

کاهش ویژگی سیستم های اطلاعاتی ناقص با تئوری مجموعه راف فازی با استفاده از الگوریتم چرخه آب

وحید نجف پور^{۱*}، ناصر سلطانی^۲ و بهروز زادمهر^۳

۱- دانشکده مهندسی برق، دانشگاه پدافند هوایی خاتم الانبیا(ص)، تهران، ایران

۲- دانشکده مهندسی برق، دانشگاه پدافند هوایی خاتم الانبیا(ص)، تهران، ایران

۳- دانشکده مهندسی برق و کامپیوتر، دانشگاه امیر کبیر، تهران، ایران

چکیده

در سال های اخیر تئوری مجموعه راف به یکی از راه حل های قدرتمند در حل مسئله هوش مصنوعی و داده کاوی تبدیل شده است. کاهش ویژگی در سیستم های اطلاعاتی به معنای حذف ویژگی های غیرضروری و کاهش بعد داده ها است. استفاده از تئوری مجموعه های فازی و الگوریتم چرخ آب می تواند در کاهش ویژگی های ناقص سیستم های اطلاعاتی مفید باشد. تئوری مجموعه فازی و نظریه مجموعه های راف دو نظریه متمایز اما مکمل است که با عدم اطمینان در داده ها مواجه می شوند. ویژگی های برجسته هر دو نظریه در محدوده تئوری تنظیم راف فازی قرار گرفته است. هدف از این مقاله ارائه یک رویکرد مجموعه ای فازی بر اساس غلبه بر سیستم های اطلاعاتی با ارزش ناقص است. در این مقاله، یک رویکرد بهینه سازی جدید، که به عنوان الگوریتم چرخه آب شناخته می شود برای حل این مسئله مورد استفاده قرار گرفته است. برای ارزیابی عملکرد الگوریتم WCA در ۹ مجموعه داده ای شناخته شده از UCI که در جدول نتایج نشان داده شده مورد آزمایش قرار گرفته شده و تست گردید. نتایج آزمایشات، نشان می دهد که راف فازی و الگوریتم پیشنهاد شده نتایج مناسبی ارائه داد که درخور تأمل است.

واژه های کلیدی: الگوریتم چرخه آب، تئوری مجموعه راف، راف فازی، سیستم اطلاعاتی ناقص، کاهش ویژگی

Attribute Reduction of Incomplete Information systems based on Fuzzy Rough Set by the Water Cycle Algorithm

Abstract

In recent years, rough set theory has become one of the powerful solutions in solving the problem of artificial intelligence and data mining. Feature reduction in information systems means removing unnecessary features and reducing data dimension. Using the theory of fuzzy sets and the water wheel algorithm can be useful in reducing the incomplete features of information systems. Fuzzy set theory and rough set theory are two distinct but complementary theories that deal with uncertainty in data. The salient features of both theories are within the scope of fuzzy rough tuning theory. The purpose of this paper is to present a fuzzy set approach based on overcoming incomplete value information systems. In this article, a new optimization approach, which is known as water cycle algorithm, is used to solve this problem. To evaluate the performance of the WCA algorithm, it was tested on 9 datasets selected from UCI, which are shown in the results table. The results of the experiments show that the fuzzy rough and the proposed algorithm gave appropriate results, which is reasonable.

Key words: Water cycle algorithm, rough set theory, fuzzy rough, incomplete information system, feature reduction

های اطلاعاتی مفید باشد. برای این منظور، ابتدا باید داده های مورد نظر را به شکل مجموعه های فازی تبدیل کنیم. سپس با استفاده از الگوریتم چرخ آب، ویژگی های ناقص را شناسایی کرده و حذف می کنیم. این الگوریتم با استفاده از خصوصیت های مجموعه های فازی، ورودی هایی را که در محدوده تعریف شده قرار ندارند را شناسایی می کند و آن ها را حذف می کند. به عنوان مثال، در صورتی که داده هایی از سیستم های اطلاعاتی را داریم که شامل ویژگی های نام، سن، آدرس، تحصیلات و شغل هستند، می توانیم با استفاده از تئوری مجموعه های فازی و الگوریتم چرخ آب، ویژگی های ناقص را شناسایی کرده و حذف کنیم. به عنوان مثال، اگر ویژگی "آدرس" در بیشتر موارد از داده های ما مقدار نداشته باشد، می توانیم این ویژگی را حذف کنیم تا اطلاعات مورد نیاز را با دقت بیشتری به دست آوریم.

کاهش ویژگی به عنوان روش کاهش تعدادی از ویژگی ها در یک سیستم اطلاعات شناخته شده است و عمل آن یک مرحله حیاتی از پردازش داده کاوی است. در فرآیند کاهش ویژگی، حداقل مجموعه ای از صفات انتخاب شده است. افزایش حجم پایگاه داده ها مسئله کاهش ویژگی را ایجاد و پراهمیت می کند. در بررسی پایگاه های داده، گاهی با سیستم های اطلاعاتی ناقص مواجه هستیم. یک سیستم اطلاعاتی ناقص به جدولی از داده ها اطلاق می شود که برخی درایه های صفات آن مقداری ندارند. یک رویکرد برای تحلیل این سیستم ها این است که سطر مربوط به مقادیر ناقص را حذف کنیم. اما این رویکرد مناسب نیست زیرا اندازه داده ها را کاهش می دهد که ممکن است اطلاعات مفیدی از بین برود. یک روش دیگر برای تحلیل این نوع سیستم های اطلاعاتی این است که داده های ناموجود را بصورت های مختلفی مقاردهی کنیم. برای مثال از تحلیل های آماری استفاده کرده و داده های ناموجود را حدس بزنیم. هر دو رویکرد برای کامل شدن سیستم های اطلاعاتی کاربرد دارند اما بکارگیری یک روش مناسب برای کامل کردن سیستم های اطلاعاتی ناقص خود یک موضوع تحقیق پیچیده است.

یکی از روشهای کنترل و مدیریت داده ها با حجم زیاد تئوری مجموعه راف است. در این تئوری با تفکیک ویژگی ها به ویژگی های شرط و تصمیم گیری، به کمک تئوری نظریه مجموعه ها و رویکردهای ریاضی اقدام به شناسایی رابطه و وابستگی بین این دو گروه ویژگی ها می نماید. تئوری مجموعه راف تلاش می کند، یک زیر مجموعه ای با حداقل ویژگی از مجموعه ویژگی های اصلی با استفاده از روش تولید تصادفی ایجاد کند. اشکال اصلی روش تولید تصادفی، این است که تنها برای مجموعه داده های کوچک مناسب است و روش بسیار وقت گیر است. یک عامل مهم در فرآیند انتخاب ویژگی ها، جستجو در تمام فضای زیر مجموعه ویژگی های ممکن است. جستجو جامع در این فضا، معمولاً امکان پذیر نیست. روش کلاسیک معمولاً منجر به بهینه

در اختیار داشتن حجم زیاد داده ها در دنیای واقعی نمایانگر این است که ما از نظر مقدار داده غنی، و از نظر مقدار دانش، فقیر هستیم و به همین دلیل متدهای تحلیل داده کاوی جهت استخراج دانش از دل داده ها معرفی شده اند [۱]. همچنین با توجه به ذخیره ده ها، صدها یا حتی هزاران ویژگی در پایگاه داده برنامه های کاربردی، بحث کاهش ویژگی در جهت تحلیل این داده ها امری ضروری و انکار ناپذیر است [۲-۳]. تئوری مجموعه های فازی (Fuzzy Set Theory) یک روش ریاضی است که برای تعامل با عدم قطعیت در داده ها و اطلاعات استفاده می شود. یک مجموعه فازی به دنباله ای از اعداد فازی گفته می شود که هر عدد فازی نشان دهنده احتمال عضویت عنصر در آن مجموعه است. کاهش ویژگی بر مبنای تئوری مجموعه های فازی (Fuzzy Set Reduction) یکی از روش های مورد استفاده در پردازش داده های فازی است که به منظور کاهش بعد داده ها و حذف ویژگی های تکراری یا غیر ضروری، استفاده می شود. در مجموعه های فازی، ویژگی های تکراری یا غیر ضروری ممکن است باعث ایجاد ابهام در تصمیم گیری و حفظ سطح دقت مدل سازی نشوند. با استفاده از روش کاهش ویژگی، امکان حذف این ویژگی ها و در نتیجه سرعت بخشیدن به فرآیند آموزش مدل های پیش بینی و تصمیم گیری فراهم می شود. تحقیقات پیشین در زمینه کاهش ویژگی بر مبنای تئوری مجموعه های فازی، شامل روش های مختلفی مانند الگوریتم جستجوی ژنتیک، الگوریتم هایبیریدی، الگوریتم تکاملی، روش های خوشه بندی و غیره می شود. هر یک از این روش ها با استفاده از مفاهیم تئوری مجموعه های فازی و الگوریتم های خاص خود، قادر به کاهش بعد داده ها و یافتن ویژگی های مهم و حذف ویژگی های تکراری هستند. در بررسی پیشینه نظری کاهش ویژگی بر مبنای تئوری مجموعه راف، دو رویکرد اساسی را در مقابل داریم. ابتدا رویکردهایی که کلاسیک هستند و از همان نسخه اصلی تئوری مجموعه راف استفاده می کنند. رویکرد دوم، توسعه یافته رویکرد قبلی است که الگوریتم های اکتشافی کاهش ویژگی را ارائه می کند. نکته مهم در رویکردهای مختلف کاهش ویژگی بر مبنای تئوری مجموعه راف، نوع سیستم های اطلاعاتی است. در اغلب پژوهش های اجرا شده بر سیستم های اطلاعاتی کامل تمرکز شده است؛ اما در دنیای واقعی، مجموعه داده های جمع آوری شده همیشه کامل نیستند. این مشکل بخصوص در سامانه های اطلاعات نظامی بیشتر به چشم می خورد. پس چالش اساسی کاهش ویژگی مبتنی بر تئوری مجموعه راف، کار روی سیستم های اطلاعاتی ناقص است.

کاهش ویژگی در سیستم های اطلاعاتی به معنای حذف ویژگی های غیر ضروری و کاهش بعد داده ها است. استفاده از تئوری مجموعه های فازی و الگوریتم چرخ آب می تواند در کاهش ویژگی های ناقص سیستم

این تئوری دانست که مهمترین آن ها به شرح زیر است: این تئوری تنها بر پایه داده های خام بنا شده و نیاز به اطلاعات خارجی ندارد. RST نه تنها برای تحلیل ویژگی های کیفی، بلکه برای مطالعه ویژگی های کمی نیز به کار گرفته می شود. این تئوری امکان کاهش معیارهای اضافی را داشته و الگوریتم هایی ساده برای این کار ارائه می دهد. نتیجه RST یکسری قوانین تصمیم گیری است که از مدل های مختلف آن استخراج می شود.

RST نیازی به اصلاح ناسازگاری ها نداشته و می تواند قوانین متناقض را به قوانین قطعی و غیر قطعی تقسیم کند؛ به طوری که قوانین ناشی از این مدل ها به راحتی قابل درک باشند. نقطه آغاز RST مجموعه داده هاست که معمولاً در قالب یک جدول، سازمان داده شده و تحت نام سیستم اطلاعاتی یا پایگاه داده معرفی می شوند. اصلی ترین عملیات در RST تقریب بالا و پایین است و این تقریبهها برای تشخیص اجزای کاملاً وابسته یا نیمه وابسته در جدول اطلاعات به کار می روند. تقلیل داده نیز یکی از مهمترین مفاهیم این تئوری است، که در ادامه هر یک از این مفاهیم اساسی تشریح می شود [۱۲]. با این حال یکی از محدودیت های اصلی تئوری راف این است که مبتنی بر داده های گسسته است. در تئوری راف کالاسیک، رسیدگی به داده های نوپز دار و داده های حقیقی و پیوسته ممکن نیست. در نتیجه بسیاری از محققان، به دنبال تعمیم تئوری راف بودند، بطوری که بتوان از آن برای رسیدگی به مجموعه داده های نوپز دار و پیوسته استفاده کرد [۱۳]. یکی از ایدئولوژی های بکار رفته در تئوری مجموعه راف این است که فقط از داده های تولید شده استفاده می کند؛ در حالی که در سایر روش های موجود برای حل مسئله علاوه بر داده های تولید شده، نیاز به یک دانش تکمیلی دیگری هم دارند.

۲-۱- فرایند کاهش راف فازی

مجموعه راف فازی که بر پایه انتخاب ویژگی است، براساس مفهوم پایین تقریب فازی، عملیات کاهش مجموعه داده ها که شامل ویژگی های حقیقی می شود، را آغاز می کند. عملگر منطقه مثبت در حالت مجموعه راف سنتی تنها بوسیله عملگر پایین تقریب تعریف می شود، در حالی که در حالت راف فازی، به همان حالت است، تنها درجه عضویت شی $X \in U$ به آن اضافه شده [۱۴]، در نتیجه منطقه مثبت در حالت فازی به این صورت است [۱۵]:

$$\mu_{POS_P(Q)}(x) = \sup_{x \in U/Q} \mu_{PX}(x). \quad (1)$$

با استفاده از تعریف منطقه مثبت، تابع وابستگی راف فازی به این صورت تعریف می شود [۱۱]:

$$\gamma_P(Q) = \frac{|\mu_{POS_P(Q)}(x)|}{|U|} = \frac{\sum_{x \in U} \mu_{POS_P(Q)}(x)}{|U|}. \quad (2)$$

محل می شود. به منظور این نقطه ضعف، استراتژی های جستجو مبتنی بر الگوریتم های فرابتکاری پیشنهاد شده است. هدف این الگوریتم ها، رسیدن به بهترین جواب یافت شده در کوتاهترین زمان ممکن می باشد. امروزه با بزرگ و پیچیده شدن مسائل، استقبال از روش های فرابتکاری بطور چشم گیری افزایش یافته است. این امر باعث شده بسیاری از محققان به سمت بررسی مشکلات کاهش ویژگی با استفاده از انواع مختلفی از الگوریتم ها، مانند: الگوریتم ژنتیک [۴]، بهینه سازی ازدحام ذرات [۵]، جستجو پراکنده [۶]، بهینه سازی کلونی مورچه ها و زنبور عسل [۷-۹] شبیه سازی حرارتی [۱۰] تشویق گردند.

اغلب روش های بهینه سازی مبتنی بر هوش مصنوعی همانند الگوریتم ژنتیک، شبیه سازی فرایندهای طبیعی می باشند. یکی از دلایل این امر ملموس بودن، سادگی فرموله کردن و درک تکامل این الگوریتم ها است. الگوریتم چرخه آب، با الهام گیری از فرایند چرخه آب در طبیعت نسبت به دیگر روش های مطرح شده در این زمینه دارای توانایی بالایی بوده و از سرعت مناسبی نیز برخوردار می باشد. در اغلب پژوهش های انجام شده، سیستم های اطلاعاتی کامل مورد هدف بوده اند. اما در دنیای واقعی، مجموعه داده های جمع آوری شده همیشه بصورت کامل نیستند. پس یک چالش اساسی کاهش ویژگی، کار بر روی سیستم های اطلاعاتی ناقص می باشد. در گذشته، روش هایی برای حل مسئله کاهش ویژگی معرفی شده اند؛ در این مقاله از تئوری مجموعه راف برای سیستم های اطلاعاتی ناقص که از جمله مسائل NP-hard بشمار می آید، استفاده می شود. رویکرد مجموعه راف را به سیستم های اطلاعات ناقص، یک سیستم اطلاعاتی ناقص به جداولی از داده ها اطلاق می شود که برخی درایه های صفات آن مقداری ندارند یعنی سیستم هایی که مقادیر مشخصه برای اشیاء ممکن است ناشناخته (null, missing) باشند، ارائه می دهیم. نگرانی اصلی ما برای یافتن قوانین از چنین سیستم هایی است [۱۱].

در ادامه این مقاله در بخش ۲، تئوری مجموعه راف، راف ناقص و راف فازی به طور دقیق تشریح شده است. ساختار الگوریتم چرخه آب نحوه سازگاری آن با مسئله کاهش ویژگی در بخش ۳ و نتایج پیاده سازی الگوریتم مذکور در بخش ۴ ارائه شده است. نهایتاً، خلاصه و نتیجه گیری و رویکرد پیشنهادی برای کاهش ویژگی و دستاوردهای آن در بخش ۵ جمع بندی شده است.

۲- تئوری مجموعه راف (RST)

در اختیار داشتن حجم زیاد داده ها در دنیای واقعی نمایانگر این است که ما از نظر مقدار داده غنی، و از نظر مقدار دانش، فقیر هستیم و به همین دلیل متدهای تحلیل داده کاوی جهت استخراج دانش از دل داده ها معرفی شده اند. تئوری مجموعه راف یکی از این روش هاست که توسط محقق لهستانی پاولاک در سال ۱۹۸۰ معرفی گردیده است [۴]. موفقیت RST و توسعه آن را می توان به دلیل قابلیت ها و مزایای

جدول ۱- نمونه‌ای از یک سیستم اطلاعاتی ناقص

افراد	قد	وزن	سن
شماره ۱	بلند	زیاد	خیلی بالا
شماره ۲	کوتاه	زیاد	بالا
شماره ۳	*	کم	پایین
شماره ۴	بلند	زیاد	خیلی بالا
شماره ۵	خیلی بلند	*	خیلی بالا
شماره ۶	بلند	متوسط	*
شماره ۷	*	کم	بالا
شماره ۸	بلند	*	خیلی پایین

۳- الگوریتم چرخه آب برای کاهش ویژگی

مجموعه راف فازی برای انتخاب ویژگی، به منظور کاهش ابعاد داده‌های حجیم و استفاده از ویژگی‌های مهم‌تر برای تحلیل و پردازش داده‌ها استفاده می‌شود. در این روش، با استفاده از مفهوم پایین تقریب فازی، ابعاد داده‌ها با حفظ دقت مورد نیاز، کاهش می‌یابد. برای این منظور، ابتدا یک مجموعه اولیه از ویژگی‌ها به عنوان ورودی در نظر گرفته می‌شود. سپس با استفاده از الگوریتم‌های مختلفی مانند الگوریتم ژنتیک، الگوریتم جستجوی محلی و غیره، فضای ویژگی‌ها جستجو شده و ویژگی‌های بی‌اهمیت حذف می‌شوند. در این روش، با توجه به مفهوم پایین تقریب فازی، برای هر ویژگی، یک مقدار عضویت در نظر گرفته می‌شود. سپس با استفاده از الگوریتم‌های کاهش ابعاد، ویژگی‌هایی که دارای مقدار عضویت پایینی هستند، حذف می‌شوند.

به عنوان مثال، فرض کنید یک مجموعه داده شامل ۱۰ ویژگی باشد. با استفاده از مجموعه راف فازی، می‌توانیم تعدادی از ویژگی‌های مهم را انتخاب کنیم و ویژگی‌های بی‌اهمیت را حذف کنیم. برای مثال، اگر ویژگی‌های ۲ و ۶ و ۹ دارای مقدار عضویت پایینی هستند، آن‌ها حذف می‌شوند و مجموعه داده به صورت ۷ ویژگی کاهش می‌یابد. در نهایت، با استفاده از مجموعه راف فازی، می‌توانیم تنها با استفاده از ویژگی‌های مهم‌تر داده‌ها را تحلیل کنیم و همچنین در زمان اجرای الگوریتم‌ها صرفه جویی شود.

الگوریتم چرخه آب، به عنوان یک روش بهینه‌سازی جدید که توسط اسکندر و همکاران در سال ۲۰۱۲ معرفی شده است، این الگوریتم همانند چرخه آب در طبیعت، با الگوبرداری از تبخیر آب از اقیانوس تشکیل ابر، تشکیل رودخانه و همچنین الگوبرداری از سرازیر شدن آب از گودال‌ها در طبیعت ارائه شده است که قابلیت بالایی در فرار از بهینه محلی و همچنین سرعت زیادی در رسیدن به بهینه سراسری دارد [۷]. مشابه سایر الگوریتم‌های متاهوریستیک، این الگوریتم هم با یک جمعیت اولیه آغاز می‌شود که قطرات باران نامیده می‌شوند. در ابتدا ما فرض می‌کنیم باران یا تگرگ داریم و بهترین شخص را به عنوان دریا در نظر می‌گیریم. پس تعدادی از قطرات باران را که موقعیت بهتری دارند، به عنوان رودخانه و بقیه قطرات باران را به عنوان جریان‌هایی که به سمت رودخانه و دریا جاری هستند، در نظر می‌گیریم. مقدار آبی

در حالت crisp مجموعه راف، وابستگی Q بر روی P یعنی، نسبت اشیائی که از کل مجموعه داده قابل تشخیص هستند. اما در حالت فازی، وابستگی یعنی تعیین کاردینالیته فازی $\mu_{POS_P(Q)}(X)$ ، تقسیم بر تعداد اشیاء موجود در مجموعه مرجع. اگر روند مجموعه راف فازی بدرستی نتیجه دهد، این امکان وجود دارد که بتوان آن را برای ویژگی‌های با خاصیت‌های متعدد استفاده کرد و وابستگی بین زیر مجموعه‌های مختلف از مجموعه ویژگی‌های اصلی را می‌توان براحتی پیدا کرد. برای مثال، می‌توان درجه وابستگی ویژگی‌های شرطی $P = \{a, b\}$ را با توجه به ویژگی تصمیم‌گیری مورد نظر بدست آورد. در حالت crisp، عملگر U/P دربردارنده گروهی از اشیاء است که با توجه به هر دو ویژگی a, b نوشته می‌شود. اما در حالت فازی اشیاء ممکن است به بسیاری از کلاس‌های هم‌ارزی تعلق داشته باشند، بنابراین حاصل ضرب دکارتی $U/IND(\{a\})$ و $U/IND(\{b\})$ باید در نظر گرفته شود. در حالت کلی:

$$U/P = \otimes \{a \in P: U/IND(\{a\})\} \quad (۳)$$

هر مجموعه‌ای که در U/P است، مشخص کننده کلاس هم‌ارزی است. برای مثال اگر

$$P = \{a, b\}, \quad U/IND(\{a\}) = \{N_a, Z_a\}, \quad U/IND(\{b\}) = \{N_b, Z_b\}, \quad (۴)$$

آنگاه:

$$U/P = \{N_a \cap N_b, N_a \cap Z_b, Z_a \cap N_b, Z_a \cap Z_b\} \quad (۵)$$

میزان اینکه یک شیء متعلق به کلاس هم‌ارزی است، با استفاده از ترکیب کلاس‌های هم‌ارزی فازی تشکیل دهنده، محاسبه می‌شود. مثلاً:

$$F_i, i = 1, 2, \dots, n$$

$$\mu_{F_1 \cap \dots \cap F_n}(x) = \min(\mu_{F_1}(x), \mu_{F_2}(x), \dots, \mu_{F_n}(x)) \quad (۶)$$

۲-۲- سیستم‌های اطلاعاتی ناقص

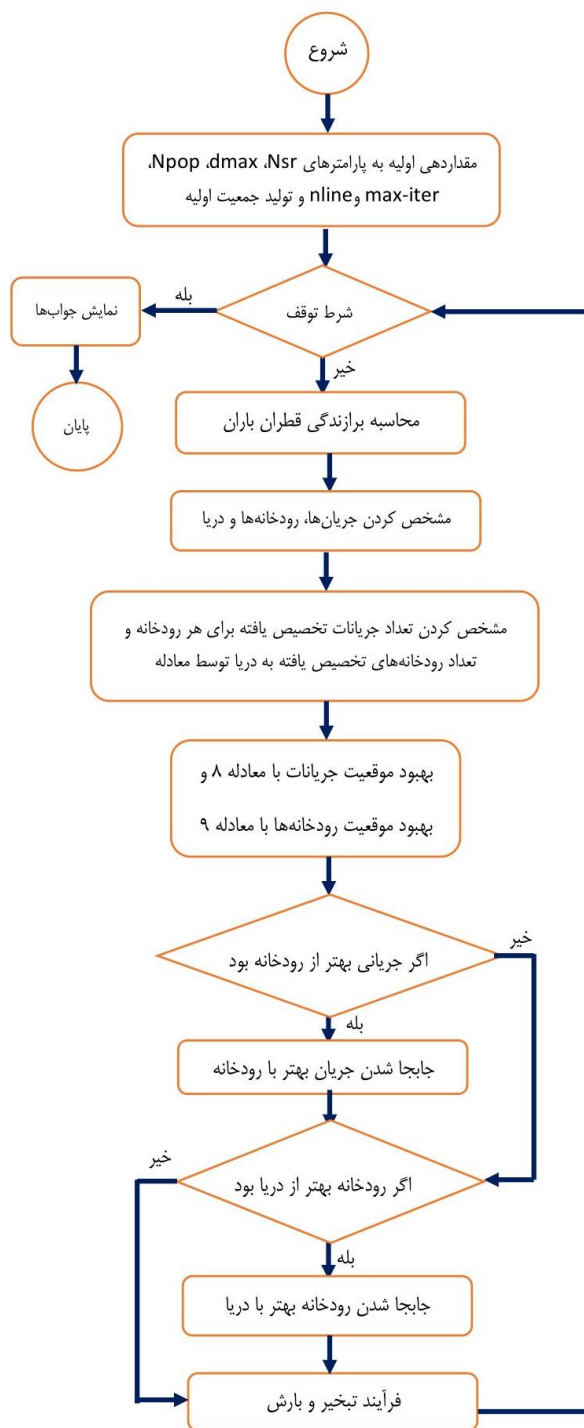
یک سیستم اطلاعاتی معمولاً به فرم زیر ارائه می‌شود:

$$IS = (U, A, \{V_a\}, f_a)_{a \in A} \quad (۷)$$

که در رابطه ۱، U، تعداد متناهی از اشیای دارای مقدار و A، تعداد محدودی از ویژگی‌های دارای مقدار می‌باشند. V_a ، مجموعه مقادیر ویژگی‌ها و f_a ، یک تابع اطلاعاتی از سیستم اطلاعات است که V_a را به U نگاشت می‌کند. اگر در سیستم اطلاعاتی یک شی x و یک ویژگی به نام a وجود داشته باشد که تابع $f_a(x)$ مقداری نداشته باشد، آنگاه به آن سیستم اطلاعاتی ناقص گویند.

مثال: جدول شماره ۱ نمونه‌ای از سیستم اطلاعاتی ناقص را نشان می‌دهد. این سیستم شامل اطلاعات مربوط به چند فرد است که ویژگی‌های آن‌ها شامل قد، وزن و سن می‌باشد. هر کدام از مقادیر مربوط به این ویژگی‌ها با مقادیر کیفی مشخص شده‌اند.

که از یک جریان به رودخانه یا دریا می‌ریزد، متفاوت است با جریانی دیگر، به علاوه رودخانه‌ها به دریا که در مکان پایین‌تری قرار دارد، سرازیر می‌شوند. در شکل ۱ فلوچارت چرخه آب را آورده ایم.



شکل ۱- فلوچارت الگوریتم چرخه آب

۳-۱-گام اول: ایجاد جمعیت اولیه و مقداردهی اولیه به پارامترها

ایجاد جمعیت اولیه: برای حل یک مسئله بهینه‌سازی با استفاده از روش‌های متاهیورستیک مبتنی بر جمعیت لازم است که مقادیر

متغیرهای مسئله همانند یک آرایه تشکیل شوند. در PSO و GA، این آرایه، موقعیت ذره و کروموزوم نامیده می‌شود. در این روش به آن قطره‌ی باران می‌گویند. در یک مسئله بهینه‌سازی N بعدی، یک قطره‌ی باران، آرایه‌ی $1 \times N$ تایی است که این آرایه به صورت معادله (۸) تعریف شده است:

$$Raindrop = [x_1, x_2, x_3, \dots, x_n] \quad (8)$$

مقداردهی اولیه: به تعداد Npop قطره‌ی باران ایجاد می‌شود. برای این کار، یک ماتریس با سائز $N_{pop} \times N_{var}$ تولید شده که سطرها تعداد اعضای جمعیت و ستون‌ها تعداد متغیرها هستند و در واقع نشان‌دهنده جمعیتی از قطرات باران است. مقادیر متغیرها با توجه به مسئله‌ی مورد بررسی می‌تواند عدد حقیقی، اعشاری یا گسسته باشد. از این رو، ماتریس X به صورت تصادفی تولید شده و به صورت معادله (۹) است:

$$Population\ of\ raindrops = \begin{bmatrix} Raindrop_1 \\ Raindrop_2 \\ Raindrop_3 \\ \vdots \\ Raindrop_{N_{pop}} \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_{N_{var}}^1 \\ x_1^2 & x_2^2 & \dots & x_{N_{var}}^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{N_{pop}} & x_2^{N_{pop}} & \dots & x_{N_{var}}^{N_{pop}} \end{bmatrix} \quad (9)$$

۳-۲-گام دوم: ارزیابی جمعیت

هزینه‌ی یک قطره‌ی باران به صورت معادله (۱۰) ارزیابی می‌شود:

$$C_i = Cost_i = f(x_1^i, x_2^i, \dots, x_{N_{var}}^i) \quad i = 1, 2, 3, N_{pop} \quad (10)$$

۳-۳-گام سوم: انتخاب دریا و رودخانه

به تعداد Nsr از بهترین افراد (که دارای کمترین مقدار Cost هستند) به عنوان دریا و رودخانه‌ها انتخاب می‌شوند. قطره‌ای که کمترین مقدار را در میان دیگران دارد، به عنوان دریا در نظر گرفته می‌شود. در حقیقت، Nsr مجموعه‌ای از رودخانه‌ها و یک دریا است که این تعداد از معادله‌ی (۱۱) محاسبه می‌شود:

$$N_{sr} = Number\ of\ Rivers + \frac{1}{sea} \quad (11)$$

باقی جمعیت (قطرات باران)، جریان‌ها را شکل می‌دهند. این جریانات به سمت رودخانه‌ها یا مستقیماً به سمت دریا جاری می‌شوند.

$$N_{Raondrops} = N_{pop} - N_{sr} \quad (12)$$

برای تخصیص قطرات باران به رودخانه‌ها یا دریا با توجه به شدت جاری شدن، از معادله‌ی (۱۳) استفاده می‌کنیم:

$$NS_n = round \left\{ \left| \frac{Cost_n}{\sum_{i=1}^{N_{sr}} Cost_i} \right| \times N_{Raondrops} \right\} \cdot n = 1, 2, \dots, N_{sr} \quad (13)$$

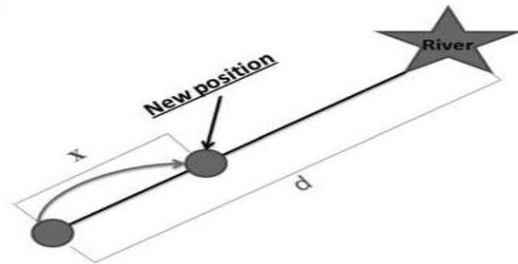
که NS_n تعداد جریاناتی است که به رودخانه‌های خاص یا دریا جاری می‌شوند.

۳-۴-گام چهارم: بهبود موقعیت

همه رودخانه‌ها و جریان‌ها به دریا ختم می‌شوند (بهترین نقطه‌ی بهینه). یک جریان در طول خط اتصال با رودخانه به سمت رودخانه جاری می‌شود، توسط فاصله‌ای که به صورت تصادفی با معادله (۱۴) انتخاب می‌شود:

$$X \in (0, C \times d), C > 1 \quad (14)$$

که C یک مقدار بین ۱ و ۲ است (نزدیک به ۲). بهترین مقدار برای C می‌تواند ۲ باشد. فاصله‌ی بین جریان (نهر) و رودخانه d در نظر گرفته شده است. مقدار X در معادله‌ی (۱۴) با یک عدد تصادفی توزیع شده بین $(0, C \times d)$ مرتبط است. مقدار C بیشتر از یک باعث می‌شود که جریان‌ها (نهرها) قادر باشند، در جهت‌های مختلف به سمت رودخانه‌ها جاری شوند.



شکل ۲- دید سمانتیک از جاری شدن نهرها به یک رودخانه

دیگر

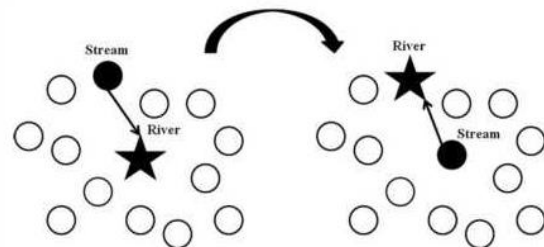
این مفهوم هم‌چنین ممکن است در جریان رودخانه‌ها به سمت دریا هم استفاده شود. بنابراین موقعیت جدید جریان‌ها و رودخانه‌ها به صورت معادله (۱۵) و (۱۶) بدست می‌آید:

(۱۵)

$$X_{stream}^{i+1} = X_{stream}^i + rand \times C \times (X_{River}^i - X_{stream}^i)$$

$$X_{River}^{i+1} = X_{River}^i + rand \times C \times (X_{sea}^i - X_{River}^i) \quad (16)$$

که $rand$ یک عدد تصادفی توزیع شده به صورت یکنواخت است بین $[0, 1]$. اگر راه‌حل داده شده توسط یک جریان بهتر از رودخانه‌ی متصل به آن باشد، موقعیت‌های رودخانه و نهر (جریان) جابجا می‌شود (یعنی رودخانه نهر می‌شود و نهر رودخانه می‌شود). چنین جابجایی‌ای می‌تواند به طور مشابه برای رودخانه‌ها و دریا نیز اتفاق بیفتد. شکل ۳، جابجایی جریانی را که در میان دیگر جریان‌ها بهترین است را با رودخانه نشان می‌دهد.



شکل ۳- مبادله موقعیت جریان‌ها و رودخانه‌ها که ستاره نماد

رودخانه و دایره با رنگ سیاه نماد بهترین جریان در بین دیگر جریان‌ها است.

۳-۵-گام پنجم: تبخیر

تبخیر یکی از فاکتورهای بسیار مهمی است که می‌تواند از همگرایی زودرس الگوریتم جلوگیری کند. آب تبخیر شده در اتمسفر برای تشکیل ابر در اتمسفر سردتر، متراکم می‌شود و ابر آب را به شکل باران به زمین برمی‌گرداند. باران جریان‌های (نهرهای) جدید را ایجاد می‌کند و جریان‌های جدید به سمت رودخانه‌ها یا دریا جاری می‌شوند. در روش ارائه شده فرآیند تبخیر، باعث می‌شود که آب دریا مانند رودخانه/نهر که به سمت دریا جاری هستند، تبخیر شود. این فرض به منظور جلوگیری از قرار گرفتن در بهینه محلی در نظر گرفته شده است. شبه کد زیر را در نظر بگیرید:

$$\text{If } |X_{sea}^i - X_{River}^i| < d_{max} \quad i = 1, 2, 3, \dots, N_{sr} - 1$$

Evaaporation and raining process end
dmax

یک عدد کوچک نزدیک به صفر است. اگر فاصله بین رودخانه و دریا کمتر از d_{max} باشد، نشان می‌دهد که رودخانه به دریا ملحق شده است. در این وضعیت، فرآیند تبخیر انجام می‌شود و همان‌طور که در طبیعت می‌بینیم بعد از میزان کافی عمل تبخیر، رطوبت آغاز خواهد شد.

یک مقدار بزرگ برای d_{max} جستجوی نزدیک به دریا را کاهش می‌دهد، در حالی که مقدار کوچک باعث می‌شود، جستجو نزدیک دریا شدت یابد. بنابراین d_{max} شدت جستجو در نزدیکی دریا را کنترل می‌کند (راه‌حل بهینه). مقدار d_{max} به صورت تطبیقی به صورت زیر کاهش می‌یابد:

$$d_{max}^{i+1} = d_{max}^i - \frac{d_{max}^i}{\max \text{ iteration}} \quad (17)$$

۳-۶-گام ششم: فرآیند بارش

بعد از عمل فرآیند تبخیر، بارش انجام می‌شود. در فرآیند بارندگی قطرات باران جدید جریانی را در مکان‌های مختلف شکل می‌دهند (شبه عملگر جهش در الگوریتم ژنتیک). برای مشخص کردن مکان‌های نهرهای جدید و تولید شده، معادله‌ی (۱۸) استفاده می‌شود:

$$x_{stream}^{new} = LB + rand \times (UB - LB) \quad (18)$$

LB و UB محدوده‌ی پایین و بالایی هستند که توسط مسئله تعریف شده‌اند. مجدداً قطره‌ی بارانی که جدیداً تشکیل شده است، به عنوان یک رودخانه که به سمت دریا می‌رود، در نظر گرفته می‌شود. بقیه‌ی قطرات جدید برای تشکیل نهرهای جدیدی که به سمت رودخانه یا دریا جاری هستند، فرض شده‌اند. به منظور افزایش نرخ همگرایی و عملکرد محاسباتی الگوریتم معادله‌ی (۱۹) فقط برای جریان‌هایی که به سمت دریا جاری هستند در نظر گرفته شده است. این معادله قصد دارد کاوش در نزدیکی دریا را بهبود بخشد.

$$X_{stream}^{new} = X_{sea} + \sqrt{\mu} \times randn(1, N_{var}) \quad (19)$$

که μ یک ضریب است که محدوده‌ی ناحیه‌ی جستجو نزدیک دریا را نشان می‌دهد. randn یک عدد تصادفی توزیع شده‌ی یکنواخت است. مقادیر بزرگ‌تر برای μ امکان خروج از ناحیه‌ی موجه را افزایش می‌دهد. به عبارت دیگر مقادیر کوچک‌تر برای μ باعث می‌شود الگوریتم در ناحیه‌ی کوچک‌تر نزدیک دریا به جستجو بپردازد. مقدار مناسب برای μ ، ۰.۱ است.

راه‌حل الگوریتم چرخه آب مجموعه‌ای از راه‌حل‌ها (موقعیت‌ها) و یا مجموعه حالاتی از کارها را معرفی می‌گویند. برای شروع کار چرخه آب تنظیم راه‌حل کار مهم و ضروری است. چگونگی نمایش و ارائه راه‌حل از روی مشکلی که باید حل شود، تعیین می‌گردد. در واقع این امر، این اجازه را می‌دهد که یک راه‌حل به راحتی ساخته و به سرعت مورد ارزیابی قرار گیرد. الگوریتم چرخه آب مبتنی بر جمعیت است، پس با مجموعه‌ای از راه‌حل‌ها کار خود را شروع می‌کند. که تعداد این راه‌حل‌ها با تعداد ویژگی‌های شرطی $|C|$ برابر است. در این الگوریتم هر راه‌حل معادل یک قطره باران است. که به حالات مختلف تولید می‌شود. تولید راه‌حل در الگوریتم چرخه آب برای حل مسئله کاهش ویژگی با توجه به اینکه از کدام حالت تئوری راف استفاده شود، متفاوت است. الگوریتم چرخه آب از راه‌حل‌های باینری استفاده می‌کند. ساخت راه‌های اولیه که در این مرحله صورت می‌گیرند، اصولاً بصورت تصادفی انجام می‌گیرد. برخی از راه‌حل‌های اولیه‌ی تولید شده، دارای درجه وابستگی یک می‌باشند. ویژگی‌ها که در اینجا همان راه‌حل‌ها هستند، با برداری به نام Y نمایش داده می‌شوند. در میان تمام ویژگی‌ها از بردار Y ، اگر خانه y_i از بردار Y که $i=1, \dots, |C|$ مقدارش برابر با یک باشد، آنگاه ویژگی i^{th} به عنوان بخشی از زیر مجموعه راه‌حل‌ها انتخاب می‌شود؛ و اگر مقدار خانه از بردار Y برابر با صفر باشد، آن ویژگی بخشی از راه‌حل محسوب نمی‌شود. در این روش ویژگی‌ها بصورت تصادفی انتخاب می‌شوند.

۳-۷- شرط خاتمه

شرط خاتمه در الگوریتم را می‌توان تعداد دفعات تکرار الگوریتم یا رسیدن به بهترین زیر مجموعه حداقلی تعریف کرد. در این پژوهش ما شرط پایان الگوریتم را مطابق با سایر پژوهش‌های انجام شده در این زمینه [20, 7]، تعداد دفعات تکرار الگوریتم و برابر با ۲۰ بار تعریف می‌کنیم.

۴- نتایج پیاده سازی

برای ارزیابی عملکرد الگوریتم WCA در ۹ مجموعه داده‌ی شناخته شده از UCI که در سایت <http://www.ics.uci.edu/ml> قابل دسترسی است که در جدول (۲) نشان داده شده مورد آزمایش قرار گرفته شده است. در ادامه نتایج بدست آمده کاهش ویژگی با استفاده از الگوریتم چرخه آب (WCA) با نتایج به دست آمده با الگوریتم‌های معروف ابتکاری ژنتیک (GA) و جستجوی هارمونی (HS) در جدول (۳)

آمده است. این نتایج بهترین از ۲۰ بار تکرار هر الگوریتم است که تعداد ویژگی‌های کاهش یافته را نشان می‌دهد. برای برنامه نویسی این الگوریتم‌ها از متلب Matlab R2017a استفاده شده است. الگوریتم‌ها بر روی یک کامپیوتر با حافظه اصلی ۴ گیگابایت و پردازنده دو هسته‌ای ۲.۲ گیگا هرتز اینتل اجرا شده اند در همه‌ی این الگوریتم‌ها مقدار جمعیت اولیه برابر ۱۰۰ و شرط توقف نیز تعداد تکرار ۲۰ در نظر گرفته شده است.

جدول ۲- داده‌های پایگاه داده UCI

تعداد سطرها	تعداد ویژگی‌ها	نام دیتاست
14	297	Cleveland
10	214	Glass
14	270	Heart
35	230	Ionosphere
26	120	Olitos
39	390	Water 2
39	390	Water3
2557	149	Web
14	178	Wine

۵- نتیجه‌گیری و بحث

نتایج برآمده از متد ارائه شده، بخوبی نمایانگر توانایی‌های الگوریتم چرخه آب در امر کاهش ویژگی سیستم‌های اطلاعاتی ناقص می‌باشد. با بررسی دیدگاه‌های مختلف میتوان پی برد که؛ زیرمجموعه‌های کاهش یافته قابل قبول، موازی سازی الگوریتم به منظور صرف زمان کمتر در رسیدن به پاسخ، کاردینالیته کمتر و کیفیت، دقت و صحت بالا در امر کاهش ویژگی همه و همه دلایلی بر برتری مزایا و راندمان مناسب این روش پیشنهادی می‌باشند. بررسی الگوریتم چرخه آب بعنوان روشی مبتنی بر جمعیت، میتواند سبب بهبود تمام راه حل‌های موجود (در جمعیت) گردد. داده‌های شناخته شده از UCI که در ۹ مجموعه داده استفاده شده‌اند و الگوریتم WCA بر روی آن‌ها آزمایش شده است، نتایج نشان می‌دهد که الگوریتم پیشنهادی بهترین عملکرد را نسبت به راف فازی و الگوریتم WCA ارائه می‌دهد. این نتایج نشان می‌دهند که الگوریتم پیشنهادی در تحلیل داده‌ها و استخراج اطلاعات، عملکرد بهتری نسبت به روش‌های موجود دارد. با این حال، برای ارزیابی بهتر و قطعی‌تر عملکرد الگوریتم پیشنهادی، می‌توان بر روی مجموعه داده‌های دیگری نیز آزمایشات جدیدی انجام داد.

جدول ۳- نتایج بدست آمده توسط هر یک از روش ها

الگوریتم	WCA		GA		HS	
	تعداد ویژگی ها	زمان اجرا(ثانیه)	تعداد ویژگی ها	زمان اجرا(ثانیه)	تعداد ویژگی ها	زمان اجرا(ثانیه)
Cleveland	$9^{12} 10^8$	521 second	$10^{19} 11^1$	490 second	$9^{10} 10^{10}$	671 second
Glass	$9^{11} 10^9$	320 second	$9^{10} 10^{10}$	501 second	$10^{19} 11^1$	452 second
Heart	8^{20}	40 second	$8^{18} 9^2$	77second	$8^{17} 9^3$	129 second
Ionosphere	$8^{19} 9^1$	373 second	$9^{17} 10^3$	800 second	$8^{15} 9^5$	326 second
Olitos	$6^{18} 9^2$	238 second	7^{20}	130 second	$6^{16} 7^4$	197 second
Water 2	7^{20}	829 second	$7^{19} 8^1$	1003 second	$7^{17} 8^3$	3123 second
Water3	$6^{13} 7^7$	3618 second	8^{20}	2321 second	$7^{19} 8^1$	1928 second
Web	$18^{12} 19^6 20^2$	2238 second	20^{20}	2520 second	$20^{18} 2^{12}$	2639 second
Wine	6^{20}	156 second	$6^{19} 7^1$	201 second	6^{20}	195 second

مراجع

[8] Palanisamy, S. and S. Kanmani, "Artificial Bee Colony Approach for Optimizing Feature Selectio," *International Journal of Computer Science Issues (IJCSI)*, vol. 9, no. 3, 2012.

[9] Ke, L., Z. Feng, and Z. Ren, "An efficient ant colony optimization approach to attribute reduction in rough set theory," *Pattern Recognition Letters*, vol. 29, no. 9, pp. 1351-1357, 2008.

[10] Abdolrazzagh-Nezhad, M., and Ali Adibiyan. "Attribute reduction based on rough set theory by soccer league competition algorithm," *Iranian Journal of Electrical and Computer Engineering*, vol. 82, no. 3, 2021.

[11] Pawlak, Z., "Rough sets," *International journal of computer & information sciences*, vol. 11, no. 5, pp. 341-356, 1982.

[12] Jain, Pankhuri, Anoop Tiwari, and Tanmoy Som. "Fuzzy rough assisted missing value imputation and feature selection." *Neural Computing and Applications*, vol. 35. no.3, pp. 2773-2793, 2023.

[13] Al Shalabi, L., Z. Shaaban, and B. Kasasbeh, "Data mining: A preprocessing engine," *Journal of Computer Science*, vol. 2, no. 9, pp. 735-739, 2006.

[14] Zadeh, L.A., "The concept of a linguistic variable and its application to approximate reasoning—I," *Information sciences*, vol. 8, no. 3, pp. 199-249, 1975.

[15] Ayesha, Shaeela, Muhammad Kashif Hanif, and Ramzan Talib. "Overview and comparative study of dimensionality reduction techniques for high dimensional data." *Information Fusion*, vol. 59, pp. 44-58, 2020.

[1] Han, J., J. Pei, and M. Kamber, *Data mining: concepts and techniques*. 2011: Elsevier.

[2] Guyon, I. and A. Elisseeff, An introduction to variable and feature selection. "*Journal of machine learning research*," vol. 3, pp. 1157-1182, 2003.

[3] Yu, J., "General C-means clustering model," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 8, pp.1197-1211, 2005.

[4] Al Imran, Abdullah, Md Rifatul Islam Rifat, and Rafeed Mohammad. "Enhancing the classification performance of lower back pain symptoms using genetic algorithm-based feature selection." *Proceedings of International Joint Conference on Computational Intelligence: IJCCI 2018*. Singapore: Springer Singapore, 2019.

[5] Tawhid, Mohamed A., and Abdelmonem M. Ibrahim. "Hybrid binary particle swarm optimization and flower pollination algorithm based on rough set approach for feature selection problem." *Nature-inspired computation in data mining and machine learning*, pp. 249-273, 2020.

[6] Som, Tanmoy, et al. "Fuzzy Rough Set Theory-Based Feature Selection: A Review." *Mathematical Methods in Interdisciplinary Sciences*, pp. 145-166, 2020.

[7] Nasir, Mohammad, et al. "A comprehensive review on water cycle algorithm and its applications." *Neural Computing and Applications*, vol. 32, pp. 17433-17488, 2020.