

یک روش خوشه‌بندی ترکیبی جدید مبتنی بر خوشه‌بند cmeans فازی با حفظ تنوع در اجماع

فاطمه نجفی^۱، حمید پروین^{۲*}، کمال میرزا^۳، صمد نجاتیان^۴ و سیده وحیده رضایی^۵

^{۱,۲}دانشکده فنی و مهندسی، واحد میبد، دانشگاه آزاد اسلامی، میبد، ایران

^۳دانشکده فنی و مهندسی، واحد مهندسی، دانشگاه آزاد اسلامی، نورآباد ممسنی، فارس، ایران

^۴دانشگاه پژوهشگران جوان و نخبگان، دانشگاه آزاد اسلامی، نورآباد ممسنی، فارس، ایران

^۵دانشگاه پژوهشگران جوان و نخبگان، دانشگاه آزاد اسلامی، یاسوج، کهگیلویه و بویراحمد، ایران

^۶دانشکده مهندسی برق، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

^۷دانشکده علوم، گروه ریاضی، واحد یاسوج، دانشگاه آزاد اسلامی، یاسوج، ایران

چکیده

به علت بدون ناظربودن مسأله خوشه‌بندی، انتخاب یک الگوریتم خاص جهت خوشه‌بندی یک مجموعه ناشناس امری پر خطر و به طور معمول شکست‌خورده است. به خاطر پیچیدگی مسأله و ضعف روش‌های خوشه‌بندی پایه، امروزه بیشتر مطالعات به سمت روش‌های خوشه‌بندی ترکیبی هدایت شده است. در خوشه‌بندی ترکیبی ابتدا چندین خوشه‌بندی پایه تولید و سپس برای تجمعی آن‌ها، از یک تابع توافقی جهت ایجاد یک خوشه‌بندی نهایی استفاده می‌شود که بیشینه شباهت را به خوشه‌بندی‌های پایه داشته باشد. خوشه‌بندی توافقی تولید شده باید با استفاده از بیشترین اجماع و توافق به دست آمده باشد. وروودی تابع یادشده همه خوشه‌بندی‌های پایه و خروجی آن یک خوشه‌بندی به نام خوشه‌بندی توافقی است. در حقیقت روش‌های خوشه‌بندی ترکیبی با این شعار که ترکیب چندین مدل ضعیف بهتر از یک مدل قوی است، به میدان آمده‌اند. با این وجود، این ادعا در صورتی درست است که برخی شرایط همانند تنوع بین اعضای موجود در اجماع و کیفیت آن‌ها رعایت شده باشند. این مقاله یک روش خوشه‌بندی ترکیبی را ارائه داده که از روش خوشه‌بندی پایه ضعیف cmeans فازی به عنوان خوشه‌بند پایه استفاده کرده است. همچنین با اتخاذ برخی تمهیدات، تنوع اجماع را بالا برده است. روش خوشه‌بندی ترکیبی پیشنهادی مزیت الگوریتم خوشه‌بندی cmeans فازی را که سرعت آن است، دارد و همچنین ضعف‌های عمده آن را که عدم قابلیت کشف خوشه‌های غیرکروی و غیریکنواخت است، ندارد. در بخش مطالعات تجربی الگوریتم خوشه‌بندی ترکیبی پیشنهادی با سایر الگوریتم‌های خوشه‌بندی مختلف به روز و قوی بر روی مجموعه داده‌های مختلف آزموده و با یکدیگر مقایسه شده است. نتایج تجربی حاکی از برتری کارایی روش پیشنهادی نسبت به سایر الگوریتم‌های خوشه‌بندی به روز و قوی است.

واژگان کلیدی: یادگیری ترکیبی، خوشه‌بندی ترکیبی، الگوریتم خوشه‌بندی cmeans فازی، اعتبار داده‌ها.

A new ensemble clustering method based on fuzzy cmeans clustering while maintaining diversity in ensemble

Fatemeh Najafi¹, Hamid Parvin^{*2}, Kamal Mirza³, Samad Nejatian⁴ & Vahideh Rezaie⁵

^{1,3}Department of Computer Engineering, Maybod Branch, Islamic Azad University,
Maybod, Iran

²Department of Computer Engineering, Mamasani Branch, Islamic Azad University,
Fars, Iran

^{4,5}Department of Electrical Engineering, Yasooj Branch, Islamic Azad University,
Yasooj, Iran

Abstract

An ensemble clustering has been considered as one of the research approaches in data mining, pattern recognition, machine learning and artificial intelligence over the last decade. In clustering, the combination first produces several bases clustering, and then, for their aggregation, a function is used to create a final cluster that is as similar as possible to all the cluster bundles. The input of this function is all base clusters

* Corresponding author

*نویسنده عهده‌دار مکاتبات

and its output is a clustering called clustering agreement. This function is called an agreement function. Ensemble clustering has been proposed to increase efficiency, strong, reliability and clustering stability. Because of the lack of cluster monitoring, and the inadequacy of general-purpose base clustering algorithms on the other, a new approach called an ensemble clustering has been proposed in which it has been attempted to find an agreed cluster with the highest Consensus and agreement. In fact, ensemble clustering techniques with this slogan, the combination of several poorer models, is better than a strong model. However, this claim is correct if certain conditions (such as the diversity between the members in the consensus and their quality) are met. This article presents an ensemble clustering method. This paper uses the weak clustering method of fuzzy cmeans as a base cluster. Also, by adopting some measures, the diversity of consensus has increased. The proposed hybrid clustering method has the benefits of the clustering algorithm of fuzzy cmeans that has its speed, as well as the major weaknesses of the inability to detect non-spherical and non-uniform clusters. In the experimental results, we have tested the proposed ensemble clustering algorithm with different, up-to-date and robust clustering algorithms on the different data sets. Experimental results indicate the superiority of the proposed ensemble clustering method compared to other clustering algorithms to up-to-date and strong.

Keywords: Ensemble Learning, Ensemble Clustering, Fuzzy Cmeans Clustering Algorithm, Data Validity.

به آهستگی با زمان تغییر کنند. اگر این تغییرات بتواند با یک ردهبندی کننده^۵ به صورت بدون ناظر^۶ رهگیری^۷ شود، عملکرد بهتری می‌تواند به دست آید.
 ۴) می‌توانیم از روش‌های بدون ناظر (خوشه‌بندی) برای پیداکردن و استخراج ویژگی‌ها استفاده کنیم.
 ۵) با خوشه‌بندی می‌توانیم یک دید و بینشی از طبیعت و ساختار داده به دست آوریم که این می‌تواند برای ما با ارزش باشد. کشف زیرداده‌های^۸ مجزا یا شباهت‌های بین الگوها ممکن است به طور چشم‌گیری در روش طراحی ردهبندی کننده به ما پیشنهاد ارایه کند.
 این تنوع در الگوریتم‌های خوشه‌بند خود یک چالش محسوب می‌شود؛ چون هر کدام نقاط ضعف و قوت گوناگونی دارد؛ پس هیچکدام برای همه مجموعه‌داده‌ها مناسب نیست. چالش این است که برای یک مجموعه‌داده در دست، چطور بهترین و مناسب‌ترین روش خوشه‌بندی را انتخاب کرد. برای مثال الگوریتم خوشه‌بند kmeans که یکی از رویکردهای مسطح است، به عنوان یک الگوریتم بسیار سریع و با کارایی به نسبه مناسب شناخته می‌شود [۴]. همچنین این الگوریتم به عنوان یک روشی که (۱) در بهینه محلی گیر می‌افتد، (۲) به انتخاب مراکز اولیه خوشه‌ها حساس است، (۳) همچنین بر این فرض که ساختار خوشه‌ها کروی و یک‌نواخت است، استوار است و (۴) توزیع داده‌ها بر کارایی آن مؤثر است، شناخته می‌شود. این الگوریتم به عنوان یک الگوریتم خوشه‌بند ضعیف، یکی از الگوریتم‌های خوشه‌بند پایه مناسب برای مشارکت در ساخت اجماع محسوب می‌شود. در نقطه مقابل الگوریتم‌های خوشه‌بند محلی

⁵ Classifier

⁶ Unsupervised

⁷ Tracking

⁸ SubClass

۱- مقدمه

یکی از مهم‌ترین مسائل در حوزه‌ی داده‌کاوی و شناسایی الگو، خوشه‌بندی است. خوشه‌بندی نوعی یادگیری بدون ناظر است که به معنی سازمان‌دهی الگوها در چند دسته است؛ به طوری که اعضای هر دسته، از جنبه‌هایی به هم شبیه باشند. در خوشه‌بندی سعی می‌شود الگوها در چند خوشه چنان تقسیم شوند که اعضای هر خوشه به هم شبیه باشند و با اعضای دیگر خوشه‌ها بیشینه تفاوت را داشته باشند. تقسیم‌بندی‌های مختلفی برای روش‌های خوشه‌بندی وجود دارد. این تقسیم‌بندی عبارتند از: سلسله‌مراتبی^۱ در برابر افزایشی^۲، انحصاری^۳ در برابر غیرانحصاری، فازی در برابر غیرفازی، جزئی در برابر کامل و

در واقع خوشه‌بندی داده‌ها یک ابزار ضروری برای یافتن گروه‌ها در داده‌های بدون برچسب است. دست کم پنج دلیل اصلی برای اهمیت خوشه‌بندی وجود دارد:
 ۱) جمع‌آوری و برچسب‌گذاری یک مجموعه بزرگ از الگوهای نمونه می‌تواند بسیار با ارزش باشد.

۲) ممکن است ما به دنبال کردن در جهت معکوس علاقمند باشیم؛ یعنی آموزش با مقدار زیاد داده‌های بدون برچسب و سپس تنها استفاده از ناظر برای برچسب‌گذاری خوشه‌های پیدا شده است. این می‌تواند برای کاربردهای داده‌کاوی بزرگ که محتويات یک پایگاه داده از قبل شناخته شده نیست، مناسب باشد.

۳) در خیلی از کاربردها، مشخصه‌های الگوها مثل ردهبندی^۴ خودکار مواد غذایی با تغییر فصل می‌توانند

¹ Hierarchical

² Partitioning

³ Exclusive

⁴ Classification

خوشه‌بندی ترکیبی هنوز هم به عنوان یک ابزار و هم به عنوان یک زمینه پژوهشی تئوری مورد مطالعه است. یک مقاله مروری در [14]، برای انواع این روش‌ها ارایه شده است. به علت آن که دقت در خوشه‌بندی معنای سرراستی همچون رده‌بندی ندارد؛ بنابراین، مفهوم جایگزینی برای آن ارایه شده است که بیان می‌دارد یک خوشه‌بندی دقیق، خوشه‌بندی‌ای است که به خوشه‌بندی‌های دیگر شکل‌گرفته بر روی داده‌های مورد نظر، بیشترین شباهت را داشته باشد؛ به عبارتی خوشه‌بندی بهتر یعنی خوشه‌بندی پایدارتر. بدليلى مشابه با علت مناسب‌بودن یک مجموع متعدد از رده‌بندها برای رده‌بندی ترکیبی، یک مجموعی از خوشه‌بندی‌ها را مناسب می‌گوییم اگر خوشه‌بندی‌های پایه آن، متعدد باشند. برای متعدد قلمدادشدن یک اجتماعی از خوشه‌بندی‌ها، باید یک الگوریتم خوشه‌بند ضعیف بر روی داده‌ها چندین بار اعمال شود. برای حل این مسأله، الگوریتم خوشه‌بند cmeans فازی بهبودیافته را به عنوان یک خوشه‌بند ضعیف به کار می‌بریم. چهار زیرمسأله در خوشه‌بندی ترکیبی در زیر آورده شده است:

۱. مسأله تشخیص برچسب‌های به نسبه صحیح در خوشه‌بندی: بر خلاف رده‌بندی، هیچ اطلاعات واقعی از برچسب‌ها در خوشه‌بندی وجود ندارد.

۲. مسأله حصول خوشه‌بندی‌های متعددی که توصیف‌گر کل داده‌ها باشد: در یادگیری ترکیبی، در حالی که چندین یادگیرنده ضعیف به عنوان یک یادگیرنده قوی ترکیب می‌شوند، هر چه یادگیرنده‌های پایه بیشتر مکمل هم‌دیگر باشند، یادگیرنده ترکیبی بهتر عمل می‌کند. یعنی هر خوشه‌بندی ضعف بقیه خوشه‌بندی‌ها را پیوшуند؛ بنابراین، برای این منظور ما بایستی چندین خوشه‌بندی مکمل با اعمال الگوریتم خوشه‌بند cmeans فازی بهبودیافته تولید کنیم.

۳. مسأله تشخیص تناسب بین خوشه‌ها: بر خلاف رده‌بندی که در آن هر برچسبی فقط به یک رده دلالت دارد، برچسب‌ها در خوشه‌بندی هیچ معنای واحدی ندارد و تنها هم‌خوشه‌بودن داده‌ها را نشان نمی‌دهد و خوشه‌های هم‌نام در دو خوشه‌بندی گوناگون به هیچ حقیقتی دلالت نمی‌کند؛ بنابراین، پیش از هر کاری در خوشه‌بندی ترکیبی، برچسب خوشه‌بندی‌های مختلف باید بر اساس تناظر باز برچسب‌گذاری شود؛ علاوه براین، حتی دو خوشه از یک خوشه‌بندی یکسان نیز احتمال دارد که به یک خوشه واقعی دلالت کند.

(local) شبیه به kmeans، الگوریتم‌های خوشه‌بند سراسری قرار می‌گیرد که برخلاف داشتن کارایی خوب، از ضعف‌های واضحی چون داشتن پیچیدگی زمانی بالا رنج می‌برند (چراکه ممکن است برای مثال نیازمند محاسبه فاصله بین همه زوج اشیای داده‌ای باشند). برخی از این دسته the الگوریتم‌ها شامل kmeans سراسری (gkmeans) [5]، DBSCAN [6] noise (DBSCAN) و find of density peaks (CFSFDP) [7] خوشه‌بندی طیفی [8-9] است. الگوریتم gkmeans از محاسبه تعدادی از فواصل بین داده‌ها برای حصول به یک خوشه‌بندی بهتر استفاده می‌کند. DBSCAN می‌تواند هر ساختاری از خوشه‌ها را در محیط‌های حتی نویه‌ای با استفاده از تعداد زیادی از محاسبات فاصله‌ای و چگالی تشخیص دهد. خروجی الگوریتم‌های خوشه‌بند سراسری اگرچه قوی و باکیفیت هست؛ اما در مقایسه با الگوریتم‌های خوشه‌بند محلی، هزینه‌های محاسباتی بسیار چشم‌گیری دارد؛ بنابراین، در همین‌واخر به جای پرداختن به ساخت یک الگوریتم خوشه‌بند سراسری قوی، توجه بیشتر به ساخت چارچوب‌هایی شده است که چندین خوشه‌بند ضعیف را یک پارچه کنند. در این راستا "خوشه‌بند ترکیبی" یا "تجمع خوشه‌بندها" [10-11]، برای بهبود استحکام و کیفیت فرایند خوشه‌بندی ارایه شده است.

در یادگیری مبتنی بر اجماع به عنوان یکی از مباحث داغ پژوهشی در بحث داده‌کاوی، شناسایی الگو، یادگیری ماشین و هوش مصنوعی، چندین یادگیرنده ساده (اغلب ضعیف) را برای حل یک مسأله واحد آموزش می‌دهیم. در این روش یادگیری، به جای یادگیری مستقیم داده‌ها به وسیله یک یادگیر قوی (که به طور معمول کند هستند)، سعی در یادگیری یک مجموعه از یادگیرهای ضعیف (که به طور معمول سریع هستند) و تلفیق نتایج آنها با یک روش تابع توافقی (شبیه رأی گیری) دارند [12]. در یادگیری با ناظر به علت وجود برچسب‌های اشیای داده‌ای، ارزیابی هر یک از یادگیرهای ساده سرراست است؛ اما در یادگیری بدون ناظر چنین نیست و بسیار سخت بتوان راه حلی بدون استفاده از اطلاعات جانبی برای ارزیابی اینکه یک الگوریتم خوشه‌بند روی یک مجموعه داده چه ضعف‌ها و نقاط قوتی دارد، یافت [13]. هم‌اکنون، روش‌های متعددی برای خوشه‌بندی ترکیبی به منظور بهبود استحکام و کیفیت نتایج خوشه‌بندی، مطرح شده‌اند. هر خوشه‌بندی موجود در خوشه‌بندی ترکیبی، به عنوان یک یادگیر پایه در نظر گرفته می‌شود.

متنوع تولید کنیم و (۲) این که چه طور از یک اجماع در دست، بهترین نتایج توافقی را تولید کنیم. البته اگرچه کارایی هر یک از این دو بر کارایی دیگری تأثیر دارد، اما این دو به عنوان دو مسئله به طور کامل مستقل شناخته می‌شوند. به همین علت، تاحدودی در مقاله‌های پژوهشی، تنها به یکی از این دو مقوله پرداخته شده و کمتر دیده شده است که هر دو توان در نظر گرفته شوند.

مسئله نخست که مسئله تولید اجماع نیز نام دارد، سعی در تولید یک مجموعه‌ای از خوشبندی‌های پایه معتبر و البته متنوع دارد. این مسئله به کمک روش‌های گوناگونی انجام شده است؛ برای مثال می‌تواند توسط اعمال یک الگوریتم خوشبند پایه ناپایدار بر روی مجموعه‌داده داده شده با تغییر در پارامترهای الگوریتم تولید [15-17] و همچنین می‌تواند توسط اعمال الگوریتم‌های خوشبند پایه گوناگون بر روی مجموعه‌داده داده شده تولید شود [11، 18-19]. راه دیگر می‌تواند با اعمال یک الگوریتم خوشبند پایه بر روی نگاشتهای گوناگون از مجموعه‌داده داده شده خوشبندی‌های پایه معتبر و متنوع را تولید کند [20-28]. در راه بعدی خوشبندی‌های پایه معتبر و متنوع می‌تواند به کمک اعمال یک الگوریتم خوشبند پایه بر روی زیرمجموعه‌های گوناگون که می‌تواند با جای‌گذاری یا بدون جای‌گذاری تولید شده باشد، از مجموعه‌داده داده شده تولید شود [18-19].

برای حل مسئله دوم نیز راه‌کارهای فراوانی مطرح شده‌اند. نخستین راه‌کار، رویکردهای مبتنی بر ماتریس هم‌رخدادی است که در این رویکردها ابتدا تعداد هم‌خوششدن‌های هر زوج داده در یک اجماع را در یک ماتریس به نام ماتریس هم‌رخدادی ذخیره می‌کنیم؛ سپس با درنظر گرفتن این ماتریس به عنوان ماتریس شباهت و اعمال یک روش خوشبندی سلسله‌مراتبی، خوشبندی‌های نهایی توافقی به دست می‌آیند. این رویکرد به طور تقریبی به عنوان مرسوم‌ترین روش قدیمی شناخته می‌شود [29-32]. رویکردی دیگر، رویکرد مبتنی بر برش (ابر) گراف است. در این رویکرد، ابتدا مسئله یافتن خوشبندی توافقی به یک مسئله افزار گراف تبدیل می‌شود؛ سپس به کمک الگوریتم‌های افزار یا برش گراف، خوشبندی‌های نهایی توافقی به دست می‌آیند [10، 33-35]. چهار الگوریتم ترکیبی ابرگراف معروف CSPA، MCLA، HGPA، و HBGF هستند. رویکردی دیگر، رویکرد رأی‌گیری است [21-22، 36، 38-39]. برای این منظور باید ابتدا عمل

۴. مسئله ترکیب نتایج خوشبندی‌های پایه همسان شده: در خوشبندی‌های گوناگون، هر شیء ممکن است برچسب‌های گوناگونی داشته باشد؛ بنابراین ما باید یک برچسب نهایی را به نام برچسب توافقی تعیین کنیم. در یادگیری ترکیبی، در حالی که چندین یادگیرنده ضعیف به عنوان یک یادگیرنده قوی ترکیب می‌شوند، هر چه عمل ترکیب مؤثرتر باشد، یادگیرنده ترکیبی بهتر عمل می‌کند.

این مقاله سعی در مرتفع‌سازی همه زیرمسائل یادشده به کمک تعریف خوشبندی‌های محل‌اعتبر خواهد داشت. در حقیقت این مقاله داده‌های اطراف یک مرکز خوشبندی‌های الگوریتم خوشبند cmeans فازی بهبودیافته را خوشبندی‌های داده‌ای محل‌اعتبر می‌نماید. برای تولید خوشبندی‌های متنوع، از یک استراتژی تکراری تولید خوشبندی‌های ضعیف (به کمک الگوریتم خوشبند cmeans فازی بهبودیافته) بر روی داده‌های ظاهرنشده در خوشبندی‌های محل‌اعتبر قبلی، استفاده می‌شود؛ سپس به کمک یک معیار شباهت بین خوشبندی‌های پایه دیگر نیز می‌تواند این گراف را اندازه‌گیری می‌کنیم. با تشکیل یک گراف وزن دار که رأس‌های آن خوشبندی‌های محل‌اعتبر است و وزن یال‌های آن میزان شباهت بین خوشبندی‌ها است، گام بعدی الگوریتم پیشنهادی این مقاله انجام می‌شود. در گام بعد الگوریتم، کمینه برش گراف برای افزار این گراف به تعدادی (که از پیش تعیین شده است) خوشبندی‌های اعمال می‌شود. در گام نهایی به کمک خروجی این خوشبندی‌ها و میزان اعتبار متوجه و بیشینه‌سازی توافق، خوشبندی‌های نهایی توافقی تولید می‌شوند. گفتنی است که هر خوشبند پایه دیگر نیز می‌تواند به عنوان خوشبندی این ماتریس به دست می‌شود؛ برای مثال می‌توان الگوریتم خوشبند پایه kmeans را به کار برد. همچنین از دیگر روش‌های تابع توافقی مرسوم نیز می‌توان به عنوان تابع توافقی جهت ترکیب نتایج خوشبندی‌های پایه بهره برد.

در ادامه مقاله، بخش دوم به کارهای مرتبط می‌پردازد. در بخش سوم روش پیشنهادی ارایه شده و در بخش چهارم نتایج تجربی ارایه شده و در بخش نهایی، نتیجه‌گیری و کارهای آینده بحث شده است.

۲- کارهای مرتبط

دو مقوله بسیار مهم در خوشبندی ترکیبی وجود دارد: (۱) این که چه طور یک اجماعی از خوشبندی‌های پایه معتبر و



در [48]، یک چارچوب جدید خوشبندی ترکیبی براساس وزن‌گیری در سطح خوش ارائه شده است. مقدار اطمینان این اجماع در مورد یک خوش، به عنوان قابلیت اطمینان آن خوش در نظر گرفته شده است. مقدار قطعیتی که یک اجماع خاص در مورد یک خوش دارد بر اساس میانگین میزان قابلیت اطمینان آن خوش توسط اجماع محاسبه و سپس با انتخاب بهترین خوشها و تعیین وزنی به هر خوش انتخابی براساس قابلیت اطمینان آن، مجموعه خوش‌های نهایی ایجاد می‌شود؛ پس از آن، مقاله به جای ماتریس توافقی سنتی، ماتریس توافقی وزنی در سطح خوش را پیشنهاد می‌کند؛ سپس دو رویکرد اجماع برای تولید افزار توافقی معرفی و مورد استفاده قرار گرفته است.

۳- روش پیشنهادی

در این بخش ابتدا تعاریف و علائم لازم ارائه می‌شود؛ سپس مسئله خوشبندی ترکیبی را تعریف می‌کنیم. در گام بعد الگوریتم پیشنهادی ارائه خواهد شد. بالاخره در گام نهایی تحلیل الگوریتم آورده خواهد شد.

۳-۱- تعاریف و علائم

تمامی علائم مورد استفاده در این مقاله در جدول (۱) ارایه شده‌اند.

مجموعه‌داده: یک مجموعه‌داده، یک مجموعه‌ای از اشیای داده‌ای است که هر شیء داده خود یک بردار عددی (یا بردار ویژگی) است. مجموعه‌داده با X و هر شیء داده با x_i نشان داده می‌شود؛ بدیهی است که $x_i \in X$ ویژگی زام از شیء داده x_i را با x_j نشان می‌دهیم. تعداد ویژگی‌های مجموعه‌داده X را با $|X|$ نشان می‌دهیم.

خوشبندی: یک مجموعه‌ای را از C زیرمجموعه از داده‌ها خوشبندی یا افزار می‌گوییم اگر اجتماع زیرمجموعه‌ها، کل مجموعه‌داده و اشتراک هر زوج از زیرمجموعه‌ها تهی باشد. به هر زیرمجموعه‌ای از یک مجموعه‌داده خوش می‌گوییم. یک خوشبندی را با $\{\pi^1, \pi^2, \dots, \pi^c\} = \pi$ نشان می‌دهیم که π^i نشان‌دهنده خوش i -ام است. بدیهی است که خوش π^i را با C^{π^i} نشان می‌دهیم و زامین ویژگی آن به شکل رابطه (۱) تعریف می‌شود [4].

بازبرچسب‌گذاری انجام شود. عمل باز برچسب‌گذاری به منظور همسان‌سازی برچسب‌های خوشبندی‌های گوناگون برای تطابق است. از رویکردهای مهم دیگر می‌توان به موارد زیر اشاره کرد [40-45]:

۱- رویکرد درنظر گرفتن خوشبندی‌های اولیه به عنوان یک فضای واسط (یا مجموعه‌داده جدید) و خوشبندی این فضای جدید به کمک یک الگوریتم خوشبند پایه شبیه الگوریتم بیشینه‌سازی انتظار [41]

۲- رویکرد استفاده از الگوریتم‌های تکاملی برای یافتن سازگارترین خوشبندی به عنوان خوشبندی توافقی [40]، رویکرد استفاده از الگوریتم خوشبندی kmods برای یافتن خوشبندی توافقی [44-45]

بسیار محتمل است که یک پارتبیشن وجود داشته باشد که با استفاده از یک اندازه‌گیری پایداری به عنوان یک پارتبیشن بد قضاوت شود؛ در حالی که دارای یک (یا بیشتر) خوش با کیفیت بالا است [46]؛ بنابراین، پژوهش‌گران با الهام از ارزیابی پارتبیشن‌ها، اقدامات لازم برای ارزیابی خوش‌ها را تعیین می‌کنند. بسیاری از سنجه‌های پایداری مانند اطلاعات متقابل نرمال برای اعتبارسنجی یک پارتبیشن پیشنهاد شده است. سنجه‌های تعریف شده مبتنی بر اطلاعات متقابل نرمال است. اشکال رویکرد متدائل در این مقاله مورد بحث قرار گرفت و ملاکی برای ارزیابی ارتباط بین یک خوش و یک پارتبیشن ارائه شد که به آن معیار ENMI گویند. معیار ENMI اشکال اندازه‌گیری معمول اطلاعات متقابل نرمال معمولی را جبران می‌کند؛ همچنین، یک روش خوشبندی ترکیبی که مبتنی بر جمع کردن زیرمجموعه‌ای از خوش‌های اولیه است، ارائه شد [46].

برخلاف برخی از تلاش‌ها برای بهبود کیفیت روش‌های خوشبندی، به نظر رسد که پژوهش‌های اندکی به رویه انتخاب در خوشبندی ترکیبی فازی اختصاص یافته است؛ علاوه براین، کیفیت و تنوع محلی دو عامل مهم در انتخاب خوشبندی‌های پایه است. تعداد کمی از مطالعات، این دو عامل را برای انتخاب بهترین خوشبندی‌های پایه فازی در اجماع در نظر گرفته‌اند. در [47] یک چارچوب خوشبندی ترکیبی فازی جدید براساس یک معیار اندازه‌گیری تنوع فازی جدید پیشنهاد شده تا خوش‌های پایه را با بهترین عملکرد پیدا کند. تنوع و کیفیت براساس اطلاعات متقابل عادی فازی بین خوشبندی‌های پایه فازی تعریف شده است.

جدید را معرفی می‌کنیم. بر اساس معیار معرفی شده در این مقاله، شباهت بین دو خوش π^i و π^j را که با (۳) نشان می‌دهیم، به شکل رابطه (۴) تعریف می‌کنیم:

(۴)

$$\text{sim}(\pi^i, \pi^j) = \begin{cases} \frac{\sum_{q=1}^9 T_q(\pi^i, \pi^j) - (\pi^i \cup \pi^j)}{\sqrt{\sum_{w=1}^{|X|} |C_w^{\pi^i} - C_w^{\pi^j}|^2}} & \text{if } \sqrt{\sum_{w=1}^{|X|} |C_w^{\pi^i} - C_w^{\pi^j}|^2} \leq 4\gamma \\ 0 & \text{otherwise.} \end{cases}$$

که (π^i, π^j) از رابطه زیر به دست می‌آید:

(4)

$$T_q(\pi^i, \pi^j) = \left\{ x_k : X \mid \sqrt{\sum_{w=1}^{|X|} |p_{qw}(\pi^i, \pi^j) - x_{kw}|^2} \leq \gamma \right\}$$

که $p_{qw}(\pi^i, \pi^j)$ یک نقطه است که ویژگی w -ام آن بنابر رابطه (۵) تعریف می‌شود:

$$p_{qw}(\pi^i, \pi^j) = \frac{(q) \times C_w^{\pi^i} + (10 - q) \times C_w^{\pi^j}}{10} \quad (5)$$

اجماع خوشبندی: یک مجموعه‌ای از B خوشبندی از داده‌ها را یک اجماع خوشبندی می‌گوییم و آن را با $\Pi = \{\pi_1, \pi_2, \dots, \pi_B\}$ نشان می‌دهیم که π_i نشان‌دهنده خوشبندی i -ام است. بدیهی است که $\pi_i = \{\pi_i^1, \pi_i^2, \dots, \pi_i^{c_i}\}$ و c_i نشان‌دهنده تعداد خوشبندی خوشبندی i -ام است. خوش i -ام از خوشبندی i -ام با π_i^j نشان داده شده است. خوشبندی هدف یا بهترین خوشبندی را با π^* نشان می‌دهیم.

گراف وزن دار متناظر با یک اجماع خوشبندی: یک گراف وزن دار متناظر با یک اجماع خوشبندی Π را با $G(\Pi) = (V(\Pi), E(\Pi))$ و به شکل (۶) تعریف می‌شود. مجموعه رئوس این گراف همان زیرخوشبندی هستند که همه خوشبندی‌های اجماع است؛ یعنی:

$$V(\Pi) = \left\{ r_{\pi_1^1}, r_{\pi_1^2}, \dots, r_{\pi_1^{c_1}}, r_{\pi_2^1}, r_{\pi_2^2}, \dots, r_{\pi_2^{c_2}}, \dots, r_{\pi_B^1}, r_{\pi_B^2}, \dots, r_{\pi_B^{c_B}} \right\}$$

وزن یال‌های بین رئوس این گراف، یا یال‌های بین خوشبندی میزان مشابهت آن‌ها است و مطابق رابطه (۶) به دست می‌آید.

$$C_j^{\pi^i} = \frac{\sum_{k \in \pi^i} x_{kj}}{|\pi^i|} \quad (1)$$

(جدول-۱): شرح علائم

(Table-1): Description of symbols

نماد	شرح
X	یک مجموعه داده
x_i	ام در آتشی داده X
x_{ij}	ام از شی داده x_i ویژگی
$ X $	اندازه یک مجموعه X
$ x_i $	تعداد ویژگی‌های مجموعه داده X
π	یک مجموعه از خوشبندی‌های اولیه یا یک خوشبندی
π^i	ام در خوشبندی نشان‌دهنده خوش π
$c = \pi $	تعداد خوشبندی‌های π
r_{π^i}	زیرخوش معتبر از یک خوش π^i
γ	پارامتر شعاع همسایگی خوش معتبر
C^{π^i}	نقطه مرکز خوش π^i
$C_w^{\pi^i}$	ام از نقطه مرکز خوش w ویژگی
π^*	خوشبندی توافقی
$ \pi^* $	تعداد خوشبندی توافقی
$\text{sim}(\pi^i, \pi^j)$	شباهت بین دو خوش π^i و π^j
$T_q(\pi^i, \pi^j)$	امین خوش فرضی بین دو خوش π^i و π^j
$p_q(\pi^i, \pi^j)$	امین خوش فرضی بین دو خوش π^i و π^j مرکز خوشبندی پایه q از اجتماعی از
Π	خوشبندی پایه Π ایک اجتماعی از
B	اندازه اجماع؛ $ B = \Pi $
π_i	امین خوشبندی در اجماع i
c_i	تعداد خوشبندی خوش π_i ؛ $c_i = \pi_i $
$G(\Pi)$	گراف تعریف شده بر روی اجماع Π
$V(\Pi)$	گره‌های گراف تعریف شده بر روی اجماع Π
$E(\Pi)$	یال‌های گراف تعریف شده بر روی اجماع Π
L	خوشبندی‌های واقعی تعریف شده بر حسب برجسب‌های داده

زیرخوش معتبر از یک خوش: زیرخوش معتبر از یک خوش π^i به شکل r_{π^i} نمایش داده می‌شود و بنابر رابطه (۲) تعریف می‌شود:

$$r_{\pi^i} = \left\{ x_k : \pi^i \mid \sqrt{\sum_{j=1}^{|X|} |C_j^{\pi^i} - x_{kj}|^2} \leq \gamma \right\} - \bigcup_{j=1}^{i-1} r_{\pi^j} \quad (2)$$

که γ یک پارامتر است. گفتگی است که یک زیرخوش می‌تواند به عنوان یک خوش در نظر گرفته شود.

شباهت بین دو خوش: معیارهای فاصله/شباهت بین دو خوشه گوناگونی وجود دارد. در حال حاضر، معیارهای اندازه‌گیری فاصله بین خوشبندی‌های گوناگونی مطرح هستند [49-52] که ما نیز در این مقاله، به فراخور نیاز یک معیار



بیشتر خوشه تولید کند (در حقیقت به اندازه ضریبی از تعداد خوشه‌های که به عنوان ورودی دریافت می‌کند، خوشه تولید می‌کند). الگوریتم ارائه شده در شکل (۲)، یعنی الگوریتم تولید خوشه‌های محلی معتبر، هر بار یک الگوریتم خوشه‌بندی فازی cmeans بهبودیافته را با تعداد خوشه‌های ورودی فراخوانی و به آن اندازه خوشه اولیه محلی معتبر تولید می‌کند؛ پس این الگوریتم یک ضریبی از تعداد خوشه‌های ورودی، خوشه اولیه محلی معتبر تولید می‌کند. این روش یادگیری افزایشی را برای حل مسئله تولید خوشه‌های محلی معتبر استفاده می‌کند. با تبدیل مسئله تولید خوشه‌های محلی معتبر به یک مسئله افزایشی، این روش به تدریج خوشه‌بندی‌های پایه را در هر مرحله تولید می‌کند. در روش یادگیری افزایشی، در هر مرحله به طور تصادفی، c شیء داده را از بین داده‌های هنوز خوشه‌بندی نشده به عنوان مراکز خوشه‌های اولیه انتخاب و از آن در الگوریتم خوشه‌بندی فازی cmeans بهبودیافته استفاده می‌کنیم. این روال تا زمانی که تعداد داده‌های هنوز خوشه‌بندی نشده کمتر از c^2 نشود، ادامه داده می‌شود [53-54]. پیچیدگی زمانی کلی الگوریتم تولید افرازهای اولیه که در شکل (۱) آورده شده، $O(|X| \sum_{i=1}^B c_i)$ است به طوری که در شکل (۳) آورده شده، تعداد تکرارها است.

$$E(v_i, v_j) = sim(v_j, v_i) \quad (6)$$

۳-۳- الگوریتم خوشه‌بندی برای تولید خوشه‌های پایه فازی متنوع

پیچیدگی زمانی الگوریتم cmeans فازی $O(|X|cI)$ است؛ به طوری که I تعداد تکرارها است. گفتنی است که الگوریتم cmeans فازی یک یادگیرنده ضعیف است که عملکرد آن تحت تأثیر عوامل بسیاری قرار دارد. به عنوان مثال، الگوریتم بسیار حساس به مراکز اولیه خوشه است. به طوری که انتخاب مراکز اولیه مختلف، اغلب منجر به نتایج خوشه‌بندی متفاوتی می‌شود؛ علاوه بر این، الگوریتم cmeans فازی تمایل به کشف خوشه‌های کروی با اندازه‌های بنسیبه یکنواخت دارد که برای دیگر توزیع داده‌ها، مناسب نیست؛ بنابراین، ما تلاش خواهیم کرد تا خوشه‌بندی‌های چندگانی تولید شده به وسیله الگوریتم cmeans فازی را برای ایجاد یک نتیجه خوشه‌بندی خوب در مجموعه داده‌ها، با توزیع داده‌های مختلف، به جای یک خوشه‌بند قوی داشته باشیم. بدلیل اهمیت موضوع تنوع در میان خوشه‌بندی‌های پایه، ما در ابتدا در مورد یک الگوریتم خوشه‌بندی محلی معتبر بحث می‌کنیم. تولید خوشه‌بندی‌های پایه بر اساس الگوریتم شکل (۱) انجام می‌پذیرد. در الگوریتم پیشنهادی که مبتنی بر شکل (۱) است، A فراز با کمک الگوریتم خوشه‌بندی پایه که در شکل (۲) آورده شده است، تولید می‌شود. هر فراز تعداد خوشه‌هایی بین c تا $\min(\sqrt{|X|}, 100)$ دارد. بعد از اجرای الگوریتم خوشه‌بندی پایه، ممکن است الگوریتم پایه تعداد

The algorithm for generating a diverse clustering ensemble	
Input:	X, B, c
Output:	Π
01.	$\Pi = \emptyset;$
02.	$\gamma = \text{rand} \times 0.5;$
03.	For $i=1$ to B
04.	$n = X ;$
05.	$c_i = \text{rand}([c, \dots, \min(\sqrt{X}, 100)])$;
06.	$[\pi_i, c_i] = \text{BaseClustering}(X, c_i, \gamma);$
07.	$\Pi = \Pi \cup \{\pi_i\};$
08.	EndFor
09.	Return Π

(شکل-۱): شبکه کد تولید اجماعی از خوشه‌بندی‌های محلی معتبر پایه

(Figure-1): A pseudocode of consensus production of local clustering - valid base

The base clustering algorithm	
Input:	X, c, γ
Output:	π, c
01.	$\pi = \emptyset;$
02.	$TempX = \emptyset;$
03.	$counter = 0;$
04.	$CN = 0;$

```

05. While  $(|X| - |TempX|) \geq c^2$ 
06.    $[\pi^{CN+1}, \dots, \pi^{CN+c}] = \text{modifiedfuzzycmeans}(X, c, \gamma, TempX);$ 
07.    $CN = CN + c;$ 
08.   For  $k=1$  to  $c$ 
09.     If  $(|\pi^{counter \times c+k}| \leq c)$ 
10.        $\pi^{counter \times c+k} = \emptyset; counter = counter + 1;$ 
11.     EndIf
12.      $TempX = TempX \cup \pi^{counter \times c+k};$ 
13.   EndFor
14. EndWhile
15.  $c = CN - counter;$ 
16. Remove all empty clusters from  $\pi$ ;
17. Return  $\pi, c;$ 

```

(شکل-۲): شبکه کد الگوریتم تولید یک خوشبندی محلی معتبر پایه

(Figure-2): Pseudo-code algorithm for producing a valid local-clustering base

The modified fuzzy c-means clustering algorithm

```

Input:  $X, c, \gamma, RDI$  % Removed Data Index
Output:  $[r_{\pi^1}, r_{\pi^2}, \dots, r_{\pi^c}]$ 
01.  $\forall i \in \{1, 2, \dots, c\}: \pi^i = \emptyset;$ 
02.  $counter = 0;$ 
03.  $Y = \emptyset;$ 
04.  $Temp = \{1, \dots, |X|\} - RDI;$ 
05. For each  $j$  in  $Temp$ 
06.    $counter = counter + 1; y_{counter} = x_j;$ 
07.   If ( $counter \leq c$ )
08.      $C_{counter} = y_{counter};$ 
09. EndFor
10. For  $j=1$  to  $MaxIteration$ 
11.   For  $p=1$  to  $|X|$ 
12.     For  $k=1$  to  $c$ 
13.        $dis_{kp} = |x_i - C_k|;$ 
14.        $DIR_{kp} = \begin{cases} 1 & dis_{kp} < \gamma; \\ 0 & \text{o. w.} \end{cases}$ 
15.       If ( $p \in RDI$ )  $DIR_{kp} = -DIR_{kp};$ 
16.     EndFor
17.      $ADIR_k = \sum_{k=1}^c |DIR_{kp}|;$ 
18.      $cln_p = \arg \min_{k \in \{1, \dots, c\}} dis_{kp};$ 
19.     If  $((ADIR_k > 1) \& (p \notin RDI))$ 
20.       For  $p=1$  to  $c$ 
21.          $DIR_{kp} = -DIR_{kp};$ 
22.          $DIR_{k, cln_p} = 1;$ 
23.       EndIf
24.     EndFor
25.     For  $k=1$  to  $c$ 
26.       
$$\Delta C = \frac{\left( DIR \times \left( X - \begin{bmatrix} C_k \\ \vdots \\ C_k \end{bmatrix}_{|X| \times f} \right) \right)}{ADIR_k};$$

27.     EndFor
28.      $C = C + \Delta C;$ 
29.   EndFor
30.   For  $p=1$  to  $|X|$ 
31.     For  $k=1$  to  $c$ 
32.       If ( $DIR_{kp} == 1$ )
33.          $r_{\pi^k} = r_{\pi^k} \cup \{p\};$ 
34.       EndIf
35.     EndFor
36.   EndFor
37. Return  $[r_{\pi^1}, r_{\pi^2}, \dots, r_{\pi^c}]$ 

```

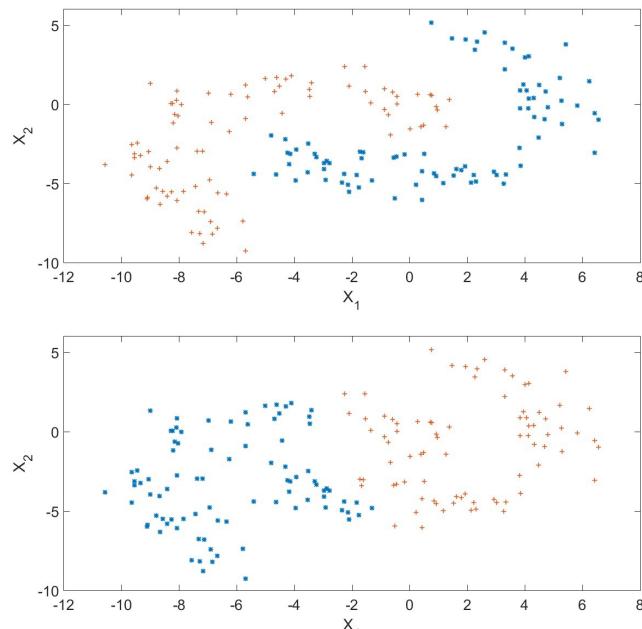
(شکل-۳): شبکه کد الگوریتم خوشبندی فازی cmeans بهبود یافته

(Figure-3): Pseudo-coding of fuzzy clustering algorithms improved cmeans



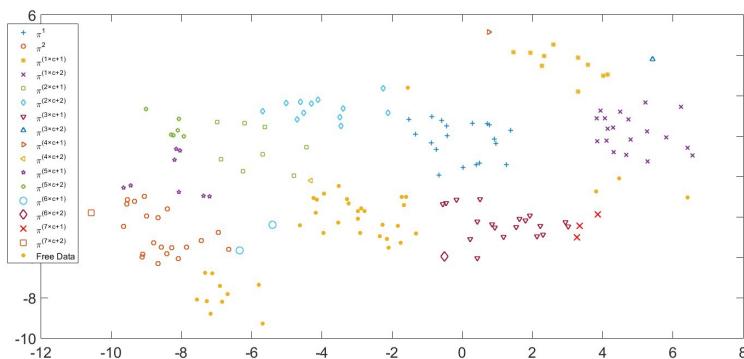
داده‌های Free Data مشخص شده‌اند؛ ولی در اجرای سوم از این الگوریتم (نشان داده شده در دو تصویر پایینی در شکل ۵)، هفت‌بار حلقه موجود در سطر پنج الگوریتم پیشنهادی (ارایه شده در شکل ۲) تکرار می‌شود که منجر به تولید هفت زوج خوشی‌یعنی چهارده خوشی می‌شود. از خروجی‌های این الگوریتم مشاهده می‌کنیم که این خوشی‌های پایه تا اندازه‌ای متنوع بوده، که این حقیقت برای خوشه‌بندی اجتماعی بسیار مفید است. توجه داشته باشید که تعداد خوشی‌های محلی‌معتبر در یک خوشه‌بندی پایه بستگی به پارامتر اندازه همسایگی ۷ دارد. هنگامی‌که این مقدار کم باشد، تعداد خوشی‌های محلی‌معتبر در یک خوشه‌بندی پایه زیاد است و بالعکس؛ بنابراین، اگر این مقدار را به مقادیر کوچک‌تر تنظیم کنیم، نیاز به خوشه‌بندی‌های پایه بیشتری در اجماع داریم. تنظیم این پارامتر بسته به نیاز می‌تواند تغییر کند.

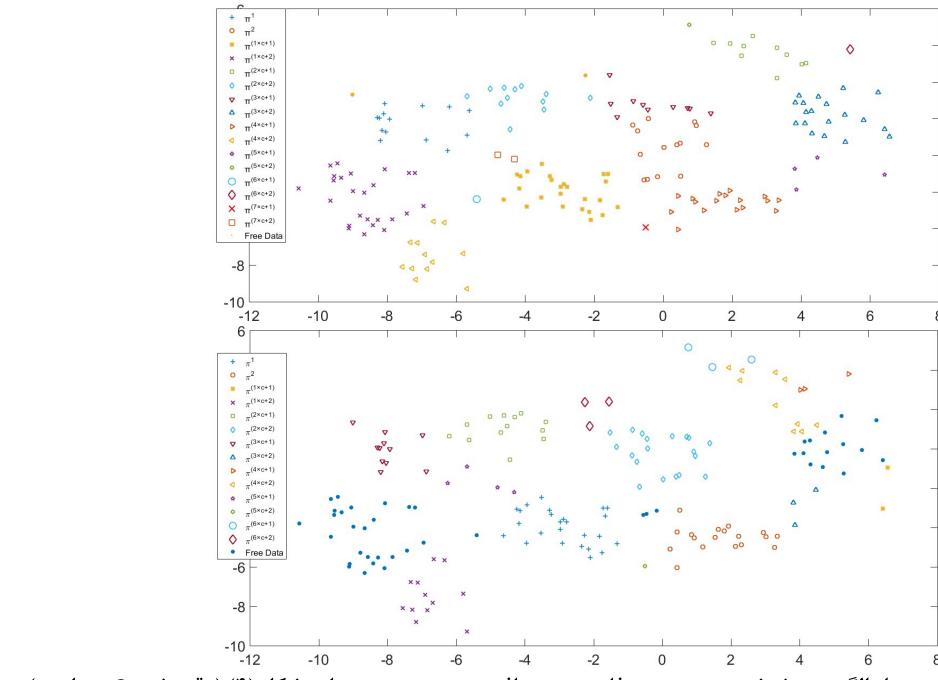
مثال زیر را در شکل (۴) بینید. یک مجموعه داده تصنیعی با دو خوشی و ۱۷۰ داده که در هر خوشی ۸۵ شیء داده موجود است، در شکل (۴-الف) نشان داده شده است. خروجی الگوریتم خوشه‌بندی فازی cmeans بر روی این داده‌ها در شکل ۴-ب قابل مشاهده است. چنان‌که در شکل (۴) مشهود است، الگوریتم خوشه‌بندی فازی cmeans بر روی این داده‌ها ناکارامد است. در شکل (۵) سه خروجی الگوریتم خوشه‌بندی فازی cmeans بهبودیافته (ارایه شده در شکل ۲) بر روی این داده‌گان به تصویر کشیده شده است. در هر یک از دو اجرای نخست و دوم از این الگوریتم (نشان داده شده در سطر پنج بالاتر در شکل ۵)، هشت بار حلقه موجود در سطر پنج الگوریتم پیشنهادی (ارایه شده در شکل ۲) تکرار می‌شود که منجر به تولید هشت زوج خوشی‌یعنی ۱۶ خوشی می‌شود؛ و درنهایت تعدادی داده که بدون خوشی مانده‌اند، به عنوان



(شکل-۴): (الف) یک مجموعه داده و برچسب‌های خوشی واقعی (بالا). (ب) نتیجه اعمال الگوریتم خوشه‌بندی cmeans فازی (پایین).

(Figure-4): (a) A dataset and true cluster labels (high). (B) The result of the fuzzy cmeans clustering algorithm (bottom).



(شکل-۵): نتیجه اعمال سه بار الگوریتم خوشبندی فازی cmeans بر روی مجموعه داده شکل (۴) (دقت شود $c = 2$ است).(Figure-5): The result of applying the fuzzy cmeans clustering algorithm improved on the data set in Figure 4. (Note $c = 2$).

۴- مطالعات تجربی

در بخش جاری، الگوریتم خوشبندی ترکیبی پیشنهادی را بر روی چهار مجموعه داده مصنوعی و پنج مجموعه داده واقعی محک می‌زنیم و بر حسب کارایی و زمان اجرا با الگوریتم‌های به روز مقایسه می‌کنیم.

(جدول-۲): شرح مجموعه داده‌ها: تعداد اشیاء داده ($|X|$), تعداد صفات ($|x_1|$), تعداد خوشبندی (c)

(Table-2): Description of the data set: The number of data objects ($|X|$), the number of traits ($|x_1|$), the number of clusters (c)

c	$ x_1 $	$ X $	مجموعه داده
3	2	600	چرخه-سه‌تایی (شکل ۷-پایین)
2	2	300	موسی-دوتایی (شبیه شکل ۴-بالا، فقط وقتی تعداد داده‌ها ۳۰۰ تا باشد)
7	2	788	توده‌ای-هفت‌تایی (شکل ۷-وسط)
2	2	300	نامتوازن-دوتایی (شکل ۷-بالا)
3	4	150	Iris
3	13	178	Wine
2	30	569	Breast
10	63	5620	Digits
2	39	1048576	KDD-CUP'99

۱-۴- توصیف مجموعه داده‌ها

ارزیابی‌های تجربی بر روی نه مجموعه داده انجام شده است. جدول (۲) جزئیات این مجموعه داده‌ها را نشان می‌دهد.

۳-۳- تابع توافقی پیشنهادی

الگوریتم خوشبندی ترکیبی پیشنهادی در شکل (۶) ارایه شده است. پیچیدگی کلی این الگوریتم برابر با:

$$O\left(|X|(I \sum_{i=1}^B c_i + \sum_{i=1}^B c_i^2 + B)\right)$$

قابل مشاهده است که رابطه پیچیدگی زمانی الگوریتم با تعداد اشیای خطی است و برای یادگیری ترکیبی، بیشتر بودن تعداد B خوشبندی‌های پایه به منزله بهتر بودن نیست؛ بنابراین، به طور کلی ما می‌توانیم فرض کنیم که عبارت:

$$\sum_{i=1}^B c_i \leq \sum_{i=1}^B \min(\sqrt{|X|}, 100) \leq 100B \ll |X|$$

(بهخصوص برای داده‌های خیلی بزرگ) صحیح است؛ به طوری این نشان‌دهنده آن است که الگوریتم پیشنهادی برای مقابله با مجموعه داده‌های با اندازه بسیار بزرگ نیز مناسب است.

۱. π^* خروجی:

۱. ایجاد یک مجموعه خوشبندی پایه به کمک الگوریتم شکل ۲ وقتی تعداد است c خوشبندی‌ها برابر یک

دست آمده از مرحله قبل ۰۲. ساخت یک گراف وزن دار از اجماع به

دست آمده از اجماع ۰۳. افزایش گراف وزن دار به

از طریق خروجی افزایش گراف $4\pi^*$. به دست آوردن یک خوشبندی نهایی

۵. بازگرداندن خوشبندی نهایی

(شکل-۶): شبکه کد الگوریتم خوشبندی ترکیبی پیشنهادی

(Figure-6): Pseudo-code of proposed ensemble clustering algorithm

فصلنامه
پژوهش‌های
دانشجویی





همچنین n_{ij} به شکل زیر تعریف می‌شود:

$$n_{ij} = |\pi_p^i \cap L_j| \quad (9)$$

n نیز تعداد داده‌های مجموعه‌داده را نشان می‌دهد. شاخص ARI تنظیم شده [52] به صورت زیر تعریف می‌شود:

$$(10)$$

$$\text{ARI}(\pi_p, L) = \frac{\binom{n}{2} \times \sum_i \binom{n_{ij}}{2} - \sum_i \binom{n_{ii}}{2} \times \sum_j \binom{n_j}{2}}{0.5 \times \binom{n}{2} \times [\sum_i \binom{n_{ii}}{2} + \sum_j \binom{n_j}{2}] - \sum_i \binom{n_{ii}}{2} \times \sum_j \binom{n_j}{2}}$$

اطلاعات متقابل نرمال شده [56] به صورت زیر تعریف می‌شود:

$$(11)$$

$$\text{NMI}(\pi_p, L) = -\frac{2 \times \sum_i \sum_j n_{ij} \times \log \frac{n_{ij} \times n}{n_{ii} \times n_{jj}}}{\sum_{ij} n_{ii} \times \log \frac{n_{ii}}{n} + \sum_{ij} n_{jj} \times \log \frac{n_{jj}}{n}}$$

هر دو معیار، عددی بین صفر و یک را بر می‌گردانند که هر چه عدد یادشده بزرگ‌تر باشد، بهمنزله این است که افزار نتیجه، نزدیک‌تر به افزار واقعی است.

۴-۳- روش‌های به روز مقایسه شده

برای این که درستی عملکرد الگوریتم پیشنهادی را بررسی کنیم، آن را با الگوریتم‌های خوشه‌بندی ترکیبی زیر مقایسه کردیم:

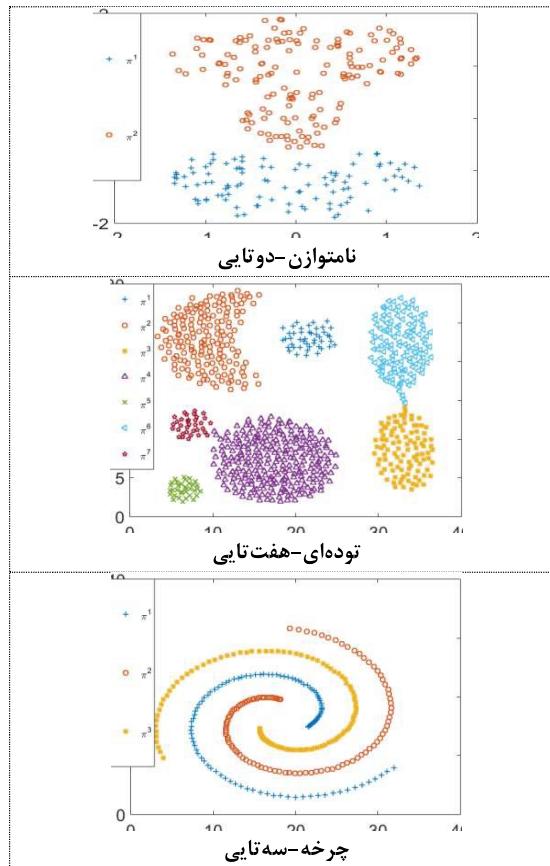
۱. الگوریتم‌های شباهت دوگانه شامل ماتریس هم‌خدادی: روش‌های مبتنی بر ماتریس هم‌خدادی [15] و سه ماتریس شباهت مبتنی بر پیوند WCT، CSM و WTQ [31] از جمله روش‌های این دسته محاسبه می‌شوند. همچنین الگوریتم‌های سلسه‌مراتبی پیوندیکی (SL) و پیوندمتوسط (AL) برای به دست آوردن خوشه‌بندی توافقی نهایی در این مورد استفاده شده است.

۲. الگوریتم‌های مبتنی بر ابرگراف: این روش‌ها شامل MCLA و HGPA [10] هستند.

۳. الگوریتم‌های مبتنی بازبرچسب‌گذاری: این روش‌ها عبارتند از الگوریتم‌های SV و SWV [22].

۴. الگوریتم‌های مبتنی بر فضای واسطه: این روش‌ها شامل الگوریتم پیشینه انتظار (EM) [41] و الگوریتم توافق رأی‌گیری تکراری (IVC) [44] است؛ علاوه‌بر این، روش

توزیع خوشه‌ای این مجموعه‌داده‌های مصنوعی در شکل (7) نشان داده شده است. مجموعه‌داده‌های واقعی برگرفته از مجموعه‌داده‌های UCI هستند [55].



(شکل-7): توزیع سه داده مصنوعی
(Figure-7): Distribution of three artificial datasets

۴-۴- معیارهای ارزیابی

ما به طور گسترده از دو معیار خارجی استفاده کردیم تا شباهت بین نتیجه خوشه‌بندی و تقسیم درست بر روی مجموعه‌داده‌ها را اندازه‌گیری کنیم. با توجه به یک مجموعه‌داده X و دو افزار (یا خوشه‌بندی) از این مجموعه‌اشیا، یعنی $\pi_p = \{\pi_p^1, \pi_p^2, \dots, \pi_p^{c_p}\}$ (که نتیجه یک الگوریتم خوشه‌بندی پایه است) و $L = \{L^1, L^2, \dots, L^c\}$ (که افزار هدف واقعی برای مجموعه‌داده است)، n_{ij} به شکل زیر تعریف می‌کنیم:

$$n_{ij} = |\pi_p^i \cap L_j| \quad (7)$$

همچنین n_i به شکل زیر تعریف می‌شود:

$$n_i = |\pi_p^i| \quad (8)$$



پیشنهادی را با دیگر الگوریتم‌های خوشبند قوی پایه مقایسه خواهیم کرد. الگوریتم‌های خوشبند قوی پایه مورد مقایسه شامل الگوریتم خوشبندی طیفی نرمال (NSC) [9], خوشبندی فضایی مبتنی بر تراکم در شرایط نوفهای (DBSCAN) [6] و خوشبندی جستجوی سریع و پیداکردن قلهای چگالی (CFSFDP) [41] هستند.

۴-۴- تنظیمات پارامترهای الگوریتم‌های

گوناگون

برای روشن‌بودن مقایسات تجربی، تمامی پارامترهای الگوریتم‌های گوناگون در این بخش ارایه شده‌اند تا بتوان نتایج روش‌های گوناگون را در صورت نیاز دوباره بازتولید کرد؛ همچنین برای اطمینان از منصفانه‌بودن نتایج به‌دست‌آمده برای روش‌های گوناگون، نتایج ارایه شده میانگین پنجاه‌بار اجرای مجزای هر الگوریتم خواهد بود. تعداد خوش‌های موجود در هر خوشبندی پایه برابر است با عددی تصادفی بین تعداد واقعی خوش‌ها در هر یک از مجموعه داده‌های مورد نظر تا دست‌کم صد یا چذر تعداد داده‌های آن مجموعه داده؛ همچنین تعداد خوش‌های خوشبندی توافقی برابر تعداد واقعی خوش‌ها در هر یک از مجموعه داده‌های مورد نظر در نظر گرفته می‌شود. همچنین الگوریتم خوشبند cmeans فازی بهبودیافته به عنوان مولد خوشبندی‌های پایه استفاده شده است.

دو روش برای خوشبندی‌های پایه وجود دارد:

(۱) اجرای kmeans به تعداد B بار مجزا، که هر کدام دارای

مقادیر موقعیت اولیه خوش‌های متفاوتی هستند.

(۲) اجرای الگوریتم پیشنهادی (شکل ۱) برای تولید خوشبندی‌های پایه.

برای ارائه نتایج روش‌های مقایسه‌شده، ما بر اساس پیشنهادها نویسنده‌گان، پارامترهای آن‌ها را تعیین می‌کنیم. کیفیت هر یک از الگوریتم‌های خوشبندی مبتنی بر اجماع با توجه به تنظیمات آن اجماع خاص، با میانگین پنجاه اجرا می‌شوند.

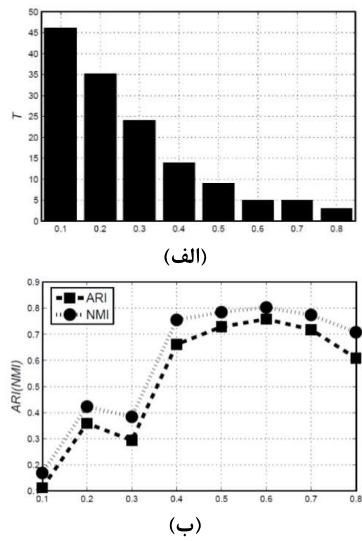
برای الگوریتم NSC، از هسته گاووسی استفاده شده است و پارامتر هسته را از مجموعه $\{0.1, 0.2, \dots, 1.9, 2.0\}$ انتخاب کرده‌ایم. در بین همه این مقادیر، مقداری که به بهترین نتیجه خوشبندی منتج شده، برای مقایسه انتخاب شده است.

الگوریتم‌های CFSFDP و DBSCAN نیز به پارامتر ورودی ϵ نیاز دارند. مقدار ϵ استفاده شده با استفاده از $\bar{d} = \frac{1}{n} \sum_{i=1}^n d(X_i, \bar{X})$ تخمین زده شده که به صورت $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$ است. با این حال، هر یک از این الگوریتم‌ها ممکن است، نیاز به مقادیر مختلف ϵ داشته باشند؛ بنابراین، هر یک از این الگوریتم‌ها با ده مقدار مختلف آزمایش شده او پارامتری که منتج به بهترین نتیجه خوشبندی شده است، برای مقایسه انتخاب شده است. مقادیر مختلف ϵ پژوهش‌شده در این مقاله شامل مجموعه زیر است:

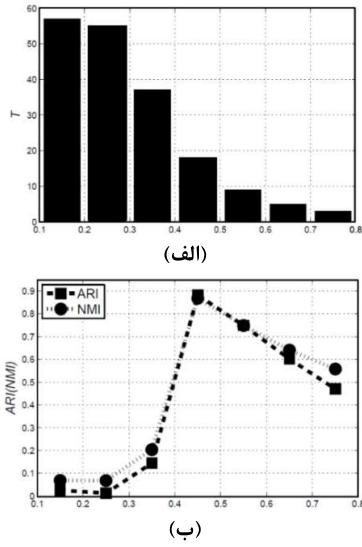
$$\{\bar{d}, \bar{d}/2, \bar{d}/3, \bar{d}/4, \bar{d}/5, \bar{d}/6, \bar{d}/7, \bar{d}/8, \bar{d}/9, \bar{d}/10\}$$

۴-۵- نتایج تجربی

مقایسه با روش‌های ترکیبی دیگر: براساس معیارهای اعتباری ARI و NMI، شکل (۱۰ و ۱۱) نشان‌دهنده مقایسه عملکرد الگوریتم‌های خوشبندی مختلف بر روی مجموعه داده‌های مصنوعی و واقعی هستند. بر طبق شکل (۱۰)، مشاهده می‌کنیم که الگوریتم خوشبندی ترکیبی پیشنهادی دارای دقت خوشبندی بسیار بالا در این مجموعه داده‌های مصنوعی نسبت به دیگر الگوریتم‌های موجود است. نتایج تجربی نشان می‌دهد که الگوریتم ترکیبی پیشنهادشده می‌تواند به طور مؤثر خوش‌های مختلف‌الشكل را کشف و عملکرد الگوریتم خوشبندی پایه را افزایش دهد؛ همچنین در شکل (۱۱) مشاهده می‌کنیم که کارایی الگوریتم خوشبندی ترکیبی پیشنهادی بهتر از الگوریتم‌های دیگر در مجموعه‌داده‌های واقعی است؛ با این حال، بهبود دقت الگوریتم خوشبندی ترکیبی پیشنهادی در مجموعه‌داده‌های واقعی از آن‌چه در مجموعه‌داده‌های مصنوعی است، کمتر است. دلیل اصلی آن، این است که پیچیدگی مجموعه‌داده‌های واقعی بزرگ‌تر از مجموعه‌داده‌های مصنوعی است؛ علاوه‌براین، با توجه به مقایسه‌ها، مشاهده می‌کنیم که بیشتر الگوریتم‌های موجود در سازوکار تولید خوشبندی‌های پایه تصادفی از روش تولید خوشبندی‌های پایه پیشنهادی بهتر عمل می‌کنند؛ زیرا هر خوشبندی پایه تولیدشده بهوسیله روش تولید خوشبندی‌های پایه پیشنهادی، قابل فهم کامل بهوسیله آن‌ها نیستند؛ بنابراین، آن‌ها نمی‌توانند یک نتایج خوشبندی خوب را در روش تولید خوشبندی‌های پایه پیشنهادی بهدست آورند. الگوریتم خوشبندی ترکیبی پیشنهادی عملکرد بهتری در



(شکل-۸): اثر پارامتر s بر روی داده‌های Iris
(Figure-8): Effect of the parameter s on the Iris data



(شکل-۹): اثر پارامتر s بر داده‌های Wine
(Figure-9): Effect of the parameter s on the Wine data

بررسی زمان محاسباتی روش پیشنهادی: در پایان بخش نتایج تجربی، کارایی الگوریتم خوشه‌بندی ترکیبی پیشنهادی را بر روی مجموعه‌داده بسیار بزرگ KDD-CUP'99 محقق می‌زنیم. ما در اینجا دو خوشه (خوشه حمله و غیرحمله) داریم و $\gamma = 0.14$ را تنظیم کردیم. تمامی الگوریتم‌ها بر روی یک دستگاه رایانه‌ای واحد اجرا شده است. جدول (۳) زمان اجرای الگوریتم پیشنهادی با تعداد شی‌عده‌های مختلف را نشان می‌دهد. همان‌طور که مشاهده می‌کنیم، با افزایش تعداد اشیا، تعداد خوشه‌های پایه، یعنی $c_i \sum_{i=1}^B$ نیز افزایش می‌یابد. با توجه به پیچیدگی

روش تولید خوشه‌بندی‌های پایه پیشنهادی در مقایسه با الگوریتم‌های دیگر دارد. توجه داشته باشید که الگوریتم خوشه‌بندی ترکیبی پیشنهادی فقط در روش تولید خوشه‌بندی‌های پایه پیشنهادی اجرا می‌شود، به این دلیل که روش تولید خوشه‌بندی‌های پایه پیشنهادی بخشی از آن است. مشاهده می‌کنیم که الگوریتم خوشه‌بندی ترکیبی پیشنهادی در روش تولید خوشه‌بندی‌های پایه پیشنهادی نیز بر اساس الگوریتم‌های دیگر در روش تولید خوشه‌بندی‌های پایه تصادفی، از نظر ARI و NMI بهتر عمل می‌کند.

مقایسه با الگوریتم‌های پایه مقاوم: شکل (۱۲) نتایج مقایسه الگوریتم خوشه‌بندی ترکیبی پیشنهادی را با سه الگوریتم خوشه‌بندی پایه مقاوم بر اساس مجموعه داده‌های مورد نظر نشان می‌دهد. این شکل، میانگین (mean) و انحراف معیار (standard deviation) اعتبار خوشه‌بندی هر الگوریتم برای این مجموعه داده‌ها نیز بیان شده است. مشاهده می‌کنیم که اعتبار خوشه‌بندی به دست آمده با الگوریتم خوشه‌بندی ترکیبی پیشنهادی برتر یا نزدیک به بهترین نتایج حاصل از سه الگوریتم دیگر است. این آزمایش‌ها به ما می‌گویند که الگوریتم پیشنهادی می‌تواند نتایج مقاوم به دست آمده به وسیله الگوریتم‌های خوشه‌بندی مقاوم را تولید کند.

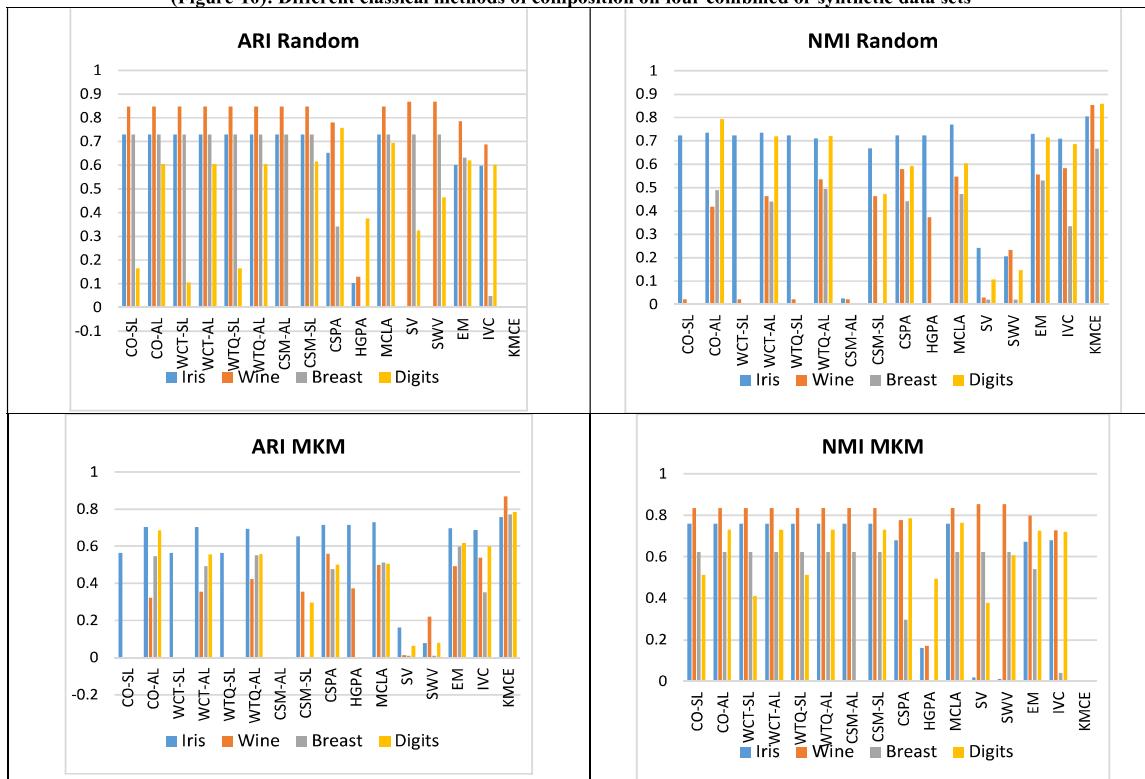
تجزیه و تحلیل پارامتر: نحوه تنظیم پارامتر شعاع همسایگی خوشه معتبر (۷) در الگوریتم خوشه‌بندی ترکیبی پیشنهادی یک چالش مهم محسوب می‌شود. ما در بخش ۳ بحث کردیم که انتخاب این پارامتر بر روی تعداد خوشه‌بندی‌های پایه‌ای تأثیر مستقیم دارد. اینک با بررسی اثر این پارامتر بر روی عملکرد الگوریتم خوشه‌بندی ترکیبی پیشنهادی با آزمایش بر روی داده‌های Iris و Wine در شکل‌های (۸-الف) و (۹-الف) مشاهده می‌کنیم که تعداد خوشه‌های پایه تولید شده به وسیله الگوریتم خوشه‌بندی ترکیبی پیشنهادی با افزایش مقدار این پارامتر کاهش می‌یابد. با این حال، شکل‌های (۸-ب) و (۹-ب) نشان می‌دهند که دقیق خوشه‌بندی افزایش نمی‌یابد؛ بنابراین مقدار این پارامتر بایستی تا حد مشخصی رشد کند. این آزمایش‌ها به ما می‌گویند که تعداد خوشه‌بندی‌های پایه برای به دست آوردن یک نتیجه خوب، داشتن اجماع‌های بزرگ تراز یک آستانه و کوچک‌تر از یک آستانه دیگر هستند؛ بنابراین، ما باید یک مقدار مناسب از این پارامتر را برای کنترل تعداد خوشه‌بندی‌های پایه بر روی هر مجموعه‌داده انتخاب کنیم.

$\sum_{i=1}^B c_i$ چشم‌گیر نیست؛ همچنین از آن جایی که می‌دانیم رابطه هزینه محاسباتی اجرای الگوریتم پیشنهادشده، با تعداد اشیا مجموعه‌داده خطی است، انتظار داریم که الگوریتم پیشنهادی بتواند به سرعت، خوشبندی نهایی را در یک مجموعه‌داده بسیار بزرگ به دست آورد. این مهم به وسیله نتایج ارایه شده در جدول (۳)، تایید می‌شود

زمانی الگوریتم پیشنهادی، رابطه هزینه زمانی با تعداد خوشبندی‌های پایه از نوع درجه دوم است؛ با این حال، با توجه به این واقعیت که می‌دانیم عبارت $\sum_{i=1}^B c_i$ در مقابل $|X|$ در مجموعه‌داده‌های بزرگ قابل چشم‌پوشی است و همچنین این که به آرامی در مقایسه با $\sum_{i=1}^B c_i$ افزایش می‌یابد، افزایش هزینه محاسباتی به خاطر این افزایش در

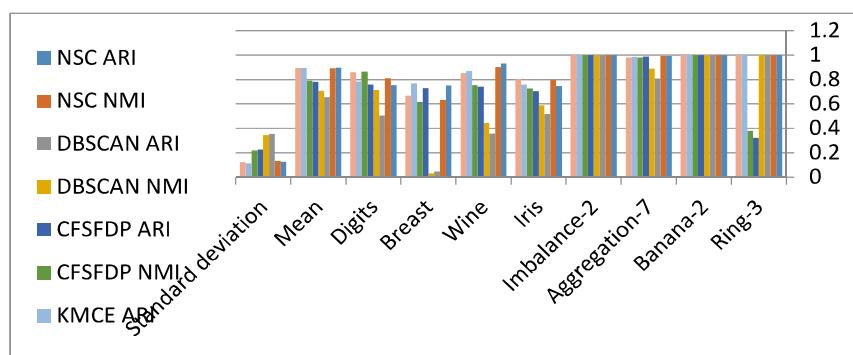
(شکل-۱۰): روش‌های مختلف کلاسیک ترکیبی بر روی چهار مجموعه‌داده‌های ترکیبی یا مصنوعی

(Figure-10): Different classical methods of composition on four combined or synthetic data sets



(شکل-۱۱): روش‌های مختلف کلاسیک ترکیبی بر روی چهار مجموعه‌داده‌ای مصنوعی یا ترکیبی

(Figure-11): Different combinations of classical methods on four synthetic or combined data sets



(شکل-۱۲): مقایسه با الگوریتم‌های خوشبندی "قوی"

(Figure-12): Comparison with Strong Clustering Algorithms



(جدول-۶): نتیجه آزمون آماری KMCE در مقابل سایر روش‌ها بر روی مجموعه داده Breast

(Table-6): The result of KMCE statistical test against other methods on the Breast data set

datasets	methods	نتیجه آزمون آماری KMCE در مقابل:	
		ARI	NMI
Breast	NSC	+	+
	DBSCAN	+	+
	CFSFDP	~	~

(جدول-۷): نتیجه آزمون آماری KMCE در مقابل سایر روش‌ها بر روی مجموعه داده Digits

(Table-7): The result of KMCE statistical test against other methods on the Digits data set

datasets	methods	نتیجه آزمون آماری KMCE در مقابل:	
		ARI	NMI
Digits	NSC	+	+
	DBSCAN	+	+
	CFSFDP	+	-

ارزیابی روش پیشنهادی در مقایسه با روش‌های رقیب

بر روی مجموعه داده‌های پیچیده:

در آزمایش‌های ما از ۱۰ مجموعه داده‌های واقعی استفاده شده است [55]:

(با ۱۵۹۳ نقطه داده، ۲۵۶ ویژگی و ۵ ده طبقه)، Semeion چندین ویژگی (MF) (با دوهزار نقطه داده، ۶۴۹ ویژگی و ۵ ده طبقه)، تقسیم‌بندی تصویر (IS) (با ۲۳۱ نقطه داده، نوزده ویژگی و هفت طبقه)، Forest-Cover-Type (FCT) (با پنج هزار ۳۷۸ نقطه داده، ۵۴ ویژگی و طبقه)، MNIST نقطه داده، ۷۸۴ ویژگی و ۵ ده طبقه)، تشخیص رقمی نوری (با ۵۶۲۰ نقطه داده، ۶۴ ویژگی و ۵ ده طبقه)، ماهواره لنست (LS) (با ۶۴۳۵ نقطه داده، ۳۶ ویژگی و شش طبقه)، SOLET (با ۷۷۹۷ نقطه داده، ۶۱۷ ویژگی، و ۲۶ طبقه)، USPS خوشه‌بندی‌های پایه توسط الگوریتم‌های خوشه‌بندی fuzzy c-means و k-means تولید می‌شوند. گفتنی است که خوشه‌های تولید شده به وسیله الگوریتم خوشه‌بندی fuzzy c-means ابتدا به خوشه‌های واضح تبدیل می‌شوند؛ سپس در ادامه، کل فرآیند برای هر دو الگوریتم، خوشه‌بندی پایه یکسان است. تعداد خوشه‌ها در هر خوشه‌بندی پایه به طور تصادفی از محدوده [۲, ۴K] انتخاب می‌شود [48].

در [57] سعی شده است یکتابع تجمعی، بهنام خوشه‌بندی جمعی مقاوم، مبتنی بر نمونه‌برداری و خوشه‌بندی خوشه‌ای (RCESCC) ارائه شود؛ سپس، یک

(جدول-۳): اثر تعداد داده‌ها بر هزینه محاسباتی الگوریتم

خوشه‌بندی ترکیبی پیشنهاد

(Table-3): Effect of number of data on computational cost of proposed hybrid clustering algorithm

X	$\sum_{i=1}^B c_i$	Time (in sec.)
10K	91	11.23
20K	213	51.29
30K	225	80.11
40K	232	114.06
50K	233	138.91
60K	242	178.71
70K	245	197.62
80K	353	331.02
90K	461	516.96
100K	472	576.58

آزمون آماری معناداری:

برای مقایسه آماری الگوریتم‌ها از آزمون Wilcoxon-Signed-Rank استفاده شده است. در این آزمون از ۹۸ درجه آزادی با سطح معنی دار برابر با $0.05/5$ استفاده کرده‌ایم. بعد از انجام آزمایش نتایج بدین صورت نمایش داده می‌شوند که:

علامت + نشان می‌دهد که KMCE به صورت معناداری از الگوریتم رقیب معین شده بهتر است. علامت - نشان دهنده برتری معنادار الگوریتم رقیب مقابل KMCE بوده و علامت ~ نشان دهنده عدم وجود تفاوت معنادار بین KMCE و رقیب آن است. جداول (۴ تا ۷) نتایج آماری را برای روش پیشنهادی و سایر روش‌ها بر روی مجموعه داده‌های استاندارد مختلف نشان می‌دهند.

(جدول-۴): نتیجه آزمون آماری KMCE در مقابل سایر روش‌ها

بر روی مجموعه داده Iris

(Table-4): The result of KMCE statistical test against other methods on the Iris data set

datasets	methods	نتیجه آزمون آماری KMCE در مقابل:	
		ARI	NMI
Iris	NSC	+	-
	DBSCAN	+	+
	CFSFDP	+	+

(جدول-۵): نتیجه آزمون آماری KMCE در مقابل سایر روش‌ها

بر روی مجموعه داده Wine

(Table-5): The result of KMCE statistical test against other methods on the Wine data set

datasets	methods	نتیجه آزمون آماری KMCE در مقابل:	
		ARI	NMI
Wine	NSC	-	-
	DBSCAN	+	+
	CFSFDP	+	+



ماتریس شباهت خوشخواهی از خوشخواهی فازی به دست می‌آورد؛ پس از آن، با استفاده از یک الگوریتم خوشبندی سلسله‌مراتبی بر روی ماتریس شباهت خوشخواهی، خوشخواهی فازی را تقسیم می‌کند. در مرحله بعدی، الگوریتم RCESCC نقاط داده را به خوشخواهی ادغام‌شده اختصاص می‌دهد.

در [58]، یک رویکرد جدید خوشبندی با استفاده از یک رویکرد وزنی ارائه شده است. در این مقاله روشی برای انجام خوشبندی جمعی با بهره‌برداری از مفهوم عدم اطمینان خوش ارائه شده است. درواقع، هر خوشخواهی عدم وابستگی آن محاسبه شده است. همه برچسب‌های خوشخواهی پیش‌بینی شده موجود در اجماع برای ارزیابی عدم وابستگی خوشخواهی، از معیار مبتنی بر تئوری اطلاعات استفاده می‌کنند. در این مقاله، دو روش مبتنی بر عدم وابستگی یا عدم اطمینان از خوشخواهی برآورد قابلیت اطمینان ارائه شده است. در این مقاله دو رویکرد ارائه شده است: جمع آوری شواهد با وزن خالص و تقسیم‌بندی گراف با وزن خوش. جدول (۸) روش پیشنهادی را با چهار روش جدید که دو روش نخست از نوع خوشبندی k-means و fuzzy c-means به عنوان الگوریتم خوشبندی پایه استفاده می‌کنند، مقایسه کرده است. نتایج مقایسه حاکی از برتری روش پیشنهادی نسبت به چهار روش دیگر است.

(جدول-۸): مقایسه عملکرد الگوریتم خوشبندی پیشنهادی با CLWGC در صورت استفاده خوشبندی fuzzy c- و k-means به عنوان الگوریتم خوشبندی پایه means

(Table-8): Comparison of the performance of the proposed clustering algorithm with CLWGC using k-means and fuzzy c-means clustering as the basic clustering algorithm.

Dataset	CLWGC with k-means	CLWGC with fuzzy c-means	π_{GND}	RCESCC	PROPOSED
Semeion	66.81	66.43	68.59	67.43	69.40
MF	68.89	69.13	69.28	70.15	72.10
IS	67.05	67.26	62.71	67.22	68.54
FCT	23.31	23.38	26.19	30.23	34.30
MNIST	65.26	64.16	66.16	67.02	68.10
ODR	83.12	82.57	85.50	84.04	85.62
LS	63.28	61.42	65.60	63.15	63.48
ISOLET	76.38	75.51	77.71	77.19	78.55
USPS	65.68	65.56	67.12	67.20	68.89
LR	41.47	41.69	47.07	43.42	45.85
Average	62.12	61.71	63.59	63.70	65.48

جدول (۸) نتایج میانگین عملکرد روش‌های مختلف را در بیش از سی اجرای مختلف از نظر NMI مقایسه کرده و نشان داده روش پیشنهادی عملکرد بهتری نسبت به سایر رقبا دارد.

6- References

- [1] J. Han, M. Kamber, *Data Mining: Concepts and Techniques*, Morgan Kaufmann, 2001.
- [2] A.K. Jain, R.C. Dubes, *Algorithms for Clustering Data*, Prentice Hall, 1988.
- [3] A.K. Jain, "Data clustering: 50 years beyond Kmeans", *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651-666, 2010.
- [4] J.B. MacQueen, "Some methods for classification and analysis of multivariate observations". *Proc. of 5-th Berkeley Symposium on Mathematical Statistics and Probability*, Berkeley, University of California Press, vol. 1, pp. 281-297, 1967.
- [5] A. Likas, M. Vlassis, J. Verbeek, "The global fc-means clustering algorithm", *Pattern Recognition*, vol. 35, no. 2, pp. 451-461, 2003.
- [6] M. Ester, H. Kriegel, J. Sander, X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise", In Evangelos Simoudis, Jiawei Han, Usama M.

- Transactions on Computational Biology and Bioinformatics*, vol. 10, no. 3, pp. 657-670, 2013.
- [20] B. Fischer, J. Buhmann, "Bagging for path-based clustering", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 11, pp. 1411-1415, 2003.
- [21] A. Topchy, B. Minaei-Bidgoli, A. Jain, "Adaptive clustering ensembles", *Proc. the 17th International Conference on Pattern Recognition*, 2004.
- [22] Z. Zhou, W. Tang, "Clusterer ensemble", *Knowledge-Based Systems*, vol. 19, no. 1, pp. 77-83, 2006.
- [23] Y. Hong, S. Kwong, H. Wang, Q. Ren, "Resampling-based selective clustering ensembles", *Pattern Recognition Letters*, vol. 41(9), pp. 2742-2756, 2009.
- [24] X. Fern, C. Brodley, "Random projection for high dimensional data clustering: A cluster ensemble approach", *Proc. International Conference on Machine Learning*, 2003.
- [25] P. Zhou, L. Du, L. Shi, H. Wang et al., "Learning a robust consensus matrix for clustering ensemble via kullback-leibler divergence minimization", *Proc. the 25th International Joint Conference on Artificial Intelligence*, 2015.
- [26] Z. Yu, L. Li, J. Liu et al., "Adaptive noise immune cluster ensemble using affinity propagation", *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 19, pp. 3176-3189, 2015.
- [27] F. Gullo, C. Domeniconi, "Metacluster-based projective clustering ensembles", *Machine Learning*, vol. 98, no. 1-2, pp. 1-36, 2013.
- [28] Y. Yang, J. Jiang, "Hybrid Sampling-Based Clustering Ensemble with Global and Local Constitutions", *Ieee Transactions on Neural Networks and Learning Systems*, vol. 27, no. 5, pp. 952-965, 2016.
- [29] A. Fred, A. K. Jain, "Data clustering using evidence accumulation", *Proc. the 16th International Conference on Pattern Recognition*, , 2002, pp. 276-280.
- [30] Y. Yang, K. Chen, "Temporal data clustering via weighted clustering ensemble with different representations", *IEEE Transactions on Knowledge and Data Engineering*, vol. 23, no. 2, pp. 307-320, 2011.
- [31] N. Iam-On, T. Boongoen, S. Garrett, C. Price, "A link-based approach to the cluster ensemble problem", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 12, pp. 2396-2409, 2011.
- [32] N. Iam-On, T. Boongoen, S. Garrett, C. Price, "A Fayyad, *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*, AAAI Press, 1996, pp. 226-231.
- [7] A. Rodriguez, A. Laio, "Clustering by fast search and find of density peaks", *Science*, vol. 344, no. 6191, pp. 1492-1496, 2014.
- [8] J. Shi, J. Malik, "Normalized cuts and image segmentation", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 888-905, 2000.
- [9] A.Y. Ng, M.I. Jordan, Y. Weiss, "On Spectral Clustering: Analysis and an Algorithm", in: T.G. Dietterich, S. Becker, Z. Ghahramani (Eds.)", *Advances in Neural Information Processing Systems*, vol. 14, MIT Press, Cambridge, MA, 2002.
- [10] A. Strehl, J. Ghosh, "Cluster ensembles: a knowledge reuse framework for combining multiple partitions", *Journal on Machine Learning Research*, vol. 3, pp. 583-617, 2002.
- [11] A. Gionis, H. Mannila, P. Tsaparas, "Clustering aggregation", *ACM Transactions on Knowledge Discovery from Data*, vol. 1, no. 1, pp. 1-30, 2007.
- [12] Z. Zhou, *Ensemble Methods: Foundations and Algorithms*, CRC Press, 2012.
- [13] E. Gonzlez, J. Turmo, "Unsupervised ensemble minority clustering", *Machine Learning*, vol. 98, pp. 217-268, 2015.
- [14] N. Iam-On, T. Boongoen, "Comparative study of matrix refinement approaches for ensemble clustering", *Machine Learning*, vol. 98, pp. 269-300, 2015.
- [15] A. Fred, A. Jain, "Combining multiple clusterings using evidence accumulation", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 6, pp. 835-850, 2005.
- [16] L. Kuncheva, D. Vetrov, "Evaluation of stability of kmeans cluster ensembles with respect to random initialization", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 11, pp. 1798-1808, 2006.
- [17] X. Zhang, L. Jiao, F. Liu, L. Bo, M. Gong, "Spectral clustering ensemble applied to SAR image segmentation", *IEEE Transactions on Geoscience and Remote Sensing*, vol. 46, no. 7, pp. 2126-2136, 2008.
- [18] M. Law, A. Topchy, A. Jain, "Multiobjective data clustering", *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2004.
- [19] Z. Yu, H. Chen, J. You, et al, "Hybrid fuzzy cluster ensemble framework for tumor clustering from bio-molecular data", *IEEE/ACM*



- Intelligence Review*, vol. 52, pp. 1311–1340, Springer Nature B.V. 2018, <https://doi.org/10.1007/s10462-018-9642-2>.
- [47] A. Bagherinia, B. Minaei-Bidgoli, M. Hossinzadeh, H. Parvin, “Elite fuzzy clustering ensemble based on clustering diversity and quality measures,” *Springer Science+Business Media, LLC, part of Springer Nature, Applied Intelligence*, vol.49 , PP. 1724–1747, 2019. <https://doi.org/10.1007/s10489-018-1332-x>.
- [48] A. Nazari, A. Dehghan, S Nejatian, V. Rezaie, H. Parvin, “A comprehensive study of clustering ensemble weighting based on cluster quality and diversity,” *Pattern Analysis and Applications*, vol. 22, pp.133–145, 2019.
- [49] S. Guha, R. Rastogi, K. Shim, “Cure: an efficient clustering algorithm for large databases”, *Proc. of the Conference on Management of Data (ACM SIGMOD)*, pp.73-84, 1998.
- [50] P.H.A. Sneath, R.R. Sokal, *Numerical Taxonomy*, Freeman, San Francisco, London, 1973.
- [51] B. King, “Step-wise clustering procedures”, *Journal of the American State Association*, vol. 69, pp. 86-101, 1967.
- [52] G. Karypis, E.-H.S. Han, V. Kumar, “Chameleon: ahierarchical clustering algorithm using dynamic modeling”, *IEEE Computer*, vol. 32, no. 8, pp. 68-75, 1999.
- [53] J.C. Bezdek, N. R. Pal, “Some new indexes of cluster validity”, *IEEE Transactions on Systems Man and Cybernetics Part B*, vol. 28, no. 3, pp. 301-15, 1998.
- [54] N.R. Pal, J.C. Bezdek, “On cluster validity for the fuzzy c-means model”, *IEEE Transactions on Fuzzy Systems*, vol. 3, no. 3, pp. 370-379, 1995.
- [55] UCI Machine Learning Repository, <http://www.ics.uci.edu/~mlearn/ML-Repository.html>, 2016.
- [56] T. S. A. V. W. T. Press, W. H. and B. P. Flannery, *Conditional Entropy and Mutual Information. Numerical Recipes: The Art of Scientific computing (3rd ed)*, New York: Cambridge University Press, 2007.
- [57] F. Rashidi, S. Nejatian, H. Parvin, V. Rezaie, “Diversity based cluster weighting in cluster ensemble: an information theory approach,” *Artificial Intelligence Review*, vol. 52, pp.1341–1368, 2019.
- [58] F. Najafi, H. Parvin, K. Mirzaie, S. Nejatian, V. Rezaie, “Dependability-based cluster weighting in clustering ensemble,” *Stat Anal Data Min: The ASA Data Sci Journal*, vol. 13, pp. 151-164, 2020.
- link-based cluster ensemble approach for categorical data clustering”, *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 3, pp. 413-425, 2012.
- [33] X. Fern, C. Brodley, “Solving cluster ensemble problems by bipartite graph partitioning”, *Proc. of the 21st International Conference on Machine Learning*, 2004.
- [34] D. Huang, J. Lai, C. D. Wang, “Ensemble clustering using factor graph”, *Pattern Recognition*, vol. 50, pp. 131-142, 2016.
- [35] M. Selim, E. Ertunc, “Combining multiple clusterings using similarity graph”, *Pattern Recognition*, vol. 44, no. 3, 694-703, 2011.
- [36] C. Boulis, M. Ostendorf, “Combining multiple clustering systems”, Proc. European Conf. Principles and Practice of Knowledge Discovery in Databases, 2004.
- [37] A. Topchy, B. Minaei-Bidgoli, A. Jain, “Adaptive clustering ensembles”, *Proc. the 17th International Conference on Pattern Recognition*, 2004.
- [38] P. Hore, L. O. Hall, B. Goldgo, “A scalable framework for cluster ensembles”, *Pattern Recognition*, vol. 42, no. 5, 676-688, 2009.
- [39] B. Long, Z. Zhang, P. S. Yu, “Combining multiple clusterings by soft correspondence”, *Proc. the 4th IEEE International Conference on Data Mining*, 2005.
- [40] D. Cristofor, D. Simovici, “Finding median partitions using information theoretical based genetic algorithms”, *J. Universal Computer Science*, vol. 8, no. 2, pp. 153-172, 2002.
- [41] A. Topchy, A. Jain, W. Punch, “Clustering ensembles: Models of consensus and weak partitions”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, 1866-1881, 2005.
- [42] H. Wang, H. Shan, A. Banerjee, “Bayesian cluster ensembles”, *Statistical Analysis and Data Mining*, vol. 4, no. 1, pp. 54-70, 2011.
- [43] Z. He, X. Xu, S. Deng, “A cluster ensemble method for clustering categorical data”, *Information Fusion*, vol. 6, no. 2, pp. 143-151, 2005.
- [44] N. Nguyen, R. Caruana, “Consensus Clusterings”, *Proc. IEEE Intl Conf. Data Mining*, 2007, pp. 607-612.
- [45] Z. Huang, “Extensions to the kmeans algorithm for clustering large data sets with categorical values”, *Data Mining and Knowledge Discovery*, vol. 2, no. 3, pp. 283-304, 1998.
- [46] S. Abbasi, S. Nejatian, H. Parvin, V. Rezaie &K. Bagherifard, “Clustering ensemble selection considering quality and diversity,” *Artificial*





سamed نجاتیان تحصیلات خود را در مقطع کارشناسی در رشته مهندسی برق (الکترونیک) دانشگاه سیستان و بلوچستان در سال ۱۳۸۲ به پایان رساند. ایشان مدرک کارشناسی ارشد خود را در رشته برق (مخابرات)، از دانشگاه مشهد در سال ۱۳۸۶ و مدرک دکترای تخصصی خود را در رشته برق (مخابرات)، در سال ۱۳۹۳ از دانشگاه UMT مالزی اخذ کردند. وی هم‌اکنون عضویت هیأت علمی دانشگاه آزاد واحد یاسوج است.

نشانی رایانامه ایشان عبارت است از:

nejatian@iauyasooj.ac.ir



سیده وحیده رضایی دارای مدرک تحصیلی در مقطع دکترای تخصصی رشته ریاضیات هستند. وی هم‌اکنون عضو هیأت علمی دانشگاه آزاد اسلامی واحد یاسوج است. زمینه‌های پژوهشی ایشان بهینه‌سازی ریاضی، متن کاوی، پردازش سیگنال، داده کاوی و خوشه‌بندی داده‌ها است.

نشانی رایانامه ایشان عبارت است از:

v.rezaie@iauyasooj.ac.ir



فاطمه نجفی تحصیلات خود را در مقطع کارشناسی ارشد در رشته مهندسی رایانه در دانشگاه علوم و تحقیقات و مقطع دکترای تخصصی را در رشته مهندسی رایانه در دانشگاه آزاد اسلامی واحد میبد به پایان رساند. وی هم‌اکنون عضو هیأت علمی دانشگاه آزاد اسلامی واحد ایده است. زمینه‌های پژوهشی ایشان طبقه‌بندی و خوشه‌بندی داده‌ها است.

نشانی رایانامه ایشان عبارت است از:

najafi.un@gmail.com



حمید پروین تحصیلات خود را در مقطع کارشناسی در دانشگاه چمران اهواز به پایان رساند. ایشان مدرک کارشناسی ارشد و دکترا را از دانشگاه علم و صنعت دریافت کردند و پس از آن به عضویت هیأت علمی دانشگاه آزاد واحد نورآباد ممتنی در آمدند. وی هم‌اکنون در چندین واحد دانشگاهی در رشته رایانه مشغول تدریس است. زمینه‌پژوهشی وی مباحثی نظریه‌گوریتم‌های بهینه‌سازی، طبقه‌بندی و خوشه‌بندی داده‌ها است.

نشانی رایانامه ایشان عبارت است از:

parvin@iust.ac.ir



کمال میرزائی مدرک کارشناسی خود را در رشته مهندسی رایانه در سال ۱۳۸۰ از دانشگاه علم و صنعت تهران و مدرک کارشناسی ارشد را در رشته مهندسی کامپیوتر در سال ۱۳۸۳ از دانشگاه اصفهان دریافت کرد. ایشان در سال ۱۳۹۰ موفق به کسب درجه دکترا در رشته مهندسی رایانه از دانشگاه علوم و تحقیقات تهران شد. زمینه‌های پژوهشی مورد علاقه ایشان، علوم شناختی، الگوریتم‌های تکاملی، شبکه‌های پیچیده پویا، محاسبات نرم، داده کاوی پزشکی، پردازش تصویر و شناسایی الگو بوده و در حال حاضر عضو هیأت علمی با مرتبه استادیار در دانشگاه آزاد اسلامی واحد میبد است.

نشانی رایانامه ایشان عبارت است از:

k.mirzaie@maybodiau.ac.ir

