

بررسی و خوشه بندی ترکیبات شیمیایی با مدل ترکیبی فازی نزدیکترین همسایگی در راستای سیاستگذاری و تحلیل تولید علم کاربردی در ایران

تاریخ دریافت: ۱۳۹۸/۵/۲

مریم عابدی^۱

تاریخ پذیرش: ۱۳۹۸/۱۱/۳۰

هادی یزدانی^۲

چکیده

ترکیبات شیمیایی مختلفی در صنعت مورد استفاده قرار می‌گیرند. بسیاری از صنایع از نتایج حاصل از ترکیبات شیمیایی نگهداری می‌کنند. در این حالت نگهداری و استفاده از داده‌های شیمیایی موجود یک چالش را بوجود می‌آورد. اگر میزان این داده‌های شیمیایی زیاد شود، به مدلی برای خوشه‌بندی داده‌ها نیاز می‌شود تا بتواند داده‌های ترکیبات مختلف را جداسازی کند. خوشه‌بندی یافتن داده‌های دارای ویژگی‌های نظیرهم، در خوشه‌های مجزا و بدون داشتن اطلاعات اولیه از داده‌های موجود است. در صنایع شیمیایی، امکان آنکه برای تمام داده‌های ترکیبات شیمیایی برچسب گذاری انجام شود، وجود ندارد زیرا هر لحظه ممکن است بوجود بیایند یا تغییر کنند. در این حالت بایستی از خوشه بندی استفاده شود که عمل تقسیم داده‌های شیمیایی به تعدادی از زیر مجموعه‌ها را انجام می‌دهد. از دیدگاه داده‌کاوی تشخیص داده‌های شیمیایی جزء مسائل خوشه‌بندی داده‌ها محسوب می‌شود. با معرفی الگوریتم‌های مناسب در این زمینه و سپس تلاش برای افزایش کارایی و میزان درستی اطلاعات شیمیایی، می‌توان به سمت ایجاد سیستم‌های مکانیزه با قابلیت اعتماد بالا با توانایی کشف الگوهای پیچیده گام برداشت. در اینجا یک سری داده‌های ترکیبات شیمیایی صنایع مختلف جمع‌آوری شده و با کمک یک مدل ترکیبی مناسب عمل خوشه‌بندی انجام می‌شود. روش پیشنهادی یک مدل ترکیبی از نزدیکترین همسایگی با کمک خوشه بندی فازی است. در این مدل داده‌های شیمیایی موجود، تحت یک عملیات پیش‌پردازش قرار می‌گیرند تا داده‌های نامناسب و تهی از سیستم خارج شوند. سپس عمل خوشه بندی با مدل نزدیکترین همسایگی انجام می‌شود. در این مدل ابتدا فاصله میان خوشه‌ها با کمک خوشه‌بندی فازی محاسبه شده و سپس عملیات خوشه‌بندی انجام می‌شود. نتایج تجربی بدست آمده از این مدل ترکیبی نشان‌دهنده بهبود نتایج نسبت به چندین مدل دیگر است که تاکنون بر روی داده‌های شیمیایی عمل خوشه‌بندی را انجام داده‌اند.

کلمات کلیدی: خوشه‌بندی فازی، اطلاعات شیمیایی، دسته‌بندی، الگوریتم نزدیکترین همسایه

^۱گروه شیمی، دانشگاه فنی و حرفه‌ای، تهران، ایران (نویسنده مسئول)

^۲گروه مهندسی کامپیوتر، دانشگاه فنی و حرفه‌ای، تهران، ایران

روشهای کلاسیک بسیاری برای بررسی وجود یا عدم وجود یک ترکیب خاص در یک آنالیز داده شده وجود دارد. همه روشهای موجود می‌توانند نتایج مختلفی را بدست آورند. شیمی تحلیلی در این زمینه بسیار مورد توجه قرار گرفته شده است که عملیات اصلی آن، بررسی داده‌های شیمیایی و خروجی آنالیز داده‌های شیمیایی است. نکته ای که در مورد شیمی تحلیلی وجود دارد این است که، تحلیل همه داده‌ها بدون داشتن یاختر منظم داده‌ای امکان‌پذیر نیست. برای جلوگیری از افزونگی داده‌ها بایستی تمام موارد تهیه‌شده دارای یک دسته‌بندی مشخص باشند. عمل دسته‌بندی بایستی به همراه برچسب‌گذاری داده‌ها باشد که بسیار وقت‌گیر و هزینه‌بر است. پس می‌توان بجای عمل دسته‌بندی از خوشه‌بندی داده‌ها استفاده نمود که یکی از روش‌های داده‌کاوی است [۱، ۲].

به مجموعه‌ای از روش‌های قابل اعمال بر پایگاه داده‌های بزرگ و پیچیده به منظور کشف الگوهای پنهان و جالب توجه نهفته در میان داده‌ها، داده‌کاوی گفته می‌شود. روش‌های داده‌کاوی [۳] تقریباً همیشه به لحاظ محاسباتی پر هزینه هستند. علم میان‌رشته‌ای داده‌کاوی، پیرامون ابزارها، متدولوژی‌ها و تئوری‌هایی است که برای آشکارسازی الگوهای موجود در داده‌ها مورد استفاده قرار می‌گیرند و گامی اساسی در راستای کشف دانش محسوب می‌شود. خوشه‌بندی یکی از شاخه‌های یادگیری بدون نظارت می‌باشد و فرآیند خودکاری [۴] است که در طی آن، نمونه‌ها به دسته‌هایی که اعضای آن مشابه یکدیگر می‌باشند تقسیم می‌شوند که به این دسته‌ها خوشه گفته می‌شود. بنابراین خوشه مجموعه‌ای از اشیاء می‌باشد که در آن اعضای مجموعه با یکدیگر مشابه بوده و با اشیاء موجود در خوشه‌های (مجموعه‌های) دیگر غیر مشابه می‌باشند [۵، ۶].

برای مشابه بودن می‌توان معیارهای مختلفی را در نظر گرفت مثلاً می‌توان معیار فاصله را برای خوشه‌بندی مورد استفاده قرار داد و اشیائی را که به یکدیگر نزدیکتر هستند را بعنوان یک خوشه در نظر گرفت که به این نوع خوشه-بندی، خوشه‌بندی مبتنی بر فاصله نیز گفته می‌شود. در الگوریتم خوشه بندی فازی [۷] FCM باید تعداد و مراکز خوشه‌ها توسط کاربر در ابتدا مشخص شوند. کیفیت این الگوریتم بشدت به تعداد اولیه خوشه‌ها و مکان اولیه مراکز خوشه‌ها بستگی دارد. هدف از خوشه‌بندی فازی استخراج مدل‌های فازی [۸] از داده هاست. در این تحقیق ما به پردازش اطلاعات شیمیایی مبتنی بر خوشه‌بندی فازی و الگوریتم k-means می‌پردازیم. به گونه‌ایی که نتایج الگوریتم k-means میتواند در خوشه‌بندی اطلاعات و خوشه‌بندی فازی اطلاعات بدست آمده تاثیر گذار باشد.

در بخش دوم، معرفی اجمالی از داده‌کاوی ارائه شده است و در بخش سوم عمل خوشه‌بندی شرح داده‌شده است. روش پیشنهادی و ارائه نتایج کاربردی مدل به ترتیب در بخش‌های چهارم و پنجم آورده شده است. در انتها یک نتیجه گیری کلی از مدل پیشنهادی ترکیبی آورده شده است.

۲- داده‌کاوی

داده‌کاوی عبارت است از فرآیند اکتشاف الگو و روندهای منظم و پنهان در داده‌های بزرگ و توزیع شده، با استفاده از مجموعه وسیعی از الگوریتم‌های مبتنی بر علوم ریاضی و آمار. این الگوریتم‌ها معمولاً بروی مقادیر عددی و غیرممتنی

اعمال می‌شوند و برای داده‌های متنی، از الگوریتم‌های متن‌کاوی استفاده می‌شود [۸]. داده‌کاوی از علمی مانند هوش مصنوعی، یادگیری ماشینی، آمار، پژوهش عملیاتی و مدیریت پایگاه‌های داده برای ساخت مدل‌ها و پاسخ به سوالات بهره می‌برد. استخراج و تحلیل اطلاعات سازمان از داده‌های در دسترس توسط کارکنان، فرایندی است که برای سال‌های متمادی انجام شده و وظیفه جدیدی در سازمان‌ها به شمار نمی‌آید. اولین الگوریتم‌های شناسایی روندهای منظم و الگوها در پایگاه داده، از علم آمار و نظریه‌های احتمال نشأت گرفته‌اند. در سال‌های اخیر، با رشد روزافزون قدرت محاسباتی رایانه‌ها و امکان دستیابی به نتایج حاصل از محاسبات پیچیده [۹] در مدت زمان کوتاه، سبب شده است تا الگوریتم‌های پیشرفته ریاضی مورد توجه قرار بگیرند. این الگوریتم‌ها با در نظر گرفتن ابعاد مختلف داده، به پالایش و تحلیل آن پرداخته و الگوهای پیچیده و غیرقابل شناسایی توسط روش‌های قدیمی را استخراج و ارائه می‌کنند. رایانه‌ها کمک کرده‌اند تا فرآیند استخراج، پالایش، پیش پردازش و مدل‌سازی داده‌ها و همچنین اعتبارسنجی یافته‌ها [۹، ۱۰] با دقت بیشتر و سرعتی بی‌نظیر انجام شود. تکنیک‌های داده‌کاوی در پردازش اطلاعات دارای انواع گوناگون می‌باشد که این تکنیک‌ها عبارتند از:

۱-۲ - دسته بندی

متداولترین تکنیک است و یکسری نمونه‌های از پیش تعیین شده را شامل می‌شود که برای توسعه مدل به کار می‌رود که بتواند انواعی از موارد ثبت شده را دسته بندی نماید. این شیوه غالباً از درخت تصمیم‌گیری یا الگوریتم‌های دسته بندی شبکه استفاده می‌کند [۹]، فرآیند شامل یادگیری و رده بندی است. در یادگیری اطلاعات آموزشی با الگوریتم دسته بندی تحلیل می‌شود و اطلاعات برای برآورد دقیق قواعد به کار می‌رود، اگر دقت آن در حد مناسبی باشد می‌توان از آن برای موارد جدید استفاده نمود. [۷].

۲-۲ - خوشه بندی

خوشه بندی به عمل تقسیم جمعیت ناهمگن به تعدادی از زیر مجموعه‌ها یا خوشه‌های همگن گفته می‌شود. نقطه تمایز خوشه بندی از دسته بندی: در دسته بندی [۱۱] بر اساس یک مدل هرکدام از داده‌ها به دست‌های از پیش تعیین شده اختصاص می‌یابد. این دسته‌ها از طریق پژوهش‌های پیشین تعیین گردیده‌اند. لیکن در روش خوشه بندی هیچ دسته از پیش تعیین شده ای وجود ندارد و داده‌ها صرفاً براساس تشابه، گروه بندی می‌شوند و عناوین هرگروه نیز توسط کاربر تعیین می‌گردد [۳].

۳-۲ - رگرسیون گیری

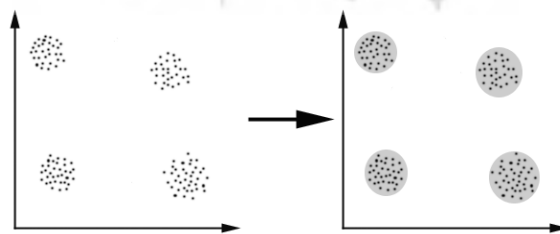
تکنیک رگرسیون گیری را میتوان برای پیش‌بینی پذیرفت. از تحلیل رگرسیون میتوان برای مدل‌سازی روابط یک یا چند متغیر مستقل و وابسته استفاده نمود. در استخراج اطلاعات متغیره ای مستقل ویژگی‌هایی هستند [۱۲] که قبلاً شناخته شده و متغیرهای وابسته مربوط به چیزی هستند که میخواهیم پیش‌بینی کنیم. انواعی از مدل‌ها غالباً برای رگرسیون گیری و دسته بندی به کار می‌روند برای نمونه CART (دسته بندی و درخت رگرسیون گیری) الگوریتم درخت تصمیم گیری را می‌توان برای ایجاد درخت دسته بندی و درخت رگرسیون گیری استفاده نمود [۱۳].

۳ - خوشه بندی

خوشه‌بندی یکی از شاخه‌های یادگیری بدون نظارت می‌باشد و فرآیند خودکاری است که در طی آن، نمونه‌ها به دسته‌هایی که اعضای آن مشابه یکدیگر می‌باشند تقسیم می‌شوند که به این دسته‌ها خوشه گفته می‌شود. بنابراین خوشه مجموعه‌ای از اشیاء می‌باشد [۱۴، ۱۵] که در آن اعضای مجموعه با یکدیگر مشابه بوده و با اشیاء موجود در خوشه-های دیگر غیر مشابه می‌باشد. برای مشابه بودن می‌توان معیارهای مختلفی را در نظر گرفت مثلاً می‌توان معیار فاصله را برای خوشه بندی مورد استفاده قرار داد و اشیائی را که به یکدیگر نزدیکتر هستند را بعنوان یک خوشه در نظر گرفت که به این نوع خوشه‌بندی، خوشه‌بندی مبتنی بر فاصله نیز گفته می‌شود [۱۶]. در کنار مفهوم خوشه‌بندی، مفهوم دسته-بندی وجود دارد. در خوشه‌بندی، هدف یافتن مجموعه متناهی از خوشه‌ها برای توصیف داده‌هاست. هدف در دسته-بندی، ایجاد یک مدل پیشگویی کننده است که این مدل اولاً توانایی دسته‌بندی داده‌های ورودی را داشته باشد و ثانياً بتوان از آن جهت پیشگویی برای تعیین دسته‌ی یک داده که تازه به سیستم اضافه شده، استفاده نمود [۱۷]. خوشه‌بندی در صنایع شیمیایی کاربردهای مختلفی مانند جداسازی نمونه، کشف تقلب [۱۸] و ارائه چارچوب منطقی تریکییات شیمیایی را برعهده دارد.

هدف خوشه‌بندی یافتن خوشه‌های مشابه از اشیاء در بین نمونه‌های ورودی می‌باشد. هیچ معیار مطلقاً برای بهترین خوشه بندی وجود ندارد بلکه این بستگی به مساله و نظر کاربر دارد که باید تصمیم بگیرد [۲۰] که آیا نمونه‌ها بدرستی خوشه بندی شده اند یا خیر. با این حال معیارهای مختلفی برای خوب بودن یک خوشه بندی ارائه شده است که می‌توانند کاربر را برای رسیدن به یک خوشه‌بندی مناسب راهنمایی کند. یکی از مسائل مهم در خوشه‌بندی انتخاب تعداد خوشه‌ها می‌باشد [۹]. در بعضی از الگوریتم‌ها تعداد خوشه‌ها از قبل مشخص شده است و در بعضی دیگر خود الگوریتم تصمیم می‌گیرد که داده‌ها به چند خوشه تقسیم شوند. الگوریتم‌های خوشه‌بندی، با توجه به کاربرد و تنوع مسائل مرتبط با آن، بسیار زیادند و علیرغم پیگیری یک هدف از لحاظ تکنیکی تفاوت‌هایی گاه متضاد دارند [۷] که برخی از آنها عبارتند از:

- ✓ مبتنی بر ناحیه الگوریتم K-means و K-medoids
- ✓ روش سلسله مراتبی
- ✓ مبتنی بر چگالی
- ✓ روش Grid-based
- ✓ روش مبتنی بر مدل

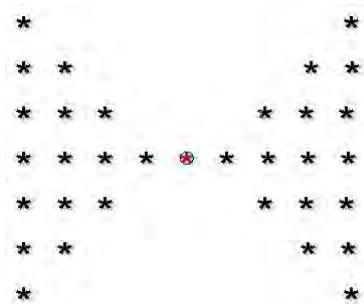


شکل ۱- خوشه‌بندی نمونه‌های ورودی

برای برای درک بهتر خوشه‌بندی فازی و الگوریتم‌های مختلف آن لازم است تا ابتدا با مفهوم مجموعه‌های فازی و تفاوت آنها با مجموعه‌های کلاسیک آشنا شویم. در مجموعه‌های کلاسیک یک عضو از مجموعه مرجع یا عضوی از مجموعه A است یا عضو مجموعه A نیست [۲]. مثلاً مجموعه مرجع اعداد حقیقی را در نظر بگیرید. عدد ۲.۵ عضو مجموعه اعداد صحیح نمی‌باشد حال آنکه عدد ۲ عضو این مجموعه است. به زبان دیگر تعلق عدد ۲.۵ به مجموعه اعداد صحیح ۰ است و تعلق عدد ۲ به این مجموعه ۱ است. در واقع می‌توان برای هر مجموعه یک تابع تعلق تعریف کرد که مقدار این تابع تعلق برای اعضای مجموعه ۱ می‌باشد و برای بقیه ۰. در مجموعه‌های کلاسیک مقدار این تابع تعلق یا ۰ است یا ۱. حال مجموعه انسان‌های جوان و پیر را در نظر بگیرید. سوالی که در اینجا مطرح می‌شود این است که آیا فردی با سن ۲۵ جزء این مجموعه است یا خیر؟ سن ۳۰ چطور؟ ۳۵؟ همانطور که حدس زدید نمی‌توان بطور قطع و یقین مرزی برای انسان‌های جوان و پیر در نظر گرفت. دلیل آن هم این است که اگر فرضاً ۳۵ جوان محسوب شود ۳۶ نیز می‌تواند جوان باشد و همینطور ۳۷ و ۳۸ و غیره. در واقع در اینجا با مفهوم عدم قطعیت مواجه هستیم. ما خودمان نیز از عدم قطعیت در زندگی روزمره بارها استفاده کرده ایم مثلاً هوای سرد، آب داغ و غیره. در واقع تمامی مثالهای بالا مثالهایی از مجموعه‌های فازی می‌باشند [۱].

تفاوت اصلی مجموعه‌های فازی و مجموعه‌های کلاسیک در این است که تابع تعلق مجموعه‌های فازی دو مقداری نیست (۰ یا ۱) بلکه می‌تواند هر مقداری بین ۰ تا ۱ را اختیار کند. حال مجموعه انسان‌های جوان و پیر را در نظر بگیرید اگر ۲۵ سال را سن جوانی در نظر بگیریم می‌توانیم به ۲۵ تعلق ۱ بدهیم و مثلاً به ۳۰ تعلق ۰.۸ و به ۳۵ تعلق ۰.۷۵ و به ۹۰ تعلق ۰.۱ را بدهیم. اگر اعضای یک مجموعه فازی تنها دارای تابع تعلق ۰ و ۱ باشند این مجموعه فازی یک مجموعه کلاسیک خواهد بود. نکته جالب توجه این است که مثلاً سن ۵۰ می‌تواند با تعلق ۰.۵ عضو مجموعه جوان باشد و با تعلق ۰.۵ عضو مجموعه پیر یعنی یک عضو مجموعه مرجع می‌تواند با درجه‌های تعلق مختلف عضو مجموعه‌های فازی تعریف شده روی مجموعه مرجع باشد [۲۱، ۲۲].

در خوشه‌بندی کلاسیک هر نمونه ورودی متعلق به یک و فقط یک خوشه می‌باشد و نمی‌تواند عضو دو خوشه و یا بیشتر باشد. مثلاً در شکل دو هر یک وسایل نقلیه عضو یک خوشه می‌باشد و نمونه ای عضو دو خوشه نیست و به زبان دیگر خوشه‌ها همپوشانی ندارند. حال حالتی را در نظر بگیرید که میزان تشابه یک نمونه با دو خوشه و یا بیشتر یکسان باشد در خوشه‌بندی کلاسیک باید تصمیم‌گیری شود که این نمونه متعلق به کدام خوشه است [۲۲]. تفاوت اصلی خوشه‌بندی کلاسیک و خوشه‌بندی فازی [۲۳] در این است که یک نمونه می‌تواند متعلق به بیش از یک خوشه باشد. برای روشن شدن مطلب شکل زیر را در نظر بگیرید:



شکل ۲ - مجموعه داده پروانه‌ای

اگر نمونه‌های ورودی مطابق شکل فوق باشند مشخص است که می‌توان داده‌ها را به دو خوشه تقسیم کرد اما مشکلی که پیش می‌آید این است که داده مشخص شده در وسط می‌تواند عضو هر دو خوشه باشد بنابراین باید تصمیم گرفت که داده مورد نظر متعلق به کدام خوشه است، خوشه سمت راست یا خوشه سمت چپ [۲۴، ۲۵]. اما اگر از خوشه - بندی فازی استفاده کنیم داده مورد نظر با تعلق ۰.۵ عضو خوشه سمت راست و با تعلق مشابه عضو خوشه سمت چپ است. تفاوت دیگر در این است که مثلاً نمونه‌های ورودی در سمت راست شکل فوق می‌توانند با یک درجه تعلق خیلی کم عضو خوشه سمت چپ نیز باشند که همین موضوع برای نمونه‌های سمت چپ نیز صادق است.

۲-۳ - روش آنالیز مولفه اصلی PCA

PCA یکی از روش‌های فروگاهی داده می‌باشد. این الگوریتم سعی در تحلیل مولفه‌های اصلی داده‌های ما دارد. آنالیز اجزای اصلی یک تکنیک مفید آماری است که کاربرد آن در زمینه‌های از قبیل: تشخیص چهره، فشرده سازی تصویر و یک تکنیک رایج برای شناسایی یک نمونه در داده‌های از بعد بالا است [۲۵]. این تبدیل که با اسامی دیگری چون هتلینگ، کارهان-لو و بردارهای ویژه نیز شناخته می‌شود، تبدیل بهینه در کارهای فشرده سازی و کاهش بعد است [۳] و خطای میانگین مربعات حاصل از فشرده سازی را کمینه می‌کند. هر چند این تبدیل به علت وابسته بودن به داده ورودی، جای خود را در الگوریتم‌های کاربردی و عملی، به تبدیل گسسته کسینوسی داده است اما در صورت کافی بودن داده ورودی می‌تواند تبدیل بهینه را استخراج نماید [۵]. یکی از کاربردهای اصلی PCA در عملیات کاهش ویژگی است. PCA همان‌طور که از نامش پیداست می‌تواند مولفه‌های اصلی را شناسایی کند و به ما کمک می‌کند تا به جای اینکه تمامی ویژگی‌ها را مورد بررسی قرار دهیم، یک سری ویژگی‌هایی را ارزش بیشتری دارند، تحلیل کنیم. در واقع PCA آن ویژگی‌هایی را که ارزش بیشتری فراهم می‌کنند برای ما استخراج می‌کند [۷].

۳-۳ - روش K-means

روش K-Means یکی از روش‌های خوشه‌بندی داده‌ها در داده‌کاوی است. این روش علی‌رغم سادگی آن یک روش پایه برای بسیاری از روش‌های خوشه‌بندی دیگر (مانند خوشه‌بندی فازی) محسوب می‌شود. این روش روشی انحصاری و مسطح محسوب می‌شود. برای این الگوریتم شکل‌های مختلفی بیان شده است [۱۵]. ولی همه آنها دارای روالی تکراری هستند که برای تعدادی ثابت از خوشه‌ها سعی در تخمین موارد زیر دارند: بدست آوردن نقاطی به عنوان مراکز خوشه‌ها این نقاط در واقع همان میانگین نقاط متعلق به هر خوشه هستند [۱۶]. نسبت دادن هر نمونه داده به یک خوشه که آن داده کمترین فاصله تا مرکز آن خوشه را دارا باشد [۹].

در نوع ساده‌ای از این روش ابتدا به تعداد خوشه‌های مورد نیاز نقاطی به صورت تصادفی انتخاب می‌شود. سپس در داده‌ها با توجه با میزان نزدیکی (شباهت) به یکی از این خوشه‌ها نسبت داده می‌شوند و بدین ترتیب خوشه‌های جدیدی حاصل می‌شود. با تکرار همین روال می‌توان در هر تکرار با میانگین‌گیری از داده‌ها مراکز جدیدی برای آنها محاسبه کرد و مجدداً داده‌ها را به خوشه‌های جدید نسبت داد [۱۰]. این روند تا زمانی ادامه پیدا می‌کند که دیگر تغییری در داده‌ها حاصل نشود. تابع زیر به عنوان تابع هدف مطرح است. در الگوریتم K-means ابتدا k عضو (که k تعداد خوشه‌ها است) بصورت تصادفی از میان n عضو به عنوان مراکز خوشه‌ها انتخاب می‌شود [۲، ۷]. سپس $n-k$ عضو باقیمانده به نزدیک‌ترین خوشه تخصیص می‌یابند. بعد از تخصیص همه اعضا مراکز خوشه مجدداً محاسبه

می‌شوند و با توجه به مراکز جدید به خوشه‌ها تخصیص می‌یابند و این کار تا زمانی که مراکز خوشه‌ها ثابت بماند ادامه می‌یابد.

بهترین خوشه‌بندی آن است که مجموع تشابه بین مرکز خوشه و همه اعضای خوشه را حداکثر و مجموع تشابه بین مراکز خوشه‌ها را حداقل کند. برای انتخاب بهترین خوشه [۲۳] ابتدا براساس نظرات خبره و مطالعات قبلی یک محدوده پیشنهادی برای تعداد خوشه‌ها مشخص می‌شود. معمولاً این محدوده بین انتخاب می‌شود. سپس مقدار p برای هر یک از مقادیر k محاسبه می‌شود. مقداری از k که در آن $\rho(k)$ حداکثر شود، به عنوان تعداد بهینه خوشه‌ها [۲۲] انتخاب می‌شود. به این ترتیب می‌توان تعداد خوشه‌ای را انتخاب نمود که به ازای آن فاصله بین مراکز خوشه‌ها و شباهت مراکز خوشه با اعضای درون هر خوشه حداکثر است.

۴- روش پیشنهادی

خود پدیده ی خوشه بندی که یکی دیگر از اهداف داده کاوی می باشد، به فرآیند تقسیم مجموعه ای از داده ها به زیرکلاس هایی با مفهوم خوشه اتلاق می شود. به این ترتیب یک خوشه ، مجموعه داده های مشابه می باشد که همانند یک گروه واحد رفتار می کنند . لازم به ذکر است خوشه بندی همان کلاسه بندی ۳ است، با این تفاوت که کلاس ها از پیش تعریف شده و معین نمی باشند. در خوشه بندی عمل گروه بندی داده ها بدون نظارت انجام می گیرد.



شکل ۳- فلوچارت روش پیشنهادی

داده‌های مورد نظر دارای مشخصات زیر است:

- ✓ نودها به طور تصادفی در محیط پخش شده اند و نودها ساکن و همگن فرض شده اند.
- ✓ ایستگاه پایه آغاز کار در مرکز محیط قرار دارد سپس برای جمع آوری داده ها در مسیری پیدا شده حرکت می کند و در انتهای کار به مرکز محیط بر میگردد.
- ✓ گره ها قادرند توان ارسال خود را با توجه به فاصله خود تا گیرنده مورد نظر، تنظیم کنند.[۹]
- ✓ گره ها همگی دارای انرژی و توانایی یکسان هستند.
- ✓ موقعیت و شناسه تمام گره ها برای ایستگاه پایه معلوم است.
- ✓ هر گره ، از همسایگانش تا شعاع ۲۲ مطلع است.

هر گره برای ارسال l بیت داده به فاصله d از خود به اندازه E_s انرژی مصرف می کند که این از رابطه ۱ بدست می آید:

$$E_s = \begin{cases} lE_{elect} + l\varepsilon_{fs}d^2 & d < d_{co} \\ lE_{elect} + l\varepsilon_{mp}d^4 & d \geq d_{co} \end{cases} \quad (1)$$

همچنین مقدار انرژی برای دریافت این l بیت، در گره گیرنده صرف می شود از رابطه ۲ بدست می آید:

$$E_r = lE_{elect} \quad (2)$$

فرض می شود $\langle U, C, D \rangle$ اطلاعات شیمیایی برای دسته بندی باشند. جدول ۱ نشان دهنده داده های چند برجسبه است.

$$U = \{x_1, x_2, \dots, x_N\}$$

$$C = \{c_1, c_2, \dots, c_N\}$$

$$D = \{d_1, d_2, \dots, d_N\}$$

U نمونه، C مجموعه ویژگی، D مجموعه برجسب هست.

جدول ۱- داده های چند برجسبه

U	c_1	c_2	...	c_M	d_1	d_2	...	d_L
x_1	$x_1^{c_1}$	$x_1^{c_2}$...	$x_1^{c_M}$	$x_1^{d_1}$	$x_1^{d_2}$...	$x_1^{d_L}$
x_2	$x_2^{c_1}$	$x_2^{c_2}$...	$x_2^{c_M}$	$x_2^{d_1}$	$x_2^{d_2}$...	$x_2^{d_L}$
...
x_N	$x_N^{c_1}$	$x_N^{c_2}$...	$x_N^{c_M}$	$x_N^{d_1}$	$x_N^{d_2}$...	$x_N^{d_L}$

در این تحقیق تسک های چند کلاسه به چندین تسک تک برجسبه (باینری) تبدیل می شود و سپس امتیاز میانگین ویژگی ها در تمام تسک های تک برجسبه محاسبه می شود. در جدول ۲ نوع برجسب دهی به کلاس ها تعیین شده است به این صورت که یک نمونه اگر در ارتباط با آن لیبیل باشد برجسب مثبت می گیرد و در غیر این صورت برجسب منفی می گیرد.

(۳)

$$\forall x \in U \quad d_i(x) = \begin{cases} -1 & x \notin d_i \\ +1 & x \in d_i \end{cases}$$

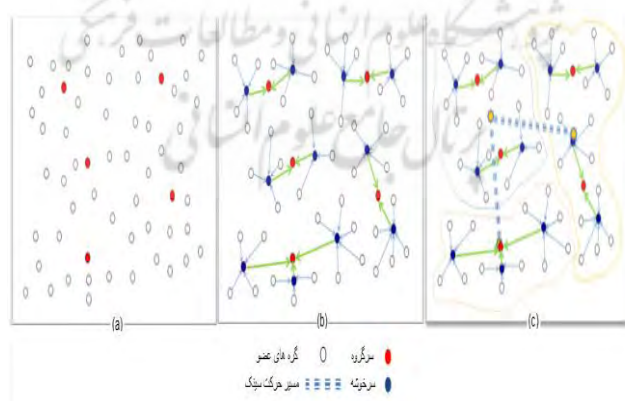
جدول ۲- نحوه تبدیل داده چند برجسبه به داده های باینری تک برجسبه

U	C	d_1	U	C	d_2	U	C	d_L
x_1	...	$x_1^{d_1}$	x_1	...	$x_1^{d_2}$	x_1	...	$x_1^{d_L}$
x_2	...	$x_2^{d_1}$	x_2	...	$x_2^{d_2}$	x_2	...	$x_2^{d_L}$
...
x_N	...	$x_N^{d_1}$	x_N	...	$x_N^{d_2}$	x_N	...	$x_N^{d_L}$

در اینجا برای انتخاب ویژگی از مجموعه‌های فازی استفاده شده است. مجموعه‌های فازی بر پایه مقابله با غیر قابل مشاهده بودن ارائه شده‌اند. موفقیت مجموعه‌های فازی به سه دلیل است:

- ✓ واقعیت پنهان شده در داده‌ها را مشخص می‌کند.
- ✓ هیچ اطلاعات دیگری مثل سطح آستانه یا دانش فرد خبره علاوه بر خود اطلاعات نیاز نیست.
- ✓ دانش حداقلی را برای داده‌ها ارائه می‌دهد.

این مجموعه‌ها فقط با داده‌های کریسپ سروکار دارد. برای حل آن می‌توان مجموعه‌های راف و فازی را با هم ترکیب کرد. چرا که در مجموعه‌های فازی میزان عضویت اعضا به مجموعه‌ها به صورت نسبی (عددی بین ۰ و ۱) بیان می‌شود. روش پیشنهادی از ۳ گام تشکیل شده است. گام اول در شروع هر راند تمام گره‌ها پارامترهای انرژی و مرکزیت را به ماژول فازی موجود در خود ارسال می‌کند. منظور از مرکزیت اینست که یک گره تا چه حد نسبت به همسایگان خود با شعاع I در مرکز واقع شده است. بر اساس خروجی ماژول فازی هر گره، تایمیری برای گره‌ها فعال می‌شود و شروع به شمارش معکوس از مقدار بدست آمده از ماژول فازی می‌کند و گره‌ای که در هر منطقه قابلیت بهتری داشته باشد تایمر آن زودتر به صفر می‌رسد و سپس یک پیام کنترل به شعاع $I+2$ به اطراف ارسال می‌کند و خود را به عنوان سرگروه منطقه معرفی می‌کند شکل ۱(a). گام دوم اینکه در هر منطقه سرگروه‌ها مدیریت خوشه بندی منطقه خود را بر عهده می‌گیرد و در اینجا هم کیفیت گره‌ها بر اساس منطق فازی با پارامترهای انرژی و تعداد همسایگان مشخص می‌شود و گره‌ای که کیفیت بهتری در شعاع I داشته باشد به عنوان سرخوشه انتخاب می‌شود و خود را به سرگروه منطقه معرفی می‌کند و گره‌هایی که کیفیت پایین تری داشته‌اند خود را به نزدیکترین سرخوشه متصل می‌کنند شکل ۱(b). گام سوم تعیین مسیر حرکت سینک می‌باشد که با استفاده الگوریتم k -means سرگروه‌های هر منطقه به k دسته تقسیم بندی می‌شوند در واقع می‌توان گفت منجر به افزای محیط شبکه بر اساس موقعیت سرگروه‌ها می‌شود. برای تعیین مسیر حرکت سینک ابتدا مرکز هر دسته را پیدا کرده و سپس سینک برای جمع‌آوری داده‌ها، کوتاهترین مسیری را که از مراکز دسته‌ها عبور کند را می‌پیماید.



شکل ۴- گام‌های روش پیشنهادی

همانطور که گفته شد در روش پیشنهادی ما تمام گره‌ها با ماژول فازی تست می‌شوند تا در شعاع‌های بخصوص بهترین گره‌ها مشخص شوند.

پارامترهای ورودی منطق فازی برای تعیین گروه‌ها شامل موارد زیر است:

(۱) انرژی گره (energy).

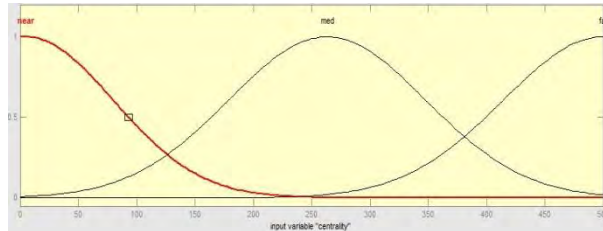
(۲) مرکزیت گره (centrality) .

پارامترهای ورودی منطق فازی برای تعیین سرخوشه ها شامل موارد زیر است

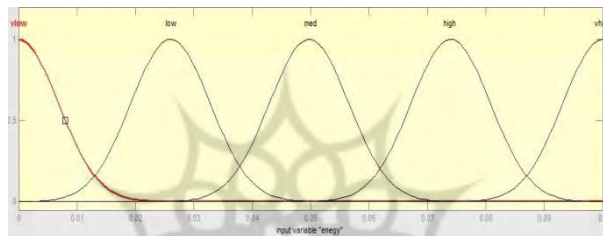
(۱) انرژی گره (energy) .

(۲) چگالی گره (density) .

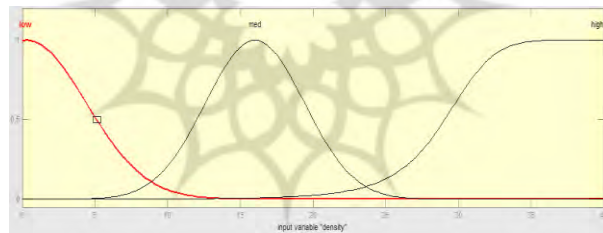
شکل ۵ تا ۸ نداشت ورودی های فازی و خروجی های تابع عضویت را بیان می کند.



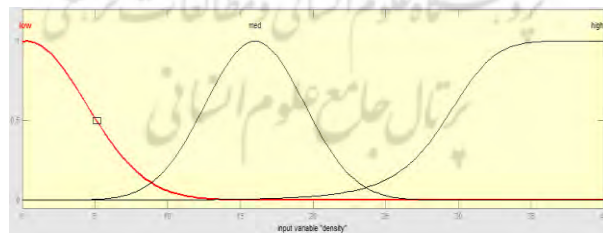
شکل ۵- مجموعه فازی برای متغیر مرکزیت



شکل ۶- مجموعه فازی برای متغیر انرژی



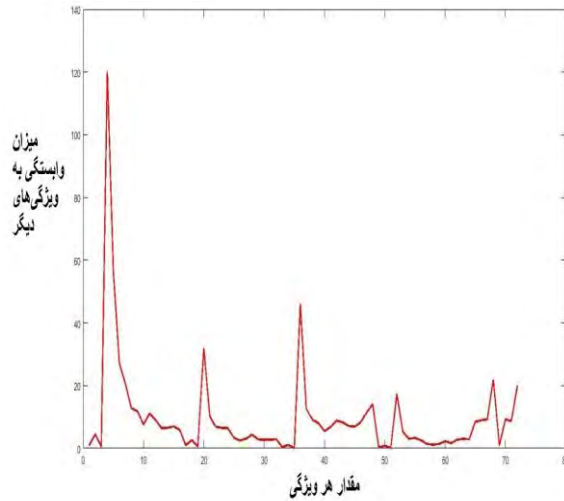
شکل ۷- مجموعه فازی برای متغیر چگالی



شکل ۸- مجموعه فازی برای خروجی تابع عضویت

۵- بررسی نتایج

شکل زیر تابع وابستگی ویژگی ها را طبق رابطه های فصل قبل نشان می دهد.



شکل ۹- تابع وابستگی ویژگی‌ها

با دستور **sort** در نرم افزار متلب، ۱۵ ویژگی از ۷۵ ویژگی که بیشترین وابستگی را دارند انتخاب می‌شود. حال با استفاده از این ویژگی‌ها باید کلاس بندی انجام شود. نوع طبقه بندی کننده‌ای که ما استفاده کرده‌ایم، **K-Means** با ۱۰ نمونه همسایه استفاده می‌گردد. **K-Means** یکی از الگوریتم‌های طبقه بندی می‌باشد. مبنای الگوریتم پیدا کردن تعداد معینی از نزدیکترین عناصر موجود در جامعه آماری به عنصر جدید وارد شده در آن جامعه است که بر اساس آن بتوان نزدیکترین داده موجود به عنصر جدید را از لحاظ ویژگی‌های مختلف پیدا کرد تا عنصر جدید را در همان طبقه ای قرار داد که عناصر نزدیک به آن قرار دارند. **K-Means** یکی از روش‌های غیر پارامتریک برای بدست آوردن تابع توزیع از روی داده‌های توزیع شده می‌باشد. همچنین این روش یکی از متداول‌ترین روش‌ها برای دسته بندی داده‌ها می‌باشد. طبقه بندی کننده **K-Means** بر اساس یادگیری مقایسه‌ای عمل می‌کند. در این جا هر کلاس بصورت جدا با ویژگی‌های انتخاب شده جدا می‌شود و عملکرد کلی کلاس بندی به ازای تمام کلاس‌ها بدست می‌آید.

این روش با سه روش دیگر مقایسه شده است:

- ✓ کلاس بندی با کل ویژگی‌ها انجام می‌گیرد.
- ✓ کلاس بندی با ۱۵ ویژگی که به صورت رندم از کل ویژگی‌ها انتخاب شده است
- ✓ کلاس بندی با ۱۵ ویژگی انتخاب شده با روش آنالیز اجزای اصلی

روش **PCA** تصویر بیشترین واریانس را به دست می‌آورد. عملکرد این روش به این صورت است که ابتدا میانگین داده‌ها بر روی هر بعد را از داده‌ها کم می‌کند و داده‌های جدید با میانگین صفر تولید می‌کند. سپس ماتریس کواریانس داده‌های جدید محاسبه می‌شود. بردارهای ویژه بکه ماتریس کواریانس را می‌توان به عنوان بردار ویژگی‌ها در نظر گرفت زیرا به نوعی پراکندگی داده‌ها را نشان می‌دهد. داده‌های نهایی با ضرب بردارهای ویژگی در داده‌هایی با میانگین صفر به دست می‌آیند. در حقیقت داده‌های نهایی از دوران داده‌های با میانگین صفر به دست می‌آیند. به صورتی که محورهای مختصات آن‌ها بردارهای ویژه ماتریس کواریانس شود. با دقت در مقادیر ویژه ماتریس کواریانس می‌توان بعضی ابعاد که متناظر با مقادیر ویژه‌ای به نسبت کوچک هستند را حذف نمود و به این وسیله ابعاد فضای ورودی را کاهش داد.

جدول زیر نتایج حاصل از مقایسه روش‌ها را نشان می‌دهد. عملکرد طبقه‌بندی‌کننده‌ها براساس روش مبتنی بر وابستگی و مبتنی بر سازگاری برای هر روش محاسبه شده است. [۶]

جدول ۳- مقایسه بین روش‌ها

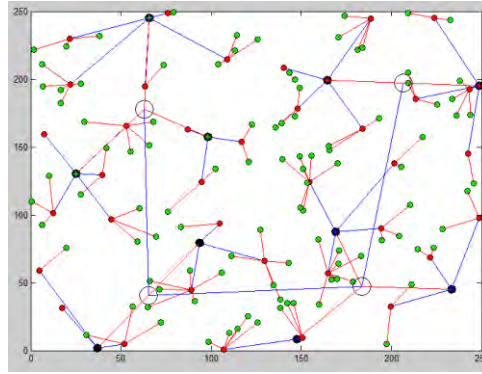
روش انتخاب ویژگی	میزان عملکرد مبتنی بر وابستگی	میزان عملکرد مبتنی بر سازگاری
کل ویژگی‌ها	۸۲۲.۷۰	۲۱۷۰.۰
k-means	۳۲۲.۷۱	۲۱۶۳.۰
Pca	۶۲.۷۱	۲۱۶۵.۰
مجموعه‌های فازی	۹۸۲.۷۱	۲۱۹۵.۰

الگوریتم ارائه شده و پروتکل‌های دیگر با نرم افزار MATLAB شبیه سازی شده‌اند. مقادیر در نظر گرفته شده در این شبیه ازی در جدول زیر آورده شده است.

جدول ۴- پارامترهای شبیه سازی

انرژی اولیه	۱۰.۰ J
مصرف انرژی	۵۰ nJ/bit
اثرگذاری پراکندگی	2 pJ/bit/m^2
مرکزیت	۵ m
اندازه بسته‌ها	۴۰۰۰ bits/nJ/bit/signal

شکل زیر یک نمونه از پاسخ‌های الگوریتم را در اندازه 250×250 نشان می‌دهد.



شکل ۱۰- یک نمونه از پاسخ های الگوریتم پیشنهادی

شکل فوق نشان می دهد که در میزان مصرف انرژی یکسان میزان تاثیرگذاری و وابستگی نود ها در الگوریتم پیشنهادی بیشتر از سایر الگوریتم ها است و ما میتوانیم خوشه بندی اطلاعات را با میزان تاثیرگذاری بیشتری نسبت به ویژگی های انتخابی انجام دهیم. همچنین در شکل زیر میانگین کاهش انرژی روش پیشنهادی را با روش های دیگر را مقایسه می - کند و همانطور که دیده می شود روش پیشنهادی مصرف انرژی بهینه تری را نسبت به روش های دیگر دارد.



شکل ۱۱- مقایسه میانگین انرژی مصرفی در هر راند شبکه

پژوهشگاه علوم انسانی و مطالعات فرهنگی
 رتال جامع علوم انسانی

روشی در داده‌کاوی از بخشی از علم آمار به نام تحلیل اکتشافی داده‌ها استفاده می‌شود که در آن بر کشف اطلاعات نهفته و ناشناخته از درون حجم انبوه داده‌ها تاکید می‌شود. خوشه‌بندی به معنای یافتن داده‌های دارای ویژگی‌های نظیرهم، در خوشه‌های مجزا بوده و کاربرد فراوان در علوم و صنایع مختلف دارد. استفاده از ایده‌ی فازی سازی یکی از روش‌های مورد توجه برای خوشه‌بندی داده‌ها است که این امکان را فراهم می‌کند. خوشه بندی همان کلاسه بندی ۳ است، با این تفاوت که کلاس‌ها از پیش تعریف شده و معین نمی‌باشند. در خوشه بندی عمل گروه بندی داده‌ها بدون نظارت انجام می‌گیرد. K-Means یکی از الگوریتم‌های طبقه بندی می‌باشد. مبنای الگوریتم پیدا کردن تعداد معینی از نزدیکترین عناصر موجود در جامعه آماری به عنصر جدید وارد شده در آن جامعه است که بر اساس آن بتوان نزدیکترین داده موجود به عنصر جدید را از لحاظ ویژگی‌های مختلف پیدا کرد تا عنصر جدید را در همان طبقه‌ای قرار داد که عناصر نزدیک به آن قرار دارند. K-Means یکی از روش‌های غیر پارامتریک برای بدست آوردن تابع توزیع از روی داده‌های توزیع شده می‌باشد. همچنین این روش یکی از متداول‌ترین روش‌ها برای دسته بندی داده‌ها می‌باشد. در این تحقیق با استفاده از الگوریتم نزدیکترین همسایه تسک‌های چند کلاسه به چندین تسک تک برچسبه (باینری) تبدیل و سپس امتیاز میانگین ویژگی‌ها در تمام تسک‌های تک برچسبه محاسبه شد و در نهایت با استفاده از روش پیشنهادی تمام گره‌ها با مازول فازی تست تا در شعاع‌های بخصوص بهترین گره‌ها مشخص شوند. در پایان برای ارزیابی نهایی نمونه‌ی تست با نمونه‌های آموزش داده شده مقایسه شد.



- [۱] ترابی، سحر و مریم خادمی، ۱۳۹۵، منطق فازی در روش‌های خوشه‌بندی، اولین کنفرانس بین‌المللی دستاوردهای نوین پژوهشی در مهندسی برق و کامپیوتر، تهران، کنفدراسیون بین‌المللی مخترعان جهان، دانشگاه جامع علمی کاربردی.
- [۲] همایون، مهدی، سیده اعظم ابوالقاسم پور و محمد رفاهی، ۱۳۹۸، استفاده از سی - مینز فازی جهت خوشه‌بندی خودکار داده‌ها برای الگوریتم ژنتیک چند هدفه، کنفرانس بین‌المللی علوم، مهندسی، تکنولوژی و کسب و کارهای فناوریانه، تهران، شرکت همایش آروین البرز.
- [۳] خسروانیا، آسیه، محمد رحمانی منش، پرویز کشاورزی و سعید مظفری، ۱۳۹۸، ارائه یک روش سطوح همتراز فازی در ناحیه‌بندی خودکار تصاویر پزشکی، اولین کنفرانس سیستم‌ها و فناوری‌های محاسباتی مراقبت از سلامت، بیرجند، دانشگاه بیرجند.
- [۴] مهرآفرید، سمیرا و محمد اکبرپورسکه، ۱۳۹۷، تشخیص تومورهای مغزی با استفاده از الگوریتم خوشه‌بندی فازی و هوش ازدحامی، فصلنامه پژوهش‌های کاربردی در فنی و مهندسی.
- [۵] حمدالهی اسکویی، سعید و مهدی هاشم زاده، ۱۳۹۷، بهبود الگوریتم خوشه‌بندی میانگین فازی با استفاده از الگوریتم سیستم ایمنی مصنوعی، چهارمین کنفرانس ملی محاسبات توزیعی و پردازش داده‌های بزرگ، تبریز، دانشگاه شهید مدنی آذربایجان.
- [۶] هادی نسب، نسیم، جواد محمدزاده و علی سلیمانی، ۱۳۹۷، پیش‌بینی لینک مبتنی بر خوشه‌بندی در شبکه‌های اجتماعی با رویکرد عصبی-فازی تکاملی، سومین کنفرانس ملی فناوری در مهندسی برق و کامپیوتر، سمنان، دانشگاه پیام نور.
- [۷] منصور، محدثه و مرضیه دادور، ۱۳۹۷، ارائه یک مدل برای پیش‌بینی بیماری‌های قلبی با به‌کارگیری خوشه‌بندی ترکیبی و رده‌بند ماشین بردار پشتیبان، پنجمین کنفرانس ملی علوم و مهندسی کامپیوتر و فناوری اطلاعات، بابل، موسسه علمی تحقیقاتی کومه علم آوران دانش.
- [۸] منصور، محدثه و مرضیه دادور، ۱۳۹۷، تشخیص حملات قلبی با استفاده از تکنیک‌های داده‌کاوی فازی و اطلاعات پزشکی، پنجمین کنفرانس ملی علوم و مهندسی کامپیوتر و فناوری اطلاعات، بابل، موسسه علمی تحقیقاتی کومه علم آوران دانش.
- [۹] قانعی رودی، نادر و حمید علیمیرزائی، ۱۳۹۳، بهبود خوشه‌بندی k -means با به‌کارگیری الگوریتم ژنتیک، اولین همایش ملی الکترونیک پیشرفت‌های تکنولوژی در مهندسی برق، الکترونیک و کامپیوتر، بصورت الکترونیک، دانشگاه خیام الکترونیک.
- [۱۰] بنده پی، فرزانه، ۱۳۹۵، بررسی کاربرد الگوریتم K -means در فرآیند خوشه‌بندی داده‌های بزرگ، اولین کنفرانس ملی مهندسی کامپیوتر و فناوری اطلاعات، قم، موسسه مدیریت کنفرانس‌های علمی اندیشوران هزاره سوم.
- [۱۱] مهدی‌زاده، اسماعیل، تیموری، محمد، زارع‌طلب، آرش. (۱۳۹۶). ارائه‌ی یک الگوریتم ترکیبی برای خوشه‌بندی داده‌ها با استفاده از الگوریتم‌های K -means و الکترومغناطیس. مهندسی صنایع و مدیریت.

- [۱۲] Nguyen, T. P. Q., & Kuo, R. J. (۲۰۱۹). Partition-and-merge based fuzzy genetic clustering algorithm for categorical data. *Applied Soft Computing*, ۷۵, ۲۵۴-۲۶۴.
- [۱۳] Ngo, L. T., Dang, T. H., & Pedrycz, W. (۲۰۱۸). Towards interval-valued fuzzy set-based collaborative fuzzy clustering algorithms. *Pattern Recognition*, ۸۱, ۴۰۴-۴۱۶.
- [۱۴] Shang, R., Zhang, W., Li, F., Jiao, L., & Stolkin, R. (۲۰۱۹). Multi-objective artificial immune algorithm for fuzzy clustering based on multiple kernels. *Swarm and Evolutionary Computation*.
- [۱۵] Wu, T., Zhou, Y., Xiao, Y., Needell, D., & Nie, F. (۲۰۱۹). Modified fuzzy clustering with segregated cluster centroids. *Neurocomputing*, ۳۶۱, ۱۰-۱۸.
- [۱۶] Mahata, N., Kahali, S., Adhikari, S. K., & Sing, J. K. (۲۰۱۸). Local contextual information and Gaussian function induced fuzzy clustering algorithm for brain MR image segmentation and intensity inhomogeneity estimation. *Applied Soft Computing*, ۶۸, ۵۸۶-۵۹۶.
- [۱۷] Li, D., Deogun, J., Spaulding, W., & Stuart, B. (۲۰۰۴, June). Towards missing data imputation: a study of fuzzy k-means clustering method. In *International Conference on Rough Sets and Current Trends in Computing* (pp. ۵۷۳-۵۷۹). Springer, Berlin, Heidelberg.
- [۱۸] Amiri, E., Mosallanejad, A., & Sheikahmadi, A. (۲۰۲۱). Copy-Move Forgery Detection by an Optimal Keypoint on SIFT (OKSIFT) Method. *Journal of Computer & Robotics*, ۱۴(۲), ۱۱-۱۹.
- [۱۹] Meng, Y., Liang, J., Cao, F., & He, Y. (۲۰۱۸). A new distance with derivative information for functional k-means clustering algorithm. *Information Sciences*, ۴۶۳, ۱۶۶-۱۸۵.
- [۲۰] Kaur, A., Pal, S. K., & Singh, A. P. (۲۰۱۹). Hybridization of chaos and flower pollination algorithm over K-means for data clustering. *Applied Soft Computing*, ۱۰۵۵۲۳.

[۲۱] Bai, L., Cheng, X., Liang, J., Shen, H., & Guo, Y. (۲۰۱۷). Fast density clustering strategies based on the k-means algorithm. *Pattern Recognition*, ۷۱, ۳۷۵-۳۸۶.

[۲۲] C.Perera, A.Zaslavsky, P.Christan & D.Geargakopoulos, "Contex aware Computing for the Internet of Things: A Survey,". *Communication Survey & Tutorials*, IEEE, vol. ۱۶, pp.۴۱۴-۴۵۴, ۲۰۱۴.

[۲۳] A. Al-Fuqaha, M. Guizani, M. Mohammadi, M. Aledhari and M. Ayyash, "Internet of things: A survey on enabling technologies, protocols and applications". *IEEE Communications Surveys & Tutorials*, vol. ۱۷, pp. ۲۳۴۷ – ۲۳۷۶, ۲۰۱۵.

[۲۴] P. Maia, T. Batista, E. Cavalcante, A. Baffa, F.C. Delicato, P.F. Pires and A. Zomaya, "A web platform for interconnecting body sensors and improving health care". *Procedia Computer Science*, vol. ۴۰, pp. ۱۳۵-۱۴۲, ۲۰۱۴.

[۲۵] J.H. Abawajy, M. Mohammad and M. Hassan, "Federated Internet of Things and Cloud Computing Pervasive Patient Health Monitoring System". *IEEE Communications Magazine*, vol. ۵۵, pp. ۴۸-۵۳, ۲۰۱۷.