

ترکیب وزن دار خوشبندی ها با هدف افزایش صحت خوشبندی نهایی

صدیقه وحیدی فردوسی و حسین امیرخانی*

گروه مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه قم، قم، ایران

چکیده

با توجه به ماهیت بدون ناظر مسائل خوشبندی و تأثیرگذاری مؤلفه های مختلف از جمله تعداد خوشبندی ها، معیار فاصله و الگوریتم انتخابی، ترکیب خوشبندی ها برای کاهش تأثیر این مؤلفه ها و افزایش صحت خوشبندی نهایی معرفی شده است. در این مقاله، روشی برای ترکیب وزن دار خوشبندی های پایه با وزن دهنی به خوشبندی ها بر اساس روش AD ارائه شده است. روش AD برای برآورد صحت انسان ها در مسائل جمع سپاری از هماهنگی یا تضاد بین آرای آنها استفاده می کند و با پیشنهاد مدلی احتمالاتی، فرآیند برآورد صحت را به کمک یک فرآیند بهینه سازی انجام می دهد. نوآوری اصلی این مقاله، تخمین صحت خوشبندی های پایه با استفاده از روش AD و استفاده از صحت های تخمین زده شده در وزن دهنی به خوشبندی های پایه در فرآیند ترکیب است. نحوه تطبیق مسئله خوشبندی به روش برآورد صحت AD و نحوه استفاده از صحت های برآورده شده در فرآیند ترکیب نهایی خوشبندی ها، از چالش هایی است که در این پژوهش به آنها پرداخته شده است. چهار روش برای تولید خوشبندی های شامل الگوریتم های متفاوت، معیار های فاصله ای متفاوت در اجرای k-means ویژگی های توزیع شده و تعداد خوشبندی های متفاوت برسی شده است. در فرآیند ترکیب، قابلیت وزن دهنی به الگوریتم های خوشبندی ترکیبی CSPA و HGPA اضافه شده است. نتایج روش پیشنهادی روی سیزده مجموعه داده مصنوعی و واقعی مختلف و بر اساس نه معیار ارزیابی متفاوت نشان می دهد که روش ترکیب وزن دار ارائه شده در بیشتر از روش ترکیب خوشبندی بدون وزن عمل می کند که این بهبود برای روش HGPA نسبت به CSPA بیشتر است.

واژگان کلیدی: خوشبندی ترکیبی وزن دار، یادگیری بدون ناظر، CSPA، HGPA

Weighted Ensemble Clustering for Increasing the Accuracy of the Final Clustering

Sedigheh Vahidi Ferdosi & Hossein Amirkhani*

Computer Engineering and Information Technology, University of Qom, Qom, Iran

Abstract

Clustering algorithms are highly dependent on different factors such as the number of clusters, the specific clustering algorithm, and the used distance measure. Inspired from ensemble classification, one approach to reduce the effect of these factors on the final clustering is ensemble clustering. Since weighting the base classifiers has been a successful idea in ensemble classification, in this paper we propose a method to use weighting in the ensemble clustering problem. The accuracies of base clusterings are estimated using an algorithm from crowdsourcing literature called agreement/disagreement method (AD). This method exploits the agreements or disagreements between different labelers for estimating their accuracies. It assumes different labelers have labeled a set of samples, so each two persons have an agreement ratio in their labeled samples. Under some independence assumptions, there is a closed-form formula for the agreement ratio between two labelers based on their accuracies. The AD method estimates the labelers' accuracies by

* Corresponding author

*نویسنده عهده دار مکاتبات

minimizing the difference between the parametric agreement ratio from the closed-form formula and the agreement ratio from the labels provided by labelers. To adapt the AD method to the clustering problem, an agreement between two clusterings are defined as having the same opinion about a pair of samples. This agreement can be as either being in the same cluster or being in different clusters. In other words, if two clusterings agree that two samples should be in the same or different clusters, this is considered as an agreement. Then, an optimization problem is solved to obtain the base clusterings' accuracies such that the difference between their available agreement ratios and the expected agreements based on their accuracies is minimized. To generate the base clusterings, we use four different settings including different clustering algorithms, different distance measures, distributed features, and different number of clusters. The used clustering algorithms are mean shift, k-means, mini-batch k-means, affinity propagation, DBSCAN, spectral, BIRCH, and agglomerative clustering with average and ward metrics. For distance measures, we use correlation, city block, cosine, and Euclidean measures. In distributed features setting, the k-means algorithm is performed for 40%, 50%, ..., and 100% of randomly selected features. Finally, for different number of clusters, we run the k-means algorithm by k equals to 2 and also 50%, 75%, 100%, 150%, and 200% of true number of clusters. We add the estimated weights by the AD algorithm to two famous ensemble clustering methods, i.e., Cluster-based Similarity Partitioning Algorithm (CSPA) and Hyper Graph Partitioning Algorithm (HGPA). In CSPA, the similarity matrix is computed by taking a weighted average of the opinions of different clusterings. In HGPA, we propose to weight the hyperedges by different values such as the estimated clustering accuracies, size of clusters, and the silhouette of clusterings. The experiments are performed on 13 real and artificial datasets. The reported evaluation measures include adjusted rand index, Fowlkes-Mallows, mutual index, adjusted mutual index, normalized mutual index, homogeneity, completeness, v-measure, and purity. The results show that in the majority of cases, the proposed weighted-based method outperforms the unweighted ensemble clustering. In addition, the weighting is more effective in improving the HGPA algorithm than CSPA. For different weighting methods proposed for HGPA algorithm, the best average results are obtained when we use the accuracies estimated by the AD method to weight the hyperedges, and the worst results are obtained when using the normalized silhouette measure for weighting. Finally, among different methods for generating base clusterings, the best results in weighted HGPA are obtained when we use different clustering algorithms to come up with different base clusterings.

Keywords: Weighted Ensemble Clustering, Unsupervised Learning, HGPA, CSPA, AD

مسایلی که می‌توان به آن اشاره کرد، انتخاب الگوریتم خوشبندی مناسب از میان روش‌های متفاوت است که معیار ارزیابی یکتایی برای تمایزدادن بین آنها وجود ندارد. همچنین نبود معیار سنجش استانداردی برای نتایج خوشبندی موجب شده بدون هیچ دانش قبلی در مورد بهترین روش، هر نتیجه‌ای معمول به نظر برسد [2].

نتایج متفاوت حاصل از اجرای الگوریتم‌های خوشبندی مختلف بر روی مجموعه یکسانی از داده‌ها، حاکی از تمرکز هر کدام از آنها بر جنبه‌ای متفاوت از بهینگی است. در دسترس نبودن خوشبندی ایده‌آل، اعتبارسنجی خوشبندی‌ها را مشکل ساخته است. علاوه بر این بسیاری از الگوریتم‌ها نیازمند تنظیم پارامترهای ورودی از جانب کاربر هستند که این مسئله کنترل کیفیت خوشبندی را دشوار می‌سازد. بعد از داده نیز از مواردی است که عملکرد خوشبندی را از نظر پیچیدگی زمانی به چالش کشیده است.

از آن جایی که پژوهش‌های انجام شده در زمینه دسته‌بندی ترکیبی نمایان گر بهبود نتایج بوده [3-10]، منطقی است که این روش در خوشبندی نیز مورد بررسی

۱- مقدمه

خوشبندی^۱ به فرآیندی گفته می‌شود که در آن مجموعه‌ای از اشیا به زیرمجموعه‌هایی به نام خوش، بخش‌بندی می‌شوند؛ به طوری‌که، اعضای هر خوش به یکدیگر شبیه و با اعضای دیگر خوش‌ها متفاوت‌اند. ویژگی متایزکننده تحلیل خوشة نسبت به دسته‌بندی^۲، گروه‌بندی بدون نظارت داده‌ها است. این فرآیند سودمند می‌تواند موجب یافتن گروه‌ها و الگوهای ناشناخته‌ای در داده‌ها شود [1].

امروزه استفاده از خوشبندی در بسیاری از حوزه‌ها همچون داده‌گاوی، یادگیری ماشین، جستجوی وب، بازیابی اطلاعات، تشخیص الگو، امنیت و غیره رو به افزایش است. با وجود حجم بالا و گونه‌های مختلفی از داده و داده‌هایی با ابعاد بالا تلاش‌های بسیاری برای یافتن روش‌های کارآمد انجام شده است.

اگرچه روش‌های خوشبندی برای خلاصه‌سازی و فهم حجم انبوی داده‌ها بسیار مناسب هستند، اما همواره پژوهش‌گران با چالش‌های آن نیز رو به رو بوده‌اند. از نخستین

¹ clustering

² classification



خوشه ها، از چالش هایی است که در این پژوهش به آنها پرداخته شده است. در ادامه ابتدا ادبیات این موضوع در بخش ۲ تعریف و پژوهش های انجام شده در این زمینه بررسی شده است. در بخش ۳ الگوریتم هایی سرح داده شدند که در روش پیشنهادی مورد استفاده قرار گرفتند. در بخش ۴ به شرح الگوریتم پیشنهادی پرداخته خواهد شد و نتایج تجربی حاصل از آزمایش های انجام شده در بخش ۵ گزارش می شود. در نهایت نتیجه گیری و پیشنهاد هایی جهت کارهای آینده بیان شده است.

۲- ادبیات پژوهش و مرور کارهای پیشین

ایده اصلی ترکیب مدل های آموزشی و منابع داده در رشته های متعددی وجود دارد. در روش های دسته بندی ترکیبی مدلی جهت دسته بندی ساخته می شود که ترکیب چندین دسته بندی است. درواقع، برچسب دسته یک نمونه جدید براساس تجمعی آرای دسته بندی های اولیه تعیین می شود. از جمله روش های رایج در دسته بندی ترکیبی می توان به روش boosting, bagging و جنگل های تصادفی اشاره کرد [1].

الگوریتم آرای بیشینه وزن دار^۱، در سال ۱۹۹۴ ارائه شده است [14] این روش با استفاده از پیش بینی های انجام شده توسط مجموعی^۲ از الگوریتم ها با خطای کمتری برچسب مورد نظر را پیش بینی می کند. روش دیگری مشابه این روش با نام voted-perceptron نیز پیشنهاد شده است [15]. در دسته بندی ترکیبی نسبت به خوشه بندی ترکیبی مطالعات بیشتری صورت گرفته است [16-19]. در میان روش هایی که در همین اواخر در این زمینه مطرح شده است، می توان به روش AD اشاره کرد [13]. این روش نیز از مفهوم رأی گیری بیشینه در آرای مجموعی از انسان ها استفاده می کند و برای وزن دهی نظرات، روشی را جهت برآوردن صحت انسان ها در مسائل جمع سپاری معرفی کرده است. بخش عمده پژوهش های انجام شده در زمینه خوشه بندی ترکیبی، طی دهه اخیر صورت گرفته که نمایان گر نوبایودن این بخش است. دو گام اساسی در روش های مطرح شده برای خوشه بندی ترکیبی، گام تولید^۳ و گام تابع توافقی^۴ است (شکل ۱).

در گام تولید که نخستین مرحله است، مجموعه ای از

قرار بگیرد. خوشه بندی ترکیبی شامل تجمیع چندین افزار است که این افزارها می توانند نتیجه اجرای یک الگوریتم با شرط های آغازین گوناگون، یا نمونه برداری های متفاوت از مجموعه داده ها، و یا نتایج اجرای الگوریتم های مختلف بر روی مجموعه داده یکسان باشند.

روش های خوشه بندی ترکیبی می توانند راه حل های قدرتمند و پایداری برای چالش های ذاتی خوشه بندی ارائه دهند و بر مشکلات خوشه بندی غلبه کنند. این روش ها با استفاده از اجماع بین نتایج چندین خوشه بندی به متعادل کردن ساختارهای ناخواسته ای ناشی از تنظیمات متفاوت هر الگوریتم، به وسیله واریانس ایجاد شده در اثر نمونه گیری های گوناگون از داده می پردازند [11]. در پایگاه داده های توزیع شده که امکان اجرای خوشه بندی به صورت یکجا بر روی داده ها وجود ندارد، می توان خوشه بندی را برای هر پایگاه داده اجرا و تنها برچسب خوشه ها را به یک محل مرکزی منتقل کرد؛ سپس با استفاده از الگوریتم های خوشه بندی ترکیبی، افزار نهایی به دست می آید [12]. در کنار مزیت های خوشه بندی ترکیبی که ذکر شد، خود این روش نیز با چالش هایی روبرو است و عمدت ترین مسئله نحوه ترکیب خوشه بندی ها است. در ترکیب خوشه ها ممکن است با مسئله سازگار نبودن نتایج مواجه شویم و همه روش ها نظر یکسانی در مورد هم خوشه بودن دو نمونه نداشته باشند.

در این مقاله این ایده استفاده می شود که با نظرات خوشه بندی های مختلف به طور یکسان برخورد نشود و برای هر کدام وزنی در نظر گرفته شود. این وزن دهی می تواند براساس صحت^۱ هر یک از الگوریتم ها باشد، در این صورت چگونگی محاسبه صحت در خوشه بندی ها که روش های بدون ناظر هستند از مشکلات پیش روست. مرجع [13] روش AD را جهت برآوردن صحت انسان ها در مسائل جمع سپاری^۲ ارائه داده است. این روش، از هماهنگی یا تضاد بین نظرات انسان ها به عنوان یک منبع مفید اطلاعات استفاده کرده است و با پیشنهاد مدلی احتمالاتی، فرآیند برآوردن صحت را به کمک یک فرآیند بهینه سازی انجام می دهد. در این مقاله، در فرآیند ترکیب خوشه بندی ها وزن دهی به خوشه ها براساس روش AD انجام می شود تا بتوانیم صحت خوشه بندی نهایی را افزایش دهیم. نحوه تطبیق مسئله خوشه بندی به روش برآوردن صحت AD و نحوه استفاده از صحت های برآورده شده در فرآیند ترکیب نهایی

¹ accuracy

² agreement/disagreement

³ crowdsourcing

مجدد پیشنهاد داده‌اند [26]. تاپچی^۹ و همکارانش نیز چارچوبی یکپارچه برای ایجاد افزایش‌های متعدد ارائه داده‌اند [27].

علیزاده و همکاران برای خوشبندی ترکیبی روشی مبتنی بر زیرمجموعه‌ای از خوشه‌های پایه پیشنهاد داده‌اند [28]. در این روش تنها خوشه‌های پایدار در ترکیب نهایی مؤثر هستند و برای ساخت تابع توافقی مبتنی بر ماتریس همبستگی، روش خوشبندی انباست مدارک را توسعه داده‌اند. مرجع [29] با تعبیه خوشبندی در فضاهای برداری، مسئله پیچیده خوشبندی ترکیبی را به مسئله شناخته‌شده میانه اقلیدسی^{۱۰} کاهش می‌دهد و بهوسیله الگوریتم Weiszfeld آن را حل می‌کند. روش‌های ذکر شده بر روی خوشبندی مسطح تمرکز دارند؛ اما میزبانی در رساله خود به معرفی ترکیب خوشبندی‌های سلسله‌مراتبی پرداخته است [30].

در میان پژوهش‌هایی که در زمینه خوشبندی ترکیبی وزن‌دار صورت گرفته است، می‌توان به روش خوشبندی زیرفضایی COSA¹¹ اشاره کرد [31]. این روش به ازای هر داده یک بردار وزن برای ویژگی‌ها تعریف می‌کند که در ابتدا وزن تمام ویژگی‌ها یکسان است. فرآیند محاسبه فواصل داده‌ها با توجه به وزن ویژگی‌ها، اجرای KNN و بهروزرسانی وزن ویژگی‌ها تا جایی تکرار می‌شود که کمینه تغییر مورد نظر رخ دهد. الگوریتم LAC نیز مشابه COSA از وزن‌دهی استفاده می‌کند، با این تفاوت که بردار وزن را به خوشها اختصاص می‌دهد [32].

الگوریتم افزایش‌بندی مشابه وزن‌دار (WSPA)¹² برای حل مشکل ابعاد بالای داده از الگوریتم LAC استفاده می‌کند و افزار توافقی با افزایش‌بندی گرافی به دست می‌آید که از ماتریس شباهتی مبتنی بر بردار احتمال نسبت داده شده به هر نمونه ساخته شده است [33]. روش WBPA در ابتدا همانند روش WSPA است، اما در ساخت گراف پایانی همان‌طور که از نام آن مشخص است، یک گراف دوبخشی از نمونه‌ها و خوشها می‌سازد [33]. WSBPA¹³ نخستین رویکرد در خوشبندی زیرفضایی با استفاده از خوشبندی WBPA است [11]. این روش توسعه‌ای از روش WBPA است که علاوه‌بر افزایش‌بندی نهایی، بردار وزن ویژگی‌ها برای هر خوشه نیز ارائه می‌شود.

⁹ Topchy

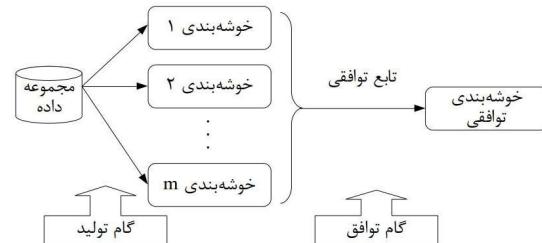
¹⁰ euclidean median

¹¹ clustering objects on subsets of attributes

¹² weighted similarity partitioning algorithm

¹³ weighted subspace bipartite partitioning algorithm

خوشبندی‌ها از مجموعه داده تولید می‌شوند؛ سپس در مرحله دوم باید نتایج حاصل از خوشبندی‌های پایه بهوسیله یک تابع توافق با یکدیگر تجمیع شوند تا خوشبندی نهایی به دست آید.



شکل-۱: نمودار فرآیند کلی خوشبندی ترکیبی [2]

(Figure-1): The general process of ensemble clustering [2]

استرل و گاش^۱ جزء نخستین کسانی بودند که ترکیب را در خوشبندی مورد پژوهش قرار دادند [12]. آنها در پژوهش خود سه روش افزایش‌بندی مبتنی بر گراف، جهت ساخت تابع توافقی پیشنهاد داده‌اند: الگوریتم افزایش‌بندی با شباهت مبتنی بر خوشة (CSPA)^۲، الگوریتم افزایش‌بندی ابرگراف (HGPA)^۳ و الگوریتم متا خوشبندی (MCLA)^۴. در این روش‌ها ابتدا مجموعه خوشبندی‌های پایه بهصورت یک ابرگراف مدل و محاسبات دیگر با استفاده از این ابرگراف انجام می‌شود.

حدود یک سال بعد فرد و جین^۵ به بررسی ایده جمع‌آوری شواهد برای ترکیب نتایج چندین خوشبندی برداختند [20]. آنها برای ساخت تابع توافق، افزایش‌های اشیا را به یک ماتریس همبستگی^۶ از الگوهای تکراری نگاشتند. روش‌های رأی‌گیری در مراجع [21] و [22] برای ساختن تابع توافقی به عنوان جایگزینی برای ماتریس همبستگی مطرح شده‌اند. مرجع [23] تابع توافقی را بر اساس اصل تنگی اطلاعات^۷ می‌سازد.

مرجع [24] دو روش برای تولید افزایش‌های گوناگون، با عنوان الگوریتم‌های خوشبندی ضعیف^۸ معرفی می‌کند. در پژوهش دیگری با استفاده از اجرای الگوریتم k-means با تعداد تصادفی خوشه‌ها، افزایش‌های گوناگون تولید و بهوسیله ماتریس تطابق بین تمام جفت اشیا افزایش‌ها ترکیب می‌شوند [25]. مینایی و همکارانش روشی مبتنی بر نمونه‌برداری

¹ Strehl and Ghosh

² cluster-based similarity partitioning algorithm

³ hypergraph partitioning algorithm

⁴ mcta-clustering algorithm

⁵ Fred and Jain

⁶ co-association matrix

⁷ information bottleneck principle

⁸ weak clustering algorithms

ترکیب کرده و روشی با عنوان K-MWO ارائه داده است. این روش از وزن دهی به نقاط داده استفاده می‌کند. در میان پژوهش‌های داخلی نیز پروین در رساله خود روش جدیدی ارائه داده که علاوه بر وزن دهی خوشه‌های پایه، از وزن دهی ویژگی‌های هر خوشه نیز استفاده کرده است [42]. در پژوهش [43] پس از گرomboندی خوشه‌بندی‌ها با استفاده از معیارهای شباهت، با کیفیت‌ترین خوشه‌بندی از هر گروه بر اساس شاخص‌های اعتبار خوشه انتخاب می‌شود؛ سپس خوشه‌بندی‌های انتخاب شده با استفاده از سه تابع MCLA و CSPA، HGPA ترکیب می‌شوند. مرجع [44] هم از توافق/عدم توافق خوشه‌بندی‌ها برای سنجش اعتبار آنها استفاده می‌کند. از میان پژوهش‌های اخیر نیز می‌توان به مراجع [45-48] اشاره کرد.

با توجه به پژوهش‌های انجام شده، برآورد صحّت خوشه‌بندی‌ها با توجه به ماهیت بدون ناظر آنها یکی از چالش‌های روش‌های خوشه‌بندی ترکیبی است. از آنجایی که در روش AD وزن دهی صرفاً با استفاده از آرای دسته‌بندها انجام می‌شود، بررسی این روش در زمینه خوشه‌بندی ترکیبی نیز می‌تواند سودمند باشد. در این مقاله، علاوه بر معرفی یک روش خوشه‌بندی ترکیبی وزن دار مبتنی بر AD، تأثیر عواملی مانند اندازه خوشه‌ها و ارزیابی خوشه‌بندی‌ها با استفاده از یک معیار داخلی نیز بررسی می‌شود.

۳- الگوریتم‌های مورد استفاده

۱-۳- الگوریتم CSPA

الگوریتم CSPA یکی از روش‌های خوشه‌بندی ترکیبی مبتنی بر گراف است که در مرجع [12] ارائه شده است. اگر n خوشه‌بندی، n داده و k خوشه داشته باشیم، برای بردار پرچسب هر یک از خوشه‌بندی‌ها $\lambda \in N^n$ یک ماتریس عضویت دودویی ($H^{(q)}$) می‌سازیم که سطرهای آن داده‌ها یا همان رئوس ابرگراف و ستون‌ها، خوشه‌ها یا ابریال‌ها هستند. مقدار یک در هر ستون بیان کننده عضویت رأس متناظر با آن سطر در این خوشه (ابریال) است. با کنار هم قراردادن ماتریس تمام خوشه‌بندی‌های پایه، بلوکی از ماتریس‌ها به دست می‌آید ($| H = h^{(1)}, \dots, h^{(r)} | = | H^{(1)}, \dots, H^{(r)} |$) که خود ماتریس مجاورت یک ابرگراف با n رأس و $(q) \sum_{q=1}^r K$ ابریال است. بدین ترتیب هر خوشه به یک ابریال و هر خوشه‌بندی به یک ابرگراف نگاشت می‌شود.

در روش ^۱ WKF گام ارزیابی کیفیت افزایش‌بندی‌ها میان گام تولید وتابع توافقی اضافه شده است که موجب تأثیر کمتر افزایش‌های نوافه با وزن کمتر، در افزایش نهایی می‌شود [34]. این روش سپس به ^۲ GWKF ارتقا پیدا کرد [35] که در آن اینکه کدام افزار دو نمونه را در یک خوشه قرار می‌دهد، حائز اهمیت است و در گام ارزیابی کیفیت افزایش‌بندی‌ها دیگر احتیاجی به دسترسی به داده‌های اصلی وجود ندارد.

در زمینه وزن دهی به خوشه‌بندی‌ها مرجع [36] ثابت می‌کند تنظیم ^۳ LASSO از مسئله معادل بهینه‌سازی وزن برای روشی است که جهت خوشه‌بندی ترکیبی وزن دار پیشنهاد می‌دهد. ^۴ GULLO و همکارانش سه طرح برای تجمیع خوشه‌بندی‌ها بر اساس دو معیار تنوع مبتنی بر F-measure و ^۵ NMI معرفی کردن؛ وزن دهی تکی ^۶، وزن دهی گروهی ^۷ و وزن دهی دندروگرام ^۸ [37].

پژوهش [38] مسئله خوشه‌بندی ترکیبی را به یک مسئله برنامه‌نویسی خطی دودویی ^۹ تبدیل کرده و راه حلی مبتنی بر گراف فاکتور ^{۱۰} پیشنهاد داده است. این پژوهش از روشی بدون ناظر برای اعتبار سنجی خوشه‌بندی‌های پایه استفاده کرده و نمایش ممتاز-شیء ^{۱۱} را برای داده‌های حاصل تجمیع پیشنهاد داده است.

مرجع [39] خوشه‌بندی ترکیبی طیفی ^{۱۲} را برای استفاده بیشینه‌ای از مزایای ماتریس همبستگی پیشنهاد می‌دهد. در این مقاله نشان داده شده که این خوشه‌بندی معادل خوشه‌بندی k-means وزن دار است که پیچیدگی آن به طور چشم‌گیری کاهش یافته است. در مقاله [40] از وزن دهی به اشیا در فرآیند تجمیع خوشه‌بندی استفاده شده است. این وزن دهی به وسیله ماتریس همبستگی صورت می‌گیرد که نتایج خوشه‌بندی پایه را خلاصه می‌کند.

مرجع [41] یک مدل هوش گروهی ^{۱۳} جدید به نام بهینه‌سازی سرگردان صدف‌ها ^{۱۴} را خوشه‌بندی k-means

¹ weighted cluster ensemble using a kernel consensus function

² generalized WKF

³ regularization

⁴ GULLO

⁵ normalized mutual information

⁶ single weighting

⁷ group weighting

⁸ dendrogram weighting

⁹ binary linear programming

¹⁰ factor graph

¹¹ super-object

¹² spectral ensemble clustering

¹³ swarm intelligence model

¹⁴ mussels wandering optimization

برآورده شده جهت برآورد برچسب‌های درست. با توجه به استفاده از گام نخست در این پژوهش، تنها به شرح این گام اکتفا می‌کنیم. برآورد صحت برچسب‌گذارها در سه بخش انجام می‌شود. در بخش نخست دوشی برای برآورد احتمال اختلاف نظر بر حسب برچسب‌های موجود مطرح می‌شود. این احتمال که برای هر جفت برچسب‌گذار محاسبه و P_i نامیده می‌شود، نسبت اختلاف‌نظرها به تعداد کل نمونه‌هایی است که هر دو برچسب‌گذار در مورد آن نظر داده‌اند.

در بخش دوم احتمال اختلاف نظر بر حسب صحت‌های فرضی محاسبه می‌شود. زمانی دو فرد i و j در مورد نمونه k اختلاف نظر پیدا می‌کنند که یکی از دو حالت زیر رخ دهد:

۱. فرد i درست و فرد j غلط برچسب‌گذاری کرده است.

۲. فرد j درست و فرد i غلط برچسب‌گذاری کرده است.

حالات نخست را E_1 و حالات دوم را E_2 می‌نامیم.

بنابراین:

$$P(O_k^i \neq O_k^j) = P(E_1) + P(E_2) = \quad (2)$$

$$P(O_k^i = y_k, O_k^j = 1 - y_k) +$$

$$P(O_k^i = 1 - y_k, O_k^j = y_k)$$

اگر صحت فرضی i و j به ترتیب $\hat{\alpha}^i$ و $\hat{\alpha}^j$ باشد، پس می‌توان احتمال E_1 و E_2 را به صورت نیز نوشت:

$$P(O_k^i = y_k, O_k^j = 1 - y_k) = \quad (3)$$

$$P(O_k^i = y_k)P(O_k^j = 1 - y_k) = (\hat{\alpha}^i)(1 - \hat{\alpha}^j)$$

$$P(O_k^i = 1 - y_k, O_k^j = y_k) = \quad (4)$$

$$P(O_k^i = 1 - y_k)P(O_k^j = y_k) = (1 - \hat{\alpha}^i)\hat{\alpha}^j$$

در نتیجه معادله (۲) به این ترتیب ساده‌سازی می‌شود:

$$P(O_k^i \neq O_k^j) = \hat{\alpha}^i + \hat{\alpha}^j - 2\hat{\alpha}^i\hat{\alpha}^j \quad (5)$$

ماتریس این احتمالات P_2 نامیده می‌شود. بخش سوم در برآورد صحت برچسب‌گذارها، حل مسئله بهینه‌سازی است. در این بخش با استفاده از یکتابع، بهینه‌سازی اختلاف بین P_1 و P_2 کمینه می‌شود. این تابع، مقدار دلخواهی برای صحت یک برچسب‌گذار در نظر می‌گیرد و با جای‌گذاری در معادله $P_2 - P_1 = 0$ صحت بقیه را برآورد می‌زند و این فرآیند را برای مقادیر مختلف صحت یک برچسب‌گذار تکرار می‌کند. خروجی تابع بهینه‌سازی بردار صحت‌های برآورده شده برای برچسب‌گذاران است.

بنابراین، اگر شباهت دو شیء را که در خوشة یکسان هستند یک و غیر این صورت صفر تعريف کنیم، می‌توانیم یک ماتریس شباهت H برای هر خوشه‌بندی ایجاد کنیم. در این صورت با محاسبه میانگین بین درایه‌های متناظر ماتریس هر خوشه‌بندی، می‌توان ماتریس شباهتی تعريف کرد که از طریق معادله (۱) قابل محاسبه است.

$$S = \frac{1}{r} HH^t \quad (1)$$

رئوس گراف شباهت حاصل‌شده از این ماتریس با اشیا و وزن یال‌ها با مقادیر شباهت دوبه‌دوی اشیا، متناظر است. درنهایت این گراف برای دوباره خوشه‌بندی کردن اشیا با استفاده از METIS [49] افزایش‌بندی می‌شود و نتیجه این افراز، تجمعی خوشه‌بندی‌های مختلف را به دست می‌دهد.

۳-۲- الگوریتم HGPA

این الگوریتم نیز جزء سه روش افزایش‌بندی مبتنی بر گراف است که در مرجع [12] ارائه شده است. HGPA از ابرگرافی که مشابه ابرگراف الگوریتم CSPA به دست می‌آید، برای ترکیب خوشه‌بندی‌ها استفاده می‌کند. این الگوریتم سعی در یافتن کمترین تعداد از ابریال‌هایی دارد که فقط آنها ابرگراف را به k بخش مجزا با اندازه‌های به طور تقریبی یکسان افزایش‌بندی می‌کند. در این روش تمام یال‌ها و تمام رأس‌ها وزن یکسانی دارند. در افزایش‌بندی ابرگراف از بسته HMETIS [50] استفاده می‌شود که بسیار مقیاس‌پذیر است و افزایه‌ای با کیفیت بالایی تولید می‌کند.

۳-۳- الگوریتم AD

این الگوریتم جهت برآورد صحت انسان‌ها در مسائل جمع‌سپاری ارائه شده است که بر روی مسائل دودسته‌ای تمرکز دارد [13]. فرض کنید N نمونه آموزشی هستند که برچسب‌های درست آنها $\{y_i\}_{i=1}^N$ نامشخص است. به جای این برچسب‌های درست، نظرات R انسان را در مورد بعضی از این نمونه‌ها داریم. نظر برچسب‌گذار ز در مورد نمونه i با $O_i^j \in \{-1, 0, 1\}$ نشان داده می‌شود که -1 یعنی نظری داده نشده است. صحت برچسب‌گذار ز با α^j نمایش داده می‌شود.

AD شامل دو گام است:

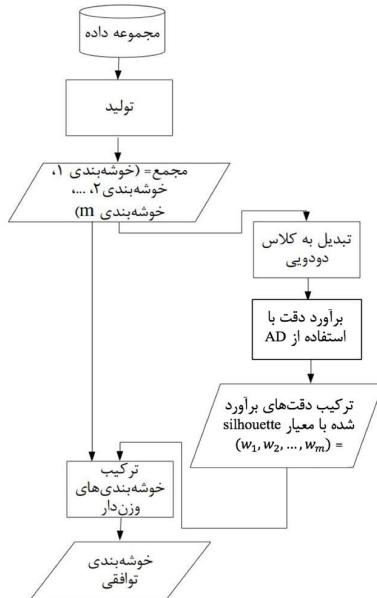
۱. برآورد صحت برچسب‌گذارها با استفاده از توافقات/اختلافات موجود در نظرات جمع‌آوری شده،

۲. یکی کردن نظرات با استفاده از صحت‌های

فصلنامه



انتقال میانگین^۱، خوشه بندی تجمعی^۲ با دو معیار ward و DBSCAN، Mini-batch k-means، average k-means، Spectral، Affinity Propagation و BIRCH به طور جداگانه بر روی داده ها اجرا می شوند؛ سپس نتایج این خوشه بندی ها برای بدست آوردن خوشه بندی نهایی ترکیب می شوند.



(شکل-۲): چارچوب روش پیشنهادی برای ترکیب وزن دار خوشه بندی ها

(Figure-2): The proposed method for weighted ensemble clustering

• خوشه بندی ویژگی ها

ای توزیع شده (FDC³): در روش ویژگی های توزیع شده خوشه بندی های پایه به عنوان اجرای k-means برای زیر مجموعه هایی از ویژگی ها تولید می شوند. این زیر مجموعه ها با انتخاب تصادفی ۴۰٪، ۵۰٪، ۶۰٪، ۷۰٪، ۸۰٪، ۹۰٪ و ۱۰۰٪ از ویژگی ها به دست می آیند. در نهایت هفت خوشه بندی ایجاد شده ترکیب و نتیجه توافقی حاصل می شود.

استفاده از معیار های فاصله متفاوت در اجرای k-means برای استفاده از معیار های فاصله متفاوت در تولید خوشه بندی های پایه، الگوریتم k-means را به ازای هر یک از فاصله های کسینوسی^۴، اقلیدسی^۵، همبستگی^۶ و بلوك

¹ mean shift

² agglomerative clustering

³ feature-distributed clustering

⁴ cosine

⁵ euclidean

⁶ correlation

با توجه به معادله (۵) برای برآورده صحت برچسب گذارها، این مسئله دو نقطه بهینه دارد:

$$\hat{\alpha}^i + \hat{\alpha}^j - 2\hat{\alpha}^i\hat{\alpha}^j = \\ (1 - \hat{\alpha}^i) + (1 - \hat{\alpha}^j) - 2(1 - \hat{\alpha}^i)(1 - \hat{\alpha}^j) \quad (6)$$

بنابراین، با جایگزین کردن بردار $(\hat{\alpha}^1, \dots, \hat{\alpha}^R)$ با بردار $(1 - \hat{\alpha}^1, \dots, 1 - \hat{\alpha}^R)$ مقدار تابع هدف تعییر نمی کند. در نتیجه، اگر خروجی الگوریتم بهینه سازی بردار $\bar{\alpha} = (1 - \hat{\alpha}^1, \dots, 1 - \hat{\alpha}^R)$ باشد، بردار $(\hat{\alpha}^1, \dots, \hat{\alpha}^R)$ برای می تواند پاسخ درست باشد. به عبارت دیگر، یکی از این نقاطه ها بهینه درست و دیگری غلط است. برای مسئله بحث شده در این مقاله، روشی جهت تشخیص بهینه درست طراحی و بررسی شده است که در بخش بعد معرفی خواهد شد.

۴- روش پیشنهادی

ایده مطرح شده در روش پیشنهادی، استفاده از صحت های خوشه بندی های پایه به عنوان وزن هر یک از آنها در فرآیند ترکیب است. در طراحی روش پیشنهادی سه گام طی می شود: تولید خوشه بندی های پایه، برآورده خوشه بندی با استفاده از روش AD و معیار silhouette و ترکیب خوشه بندی های پایه مبتنی بر وزن های برآورده شده. شکل (۲) این مراحل را به تصویر کشیده است.

۴-۱- گام نخست: تولید خوشه بندی های پایه

در گام نخست روش های متفاوتی برای تولید خوشه بندی های پایه مورد بررسی قرار داده شده است. با اجرای هر روش بر روی مجموعه ای از داده ها، خوشه بندی های اولیه به دست می آیند. در ادامه روش های مورد استفاده در آزمایش ها برای تولید خوشه بندی های پایه بیان می شوند:

- اجرای k-means با k های متفاوت: در این روش، ترکیب بر روی خوشه بندی های حاصل از اجرای الگوریتم k-means برای ۲ خوشه و ۵۰٪، ۷۵٪، ۱۰۰٪، ۱۵۰٪ و ۲۰۰٪ تعداد واقعی خوشه ها صورت می گیرد. خوشه بندی توافقی برای تعداد خوشه هایی معادل تعداد واقعی به دست می آید.

- استفاده از الگوریتم های متفاوت: در این روش مجموعه ای شامل ۹ الگوریتم خوشه بندی داریم. هر یک از الگوریتم های

Algorithm 1: Convert ensemble to binary opinions

Parameters: N (number of instances), R (number of clusterings)

Input: Dataset, Ensemble ($R \times N$)

Output: ADinput ($R \times N^2$)

```

(1) for clusterings in Ensemble do
(2)   Index← 0
(3)   for (i,j) in Dataset2 do
(4)     If clusterings (i) = clusterings (j) then
(5)       ADinput (labeler)(index) ← 1
(6)     else do
(7)       ADinput (labeler)(index) ← 0
(8)     Index← Index+1
(9)   end for
(10) end for

```

(شکل-۳): الگوریتم تبدیل مجموعه‌ای از خوشبندی‌ها به نظرات

دودویی قابل استفاده در الگوریتم AD

(Figure-3): Converting an ensemble of clusterings to binary opinions which can be used in AD algorithm

۱-۲-۴-۴-۲-۱- انتخاب جواب درست از میان دو بهینه روش AD

برای تشخیص بهینه درست از غلط آخرين نمونه موجود در مجموعه داده را بیست مرتبه تکرار و سپس الگوریتم پیشنهادی را یک مرتبه با استفاده از صحت‌های α و بار دیگر با استفاده از $\alpha - 1$ مطابق توضیحات بخش قبل بر روی این مجموعه داده اجرا می‌کیم. با توجه به یکسان‌بودن بیست نمونه آخر، نتایج به دست‌آمده تنها برای این بیست نمونه با استفاده از معیار Fowlkes-Mallows³ مورد سنجش قرار می‌گیرد.

برای تمام مجموعه‌داده‌ها، روش‌های تولید خوشبندی پایه و معیارهای ارزیابی مطرح شده، با مقایسه نتایج خوشبندی نهایی و نتایج بیست نمونه تکراری، مشاهده شد هر یک از وزن‌دهی‌ها که برای بیست نمونه تکراری بهتر عمل کرده، در خوشبندی نهایی نیز عملکرد بهتری داشته است. بنابراین بیشترین مقدار معیار Fowlkes-Mallows برای این بیست نمونه نشان‌دهنده استفاده از وزن‌های برآورده درست است. تمام مواردی که در این فصل به عنوان نتیجه وزن‌دهی با استفاده از صحت‌های برآورده AD گزارش شده است، وزن‌دهی با جواب درست برمبنای این روش است.

۱-۲-۴-۳- گام سوم: ترکیب خوشبندی‌های پایه مبتنی بر وزن‌های برآورده شده

در این گام خوشبندی‌های پایه با استفاده از روش CSPA و HGPA وزن‌دار ترکیب می‌شوند که در ادامه وزن‌دار کردن آنها را شرح می‌دهیم. در روش CSPA وزن‌دار از صحت‌های برآورده شده توسط AD به عنوان وزن خوشبندی‌ها استفاده می‌شود. در روش HGPA وزن‌دار، برای وزن ابریال‌ها در

شهری^۱ دو بار به طور مستقل اجرا می‌کنیم (جدول ۱). سپس الگوریتم ترکیب بر روی ده خوشبندی به دست‌آمده اجرا می‌شود. مستند جزئیات پیاده‌سازی‌های این معیارها در پایتون از طریق پایگاه^۲ scipy قابل دسترسی است.

(جدول-۱): معیارهای فاصله مورد استفاده

(Table-1): The used distance metrics

نام معیار	نحوه محاسبه (برای دو شی و p با $j = (x_{j1}, \dots, x_{jp})$ ویژگی عددی)
اقلیدسی [1]	$\sqrt{(x_{i1} - x_{j1})^2 + \dots + (x_{ip} - x_{jp})^2}$
بلوک شهری [1]	$ x_{i1} - x_{j1} + \dots + x_{ip} - x_{jp} $
کسینوسی [51]	$d(i, j) = 1 - \frac{i^T \cdot j}{\ i\ \cdot \ j\ }$ $\ i\ = \sqrt{x_{i1}^2 + \dots + x_{ip}^2}$
همبستگی [51]	$d(i, j) = 1 - \frac{(i - \bar{i})^T \cdot (j - \bar{j})}{\ \bar{i} - \bar{j}\ \cdot \ \bar{j}\ }$ $\bar{i} = \frac{x_{i1} + \dots + x_{ip}}{p}$

۱-۲-۴-۲- گام دوم: برآورد صحت هر خوشبندی

با استفاده از روش AD

از آنجایی که روش AD بر روی مسائل دودسته‌ای کار می‌کند باید نتایج خوشبندی‌های پایه را به شکل مناسبی تبدیل کنیم. شکل (۳) الگوریتم این تبدیل را نمایش می‌دهد که در آن ورودی، مجموعه داده (Dataset) و ماتریس خوشبندی‌های پایه در گام تولید (ADinput) در نظر است. در این الگوریتم ماتریسی دودویی (ADinput) در نظر می‌گیریم که سطرها متناظر با خوشبندی‌ها (برچسب‌گذاران) هستند. ستون‌ها نیز، زوج مرتب‌هایی هستند که از مجموعه داده‌ها ایجاد شده‌اند. مقدار یک برای یک زوج مرتب (j, i) در سطر k به این معنی است که خوشبندی k دو نمونه i و j را در خوشبندی‌های یکسان قرار داده است و صفر یعنی در خوشبندی‌های یکسان نیستند. به سبب متقارن‌بودن این ماتریس تنها یک نیمه ماتریس محاسبه می‌شود و عناصر روی قطر اصلی که همگی مقدار یک دارند، به منظور کاهش هزینه محاسباتی حذف می‌شوند؛ سپس این ماتریس به عنوان ورودی الگوریتم AD قرار داده می‌شود تا صحت‌های هر خوشبندی برآورد شود.

³ <http://scikit-learn.org/stable/modules/clustering.html#fowlkes-mallows-scores>

¹ city block

² <http://docs.scipy.org/doc/scipy/reference/spatial.distance.html>

$$\text{score}(\lambda^t) = \begin{cases} w_i & \text{if } \text{label}(x_i) = \text{label}(x_j) \\ 1 - w_i & \text{otherwise.} \end{cases} \quad (7)$$

به عبارت دیگر اگر خوشه بندی با وزن w دو نمونه را در یک خوشه قرار نداده است، پس با وزن $w - 1$ به هم خوشه بودن آنها رأی می دهد. در این صورت درایه λ و زاز ماتریس شباهت از معادله (8) بدست می آید.

$$S(i, j) = \sum_{l=1}^r \frac{\text{score}(\lambda^l)}{r} \quad (8)$$

برای وزن دار کردن HGPA کافی است به جای یکسان در نظر گرفتن وزن ابریال ها، وزن های محاسبه شده برای هر خوشه بندی (طبق جدول ۲) را به عنوان وزن ابریال های معادل آن در ابر گراف در نظر بگیریم.

حالتهای مختلف ترکیب، سه پارامتر صحت های برآورده شده به وسیله AD، معیار silhouette و اندازه خوشه ها بررسی می شود. جدول (۲) انواع وزن دهی های مورد آزمایش برای HGPA را نشان می دهد.

(جدول ۲): شیوه های مختلف وزن دهی ابریال های در HGPA

(Table-2): Different ways of weighting hyperedges in HGPA algorithm

نماد	وزن ابریال ها
w_1	$\text{int}(\text{accuracy} \times 1000)$
w_2	$\text{int}(\text{accuracy} \times 1000) \times \text{Vertices}$
w_3	$\text{int}(\text{accuracy} \times 1000) \times \text{int}(\frac{\text{Vertices}}{N} \times 1000)$
w_4	$\text{int}(\text{silhouette} \times 1000) \times \text{int}(\frac{\text{Vertices}}{N} \times 1000)$
w_5	$\text{int}(\text{silhouette} \times 1000)$
w_6	$\text{int}(\text{accuracy} \times \text{silhouette} \times 1000)$
w_7	$\text{int}(\text{accuracy} \times \text{silhouette} \times 1000) \times \text{int}(\frac{\text{Vertices}}{N} \times 1000)$

پارامتر accuracy در جدول (۲) با جای گذاری صحت های برآورده شده توسط AD به طور جداگانه آزمایش شده است. از آنجایی که HMETIS وزن ها را به صورت عدد صحیح می پذیرد، صحت را که عددی بین صفر و یک است در هزار ضرب و سپس به عدد صحیح تبدیل کرده ایم. پارامتر Vertices تعداد رأس های هر ابریال یا به عبارتی دیگر تعداد نمونه های عضو هر خوشه است؛ زیرا در این الگوریتم هر خوشه به یک ابریال نگاشت می شود. N تعداد کل نمونه ها است که برای نرم افزار silhouette وزن های w_3 ، w_4 و w_7 استفاده شده است. با توجه به این که این پارامتر پس از نرم افزار silhouette بین صفر و یک قرار می گیرد، مانند صحت آن را به عدد صحیح تبدیل می کنیم. یک معیار ارزیابی silhouette داخلی خوشه بندی در بازه [-1, 1] است [52] و همان طور که پیش تر توضیح داده شد به عدد صحیح تبدیل می شود. در وزن دهی به وسیله w_1 ، w_5 و w_6 وزن تمام ابریال ها یا خوشه های یک خوشه بندی یکسان است؛ زیرا صحت و معیار silhouette برای کل خوشه بندی محاسبه می شود و برای خوشه های داخل خوشه بندی متفاوت نیست.

۱-۳-۴- وزن دار کردن HGPA و CSPA

با اندک تغییری در CSPA، امکان ترکیب وزن دار خوشه بندی ها را به این الگوریتم خوشه بندی ترکیبی اضافه کرده ایم. اگر r خوشه بندی داشته باشیم، برای ساخت ماتریس شباهت، ابتدا برای هر زوج نمونه (x_i, x_j) امتیاز هر خوشه بندی λ^l با وزن w_l ($0 \leq w_l \leq 1$) را مطابق (7) محاسبه می کنیم.

۵- نتایج تجربی

در این بخش به بیان نتایج حاصل از آزمایش روش پیشنهادی می پردازیم. به منظور پیاده سازی این روش از زبان برنامه نویسی Python استفاده شده است. از جمله کتابخانه های مورد استفاده در این پژوهش می توان به `Scipy`¹ و `Numpy`² اشاره کرد که برای انواع محاسبات عددی و مهندسی ارائه شده اند. همچنین `Matplotlib` یادگیری ماشین `Scikit-learn`³ ابزارهای ساده و کارایی جهت تحلیل داده و داده کاوی فراهم کرده است. در ادامه ابتدا مجموعه داده های مورد بررسی و معیارهای ارزیابی مورد استفاده بیان می شوند و سپس به بیان نتایج می پردازیم.

۱-۵- مجموعه داده های مورد آزمایش

کارایی الگوریتم پیشنهادی در مقایسه با الگوریتم های دیگر در چندین مجموعه داده واقعی و دست ساز بررسی می شود. مجموعه داده دست ساز داده هایی هستند که توسط استرل و گاش تهیه شده اند و از طریق (<http://strchl.com>) قابل دسترسی است. مجموعه داده های واقعی مورد استفاده، تعدادی از مجموعه داده های استاندارد UCI⁴ و دو مجموعه داده از scikitlearn⁵ هستند. بسیاری از مطالعات اخیر با استفاده از این مجموعه داده ها گزارش شده اند. مشخصات داده های مورد آزمایش در جدول (۳) ارائه شده است.

¹ <http://www.scipy.org>

² <http://www.numpy.org>

³ <http://www.scikit-learn.org>

⁴ <https://archive.ics.uci.edu/ml/datasets.html>

⁵ <http://scikit-learn.org/stable/datasets/#toy-datasets>



۲-۵- معیارهای ارزیابی نتایج

برای ارزیابی روش پیشنهادی از ۹ معیار استفاده شده است. جدول (۴) خلاصه‌ای از معیارهای ارزیابی استفاده شده را نشان می‌دهد. مستند جزئیات این معیارها در پایتون از طریق پایگاه scikitlearn¹ قابل دسترسی است.

(جدول-۳): مجموعه داده‌های مورد آزمایش

(Table-3): Datasets used in the experiments

نام مجموعه داده	تعداد خوش‌ها	تعداد ویژگی‌ها (ابعاد)	تعداد نمونه‌ها
(دسترساز) 8d5k	5	8	1000
Iris (scikitlearn)	3	4	150
Digits (scikitlearn)	10	64	1797
Seeds (UCI)	3	6	210
Data User Modeling Int (UCI)	4	4	258
Wholesale customers data (UCI)	3	6	440
Movement libras (UCI)	15	90	360
Zoo (UCI)	7	16	110
Wine (UCI)	3	13	178
Ionosphere (UCI)	2	34	351
Ecoli (UCI)	8	7	336
Image segmentation (UCI)	7	19	210
Glass (UCI)	7	10	214

(جدول-۴): معیارهای ارزیابی روش پیشنهادی

(Table-4): Evaluation measures

معیار	اطلاعات متقابل (MI)	اطلاعات متقابل (AMI)	اطلاعات متقابل (NMI)	همگوئی (Homogeneity)	جامیعت (Completeness)	V-measure	خلوص ² (Purity) ²
نُرمال شده در برابر برچسب‌های تصادفی	متقارن بازه	متقارن	بازه	در برابر برچسب‌های تصادفی	نُرمال شده		
✓	✓	[−1, 1]				adjusted rand index (ARI)	
✓	✗	[0, 1]				Fowlkes-Mallows (FMI)	
✓	✓	✗				اطلاعات متقابل (MI)	
✗	✓	Max=1				اطلاعات متقابل تطبیق داده شده (AMI)	میتنی بر اطلاعات متقابل
✗	✓	[0, 1]				اطلاعات متقابل نُرمال شده (NMI)	
✗	✗	[0, 1]					
✗	✗	[0, 1]					
✗	✗	[0, 1]					

¹ <http://scikit-learn.org/stable/modules/clustering.html#clustering-performance-evaluation>

² <http://nlp.stanford.edu/IR-book/html/htmledition/evaluation-of-clustering-1.html>

۳-۵- نتایج آزمایش‌ها

در این بخش نتایج هر یک از روش‌های تولید خوشبندی‌های پایه و ترکیب وزن‌دار آنها با استفاده از CSPA و HGPA وزن‌دار گزارش شده است. برای اینکه نتایج بیشتر قابل اعتماد باشد، میانگینی از ۵ بار اجرای مجزای فرآیند تولید خوشبندی‌های پایه و سپس ترکیب وزن‌دار ارائه شده است. در نمودارهای این بخش خط افقی معیارهای ارزیابی را نشان می‌دهد و چهار میله رنگی روش‌های تولید خوشبندی‌های پایه را مشخص کرده است که میله آبی رنگ استفاده از الگوریتم‌های مختلف، میله قرمز k-means، میله خاکستری خوشبندی ویژگی‌های توزیع شده و میله زرد اجرای k-means با k متفاوت است.

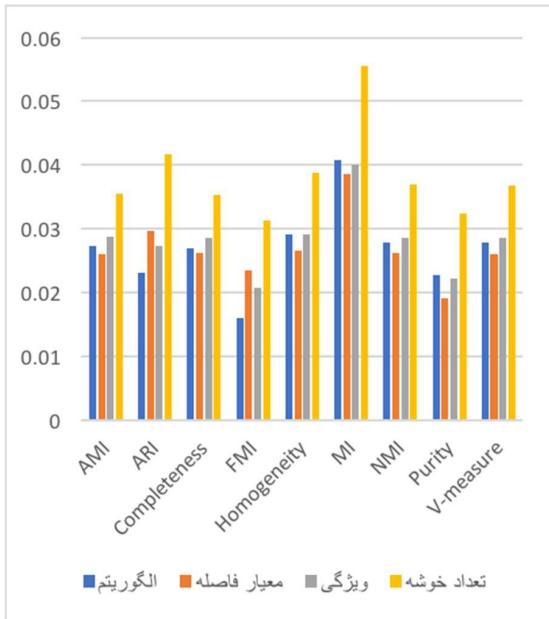
هر میله در نمودار، میانگین تفاضل CSPA (یا HGPA) بدون وزن از CSPA (یا HGPA) وزن‌دار برای هر یک از معیارهای ارزیابی گزارش شده است. بدین ترتیب اعداد مشبت نمایان گر بهبود عملکرد CSPA (یا HGPA) با استفاده از وزن دهی است. بهدلیل حجم بالای آزمایش‌ها، در این بخش خلاصه‌ای از نتایج به صورت میانگین نتایج تمام مجموعه داده‌ها گزارش شده است. گزارش آزمایش تمام مجموعه داده‌ها از طریق فایل موجود در وب‌سایت آزمایشگاه داده‌کاوی و پادگیری ماشین دانشگاه قم³ در دسترس است. در نمودار (۱)، نتایج برای تولید خوشبندی‌های پایه بر مبنای ویژگی‌های مختلف با استفاده از الگوریتم CSPA نمایش داده شده است. همان‌طور که مشاهده می‌شود، در این آزمایش بهبود قابل توجهی مشاهده نمی‌شود؛ به طوری که بهترین بهبود ایجاد شده کمتر از ۰.۰۱ است. به علاوه، استفاده از ویژگی‌های توزیع شده برای تولید خوشبندی‌های پایه باعث کاهش دقت در روش CSPA وزن‌دار شده است.

نمودارهای (نمودار تا نمودار تا) نتایج تفاضل الگوریتم HGPA بدون وزن از HGPA وزن‌دار را با وزن دهی توسط w_7 تا w_1 (طبق (جدول) زمانی که accuracy برابر صحّت‌های برآورده شده توسط AD است نشان می‌دهند).

برای مقایسه روش‌های تولید خوشبندی‌های پایه توجه کنید که عملکرد روش تولید به روش وزن دهی وابسته است. به عنوان مثال، روش استفاده از تعداد متفاوت خوشبندی (میله زرد) هرچند برای روش وزن دهی w_1 و w_2 و w_7

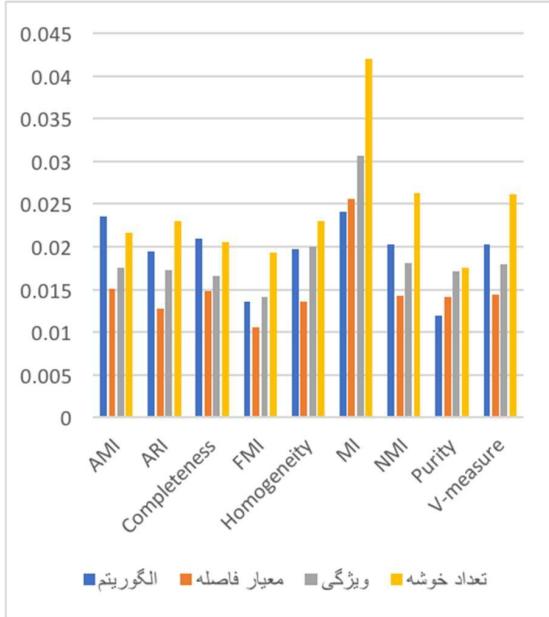
³ http://dml.qom.ac.ir/wp-content/uploads/2017/05/Results_Vahidi.pdf

(نمودار-۱): تفاضل الگوریتم CSPA بدون وزن از CSPA با وزن دهنده به وسیله صحت های برآورد شده توسط AD
(Chart-1): Subtraction of unweighted CSPA from weighted CSPA using accuracies estimated by AD method



(نمودار-۲): تفاضل الگوریتم HGPA بدون وزن از HGPA با وزن دهنده w_1 و وزن دار AD

(Chart-2): Subtraction of HGPA from weighted HGPA by using w_1 and AD estimated



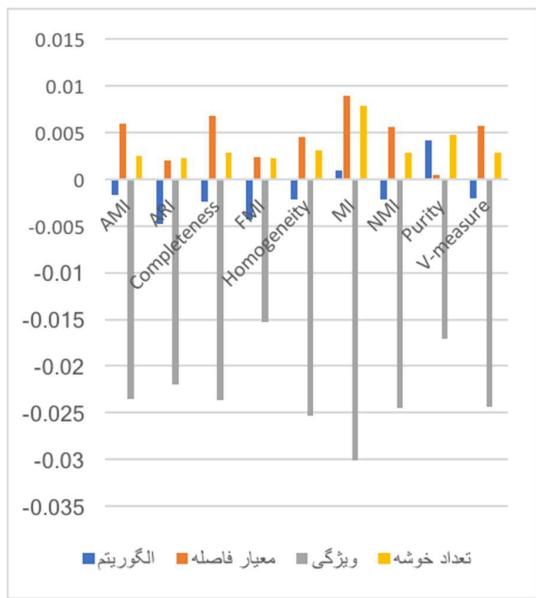
(نمودار-۳): تفاضل الگوریتم HGPA بدون وزن از HGPA با وزن دهنده w_2 و وزن دار AD

(Chart-3): Subtraction of HGPA from weighted HGPA by using w_2 and AD estimated accuracies

بهترین نتایج را حاصل کرده است، برای روش های w_3 و w_4 ضعیفتر از بقیه روش ها عمل کرده است. به همین صورت، روش استفاده از الگوریتم های مختلف (میله آبی) برای روش وزن دهنده w_3 و w_6 عمل کرد بسیار خوبی داشته است، در حالی که در روش w_5 باعث کاهش دقت روش تحمیعی شده است. روش استفاده از معیارهای فاصله متفاوت برای تولید خوشبندی های پایه نیز در روش وزن دهنده w_5 بهترین نتایج را به دست آورده است؛ درنهایت روش ویژگی های توزیع شده که در بسیاری از روش های وزن دهنده عمل کرد قابل قبولی نشان داده است، در روش w_5 ضعیفتر از بقیه روش ها عمل کرده و باعث کاهش دقت الگوریتم ترکیب شده است. مقایسه روش های وزن دهنده مختلف نشان می دهد که در مجموع استفاده از w_1 , w_2 , w_6 و w_7 نتایج را بهبود داده است و به طور متوسط بهترین نتایج مربوط به زمانی است که وزن دهنده به وسیله w_1 یعنی تنها با استفاده از صحت های برآورده AD است. از مقایسه نمودار (۳) و نمودار (۴) این نتیجه حاصل می شود که نرمال سازی کارایی را کاهش داده است. نمودار (۶) نیز نشان می دهد استفاده از نیز کارایی پایینی دارد.

در نمودار (۵) ضعیفترین عملکرد مشاهده می شود که وزن دهنده w_4 که وزن دهنده تحت تأثیر silhouette و نرمال سازی است و در آن تأثیری ندارد. در حالی که در دیگر وزن دهنده ها که AD حضور داشته است، دست کم برای دو مورد از چهار روش تولید خوشبندی های پایه بهبود مشاهده می شود که نشان دهنده مؤثر بودن روش ارائه شده در این مقاله است.

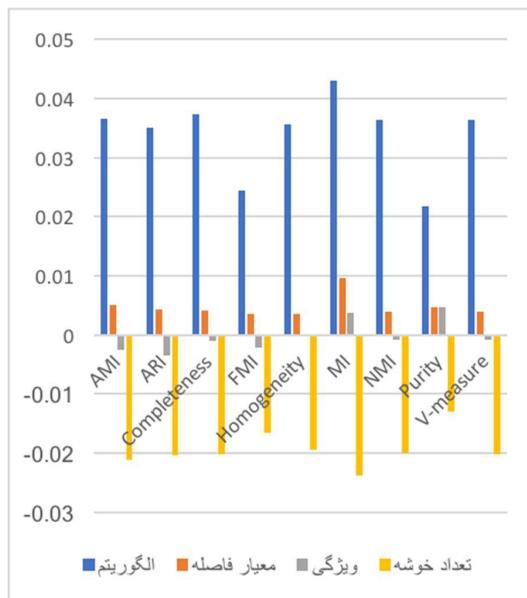




(نمودار-۶): تفاضل الگوریتم HGPA بدون وزن از HGPA وزن دار

با وزن دهی توسط w_5 و صحت های برآورد شده توسط AD

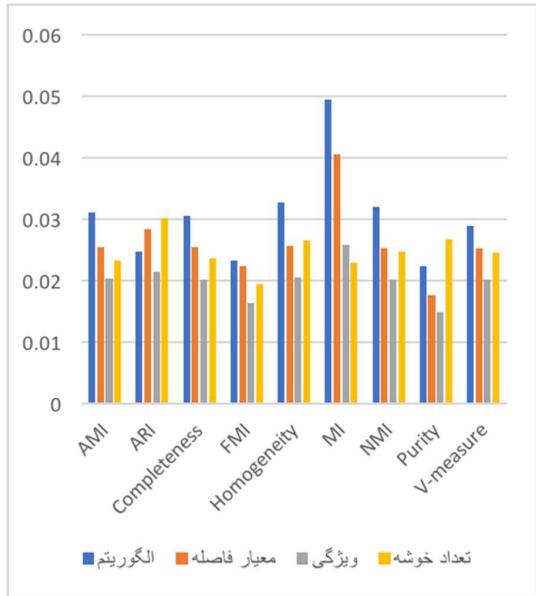
(Chart-6): Subtraction of HGPA from weighted HGPA by using w_5 and AD estimated accuracies



(نمودار-۴): تفاضل الگوریتم HGPA بدون وزن از HGPA وزن دار

با وزن دهی توسط w_3 و صحت های برآورد شده توسط AD

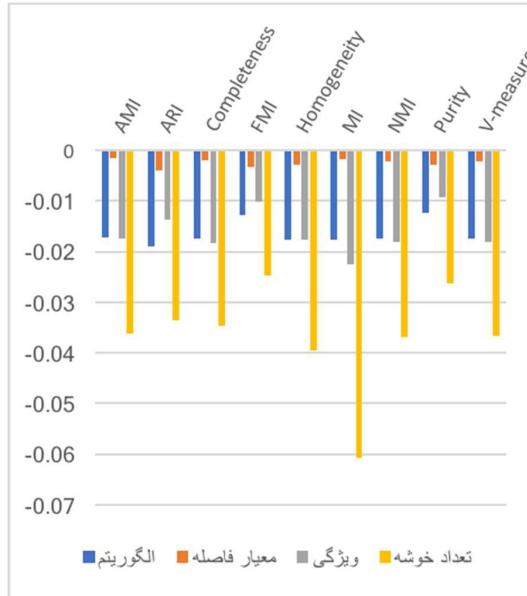
(Chart-4): Subtraction of HGPA from weighted HGPA by using w_3 and AD estimated accuracies



(نمودار-۷): تفاضل الگوریتم HGPA بدون وزن از HGPA وزن دار

با وزن دهی توسط w_6 و صحت های برآورد شده توسط AD

(Chart-7): Subtraction of HGPA from weighted HGPA by using w_6 and AD estimated accuracies



(نمودار-۵): تفاضل الگوریتم HGPA بدون وزن از HGPA وزن دار

با وزن دهی توسط w_4 و صحت های برآورد شده توسط AD

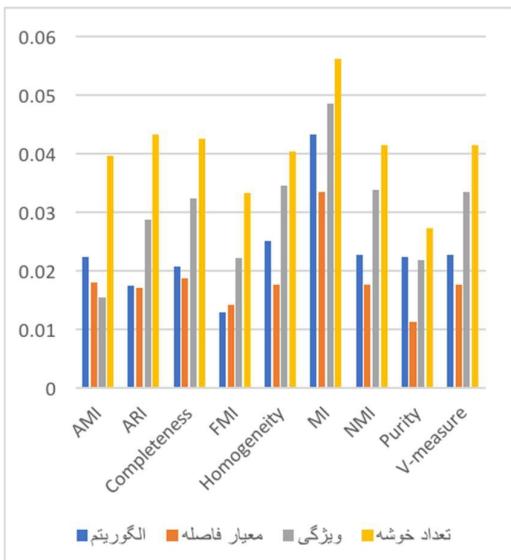
(Chart-5): Subtraction of HGPA from weighted HGPA by using w_4 and AD estimated accuracies

با توجه به بهبود صحت خوش بندی نهایی با وزن دهنده توسط صحت های برآورده در شرایط خاصی، برای پژوهش های آینده می توان به ساخت مارژولی جهت ابریادگیری پرداخت. به طوری که به وسیله ابریوزگری های^۱ استخراج شده از داده ها و خوش بندی های پایه، میزان بهبود خوش بندی نهایی را پیش بینی کند. برای این منظور باید روی مجموعه داده های بیشتری آزمایش صورت بگیرد تا ابریوزگری های مؤثر متفاوت، مانند ابعاد داده، تعداد داده، صحت خوش های پایه، نوع خوش بندی و غیره شناسایی شوند.

7- References

- [1] J. Han, M. Kamber, and J. Pei, *Data mining: concepts and techniques*. Elsevier, 2011.
- [2] S. Vega-Pons and J. Ruiz-Shulcloper, “a Survey of Clustering Ensemble Algorithms,” *Int. J. Pattern Recognit. Artif. Intell.*, vol. 25, no. 03, pp. 337–372, 2011.
- [3] J. Kittler, M. Hatef, R. P. W. Duin, and J. Matas, “On combining classifiers,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 3, pp. 226–239, 1998.
- [4] S.-B. Cho and J. H. Kim, “Combining multiple neural networks by fuzzy integral for robust classification,” *Syst. Man Cybern. IEEE Trans.*, vol. 25, no. 2, pp. 380–384, 1995.
- [5] J. Franke and E. Mandl, “A comparison of two approaches for combining the votes of cooperating classifiers,” in *Pattern Recognition, 1992. Vol. II. Conference B: Pattern Recognition Methodology and Systems, Proceedings., 11th IAPR International Conference on*, 1992, pp. 611–614.
- [6] L. K. Hansen and P. Salamon, “Neural network ensembles,” *IEEE Trans. Pattern Anal. Mach. Intell.*, no. 10, pp. 993–1001, 1990.
- [7] S. Hashem and B. Schmeiser, “Improving model accuracy using optimal linear combinations of trained neural networks,” *Neural Networks, IEEE Trans.*, vol. 6, no. 3, pp. 792–794, 1995.
- [8] J. Kittler, “Improving recognition rates by classifier combination: A theoretical framework,” *DAC and IS, editors, Progress in Handwriting Recognition*, pp. 231–248, 1997.
- [9] D. J. Miller and L. Yan, “Critic-driven ensemble classification,” *Signal Process. IEEE Trans.*, vol. 47, no. 10, pp. 2833–2844, 1999.

۷- مراجع



(نمودار-۸): تفاضل الگوریتم HGPA بدون وزن از HGPA وزن دار

با وزن دهنده توسط w_7 و صحت های برآورده شده توسط AD(Chart-8): Subtraction of HGPA from weighted HGPA by using w_7 and AD estimated accuracies

۶- نتیجه گیری و کارهای آینده

در این مقاله روشی برای ترکیب خوش بندی های پایه با وزن دهنده مبتنی بر الگوریتم AD پیشنهاد شد و براساس نه معیار ارزیابی مختلف روی سیزده مجموعه داده متفاوت نشان داده شد که روش ترکیب وزن دار پیشنهادی بهتر از روش های ترکیب بدون وزن قادر به خوش بندی داده ها است. از نتایج به دست آمده مشاهده می شود که استفاده از AD در برآورد صحت خوش بندی ها مناسب تر است. از میان روش های وزن دهنده مختلفی که برای ابریال ها پیشنهاد شد، نتایج هر یک از وزن دهنده ها بر روی مجموعه داده خاصی توانستند بهبود قابل توجهی ایجاد کنند که حاکی از واپسگی نوع وزن دهنده به ویژگی های ساختاری هر مجموعه داده است؛ اما به طور متوسط بهترین نتایج مربوط به زمانی است که وزن دهنده به وسیله w_1 یعنی تنها با استفاده از صحت های برآورده AD است. ضعیفترین عملکرد را زمانی مشاهده می کنیم که وزن دهنده تحت تأثیر silhouette نرمال سازی می باشد و AD در آن تأثیری ندارد.

در این پژوهش HGPA وزن دار با روش های متفاوت تولید خوش بندی های پایه مورد بررسی قرار گرفت. نتایج آزمایش HGPA با هفت نوع وزن دهنده متفاوت با صحت های AD HGPA نشان داد که در بیشتر موارد بیشترین بهبود برای تولید خوش بندی های پایه با استفاده از الگوریتم های متفاوت رخ داده است.

- فصلنامه
پژوهش در زمینه حسابات
و مدیریت اطلاعات
- سال ۱۳۹۹ شماره ۲ پیاپی ۴۴
- ۹۸
- international conference on Knowledge discovery in data mining*, 2005, pp. 70–77.
- [24] A. Topchy, A. K. Jain, and W. Punch, “Combining multiple weak clusterings,” in *Data Mining, 2003. ICDM 2003. Third IEEE International Conference on*, 2003, pp. 331–338.
- [25] L. I. Kuncheva, S. T. Hadjitodorov, and Others, “Using diversity in cluster ensembles,” in *Systems, man and cybernetics, 2004 IEEE international conference on*, 2004, vol. 2, pp. 1214–1219.
- [26] B. Minaei-Bidgoli, A. Topchy, and W. F. Punch, “Ensembles of partitions via data resampling,” in *Information Technology: Coding and Computing, 2004. Proceedings. ITCC 2004. International Conference on*, 2004, vol. 2, pp. 188–192.
- [27] A. Topchy, A. K. Jain, and W. Punch, “Clustering ensembles: Models of consensus and weak partitions,” *Pattern Anal. Mach. Intell. IEEE Trans.*, vol. 27, no. 12, pp. 1866–1881, 2005.
- [۲۸] علیزاده حسین، مشکی محسن، پروین حمید و مینایی بیدگلی بهروز، ”خوشه‌بندی ترکیبی مبتنی بر زیرمجموعه‌ای از خوشه‌های اولیه“، پژوهش علوم و تاریخ، ۱۳۸۹، ۷(۱) : ۱۹-۳۲.
- [28] H. Alizadeh, M. Moshki, H. Parvin, B. Minaei Bidgoli, “Clustering ensemble based on combination of subset of primary clusters,” *Signal and Data Processing*, vol. 7, no. 1, pp. 19-32, 2010.
- [29] L. Franek and X. Jiang, “Ensemble clustering by means of clustering embedding in vector spaces,” *Pattern Recognit.*, vol. 47, no. 2, pp. 833–842, 2014.
- [۳۰] میرزایی عبدالرضا، ”ترکیب خوشه‌بندی‌های سلسه‌مراتبی با تأکید بر حفظ اطلاعات ساختاری خوشه‌بندی‌های پایه“، دانشگاه صنعتی امیر کبیر، ۱۳۸۸.
- [30] A. Mirzaei, "Combining hierarchical clusterings with emphasis on retaining the structural contents of the base clusterings," *PhD dissertation*, Amirkabir University of Technology, 2009.
- [31] J. H. Friedman and J. J. Meulman, “Clustering objects on subsets of attributes,” *J. R. Stat. Soc. Ser. B (Statistical Methodol.)*, vol. 66, no. 4, pp. 815–849, 2004.
- [32] C. Domeniconi, D. Gunopoulos, S. Ma, B. Yan, M. Al-Razgan, and D. Papadopoulos, “Locally adaptive metrics for clustering high dimensional
- [10] K. W. De Bock, K. Coussement, and D. Van den Poel, “Ensemble classification based on generalized additive models,” *Comput. Stat. Data Anal.*, vol. 54, no. 6, pp. 1535–1546, 2010.
- [11] C. Domeniconi and M. Al-Razgan, “Weighted cluster ensembles: Methods and analysis,” *ACM Trans. Knowl. Discov. from Data*, vol. 2, no. 4, pp. 17, 2009.
- [12] A. Strchl and J. Ghosh, “Cluster ensembles---a knowledge reuse framework for combining multiple partitions,” *J. Mach. Learn. Res.*, vol. 3, no. Dec, pp. 583–617, 2002.
- [13] H. Amirkhani and M. Rahmati, “Agreement/disagreement based crowd labeling,” *Appl. Intell.*, vol. 41, no. 1, pp. 212–222, Jul. 2014.
- [14] N. Littlestone and M. K. Warmuth, “The weighted majority algorithm,” in *Foundations of Computer Science, 1989., 30th Annual Symposium on*, 1989, pp. 256–261.
- [15] S. B. Kotsiantis, I. Zaharakis, and P. Pintelas, “Supervised machine learning: A review of classification techniques.” *Emerging Artificial Intelligence Applications in Computer Engineering*, vol. 160, pp. 3-24, 2007.
- [16] T. G. Dietterich, “Ensemble methods in machine learning,” in *International workshop on multiple classifier systems*, 2000, pp. 1–15.
- [17] T. Windeatt, “Vote counting measures for ensemble classifiers,” *Pattern Recognit.*, vol. 36, no. 12, pp. 2743–2756, 2003.
- [18] S. Wang, A. Mathew, Y. Chen, L. Xi, L. Ma, and J. Lee, “Empirical analysis of support vector machine ensemble classifiers,” *Expert Syst. Appl.*, vol. 36, no. 3, pp. 6466–6476, 2009.
- [19] A. J. C. Sharkey, *Combining artificial neural nets: ensemble and modular multi-net systems*. Springer Science & Business Media, 2012.
- [20] A. L. N. Fred and A. K. Jain, “Data clustering using evidence accumulation,” in *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, 2002, vol. 4, pp. 276–280.
- [21] A. Topchy, A. K. Jain, and W. Punch, “A mixture model for clustering ensembles,” in *Society for Industrial and Applied Mathematics. Proceedings of the SIAM International Conference on Data Mining*, 2004, pp. 379.
- [22] S. Dudoit and J. Fridlyand, “Bagging to improve the accuracy of a clustering procedure,” *Bioinformatics*, vol. 19, no. 9, pp. 1090–1099, 2003.
- [23] D. Gondek and T. Hofmann, “Non-redundant clustering with conditional ensembles,” in *Proceedings of the eleventh ACM SIGKDD*



- [43] A. Litifi Pakdehi, N. Daneshpour, "Cluster ensemble selection using voting," *Signal and Data Processing*, vol. 15, no. 4, pp. 17-30, 2019.
- [44] Z. Yu and H. S. Wong, "Class discovery from gene expression data based on perturbation and cluster ensemble," *IEEE Trans. Nanobioscience*, vol. 8, no. 2, pp. 147-160, 2009.
- [45] H. Liu, M. Shao, S. Li, and Y. Fu, "Infinite ensemble clustering," *Data Min. Knowl. Discov.*, vol. 32, no. 2, pp. 385-416, 2018.
- [46] D. Huang, C.-D. Wang, and J.-H. Lai, "Locally weighted ensemble clustering," *IEEE Trans. Cybern.*, vol. 48, no. 5, pp. 1460-1473, 2018.
- [47] C. Zhong, L. Hu, X. Yue, T. Luo, Q. Fu, and H. Xu, "Ensemble clustering based on evidence extracted from the co-association matrix," *Pattern Recognit.*, vol. 92, pp. 93-106, 2019.
- [48] F. Li, Y. Qian, J. Wang, C. Dang, and L. Jing, "Clustering ensemble based on sample's stability," *Artif. Intell.*, vol. 273, pp. 37-55, 2019.
- [49] G. Karypis and V. Kumar, "A fast and high quality multilevel scheme for partitioning irregular graphs," *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 359-392, 1998.
- [50] G. Karypis, R. Aggarwal, V. Kumar, and S. Shekhar, "Multilevel hypergraph partitioning: applications in VLSI domain," *IEEE Trans. Very Large Scale Integr. Syst.*, vol. 7, no. 1, pp. 69-79, 1999.
- [51] L. Rokach and O. Maimon, "Clustering methods," in *Data mining and knowledge discovery handbook*, Springer, 2005, pp. 321-352.
- [52] P. J. Rousseeuw, "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis," *J. Comput. Appl. Math.*, vol. 20, pp. 53-65, 1987.



صدیقه وحیدی فردوسی، در سال ۱۳۹۳ مدرک کارشناسی خود را در رشته مهندسی کامپیوتر با گرایش نرم افزار از دانشگاه قم دریافت و مدرک کارشناسی ارشد خود را نیز از همان دانشگاه در رشته مهندسی فناوری اطلاعات با گرایش تجارت الکترونیک در سال ۱۳۹۵ اخذ کرد. موضوع پایان نامه کارشناسی ارشد ایشان ترکیب وزن دار خوش بندی ها با هدف افزایش دقت خوش بندی نهایی بوده است. زمینه های پژوهشی مورد علاقه وی یادگیری ماشین و داده کاوی است.

نشانی رایانامه ایشان عبارت است از:

s.vahidi@stu.qom.ac.ir

- [42] پروین حمید، "خوش بندی ترکیبی با وزن دهنده توأم خوش ها و ویژگی ها"، دانشگاه علم و صنعت ایران، ۱۳۹۲.
- [43] H. Parvin, "Clustering ensembles with weighting clusters and features," *PhD dissertation*, Iran University of Science & Technology, 2013.
- [44] لطیفی پاکدهی علیرضا، دانشپور نگین، "انتخاب اعضای ترکیب در خوش بندی ترکیبی با استفاده از رأی گیری." پردازش علائم و داده ها، ۱۳۹۷، ۱۵ (۴): ۳۰-۱۷.

حسین امیرخانی، مدرک کارشناسی خود را در رشته مهندسی کامپیوتر با گرایش نرمافزار در سال ۱۳۸۶ از دانشگاه اصفهان دریافت کرده است. ایشان در سال ۱۳۸۸ مدرک کارشناسی ارشد خود

را در رشته مهندسی کامپیوتر با گرایش هوش مصنوعی از دانشگاه صنعتی امیرکبیر و در سال ۱۳۹۳ نیز مدرک دکترای خود را از همان دانشگاه و در همان رشته کسب کرده است. وی هم‌اکنون استادیار گروه مهندسی کامپیوتر و فناوری اطلاعات دانشگاه قم و زمینه‌های پژوهشی مورد علاقه ایشان یادگیری ماشین و پردازش زبان‌های طبیعی است.

نشانی رایانمه ایشان عبارت است از:

amirkhani@qom.ac.ir



فصلنامه
پژوهش‌های
دانشگاه قم

سال ۱۳۹۹ شماره ۲ پیاپی ۴۴

۱۰۰