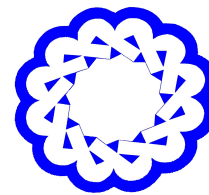


## موجک‌ها و جبرخطی

<http://wala.vru.ac.ir>



دانشگاه ولیعصر (عج)  
رفسنجان

### خوشه‌بندی همزمان دیدگاه‌های مختلف در بررسی داده‌های مالیاتی

ام‌حکیمه اربابی<sup>آ</sup>، مینا جمشیدی<sup>\*آ</sup>، محمدعلی عبدالرزاقی<sup>ب</sup>

<sup>آ</sup>گروه ریاضی کاربردی، دانشگاه تحصیلات تکمیلی صنعتی و فناوری پیشرفته، کرمان، ایران  
<sup>ب</sup>دانشجوی دکتری حسابداری، کارشناس ارشد اداره امور مالیاتی استان کرمان، کرمان، ایران

#### چکیده

استفاده از روشهای داده کاوی در تحلیل داده‌های مالیاتی از جمله تقلب و یا فرار مالیاتی از راهکارهای جدید و مورد توجه می‌باشد. در روشهای ارائه شده، اغلب از الگوریتمهای بانظارت استفاده شده است که ممکن است برای تحلیل داده‌های زیاد هزینه‌بر باشد. همچنین تحلیل داده‌های مالیاتی کمتر از دیدگاه انتخاب ویژگی و یا خوشه‌بندی طیفی مورد بررسی قرار گرفته است. در این مقاله ابتدا مفاهیم خوشه‌بندی طیفی، خوشه‌بندی  $k$ -میانگین و انتخاب ویژگی را مورد مطالعه قرار می‌گیرد و سپس روش جدیدی بر اساس در نظر گرفتن حالت‌های متفاوتی از انتخاب ویژگی معرفی می‌گردد. در این روش ماتریس مشخصه خوشه‌بندی با در نظر گرفتن همزمان همه ماتریسهای مشخصه حاصل از انتخاب تعداد متفاوت ویژگیها، به همراه ضریب وزن تاثیر هر حالت، به دست می‌آید. سپس روش‌های ذکر شده روی داده‌های مالیاتی مورد بررسی قرار می‌گیرند و ویژگیهای مهم مرتبط با فرار مالیاتی با استفاده از انتخاب ویژگی به دست می‌آید. همچنین نتایج حاصل از انتخاب ویژگی و سه روش خوشه‌بندی  $k$ -میانگین، خوشه‌بندی طیفی و روش خوشه‌بندی پیشنهادی جدید ارائه و مقایسه می‌گردد. ملاحظه می‌شود که بر اساس معیارهای ارزیابی خوشه‌بندی، این خوشه‌بندی ها نتایج نسبتاً خوبی ارایه می‌نمایند و در نتیجه راهکاری مناسب برای تحلیل داده‌های مالیاتی می‌باشند و با استفاده از این روشها می‌توان فرار مالیاتی را، در داده‌هایی به صورت ارائه شده، مورد بررسی یا پیش بینی قرار داد.

موجک‌ها و جبرخطی (۱۴۰۴) ©

#### اطلاعات مقاله

تاریخچه مقاله:

دریافت شده: ۲۷ بهمن ۱۴۰۲  
پذیرفته شده: ۳۱ فروردین ۱۴۰۴  
دسترسی آنلاین: ۳۱ فروردین ۱۴۰۴

کلمات کلیدی:

خوشه‌بندی، انتخاب ویژگی،  
داده مالیاتی.

\*نویسنده مسئول

آدرس ایمیلها: omhakimeh75@gmail.com (ام‌حکیمه اربابی)، m.jamshidi@kgut.ac.ir (مینا جمشیدی)

<http://doi.org/10.22072/wala.2025.2023159.1445>

موجک‌ها و جبرخطی (۱۴۰۴) ©

## ۱. مقدمه

یکی از مهم‌ترین منابع درآمدی دولت و کشورهای در حال توسعه مالیات محسوب می‌شود. مقدار قابل توجهی از درآمدهای مالیاتی دولت‌ها به دلیل استفاده کارشناسان مالیاتی از روشهای حسابرسی سنتی از دست می‌رود. از آنجا که نرخ تقلب و فرار مالیاتی رو به رشد است، نیازمند روشهایی می‌باشیم تا بتوانیم این موارد را تشخیص دهیم. با پیشرفت‌های به وجود آمده در جمع‌آوری داده‌ها و قابلیت‌های ذخیره‌سازی در طی دهه‌های اخیر مجموعه‌های داده‌ای با ابعاد بالا در علوم مختلف به سرعت در حال افزایش هستند. مسائل مربوط به امور مالیاتی نیز از این امر مستثنی نمی‌باشند و سیستم‌های سنتی پاسخگوی مدیریت این حجم اطلاعات نیستند. بنابراین در جهت مدیریت این حجم داده‌های بزرگ نیازمند به سیستم‌های هوشمند است که در این سیستم‌های هوشمند می‌توان از انواع الگوریتم‌های هوش مصنوعی و فرآیند داده‌کاوی در جهت کشف الگوی مورد نیاز و پنهان استفاده کرد. در واقع حجم بالای داده‌های ذخیره شده در سیستم‌های مالیاتی نیازمند به ابزاری است تا داده‌های ذخیره شده را پالایش و پردازش کرده و اطلاعات و دانش مورد نظر را استخراج کند و روابط موجود بین داده را بررسی و استخراج کند. بنابراین محققین به پیش بینی و بررسی عوامل مؤثر بر شناسایی و تشخیص فرارها یا تقلبهای مالیاتی با استفاده از روشهای داده‌کاوی می‌پردازند. در این راستا از روش‌های مختلفی از قبیل قوانین همبستگی، خوشه‌بندی، شبکه‌های عصبی، درخت‌های تصمیم، شبکه‌های بیزین، رگرسیون و روشهای ژنتیک استفاده شده است.

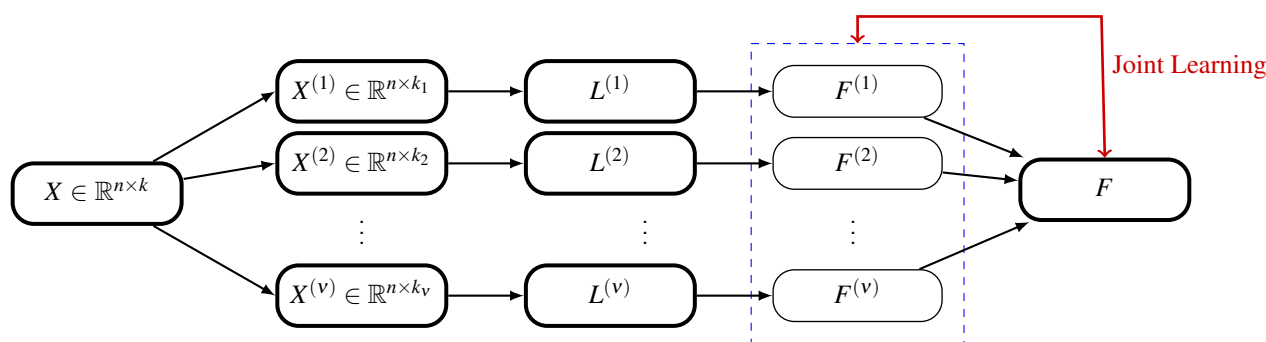
انتخاب ویژگی و خوشه‌بندی از جمله روشهایی می‌باشند که برای بررسی داده‌ها در داده‌کاوی به کار می‌رود. با اینکه این روشها می‌توانند برای داده‌های مالیاتی نیز مورد استفاده قرار گیرند، اما از روشهای انتخاب ویژگی کمتر در این امر استفاده شده‌است. در این مقاله سعی شده است که تحلیل داده‌های مالیاتی با ابزاری خوشه‌بندی و انتخاب ویژگی به صورت همزمان انجام شود. برای این کار برخی از اطلاعات مؤدیان مورد بررسی قرار می‌گیرد و هر مؤدی به عنوان یک داده و اطلاعات مؤدی به عنوان ویژگیهای داده در نظر گرفته می‌شوند. انتخاب ویژگی، یکی از روشهای داده‌کاوی است که موجب صرفه‌جویی در ذخیره‌سازی و محاسبه زمان می‌شود. انتخاب ویژگی در مسائلی که تعداد نمونه‌ها نسبت به تعداد ویژگی‌ها بسیار کم می‌باشد مفید است. اما انتخاب ویژگی در برخی مسائل در تشخیص ویژگیهایی که تأثیر بیشتری در فرآیند داده‌کاوی دارند، نیز مؤثر می‌باشد و نشان می‌دهد که جمع‌آوری اطلاعات این ویژگیها یا تمرکز بر آنها ضروری می‌باشد. بنابراین در تشخیص برخی عوامل مؤثر در تشخیص تقلبهای مالیاتی می‌تواند مفید باشد. خوشه‌بندی نیز در راستای تشخیص رفتار مؤدیان می‌تواند به کار گرفته شود.

الگوریتم‌های متفاوتی هم برای انتخاب ویژگی و هم برای خوشه‌بندی ارائه شده است. از متداولترین الگوریتمهای انتخاب ویژگی استفاده از امتیاز لاپلاسین می‌باشد که بر پایه گراف متناظر داده‌ها محاسبه می‌شود. همچنین روشهای دیگری در سالهای اخیر ارائه شده است که اهداف مختلفی از جمله نداشتن پارامتر، کارایی بر روی داده‌های بزرگ و ... را دنبال می‌نمایند [۱۶، ۱۸، ۲۸، ۲۹]. همچنین مشهورترین و پرکاربردترین روشهایی که برای خوشه‌بندی ارائه شده‌است روش  $k$ -میانگین و خوشه‌بندی طیفی می‌باشند. در بحث خوشه‌بندی در سالهای اخیر موضوع خوشه‌بندی چنددیدگاه بسیار مورد توجه قرار گرفته است. در این روشها داده‌هایی وجود دارند که اطلاعات موجود برای آنها از دیدگاههای مختلف وجود دارد؛ به عنوان مثال یک خبر می‌تواند به زبانهای مختلف انتشار یابد. بنابراین برای هر مجموعه داده چندین ماتریس ویژگی وجود دارد و هدف این است که با در نظر گرفتن همزمان همه این ماتریسهای ویژگی خوشه‌بندی بهتری برای مجموعه داده ارائه گردد. در برخی روشها، که داده به صورت چند دیدگاه نمی‌باشد، جهت دسترسی به اطلاعات بهتری از داده سعی می‌شود که با روشهای ابتکاری ماتریس ویژگیهای جدیدی برای داده به دست آورند، به عنوان مثال از توان‌های ماتریس مشابهت استفاده می‌شود [۲۲]. در این راستا تحقیقات زیادی بر اساس معیارهای مختلف انجام شده است [۶، ۸، ۱۲، ۲۲]، ولی هیچ‌گاه مساله از این منظر مورد بررسی قرار نگرفته است که اگر انتخاب ویژگیهایی به تعداد متفاوت انجام پذیرد، آیا با در نظر گرفتن همه این انتخاب ویژگیها به صورت همزمان به عنوان اطلاعات جدید، در خوشه‌بندی نهایی بهبودی حاصل می‌شود یا خیر؟ به بیان دیگر آیا صرف انتخاب ویژگی به نتیجه بهتری دست می‌یابیم و یا اینکه می‌توان از انتخاب ویژگی به شیوه‌های مختلف یا تعداد مختلف به عنوان ابزاری جهت دستیابی به اطلاعات پنهان داده‌ها استفاده کرد و احاطه بهتری جهت بررسی و تحلیل داده‌ها داشت؟ به دست آوردن ماتریسهای ویژگی جدید با تعداد ویژگیهای متفاوت در واقع به دیدگاهی جدید از داده منجر

می‌شود و می‌توان تحلیل و خوشه‌بندی داده را با در نظر گرفتن آن به صورت داده ای چنددیدگاه انجام داد. در این راستا، در این مقاله روش جدیدی برای خوشه‌بندی بر اساس در نظر گرفتن انتخاب ویژگی های متفاوت به صورت همزمان ارائه می‌شود که کارایی این روش با استفاده از داده‌های مالیاتی و در مقایسه با برخی روشهای دیگر، مورد ارزیابی قرار می‌گیرد و طبق معیارهای ارزیابی می‌توان نتیجه گرفت که استفاده از این روش می‌توان راهکارهایی برای بررسی یا پیش‌بینی مؤدیانی که فرار مالیاتی انجام می‌دهند داشت. به صورت کلی در این مقاله موارد ذیل ارائه می‌گردد:

- راهکاری جدید برای به کارگیری روش انتخاب ویژگی در تحلیل داده های مالیاتی ارائه می‌گردد که در آن از انتخاب تعداد ویژگیهای مختلف برای به دست آوردن دیدگاههای جدید برای داده اصلی استفاده می‌گردد.
- ماتریسهای جدید داده با تعداد ویژگیهای مختلف، به عنوان اطلاعاتی جدید برای آن در نظر گرفته می‌شود و با استفاده از روش خوشه‌بندی طیفی ماتریسهای مشخصه خوشه‌بندی برای هر منظر و نیز ماتریس کلی مشخصه خوشه‌بندی به صورت همزمان مورد یادگیری قرار می‌گیرند.
- با استفاده از ماتریس مشخصه کلی، خوشه‌بندی بر روی داده های مالیاتی انجام می‌شود و با برخی روشهای متداول و اخیر مورد مقایسه قرار می‌گیرد.

می‌توان نمودار ۱ را به عنوان شکل روش پیشنهادی ارائه نمود.



در ادامه در بخش دوم مروری بر تحقیقات گذشته در مورد حل مسائل مالیاتی از طریق داده‌کاوی داریم. سپس در بخش سوم به توضیحاتی در مورد انتخاب ویژگی، خوشه‌بندی  $k$ -میانگین، خوشه‌بندی طیفی می‌پردازیم. در بخش چهارم روش جدیدی بر پایه در نظر گرفتن حالت‌های مختلف از انتخاب ویژگی داده‌ها ارائه می‌نماییم. بخش پنجم به بررسی روی داده‌های مالیاتی اختصاص دارد و نتایج را ارزیابی و ارائه می‌نماییم. در بخش پنجم نتیجه‌گیری این تحقیق را ارائه می‌نماییم.

## ۲. تاریخچه حل مسائل مالیاتی از طریق داده‌کاوی

در این بخش مروری کوتاه از کاربردهای علم داده‌کاوی در مسائل مالی و امور مالیاتی ارائه می‌گردد. در سالهای اخیر در زمینه کشف فرارها و تقلب‌های مالی و مالیاتی، روشهایی با استفاده از الگوریتمهای مختلف داده‌کاوی ارائه شده است. شبکه‌های عصبی چند لایه، ماشین بردار پشتیبان، الگوریتمهای ژنتیک و رگرسیون لجستیک از آن دسته اند. کرکوس، اسپاتیس و مانولوپس در سال ۲۰۰۷ از روش درخت تصمیم‌گیری، شبکه‌های عصبی مصنوعی و شبکه‌های بیزی برای پژوهش‌های روشهای داده‌کاوی برای تشخیص صورتهای مالی جعلی به این نتیجه رسیدند که روشهای بیزی نسبت به بقیه کارایی بهتری دارند [۲۰]. ژو و کاپور در سال ۲۰۱۱ از تکنیک‌های شامل مجموعه مدل‌های رگرسیونی، درخت تصمیم‌گیری و شبکه‌های بیزی برای تشخیص تقلب در صورتهای مالی استفاده کردند و به این نتیجه رسیدند که ترکیب این مدل‌ها نتایج بهتری خواهد

داشت [۳۱]. ووا و همکارانش در سال ۲۰۱۲ منظور افزایش عملکرد تشخیص فرار پرداخت از مالیات روشهای مبتنی بر داده کاوی استفاده کردند که استفاده از سیستم های مبتنی بر داده کاوی موجب بهبود چشمگیر در تشخیص شیوه‌های فرار در حوزه مالیات شد [۲۷].

باقرپور و همکاران از جمله محققان داخلی بودند که در سال ۲۰۱۲ با استفاده از تکنیک های داده کاوی مانند درخت تصمیم و شبکه‌های عصبی مصنوعی به بررسی عوامل مالی و غیر مالی برگریز مالیاتی پرداختند. [۱]. در همین سال برای پیش بینی گزارش حسابرس مستقل در ایران از روش های درخت تصمیم، شبکه‌های عصبی مصنوعی و رگرسیون لجستیک استفاده شد و این نتیجه حاصل شد که میانگین دقت مدل حاصل از روش درخت تصمیم قابل قبول تر می‌باشد [۲]. سهرابی و همکاران در سال ۲۰۱۵ برای ارزیابی مالیات عملکرد شرکت ها و تحلیل روندهای مالیاتی با استفاده از الگوریتم های داده کاوی، که شامل روش های خوشه‌بندی و طبقه بندی بود، استفاده کردند که برترین این روشها خوشه‌بندی مبتنی بر چگالی است [۴]. وانهلودا و همکارانش نیز در سال ۲۰۲۰ تحقیقی درباره شیوه کشف تقلب مالیات بر ارزش افزوده با تکنیکهای خاص مانند قابل تشخیص ناهنجاری بدون نظارت انجام داده‌اند و یک روش ارزیابی جدید معرفی نموده‌اند که با استفاده از آن نشانه‌های رفتار قابل اطمینان را می‌توان بررسی نمود [۲۴]. یان و همکارانش در سال ۲۰۲۰ در یک کار پژوهشی از الگوریتم ژنتیک تطبیقی بهبودیافته برای ارائه مدلی جهت شناسایی تقلب در بیمه خودرو بر اساس شبکه عصبی استفاده کردند. نتایج تجربی نشان داد که الگوریتم ژنتیک بهبودیافته دارای دقت بهتر و مؤثر تری نسبت به الگوریتمهای ژنتیک متداول و سنتی می‌باشد [۳۰]. در سال ۲۰۲۰، دونگ و همکارانش تشخیص تقلب از طریق جنگل تصمیم گیری که ویژگی خود رمزگذار عصبی را دارد، مورد بررسی قرار دادند. آنها یک روش و الگوریتم قابل آموزش را پیشنهاد دادند و کارایی این روش را با بررسی روی داده‌ها نشان دادند [۱۱]. شا و کومار در سال ۲۰۲۱ از الگوریتمهای متعدد یادگیری ماشین مانند ماشین بردار پشتیبان،  $k$  - نزدیکترین همسایه و شبکه عصبی مصنوعی برای پیش بینی وقوع تقلب استفاده کرده اند [۵]. در جدول ۱ مزایا و معایب برخی روشهای ذکر شده و نیز روش پیشنهادی جدید ارائه می‌گردد.

### ۳. مروری بر انتخاب ویژگی، خوشه‌بندی $k$ - میانگین، خوشه‌بندی طیفی

در این بخش ابتدا یک روش متداول انتخاب ویژگی بیان می‌شود و سپس روشهای خوشه‌بندی  $k$  - میانگین و خوشه‌بندی طیفی به صورت خلاصه ارائه می‌گردد. اما قبل از آن لازم است گراف متناظر با داده‌ها توضیح داده شود. فرض کنیم که مجموعه داده  $X = \{x_1, \dots, x_n\}$  را داریم. در گراف متناظر با داده‌ها رأس  $i$  ام متناظر داده  $x_i$  می‌باشد. این گراف با یکی از سه روش  $\epsilon$  - همسایگی،  $k$  - نزدیکترین همسایگی و گراف کامل به دست می‌آید چنانچه رئوس متناظر داده‌های  $x_i$  و  $x_j$  به هم متصل باشند، اغلب از تابع کرنل گاوسی برای تعیین میزان شباهت آنها یا همان وزن یالها استفاده می‌شود. با داشتن گراف متناظر داده‌ها می‌توان ماتریس مجاورت گراف،  $W$ ، که ماتریسی متقارن است و نیز ماتریس لاپلاسیان را به صورت  $L = D - W$  که  $D$  ماتریس قطری متناظر با گراف است و درایه‌های قطر درجه وزن دار رئوس می‌باشند. مرجع [۲۵] نگاهی اجمالی به تعدادی از خواص ماتریس لاپلاسیان دارند.

#### ۳.۱. انتخاب ویژگی به روش امتیاز لاپلاسیان

در روشهای انتخاب ویژگی، هدف انتخاب ویژگی‌های مرتبط است که غالباً بیشترین خاصیت جداکنندگی داده‌های کلاس‌های مختلف را دارند و در نتیجه بیشترین تأثیر را در نتیجه مطلوب دارند [۱۹]. امتیاز لاپلاسیان یکی از روشهای انتخاب ویژگی می‌باشد که اهمیت ویژگیها با توجه به قدرت آنها برای حفظ همسایگی داده‌ها مورد توجه قرار می‌گیرد. بنابراین ویژگیهای مرتبط به گونه‌ای انتخاب می‌شوند که این خاصیت حفظ شود [۱۴]. فرض کنید  $l_r$  نمایانگر امتیاز لاپلاسیان ویژگی  $r$  ام باشد و نیز  $x_{ri}$  نشان دهنده میزان ویژگی  $r$  در داده  $i$  ام باشد. برای ویژگی  $r$  ام بردار  $f_r = [x_{r1}, \dots, x_{rm}]^T$  معرفی می‌شود.  $\bar{f}_r$  به صورت

$$\bar{f}_r = f_r - \frac{f_r^T D \mathbf{1}}{\mathbf{1}^T D \mathbf{1}}$$

محدودیتها	مزایا	نوع یادگیری	روش
حساس به شرایط اولیه عدم توانایی تشخیص شکل پیچیده	کم هزینه بودن قابلیت اجرا روی داده بزرگ	بدون نظارت	$k$ - میانگین
عدم توانایی تشخیص شکل پیچیده بیش برآزش هنگام استفاده تعداد داده زیاد عدم تشخیص روابط بین ویژگیها	تعداد پارامتر کم قدرت تشخیص ویژگیهای مهم نتایج خوب روی تعداد داده کم	با نظارت	رگرسیون لجستیک
احتمال تولید روابط نادرست	کار با هر نوع ویژگی بدون نیاز به پارامتر	با نظارت	درخت تصمیم‌گیری
حساس به شکل داده ورودی	حساسیت کم به نقاط گم شده عملکرد خوب روی داده کم	با نظارت	شبکه بیزی
نیاز به تعداد زیاد داده تست تعداد پارامتر زیاد	دقت بالا قدرت تشخیص شکل‌های پیچیده	با نظارت	شبکه عصبی
حساسیت به نویز	تعداد پارامتر کم قدرت کار با داده‌های با بعد بالا	با نظارت	ماشین بردار پشتیبان
نیاز به ساخت گراف داده‌ها عدم تشخیص روابط بین ویژگیها	قدرت تشخیص شکل پیچیده تحلیل داده با دیدگاههای متفاوت تعداد پارامتر اولیه کم قدرت تشخیص ویژگیهای مهم	بدون نظارت	روش پیشنهادی

جدول ۱: مقایسه برخی روشهای استفاده‌شده برای تحلیل داده‌های مالیاتی

تعریف می‌شود و امتیاز لاپلاسی  $r$  ام ویژگی به صورت

$$l_r = \frac{\bar{f}_r^T L \bar{f}_r}{\bar{f}_r^T D \bar{f}_r}$$

محاسبه می‌شود.

### ۲.۳. خوشه‌بندی $k$ - میانگین، خوشه‌بندی طیفی

روش خوشه‌بندی  $k$  - میانگین یک روش پایه برای بسیاری از روش‌های خوشه‌بندی دیگر محسوب می‌شود. روال کار در این خوشه‌بندی به این صورت است که ابتدا به تعداد خوشه‌های مورد نیاز نقاطی به صورت تصادفی انتخاب می‌شود. سپس داده‌ها با توجه به میزان نزدیکی (شبهت)، به یکی از خوشه‌ها نسبت داده می‌شوند و بدین ترتیب خوشه‌های جدیدی حاصل می‌شود. با تکرار همین روال می‌توان در هر تکرار با میانگین‌گیری از داده‌ها مراکز جدیدی برای آن‌ها محاسبه کرد و مجدداً داده‌ها را به خوشه‌های جدید نسبت داد. این روند تا زمانی ادامه پیدا می‌کند که شرط توقف برقرار گردد.

یکی از مشکلات اصلی خوشه‌بندی  $k$  - میانگین این است که خوشه‌هایی با شکل‌های پیچیده را نمی‌تواند تشخیص دهد و همین موضوع باعث شده است که روش خوشه‌بندی طیفی برای تعیین خوشه‌هایی با شکل‌های پیچیده تر به کار رود. در حل مسأله مربوط به خوشه‌بندی طیفی در گراف متناظر داده‌ها به حذف یالهای کم ارزش می‌پردازیم به طوری که گراف تقریباً به چند مولفه همبندی که همان خوشه‌ها می‌باشند، برسد. ثابت شده است، [۲۵]، که این راهکار معادل است با حل مسأله مینم

سازی زیر:

$$\begin{aligned} \min_F \quad & Tr(F^T L F), \\ \text{s.t.} \quad & F^T F = I_k. \end{aligned} \quad (۱.۳)$$

که  $k$  تعداد خوشه‌ها می‌باشد. با استفاده از قضیه ریلی ریتز [۲۵] جواب این مسأله به یافتن مقادیر ویژه ماتریس لاپلاسیان می‌رسد و در واقع  $F$  ماتریسی است که ستونهای آن بردارهای ویژه  $k$  کوچکترین مقدار ویژه ماتریس  $L$  می‌باشد. سطرهای  $F$  به عنوان بازنمایشی از داده‌ها در فضای با بعد کمتر در نظر گرفته می‌شوند. حال خوشه‌ها با استفاده از خوشه‌بندی  $k$ - میانگین روی سطرهای  $F$  مشخص می‌شوند. قبل از ورود به بخش اصلی مقاله لازم است که تعریف و نکات زیر را نیز ارائه نماییم.

**تعریف ۱.۳.** مجموعه داده  $\mathcal{X} = \{x_1, \dots, x_n\}$  که دارای  $k$  خوشه  $C_1, \dots, C_k$  است، را در نظر بگیرید. ماتریس  $F \in \{0, 1\}^{n \times k}$  ماتریس مشخصه این خوشه‌بندی نامیده می‌شود چنانچه درایه  $(i, j)$  برابر یک باشد اگر  $x_i \in C_j$  و در غیر اینصورت صفر باشد.

در حل مسأله خوشه‌بندی طیفی، چنانچه هر مولفه گراف دقیقاً مشخص کننده یک خوشه باشد، ماتریس حاصل از حل مسأله مینیم سازی ۱.۳ در شرایط تعریف فوق صدق می‌کند. اما چنانچه هر مولفه بیانگر یک خوشه نباشد ماتریس  $F$  حاصل از مسأله مینیم سازی این ویژگی را نخواهد داشت با این حال از آنجا که در روند خوشه بندی طیفی یالهای کم اهمیت می‌شوند، گراف ایده آل داده‌ها (منظور گرافی که هر مولفه یک خوشه است) با گراف اصلی داده تنها در تعدادی یال کم اهمیت تفاوت دارد. بنابراین ماتریس مشابهت و به تبع آن ماتریس لاپلاسیان گراف ایده آل داده‌ها،  $L_1$ ، و ماتریس اصلی و ماتریس لاپلاسیان متناظر با داده‌ها،  $L$ ، اختلاف کمی دارند. به عبارت دیگر ماتریس  $E$  با درایه های کوچک وجود دارد به طوری که  $L = L_1 + E$ . حال طبق قضیه دیویس-کاهان [۲۵] بردارهای ویژه این ماتریسها به هم نزدیک می‌باشند. بنابراین می‌توان نتیجه گرفت که ماتریس حاصل از مسأله بهینه سازی، فارغ از ترتیب قرار گرفتن ستونها، تا حد زیادی به ماتریس حاصل از یک گراف ایده آل شبیه می‌باشد، بنابراین با اغماض می‌توان آن را به صورت یک ماتریس مشخصه در نظر گرفت و به این ماتریس نیز ماتریس مشخصه گوئیم.

#### ۴. خوشه‌بندی طیفی با در نظر گرفتن حالت‌های مختلف انتخاب ویژگی به صورت همزمان

یکی از مباحثی که در سالهای اخیر مورد توجه قرار گرفته است یادگیری چند دیدگاه می‌باشد. در یادگیری چند دیدگاه برای یک داده چند بردار ویژگی وجود دارد. به عنوان مثال اگر داده یک انسان باشد، تصویر شخص، امضای شخص و ویژگیهای اخلاقی شخص می‌توانند سه بردار ویژگی برای شخص ارائه دهند. بنابراین اگر مجموعه ای از داده‌ها (انسانها) داشته باشیم می‌توان سه نوع ماتریس ویژگی (مربوط به تصویر اشخاص، امضا و ویژگیهای اخلاقی) در نظر گرفت و یا یک خبر ممکن است به زبانهای مختلف در وبسایتها موجود باشد و برای خوشه‌بندی یا دسته بندی اخبار می‌توان همه این زبانها را به صورت همزمان در نظر گرفت. تهیه اطلاعات برای یک مجموعه داده از منظرهای مختلف می‌تواند توسط تحلیل گر نیز انجام شود. به طور کلی می‌توان داده چند دیدگاه را به صورت زیر تعریف نمود:

**تعریف ۱.۴.** فرض کنید مجموعه پایه  $\mathcal{X} = \{x_1, \dots, x_n\}$  را داریم. این مجموعه داده را داده چند دیدگاه گوئیم چنانچه ماتریسهای ویژگی متفاوت برای این داده به صورت  $X_1, \dots, X_r$  موجود باشد که  $X_r \in \mathbb{R}^{n \times d_r}$ .

در مسائل یادگیری چند دیدگاه هدف ارائه راهکاری است که بتوان دیدگاههای مختلف را به صورت همزمان در نظر گرفت و داده پایه که همان اشخاص در این مثال می‌باشند را خوشه بندی یا دسته بندی کرد. یکی از رویکردهای دیگری که در راستای رویکرد یادگیری چند دیدگاه می‌توان در نظر گرفت، این است که اگر داده ای به صورت چند دیدگاه در دسترس

نیست به منظور در نظر گرفتن روابط پنهان بین داده‌ها دیدگاه‌های مختلفی برای داده به وجود آورد و یا ماتریسهای شباهت متفاوتی، با توجه به اینکه چه موضوعی در تحلیل داده مربوطه اهمیت دارد، به وجود آورد. به عنوان مثال چنانچه برای یک مجموعه داده ماتریس شباهت  $A$  وجود دارد، می‌توان دیگر دیدگاهها را ماتریس‌های شباهت  $A^1, A^2, \dots, A^p$  در نظر گرفت که به نوعی روابط بین رئوس برای در نظر گرفتن ویژگیهای پنهان و به دست آوردن دیدگاههای جدید از داده به کار می‌آید [۲۳]. در این بخش به معرفی روشی جدید برای خوشه‌بندی می‌پردازیم که در آن انتخاب ویژگیهایی به تعداد متفاوت انجام می‌شود و سپس با در نظر گرفتن همه این انتخاب ویژگیها خوشه‌بندی را انجام می‌دهیم. در واقع بررسی داده‌ها با در نظر گرفتن منظرها یا دیدگاههای مختلف راهکاری در جهت بهبود نتایج یادگیری داده‌ها می‌باشد. در این بخش قصد داریم که با داشتن مجموعه داده  $X$  و انجام انتخاب ویژگی‌های متفاوت، از دیدگاه‌های مختلف به این مجموعه داده بنگریم و با داشتن همه این دیدگاه‌ها خوشه‌بندی را انجام دهیم. در ادامه ابتدا مدل پیشنهادی را ارائه می‌نماییم. سپس با استفاده از روش تکراری به حل این مسأله بهینه سازی این روش می‌پردازیم.

#### ۱.۴. روش پیشنهادی

همانطور که در مقدمه این بخش بیان گردید، راهکار پیشنهادی بر این مبنا است که برای داده‌های مالیاتی با توجه به اینکه چگونه روابط یا اطلاعاتی برای تحلیل داده‌ها مهم تر است، دیدگاههای جدید به دست آورده شود. از آنجا که تشخیص ویژگیهای مهم داده‌ها در تحلیل این نوع داده‌ها مهم است، در این روش دیدگاههای مختلف با استفاده از در نظر گرفتن ماتریس‌های ویژگی که با فرآیند انتخاب ویژگی به دست می‌آید، تعریف می‌شوند. در اینجا می‌توان فرآیند انتخاب ویژگی را انجام داد؛ بدین گونه که می‌توان انتخاب ویژگی  $r$  ویژگی را به عنوان ماتریس ویژگی جدید در نظر گرفت،  $X_r$ . یعنی هر  $X_r$  به عنوان دیدگاهی جدید برای داده پایه در نظر گرفته می‌شود. سپس همه این ماتریسهای انتخاب ویژگیهای مختلف به صورت همزمان مورد بررسی قرار می‌گیرد. در ادامه به بیان مسأله اصلی می‌پردازیم.

فرض می‌کنیم که برای مجموعه داده  $X = \{x_1, \dots, x_n\}$  را در نظر بگیرید. راهکار برای به دست آوردن دیدگاههای جدید این است که برای این مجموعه انتخاب ویژگیها به تعداد متفاوت انجام دهیم. برای این کار روند انتخاب ویژگی با یک یا با چند روش انتخاب ویژگی به تعداد  $n_1, \dots, n_r$  ویژگی به صورت مستقل انجام می‌شود. منظور با چند روش انتخاب ویژگی اینست که به عنوان مثال می‌توان  $n_1$  ویژگی را با روش شماره یک و  $n_2$  ویژگی را با روش شماره دو انتخاب کرد. حال ماتریسهای حاصل از این انتخاب ویژگیها، یعنی  $X_1, \dots, X_r$  را داریم. بنابراین برای مجموعه داده ذکر شده یک مجموعه از ماتریسها، که می‌توان آنها را دیدگاههای مختلف مجموعه داده اصلی در نظر بگیریم، داریم. حال برای دیدگاه  $t$ ، با استفاده از حل مسأله

$$\begin{aligned} \min_{F_t} & \quad Tr(F_t^T L_t F_t), \\ \text{s.t.} & \quad F_t^T F_t = I_k. \end{aligned} \quad (1.4)$$

می‌توان یک ماتریس مشخصه برای خوشه‌بندی طیفی آن دیدگاه پیدا کرد. در اینجا به دنبال ماتریس مشخصه مشترکی هستیم که تا حد امکان برای خوشه‌بندی همه این دیدگاه‌های مختلف مناسب باشد. در نگاه اول شاید به نظر برسد که بهینه کردن فاصله همه ماتریسهای مشخصه دیدگاههای مختلف از ماتریس مشخصه مشترک راهکار مناسبی باشد، به عبارت دیگر از عبارت  $\min_F \|F - F_t\|$  برای به دست آوردن یک ماتریس مشخصه  $F$  که تا حد امکان به همه این ماتریسهای مشخصه نزدیک باشد، استفاده نماییم. اما نکته حائز اهمیت این است که در به دست آوردن ماتریسهای مشخصه دیدگاههای مختلف ترتیب ستونها که همان بردارهای مشخصه خوشه هستند لزوما رعایت نمی‌شود. بنابراین برای حل این مشکل از این خاصیت استفاده می‌نماییم که در ماتریس  $F_t F_t^T$  در حالت ایده‌آل درایه  $(i, j)$  برابر ۱ است اگر داده‌های  $i$  داده  $j$  ام داخل یک خوشه قرار بگیرند و در غیر اینصورت صفر است. حال برای به دست آوردن ماتریس مشخصه حاصل از در نظر گرفتن انتخاب ویژگیهای با تعداد متفاوت به صورت همزمان، ماتریسی در نظر می‌گیریم که تا حد امکان این خاصیت را حفظ نماید، یعنی

عبارت

$$\|FF^T - F_t F_t^T\|_F^2 \quad (۲.۴)$$

را برای همه دیدگاه‌ها باید به حداقل برسانیم [۲۱]. با توجه به روابط ۱.۴ و ۲.۴ مسأله بهینه سازی زیر را در نظر می‌گیریم:

$$\begin{aligned} \min_{F_t, F} \quad & \sum_{t=1}^v \beta Tr(F_t^T L_t F_t) + \alpha_t \|FF^T - F_t F_t^T\|_F^2, \\ \text{s.t.} \quad & F^T F = I_k, F_t^T F_t = I_k. \end{aligned} \quad (۳.۴)$$

که در اینجا  $\alpha_t$  وزنی است که برای میزان تأثیر هر جمله در نظر گرفته می‌شود که مشابه آنچه که در [۲۱] انجام شده است؛ می‌توان در نظر گرفت  $\alpha_t = \frac{1}{\sqrt{\|FF^T - F_t F_t^T\|_F}}$ . در واقع هر چه  $F_t$  به  $F$  نزدیکتر باشد، مقدار  $\alpha_t$  بیشتر است و تأثیر بیشتری دارد.

۲.۴. حل مسأله

برای به دست آوردن هر یک از مجهولات مسأله ۳.۴ دیگر متغیرها را ثابت در نظر می‌گیریم و مسأله را نسبت به آن مجهول خاص حل می‌نماییم و سپس با استفاده از روش بازگشتی الگوریتمی جهت به دست آوردن همه مجهولات ارائه می‌نماییم. بنابراین برای به دست آوردن  $F$  مقادیر  $F_t$  را ثابت می‌گیریم و از رابطه زیر استفاده می‌نماییم:

$$\sum_{t=1}^v \alpha_t \|FF^T - F_t F_t^T\|_F^2 = \sum_{t=1}^v \alpha_t Tr[(FF^T - F_t F_t^T)^T (FF^T - F_t F_t^T)].$$

بنابراین باید مسأله زیر را حل نماییم:

$$\begin{aligned} \min_F \quad & Tr(F^T M F), \\ \text{s.t.} \quad & F^T F = I_k. \end{aligned} \quad (۴.۴)$$

که در آن

$$M = \sum_{t=1}^v \alpha_t (-2F_t F_t^T + I)$$

که این مسأله مشابه آنچه که در روش خوشه‌بندی طیفی بیان شد با استفاده از قضیه ریلی ریتز قابل حل می‌باشد و  $F$  ماتریسی است که ستونهای آن بردارهای ویژه متناظر با  $k$  کوچکترین مقادیر ویژه ماتریس  $M$  می‌باشد. برای به دست آوردن هر  $F_t$  نیز دیگر مقادیر را ثابت می‌گیریم؛ بنابراین باید مسأله زیر را حل نماییم:

$$\begin{aligned} \min_{F_t} \quad & \beta Tr(F_t^T L_t F_t) + \alpha_t \|FF^T - F_t F_t^T\|_F^2, \\ \text{s.t.} \quad & F_t^T F_t = I_k. \end{aligned} \quad (۵.۴)$$

که معادل است با حل مسأله

$$\begin{aligned} \min_{F_t} \quad & Tr(F_t^T N F_t), \\ \text{s.t.} \quad & F_t^T F_t = I_k. \end{aligned} \quad (۶.۴)$$

که در آن

$$N = \beta L_t - \gamma \alpha_t F F^T + \alpha_t I.$$

از آنجا که  $F$  و  $F_t$ ،  $\alpha_t$  به یکدیگر وابسته اند، این مقادیر را نمی‌توان به صورت مستقل به دست آورد و با یک مقدار دهی اولیه مقادیر ذکر شده به صورت تکراری به روزرسانی می‌شوند. الگوریتم این روش را می‌توان به صورت الگوریتم ۱ ارائه نمود.

### الگوریتم ۱ روش پیشنهادی

ورودی: مجموعه داده  $X$ ، مقادیر اولیه  $\alpha_t$ ، پارامتر  $\beta$ ، تعداد خوشه‌ها  $k$  و تعداد تکرار  $r$   
خروجی:  $F$

- ۱: دیدگاه‌های مختلف  $\{X_1, \dots, X_V\}$  را با استفاده از انتخاب ویژگی به تعداد مختلف به دست آور.
- ۲: مجموعه ماتریس‌های مشخصه  $\{F_1, \dots, F_V\}$ ، را با حل رابطه ۱.۴ محاسبه کن.
- ۳: مراحل زیر را به تعداد تکرار  $r$  انجام بده.
- ۴:  $F$  را با توجه به رابطه ۴.۴ به روز رسانی کن.
- ۵:  $F_t$  را با توجه به رابطه ۶.۴ به روز رسانی کن.
- ۶:  $\alpha_t$  را با استفاده از  $\alpha_t = \frac{1}{\gamma \|F F^T - F_t F_t^T\|_F}$  به روز رسانی کن.

### ۵. ارزیابی روش پیشنهادی برای مجموعه داده‌های مالیاتی

در این بخش به ارزیابی روش ارائه شده برای دو مجموعه داده مالیاتی می‌پردازیم. ابتدا این داده‌ها را توضیح می‌دهیم و سپس روشهای ارزیابی را به صورت مختصر معرفی می‌نماییم. بعد از آن نتایج روشها را روی داده‌ها ارائه می‌دهیم.

#### ۱.۵. معرفی مجموعه داده مالیاتی

مجموعه  $X = \{x_1, \dots, x_n\}$  دارای  $n$  داده که در اینجا داده‌ها مؤدیان مالیاتی می‌باشند را در نظر می‌گیریم. اطلاعات هر مؤدی نمایانگر ویژگیهای هر داده است که در ارتباط با انجام یا عدم انجام تکالیف مالیاتی مالیات دهنده می‌باشد. این ویژگیها شامل موارد زیر می‌باشد: اعلامیه متمم، تسلیم اظهارنامه (صورت مجلس)، تسلیم دفاتر، عدم تسلیم دفاتر یا اظهارنامه، صدور برگ تشخیص اولیه، صدور برگ متمم، نوع مؤدی (کوچک، متوسط، بزرگ)، تمکین مؤدی، اعتراض و توافق، اعتراض و عدم توافق، برگ قطعی و واریز، برگ قطعی و عدم واریز.

بنابراین جدولی از مؤدیان داریم که در هر سطر جدول اطلاعات هر مؤدی بر اساس بله یا خیر نمایش داده شده است. قسمتی از این جدول برای ده مؤدی شهر کرمان در شکل زیر آمده است. حال این اطلاعات را به ماتریس داده‌ها تبدیل

مؤدی	تسلیم اظهارنامه	صورت مجلس	عدم صورتمجلس یا اظهارنامه	صدور برگ تشخیص اولیه	صدور برگ متمم	مؤدی کوچک	مؤدی متوسط	مؤدی بزرگ	تمکین مؤدی	اعتراض و توافق ماده ۲۳۸	اعتراض و عدم توافق	برگ قطعی و عدم واریز	برگ قطعی
1	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
2	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
3	بلی	بلی	خیر	بلی	خیر	بلی	خیر	خیر	بلی	بلی	خیر	خیر	بلی
4	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
5	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
6	بلی	بلی	خیر	بلی	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
7	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
8	بلی	بلی	خیر	بلی	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
9	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی
10	بلی	خیر	بلی	خیر	خیر	بلی	خیر	خیر	خیر	خیر	خیر	خیر	بلی

شکل ۱: نمونه داده‌های مالیاتی

می‌نماییم. بدین گونه که در ماتریس چنانچه این اطلاعات در مورد مؤدی برقرار باشد عدد ۱ و در غیر اینصورت عدد صفر قرار دارد. بنابراین ماتریس ویژگیها را به این صورت می‌سازیم. در اینجا از دو جدول استفاده می‌شود یکی اطلاعاتی که در مقاله [۳] برای ۴۰ مؤدی ارائه شده است و دیگری اطلاعات مربوط به ۱۵۰ مؤدی سازمان مالیاتی کرمان در سال ۱۴۰۰.

### ۲.۵. معیارهای ارزیابی

بعد از انتخاب ویژگی و خوشه‌بندی داده‌ها، با معیارهای ارزیابی روند کار را می‌سنجند تا کارایی آن مورد تحلیل قرار گیرد. روشهای متعددی برای ارزیابی روشهای خوشه‌بندی وجود دارد که از بین آنها دو روش بسیار متداول را انتخاب نموده ایم که توضیحات آنها در زیر آمده است. این روشها در مراجع [۷، ۱۵] به تفصیل بیشتری توضیح داده شده اند.

- شاخص دقت<sup>۱</sup>: ACC به بررسی نقطه به نقطه داده‌ها می‌پردازد، میزان قرار گرفتن داده‌ها در خوشه صحیح را محاسبه می‌نماید و به صورت زیر تعریف می‌شود:

$$ACC = \frac{N_{cor}}{N}$$

که  $N$  تعداد داده‌ها و  $N_{cor}$  تعداد داده‌های قرار گرفته در خوشه درست می‌باشد.

### - شاخص خلوص<sup>۲</sup>:

یکی از متداولترین شاخص‌های ارزیابی در خوشه‌بندی، شاخص خلوص است که درصد مطابقت بین برچسب‌های خوشه‌بندی و برچسب‌های واقعی را می‌سنجد. در این حالت برچسب هر خوشه با برچسب واقعی دسته‌ای که بیشترین اشتراک را دارد مطابقت پیدا کرده و تعداد نقاطی از خوشه که در دسته صحیح طبقه‌بندی شده‌اند شمارش می‌شوند. نسبت این تعداد به تعداد کل نقاط، شاخص خلوص را می‌سازد و شکل محاسباتی آن به صورت زیر است:

$$Purity = \frac{\sum_{i=1}^k \max_j (|S_i \cap C_j|)}{N} = \frac{\sum_{i=1}^k \max_j (m_{ij})}{n}$$

که در اینجا  $S_i$  ها دسته های واقعی می‌باشند و  $C_j$  ها خوشه هایی می‌باشند که با استفاده از یک روش خوشه‌بندی به دست آمده اند.

- اطلاعات متقابل نرمال شده<sup>۳</sup>: اطلاعات متقابل نرمال شده یکی از روشهای آماری ارزیابی نتایج خوشه‌بندی داده‌ها می‌باشد که در آن از میزان همبستگی بین متغیرها استفاده می‌شود. در واقع اطلاعات متقابل میزان اطلاعاتی که بین دو متغیر مشترک است را با استفاده از آنتروپی متغیرهای تصادفی به دست آورد. این معیار به صورت زیر فرمول بندی شده است:

$$Nmi = 2 \frac{M(X_i, X_j)}{E(X_i) + E(X_j)}$$

که  $M(X_i, X_j)$  اطلاعات متقابل بین دو متغیر را برآورد می‌نماید و  $E(X_i)$  آنتروپی هر متغیر می‌باشد.

<sup>1</sup> Accuracy

<sup>2</sup> Purity Index

<sup>3</sup> Normalized Mutual Information

- شاخص رند اصلاح شده<sup>۴</sup>: شاخص رند اصلاح شده به بررسی قرار گرفتن زوج داده‌های هم دسته یا جدا، داخل خوشه‌ها می‌پردازد و به صورت زیر تعریف شده است:

$$ARI(S, C) = \frac{\sum_{i=1}^l \sum_{j=1}^k \binom{m_{ij}}{2} - t_3}{\frac{1}{4}(t_1 + t_2) - t_3}$$

$$t_3 = \frac{2t_1 t_2}{n(n-1)} \text{ و } t_2 = \sum_{j=1}^k \binom{S_{.j}}{2}, t_1 = \sum_{i=1}^l \binom{C_i}{2} \text{ که}$$

۳.۵. بررسی و تحلیل نتایج برای داده‌های مالیاتی

در این قسمت داده‌های معرفی شده در زیربخش ۱.۵ را مورد بررسی قرار می‌دهیم و این داده‌ها را با روشهای داده‌کاوی با دقت بیشتری نگاه می‌کنیم. هدف این است که بررسی کنیم، که آیا با روش خوشه‌بندی، داده‌های دارای واریز و عدم واریز می‌توان تفکیک نمود. ابتدا با استفاده از روش انتخاب ویژگی، ویژگیهای برتر داده‌ها را انتخاب می‌نماییم. سپس روشهای خوشه‌بندی را برای معین کردن خوشه‌ها به کار می‌بریم.

برای انتخاب ویژگی از روش روش امتیاز لاپلاسین استفاده می‌نماییم. بعد از آن از روش خوشه‌بندی طیفی استفاده می‌نماییم و این داده‌ها را خوشه‌بندی می‌نماییم. بعد از انجام روش امتیاز لاپلاسین ویژگیها با ترتیبی که در ادامه آمده مرتب می‌شوند: اعلامیه متمم، صدور برگ متمم، تمکین مؤدی، اعتراض و توافق ماده ۲۳۸، تسلیم اظهارنامه، تسلیم دفاتر، عدم تسلیم دفاتر یا اظهارنامه، صدور برگ تشخیص اولیه، اعتراض و عدم توافق، مؤدی کوچک، مؤدی بزرگ، مؤدی متوسط. ملاحظه می‌شود ارتباط خوبی بین ویژگیهای به دست آمده در این روش و ویژگیهای به دست آمده در [۳] که عبارت بودند از اصل شدن اعلامیه، تسلیم اظهارنامه، تسلیم دفاتر و یا عدم تسلیم یکی از آنها وجود دارد. لازم به توضیح است برخی از ویژگیهای ارائه شده در ماتریس داده به یکدیگر وابستگی خاصی دارند به عنوان مثال بردار ویژگی مربوط به اعلامیه متمم و صدور متمم کاملاً یکسان می‌باشد.

در ادامه با نگر داشتن تعدادی از ویژگیهای مهم، با استفاده از روشهای  $k$ - میانگین و خوشه‌بندی طیفی داده‌ها را خوشه‌بندی می‌نماییم. همچنین از دو روش اخیر انتخاب ویژگی با ایده‌های بر پایه کاهش مختصات [۲۸] و بر اساس واریانس کواریانس فاصله زیرفضایی [۱۸] نیز برای مقایسه روش استفاده می‌نماییم. در جدولهایی که در ادامه می‌آید به ترتیب سه تا نه ویژگی را مورد استفاده قرار داده‌ایم و نتایج ارزیابی را با روش جدید ارائه شده مقایسه کرده‌ایم. برای روش پیشنهادی، تعداد دیدگاههای انتخاب شده در جدول متفاوت می‌باشد. به عنوان مثال در ستون اول دیدگاههایی انتخاب شده‌اند که در آنها یک، دو و سه ویژگی انتخاب شده‌اند و این دیدگاه‌ها را با استفاده از روش پیشنهادی به صورت همزمان در نظر گرفته ایم. در جدولها بهترین نتایج را با هایلایت مشخص نموده‌ایم. با توجه به نتایج ارائه شده، مشخص می‌شود که همه روش‌های مطرح شده قابلیت تشخیص نسبتاً خوبی برای بررسی این نوع داده‌های مالیاتی دارند. اما از بین روشها، روش خوشه‌بندی چند دیدگاه در مجموع می‌تواند بیشتر مورد توجه قرار گیرد.

<sup>4</sup> Adjust Rand Index

## اربابی، جمشیدی، عبدالرزاقی / موجک‌ها و جبر خطی ۱۲ (۱) (۱۴۰۴) ۱۷-۱

معیار	روش	$n = 9$	$n = 8$	$n = 7$	$n = 6$	$n = 5$	$n = 4$	$n = 3$
دقت	$k$ - میانگین	۸۱/۰۹	۸۱/۱۵	۸۴/۱۱	۸۳/۳۲	۸۲/۴۳	۸۲/۴۶	۸۲/۵۰
	طیفی	۷۲/۲۳	۷۲/۵۳	۷۲/۵	۷۲/۴۱	۷۲/۵۰	۷۲/۵۰	۷۲/۵۰
	کاهش مختصات	۷۹/۶۲	۸۲/۵	۷۹/۲۵	۸۲/۵	۸۱/۷۵	۸۰	۸۲/۵
	فاصله زیرفضایی	۷۸/۲۵	۶۵/۵	۸۰	۸۲/۳۷	۸۰	۸۰	۶۵
خلوص	چنددیدگاه	۸۵/۵	۸۵/۷۵	۸۵/۱	۸۶	۸۵/۳۵	۸۵/۸۵	۸۵/۸۵
	$k$ - میانگین	۹۰/۱۰	۹۰/۱۰	۹۱/۰۵	۹۰/۶۵	۹۰/۰۰	۹۰/۰۰	۹۰/۰۰
	طیفی	۹۰/۰۳	۹۰/۱۲	۹۰	۹۰	۹۰	۹۰/۰۲	۹۰
	کاهش مختصات	۹۰	۹۰	۹۰	۹۰	۹۰/۷۵	۹۰	۹۰
اطلاعات متقابل نرمال شده	فاصله زیرفضایی	۹۰	۹۰	۹۰	۹۰	۹۰	۹۰	۹۰
	چنددیدگاه	۹۰	۹۰	۹۰	۹۰	۹۰	۹۰	۹۰
	$k$ - میانگین	۲۲/۸۶	۲۲/۴۶	۳۲/۲۹	۲۹/۱۹	۲۴/۵۵	۲۴/۵۸	۲۴/۶۳
	طیفی	۴/۳۹	۴/۳۹	۴/۳۹	۴/۳۹	۴/۳۹	۴/۳۹	۴/۳۹
شاخص رند اصلاح شده	کاهش مختصات	۲۱/۰۸	۲۴/۶۲	۲۰/۲۰	۲۴/۶۲	۲۶/۷۳	۲۱/۹۵	۲۴/۶۲
	فاصله زیرفضایی	۱۷/۰۵	۰/۱۱	۲۱/۹۶	۲۴/۴۹	۲۱/۹۶	۲۱/۹۵	۰/۱۱
	چنددیدگاه	۲۶/۳۳	۲۷/۰۴	۲۵/۲۱	۲۷/۷۵	۲۵/۹۱	۲۷/۳۲	۲۷/۳۲
	$k$ - میانگین	۷۰/۳۸	۷۰/۳۸	۷۰/۳۸	۷۰/۳۸	۷۰/۳۸	۷۰/۳۸	۷۰/۳۸
شاخص رند اصلاح شده	طیفی	۵۹/۱۰	۵۹/۱۰	۵۹/۱۰	۵۹/۱۰	۵۹/۱۰	۵۹/۱۰	۵۹/۱۰
	کاهش مختصات	۶۷/۱۸	۷۰/۳۸	۶۷/۱۹	۷۰/۳۸	۶۷/۱۸	۶۷/۱۸	۷۰/۳۸
	فاصله زیرفضایی	۶۷/۱۸	۵۳/۳۳	۵۳/۳۳	۷۰/۳۸	۶۷/۱۸	۵۳/۳۳	۵۳/۳۳
	چنددیدگاه	۷۷/۵۶	۷۷/۵۶	۷۷/۵۶	۷۷/۵۶	۷۷/۵۶	۷۷/۵۶	۷۷/۵۶

جدول ۲: نتایج خوشه‌بندی برای ۴۰ داده

معیار	روش	$n = 9$	$n = 8$	$n = 7$	$n = 6$	$n = 5$	$n = 4$	$n = 3$
دقت	$k$ - میانگین	۵۸/۶۸	۵۸/۹۲	۵۸/۶۶	۵۸/۶۶	۵۸/۸۷	۵۸/۶۶	۵۸/۷۸
	طیفی	۷۹/۴۱	۷۹/۹۴	۷۹/۹۴	۷۹/۲۰	۷۸/۹۸	۷۹/۴۱	۷۹/۴۱
	کاهش مختصات	۵۹/۳۷	۵۸/۶۶	۵۹/۳۳	۶۰	۵۹/۳۳	۶۰/۶۶	۵۸/۶۸
	فاصله زیرفضایی	۵۸/۶۶	۵۸/۶۶	۵۸/۶۶	۵۸/۶۶	۵۹/۳۳	۵۸/۶۶	۵۸/۶۶
خلوص	چنددیدگاه	۸۴/۳۳	۸۶/۱۷	۸۶/۵	۸۷/۱۷	۹۳/۳۳	۸۸/۳۳	۸۷/۱۷
	$k$ - میانگین	۷۲/۱۲	۷۲/۱۴	۷۲/۱۴	۷۲/۱۲	۷۲/۱۴	۷۲/۱۱	۷۲/۰۵
	طیفی	۸۰/۴۹	۸۰/۷۲	۸۰/۷۲	۸۰/۴۰	۸۰/۳۰	۸۰/۴۹	۸۰/۴۹
	کاهش مختصات	۷۲	۷۲	۷۲	۷۲	۷۲	۷۲	۷۲
اطلاعات متقابل نرمال شده	فاصله زیرفضایی	۷۲	۷۲	۷۲	۷۲	۷۲	۷۲	۷۲
	چنددیدگاه	۸۷/۵۷	۸۸/۷۵	۸۸/۹۶	۸۹/۳۹	۹۳/۳۳	۹۰/۱۳	۸۹/۳۹
	$k$ - میانگین	۰/۴۴	۰/۴۳	۰/۴۵	۰/۴۵	۰/۹۵	۲/۴۵	۰/۶۵
	طیفی	۳۶/۷۳	۳۷/۷۴	۳۷/۷۴	۳۶/۳۳	۳۵/۹۲	۳۶/۷۳	۳۶/۷۳
شاخص رند اصلاح شده	کاهش مختصات	۰/۰۸	۰/۴۵	۰/۰۸	۰/۰۴	۰/۰۸	۰/۰۱	۰/۴۵
	فاصله زیرفضایی	۰/۴۵	۰/۴۵	۰/۴۵	۰/۴۵	۰/۴۵	۰/۴۵	۰/۴۵
	چنددیدگاه	۴۴/۹۰	۴۸/۲۸	۴۸/۹۰	۵۰/۱۳	۶۱/۵۰	۵۲/۲۸	۵۰/۱۳
	$k$ - میانگین	۵۱/۶۷	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۷
شاخص رند اصلاح شده	طیفی	۶۹/۴۳	۶۹/۴۳	۶۹/۴۳	۵۱/۶۷	۶۹/۴۳	۶۹/۴۳	۶۹/۴۳
	کاهش مختصات	۵۱/۴۱	۵۱/۱۸	۵۱/۴۲	۵۱/۶۷	۵۱/۴۲	۵۱/۹۵	۵۱/۱۷
	فاصله زیرفضایی	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸	۵۱/۱۸
	چنددیدگاه	۸۷/۴۷	۸۷/۴۷	۵۱/۶۸	۸۷/۴۷	۸۷/۴۷	۸۷/۴۷	۵۱/۶۸

جدول ۳: نتایج خوشه‌بندی برای ۱۵۰ داده

داده	تعداد نمونه‌ها	تعداد ویژگی‌ها	تعداد دسته‌ها
Wine	۱۷۸	۱۳	۳
Glass	۲۱۴	۹	۶
Control	۶۰۰	۶۰	۶
۵Isolet	۱۵۶۰	۶۱۷	۲۶

جدول ۴: معرفی داده‌های مورد آزمایش

## ۶. ارزیابی روش پیشنهادی روی داده‌های غیر مالیاتی

در این بخش به ارزیابی عملکرد روش پیشنهادی روی چهار مدل داده حقیقی متداول می‌پردازیم. مشخصات این داده‌ها در جدول ۴ آورده شده است. معیارهای ارزیابی و روشهای پیشنهادی مطابق با بخش قبل انتخاب شده‌اند.

### ۱.۶. تحلیل نتایج

نتایج آزمایشات روی این داده‌های ذکر شده در جدولهای ۵، ۶، ۷ و ۸ آورده شده است. در سطر اول تعداد ویژگیهای انتخاب شده مشخص شده است و میانگین نتایج برای هر مورد در ستون آخر آورده شده است. همچنین بهترین نتایج را هایلایت نموده ایم.

با بررسی میانگین نتایج می‌توان به این نتیجه رسید که غیر از داده Wine که در آن روش کاهش مختصات بهترین نتایج را دارد در بقیه داده‌ها روش پیشنهادی چند دیدگاه به طور کلی عملکرد بهتری دارد. به طور دقیق‌تر در داده Wine روش پیشنهادی برای معیارهای دقت و خلوص با اختلاف خیلی کم در جایگاه دوم قرار دارد. همچنین برای معیارهای اطلاعات متقابل نرمال شده و شاخص رند اصلاح شده، نتایج به دست آمده کاملاً قابل قبول و قابل رقابت با روشهای جایگاه اول و دوم می‌باشند. همچنین در مورد داده Glass روش چند دیدگاه در شاخص‌های دقت، خلوص و اطلاعات متقابل نرمال شده بهترین نتایج را دارد و در شاخص رند اصلاح شده با اختلاف کم بعد از روش طیفی در جایگاه دوم قرار دارد. در داده Control برای همه معیارها روش پیشنهادی به طور میانگین بهترین نتایج را دارد و در داده ۵Isolet غیر از معیار شاخص رند اصلاح شده باز هم روش پیشنهادی در سایر معیارهای بهترین نتیجه را به طور میانگین به دست آورده است. همچنین در مورد شاخص رند اصلاح شده نتیجه به دست آمده نزدیک به بهترین نتیجه و قابل قبول می‌باشد.

### ۲.۶. بررسی اثر پارامتر روی داده‌های غیر مالیاتی

در روش پیشنهادی یک پارامتر در بیان مساله بهینه‌سازی مطرح شده است که پارامتر  $\beta$  می‌باشد. با توجه به نمودارهای ارائه شده در ذیل تغییر این پارامتر می‌تواند در نتیجه نهایی موثر باشد. در این زیربخش اثر این پارامتر را برای داده‌های غیر مالیاتی در مورد شاخص دقت بررسی می‌نماییم. در شکل ۲ مشخص است که شاخص دقت با تغییر این پارامتر تغییر می‌کند و این تغییرات برای داده‌های Control و Glass واضح‌تر می‌باشد.

## ۷. نتیجه‌گیری و پیشنهادات

در این مقاله روش جدیدی با توجه به در نظر گرفتن دیدگاه‌های مختلف حاصل از انتخاب ویژگیهای متفاوت ارائه گردید. این روش به همراه دو روش خوشه‌بندی طیفی و خوشه‌بندی  $k$ -میانگین پس از انتخاب ویژگی بر روی داده‌های مالیاتی و نیز دو روش انتخاب ویژگی دیگر مورد ارزیابی قرار گرفتند. نتایج حاصل، نشان می‌دهد که روش ارائه شده نسبت به سایر روشها روش کارایی بهتری برای تحلیل داده‌های مالیاتی دارد و نتایج بهتری را ارائه می‌دهد. بررسی روشهای ریاضی و داده‌کاوی بر روی داده‌های مالیاتی می‌تواند کمکی به بهبود تشخیص رفتارهای مالیاتی داشته باشد و می‌توان برای انواع دیگری از داده‌های

میانگین	$n = ۱۱$	$n = ۱۰$	$n = ۹$	$n = ۸$	$n = ۷$	$n = ۶$	روش	معیار
۶۹/۶۰	۶۹/۶۶	۶۹/۶۶	۶۹/۵۵	۶۹/۵۵	۶۹/۶۱	۶۹/۵۸	$k$ -میانگین	دقت
۶۷/۵۲	۶۷/۶۱	۶۷/۴۶	۶۷/۵۰	۶۷/۴۷	۶۷/۵۷	۶۷/۵۴	طیفی	
۶۹/۸۴	۶۹/۶۶	۶۹/۶۶	۷۰/۶۷	۶۹/۶۶	۷۴/۷۵	۶۴/۶۱	کاهش مختصات	
۶۹/۱	۶۹/۱۱	۶۹/۱۰	۶۹/۶۶	۶۶/۸۵	۶۹/۶۶	۷۰/۲۲	فاصله زیرفضایی	
۶۹/۶۶	۶۹/۶۹	۶۹/۶۹	۶۹/۶۴	۶۹/۶۳	۶۹/۶۳	۶۹/۶۹	چنددیدگاه	
۳۷/۴۷	۳۷/۵۶	۳۷/۵۶	۳۷/۳۷	۳۷/۳۷	۳۷/۴۹	۳۷/۴۴	$k$ -میانگین	خلوص
۴۰/۳۸	۴۰/۲۸	۴۰/۴۸	۴۰/۳۶	۴۰/۳۶	۴۰/۴۱	۴۰/۴۱	طیفی	
۶۹/۸۳	۶۹/۶۶	۶۹/۶۶	۷۰/۶۷	۶۹/۶۶	۷۴/۷۵	۶۴/۶۱	کاهش مختصات	
۳۶/۴۹	۳۵/۷۴	۳۵/۷۴	۳۷/۶۲	۳۱/۹۳	۳۷/۶۲	۴۰/۳	فاصله زیرفضایی	
۴۲/۹۵	۴۲/۸۷	۴۳/۰۳	۴۲/۹۵	۴۲/۹۵	۴۲/۸۷	۴۳/۰۳	چنددیدگاه	
۶۹/۶۰	۶۹/۶۶	۶۹/۶۶	۶۹/۵۵	۶۹/۵۵	۶۹/۶۰	۶۹/۵۸	$k$ -میانگین	اطلاعات متقابل نرمال شده
۶۷/۵۲	۶۷/۶۰	۶۷/۴۵	۶۷/۵۰	۶۷/۴۷	۶۷/۵۶	۶۷/۵۴	طیفی	
۷۰/۱۱	۶۹/۶۶	۶۹/۶۶	۷۰/۶۷	۶۹/۶۶	۷۴/۷۵	۶۶/۲۹	کاهش مختصات	
۶۹/۱	۶۹/۱	۶۹/۱	۶۹/۶۶	۶۶/۸۵	۶۹/۶۶	۷۰/۲۲	فاصله زیرفضایی	
۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	چنددیدگاه	
۷۰/۳	۷۰/۳	۷۰/۳۱	۷۰/۳۰	۷۰/۳۰	۷۰/۳۱	۷۰/۳۰	$k$ -میانگین	شاخص رند اصلاح شده
۶۸/۵۳	۶۸/۸۲	۶۸/۸۱	۶۸/۳۷	۶۸/۱۶	۶۸/۱۶	۶۸/۸۲	طیفی	
۷۰/۸۵	۷۱/۴۲	۷۰/۳	۷۲/۰۶	۷۰/۰۹	۷۱/۷	۶۹/۵۴	کاهش مختصات	
۶۹/۹۸	۶۹/۶۵	۶۹/۶	۷۰/۳	۶۷/۹۷	۷۰/۳	۷۲/۱۳	فاصله زیرفضایی	
۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	۶۹/۵۴	چنددیدگاه	

جدول ۵: نتایج بر روی داده Wine

میانگین	$n = ۷$	$n = ۶$	$n = ۵$	$n = ۴$	$n = ۳$	روش	معیار
۴۶/۸۵	۴۶/۵۲	۴۶/۵۹	۴۶/۵۴	۴۷/۲۰	۴۷/۳۷	$k$ -میانگین	دقت
۵۸/۸۵	۵۸/۴۱	۵۸/۷۲	۵۸/۷۴	۵۹/۲۷	۵۹/۰۹	طیفی	
۴۷/۹۳	۴۶/۰۷	۴۶/۹۶	۵۰/۸۴	۴۸/۶۹	۴۷/۰۸	کاهش مختصات	
۴۶/۷۳	۴۹/۰۶	۴۵/۷۹	۴۸/۶۰	۴۷/۶۶	۴۲/۵۲	فاصله زیرفضایی	
۵۹/۵۲	۵۹/۲۱	۶۰/۱۳	۶۰/۰۱	۵۹/۰۳	۵۹/۲۳	چنددیدگاه	
۵۸/۶۴۳	۵۹/۶۱	۵۹/۴۱	۵۸/۹۴	۵۷/۶۳	۵۷/۶۲	$k$ -میانگین	خلوص
۶۴/۰۱	۶۳/۷۴	۶۳/۹۲	۶۳/۸۱	۶۴/۳۲	۶۴/۲۳	طیفی	
۳۸/۵۱	۴۶/۰۷	۴۶/۹۶	۵۰/۸۴	۴۸/۶۹	۴۷/۰۸	کاهش مختصات	
۵۶/۲۵	۵۷/۰۷	۶۱/۰۶	۵۶/۹۳	۵۷/۰۶	۴۹/۱۴	فاصله زیرفضایی	
۶۴/۷۳	۶۴/۴۹	۶۵/۱۸	۶۵/۰۶	۶۴/۶۰	۶۴/۳۴	چنددیدگاه	
۷۷/۹۶	۷۹/۱۲	۷۸/۸۷	۷۸/۲۸	۷۶/۸۰	۷۶/۷۲	$k$ -میانگین	اطلاعات متقابل نرمال شده
۸۱/۱۷	۸۰/۸۱	۸۰/۸۴	۸۱/۱۲	۸۱/۷۷	۸۱/۲۹	طیفی	
۷۷/۹۱	۸۰/۳۵	۷۷/۸۵	۷۷/۵۵	۷۷/۲۲	۷۶/۵۹	کاهش مختصات	
۷۵/۶۱	۷۷/۵۷	۸۱/۳۱	۷۷/۵۷	۷۶/۱۷	۶۵/۴۲	فاصله زیرفضایی	
۸۲/۳۵	۸۲/۲۲	۸۲/۷۸	۸۲/۵۹	۸۲/۴۰	۸۱/۷۸	چنددیدگاه	
۷۷/۳۱	۷۹/۱۱	۷۷	۷۶/۵۱	۷۶/۹۵	۷۷	$k$ -میانگین	شاخص رند اصلاح شده
۸۲/۳۶	۸۱/۸۴	۸۱/۷۷	۸۱/۸۴	۸۳/۱۴	۸۳/۲۱	طیفی	
۷۸/۵۵	۷۹/۲	۷۸/۹۹	۷۹/۳۹	۷۸/۲۲	۷۶/۹۵	کاهش مختصات	
۷۶/۵۲	۷۸/۳۶	۷۹/۲	۷۸/۲۷	۷۷	۶۹/۷۷	فاصله زیرفضایی	
۸۰/۹۱	۷۵/۶۹	۸۳/۱۵	۸۲/۱۳	۸۱/۷۸	۸۱/۷۸	چنددیدگاه	

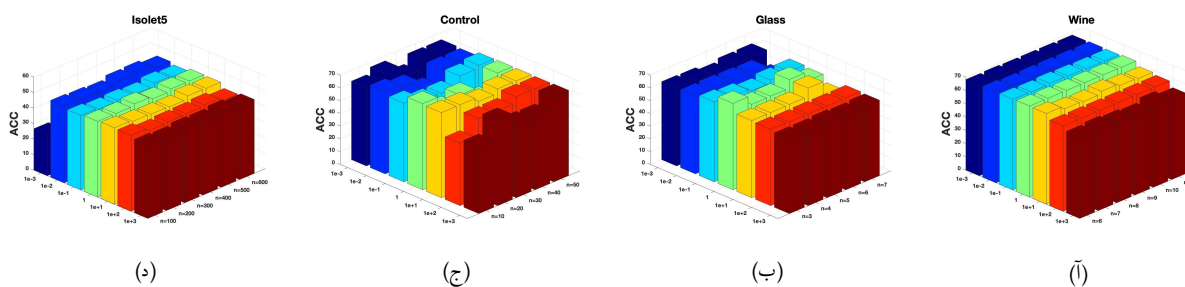
جدول ۶: نتایج بر روی داده Glass

میانگین	$n = ۵۰$	$n = ۴۰$	$n = ۳۰$	$n = ۲۰$	$n = ۱۰$	روش	معیار
۴۷/۴۹	۵۷/۸۹	۵۱/۸۴	۴۶/۵۸	۴۲/۳۷	۳۸/۷۷	$k$ -میانگین	دقت
۶۰/۸۲	۶۰/۴۸	۶۰/۵۲	۶۱/۳۲	۶۰/۸۳	۶۰/۹۲	طیفی	
۵۹/۱۳	۵۸/۹۱	۵۶/۴۲	۵۹/۶۲	۵۸/۴۸	۶۲/۲۶	کاهش مختصات	
۵۷/۷۳	۶۱	۵۹/۶۷	۵۵/۵	۶۰/۱۷	۵۲/۳۳	فاصله زیرفضایی	
۶۴/۴۵	۶۴/۹۲	۵۷/۳۳	۶۶/۶۷	۶۶/۶۷	۶۶/۶۷	چنددیدگاه	
۴۳/۳۷	۵۳/۵۷	۴۸/۰۷	۴۱/۹۱	۳۹/۱۲	۳۴/۱۵	$k$ -میانگین	خلوص
۶۹/۸۱	۶۹/۵۳	۶۹/۶۷	۷۰/۱۸	۶۹/۵۴	۷۰/۱۲	طیفی	
۵۹/۱۳	۵۸/۹۱	۵۶/۴۲	۵۹/۶۲	۵۸/۴۸	۶۲/۲۶	کاهش مختصات	
۶۵/۳۳	۶۷/۰۱	۶۶/۳۹	۶۷/۹۷	۶۶/۷۴	۵۸/۵۶	فاصله زیرفضایی	
۷۳/۳۶	۷۱/۰۳	۶۵/۵۵	۷۶/۷۵	۷۶/۷۵	۷۶/۷۵	چنددیدگاه	
۶۸/۱۲	۶۲/۴۸	۵۷/۸۴	۵۳/۸۵	۴۹/۲۸	۴۲/۹۶	$k$ -میانگین	اطلاعات متقابل نرمال شده
۵۳/۲۸	۶۷/۸۷	۶۷/۹۴	۶۸/۴۱	۶۸/۰۸	۶۸/۳۰	طیفی	
۶۵/۹۳	۶۵/۲۴	۶۴/۴۵	۶۵/۸۲	۶۵/۸۸	۶۸/۲۶	کاهش مختصات	
۶۴/۱۶	۶۵/۳۳	۶۴/۵	۶۵/۳۳	۶۴/۵	۶۱/۱۷	فاصله زیرفضایی	
۷۱/۶۸	۷۰/۶۷	۶۵/۷۵	۷۴	۷۴	۷۴	چنددیدگاه	
۷۹/۰۸	۸۲/۲۱	۸۰/۴۷	۷۹/۱۸	۷۷/۷۷	۷۵/۷۹	$k$ -میانگین	شاخص رند اصلاح شده
۸۴/۲۹	۷۷/۰۱	۸۶/۹۴	۸۶/۹۲	۸۷/۰۶	۸۳/۵۳	طیفی	
۸۵/۵۹	۸۵/۸۵	۸۵/۴۹	۸۴/۱۷	۸۵/۹۱	۸۶/۵۴	کاهش مختصات	
۸۴/۸۳	۸۵/۳۱	۸۵/۱۵	۸۵/۵۹	۸۵/۳	۸۲/۸۱	فاصله زیرفضایی	
۸۷/۷۴	۸۴/۷۶	۸۸/۴۸	۸۸/۴۸	۸۸/۴۸	۸۸/۴۸	چنددیدگاه	

جدول ۷: نتایج بر روی داده Control

میانگین	$n = ۶۰۰$	$n = ۵۰۰$	$n = ۴۰۰$	$n = ۳۰۰$	$n = ۲۰۰$	$n = ۱۰۰$	روش	معیار
۱۹/۵۸	۲۳/۳۶	۲۳/۰۱	۱۹/۷۷	۱۸/۴۳	۱۷/۸۲	۱۵/۰۹	$k$ -میانگین	دقت
۲۷/۷۱	۲۷/۸۰	۲۷/۵۷	۲۷/۶۹	۲۷/۹۱	۲۷/۶۲	۲۷/۶۷	طیفی	
۵۰/۲۱	۵۰/۱۵	۵۰/۳۶	۵۲/۴۰	۵۱/۵۵	۴۸/۴۳	۴۸/۳۵	کاهش مختصات	
۴۹/۱۷	۵۱/۸۳	۵۰/۱۶	۴۹/۵۸	۴۵/۳۵	۴۸/۲۴	۴۹/۸۴	فاصله زیرفضایی	
۵۰/۵۷	۴۹/۶۲	۴۹/۲۱	۴۸/۹۹	۵۰/۷۲	۵۱/۰۳۹	۴۹/۳۵۲	چنددیدگاه	
۳۵/۵۹	۳۹/۶۹	۳۹/۰۶	۳۶/۰۱	۳۵/۴۳	۳۴/۷۵	۲۸/۶۱	$k$ -میانگین	خلوص
۴۴/۲۶	۴۴/۲۸	۴۴/۰۶	۴۴/۲۷	۴۴/۵۵	۴۴/۲۷	۴۴/۱۵	طیفی	
۵۰/۲۱	۵۰/۱۵	۵۰/۳۶	۵۲/۴۰	۵۱/۵۵	۴۸/۴۳	۴۸/۳۵	کاهش مختصات	
۶۷/۷۵	۶۹/۴۹	۶۸/۷۱	۶۷/۵۸	۶۴/۳۹	۶۸/۴۹	۶۷/۸۳	فاصله زیرفضایی	
۶۸/۰۳	۶۸/۱۷	۶۷/۴۵	۶۷/۷	۶۸/۶	۶۸/۶۳	۶۷/۶۲	چنددیدگاه	
۲۲/۹۴	۲۶/۷۲	۲۶/۲۶	۲۳/۵۷	۲۲/۴۰	۲۱/۴۷	۱۷/۲۱	$k$ -میانگین	اطلاعات متقابل نرمال شده
۳۰/۳۸	۳۰/۴۵	۳۰/۲۶	۳۰/۳۵	۳۰/۵۸	۳۰/۳۰	۳۰/۳۶	طیفی	
۵۴/۳۲	۵۴/۷۸	۵۴/۲۴	۵۶/۸۳	۵۵/۲۴	۵۲/۴۹	۵۲/۳۴	کاهش مختصات	
۵۳/۰۸	۵۳/۶۹	۵۲/۷۳	۵۳/۰۵	۵۰/۷۴	۵۳/۴۳	۵۴/۸۴	فاصله زیرفضایی	
۵۴/۴۹	۵۴/۰۳	۵۳/۶۷	۵۳/۰۳	۵۴/۵۶	۵۴/۹	۵۴/۱۸	چنددیدگاه	
۹۰/۹۲	۹۲/۷۴	۹۲/۳۴	۹۰/۰۶	۸۹/۸۱	۹۰/۰۸	۹۰/۴۶	$k$ -میانگین	شاخص رند اصلاح شده
۶۹/۴۱	۷۱/۴۶	۶۰/۱۷	۷۴/۳۲	۷۱/۴۲	۷۲/۳۰	۶۶/۸۳	طیفی	
۹۵/۲۵	۹۴/۹۷	۹۵/۵۲	۹۵/۳۲	۹۵/۵۱	۹۴/۹۵	۹۵/۲۶	کاهش مختصات	
۹۵/۳۲	۹۵/۵۵	۹۵/۵۱	۹۵/۳۶	۹۴/۷۶	۹۵/۳۲	۹۵/۴۴	فاصله زیرفضایی	
۹۴/۳۵	۹۲/۹۵	۹۴/۷۷	۹۴/۳۲	۹۵/۳	۹۴/۱۸	۹۴/۶۱	چنددیدگاه	

جدول ۸: نتایج بر روی داده Isolet



شکل ۲: تاثیر پارامتر  $\beta$  روی شاخص دقت

مالیاتی نیز روشهای جدیدی ارائه نمود. اغلب روشهایی که قبلا ارائه شده است روشهای بانظارت می‌باشد، اما برای حجم زیادی داده مالیاتی این روشها ممکن است کارا نباشد و روش ارائه شده، به عنوان يك روش بدون نظارت می‌تواند ارزش بیشتری داشته باشد.

## مراجع

- [۱] م. باقرپور ولاشانی، م. باقری، ح. خادم و ر. حسینی‌پور، بررسی عوامل مالی و غیر مالی مؤثر برگزید مالیاتی با استفاده از تکنیک‌های داده‌کاوی: صنعت خودرو و ساخت قطعات، مطالعات تجربی حسابداری مالی، ۱ (۳۴) (۱۳۹۱) ۱۰۳-۱۲۸.
- [۲] م. باقرپور ولاشانی، م. ساعدی، ع. مشکانی و م. باقری، پیش‌بینی گزارش حسابرس مستقل در ایران: رویکرد داده‌کاوی، دهمین همایش حسابداری ایران، دانشگاه الزهراء، ۱۳۹۱.
- [۳] م. سامعی‌راد و ا. شاه‌بهرامی، بهبود کارایی الگوریتم‌های تشخیص تقلب مالیاتی با استفاده از الگوهای پردازش موازی، پژوهشنامه مالیات، ۲۴ (۲۹) (۱۳۹۵) ۱۱-۳۲.
- [۴] ب. سهرابی، ا. رئیسی وانانی و و. قانونی شیشوان، ارزیابی مالیات عملکرد شرکت‌ها و تحلیل روندهای مالیاتی با استفاده از الگوریتم‌های داده‌کاوی، تحقیقات مالی، ۱۷ (۴۰) (۱۳۹۴) ۲۱۹-۲۳۸.
- [5] R. B. Asha, K. R. Suresh Kumar, Credit card fraud detection using artificial neural network, *Global Transitions Proceedings*, (2021). doi: <https://doi.org/10.1016/j.gltp.2021.01.006>.
- [6] Y. Cai, C. Hangjun, P. Baicheng, L. Man-Fai, L. Cheng, W. Shiping, Projected cross-view learning for unbalanced incomplete multi-view clustering, *Information Fusion*, **105**, (2024), 102245.
- [7] D. Cai, X. He, J. Han, Document clustering using locality preserving indexing, *IEEE Transactions on Knowledge and Data Engineering*, **17** (12), (2005), 1624–1637.
- [8] C. Cui, R. Yazhou, P. Jingyu, L. Jiawei, P. Xiaorong, W. Tianyi, S. Yutao, H. Lifang, A novel approach for effective multi-view clustering with information-theoretic perspective, *Advances in Neural Information Processing Systems*, **36**, (2024).
- [9] C. H. Chen, Unsupervised margin-based feature selection using linear transformations with neighbor preservation, *Neurocomputing*, **171**, (2016), 1354–1366.
- [10] G. Chandrashekar, F. Sahin, A survey on feature selection methods, *Computers and Electrical Engineering*, **40**, (2014), 16–28.
- [11] M. Dong, L. Yao, X. Wang, B. Benatallah, Ch. Huang, X. Ning, Opinion fraud detection via neural autoencoder decision forest, *Pattern Recognition Letters*, **132**, (2020), 21–29.
- [12] U. Fang, L. Man, L. Jianxin, G. Longxiang, J. Tao, Z. Yanchun, A comprehensive survey on multi-view clustering, *IEEE Transactions on Knowledge and Data Engineering*, **35**(12), (2023), 12350–12368.
- [13] X. Fang, Y. Xu, X. Li, Z. Fan, H. Li, Y. Chen, Locality and similarity preserving embedding for feature selection, *Neurocomputing*, **128**, (2014), 304–315.
- [14] X. He, D. Cai, P. Niyogi, Laplacian score for feature selection, *Advances in Neural Information Processing Systems*, **18**, (2006), 507–514.
- [15] L. Hubert, P. Arabie, Comparing partitions, *Journal of Classification*, **2** (1), (1985), 193–218.

- [16] H. Htun, M. Biehl, N. Petkov, Survey of feature selection and extraction techniques for stock market prediction, *Financial Innovation*, **9**(1), (2023), 26.
- [17] Z. Kang, Z. Xinjia, P. Chong, Z. Hongyuan, Z. Joey, P. Xi, C. Wenyu, X. Zenglin, Partition level multiview subspace clustering, *Neural Networks*, **122**, (2020), 279–288.
- [18] S. Karami, F. Saberi-Movahed, P. Tiwari, P. Marttinen, S. Vahdati, Unsupervised feature selection based on variance–covariance subspace distance, *Neural Networks*, **166**, (2023), 188–203.
- [19] N. N. Karnik, J. M. Mendel, Operations on type-2 fuzzy sets, *Fuzzy Sets and Systems*, **122**, (2001), 327–348.
- [20] Y. Manolopoulos, Data mining techniques for the detection of fraudulent financial statements, *Expert Systems with Applications*, **32**(4), (2007), 995–1003.
- [21] F. Nie, J. Li, X. Li, Self-weighted multiview clustering with multiple graphs, *International Joint Conferences on Artificial Intelligence* (2017), 2564–2570.
- [22] F. Sadjadi, M. Jamshidi, Z. Kang, Multi-view subspace clustering using drop out technique on points, *International Journal of Machine Learning and Cybernetics*, **15**(5), (2024), 1841–1854.
- [23] F. Sadjadi, V. Torra, M. Jamshidi, Preprocessed Spectral Clustering with Higher Connectivity for Robustness in Real-World Applications, *International Journal of Computational Intelligence Systems*, **17**(1), (2024), 86.
- [24] J. Vanhoeyvelde, D. Martensa, B. Peetersb, Value-Added tax fraud detection with scalable anomaly detection techniques, *Applied Soft Computing*, **86**, (2020). doi:<https://doi.org/10.1016/j.asoc.2019.105895>.
- [25] U. Von Luxburg, A tutorial on spectering, *Statistics and Computing*, **17**(4), (2007), 395–416.
- [26] S. Wang, W. Pedrycz, Q. Zhu, W. Zhu, Unsupervised feature selection via maximum projection and minimum redundancy, *Knowledge-Based Systems*, **75**, (2015), 19–29.
- [27] R.S. Wu, C.S. Ou, S.I. Chang, D.C. Yen, Using data mining technique to enhance tax evasion detection performance, *Expert Sestems With Applications*, **39**, (2012), 8769–8777.
- [28] L. Xu, R. Wang, F. Nie, X. Li, Efficient top-k feature selection using coordinate descent method, *Proceedings of the AAAI Conference on Artificial Intelligence*, **37**, (2023). 10594–10601.
- [29] C. Tang, Z. Zheng, Z. Wei, L. Xinwang, Z. Xinzhong, Z. En, Unsupervised feature selection via multiple graph fusion and feature weight learning, *Science China Information Sciences*, **66**(5), (2023), 152101.
- [30] Ch. Yan, M. Li, W. Liu, M. Qi, Improved adaptive genetic algorithm for the vehicle insurance fraud identification model based on a BP neural network, *Theoretical Computer Science*, **817**, (2020), 12–23.
- [31] W. Zhou, G. Kapoor, Detecting evolutionary financial statement fraud, *Decision Support Systems*, **50**(3), (2011), 570–575.