

Research Paper

مقاله پژوهشی

Investigating the Performance of the Combined Dagging Method with the Hoeffding Tree Base Algorithm in the Qualitative Classification of Drinking Water

بررسی عملکرد روش ترکیبی دگینگ با الگوریتم پایه درخت هوفدینگ در طبقه‌بندی کیفی آب شرب

Mohammad Taghi Sattari^{1*} and Sahar Javidan²

محمدتقی ستاری^{۱*} و سحر جاویدان^۲

1- Associate Professor, Department of Water Engineering, Faculty of Agriculture, University of Tabriz, Tabriz, Iran.

۱- دانشیار گروه علوم و مهندسی آب، دانشکده کشاورزی، دانشگاه تبریز، تبریز، ایران.

2- M.Sc. Student, Department of Water Engineering, Faculty of Agriculture, University of Tabriz, Tabriz, Iran.

۲- دانشجوی کارشناسی ارشد، گروه علوم و مهندسی آب، دانشکده کشاورزی، دانشگاه تبریز، تبریز، ایران.

* Corresponding Author, Email: mtsattar@tabrizu.ac.ir

* نویسنده مسئول، ایمیل: mtsattar@tabrizu.ac.ir

Received: 11/09/2022

تاریخ دریافت: ۱۴۰۱/۰۶/۲۰

Revised: 11/04/2023

تاریخ اصلاح: ۱۴۰۲/۰۱/۲۲

Accepted: 26/04/2023

تاریخ پذیرش: ۱۴۰۲/۰۲/۰۶

© IWWA

© انجمن آب و فاضلاب ایران

Abstract

چکیده

For the effective qualitative management of drinking water, it is necessary to estimate the level of water pollution. In this research, to calculate the quality index of drinking water from the chemical parameters of Total Hardness, Alkalinity, Electrical Conductivity, Total Dissolved Solids, Calcium, Sodium, Magnesium, Potassium, Chlorine, Carbonate, Bicarbonate, and Sulfate in the hydrometric station of Bagh Kelayeh, Qazvin province used in the statistical period of 23 years (1998-2020). According to the calculated numerical values and existing standards, water quality classified into two classes, good and excellent. To predict the quality class of drinking water based on chemical parameters, different combinations of parameters were considered in the form of several scenarios. In this regard, correlation and relief algorithms were used to select different scenarios. Hoeffding tree was used as a basic model for classifying water quality based on different combinations of parameters. Also, the performance of the combined Dagging approach in improving the results was evaluated. The results showed that the combined Dagging improves the water quality classification results. Scenario 6 Dagging with Hoeffding tree base algorithm, including HCO_3^- , Ca, SO_3 , TDS, EC and TH parameters, with $\text{Kappa} = 1$, was introduced as the best method which is able to classify test samples correctly.

برای مدیریت مؤثر کیفی آب شرب، برآورد سطح آلودگی آب‌های سطحی ضروری است. در پژوهش حاضر، برای محاسبه شاخص کیفی آب شرب از پارامترهای شیمیایی سختی کل، قلیائیت، هدایت الکتریکی، کل مواد جامد محلول، کلسیم، سدیم، منیزیم، پتاسیم، کلر، کربنات، بی‌کربنات و سولفات ایستگاه هیدرومتری باغ کلاهی استان قزوین، در دوره آماری ۲۳ ساله (۱۹۹۸-۲۰۲۰) استفاده شد. با توجه به مقادیر عددی محاسبه شده و استانداردهای موجود، کیفیت آب در دو کلاس خوب و عالی طبقه‌بندی شد. برای طبقه‌بندی کلاس کیفی آب شرب براساس پارامترهای شیمیایی، ترکیب‌های مختلفی از پارامترها در قالب چندین سناریو در نظر گرفته شد. در این راستا، برای انتخاب سناریوهای مختلف، از دو روش همبستگی و الگوریتم رلیف استفاده شد. درخت هوفدینگ به‌عنوان مدل پایه برای طبقه‌بندی کلاس کیفی آب براساس ترکیب‌های مختلفی از پارامترهای شیمیایی به‌کار برده شد. هم‌چنین عملکرد روش ترکیبی Dagging در بهبود نتایج، مورد ارزیابی قرار گرفت. نتایج نشان داد که روش ترکیبی Dagging باعث بهبود نتایج طبقه‌بندی کلاس کیفی آب می‌شود. سناریوی ۶ روش Dagging با الگوریتم پایه درخت هوفدینگ، شامل پارامترهای HCO_3^- ، Ca، SO_3 ، TDS، EC و TH، با $\text{Kappa} = 1$ ، به‌عنوان بهترین روش معرفی شد. این روش توانست تمام نمونه‌های آزمایشی را به‌صورت صحیح، طبقه‌بندی کند.

Keywords: Drinking Water Quality Index, Hoeffding Tree, Kappa Statistic, Rock Curve.

کلمات کلیدی: آماره کاپا، روش ترکیبی Dagging، درخت هوفدینگ، شاخص کیفی آب شرب، منحنی راک.

جذب سدیم^۵ و سدیم قادر است با دقت زیادی کیفیت آب را طبقه‌بندی کند. دزفولی و همکاران (۱۳۹۶) طبقه‌بندی کیفی آب رودخانه کارون را براساس حداقل پارامترهای کیفی انجام دادند. نتایج مطالعه آن‌ها نشان داد که روش شبکه عصبی احتمالی^۶ با استفاده از پارامترهای کیفی کدورت، کلی‌فرم مدفوعی و کل مواد جامد با دقت ۹۰/۷۸ می‌تواند به طبقه‌بندی کیفی آب بپردازد. (Gakii and Jepkoech, 2019) با استفاده از مدل درخت تصمیم کیفیت آب در کنیا را طبقه‌بندی و تجزیه و تحلیل کردند. آن‌ها درخت تصمیم J48 و Decision Stump را به ترتیب به‌عنوان مدلی با بیشترین و کمترین دقت معرفی کردند. آن‌ها دریافتند که تجزیه و تحلیل قلیائیت آب، سطح pH و هدایت الکتریکی می‌تواند نقش مهمی در ارزیابی کیفیت آب داشته باشد.

(Meddouri et al., 2021) با استفاده از مجموعه داده‌های شناخته شده از مخزن یادگیری ماشین UCI عملکرد الگوریتم دگینگ را ارزیابی کردند. برای هر مجموعه داده، به ترتیب تعداد نمونه‌ها، تعداد ویژگی‌های عددی (قبل از گسسته‌سازی)، تعداد ویژگی‌های اسمی (پس از گسسته‌سازی)، تعداد کلاس‌ها و تنوع داده‌ها بررسی شد. نتایج نشان داد به‌علت این‌که روش‌های دگینگ از اکثریت رأی برای ترکیب طبقه‌بندی‌کننده استفاده می‌کنند، بنابراین از دقت بیشتری برخوردار هستند. ایشان پیشنهاد نمودند که استفاده از الگوریتم دگینگ می‌تواند برای بهبود عملکرد طبقه‌بندی مجموعه داده مفید باشد. (Islam Khan et al., 2021) براساس رگرسیون مؤلفه اصلی^۷ و رویکرد طبقه‌بندی‌کننده تقویت گرادیان^۸، کیفیت آب را طبقه‌بندی و پیش‌بینی نمودند. آن‌ها رویکرد طبقه‌بندی‌کننده تقویت گرادیان را با دقت طبقه‌بندی ۱۰۰ درصد، به‌عنوان روش برتر معرفی کردند. همچنین دریافتند که این روش، در مقایسه با مدل‌های پیشرفته، عملکرد مطلوبی داشته است.

(Sattari et al., 2021) کیفیت آب رودخانه آلاداغ در ترکیه را به‌منظور مصارف آبیاری و کشاورزی طبقه‌بندی کردند. پس از ارزیابی روش‌های داده‌کاوی، نتیجه گرفتند که این روش‌ها، دقت و عملکرد خوبی در طبقه‌بندی کیفیت آب داشتند. همچنین، به‌طور کلی دریافتند که از میان هسته‌های مورد استفاده برای طبقه‌بندی در روش بردار پشتیبان، هسته جهانی پیرسون و از درخت‌های تصمیم، درخت REP^۹ بهترین نتایج را برای طبقه‌بندی کیفیت آب رودخانه آلاداغ داشت. آن‌ها روش درخت هوفدینگ را هم به کار بردند که نتایج مطلوبی را به‌دنبال نداشت. (Yusri et al., 2022) طبقه‌بندی کیفیت آب را با استفاده از روش ماشین بردار پشتیبان و افزایش گرادیان شدید^{۱۰} (XGBoost) انجام دادند.

با توجه به این‌که نیازهای کشاورزی، صنعتی و خانگی افزایش یافته، در نتیجه رقابت فزاینده‌ای برای تأمین آب پاک نیز به‌وجود آمده است. سیستم‌های آب سطحی به‌دلیل افزایش فشارهای ناشی از گسترش شهری و تغییر کاربری اراضی شهری با اختلالات قابل‌توجهی از طریق احیا، تغییر و آلودگی مواجه هستند. به‌هر طریقی، بسیاری از مشکلات مربوط به آب (مانند سیل یا خشکسالی و آلودگی جدی آب) نتیجه توسعه نامنظم یا نادرست کاربری زمین است (Kavita and Jagdish, 2012). در دو دهه اخیر، آلودگی رودخانه‌های آب به‌عنوان یکی از دغدغه‌های جهانی جوامع، مطرح شده است که توجه کامل محققان محیط‌زیست را می‌طلبد (Kotti et al., 2005). مطالعات بسیاری در زمینه برآورد کیفیت آب سطحی انجام شده است. بگینگ، تقویت و دگینگ، روش‌های نمونه‌گیری مجدد شناخته شده‌ای هستند که با استفاده از الگوریتم یادگیری یکسان برای طبقه‌بندی‌کننده‌های پایه، تنوعی از طبقه‌بندی‌کننده‌ها را مخصوصاً در شرایطی که داده‌ها دارای نویز هستند، تولید و ترکیب می‌کنند (Kotsianti and Kanellopoulos, 2007).

(Kotsianti and Kanellopoulos, 2007) بر روی یک سری از داده‌های معیار استاندارد با استفاده از روش رأی‌گیری از گروه‌های بگینگ، تقویت و دگینگ با ۸ طبقه‌بندی فرعی مدل را پیاده نمودند. در مقایسه‌ای که با گروه‌های بسته‌بندی، تقویت‌کننده و دگینگ ساده با ۲۵ طبقه‌بندی فرعی و همچنین سایر روش‌های ترکیبی معروف، بر روی مجموعه داده‌های معیار استاندارد انجام دادند، به این نتیجه رسیدند که الگوریتم پیشنهادی در اکثر موارد از دقت بیشتری برخوردار بوده است. (Babar and Babar, 2017) شاخص کیفیت آب رودخانه یامونا را با استفاده از تکنیک‌های داده‌کاوی^۱، پیش‌بینی کردند. نتایج مطالعه آن‌ها نشان داد که طبقه‌بندی‌کننده‌های درخت تصمیم^۲ و ماشین‌های بردار پشتیبان^۳، بهترین مدل‌های پیش‌گوینده با نرخ خطای ۰/۰٪، در طول استخراج خودکار کلاس کیفیت آب هستند. در صورتی‌که داده‌های ارائه‌شده نمایشی واقعی از دانش حوزه باشند، تکنیک‌های داده‌کاوی پتانسیل پیش‌بینی سریع کلاس کیفیت آب را خواهند داشت.

ستاری و همکاران (۱۳۹۶) از روش‌های داده‌کاوی برای پیش‌بینی کیفیت آب‌های سطحی رودخانه‌های دامنه شمالی سه‌سند استفاده کردند. آن‌ها نتیجه گرفتند که مدل تصمیم‌گیری درختی با استفاده از چهار پارامتر هدایت الکتریکی^۴، PH، نسبت

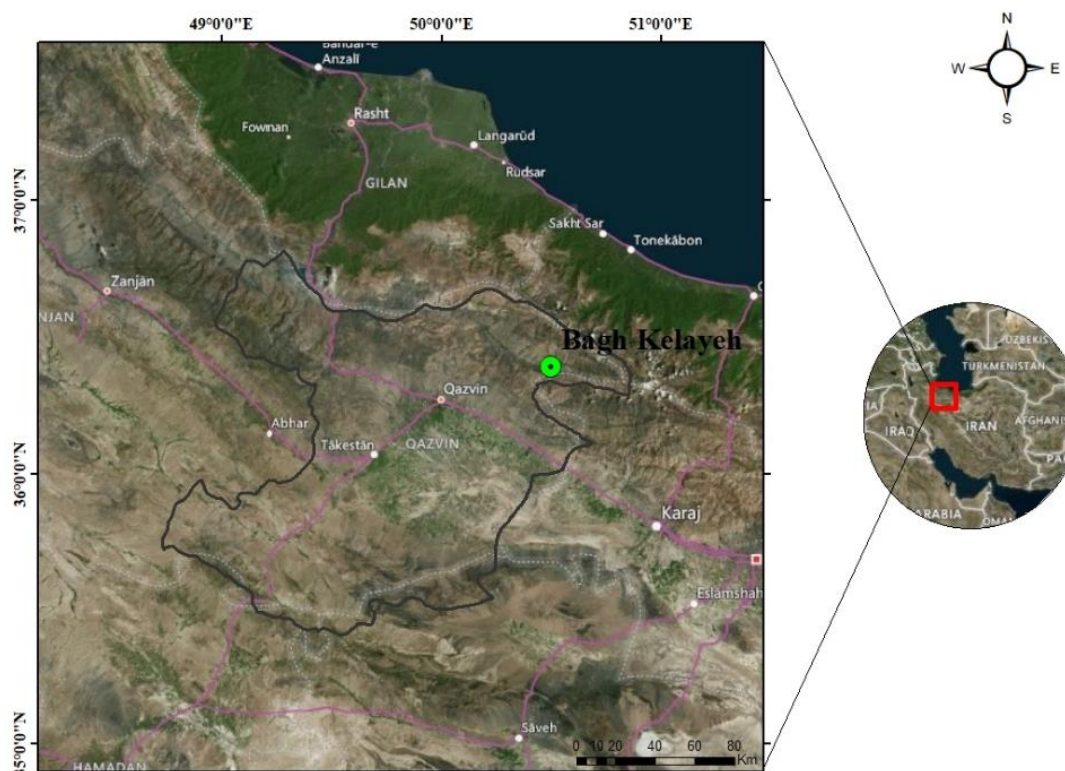
۲- مواد و روش‌ها

۲-۱- منطقه مورد مطالعه و داده‌های مورد استفاده

باغ کلايه، روستایی از توابع بخش رودبار الموت شهرستان قزوین در استان قزوین است. ایستگاه هیدرومتری باغ کلايه دارای عرض جغرافیایی "۳۸' ۲۳" ۳۶°، طول جغرافیایی "۵۱' ۲۹" ۵۰° بوده و در ارتفاع ۱۲۸۷ متر از سطح دریا واقع شده است. استان قزوین در بخش شمال غربی کشور ایران واقع شده و مساحت آن حدود ۱۵۸۲۰۰ km² است. میانگین بارش سالانه در ۲۰ سال اخیر ۳۰۲/۵۱ mm و متوسط دمای سالانه ۱۵/۲۱ °C بوده که براساس طبقه‌بندی دومارتن دارای اقلیم نیمه‌خشک و براساس اقلیم‌نمای آمبروزه دارای اقلیم خشک سرد است (جاویدان و همکاران، ۱۴۰۱). شکل ۱ موقعیت مکانی ایستگاه مورد مطالعه را نشان می‌دهد.

نتایج تحقیق آن‌ها نشان داد که مدل XGBoost با دقت ۹۴ درصد، بهتر از مدل SVM عمل کرده و تنها ۶ درصد خطا در طبقه‌بندی کیفیت آب داشته است.

با توجه به هزینه‌های سرسام‌آور مینیورینگ و نمونه‌برداری از رودخانه‌ها برای تعیین کیفیت آب شرب، ارائه مدل‌هایی که بتواند با حداقل تعداد پارامتر هیدروشیمیایی با دقت نسبتاً قابل قبولی کلاس کیفی آب را مشخص کند، یکی از اولویت‌های مدیران کیفی آب است. در این راستا هدف پژوهش حاضر، طبقه‌بندی کلاس کیفی آب شرب، با استفاده از مدل‌های داده‌کاوی براساس داده‌های مشاهداتی مربوط به پارامترهای کیفی آب ایستگاه هیدرومتری باغ کلايه در استان قزوین است. با توجه به این‌که تاکنون از روش دگینگ در زمینه مدل‌سازی کیفی آب استفاده نشده است، بررسی تأثیر استفاده از این روش در بهبود نتایج، یکی از نوآوری‌های پژوهش حاضر است.



شکل ۱- موقعیت مکانی ایستگاه هیدرومتری باغ کلايه

شد. کلاس کیفی به‌دست آمده از این روش، به‌عنوان خروجی‌های هدف در فرآیند مدل‌سازی مورد استفاده قرار گرفتند. دوره آماری در مطالعه حاضر، ۲۳ ساله (۱۹۹۸-۲۰۲۰) در نظر گرفته شد. مشخصات آماری متغیرهای مورد استفاده در جدول ۱ ارائه شده است.

در پژوهش حاضر برای محاسبه شاخص کیفی آب^{۱۱} (WQI) و طبقه‌بندی کلاس کیفی آب شرب ایستگاه هیدرومتری باغ کلايه، از پارامترهای کیفی شامل سختی کل (TH)، قلیائیت (pH)، هدایت الکتریکی (EC)، کل مواد جامد محلول^{۱۲} (TDS)، کلسیم (Ca)، سدیم (Na)، منیزیم (Mg)، پتاسیم (K)، کلر (Cl)، کربنات (CO₃)، بی‌کربنات (HCO₃) و سولفات (SO₄) استفاده

جدول ۱- مشخصات آماری پارامترهای مورد استفاده

Statistic	حداقل	حداکثر	میانگین	واریانس
سختی کل (میلی گرم بر لیتر)	۹۵/۰۰	۴۸۱/۵۰	۲۷۸/۱۵	۳۴۰۷/۳۸
قلیابیت	۴/۵۰	۸/۴۰	۷/۸۳	۰/۱۱
هدایت الکتریکی (میکروموس بر سانتی‌متر)	۲۷۹/۰۰	۱۰۴۸/۰۰	۶۲۷/۹۴	۱۷۴۸۰/۸۷
کل مواد جامد محلول (میلی گرم بر لیتر)	۱۸۶/۰۰	۶۶۳/۰۰	۳۸۸/۰۵	۵۹۸۳/۸۱
کلسیم (میلی اکی‌والان بر لیتر)	۰/۰۰	۱۵۹/۸۰	۷۳/۹۴	۳۸۰/۳۴
سدیم (میلی اکی‌والان بر لیتر)	۰/۴۶	۶۰/۴۹	۱۷/۲۶	۸۲/۹۸
منیزیم (میلی اکی‌والان بر لیتر)	۲/۷۶	۵۸/۲۰	۲۲/۱۸	۰/۹۳۷۰
پتاسیم (میلی اکی‌والان بر لیتر)	۰/۳۹	۱۹/۵۰	۱/۹۹	۲/۷۹
کلر (میلی اکی‌والان بر لیتر)	۰/۰۰	۸۹/۶۰	۲۷/۳۶	۱۹۸/۶۱
کربنات (میلی اکی‌والان بر لیتر)	۰/۰۰	۳۳/۰۰	۰/۲۱	۴/۱۷
سولفات (میلی اکی‌والان بر لیتر)	۲۲/۰۸	۳۷۴/۸۸	۱۴۰/۲۱	۲۱۲۹/۱۹
بی کربنات (میلی اکی‌والان بر لیتر)	۵۰/۰۲	۳۹۱/۶۲	۱۶۳/۶۰	۱۵۶۵/۲۱

۲-۲- سناریوهای مورد استفاده

در این پژوهش برای انجام مدل‌سازی با روش‌های داده‌کاوی از سناریوهای مختلفی استفاده شد تا بهترین ترکیب ورودی مشخص شود و تأثیر انتخاب پارامترهای شیمیایی مختلف مورد بررسی قرارگیرد. براین اساس سناریوهای استفاده شده در پژوهش حاضر در جدول ۲ ارائه شده است. انتخاب ترکیب‌های ورودی مختلف، با استفاده از روش همبستگی و الگوریتم رلیف^{۱۳} انجام شد. استفاده از الگوریتم رلیف که برای کاهش ابعاد مسئله مورد استفاده قرار می‌گیرد؛ توسط (Kira and Rendell (1992) پیشنهاد شده است. مزیت‌های استفاده از این روش، ساده بودن اصول و عدم پیچیدگی آن، قابل‌حل بودن با توابع چندجمله‌ای مرتبه پایین، قابل استفاده بودن برای داده‌های پیوسته و نیاز به تعداد کم داده‌های آموزشی در نظر گرفته می‌شود. در یک مجموعه داده با تعداد N نمونه (داده مشاهده‌ای) و تعداد P ویژگی که مربوط به دو طبقه مختلف هستند، هر ویژگی باید در بازه (۰, ۱) قرارگیرد. الگوریتم مذکور، m بار تکرار شده و در هر مرتبه از یک بردار وزنی متفاوت که از صفر شروع می‌شود، استفاده می‌کند. در هر تکرار، الگوریتم مذکور بردار ویژگی X را که متعلق به یک نمونه تصادفی است و بردارهای ویژگی نزدیک‌ترین نمونه به نمونه X در طبقه موردنظر را توسط تابع فاصله اقلیدسی انتخاب می‌کند. پس از m تکرار، هر یک از عناصر بردار وزن توسط m تقسیم‌بندی می‌شوند. نتیجه این عمل این است که یک بردار مرتبط به دست می‌آید. چنان‌چه مقدار بردار مرتبط یک ویژگی، از آستانه تعریف شده بیشتر شود، آن ویژگی انتخاب می‌شود. این الگوریتم نه به‌عنوان یک روش پیش‌پردازش است نه روش طبقه‌بندی.

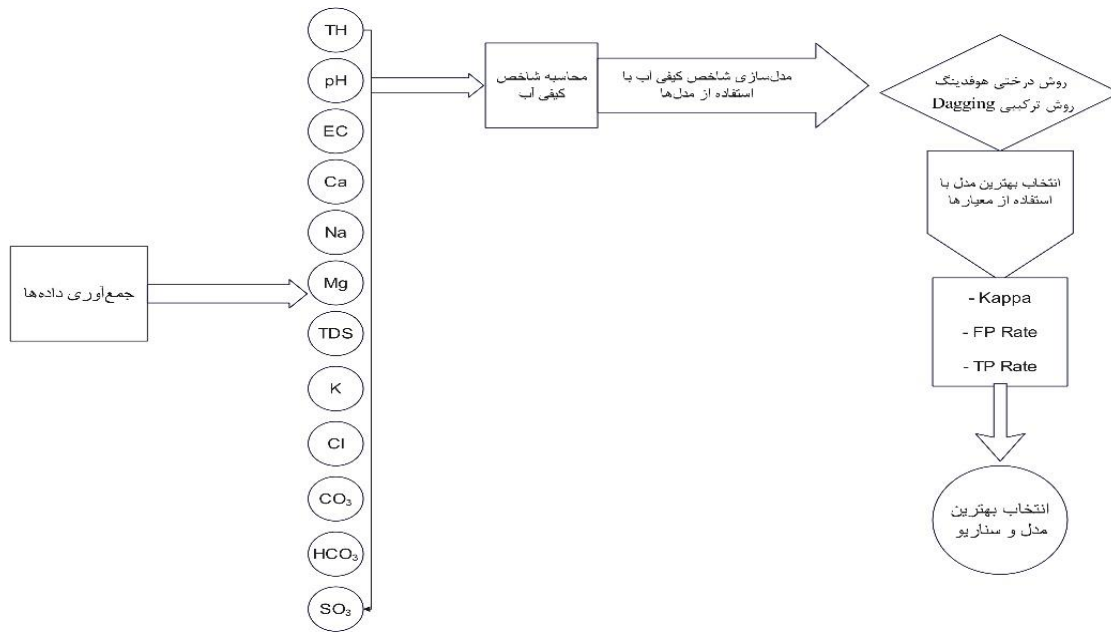
الگوریتم رلیف به‌منظور مشخص نمودن تأثیر متغیرهای مستقل بر روی متغیر پاسخ و یا هدف به‌کار می‌رود. الگوریتم

رلیف ترکیب‌های متفاوتی از متغیرهای ورودی در راستای ایجاد سناریوهای مختلف مدل‌سازی را ارائه می‌کند (Kira and Rendell, 1992). انتخاب ویژگی بر اساس همبستگی، یک روش معمول و پرکاربرد برای انتخاب متغیرهای ورودی و کاهش ابعاد مسئله است. روش همبستگی به زیرمجموعه‌هایی که دارای ویژگی‌هایی با بیشترین ضریب همبستگی با کلاس نمونه مورد نظر هستند، امتیاز می‌دهد و متغیرهایی که بیشترین امتیاز را دارا باشند، به‌عنوان متغیر اصلی در نظر می‌گیرد. این الگوریتم توانایی بالایی در تشخیص سریع داده‌های نامربوط، اضافی و دارای خطا دارد که عموماً منجر به حذف نیمی از داده‌ها می‌شود. این ویژگی با کاهش ابعاد مسئله سبب افزایش بهره‌وری مدل‌ها می‌شود (Hall, 1999). از بین داده‌های مورد مطالعه، ۷۰٪ برای واسنجی و ۳۰٪ برای صحت‌سنجی انتخاب شدند. برای طبقه‌بندی کلاس کیفی آب شرب، از روش هوفدینگ^۹ استفاده شد و عملکرد روش Dagging با الگوریتم پایه درخت هوفدینگ مورد بررسی قرار گرفت.

۲-۳- روش‌های مورد مطالعه

در این پژوهش برای تعیین کیفیت آب از دیدگاه شرب، ابتدا شاخص کیفی آب به‌عنوان یک شاخص عددی با استفاده از پارامترهای شیمیایی مؤثر در کیفیت آب محاسبه شد. سپس با استفاده از استانداردهای موجود کلاس و طبقه کیفی آب از نقطه‌نظر مساعد بودن و یا نبودن برای شرب تعیین شد. از آن‌جایی که نمونه‌برداری از منبع آب برای تعیین کیفیت آب زمان‌بر و هزینه‌بر است، لذا به‌کمک روش‌های داده‌کاوی سعی شد با استفاده از حداقل تعداد پارامتر شیمیایی شاخص کیفی آب و به‌تبع آن کلاس کیفی آب مدل‌سازی شود. در همین راستا

پژوهش معطوف به افزایش دقت مدل‌ها در پیش‌بینی شاخص کیفی آب بود. با استفاده از معیارهای ارزیابی نتایج هریک از مدل‌ها و الگوریتم‌ها برای سناریوهای معرفی شده مورد مقایسه قرار گرفت و نهایتاً بهترین سناریو (ترکیب پارامترهای ورودی) و دقیق‌ترین مدل و الگوریتم داده‌کاوی انتخاب شد. در شکل ۲ روند انجام تحقیق ارائه شده است.



شکل ۲- روندنمای مراحل انجام تحقیق

جدول ۲- طبقه بندی کیفیت آب بر اساس ارزش WQI

طبقه بندی کیفیت آب آشامیدنی		
محدوده شاخص WQI	کلاس	نوع آب
کمتر از ۵۰	I	آب عالی
۵۰-۱۰۰	II	آب خوب
۱۰۰-۲۰۰	III	آب بد
۲۰۰-۳۰۰	IV	آب بسیار بد
بیشتر از ۳۰۰	V	آب نامناسب برای آشامیدن

۲-۳-۲ روش درختی هوفدینگ

مدل درختی هوفدینگ^{۱۴} که به‌عنوان یک الگوریتم افزایشی معروف درخت تصمیم شناخته می‌شود، در ابتدا برای رفع مشکلات طبقه‌بندی در استخراج جریان داده در مقیاس بزرگ پیشنهاد شد (Domingos and Hulten, 2003). در روش مدل HT، هر نمونه در پایگاه داده آموزشی تنها یک‌بار اسکن می‌شود. از این‌رو، این مدل دارای راندمان محاسباتی برجسته با رم نسبتاً پایین است. علاوه بر این، عملیات طبقه‌بندی را می‌توان در زمانی که مدل HT در حال رشد است، اجرا کرد که به‌وضوح با مدل‌های DT معمولی متفاوت است (Mehta and Sanghavi, 2019).

۲-۳-۱ شاخص کیفیت آب شرب

از فرمول‌های (۱) تا (۳) برای محاسبه شاخص کیفی آب و طبقه‌بندی کلاس کیفی، استفاده شد (Singh, 1992).

$$W_i = \frac{w_i}{\sum_{i=1}^n w_i} \quad (1)$$

$$q_i = \left(\frac{C_i}{S_i}\right) \times 100 \quad (2)$$

$$WQI = \sum_{i=1}^n W_i q_i \quad (3)$$

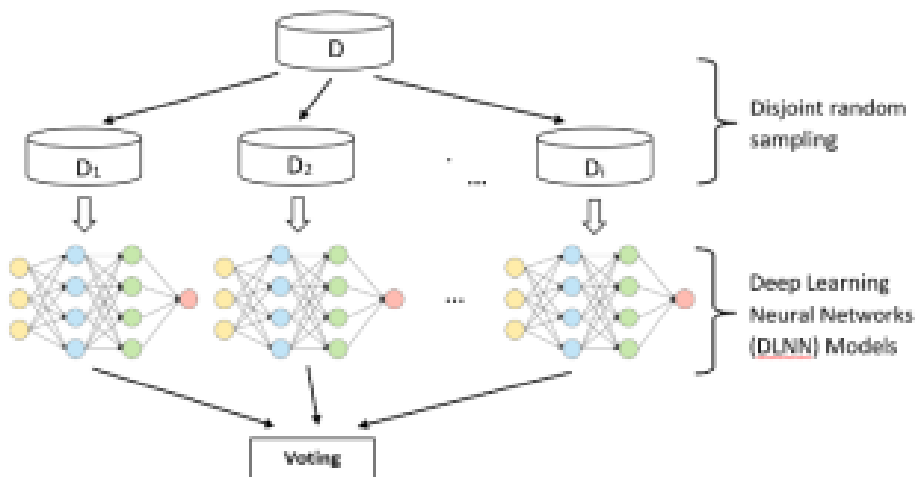
که w : وزن مربوط به هر پارامتر با توجه به اهمیت آن در شرب، W : وزن نسبی هر پارامتر، C : غلظت هر پارامتر، S : غلظت استاندارد هر پارامتر، q : رتبه کیفی هر پارامتر و WQI : نیز شاخص کیفی آب شرب است.

پس از محاسبه مقادیر WQI ، کلاس کیفی آب شرب به پنج دسته آب عالی، خوب، بد، بسیار ضعیف و نامناسب برای آشامیدن تقسیم می‌شود (جدول ۲).

۲-۳-۳- روش Dagging

این طبقه‌بندی‌کننده، تعدادی چین‌های منفصل و طبقه‌بندی‌شده از داده‌ها ایجاد می‌کند و هر تکه از داده‌ها را به یک کپی از طبقه‌بندی‌کننده پایه ارائه‌شده، تغذیه می‌کند. پیش‌بینی‌ها از طریق اکثریت آرا انجام می‌شود. نمونه‌گیری غیرمجاز تصادفی و ناهمگن در مجموعه داده آموزشی اعمال

می‌شود. پس از نمونه‌برداری مجدد از مجموعه داده آموزشی، یک طبقه‌بندی پایه برای هر نمونه ایجاد می‌شود. سپس نتایج این طبقه‌بندی‌کننده‌های چندگانه با استفاده از رأی اکثریت ترکیب می‌شوند. بنابراین روش ترکیبی Dagging یک مجموعه موازی در نظر گرفته می‌شود (Ting and Witten, 1997). شکل ۳ نمای کلی از این روش را ارائه می‌دهد.



شکل ۳- نمای کلی روش ترکیبی Dagging (Elish and Elish, 2021)

۲-۴- معیارهای ارزیابی

در هر پژوهش مبتنی بر روش‌های داده‌کاوی، از آنجایی که روش‌ها در رقابت بایکدیگر عمل می‌کنند، ضروری است تا با به کار بردن معیارهایی دقت مدل‌های مورد استفاده را تعیین نموده و بر این اساس نسبت به اولویت‌بندی مدل‌ها و انتخاب مدل برتر اقدام نمود. اساساً در شرایطی که خروجی مدل‌ها کمی (عددی) و یا کیفی (غیرعددی) باشند، معیارهای ارزیابی متفاوتی برای سنجش دقت مدل‌ها به کار می‌رود. در این پژوهش خروجی مدل‌های داده‌کاوی کیفی (غیرعددی) هستند و بر این اساس معیارهای ارزیابی، معطوف به مقایسه تعداد صحیح موارد پیش‌بینی شده در برابر تعداد غیرصحیح موارد پیش‌بینی شده در فرآیند مدل‌سازی و طبقه‌بندی است.

برای مقایسه کلاس کیفی به‌دست آمده از روش‌های داده‌کاوی با کلاس کیفی محاسبه شده از WQI از آماره کاپا، نرخ مثبت صحیح^{۱۵} (TPR) و نرخ مثبت غلط^{۱۶} (FPR) استفاده شد. فرمول‌های آماره‌های فوق به ترتیب در روابط (۴) تا (۶) ارائه شده است.

$$\text{kappa} = p_i = (PA_0 - PA_E) / (1 - PA_E) \quad (4)$$

$$\text{Sensitivity} = TPR = \frac{TP}{TP + FN} \quad (5)$$

$$1 - \text{Specificity} = FPR = \frac{TN}{TN + FP} \quad (6)$$

که PA_0 : میزان توافق دو ارزیاب، PA_E : میزان توافق مورد انتظار، TP : مثبت صحیح، FN : منفی غلط، FP : مثبت غلط و TN : منفی صحیح هستند.

منحنی راک برای مقایسه دو ویژگی عملیاتی نرخ مثبت واقعی (TPR) یا حساسیت و نرخ مثبت کاذب (FPR) یا ویژگی به کار می‌رود. هم‌چنین به‌عنوان منحنی مشخصه عملکرد گیرنده شناخته می‌شود. منحنی مشخصه عملکرد گیرنده، یک نمایش گرافیکی است که عملکرد الگوریتم‌های طبقه‌بندی‌کننده داده‌کاوی را تفسیر می‌کند. این نتیجه ترسیم نرخ مثبت واقعی در مقابل نرخ مثبت کاذب در تنظیمات آستانه‌های مختلف است. TPR در امتداد محور y و FPR در امتداد محور x رسم می‌شود. عملکرد هر الگوریتم به‌عنوان یک نقطه در منحنی ROC نشان داده می‌شود.

۳- نتایج و بحث

در پژوهش حاضر برای برآورد کلاس کیفی آب با استفاده از شاخص WQI، از پارامترهای TH، pH، EC، TDS، Ca، Na، Mg، SO_4 ، HCO_3 ، CO_3 ، Cl، K، دوره آماری ۲۳ ساله استفاده شد. روش رلیف و همبستگی برای انتخاب ترکیب‌های ورودی مختلف به کار برده شد. پارامترهای ورودی سناریوهای مورد مطالعه در جدول ۳ ارائه شده است.

جدول ۲- پارامترهای دخیل در هر سناریو و روش انتخاب سناریوها

روش انتخاب سناریو	ورودی‌ها	شماره سناریو
ماتریس همبستگی	TH	۱
	TH, EC	۲
	TH, EC, TDS	۳
	TH, EC, TDS, SO ₄	۴
	TH, EC, TDS, SO ₄ , Ca	۵
	TH, EC, TDS, SO ₄ , Ca, HCO ₃	۶
	TH, EC, TDS, SO ₄ , Ca, HCO ₃ , Mg	۷
	TH, EC, TDS, SO ₄ , Ca, HCO ₃ , Mg, Cl	۸
	TH, EC, TDS, SO ₄ , Ca, HCO ₃ , Mg, Cl, Na	۹
	TH, EC, TDS, SO ₄ , Ca, HCO ₃ , Mg, Cl, Na, K	۱۰
	TH, EC, TDS, SO ₄ , Ca, HCO ₃ , Mg, Cl, Na, K, CO ₃	۱۱
	TH, EC, TDS, SO ₄ , Ca, HCO ₃ , Mg, Cl, Na, K, CO ₃ , PH	۱۲
الگوریتم رلیف	TH, K	۱۳
	TH, K, SO ₄	۱۴
	TH, K, SO ₄ , TDS	۱۵
	TH, K, SO ₄ , TDS, EC	۱۶

نمونه‌های برداشت شده در کلاس کیفی یک و دو، یعنی کیفیت عالی (۱۳/۲۷ درصد) و خوب (۸۶/۷۳ درصد) قرار گرفتند. نتایج مدل‌سازی کلاس کیفی آب در جدول ۴ ارائه شده است.

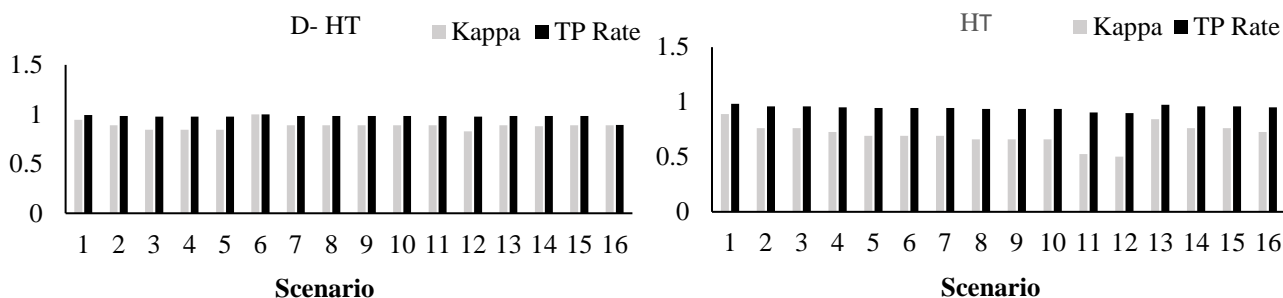
برای برآورد مقادیر کیفی WQI از روش درختی هوفدینگ و روش ترکیبی Dagging استفاده شد. پس از طبقه‌بندی کلاس کیفی آب در ایستگاه مورد مطالعه، مشخص شد که تمام

جدول ۴- معیارهای ارزیابی برای برآورد مقادیر کیفی WQI

روش										سناریو
روش ترکیبی Dagging					روش درختی هوفدینگ					
ROC Area	FP Rate	TP Rate	RMSE	Kappa	ROC Area	FP Rate	TP Rate	RMSE	Kappa	
۱/۰۰۰	۰/۰۰۱	۰/۹۹۲	۰/۱۰۰	۰/۹۴۳	۱/۰۰۰	۰/۰۰۱	۰/۹۸۴	۰/۱۱۷	۰/۸۹۱	۱
۱/۰۰۰	۰/۰۰۱	۰/۹۸۴	۰/۱۰۴	۰/۸۹۱	۱/۰۰۰	۰/۰۰۳	۰/۹۶۱	۰/۱۶۷	۰/۷۶۲	۲
۱/۰۰۰	۰/۰۰۲	۰/۹۷۶	۰/۱۰۹	۰/۸۴۵	۱/۰۰۰	۰/۰۰۳	۰/۹۶۱	۰/۱۸۶	۰/۷۶۲	۳
۱/۰۰۰	۰/۰۰۲	۰/۹۷۶	۰/۱۱۵	۰/۸۴۵	۱/۰۰۰	۰/۰۰۴	۰/۹۵۳	۰/۱۹۸	۰/۷۲۶	۴
۱/۰۰۰	۰/۰۰۲	۰/۹۷۶	۰/۱۲۲	۰/۸۴۵	۱/۰۰۰	۰/۰۰۴	۰/۹۴۵	۰/۲۱۰	۰/۶۹۲	۵
۱/۰۰۰	۰/۰۰۰	۱/۰۰۰	۰/۰۹۳	۱/۰۰۰	۱/۰۰۰	۰/۰۰۴	۰/۹۴۵	۰/۲۱۱	۰/۶۹۲	۶
۱/۰۰۰	۰/۰۰۱	۰/۹۸۴	۰/۱۰۱	۰/۸۹۱	۱/۰۰۰	۰/۰۰۴	۰/۹۴۵	۰/۲۰۱	۰/۶۹۲	۷
۰/۹۹۷	۰/۰۰۱	۰/۹۸۴	۰/۱۳۲	۰/۸۹۱	۰/۹۹۸	۰/۰۰۵	۰/۹۳۷	۰/۲۲۷	۰/۶۶۱	۸
۰/۹۹۵	۰/۰۰۱	۰/۹۸۴	۰/۱۳۰	۰/۸۹۱	۰/۹۹۸	۰/۰۰۵	۰/۹۳۷	۰/۲۳۸	۰/۶۶۱	۹
۰/۹۹۶	۰/۰۰۱	۰/۹۸۴	۰/۱۳۶	۰/۸۹۱	۰/۹۹۸	۰/۰۰۵	۰/۹۳۷	۰/۲۴۸	۰/۶۶۱	۱۰
۰/۹۹۸	۰/۰۰۱	۰/۹۸۴	۰/۱۱۳	۰/۸۹۱	۰/۹۲۱	۰/۱۱۰	۰/۹۰۶	۰/۲۹۱	۰/۵۲۶	۱۱
۰/۹۹۵	۰/۱۰۴	۰/۹۷۶	۰/۱۳۰	۰/۸۲۹	۰/۹۲۲	۰/۱۱۰	۰/۸۹۸	۰/۲۹۸	۰/۵۰۳	۱۲
۰/۹۹۷	۰/۰۰۱	۰/۹۸۴	۰/۱۱۳	۰/۸۹۱	۰/۹۹۹	۰/۰۰۲	۰/۹۷۶	۰/۱۳۸	۰/۸۴۵	۱۳
۰/۹۹۶	۰/۱۰۴	۰/۹۸۴	۰/۱۰۵	۰/۸۸۰	۰/۹۹۷	۰/۰۰۳	۰/۹۶۱	۰/۱۵۴	۰/۷۶۲	۱۴
۱/۰۰۰	۰/۰۰۱	۰/۹۸۴	۰/۰۹۱	۰/۸۹۱	۱/۰۰۰	۰/۰۰۳	۰/۹۶۱	۰/۱۷۶	۰/۷۶۲	۱۵
۱/۰۰۰	۰/۰۰۱	۰/۸۹۴	۰/۰۹۰	۰/۸۹۱	۱/۰۰۰	۰/۰۰۴	۰/۹۵۳	۰/۱۹۸	۰/۷۲۶	۱۶

روش و ۱۶ سناریوی مورد مطالعه، سناریوی ۶ روش Dagging با الگوریتم پایه HT شامل $TH, EC, TDS, SO_4, Ca, HCO_3$ و $ROC = 1, FPR = 0, TPR = 1, RMSE = 0.093, Kappa = 1$ بهترین عملکرد را داشته است. شکل ۴ نمودار ستونی مقادیر Kappa و TPR روش‌های HT و D- HT را نشان می‌دهد.

همان‌گونه که از جدول ۴ مشخص است، استفاده از روش ترکیبی Dagging باعث بهبود نتایج درخت هوفدینگ شده است. در روش HT سناریو اول شامل TH بالاترین دقت را دارد. با بهره‌گیری از روش Dagging مقدار کاپا در سناریوی اول از ۰/۸۹۱ به ۰/۹۴۳ افزایش یافته است. در حالت کلی از بین دو



شکل ۴- نمودار ستونی مقادیر کاپا و نرخ مثبت صحیح روش‌های مورد مطالعه

نمونه‌ای در این جدول‌ها خارج از قطر اصلی قرارگیرد، به این معنی است که در کلاس اشتباهی طبقه‌بندی شده است. به عبارت دیگر، هر چه تعداد نمونه‌ها در قطر اصلی بیشتر باشد، دقت و کارایی مدل در تخمین کلاس موردنظر، بیشتر می‌شود.

برتر بودن سناریوی ۱ در روش HT و سناریوی ۶ در روش D- HT از نمودارهای شکل ۴ نیز مشهود است. جدول‌های ۵ و ۶ ماتریس اغتشاش به‌دست آمده از نتایج سناریوهای برتر هر دو روش مورد مطالعه را نشان می‌دهد. این جدول‌ها نحوه توزیع نمونه‌ها در بین کلاس‌های مختلف را نشان می‌دهد. چنان‌چه

جدول ۶- ماتریس اغتشاش سناریو برتر روش D- HT

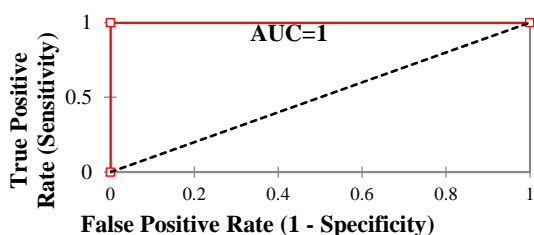
		طبقه‌بندی شده با مدل	
		خوب	عالی
واقعی	خوب	۱۱۸	۰
	عالی	۰	۹

جدول ۵- ماتریس اغتشاش سناریو برتر روش HT

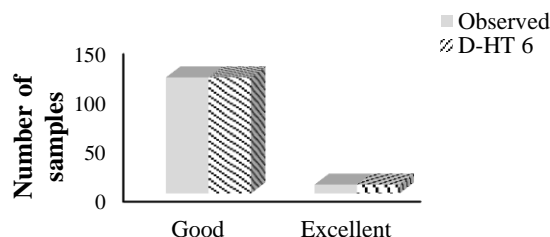
		طبقه‌بندی شده با مدل	
		خوب	عالی
واقعی	خوب	۱۱۶	۲
	عالی	۰	۹

HT6 برای دوره آزمون را نشان می‌دهد. بنابراین می‌توان نتیجه گرفت که روش Dagging عملکرد خوبی در برآورد مقادیر کیفی WQI داشته و باعث بهبود نتایج شده است. منحنی راک سناریو و روش برتر معرفی شده در پژوهش حاضر در شکل ۶ مشخص شد.

با توجه به جدول ۵، از بین ۱۲۷ نمونه آزمایشی، تنها ۲ نمونه به اشتباه طبقه‌بندی شده است. یعنی به‌جای کلاس کیفی خوب (II)، در کلاس عالی (I) قرار گرفته است. اما در روش D-HT6 که ماتریس اغتشاش آن در جدول ۶ مشخص شده است، تمام نمونه‌ها در کلاس صحیح، طبقه‌بندی شده‌اند. شکل ۵ نیز توزیع سه‌بعدی کلاس کیفی آب برای داده‌های مشاهده‌ای و مدل D-



شکل ۶- منحنی راک روش D- HT6



شکل ۵- توزیع سه‌بعدی کلاس کیفی آب برای داده‌های مشاهده‌ای و مدل D- HT6 برای دوره آزمون

۵- پی‌نوشت‌ها

- 1- Data Mining
- 2- Decision Tree
- 3- Support Vector Machine
- 4- Electrical Conductivity
- 5- Sodium Absorption Ratio
- 6- Probabilistic Neural Network
- 7- Principal Component Regression
- 8- Gradient Boosting Classifier Approach
- 9- Reduced Error Pruning Tree
- 10- Extreme Gradient Boosting
- 11- Water Quality Index
- 12- Total Dissolved Solids
- 13- Relief
- 14- Hoeffding Tree
- 15- True Positive Rate
- 16- False Positive Rate

۶- مراجع

- جاویدان، س.، ستاری، م. ت.، کریم‌زاده، پ.، و مهرابی، ا.، (۱۴۰۱)، "تحلیل عملکرد روش‌های هیدرولوژیکی و داده-مبنا در برآورد میزان رسوب معلق"، *مجله محیط‌زیست و مهندسی آب*، ۲۸(۲)، ۴۶۸-۴۸۰،
<https://doi.org/10.22034/jewe.2021.305599.1632>.
- دزفولی، د.، موغاری، م. ح.، ابراهیمی، ک.، و عراقی‌نژاد، ش.، (۱۳۹۶)، "تعیین طبقه‌بندی کیفی آب بر اساس حداقل پارامترهای کیفی (مطالعه موردی: رودخانه کارون)"، *مجله محیط‌زیست طبیعی، مجله منابع طبیعی ایران*، ۷۰(۳)، ۵۸۳-۵۹۵،
<https://doi.org/10.22059/jne.2017.213338.1224>.
- ستاری، م. ت.، میرعباسی، ر.، و عباسقلی نایب‌زاد، م.، (۱۳۹۶)، "استفاده از داده‌کاوی در پیش‌بینی کیفیت آب‌های سطحی (مطالعه موردی: رودخانه‌های دامنه شمالی سهند)"، *اکوهیدرولوژی*، ۲۴(۲)، ۴۰۷-۴۱۹،
<https://doi.org/10.22059/ije.2017.61477>.
- Babar, R., and Babar, S., (2017), "Predicting river water quality index using data mining techniques", *Environmental Earth Sciences*, 76(504), 1-15, <https://doi.org/10.1007/s12665-017-6845-9>.
- Domingos, P., and Hulten, G., (2003), "A general framework for mining massive data streams", *Journal of Computational and Graphical Statistics*, 12(4), 945-949, <https://doi.org/10.1198/1061860032544>.
- Elish, M., and Elish, K., (2021), "An empirical comparison of resampling ensemble methods of deep learning Neural Networks for cross-project software defect prediction", *International Journal of Intelligent Engineering and Systems*, 14(3), 201-209, <https://doi.org/10.22266/ijies2021.0630.18>.

هنگامی که حساسیت افزایش پیدا کند، نرخ مثبت کاذب نیز افزایش خواهد یافت. بنابراین در منحنی راک امکان مقایسه نرخ مثبت صحیح و نرخ مثبت کاذب، در هر نقطه بر روی منحنی وجود دارد. سطح زیر منحنی نشان‌دهنده کیفیت مدل‌های مورد مطالعه است. هرچه نرخ مثبت صحیح یا حساسیت بیشتر باشد، بدین معنی است که نمونه‌ها به‌درستی طبقه‌بندی شده‌اند. در نتیجه سطح زیر منحنی بیشتر شده و به ۱۰۰ درصد خواهد رسید. با توجه به توضیحات فوق و شکل ۶، چون سطح زیر منحنی برابر با یک شده است، می‌توان نتیجه گرفت که روش مورد استفاده دقت بالایی داشته است.

۴- نتیجه‌گیری

در پژوهش حاضر، کلاس کیفی آب ایستگاه هیدرومتری باغ کلابه مورد بررسی قرار گرفت. ابتدا مقادیر عددی کیفی شاخص WQI آب براساس داده‌های مشاهداتی شده پارامترهای شیمیایی محاسبه و براساس استانداردهای موجود، کلاس کیفیت آب شرب مشخص شد. برای طبقه‌بندی کیفی آب از ۱۲ پارامتر شیمیایی شامل pH، EC، TDS، Ca، Na، Mg، K، Cl، CO₃، HCO₃ و SO₄ در دوره آماری ۲۳ ساله در سناریوهای مختلف استفاده شد. طبقه‌بندی کلاس کیفی آب با روش هوفدینگ انجام شد. در مدل‌سازی کیفی، تاثیر روش ترکیبی Dagging با الگوریتم پایه درخت هوفدینگ در بهبود نتایج، مورد ارزیابی قرار گرفت. برای انتخاب روش و سناریو برتر از آماره کاپا و منحنی راک استفاده شد. در مدل‌سازی کیفی بهره‌گیری از روش Dagging باعث افزایش دقت و کاهش خطا شد. در طبقه‌بندی کلاس کیفی آب شرب، سناریوی ۶ روش D-HT شامل پارامترهای TH، EC، SO₄، Ca، HCO₃ و TDS از دقت قابل توجهی برخوردار بود و توانست تمام نمونه‌ها را در کلاس کیفی صحیح، طبقه‌بندی کند. بنابراین این روش توانست به‌جای استفاده از ۱۲ پارامتر شیمیایی، تنها با استفاده از ۶ پارامتر، با دقت مناسبی کلاس کیفی تمام نمونه‌های ایستگاه مورد مطالعه را به‌درستی، طبقه‌بندی کند. با توجه به این‌که تمام روش‌های مورد استفاده در طبقه‌بندی کلاس کیفی، دقت قابل قبولی داشتند، لذا در صورت کمبود داده و عدم دسترسی به تمام پارامترهای شیمیایی، می‌توان با استفاده از تعداد محدودی از پارامترها و با بهره‌گیری از روش‌های داده‌کاوی، نتایج مناسب و قابل قبولی را به‌دست آورده و کلاس کیفیت آب شرب را طبقه‌بندی نمود. نتایج به‌دست آمده نشان داد؛ آب منطقه مورد مطالعه از کیفیت قابل قبولی برای شرب برخوردار است.

- Ting, K.M., and Witten, I.H., (1997), "Stacking bagged and dagged models", In: *Fourteenth International Conference on Machine Learning*, San Francisco, CA, pp. 367-375.
- Yusri, H., Ab Rahim, A., Hassan, S., Halim, I., and Abdullah, N., (2022), "Water quality classification using SVM and XGBoost method", *IEEE 13th Control and System Graduate Research Colloquium (ICSGRC)*, pp. 231-236, <https://doi.org/10.1109/ICSGRC55096.2022.9845143>.
- Gakii, C., and Jepkoech, J., (2019), "A classification model for water quality analysis using desision tree", *European Journal of Computer Science and Information Technology*, 7(3), 1-8.
- Hall, M.A., (1999), "Correlation-based feature selection for machine learning", Ph.D. Thesis, University of Waikato.
- Khan, M.S.I., Islam, N., Uddin, J., Islam, S., and Nasir, M.K. (2022). Water quality prediction and classification based on principal component regression and gradient boosting classifier approach. *Journal of King Saud University-Computer and Information Sciences*, 34(8), 4773-4781, <https://doi.org/10.1016/j.jksuci.2021.06.003>
- Kavita, D., and Jagdish, S., (2012), "Water resources management and water quality, case of Bhopal", *International Conference on Chemical, Ecology and Environmental Sciences*, Bangkok.
- Khalil, B., Ouarda, T., and St-Hilaire, A., (2011), "Estimation of water quality characteristics at ungauged sites using artificial neural networks and canonical correlation analysis", *Journal of Hydrology*, 405, 277-287, <https://doi.org/10.1016/j.jhydrol.2011.05.024>.
- Kira, K., and Rendell, L. A., (1992), "The Feature Selection Problem: Traditional methods and a new algorithm", *Proceedings of the 10th National Conference on Artificial intelligence*, 129-134.
- Kotsianti, S.B., and Kanellopoulos, D., (2007), "Combining bagging, boosting and dagging for classification problems", In: Apolloni, B., Howlett, R.J., Jain, L. (eds.), *Knowledge-Based Intelligent Information and Engineering Systems, KES 2007*, Lecture Notes in Computer Science, Vol. 4693, Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-540-74827-4_62.
- Kotti, M.E., Vlessidis, A.G., Thanasoulas, N.C., and Evmiridis, N.P., (2005), "Assessment of river water quality in Northwestern Greece", *Water Resources Management*, 19, 77-94, <https://doi.org/10.1007/s11269-005-0294-z>.
- Meddouri, N., Khoufi, H., and Maddouri, M., (2021), "A performant dagging approach of classification based on formal concept", *International Journal of Artificial Intelligence and Machine Learning (IJAIML)*, 11(2), 38-62, <http://doi.org/10.4018/IJAIML.20210701.0a3>.
- Mehta, V., and Sanghavi, V., (2019), "Comparative study of various decision tree methods for data stream mining", In: *3rd International Conference on Information and Communication Technology (ICICT)*, Springer International Publishing, pp. 371-379, https://doi.org/10.1007/978-981-13-1165-9_34.
- Sattari, M.T., Feizi, H., Colak, M., Ozturk, A., Ozturk, F., and Apaydin, H., (2021), "Surface water quality classification using data mining approaches: Irrigation along the Aladag River", *Irrigation and Drainage*, 70(5), 1227-1246, <https://doi.org/10.1002/ird.2594>.
- Singh, D.F. (1992), "Studies on the water quality index of some major rivers of Pune, Maharashtra", *Proceedings of the Academy of Environmental Biology* 1, 61-66.



This article is an open-access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) license.