظت عناصر موجود در نمونه‌های خاک پرداخته شد. شکل (۱) موقعیت جغرافیایی حوزه آبخیز مورد مطالعه به همراه مکان‌های نمونه‌برداری در واحدهای زمین‌شناسی را نشان می‌دهد.



شکل (۱): موقعیت جغرافیایی حوزه و مکان‌های نمونه‌برداری

Figure (1): Geographic location of the watershed and sampling points

در این پژوهش، نمونه‌برداری از واحدهای زمین‌شناسی درون آبراهه‌هایی که تنها از همان واحد سرچشمه می‌گرفتند، انجام شد. بعد از مشخص کردن واحدهای زمین‌شناسی، برای هر تکرار بین ۲۵۰ تا ۵۰۰ گرم از رسوبات رودخانه‌ای تهیه شد. نمونه‌برداری از عمق ۰ تا ۵ سانتی‌متری (فیض‌نیا، ۲۰۰۸؛ روتون و همکاران، ۲۰۱۱) نقاط منتخب و همچنین رسوبات خروجی حوزه به صورت سیستماتیک تصادفی و در هر نقطه ۵ تکرار برای هر نمونه انجام شد. در مجموع از ۸ نقطه نمونه‌برداری (۷ واحد و خروجی) ۴۰ نمونه برداشت شد که ۳۵ نمونه مربوط به واحدهای زمین‌شناسی (۵ تکرار برای ۷ واحد) و ۵ نمونه مربوط به خروجی حوزه است.

2. Udden-Wentworth Sieve Series
3. American Society for Testing and Materials
4. Binocular
5. Yap
6. Inductively Coupled Plasma Optical Emission Spectrometry

1. Siegel

می‌رسد. درخت تصمیم شامل تعدادی از الگوریتم‌ها مثل CHAID و ID 3.C4.C5 در طبقه‌بندی است (بريمن^{۱۰}، ۲۰۰۱).
۲. روش جنگل تصادفی^{۱۱}: این روش یک تکنیک مدرن ناپارامتری^{۱۲} و متعلق به خانواده روش‌های دسته‌جمعی است. در حال حاضر یکی از بهترین الگوریتم‌های یادگیری است و برای بسیاری از مجموعه داده‌ها، دسته‌بندی را با صحت بالایی انجام می‌دهد و برخلاف مدل‌های کلاسیک چون رگرسیون^{۱۳} که تنها بر پایه یک مدل تکیه دارند با استفاده از صدها و هزاران درخت از اطلاعات بیشتری در داده‌ها استفاده می‌کند تا بتوان استنباط بهتری از متغیرها داشت. این الگوریتم از جمله دسته‌بندهایی است که متد Bagging را به کار می‌گیرد و حاوی چندین درخت تصمیم است که خروجی آن از خروجی‌های درخت‌های انفرادی به دست می‌آید (همان).

۳. روش‌های رگرسیونی^{۱۴}: مدل‌های رگرسیون بر اساس این نظریه ساخته شده‌اند که اگر دو عامل به یکدیگر بستگی داشته باشند، تغییر یکی با تغییر دیگری قرین خواهد شد. هر قدر ارتباط دو عامل به یکدیگر نزدیک‌تر و قوی‌تر باشد، ضریب همبستگی تغییرات آن‌ها بزرگ‌تر و به یک نزدیک‌تر خواهد شد. این روش انواع مختلفی دارد که مهم‌ترین آن‌ها روش رگرسیون خطی^{۱۵}، روش رگرسیون درجه دوم خالص^{۱۶}، و روش رگرسیون اثر متقابل^{۱۷} است (بريمن و همکاران، ۱۹۸۴).

۴. روش تحلیل ممیزی (آنالیز تشخیصی)^{۱۸}: تحلیل ممیزی تکنیکی چندمتغیره است که با جدا کردن مجموعه‌های متمایز مشاهده‌ها و با تخصیص مشاهده‌های جدید به دسته‌های از پیش تعریف‌شده سروکار دارد. آنالیز تشخیصی برای طبقه‌بندی پاسخگویان بر اساس کدهای یک متغیر وابسته اسمی دو یا چندوجهی به کار می‌رود. تحلیل تشخیصی ترکیب دو یا چند متغیر مستقل را که به بهترین وجه تفاوت بین دو

رسوبات مخزن فاقد مقدار بود حذف شد و در نهایت برای تمامی نمونه‌ها ۲۱ عنصر شامل Li، K، Cu، Cr، Co، Cr، B، Fe و V، Y، Zn، Ba، Mn، Sr، Mg، Al، Ca، Ti، Pb، P، Na اندازه‌گیری گردید. بعد از جمع‌آوری و ثبت داده‌ها، تحلیل‌های آماری پژوهش با استفاده از نرم‌افزارهای MATLAB و EXCEL ۲۰۱۳ صورت گرفت. آنالیزهای آماری تحقیق شامل ورود داده‌های حاصل به محیط نرم‌افزار MATLAB و سپس اجرا و ارزیابی الگوریتم‌های به کار رفته است. داده‌های ورودی به الگوریتم‌ها شامل پارامترهای دانه‌بندی و بافت‌سنجی رسوبات، غلظت عناصر شیمیایی و نهایتاً فراوانی سنگ‌ها می‌باشند که همگی به صورت کمی وارد آنالیز گردیدند. شایان ذکر است که علاوه بر نقطه خروجی، در ۲۱ نقطه از سطح حوضه، تعداد ۲۱ عنصر شیمیایی، ۱۶ پارامتر بافت‌سنجی و دانه‌بندی و نیز درصد فراوانی ۱۱ واحد سنگ‌شناسی شاخص برآورد شدند. لذا مجموعاً ۱۰۵۶ داده (۴۸ متغیر در ۲۲ نقطه) در آنالیز مورد استفاده قرار گرفت. الگوریتم‌های مورد استفاده در این تحقیق شامل جنگل تصادفی^۱، درخت‌های طبقه‌بندی^۲، شبکه عصبی مصنوعی، تحلیل ممیزی^۳، رگرسیون لجستیک^۴، انواع ماشین بردار پشتیبان، روش‌های رگرسیونی، روش گروهی مدل‌سازی داده‌ها^۵، درخت خطی محلی^۶، روش تشخیص الگو^۷ و k نزدیک‌ترین همسایه^۸ بودند که برای معرفی هریک از این الگوریتم‌ها در ادامه مختصر توضیحاتی آورده شده است.

۱. روش درخت تصمیم: درخت‌ها در هوش مصنوعی^۹ برای نمایش مفاهیم مختلفی نظیر ساختار، جملات، معادلات و... استفاده می‌شوند. درخت‌های تصمیم روشی برای نمایش یک سری از قوانین هستند که منتهی به یک رده می‌شوند. این درختان نمونه‌ها را به نحوی دسته‌بندی می‌کنند که از ریشه به سمت پایین رشد می‌کند و در نهایت به گره‌های برگ

10. Breiman
 11. Random Forest
 12. Nonlinear
 13. Regression
 14. Regression Method
 15. linear Regression
 16. Purequadratic Regression
 17. Interaction Regression
 18. Discriminant analysis

1. Random Forest
 2. Classification Trees
 3. Discriminant Analysis
 4. Logistic Regression
 5. Group Method of Data Handling
 6. Local Linear Model Tree
 7. Pattern Recognition Method
 8. K Nearest Neighbor
 9. Artificial Intelligence

مبنای آنوا^{۱۲} و اسپلاین^{۱۳} هستند. از مزایای ماشین بردار پشتیبان، توانایی حل مسائل طبقه‌بندی پیچیده با تعداد لایه‌های زیاد و نمونه‌های آموزشی کم است (اورال بنا^{۱۴} و همکاران، ۲۰۱۰).

۷. روش آنالیز k نزدیک‌ترین همسایه (KNN):^{۱۵} آنالیز

نزدیک‌ترین همسایه^{۱۶} یکی از روش‌های طبقه‌بندی متغیرها بر اساس تشابه آن‌ها با یکدیگر است که الگوی داده‌ها را بدون نیاز به الگوهای از پیش مشخص، طبقه‌بندی می‌نماید. در این روش سعی می‌شود تا ویژگی‌های نقاط داده از روی ویژگی‌های نزدیک‌ترین همسایگانشان تعیین شوند. یکی از بهترین طبقه‌بندی‌ها، طبقه‌بندی K نزدیک‌ترین همسایه است که این طبقه‌بندی، نمونه‌آزمون را متعلق به کلاسی می‌داند که بیشترین آرا را در بین k نزدیک‌ترین همسایگان آن داشته باشد. طبقه‌بندی KNN به دلیل قابلیت درک بالا و عدم نیاز به ایجاد فرضیه روی داده‌ها، روشی ساده و پرکاربرد می‌باشد.

۸. روش شبکه عصبی مصنوعی^{۱۷}: علی‌رغم به‌کارگیری

یک ساختمان ساده در این روش، سرعت و قدرت محاسباتی آن به شدت مورد توجه قرار گرفت. روش شبکه عصبی از سیستم عصبی الهام گرفته و از ساختار مغز و اعصاب انسان پیروی می‌کند (ریچارد^{۱۸}، ۲۰۱۳). امروزه مدل شبکه عصبی در بخش‌های مختلف علوم به‌منظور مدل‌سازی روابط پیچیده غیرخطی به‌کاررفته گرفته می‌شود و تا حدودی جایگزین مدل‌های آماری شده است؛ زیرا شبکه‌های عصبی مصنوعی بدون نیاز به حل معادلات دیفرانسیل جزئی، غیرخطی بودن فرایند موردنظر را شبیه‌سازی می‌کنند و حتی زمانی که مجموع داده‌های آموزشی حاوی داده‌های خطا دار باشد، عملکرد مناسبی را نشان می‌دهند. شبکه‌های چندلایه دارای توانایی بیشتری بوده که در آن یک لایه ورودی وجود دارد که اطلاعات را دریافت می‌کند و تعدادی لایه مخفی دارد که

گروه را تبیین می‌کند، نشان می‌دهد. هدف کلی تحلیل ممیزی به وجود آوردن یک ترکیب خطی بین متغیرهاست که از آن برای گروه‌بندی افراد استفاده شود. این ترکیب خطی، یک مسئله پیچیده و چندمتغیره را به یک مسئله آماری ساده و یک متغیره تبدیل می‌کند (پیتو^۱ و همکاران، ۲۰۱۲).

۵. مدل درخت خطی محلی^۲: مدل درخت خطی محلی

که در آن از نوعی مدل فازی عصبی خطی محلی استفاده شده است، الگوریتمی بر اساس استراتژی تقسیم و حل می‌باشد که در آن حل مسئله پیچیده از طریق تقسیم به تعدادی زیرمسئله کوچک‌تر صورت می‌پذیرد. بنابراین مشخصات این مدل‌های فازی عصبی به مقدار زیادی، به ساختار الگوریتم به‌کاربرده شده جهت تقسیم‌بندی وابسته است. این الگوریتم برای رسیدن به خروجی بهتر (خروجی با خطای کمتر) فضای مسئله را به تعدادی مدل خطی محلی یا LLM تقسیم می‌کند و پس از پیدا کردن بدترین LLM (LLM با خطای بیشتر) با تقسیم آن به دو LLM، الگوریتم را ادامه می‌دهد تا به انتها برسد (کشاورز امامی و کارولوکس، ۲۰۰۷).

۶. روش ماشین‌های بردار پشتیبان^۳: ماشین بردار پشتیبان

یکی از روش‌های یادگیری با نظارت^۴ است که از آن برای طبقه‌بندی و رگرسیون استفاده می‌کنند. این روش نسبتاً جدید است و در سال‌های اخیر کارایی خوبی نسبت به روش‌های قدیمی‌تر طبقه‌بندی نشان داده است. در این روش با استفاده از همه باندها و یک الگوریتم بهینه‌سازی، نمونه‌هایی که مرز کلاس‌ها را تشکیل می‌دهند، به دست می‌آید و با استفاده از آن‌ها یک مرز تصمیم‌گیری خطی بهینه^۵ برای جدا کردن کلاس‌ها محاسبه می‌شود؛ این نمونه‌ها را بردارهای پشتیبان می‌گویند. این روش چند هسته مختلف را به‌طور پیش‌فرض پشتیبانی می‌کند که شامل هسته‌های خطی^۶، چندجمله‌ای^۷، شعاع مبنای^۸، تانژانت هیپربولیک^۹، لاپلاسی^{۱۰}، بسل^{۱۱}، تابع شعاع

10. Laplacian
11. Bessel
12. Anova RBF
13. Spline
14. Ouralbena
15. K Nearest Neighbor
16. Nearest Neighbor Analysis
17. Artificial Neural Network
18. Richard

1. Pinto
2. LOLIMOT (Local Linear Model Tree)
3. Support vector machines
4. Supervised Learning
5. Optimized
6. Linear
7. Polynomial
8. Radial Basis
9. Hyperbolic Tangent

دست آورد (کاکائی لعدانی، ۲۰۱۴).

برازش الگوریتم‌ها

همان‌طور که گفته شد، برنامه‌نویسی هریک از این الگوریتم‌ها در نرم‌افزار MATLAB صورت گرفت. بدین صورت که در تمامی الگوریتم‌ها، تعداد ۷۰ درصد داده‌ها به منظور آموزش و ۳۰ درصد باقی‌مانده برای آزمون (صحت‌سنجی) انتخاب شدند و خروجی‌های مورد نظر استخراج گردید تا بهترین الگوریتم‌ها در تفکیک منابع رسوبی حوضه (هفت واحد سنگ‌شناسی به‌عنوان متغیر هدف)، در مراحل آموزش و آزمون تعیین گردد.

ارزیابی الگوریتم‌ها بر اساس چهار گروه ورودی مختلف اجرا شد. مرحله اول بر اساس ۲۱ متغیرهای ژئوشیمیایی شامل عناصر B, Ce, Co, Cr, Cu, K, Li, Na, P, Pb, Ti, V, Y, Zn, Fe, Ba, Mn, Sr, Mg, Al, Ca، مرحله دوم بر اساس ۱۱ گروه سنگ‌شناسی همگن شامل کوارتز، توف، لاتیت، داسیت، آندزیت، دولومیت، کلسیت، توف آندزیتی، آندزیت‌لیتیک‌دار، پیروژنیک و نمک با استفاده از نمونه‌های خرده‌سنگ‌ها و ذرات با قطر بیش از ۴۰۰۰، ۲۰۰۰، ۱۰۰۰، ۵۰۰، ۲۵۰، ۱۲۵، ۶۳ و کوچک‌تر از ۶۳ میکرون و مرحله سوم بر اساس دانه‌بندی ذرات شامل پارامترهای D10, D50, D90، درصد شن، درصد سیلت، درصد رس، چولگی و کشیدگی ذرات انجام شد. در مرحله چهارم نیز با ورود تمامی متغیرهای ژئوشیمیایی، دانه‌بندی و سنگ‌شناسی، الگوریتم‌های طبقه‌بندی اجرا شد.

در انتها برای حصول اطمینان از روند مدل‌سازی، اعتبارسنجی و ارزیابی دقت مدل‌های مورد استفاده، از معیارهای آماری ضریب تبیین^۶ (R^2) و میانگین مربع خطا^۷ (MSE) استفاده شد که به ترتیب از روابط (۱) و (۲) به دست می‌آیند که در این روابط مؤلفه n تعداد داده‌های مورد ارزیابی و x_i و y_i ، آمین داده برآوردی و اندازه‌گیری شده است.

$$R^2 = \frac{\sum (y_i - x_i)^2}{\sum (y_i - x_i)^2} \quad (1)$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2 \quad (2)$$

اطلاعات را از لایه‌های قبلی می‌گیرد و در نهایت، یک خروجی وجود دارد که نتیجه محاسبات به آن‌ها رفته و خروجی آن، خروجی نهایی شبکه است (آبراهام، ۲۰۰۵).

۹. روش رگرسیون لجستیک^۱: رگرسیون لجستیک یکی از انواع مدل‌های خطی تعمیم‌یافته است که برای تحلیل وجود یا عدم وجود متغیر وابسته بسیار مناسب می‌باشد. این مدل از تابع لجوجیت به‌عنوان تابع پیوند استفاده می‌کند و خطایش از توزیع چندجمله‌ای پیروی می‌کند. هدف از تحلیل رگرسیونی لجستیک، دستیابی به مدلی مناسب و در عین حال ساده برای بررسی ارتباط بین متغیر وابسته با یک یا مجموعه‌ای از متغیرهای مستقل است (چن و وانگ^۲، ۲۰۰۷).

۱۰. روش تشخیص الگو^۳: این روش شاخه‌ای از هوش مصنوعی است که با طبقه‌بندی و توصیف مشاهدات سروکار دارد. شناسایی الگو به ما کمک می‌کند داده‌ها را با تکیه بر دانش قبلی یا اطلاعات آماری استخراج‌شده از الگوها، طبقه‌بندی نماییم. تشخیص الگو عبارت است از دسته‌بندی و تفکیک الگوهای خاص بر اساس ویژگی‌های از پیش تعریف‌شده از مجموعه‌ای از داده‌های در دسترس. این روش به‌اجمال، با مسائل خوشه‌بندی و طبقه‌بندی سروکار دارد و دربرگیرنده طیف گسترده‌ای از روش‌های آماری کلاسیک الگوریتم‌های هوشمند شبکه‌های عصبی و منطق فازی است (هان و کامبر^۴، ۲۰۰۱).

۱۱. روش گروهی مدل‌سازی داده‌ها^۵: این الگوریتم قابلیت استفاده در موضوعات متنوعی چون کشف روابط، پیش‌بینی، مدل‌سازی، بهینه‌سازی و شناخت الگوریتم‌های غیرخطی را دارد. ویژگی خاص این الگوریتم استنتاجی، قابلیت شناسایی و غربال کردن متغیرهای کم‌اثر ورودی در دوره آموزش شبکه و حذف آن‌ها از روند شبیه‌سازی در دوره آزمون است. بدین ترتیب می‌توان با انجام یک فرایند قیاسی، در چند مرحله تکرار، متغیرهای کم‌اثرتر را حذف کرد و نهایتاً مدل بهینه و پیش‌بینی را بر اساس معیارهای متداول خطا به

1. Logistic Regression
2. Chen and Wang
3. Pattern Recognition
4. Han and Kamber
5. Group Method of Data Handling

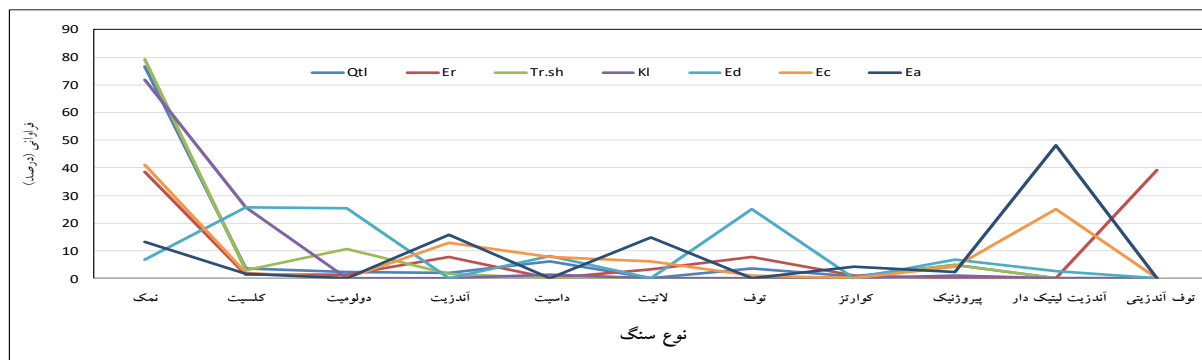
6. Coefficient of Determination
7. Mean Squared Error

نتایج

وجود دارند. جدول (۱) نشان می‌دهد تا D50، متغیر دانه‌بندی روند ثابتی داشته و سپس به‌صورت افزایشی عمل کرده و D90، بیشترین مقدار غلظت را در واحدهای زمین‌شناسی نشان می‌دهد به این دلیل که قطری است که ۹۰ درصد از ذرات از آن کوچک‌ترند و مجدد روند نزولی را ارائه کرده است. نتایج مربوط به بررسی غلظت عناصر ژئوشیمیایی نشان داد غلظت عناصر ژئوشیمیایی در هفت گروه زمین‌شناسی مورد بررسی از میزان متوسط تا زیاد برخوردار است. به‌طور کلی به‌ترتیب عناصر Ca، Fe، Mg، Al دارای بیشترین غلظت در نمونه‌های خاک بودند و سایر عناصر به میزان کمتر وجود داشتند. به‌ویژه B و Co کمترین مقدار را به خود اختصاص دادند. در بین واحدهای زمین‌شناسی، آهک اوریتولین‌دار بیشترین سهم را در رسوبات مخزن سد دارد.

همچنین نتایج مربوط به طبقه‌بندی منابع رسوبی پس از ورود داده‌های آزمایشگاهی مورد نظر به محیط نرم‌افزار MATLAB و برازش الگوریتم‌ها، در قالب جدول‌های (۳) تا (۶) در مراحل آموزش و آزمون ارائه شده‌اند.

نمودار متوسط مقادیر فراوانی نمونه سنگ‌ها، متغیرهای دانه‌بندی و متوسط غلظت عناصر ژئوشیمیایی در شکل (۲) و جداول (۱) و (۲) نمایش داده شده است. در بحث سنگ‌شناسی سنگریزه‌های قشر بالایی رسوبات پشت مخزن سد بیشتر از جنس کلسیت و نمک بوده، در رسوبات آبرفتی قدیم و گدازه‌های آتشفشانی بیشترین نوسانات فراوانی سنگ‌ها را داریم. در واحد نهشته‌های آبرفتی قدیم نیز پس از نمک که بیشترین مقدار بوده، سایر سنگ‌ها به‌دلیل خاصیت رسی و سیلتی کاهش یافته تا مواردی مانند پیروژنیک، آندزیت لیتیک‌دار تقریباً به صفر رسیده است. به این دلیل که نهشته‌های قدیمی اغلب کنگلومرای هستند و مقادیر رسی و سیلتی کمی را دارند. کاهش یافته مقادیر کلسیت، دولومیت و توف مقدار زیادی را به خود اختصاص داده است؛ زیرا در نهشته‌های قدیم کمتر رسی و سیلتی وجود دارند و اغلب کنگلومرای می‌باشند. بنابراین مقدار آندزیت لیتیک دارو توف آندزیتی در آن‌ها کم است و در نهشته‌های آبرفتی حاضر به مقدار زیاد



شکل (۲): متوسط فراوانی سنگ‌ها در واحدهای مختلف زمین‌شناسی حوضه با تحلیل مورفوسکوپی نمونه‌ها

Figure (2): The average frequency of rocks in different geologic units of watershed using morphoscopic analysis of samples

جدول (۱): متوسط متغیرهای دانه‌بندی رسوبات در واحدهای مختلف زمین‌شناسی حوضه

Table (1): The average of granulometric variables of sediments in different geologic units of watershed

کشیده‌گی	تولگی	D90	D50	D10	۶۳ >	۶۳ <	۱۲۵ <	۲۵۰ <	۵۰۰ <	۱ <	۲ <	۴ <	درصد ریز	درصد سیلت	درصد رسی	واحد	زمین‌شناسی
					μm	μm	μm	μm	μm	mm	mm	mm					
۲/۴	۰/۷	۳۳۳/۱	۶۱/۷	۷/۷	۸۵/۰	۴۷/۳	۶/۷	۳/۴	۴/۸	۷/۱	۶/۴۱	۲۷/۸	۳۸/۵	۴۰/۹	۲۰/۵	Qtl	
۱/۹	-۰/۵	۴۵۰۵/۶	۷۸۲/۹	۱۲/۷	۵۰/۶	۱۷/۶	۶/۸	۸/۲	۱۸/۳	۳۳/۱	۲۷/۲	۳۲/۸	۵۷/۸	۲۲/۷	۱۹/۵	Er	
۲/۶	۱/۰	۴۲۲۳/۱	۳۹/۳	۶/۲	۱۶۶/۰	۳۹/۰	۷/۸	۷/۰	۴/۸	۶/۷	۱۰/۰	۳۳/۰	۳۳/۸	۴۲/۶	۲۳/۶	Tr.sh	
۱/۶	۰/۱	۴۹۷۹/۷	۷۶۸/۴	۸/۸	۷۲/۶	۳۲/۰	۳/۶	۱/۹	۲/۳	۶/۲	۱۴/۰	۶۰/۷	۳۲/۱	۴۹/۸	۱۸/۱	Kl	
۲/۷	-۰/۵	۱۱۱۸/۹	۱۵۶/۵	۲۸/۲	۲۸/۳	۴۷/۱	۳۴/۰	۳۳/۷	۲۳/۸	۲۱/۴	۳/۹	۱/۳	۶۸/۹	۲۰/۶	۱۰/۵	Ed	
۲/۴	-۰/۴	۴۹۲۵/۵	۸۹۳/۱	۲۷/۶	۴۱/۲	۳۳/۶	۱۱/۴	۷/۱	۸/۶	۱۶/۷	۱۶/۸	۶۱/۴	۵۷/۸	۳۱/۹	۱۰/۳	Ec	
۲/۹	-۰/۶	۳۱۸۲/۱	۴۰۴/۶	۳۹/۵	۲۵/۸	۲۱/۸	۲۲/۴	۳۳/۸	۳۰/۴	۲۶/۹	۱۰/۳	۲۰/۷	۷۹/۸	۱۱/۳	۹/۰	Ea	

جدول (۲): متوسط غلظت عناصر ژئوشیمیایی در واحدهای مختلف زمین‌شناسی حوضه

Tablel (2): The average concentration of geochemical elements sediments in different geologic units of watershed

Fe	Ca	Al	Mg	Sr	Mn	Ba	Zn	Y	V	Ti	Pb	P	Na	Li	k	Cu	Cr	Co	Ce	B	واحد زمین‌شناسی
۴۰۵/۸۲	۱۱۴۸/۷۱	۱۳۴/۳۷	۳۶۲/۲۵	۳/۳۴	۱۰/۹۴	۲/۷۳	۱/۳۴	۰/۱۹	۰/۴۴	۰/۵۷	۰/۲۲	۶/۲۵	۲۸/۰۵	۰/۵۸	۴/۶۰	۰/۵۵	۰/۷۱	۰/۲۲	۰/۲۲	۰/۰۸	Qtl
۲۷۳/۶۱	۶۲۲/۸۷	۹۶/۶۷	۲۸۶/۴۴	۲/۰۲	۹/۴۱	۴/۷۱	۱/۱۶	۰/۲۷	۰/۲۸	۰/۵۶	۰/۲۸	۳/۶۷	۲۲/۸۹	۰/۴۳	۴/۸۲	۰/۴۹	۰/۵۰	۰/۱۸	۰/۸۶	۰/۰۶	Er
۴۲۳/۴۹	۱۸۸۴/۶۰	۱۲۸/۹۲	۳۴۲/۰۲	۵/۱۶	۹/۳۳	۲/۸۲	۱/۲۷	۰/۱۸	۰/۴۴	۰/۴۷	۰/۱۸	۴/۹۱	۲۷/۴۴	۰/۶۰	۴/۸۷	۰/۵۷	۰/۷۱	۰/۲۲	۰/۲۱	۰/۰۶	Tr.sh
۴۲۸/۵۴	۱۳۷۴/۰۸	۱۰۹/۶۳	۴۲۹/۹۷	۳/۱۷	۱۱/۸۳	۲/۴۶	۱/۳۸	۰/۱۸	۰/۴۹	۰/۶۶	۰/۲۵	۵/۴۹	۳۱/۶۲	۰/۵۹	۳/۹۳	۰/۵۵	۰/۷۶	۰/۲۵	۰/۱۷	۰/۰۸	Kl
۱۲۱/۶۴	۱۳۳۴/۲۰	۵۱/۴۶	۴۰۳/۵۴	۲/۳۴	۸/۸۴	۲/۴۳	۱/۰۱	۰/۱۴	۰/۱۷	۰/۳۲	۰/۱۲	۱/۶۱	۲۵/۷۳	۰/۳۹	۲/۴۰	۰/۳۱	۰/۲۵	۰/۱۶	۰/۲۸	۰/۰۸	Ed
۲۰۴/۷۷	۶۳۹/۶۷	۷۹/۳۲	۲۵۸/۲۱	۱/۷۱	۷/۴۹	۱/۷۴	۰/۷۱	۰/۱۵	۰/۳۷	۰/۲۳	۰/۲۵	۵/۳۸	۲۰/۴۱	۰/۳۴	۲/۸۴	۰/۲۶	۰/۵۵	۰/۱۲	۰/۴۸	۰/۰۶	Ec
۴۴۷/۲۲	۶۴۲/۷۵	۹۷/۱۶	۲۹۹/۱۱	۲/۰۷	۱۰/۹۸	۲/۰۸	۱/۲۳	۰/۲۴	۰/۵۰	۱/۳۳	۰/۲۶	۵/۹۰	۳۴/۰۰	۰/۵۰	۳/۰۴	۰/۴۳	۰/۶۶	۰/۱۸	۰/۴۸	۰/۰۷	Ea

۹/۶۲	۰	۰/۱۴	۱	تحلیل ممیزی
۶/۲۸	۰/۹۰	۰	۰/۷۵	درخت طبقه‌بندی
۲/۸۶	۰/۹۴	۰/۱۴	۰/۳۵	درخت رگرسیونی
۶/۱۰	۰	۰/۴۳	۱	جنگل تصادفی
۵/۹۰	۰	۰/۳۸	۱	نزدیک‌ترین همسایه
۵/۲۸	۰	۰/۴۸	۱	ماشین بردار پشتیبان خطی
۴/۹۵	۰	۰/۵۱	۱	ماشین بردار پشتیبان چندجمله‌ای
۵/۵۲	۰	۰/۴۸	۱	ماشین بردار پشتیبان لاپلاسیان
۶/۲۴	۱/۶۲	۰	۰/۶۶	ماشین بردار پشتیبان تابع شعاع مبنا
۴/۸۱	۰	۰/۴۳	۱	ماشین بردار پشتیبان چندگانه
۴/۹۸	۰/۵۱	۰/۲۷	۰/۶۵	روش گروهی مدل‌سازی داده‌ها

جدول (۵): نتایج طبقه‌بندی منابع رسوبی با ورود متغیرهای سنگ‌شناسی

Tablel (5): The results of sediment sources classification using geologic variables

میانگین مربع خطا	ضریب تبیین	نوع الگوریتم		
(آموزش)	(آموزش)	(آموزش)	(آموزش)	
۵/۰۵	۱/۰۴	۰/۲۴	۰/۳۹	رگرسیون خطی
۹/۷۶	۰/۳۵	۰/۱۴	۰/۶۵	رگرسیون درجه دوم خالص
۴/۵۲	۰/۳۵	۰/۱۹	۰/۶۵	رگرسیون اثر متقابل
۴/۸۱	۰/۳۹	۰/۲۸	۰/۶۱	رگرسیون درجه دوم
۴/۸۷	۲/۵۰	۰/۳۰	۰/۳	شبکه عصبی
۵/۰۸	۳/۴۰	۰/۱۳	۰/۲۴	روش تشخیص الگو
۵/۱۰	۰	۰/۶۳	۱	تحلیل ممیزی
۱۴/۸۱	۴/۳۲	۰	۰/۷۵	درخت طبقه‌بندی
۶/۳۳	۱/۱۷	۰/۱۹	۰/۳۴	درخت رگرسیونی
۱/۹۲	۰	۰/۷۵	۱	جنگل تصادفی
۳/۷۱	۰	۰/۶۲	۱	نزدیک‌ترین همسایه
۲/۰۵	۰	۰/۶۲	۱	ماشین بردار پشتیبان خطی
۲/۳۰	۰	۰/۶۰	۱	ماشین بردار پشتیبان چندجمله‌ای
۲/۱۹	۰	۰/۶۲	۱	ماشین بردار پشتیبان لاپلاسیان
۴/۳۳	۱/۳۶	۰/۲۸	۰/۸۷	ماشین بردار پشتیبان تابع شعاع مبنا
۵/۸۱	۰	۰/۶۲	۱	ماشین بردار پشتیبان چندگانه
۷/۳۳	۰/۴۸	۰/۴۱	۰/۶۶	روش گروهی مدل‌سازی داده‌ها

جدول (۳): نتایج طبقه‌بندی منابع رسوبی با ورود متغیرهای ژئوشیمیایی

Tablel (3): The results of sediment sources classification using geochemical variables

میانگین مربع خطا	ضریب تبیین	نوع الگوریتم		
(آموزش)	(آموزش)	(آموزش)	(آموزش)	
۷/۸۶	۰/۶۱	۰/۱۹	۰/۴۹	رگرسیون خطی
۳/۴۸	۰/۴۷	۰/۲۸	۰/۵۳	رگرسیون درجه دوم خالص
۴/۷۱	۰/۵۰	۰/۲۴	۰/۵۰	رگرسیون اثر متقابل
۵/۷۱	۰/۳۸	۰/۱۴	۰/۶۲	رگرسیون درجه دوم
۳/۱۳	۲/۱۱	۰/۱۴	۰/۳۱	شبکه عصبی
۳/۲۲	۲/۱۶	۰/۱۶	۰/۲۷	روش تشخیص الگو
۸/۰۵	۰	۰/۴۸	۱	تحلیل ممیزی
۵/۱۴	۱/۷۰	۰/۲۴	۰/۷۰	درخت طبقه‌بندی
۲/۰۰	۰/۷۹	۰/۲۸	۰/۳۸	درخت رگرسیونی
۳/۷۸	۰	۰/۴۰	۱	جنگل تصادفی
۲/۷۶	۰	۰/۵۲	۱	نزدیک‌ترین همسایه
۳/۰۵	۰	۰/۶۲	۱	ماشین بردار پشتیبان خطی
۴/۶۵	۰	۰/۴۳	۱	ماشین بردار پشتیبان چندجمله‌ای
۱/۴۸	۰	۰/۶۷	۱	ماشین بردار پشتیبان لاپلاسیان
۳/۹۵	۱/۳۸	۰	۰/۶۸	ماشین بردار پشتیبان تابع شعاع مبنا
۵/۸۶	۰	۰/۳۳	۱	ماشین بردار پشتیبان چندگانه
۵/۲۱	۰/۴۶	۰/۲۲	۰/۶۷	روش گروهی مدل‌سازی داده‌ها

جدول (۴): نتایج طبقه‌بندی منابع رسوبی با ورود متغیرهای دانه‌بندی

Tablel (4): The results of sediment sources classification using granulometric variables

میانگین مربع خطا	ضریب تبیین	نوع الگوریتم		
(آموزش)	(آموزش)	(آموزش)	(آموزش)	
۷/۷۶	۰/۷۱	۰/۱۴	۰/۴۶	رگرسیون خطی
۶/۱۰	۰/۳۸	۰/۲۸	۰/۶۲	رگرسیون درجه دوم خالص
۹/۵۷	۰/۴۶	۰/۱۴	۰/۵۴	رگرسیون اثر متقابل
۳/۴۸	۰/۴۶	۰/۱۹	۰/۵۴	رگرسیون درجه دوم
۴/۲۱	۲/۷۱	۰/۱۳	۰/۲۷	شبکه عصبی
۴/۱۷	۳/۳۱	۰/۱۹	۰/۲۴	روش تشخیص الگو

جدول (۶): نتایج طبقه‌بندی منابع رسوبی با ورود تمامی متغیرها

Tablel (6): The results of sediment sources classification using whole variables

میانگین مربع خطا		ضریب تبیین		نوع الگوریتم
(آزمون)	(آموزش)	(آزمون)	(آموزش)	
۱۰/۱۹	۰/۴۱	۰/۰۹	۰/۵۸	رگرسیون خطی
۶/۴۲	۰/۴۴	۰/۲۳	۰/۵۵	رگرسیون درجه دوم خالص
۹/۴۷	۰/۴۸	۰/۱۴	۰/۵۱	رگرسیون اثر متقابل
۸/۶۱	۰/۴۳	۰/۰۹	۰/۵۶	رگرسیون درجه دوم
۳/۵۳	۲/۲۶	۰/۲۳	۰/۳۳	شبکه عصبی
۴/۴۴	۲/۷۳	۰/۱۵	۰/۲۸	روش تشخیص الگو
۲/۹۰	۰	۰/۷۶	۱	تحلیل ممیزی
۱۴/۶۶	۳/۴۷	۰	۰/۷۵	درخت طبقه‌بندی
۴/۴۷	۰/۹۱	۰/۱۴	۰/۲۵	درخت رگرسیونی
۲/۸۵	۰	۰/۶۰	۱	جنگل تصادفی
۱/۱۹	۰	۰/۷۱	۱	نزدیک‌ترین همسایه
۳/۸۵	۰	۰/۷۱	۱	ماشین بردار پشتیبان خطی
۳/۴۲	۰	۰/۵۵	۱	ماشین بردار پشتیبان چندجمله‌ای
۲/۸۰	۰/۱۳	۰/۴۲	۰/۹۵	ماشین بردار پشتیبان لاپلاسیان
۱/۵۲	۰	۰/۷۶	۱	ماشین بردار پشتیبان تابع شعاع مینا
۶/۶۱	۰	۰/۴۷	۱	ماشین بردار پشتیبان چندگانه
۶/۷۱	۰/۴۰	۰/۲۶	۰/۶۸	روش گروهی مدل‌سازی داده‌ها

بحث و نتیجه‌گیری

به‌دست آمده و پس از آن ورود توأم سه دسته متغیر، بیشترین دقت را داشته است. متغیرهای سنگ‌شناسی و دانه‌بندی نیز به‌ترتیب منجر به کمترین دقت طبقه‌بندی توسط الگوریتم‌ها شده‌اند.

نتایج تفسیر و مقایسه خروجی‌های الگوریتم‌های به‌کاررفته نشان داد که در مرحله آموزش، الگوریتم‌های تحلیل ممیزی، جنگل تصادفی، k نزدیک‌ترین همسایه، ماشین بردار پشتیبان خطی^۱، ماشین بردار پشتیبان چندجمله‌ای^۲، ماشین بردار چندگانه^۳ و ماشین بردار با تابع شعاع مینا^۴ با حداکثر مقدار ضریب تبیین ($R^2=1$) و مقدار میانگین مربع خطا برابر صفر، بهترین عملکرد را در تفکیک منابع رسوبی حوضه نشان دادند و به‌عنوان بهترین الگوریتم‌های پژوهش معرفی شدند. لذا الگوریتم‌های به‌کاررفته در پژوهش در تفکیک منابع رسوبی دقت بیشتری را از خود نشان دادند که با نتایج حدادچی و همکاران (۲۰۱۳) در انگشت‌نگاری رسوبات آبرفتی همخوانی

برای تعیین الگوریتم‌های مناسب جهت تفکیک منابع رسوبی حوضه مطالعاتی، از مقایسه مقدار ضریب تبیین و میانگین مربع خطا به‌صورت جداگانه برای داده‌های آموزش و آزمون استفاده شد. به این صورت که هرچه مقدار ضریب تبیین بیشتر و میانگین مربع خطا کمتر باشد، نشان‌دهنده دقت بیشتر و در نتیجه عملکرد بهتر الگوریتم است. طبق نتایج ارائه‌شده، به‌طور کلی دقت الگوریتم‌های طبقه‌بندی در مرحله آموزش بهتر از مرحله آزمون بوده و خطای کمتری را نیز به همراه داشته است. نکته قابل ذکر اینکه در اجرای الگوریتم‌ها با چهار نوع ورودی ژئوشیمیایی، سنگ‌شناسی، دانه‌بندی و ورود توأم سه گروه متغیر، ماشین‌های بردار پشتیبان بالاترین دقت و کمترین خطا را در طبقه‌بندی منابع رسوبی داشته‌اند. بنابراین با توجه به دقت مناسب روش‌های به‌کاررفته و ارائه نتایج موفقیت‌آمیز مدل‌ها برای تفکیک منابع رسوبی حوضه بهتر است تمامی عوامل فوق دخالت داده شوند. طبق جداول فوق بیشترین دقت طبقه‌بندی منابع رسوبی از ورود متغیرهای ژئوشیمیایی

1. Linear Support Vector Machine
2. Polynomial Support Vector Machine
3. Multiple Support Vector Machine
4. Radial Basis Function (RBF)

سطحی اراضی انجام گرفت، دلیل دقت بالای الگوریتم SVM نسبت به دیگر روش‌های داده‌کاوی این‌طور بیان شد که این روش از لحاظ محاسباتی بسیار سریع بوده و از قوانین بهینه‌سازی برای مکان‌یابی مرزهای بهینه بین کلاسه‌های کاربری استفاده می‌کند. در نتیجه می‌توان آن‌ها را به‌عنوان جایگزین مناسبی برای سایر الگوریتم‌های طبقه‌بندی نیز معرفی کرد. البته باید اذعان نمود برای همهٔ مسائل الگوی تعیین‌شده‌ای وجود ندارد و تعیین بهترین الگو برای هر مسئله نیازمند انجام سعی و خطاهای مکرر و آزمون روش‌ها در شرایط مد نظر است. بنابراین توصیه می‌شود در مطالعات دیگر نیز قابلیت‌های این موارد مورد بررسی قرار گیرد تا نتایج دقیق‌تری به دست آید. در مجموع می‌توان گفت پژوهش‌های انجام‌شده در سال‌های اخیر مانند حیات‌زاده و همکاران (۲۰۱۵)، کاکائی لخدانی (۲۰۱۳)، اولکی^۲ و همکاران (۲۰۱۶) و سایر پژوهشگران حاکی از برتری روش‌های هوش مصنوعی نسبت به سایر روش‌ها در محاسبات مربوط به فرسایش و رسوب است. پس می‌توان اذعان کرد داده‌کاوی روشی تحلیلی است که امروزه به‌طور وسیعی در تحقیقات پژوهشی کاربرد ویژه‌ای دارد و با گسترش و تنوع روزافزون الگوریتم‌ها، لزوم تحقیق و بررسی در زمینه‌های گوناگون مدل‌ها را ایجاد می‌کند. بنابراین استفاده از روش الگوریتم‌های طبقه‌بندی با توجه به محاسبات دقیق و عدم نیاز به صرف هزینه و وقت زیاد و ابزارهای پیشرفته توصیه می‌گردد.

دارد. اما مغایر با مطالعات پیشین، الگوریتم شبکهٔ عصبی نسبت به سایر الگوریتم‌ها عملکرد بهتری را در تفکیک منابع رسوبی حوضه نشان نداد. همچنین روش درخت رگرسیونی با مقادیر ضریب تبیین ۰/۲۸ و خطای ۰/۹۱، دارای ضعیف‌ترین عملکرد در این زمینه است. در میان داده‌های آزمون نیز بیشترین مقدار ضریب تبیین به دو الگوریتم تحلیل ممیزی و ماشین بردار پشتیبان با تابع شعاع مبنا با مقدار برابر ۰/۷۶ اختصاص یافته که با مقایسهٔ مقادیر خطای آن‌ها مشخص می‌شود الگوریتم ماشین بردار با تابع شعاع مبنا (RBF) با میزان خطای کمتر، دقیق‌ترین عملکرد را ارائه کرد و ضعیف‌ترین عملکرد مربوط به الگوریتم درخت طبقه‌بندی است.

همچنین در مرحلهٔ آزمون نیز، دقیق‌ترین الگوریتم باز هم مربوط به خانوادهٔ ماشین‌های بردار پشتیبان با هستهٔ تابع شعاع مبنا بوده که با مطالعهٔ میسرا و همکاران (۲۰۰۹) در شبیه‌سازی رواناب و رسوب همخوانی دارد. همچنین الگوریتم درخت طبقه‌بندی با بیشترین میزان خطا، ناکارآمدترین روش در این مرحله معرفی شده است. اما در پژوهش ستاری و همکاران (۲۰۱۶) بر روی رسوبات معلق رودخانه‌ای و کومارگویال (۲۰۱۴) بر روی داده‌های تولید رسوب، الگوریتم درخت M5 به‌عنوان روش برتر گزارش شد. بررسی پژوهش‌های انجام‌شده در زمینهٔ فرسایش و رسوب می‌توان نتیجه گرفت الگوریتم شبکهٔ عصبی نسبت به سایر الگوریتم‌ها نتایج دقیق‌تری را ارائه کرده و الگوریتم برتر معرفی شده است؛ از جمله مطالعهٔ حرما و همکاران (۲۰۱۵) در مطالعهٔ داده‌های رواناب و رسوب، ملسی و همکاران (۲۰۱۱) در بررسی رسوبات معلق رودخانه‌ای، زو و همکاران (۲۰۰۷) در مطالعهٔ رسوبات معلق و نائینی و همکاران (۲۰۰۸) در بررسی داده‌های غلظت رسوب. اما در پژوهش حاضر این موضوع تأیید نشد و روش شبکهٔ عصبی کمترین دقت و ضعیف‌ترین عملکرد را ارائه نمود که می‌تواند ناشی از تفاوت در نوع داده‌های ورودی و همچنین استفاده از ۱۰ الگوریتم برتر داده‌کاوی در کنار شبکهٔ عصبی مصنوعی باشد. همچنین در مطالعه‌ای که توسط کازوگلو و سالکسن^۱ در سال ۲۰۰۹ برای طبقه‌بندی پوشش

منابع

1. Abraham, A., 2005. Artificial neural networks, Oklahoma State University, Stillwater, USA. 908 PP.
2. Arabkhedri, M., 2014. An overview of the effective factors on the water erosion in Iran. *Journal of Land Management*, 2 (1): 23-35.
3. American Society for Testing and Materials (ASTM), 2008. Standard test method for particle-size analysis of soils. In: *Annual Book of ASTM Standards*. Philadelphia.
4. Besalatpour, A.A., Ayoubi, S.A., Hajabasi, D.A., 2016. Gamma test to select the optimal inputs in modeling soil shear strength using artificial neural network. *Journal of Soil and Water Conservation Researches (Agricultural Sciences and Natural Resources)*. 20 (1): 97-114.
5. Breiman, L., 2001. Application and analysis of random forests and machine learning. *Journal of Water Management*, 15(1): 5-32.
6. Breiman, L., Friedman J., Olshen R., and Stone, C., 1984. *Classification and Regression Trees*, Chapman & Hall/CRC Press, Boca Raton, FL.
7. Chan, C., Lewis, B., 2002. A basic primer on data mining, *Information Systems Management. Journal information System Management*, 19(4): 56-69.
8. Chen, Zh. And Wang, J., 2007. Landslide hazard mapping using logistic regression model in Mackenzie Valley, Canada. *Geomorphology*, Vol.42.
9. Demirci, M., Baltaci, A., 2013. Prediction of suspended sediment in river using fuzzy logic and multi linear regression approaches. *Neural Computing and Applications*, 23 (1): 145-151.
10. Feyznia, S., 2008. *Applied sedimentology with emphasis on soil erosion and sediment production*. Gorgan University of agricultural sciences and natural resources press, 356 pp.
11. Haddadchi, A., Ryder, D.S., Evrard, O. and Olley, J., 2013. Sediment fingerprinting in fluvial systems: review of tracers, sediment sources and mixing models. *International Journal of Sediment Research*, 28(4): 560-578.
12. Harma, N., Zakallah, M.D., Tiwari, H., Kumar, D., 2015. Runoff and sediment yield modeling using ANN, and support vector machines (case study: from Nepal watershed). *Ore Geology Reviews*, 17(9): 63-89.
13. Han, D. and Kamber, M., 2001. *Data Mining: Concepts and Techniques*. San Diego Academic Press.
14. Hayatzadeh, m., Chezgi, G., Dastorani, M.T., 2015. Evaluation of sediment rating curve and neural network using a combination of morphological parameters Baghabas area. *Journal of Agricultural Sciences and Natural Resources*, 19 (70): 101-119.
15. Joudi, A. R. & Sattari, M. T., 2017. Evaluation of the performance of queneil based methods in estimating the suspended rainfall of river (case Study: Sufy Chay, Maragheh). *Journal of Research in Natural Geography Vol 38(33)*: 413-429.
16. Kakaei-Lafdani, E., Moghaddamnia, A., Ahmadi, A. Ebrahimi, C., 2013. Daily suspended sediment load prediction using artificial neural networks and support vector machines. *Journal of Hydrology*, 478(25): 50-62.
17. Kavzoglu, T. and Colkesen. I., 2009. A kernel function analysis for support vector machines for land cover classification. *Journal of applied Earth Observation and Geoinformation*, 11(5):352-359.
18. Kakaei Lafdani, E., Pournemat Roudsari, A., Qaderi, K. and Moghaddam-Nia, A., 2014. Predicting the Volume of Suspended Sediments using GMDH and SVM Models Based on Principal Component Analysis. 9th International River Engineering Conference Shahid Chamran University, Ahwaz, pp: 22-24.
19. Keshavarz-Emami, R., karolouks, A., 2007. Local linear tree algorithm development (LOLIMOT) using a fuzzy validity function and credit for time series prediction. 1st Joint Congress on Fuzzy and Intelligent Systems Ferdowsi University of Mashhad, Iran, 29-31 Aug.
20. Kumar Goyal, M., 2014. Modeling of Sediment Yield prediction Using M5 Model Tree Algorithm and Wavelet Regression *Journal of Water Resources Management*, 28, 1991-2003.
21. Melesse, A. M., Ahmad, M. E., McClain, X., Wang, F. and Lim, Y.H., 2011. Suspended sediment load prediction of river systems: An artificial neural network approach. *Journal of Agricultural Water Management*, 98(5): 855-866.
22. Milhouse, R.T., 1998. *Modeling of instream flow needs: the link between sediment and aquatic habitat* Soil Sciences. Yazd University publication, Yazd, Iran, 516 pp.
23. Misra, D., Oommenb, T., Agarwal, A., Mishra, A. and Thompson, M., 2009. Application and analysis of support vector machine based simulation for runoff and sediment yield. *Biosystems engineering*, 6(2): 527- 535.
24. Naeni, S.T., Montazeri, M., Zamani, M.M. and Soltani, F., 2008. Sensitivity analysis of stimulus function of artificial neural network

- model in estimating sediment concentration. Proceeding of the 4th National congress of Civil engineering, 17-19 May, Tehran.
25. Oralbona, C., Castellini, B., Caputo, L. and sandini, G., 2010. On-line independent support vector machines pattern Recognition Application. *Journal of the International Society for the Prevention and Mitigation of Natural Hazard*, 10(6): 127-152.
26. Pinto, U., Maheshwar, B., Shrestha, S. and Morris, C., 2012. Modeling eutrophication and microbial risks in peri-urban river systems using discriminant function analysis, *Journal of water research*, 46(21): 6476- 6488.
27. Richards, J.A., 2013. Remote Sensing digital image analysis, fifth edition, Springer, 494 pp.
28. Rhoton, F.E., Emmerich, W.E., Nearing, M.A., Mc Chesney, D.S. and Ritchie, J.C., 2011. Sediment source identification in a semiarid Watershed at soil mapping unit scales. *Catena*, 87: 12-181.
29. Siegel. F.R., 2002, Environmental geochemistry of potentially toxic metals. Springer. Berlin Heidelberg New York, 212 pp.
30. Sattari, M.T., Rezazadehjudi, A., Safdari, F., Ghahramanzadeh, F., 2016. Performance evaluation methods, support vector regression modeling M5 model tree and suspended sediment Ahar Chai River. *Journal of Soil and Water Conservation*, 6 (1): 109-124.
31. Turnbull, L., Wainwright, J. and Brazier, R. E., 2008. A conceptual framework for understanding semi-arid land degradation: Eco hydrological interaction across multiple-space and time scales. *Journal of Ecohydrology*, 1(1): 23-34.
32. Ulke, A.G., Tayfur, R., Ozkul, S., 2009. Predicting suspended sediment loads and missing data for Gediz River, Turkey. *Journal of Hydrologic Engineering*, 14(9): 954-965.
33. Yap, C.A., Esmaeili, A., Tan, S., Omar, H., 2002. Correlations between speciation of Cd, Pb and Zn in sediment and their concentrations in total soft tissue of green- lipped mussel *Perna viridis* from the west coast of Peninsular Malaysia. *Environment International*, 28(1-2): 117-126.
34. Zhu, Y. M., Lu, X. X. and Zhou, Y., 2007. Suspended sediment flux modeling with artificial neural network: An example of the Longchuanjiang River in the Upper Yangtze Catchment. China. *Journal of Geomorphology*, 84(1): 111-125.

Using Data Mining Algorithms in Separation of Sediment Sources in Nodeh Watershed, Gonabad

Mehdi bashiri¹, Mahsa Ariapour², Ali Golkarian³

Received: 5/03/2018

Accepted: 2/07/2018

Extended Abstract

Introduction: Reduction of sediment supply requires the implementation of soil conservation and sediment control programs in the form of watershed management plans. Sediment control programs require identifying the relative importance of sediment sources, their quantitative *ascription* and identification of critical areas within the watersheds. The sediment source ascription involves two main steps so that in the first, several diagnostic tracers are selected for obvious and significant separation of potential sources of sediment and in the second step selected tracers for potential sources of sediment are compared, with corresponding values extracted from the sediment samples taken in the watershed outlet. Also, due to the large amount and complexity of data available, nowadays in geo- and environmental sciences, we face the need to develop and incorporate more robust and efficient methods for their analysis and modelling. Therefore recent fundamental progress in data mining algorithms can considerably contribute to the development of the emerging field - environmental data science.

Methodology: According to what was said, in this research, the data mining algorithms used to separate sediment sources in the Nodeh watershed of Gonabad located in Razavi-Khorasan province by using the geochemical (includes the 21 elements of Mg, Sr, Mn, Ba, Zn, Y, V, Ti, Pb, P, Na, Li, K, Cu, Cr, Co, Ce, B, Ca, Al and Fe), granulometric (includes the D_{90} , D_{50} , D_{10} , percent of sand, percent of silt, percent of clay, skewness and kurtosis and the diameters less than 1, 2 and 4 millimeters and less than 500, 250, 125 and 63 microns) and lithological variables (includes the quartz, tuff, laterite, dacite, andesite, dolomite, calcite, andesitic tuff, lithic andesite and salt). A set of 11 classification algorithms includes the decision tree, random forest, regression methods, discriminant analysis, local linear model tree, nearest neighbor analysis, support vector machine, logistic regression, artificial neural network, pattern recognition and group method of data handling programmed in the MATLAB software and the results compared based on the coefficient of determination and mean squared error.

Results and Discussion: Study of geochemical element concentrations in 7 geological units showed that the Ca, Fe, Mg and Al elements have the highest and B and Co have the lowest concentrations within the soil samples. Overall evaluation of classification algorithms in training stage showed that the discriminant analysis, random forest, k nearest neighbor and support vector machines with linear, polynomial, multiple and RBF kernels with maximum values of the coefficient of determination ($R^2=1$) and minimum values of the mean squared error ($RMSE=0$) are the most accurate algorithms in sediment source separation but the regression trees method has the worst performance. Also, at testing stage, the support vector machines with RBF kernel was the most accurate and the classification trees with maximum error rate was the most inaccurate algorithm. Also, entrance of geochemical and granulometric variables lead to the highest and lowest accuracy in the sediment source separation, respectively. Using the geochemical variables for the separation of sediment sources, types of support vector machines, nearest neighbor analysis, discriminant analysis and the random forest algorithm had the highest coefficients of determination and lowest error values in the training and testing stages. By entering the lithological variables, the random forest algorithm had the highest accuracy for the sediment sources classification in the training and testing

1. Assistant Professor University of torbat heydarieh - - Razavi-khorasan, University of torbat heydarieh - me.bashiri@yahoo.com

2 M.Sc. Student University of torbat heydarieh - - Razavi-khorasan, University of torbat heydarieh

3 Assistant Professor Ferdowsi University of mashhad - - Razavi-khorasan, Ferdowsi University of Mashhad

DOI: 10.22052/deej.2018.7.19.49

stages and the discriminant analysis and support vector machines were located thereafter. Finally, fitting the classification algorithms using granulometric variables showed that the support vector machines had highest accuracy in the training and testing stages of models and the random forest and nearest neighbor analysis were ranked thereafter.

Conclusion: Totally, due to the proper accuracy and performance of data mining classifier algorithms, application of these methods in the natural sciences is suggested especially in the large amounts of data. These algorithms are used to find patterns in large sets of data and help classify new information. Especially, the support vector machines that are supervised classifier algorithms and besides that, in the natural sciences have successful results. In the watershed management considering the time and cost, sediment source ascriptions are difficult to obtain using monitoring techniques, but data mining procedures, have emerged as a potentially valuable alternative. Therefore, application and evaluation of these methods are suggested for further studies and natural sciences data.

Keywords: Classification algorithms, Element density, Nodeh watershed, Sediment source ascription.