

دریافت مقاله: ۱۳۹۴/۲/۱۵

پذیرش مقاله: ۱۳۹۴/۱۰/۲۵

## ارائه یک الگوریتم خوشه‌بندی برای داده‌های دسته‌ای با ترکیب معیارها

مریم نبی‌لو،<sup>۱</sup> نگین دانشپور<sup>۲\*</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد، دانشکده مهندسی کامپیوتر، دانشگاه تربیت دبیر شهید رجایی

mnabiloo13@yahoo.com

<sup>۲</sup> استادیار دانشکده مهندسی کامپیوتر، دانشگاه تربیت دبیر شهید رجایی

ndaneshpour@srttu.edu

### چکیده

خوشه‌بندی یکی از تکنیک‌های اصلی داده‌کاوی است. خوشه‌بندی فرایندی است که مجموعه داده‌ها را داخل گروه‌هایی طبقه‌بندی می‌کند. در خوشه‌بندی، داده‌های موجود در یک خوشه بیشترین شباهت را به هم دارند و داده‌های موجود در دو خوشه متفاوت، بیشترین تفاوت را با هم دارند. الگوریتم‌های خوشه‌بندی با توجه به نوع داده‌ها به دو دسته تقسیم می‌شوند: الگوریتم‌های خوشه‌بندی داده‌های عددی و الگوریتم‌های خوشه‌بندی دسته‌ای. الگوریتم‌های خوشه‌بندی داده‌های دسته‌ای به دلیل ماهیت و کاربرد این داده‌ها نسبت به الگوریتم‌های خوشه‌بندی داده‌های عددی، اهمیت بیشتری دارند. هر یک از این الگوریتم‌ها با توجه به نوع داده (عددی یا دسته‌ای) از معیارهای شباهت متفاوتی در خوشه‌بندی استفاده می‌کنند. در این مقاله، ابتدا به بررسی ماهیت این نوع داده‌ها پرداخته شده و سپس معیارهای شباهت و الگوریتم‌های خوشه‌بندی مطرح شده در این حوزه را بررسی کرده و در نهایت، یک روش جدید در خوشه‌بندی با استفاده از ترکیب معیارهای شباهت برای داده‌های دسته‌ای ارائه می‌کنیم. آزمایش‌ها نشان می‌دهد که روش ارائه‌شده در این مقاله، توانسته است که نتایج حاصل از خوشه‌بندی را بهبود ببخشد.

واژه‌های کلیدی: داده‌کاوی، خوشه‌بندی، داده‌های دسته‌ای، معیار فاصله، معیار چگالی.

## ۱. مقدمه

امروزه با توجه به حجم بالای داده‌ها در علوم مختلف، می‌باید اطلاعات مفید از داخل مجموعه داده‌ها استخراج شود [۱]. زمانی که حجم داده‌ها افزایش می‌یابد، کاربران هرچند زبردست و باتجربه باشند، نمی‌توانند الگوهای مفید را در میان حجم انبوه داده‌ها تشخیص دهند [۲]. یکی از تکنیک‌هایی که طی سالیان اخیر برای انجام این امر مطرح شده، تکنیک داده‌کاوی است [۳].

داده‌کاوی به‌عنوان یک ابزار قوی برای تولید اطلاعات و دانش از داده‌های خام، شناخته شده است. برای داده‌کاوی تعاریف مختلفی وجود دارد که معمول‌ترین آن عبارت از: استخراج دانش و اطلاعات از یک پایگاه داده بسیار بزرگ و پیچیده (تبدیل داده‌ها به دانش) و همچنین کشف الگوهای پنهان بین داده‌ها [۴-۶]. بنابراین، ابزارهای داده‌کاوی الگوهای پنهانی را کشف و پیش‌بینی می‌کنند که متخصصان ممکن است به‌دلیل اینکه این اطلاعات و الگوها خارج از انتظار آن‌ها باشد، آن‌ها را مدنظر قرار ندهند و به آن‌ها دست نیابند [۷].

تکنیک‌های مختلفی نظیر دسته‌بندی<sup>۱</sup> و رگرسیون<sup>۲</sup> و خوشه‌بندی<sup>۳</sup> برای داده‌کاوی ارائه شده است. خوشه‌بندی داده‌ها [۸-۱۰] به معنی گروه‌بندی نمونه‌ها در گروه‌های شبیه‌به‌هم است؛ به طوری که نمونه‌های هر گروه (خوشه) حداکثر تشابه را با یکدیگر و حداکثر فاصله را با نمونه‌های گروه‌های دیگر داشته باشند [۱۱].

روش‌های بسیاری برای خوشه‌بندی داده‌ها پیشنهاد شده است. اکثر این روش‌ها روی جنبه‌های خاصی از داده‌ها تأکید کرده و بنابراین، روی مجموعه داده‌های خاصی کارایی خوبی دارند. الگوریتم‌های زیادی در این زمینه معرفی شده است. دسته‌بندی متفاوتی برای الگوریتم‌های خوشه‌بندی وجود دارد که رویکرد زیر جامع‌ترین دسته‌بندی این روش‌هاست [۱۱]:

- الگوریتم‌های تفکیکی<sup>۴</sup>

- الگوریتم‌های سلسله‌مراتبی<sup>۵</sup>

- الگوریتم‌های مبتنی بر چگالی<sup>۶</sup>

- الگوریتم‌های مبتنی بر گرید<sup>۷</sup>

- الگوریتم‌های مبتنی بر مدل<sup>۸</sup>

همچنین، الگوریتم‌های خوشه‌بندی با توجه به نوع داده‌ها تقسیم‌بندی دیگری نیز دارند: الگوریتم‌های خوشه‌بندی داده‌های عددی و الگوریتم‌های خوشه‌بندی داده‌های دسته‌ای.

الگوریتم‌های خوشه‌بندی داده‌های دسته‌ای به دلیل ماهیت و کاربرد این داده‌ها نسبت به الگوریتم‌های خوشه‌بندی داده‌های عددی، اهمیت بیشتری دارند. داده‌های دسته‌ای شامل حداقل دو مقدار متمایزند که هیچ‌گونه تقدم یا تأخر در آن‌ها وجود ندارد. به عبارت دیگر، میان مقادیر داده‌های دسته‌ای نمی‌توان ترتیب خاصی در نظر گرفت. به همین سبب، چالش‌هایی در زمینه خوشه‌بندی داده‌های دسته‌ای مطرح است که شامل موارد زیر است [۱۲ و ۱۳].

- عدم امکان تفکیک‌پذیری روی دامنه برای داده‌های دسته‌ای: برای داده‌های عددی نماینده هر خوشه اغلب شامل میانگین داده‌ها در دامنه صفات هر خوشه است. محاسبه میانگین برای داده‌های دسته‌ای، غیرعملی است [۱۴].
- عدم امکان اندازه‌گیری شباهت برای داده‌های دسته‌ای به روش معمول داده‌های عددی: توابع فاصله نظیر فاصله منتهن و فاصله اقلیدسی برای داده‌های دسته‌ای به دلیل اینکه هیچ ترتیبی بین مقادیر داده وجود ندارد، کارایی ندارد [۱۵].

تابه‌حال، الگوریتم‌های بسیاری برای خوشه‌بندی داده‌های دسته‌ای ارائه شده است که هر یک از آن‌ها معیار شباهت جدیدی را برای این نوع داده‌ها معرفی می‌کند. این معیارهای شباهت ارائه‌شده، هر یک مزایا و معایب خاص خود را دارد. به همین منظور، یعنی برای کاهش معایب معیارهای شباهت در این مقاله، دو تا از معروف‌ترین این معیارها را با هم ترکیب کرده‌ایم.

در این مقاله، ابتدا به تعریف و تبیین مفهوم خوشه‌بندی،

به‌خصوص خوشه‌بندی داده‌های پرداخته شده است. در ادامه، روش ترکیبی جدیدی بر پایه ترکیب دو الگوریتم خوشه‌بندی سلسله‌مراتبی و خوشه‌بندی تفکیکی برای خوشه‌بندی داده‌های ارائه شده است. آزمایش‌ها نشان می‌دهد که روش یادشده موجب بهبود در نتایج خوشه‌بندی داده‌ها می‌شود. تمام الگوریتم‌های خوشه‌بندی ترکیبی در حوزه داده‌های عددی قرار دارند و از روش‌های ترکیبی در زمینه خوشه‌بندی داده‌های دسته‌ای استفاده نشده است. بنابراین، در این مقاله یک رویکرد جدید بر مبنای ترکیب الگوریتم‌های موجود برای خوشه‌بندی داده‌های دسته‌ای ارائه می‌شود.

## ۲. پیشینه پژوهش

خوشه‌بندی فرایندی است که مجموعه داده‌ها را داخل گروه‌هایی طبقه‌بندی می‌کند. الگوریتم‌های متنوعی برای خوشه‌بندی داده‌ها ارائه شده است. هریک از این الگوریتم‌ها معیار شباهت خاص خود را دارد. این الگوریتم‌ها بسته به معیار شباهتی که استفاده می‌کنند و روشی که در خوشه‌بندی داده‌ها دارند، توانایی متفاوتی در مدیریت داده‌های پرت و دارای نویز دارند. الگوریتم  $k$ -means یکی از اولین الگوریتم‌های خوشه‌بندی تفکیکی بود که به‌صورت وسیع و مؤثر، در خوشه‌بندی مجموعه داده‌های مورد استفاده قرار می‌گرفت [۱۶]. هرچند استفاده از انواع الگوریتم‌های  $k$ -means فقط محدود به داده‌های عددی است، بیشتر داده‌های جهان واقعی از نوع داده‌های دسته‌ای هستند. خوشه‌بندی داده‌های دسته‌ای به دلیل ماهیت و کاربرد این داده‌ها نسبت به خوشه‌بندی داده‌های عددی، اهمیت بیشتری دارند. برای خوشه‌بندی داده‌های دسته‌ای الگوریتم  $k$ -means به الگوریتم  $k$ -modes [۱۷] گسترش یافت. در این الگوریتم، محدودیت در خصوص داده‌های عددی از بین رفت و این الگوریتم قادر بود که به‌طور مؤثر روی داده‌های دسته‌ای اجرا شود. در این الگوریتم، معیار شباهت جدیدی برای داده‌های دسته‌ای معرفی شد. معیار شباهت بین دو داده دسته‌ای برابر تعداد مقادیر صفات غیرمشابه میان آن دو داده است. الگوریتم خوشه‌بندی  $k$ -modes میانگین خوشه‌ها را با مفهوم  $mode$  جایگزین کرد و از روش‌های مبتنی بر فرکانس برای به‌روزرسانی

این  $mode$ ها در فرایند خوشه‌بندی استفاده کرد تا به این وسیله، تابع هدف خوشه‌بندی را مینیمم کند. الگوریتم خوشه‌بندی  $k$ -modes چندین بار با مقادیر متفاوت  $mode$ های اولیه اجرا می‌شود تا پایداری راه حل خوشه‌بندی به‌دست آید. اما این الگوریتم همانند الگوریتم  $k$ -means توانایی مدیریت داده‌های پرت و نویز را در فرایند خوشه‌بندی ندارد و همچنین جواب نهایی به انتخاب مراکز اولیه خوشه‌ها وابسته است. علاوه بر این، الگوریتم  $fuzzy\ k$ -modes نیز ارائه شد که در آن، داده‌ها می‌توانند با یک درجه عضویت، عضوی از تمام خوشه‌ها باشند [۱۸]. این امر مهم‌ترین مزیت این الگوریتم است و موجب شده این الگوریتم برای داده‌های دارای نویز مناسب باشد. اما این الگوریتم نیز مانند الگوریتم  $k$ -modes توانایی مدیریت داده‌های پرت را ندارد. در ادامه، تعدادی از الگوریتم‌های خوشه‌بندی مختص داده‌های دسته‌ای را بررسی می‌کنیم:

الگوریتم خوشه‌بندی  $weighting\ k$ -modes [۳]، در واقع همان الگوریتم  $k$ -modes است با این تفاوت که به هر صفت، وزنی اختصاص داده است. دلیل وزن دادن به صفات در این الگوریتم، این است که ممکن است بعضی از صفات در مجموعه داده‌ها اهمیت بیشتری داشته باشند. این‌گونه صفات در فرایند خوشه‌بندی، نقش بیشتری را ایفا می‌کنند. وزن یک صفت معمولاً عددی بین ۰ تا ۱ است و صفتی که اهمیت بیشتری داشته باشد، وزن بیشتری دارد. مکانیزم وزندهی به صفات در این الگوریتم، موجب بهبود کیفیت خوشه‌بندی می‌گردد. اما این الگوریتم توانایی مدیریت داده‌های پرت و دارای نویز را ندارد.

الگوریتم خوشه‌بندی  $weighting\ fuzzy\ k$ -modes [۱۹]، مدل فازی الگوریتم  $weighting\ k$ -modes است. این الگوریتم در محاسبه شباهت به هر صفت، بسته به اهمیت آن یک وزن در بازه ۰ و ۱ می‌دهد. در این الگوریتم، هر داده می‌تواند با یک درجه عضویت، عضوی از تمام خوشه‌ها باشد. مزیت این الگوریتم نسبت به الگوریتم  $weighting\ k$ -modes این است که این الگوریتم توانایی مدیریت داده‌های دارای نویز را دارد. اما همچنان توانایی مدیریت داده‌های پرت را ندارد و مشکل انتخاب مراکز اولیه خوشه‌ها در آن، به قوت خود باقی است.

با یک بار بررسی مجموعه داده‌ها می‌تواند به نتیجه خوشه‌بندی خوبی دست یابد و به منظور افزایش کیفیت خوشه‌ها می‌توان این عمل را چندین بار تکرار کرد. این الگوریتم داده‌های پرت و نویزدار را نیز به خوبی مدیریت می‌کند. عیب این الگوریتم تعیین مقدار آستانه برای مقایسه شباهت است که این مقدار به صورت تجربی به دست می‌آید.

الگوریتم خوشه‌بندی COOLCAT [۲۳] از مفهوم آنتروپی برای گروه‌بندی رکوردها استفاده می‌کند. این الگوریتم، یک الگوریتم افزایشی با هدف مینیمم‌سازی آنتروپی موردانتظار برای خوشه‌هاست. با داشتن یک مجموعه از خوشه، الگوریتم COOLCAT نقطه بعدی در مجموعه نقاط داده را با مینیمم‌سازی آنتروپی کلی موردانتظار، خوشه‌بندی می‌کند. الگوریتم خوشه‌بندی COOLCAT بدون هیچ پیش‌پردازشی روی مجموعه داده‌ها، خوشه‌بندی را انجام می‌دهد. بنابراین الگوریتم COOLCAT برای داده‌های جریانی مناسب است، اما توانایی مدیریت داده‌های پرت را ندارد.

الگوریتم خوشه‌بندی CLOPE [۲۴] برای هر خوشه، یک نمودار ستونی ترسیم می‌کند و بر مبنای مینیمم‌سازی نسبت طول به عرض مربوط به آن خوشه، خوشه‌بندی را انجام می‌دهد. الگوریتم خوشه‌بندی CLOPE سریع و مقیاس‌پذیر است و برای جداسازی تراکشن‌های با ابعاد بالا مناسب است. این الگوریتم به ترتیب داده‌های ورودی حساس نیست، اما توانایی مدیریت داده‌های پرت را ندارد.

در این بخش، الگوریتم‌های مختلفی که برای خوشه‌بندی داده‌های دسته‌ای معرفی شده بود، بررسی شد. هر یک از این الگوریتم‌ها از یک معیار شباهت برای خوشه‌بندی داده‌ها استفاده می‌کند و برای موقعیت خاصی ارائه شده است. همان‌طور که بیان شد، هر یک از الگوریتم‌ها توانایی متفاوتی در مدیریت داده‌های پرت دارد. در این مقاله، روشی ترکیبی برای خوشه‌بندی بهتر داده‌های دسته‌ای بر مبنای الگوریتم سلسله‌مراتبی تجمعی و بهبود این روش ارائه می‌شود. در روش پیشنهادی، مشکل داده‌های پرت و دارای نویز در فرایند خوشه‌بندی، تا حد زیادی مرتفع شد.

الگوریتم خوشه‌بندی k-modes based on entropy [۲۰]، مشابه الگوریتم خوشه‌بندی k-modes است، با این تفاوت که این الگوریتم به جای استفاده از معیار شباهت overlap برای محاسبه شباهت میان اشیاء داده و مراکز خوشه، از معیار آنتروپی برای محاسبه شباهت استفاده می‌کند و با مینیمم‌سازی آنتروپی کلی موردانتظار، خوشه‌بندی را انجام می‌دهد. این الگوریتم که از جدیدترین نسخه‌های ارائه‌شده از الگوریتم k-modes است، قابلیت فراوانی در فرایند خوشه‌بندی دارد. این الگوریتم همچنین توانسته است مشکل مدیریت داده‌های پرت و دارای نویز را که الگوریتم k-modes در آن ناتوان بود، به خوبی حل کند. اما همچنان مشکل انتخاب مراکز اولیه خوشه‌ها در این الگوریتم نیز به قوت خود باقی است.

در الگوریتم خوشه‌بندی fuzzy k-prototype [۲۱]، k عدد prototype از میان اشیاء داده به وسیله مینیمم‌سازی تابع فاصله انتخاب می‌شود. این الگوریتم مانند الگوریتم خوشه‌بندی fuzzy k-modes عمل می‌کند، با این تفاوت که مفهوم mode با مفهوم prototype جایگزین شده است. انگیزه این الگوریتم بررسی اهمیت صفات مختلف در فرایند خوشه‌بندی است. این الگوریتم هم صفات عددی و هم صفات دسته‌ای را با وزنی متفاوت در فرایند خوشه‌بندی سهم می‌کند. این امر موجب می‌شود که این الگوریتم برای داده‌ها با صفات ترکیبی (صفات عددی و دسته‌ای) مناسب عمل کند. اما همچنان توانایی مدیریت داده‌های پرت را ندارد و همچنان مشکل انتخاب مراکز اولیه خوشه‌ها در این الگوریتم، به قوت خود باقی است.

الگوریتم خوشه‌بندی Squeezer [۲۲] شباهت بین داده‌ها را با یک معیار شباهت ویژه بررسی می‌کند تا داده موردنظر را در یکی از خوشه‌های موجود یا در یک خوشه جدید قرار دهد؛ که برای این کار از یک مقدار آستانه استفاده می‌کند. اگر میزان شباهت یک شیء به یک خوشه، بیشتر از مقدار آستانه بود، آن شیء در آن خوشه قرار می‌گیرد و در غیر این صورت، در یک خوشه جدید مستقر می‌شود. الگوریتم خوشه‌بندی Squeezer برای داده‌های دسته‌ای استفاده می‌شود و با دادن وزن بیشتر به صفات غیرمشابه، خوشه‌بندی را انجام می‌دهد. این الگوریتم تنها

## ۳. روش پیشنهادی

The Agglomerative Hierarchical Clustering  
 1: Initial Clustering, Put each data in a cluster  
 2: for each cluster do  
 3: Compute similarity measure for all pairs  
 4: Merge clusters  $C_i$  and  $C_j$ , in which  $\text{Similarity}(C_i, C_j)$  is max  
 5: end for  
 6: Determine # of clusters

شکل (۱): الگوریتم خوشه‌بندی سلسله‌مراتبی تجمعی

## ۲.۳. معیار فاصله و چگالی

معیار شباهت یا فاصله بین دو داده، یک مسئله چالش‌برانگیز در داده‌کاوی و تکنولوژی اطلاعات است. اندازه‌گیری شباهت می‌تواند به دو بخش تقسیم شود: اندازه‌گیری شباهت برای داده‌های دسته‌ای و داده‌های عددی. اندازه‌گیری شباهت برای داده‌های دسته‌ای به راحتی داده‌های عددی نیست. انواع فاصله نظیر فاصله منهن<sup>۳</sup> و فاصله اقلیدسی<sup>۴</sup> دو نوع از مهم‌ترین معیارهای فاصله برای داده‌های عددی هستند. شباهت یا فاصله برای داده‌های دسته‌ای، به راحتی داده‌های عددی محاسبه نمی‌شود و این مسئله یک موضوع چالش‌برانگیز است. این حقیقت که مقادیر مختلف در داده‌های دسته‌ای هیچ ترتیب خاصی ندارند، موجب می‌شود که مقایسه بین داده‌های دسته‌ای ممکن نباشد. به‌علاوه معیار شباهت مورداستفاده ممکن است به دامنه ویژه‌ای وابسته باشد. معیارهای مختلفی برای داده‌های دسته‌ای به کار می‌رود، اما بیشتر این معیارها در به‌کارگیری سایر ویژگی‌ها در مجموعه داده با شکست روبه‌رو می‌شوند. یک راه حل برای این مشکل ساخت یک منبع مشترک برای اندازه‌گیری شباهت برای همه مفاهیم مشترک است. اندازه‌گیری شباهت برای یک دامنه نسبت به دیگری متفاوت است. بنابراین برای تعیین شباهت داده‌ها باید درک صحیحی از دامنه داده‌ها به دست آورد. بنابراین اندازه‌گیری شباهت توسط یک خبره دامنه که به خوبی مفهوم دامنه را درک کرده، تعریف می‌شود. در بیشتر کاربردها کارشناس دامنه موجود نیست و کاربر ارتباط درونی بین داده‌ها را به خوبی درک نمی‌کند تا براین اساس شباهت و فاصله بین داده‌ها را محاسبه کند [۲۸ و ۲۹].

در اینجا از یک الگوریتم سلسله‌مراتبی با ترکیب خطی دو معیار شباهت، یعنی معیار همگرایی Jaccard و معیار چگالی استفاده می‌کنیم. به منظور تست الگوریتم پیشنهادی، دو مجموعه داده soybean data و zoo data انتخاب شده و الگوریتم پیشنهادی را روی آن‌ها به اجرا رسانده‌ایم. در ادامه، به شرح روش خوشه‌بندی سلسله‌مراتبی و هر یک از معیارهای شباهت همگرایی Jaccard و چگالی پرداخته شده و سپس روش ترکیبی پیشنهادی بیان می‌شود.

## ۱.۳. خوشه‌بندی سلسله‌مراتبی

روش خوشه‌بندی سلسله‌مراتبی یک روش تحلیل خوشه‌ای است که برای داده‌های کم حجم و داده‌های دسته‌ای به کار می‌رود. روش خوشه‌بندی سلسله‌مراتبی با یک سری از تقسیم‌بندی‌های متوالی انجام می‌شود. روش‌های خوشه‌بندی سلسله‌مراتبی بر دو دسته‌اند: سلسله‌مراتبی تجمعی<sup>۱</sup> و سلسله‌مراتبی تقسیمی<sup>۲</sup> [۲۵ و ۲۶]. روشی که در این مقاله از آن استفاده شده، خوشه‌بندی سلسله‌مراتبی تجمعی است. در مکانیزم اجرایی این روش، ابتدا با استفاده از یک معیار، فواصل دو خوشه تعریف می‌شود و سپس روش مناسب برای تشکیل خوشه‌ها و پیوند آن‌ها با یکدیگر انتخاب می‌گردد. در نهایت نیز تعداد خوشه‌های مناسب برای داده‌ها تعیین شده و خوشه‌بندی انجام می‌گیرد. خوشه‌بندی سلسله‌مراتبی با جداسازی هر مورد در یک خوشه جداگانه شروع می‌شود. در هر مرحله از تحلیل، جداسازی موارد تا جایی انجام می‌گیرد که شبیه‌ترین دو خوشه در هم ادغام شوند و در نهایت نیز تمامی موارد در یک درخت طبقه‌بندی کامل ادغام گردند. معیاری که خوشه‌بندی براساس آن انجام می‌پذیرد، فاصله است. مواردی که نزدیک یکدیگرند، در یک خوشه ادغام شده و مواردی که نسبت به یکدیگر فاصله بیشتری دارند، در خوشه‌های متفاوت قرار می‌گیرند [۲۷]. شکل (۱) الگوریتم خوشه‌بندی سلسله‌مراتبی تجمعی را بیان می‌کند.

3. Manhattan distance  
 4. Euclidean distance

1. Agglomerative  
 2. Divisive

برابر تعداد داده‌ها در مجموعه داده است، تقسیم می‌کند [۳۲]. در این مقاله، ترکیبی از دو معیار چگالی و فاصله جهت بهره‌وری از محاسن هر دو معیار استفاده شده است. در ادامه، دو معیار فاصله Jaccard و چگالی مبتنی بر فرکانس تکرار صفات، با یکدیگر ترکیب شده‌اند.

### ۳.۳. روش ترکیبی

هریک از دو معیار فاصله و چگالی، دارای مزایا و معایب مربوط به خود هستند. برای محاسبه شباهت دو داده، علاوه بر چگالی، معیار فاصله را هم در نظر می‌گیریم. اگر تنها معیار فاصله را در نظر بگیریم، ممکن است که داده‌های پرت<sup>۱</sup> در تصمیم‌گیری کیفیت خوشه‌بندی ما دخالت کند و اگر تنها معیار چگالی را در نظر بگیریم، بیشتر خوشه‌ها در یک ناحیه از مجموعه داده‌ها انتخاب می‌شوند. بنابراین برای اجتناب از بروز این گونه مشکلات در اینجا معیار چگالی با معیار فاصله برای پیدا کردن خوشه‌ها با هم ادغام می‌شوند [۳۲]. برای ترکیب معیارهای فاصله و چگالی، دو ترکیب بیان شده است که در هر دو، از یک فاکتور وزن‌دهی به این دو معیار به نام  $\eta$  استفاده می‌شود:

- ترکیب خطی  $d(1-\eta) + \eta \text{Dens}$  که در آن  $\eta \in [0,1]$  است.

- ترکیب غیرخطی  $d(1-\eta) \times \text{Dens} \eta$  که در آن  $\eta \in [0,1]$  است.

در اینجا به سه دلیل، از ترکیب خطی استفاده می‌کنیم:

- ترکیب غیرخطی ارائه شده با Log‌گیری، به رابطه خطی تبدیل می‌شود.

- ترکیب غیرخطی به مقدار فاکتور وزن‌دهی بسیار حساس است.

- در مطالعات انجام شده در این زمینه، اغلب از ترکیب خطی استفاده می‌شود.

شکل (۲) الگوریتم خوشه‌بندی سلسله‌مراتبی تجمیعی را بیان می‌کند.

این ترکیب در تحقیقات [۵ و ۳۲] برای یافتن مراکز اولیه خوشه‌ها استفاده شده است و در این مقاله، برای ترکیب دو معیار خوشه‌بندی استفاده می‌شود. در بخش بعد، ملاحظه می‌شود که ترکیب دو معیار چگالی و فاصله به‌عنوان معیار

ساده‌ترین راه برای پیدا کردن شباهت بین دو صفت دسته‌ای بر این اساس است که دو داده اگر یکسان باشند، شباهتشان ۱ است و در غیر این صورت، شباهتشان صفر است. برای داده‌های دسته‌ای چندمتغیره، شباهت بین آن‌ها برابر تعداد صفات یکسان میان آن‌هاست. این روش overlap نام دارد. میزان تشابه هر صفت برای اندازه‌گیری overlap بین ۰ و ۱ است؛ مقدار صفر به این معناست که دو داده مشابه نیستند و مقدار ۱ به این معناست که دو داده کاملاً یکسان‌اند. عیب بزرگ روش overlap این است که فرقی بین مقادیر مختلف یک صفت قائل نیست. در واقع همه صفات یکسان و غیریکسان میان دو داده، رفتار برابر دارند [۲۹]. رابطه (۱) این روش را نشان می‌دهد:

$$S_k(X_k, Y_k) = \begin{cases} 1 & \text{if } X_k = Y_k \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

که  $X_k$  و  $Y_k$  اشیاء داده در مجموعه عناصر داده هستند و  $S_k(X_k, Y_k)$  تابع شباهت بین این دو شیء را محاسبه می‌کند.

معیار فاصله دیگر در مورد داده‌های دسته‌ای، معیار همگرایی Jaccard است که با الهام از معیار overlap و بهبود آن توسعه یافته است. معیار همگرایی Jaccard تشابه بین داده‌ها را به صورت رابطه (۲) بررسی می‌کند:

$$\text{Sim}(X_k, Y_k) = \frac{|X_k \cap Y_k|}{|X_k \cup Y_k|} \quad (2)$$

که  $|X_k|$  تعداد عناصر  $X_k$  است و  $X_k$  یک شیء داده در مجموعه عناصر داده هستند. رابطه (۲) نزدیکی بین عناصر دو داده را محاسبه می‌کند و عناصری که در هر دو داده  $X_k, Y_k$  مشترک‌اند، به صورت  $|X_k \cap Y_k|$  محاسبه می‌شوند و بعد از نرمال‌سازی با پارامتر  $\theta$  مقایسه می‌شوند [۳۰ و ۳۱].

معیار دیگری که برای تعیین تشابه داده‌های دسته‌ای به‌کار می‌رود، معیار چگالی است. معیار چگالی فرکانس تکرار صفات را برای هر داده محاسبه می‌کند. برای محاسبه چگالی از رابطه (۳) استفاده می‌کنیم.

$$\text{Dens}(x) = \frac{1}{|U|} \sum_{y \in U} d(x, y) \quad (3)$$

که  $d(x, y)$  فرکانس تکرار دو داده را محاسبه می‌کند و در نهایت، برای نرمال‌سازی چگالی به‌دست‌آمده، آن را به  $|U|$  که

شبهات در خوشه‌بندی سلسله‌مراتبی، به نتایج بهتری نسبت به زمانی که تنها از یک معیار استفاده می‌شد، منجر می‌شود.

**The Agglomerative Hierarchical Algorithm.**

**Input:** Data set D, Number of Clusters k, Dimensions d:

**Output:** results.

**Begin**

Initial Clustering, Put each data in a cluster

**for** each cluster do

Find two cluster that  $\text{sim}(C_i, C_j) =$

$\max_{1 \leq i, j \leq d} (\eta * \text{Jaccard\_sim}(x_i, x_j) +$

$(1 - \eta) * \text{Overlay\_sim}(x_i, x_j))$

Merge clusters  $C_i$  and  $C_j$ , in which

$\text{Similarity}(C_i, C_j)$  is max

Compute Center new cluster

**end for**

Determine # of clusters

**End**

شکل (۲): الگوریتم خوشه‌بندی سلسله‌مراتبی تجمیعی

**۴.۳. نتایج پژوهش**

در این بخش، ابتدا خوشه‌بندی انجام شده برای الگوریتم سلسله‌مراتبی با ترکیب معیارها (فاصله و چگالی) بررسی و سپس کارایی این الگوریتم با حالت بدون ترکیب هر یک از معیارهای فاصله و چگالی با هم مقایسه می‌شود. این پیاده‌سازی روی دو مجموعه داده soybean data و zoo data انجام شده که شرح آن‌ها به صورت زیر است [۳۳]:

soybean data این مجموعه داده شامل ۴۷ رکورد است که هر کدام توسط ۳۶ صفت توصیف می‌شوند. و هر رکورد با چهار ویژگی برچسب زده شده Diaporthe Stem Canker, Phytophthora Rot Charcoal Rot, Rhizoctonia Root Rot, به غیر از Phytophthora Rot شامل ۱۷ رکورد است.

zoo data این مجموعه داده شامل ۱۰۱ نمونه است که هر کدام توسط ۱۷ صفت باینری توصیف می‌شوند. داده‌ها در ۷ خوشه طبقه‌بندی می‌شوند.

برای ارزیابی کارایی یک الگوریتم خوشه‌بندی، فاکتورهایی

بدین صورت ارائه شده است:

**Accuracy:** با فرض اینکه  $C = \{C_1, \dots, C_k\}$  مجموعه خوشه‌ها و  $a_i$  برابر تعداد داده‌هایی است که در خوشه صحیح خود ( $C_i$ ) قرار گرفته‌اند و مجموعه داده دارای  $n$  شیء و  $K$  تعداد خوشه‌ها باشد، میزان Accuracy با رابطه (۴) اندازه‌گیری می‌شود [۳۴]:

$$AC = \frac{\sum_{i=1}^K a_i}{n} \quad (4)$$

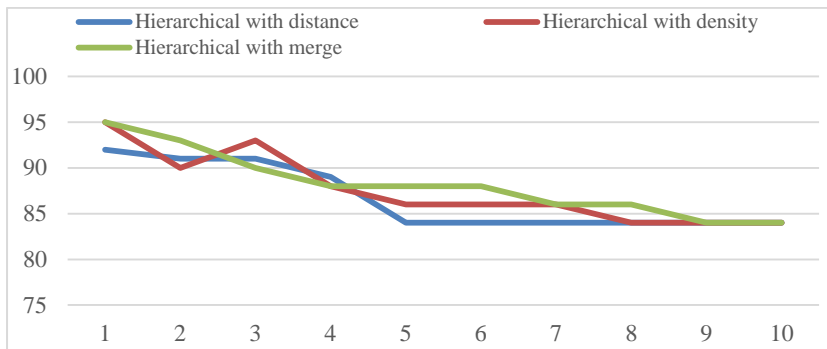
**Precision:** با فرض اینکه  $C = \{C_1, \dots, C_k\}$  مجموعه خوشه‌ها و  $a_i$  برابر تعداد داده‌هایی است که در خوشه صحیح خود ( $C_i$ ) قرار گرفته و  $b_i$  برابر تعداد داده‌هایی است که در خوشه نادرست  $C_i$  قرار گرفته‌اند و مجموعه داده دارای  $n$  شیء و  $K$  تعداد خوشه‌هاست، میزان Precision با رابطه (۵) اندازه‌گیری می‌شود [۳۴]:

$$PR = \frac{\sum_{i=1}^K \left( \frac{a_i}{a_i + b_i} \right)}{K} \quad (5)$$

جدول (۱) مقایسه بین الگوریتم‌های سلسله‌مراتبی با سه معیار مختلف روی مجموعه داده soybean data را با تعداد داده ورودی متفاوت نشان می‌دهد. جدول (۲) مقایسه بین الگوریتم‌های سلسله‌مراتبی با سه معیار مختلف روی مجموعه داده zoo data را با تعداد داده ورودی متفاوت نشان می‌دهد. اعداد جدول نتایج حاصل از خوشه‌بندی الگوریتم‌ها با استفاده از معیار ارزیابی Accuracy را نشان می‌دهد و عدد بزرگ‌تر نشان‌دهنده خوشه‌بندی بهتر است. همان‌طور که در جدول (۱) و (۲) و نمودارهای (۱) و (۲) نشان داده شده است، با افزایش تعداد داده‌ها، کارایی خوشه‌بندی افت پیدا می‌کند؛ این امر نشان‌دهنده تعداد بهینه خوشه‌ها برای این نوع داده‌هاست. مثلاً برای داده soybean تعداد بهینه خوشه‌ها ۱ یا ۲ است. نتایج نشان می‌دهد که روش پیشنهادی در مقایسه با دو روش دیگر رشد چشمگیری دارد.

جدول (۱): نتایج خوشه‌بندی داده‌های soybean با استفاده از معیار Accuracy بر حسب درصد

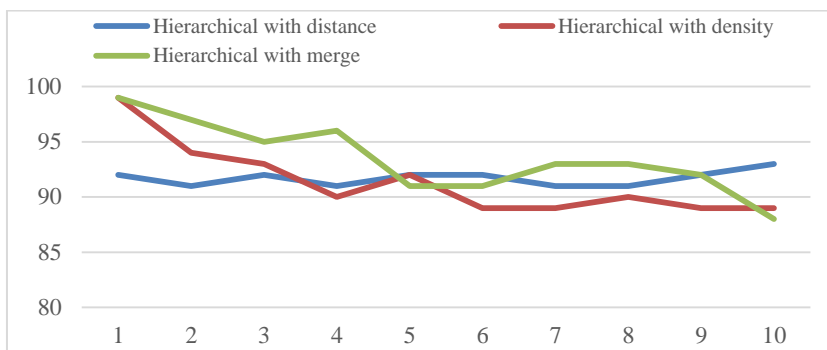
	تعداد خوشه									
	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
Hierarchical with distance	۹۲	۹۱	۹۱	۸۹	۸۴	۸۴	۸۴	۸۴	۸۴	۸۴
Hierarchical with density	۹۵	۹۰	۹۳	۸۸	۸۶	۸۶	۸۶	۸۴	۸۴	۸۴
Hierarchical with merge	۹۵	۹۳	۹۰	۸۸	۸۸	۸۸	۸۶	۸۶	۸۴	۸۴



نمودار (۱): نتایج خوشه‌بندی داده‌های soybean با استفاده از معیار Accuracy برحسب درصد

جدول (۲): نتایج خوشه‌بندی داده‌های zoo با استفاده از معیار Accuracy برحسب درصد

تعداد خوشه										
	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
<b>Hierarchical with distance</b>	۹۲	۹۱	۹۲	۹۱	۹۲	۹۲	۹۱	۹۱	۹۲	۹۳
<b>Hierarchical with density</b>	۹۹	۹۴	۹۳	۹۰	۹۲	۸۹	۸۹	۹۰	۸۹	۸۹
<b>Hierarchical with merge</b>	۹۹	۹۷	۹۵	۹۶	۹۱	۹۱	۹۳	۹۳	۹۲	۸۸



نمودار (۲): نتایج خوشه‌بندی داده‌های zoo با استفاده از معیار Accuracy برحسب درصد

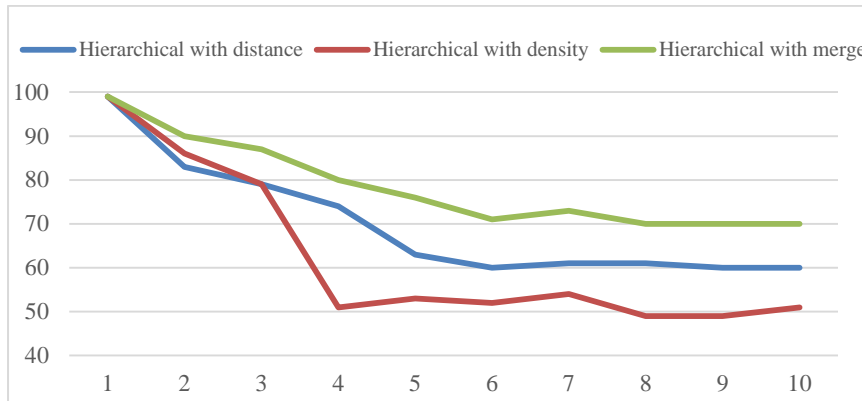
بزرگ‌تر نشان‌دهنده خوشه‌بندی بهتر است. همان‌طور که در جدول (۳) و (۴) و نمودار (۳) و (۴) نشان داده شده است، با افزایش تعداد خوشه‌ها، کارایی خوشه‌بندی افت پیدا می‌کند. همان‌طور که در جدول (۳) و (۴) دیده می‌شود، میزان Precision خوشه‌بندی الگوریتم پیشنهادی نسبت به روش دیگر رشد چشمگیری داشته است.

جدول (۳) و نمودار (۳) مقایسه بین الگوریتم‌های سلسله‌مراتبی با سه معیار مختلف با معیار Precision روی مجموعه داده soybean data با تعداد خوشه‌های متفاوت را نشان می‌دهد. جدول (۴) و نمودار (۴) همین مقایسه را روی مجموعه داده zoo data با تعداد خوشه‌های متفاوت نشان می‌دهد. اعداد جدول نتایج حاصل از خوشه‌بندی الگوریتم‌ها با استفاده از معیار ارزیابی Precision را نشان می‌دهد و عدد

جدول (۳): نتایج خوشه‌بندی داده‌های soybean با استفاده از معیار Precision برحسب درصد

تعداد خوشه										
	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
<b>Hierarchical with distance</b>	۹۹	۸۳	۷۹	۷۴	۶۳	۶۰	۶۱	۶۱	۶۰	۶۰
<b>Hierarchical with density</b>	۹۹	۸۶	۷۹	۵۱	۵۳	۵۲	۵۴	۴۹	۴۹	۵۱
<b>Hierarchical with merge</b>	۹۹	۹۰	۸۷	۸۰	۷۶	۷۱	۷۳	۷۰	۷۰	۷۰

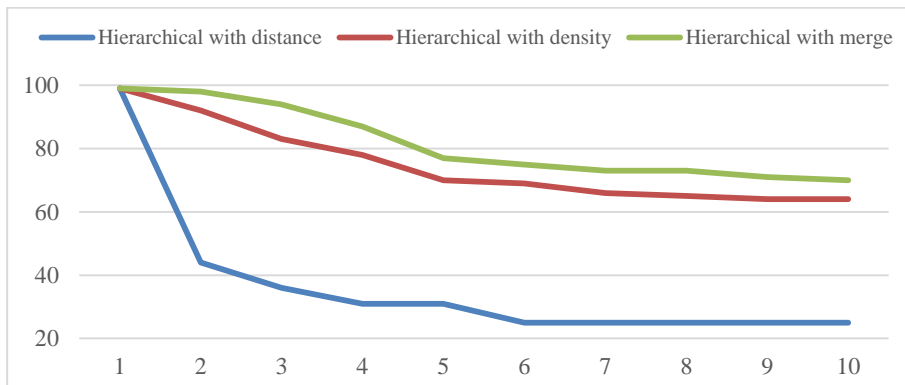




نمودار (۳): نتایج خوشه‌بندی داده‌های soybean با استفاده از معیار Precision برحسب درصد

جدول (۴): نتایج خوشه‌بندی داده‌های zoo با استفاده از معیار Precision برحسب درصد

	تعداد خوشه									
	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
Hierarchical with distance	۹۹	۴۴	۳۶	۳۱	۳۱	۲۵	۲۵	۲۵	۲۵	۲۵
Hierarchical with density	۹۹	۹۲	۸۳	۷۸	۷۰	۶۹	۶۶	۶۵	۶۴	۶۴
Hierarchical with merge	۹۹	۹۸	۹۴	۸۷	۷۷	۷۵	۷۳	۷۳	۷۱	۷۰



نمودار (۴): نتایج خوشه‌بندی داده‌های zoo با استفاده از معیار Precision برحسب درصد

در نتایج به دست آمده در نتایج خوشه‌بندی، ضریب وزنی معیارهای شباهت ثابت در نظر گرفته شده بود. جدول (۵) ارزیابی کیفیت خوشه‌بندی ترکیبی را با ضریب وزنی متفاوت در فرایند خوشه‌بندی نشان می‌دهد. همان‌طور که در

(۵) نشان داده شده است، نتایج حاصل نشان می‌دهد که هرچه ضریب وزنی برای معیار Jaccard بیشتر باشد، نتایج بهتری به دست می‌آید.

جدول (۵): نتایج خوشه‌بندی با ضرایب متفاوت برای ترکیب معیارهای شباهت

ضرایب متفاوت برای معیار Jaccard							
معیار ارزیابی	۰,۹	۰,۸	۰,۷	۰,۶	۰,۵	۰,۴	۰,۳
Accuracy	۹۲٪	۹۰٪	۹۰٪	۸۹٪	۸۸٪	۸۶٪	۸۴٪

جدول (۶) مقایسه‌ای بین الگوریتم پیشنهادی در این مقاله با دو الگوریتم بیان شده در پیشینه پژوهش (الگوریتم Squeezer و الگوریتم fuzzy k-modes) را روی سه مجموعه داده ذکر شده، با دو معیار Accuracy و Precision با تعداد خوشه ثابت ۴ بیان می‌کند. همان‌طور که در جدول (۶) نشان داده شده است، میزان دقت الگوریتم پیشنهادی در مقایسه با

می‌شود. به‌کارگیری همزمان دو معیار شباهت سبب می‌شود که مشکل تجمع خوشه‌ها در یک نقطه از مجموعه داده از بین برود و تجمع خوشه‌ها حالت متوازن‌تری به خود بگیرد.

دو الگوریتم Squeezer و fuzzy k-modes رشد قابل مقایسه‌ای را به همراه دارد. این امر به دلیل ترکیب دو معیار شباهت است؛ زیرا ترکیب این دو معیار موجب از بین بردن نقاط ضعف هریک از آن‌ها و نیز موجب افزایش درصد دقت خوشه‌بندی

جدول (۶): نتایج خوشه‌بندی سه الگوریتم روی سه مجموعه داده

معیار ارزیابی	تعداد خوشه‌ها	Squeezer algorithm	Weighting K-modes	Hierarchical with merge
Accuracy	Zoo data	۹۰٪	۸۲٪	۹۱٪
	Soy data	۸۶٪	۹۰٪	۸۹٪
	Mushroom date	۷۸٪	۷۹٪	۸۳٪
Precision	Zoo data	۸۰٪	۸۵٪	۸۸٪
	Soy data	۸۶٪	۸۲٪	۹۰٪
	Mushroom date	۷۳٪	۷۸٪	۸۳٪

مقاله، یک الگوریتم خوشه‌بندی سلسله‌مراتبی با مفهوم ترکیب معیارهای شباهت برای داده‌های دسته‌ای ارائه شد. در اینجا دو معیار فاصله Jaccard و چگالی مبتنی بر فرکانس تکرار صفات با یکدیگر ترکیب شدند. هریک از این معیارها در الگوریتم‌های جداگانه‌ای برای خوشه‌بندی داده‌های دسته‌ای استفاده شده‌اند. نتایج حاصل از خوشه‌بندی این الگوریتم نشان می‌دهد که خوشه‌های ایجادشده توسط ترکیب معیارها نتایج بهتری را نسبت به استفاده از هر معیار به‌طور جداگانه برای خوشه‌بندی به‌دست می‌دهد؛ این امر نشان‌دهنده تأثیر کمتر داده‌های پرت و دارای نویز در فرایند خوشه‌بندی است.

#### ۴. نتیجه‌گیری

هدف از آنالیز خوشه‌بندی، تقسیم داده‌ها داخل خوشه‌های بامعنی است. در مجموعه داده‌های دسته‌ای، به دلیل ماهیت درونی آن‌ها، خوشه‌بندی آن‌ها با خوشه‌بندی داده‌های عددی متفاوت است. چالش‌هایی در زمینه خوشه‌بندی داده‌های دسته‌ای مطرح است که شامل عدم امکان تفکیک‌پذیری روی دامنه و عدم امکان اندازه‌گیری شباهت برای داده‌های دسته‌ای است. مهم‌ترین ویژگی داده‌های دسته‌ای این است که معیار فاصله طبیعی درباره آن‌ها کاربردی ندارد. تغییر شکل داده‌های دسته‌ای به داده‌های عددی، معنای داده را از بین می‌برد. در این

#### مراجع

- [1] Hosseininezhad, F., Salajegheh, A., *Study and Comparison of Partitioning Clustering Algorithms*, Iranian Journal of Medical Informatics, Vol. 2, No. 1, pp. 32-44, 2012.
- [2] Babcock, B., Babu, S., Datar, M., Motwani, R., Widom, J., *Models and issues in data stream systems*, *Proceedings of Principles of DataBase Systems*, pp. 112-118, 2002.
- [3] Cao, F., Liang, J., Li, D., Zhao, X., *A weighting k-modes algorithm for subspace clustering of categorical data*, *International Journal of Neurocomputing*. Vol. 108, No. 13, pp. 23-30, 2013.
- [4] Berry, M., Linoff, G., *Mastering Data Mining: The Art and Science of Customer Relationship Management*, New York, NY, USA: John Wiley and Sons, pp. 987-998, 1999.
- [5] Bai, L., Liang, J., Dang, C., *An initialization method to simultaneously find initial cluster centers and the number of clusters for clustering categorical data*. *Knowledge-Based Systems*, Vol. 24, No. 6, pp. 785-795, 2011.
- [6] Wiederhold, G., *Advances in Knowledge Discovery in Databases*, California: *Proceedings of National Conference on Artificial Intelligence*, Vol. 32, No. 2, pp. 322-330, 1996.
- [7] Pendharkar, P., *Managing Data Mining Technologies in Organizations: Techniques and*

- Applications. Idea Group Publishing. pp. 123-154, 2003.
- [8] Dubes, R., Jain, K., Algorithms for Clustering Data. Englewood Cliffs: NJ: Prentice- Hall. pp. 432-453, 1988.
- [9] Parvin, H., Aliizadeh, H., Minaei-Bidgoli, B., Analoui, M., CCHR: Combination of Classifiers using Heuristic Retraining. *Proceedings of International Conference on Networked Computing and advanced Information Management*. pp. 231-240, 2008.
- [10] Yun, C., Chuang, K., Chen, M., Adherence clustering: an efficient method for mining market-basket clusters, *International Journal of Information Systems*, Vol. 31, No. 2, pp. 170-186, 2006.
- [11] Mesakar, S., Chaudhari, S., Review Paper On Data Clustering Of Categorical Data, *International Journal of Engineering Research & Technology (IJERT)*, Vol. 1, No. 10, pp. 1-18, 2012.
- [12] Chen, H., Chuang, K., Chen, M., On Data Labeling for Clustering Categorical Data, *IEEE Transaction on Knowledge and Data Engineering*, Vol. 20, No. 11, pp. 1458-1472, 2008.
- [13] Li, Y., Li, D., Wang, S., Zhai, Y., Incremental entropy-based clustering on categorical data streams with concept drift, *International Journal of Knowledge-Based Systems*, Vol. 59, No. 12, pp. 33-47, 2014.
- [14] Han, J., Kamber, M., Data Mining Concepts And Techniques, Harcourt India Private Limited, pp. 398-450, 2010.
- [15] Kim, M., Ramakrishna, S., Projected clustering for categorical datasets, *Pattern Recognition Letters*, Vol. 27, No. 12, pp. 1405-1417, 2006.
- [16] Zhexue, H., (ReviewPaper) Extensions to the k-Means Algorithm for Clustering Large Data Sets with Categorical Values, *Data Mining and Knowledge Discovery*, Vol. 2, No. 3, pp. 283-304, 1998.
- [17] Huang, Z., A Fast Clustering Algorithm to Cluster Very Large Categorical Data Sets in Data Mining, *Proceedings of ACM Sigmod on Research Issues on data Mining and knowledge Discovery*, pp. 56-70, 1997.
- [18] Bai, L., Liang, J., Dang, C., Cao, F., A novel fuzzy clustering algorithm with between-cluster information for categorical data, *International Journal of Fuzzy Sets and Systems*, Vol. 215, No. 1, pp. 55-73, 2013.
- [19] Saha, A. Das, S., Categorical fuzzy k-modes clustering with automated feature weight learning. *Neurocomputing*, Vol. 166, No. 1, pp. 422-435, 2015.
- [20] Sankar Sangam, R., Om, H., The k-modes algorithm with entropy based similarity coefficient, *International Journal of Procedia Computer Science*, Vol. 50, No. 1, pp. 93-98, 2015.
- [21] Ji, J., Pang, W., Zhou, C., Han, X., Wang, Z., A fuzzy k-prototype clustering algorithm for mixed numeric and categorical data, *International Journal of Knowledge-Based Systems*, Vol. 30, No. 12, pp. 129-135, 2012.
- [22] Zengyou, H., Xiaofei, X., Shengchun, D., Squeezer: An Efficient Algorithm for Clustering Categorical Data, *International Journal of Computer Science and Technology*, Vol. 17, No. 5, pp. 611-624, 2002.
- [23] Barbar'a, D., Couto, J., Li, Y., COOLCAT: An entropy-based algorithm for categorical clustering, *Proceedings of International Conference ACM CIKM on Information and Knowledge Management*, pp. 590-599, 2002.
- [24] Yang, Y., Guan, X., You, J., CLOPE: a fast and effective clustering algorithm for transactional data, *Proceedings of international conference ACM Sigkdd on Knowledge discovery and data mining*, pp. 682-687, 2002.
- [25] Johnson, A., Wichern, D., Applied Multivariate Statistical Analysis, New Jersey: Prentice Hall, Englewood Cliffs, pp. 34-67, 1988.
- [26] Kaufman, L., J. Rousseeuw, Finding Groups in Data: An Introduction to Cluster Analysis, New York: Wiley, pp. 324-367, 1990.
- [27] Zhang, X., Xu, Z., Hesitant fuzzy agglomerative hierarchical clustering algorithms, *International Journal of Systems Science*, Vol. 46, No. 3, pp. 562-576, 2015.
- [28] Cesario, E., Manco, G., Ortale, R., Top-down parameter-free clustering of high-dimensional categorical data, *IEEE Transactions on Knowledge and Data Engineering*, Vol. 19, No. 12, pp. 1607-1624, 2007.
- [29] Boriah, S., Chandola, V., Kumar, V., Similarity Measures for Categorical Data: A Comparative Evaluation, *Proceedings of International Conference SIAM on Data Mining*, pp. 212-220, 2008.
- [30] Kou, G., Peng, Y., Wang, G., Evaluation of clustering algorithms for financial risk analysis using MCDM methods, *International Journal of Information Sciences*, Vol. 275, No. 1, pp. 1-12, 2014.
- [31] Guha, S., Rastogi, R., Shim, K., ROCK: a robust clustering algorithm for categorical attributes,

- International Journal of Information Systems*, Vol. 25, No. 5, pp. 345-366, 2000.
- [32] Bai, L., Liang, J., Dang, C., A cluster centers initialization method for clustering categorical data, *International Journal of Expert Systems with Applications*, Vol. 39, No. 8, pp. 8022–8029, 2012.
- [33] "<http://archive.ics.uci.edu/ml/datasets.html>," [Online].
- [34] Bai, L., Liang, J., The k-modes type clustering plus between-cluster information for categorical data, *International Journal of Neurocomputing*, Vol. 133, No. 20, pp. 111–121, 2014.