



کاهش بعد داده‌ها ضمن حفظ خوشه‌های داده

محبوبه حقایقی پور^(۱) یحیی فرقانی*^(۲) سیدمحمدحسین معطر^(۳)

(۱) گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

(۲) گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران*

(۳) گروه مهندسی کامپیوتر، واحد مشهد، دانشگاه آزاد اسلامی، مشهد، ایران

چکیده

در این مقاله، با توجه به موفقیت روش‌های خوشه‌بندی مبتنی بر k-means، یک روش کاهش ویژگی بر پایه k-means وزن دار ارائه می‌شود. در روش پیشنهادی، نخست با استفاده از روش k-means وزن دار، به ویژگی‌های داده‌ها وزن داده می‌شود. ویژگی‌های وزین تر ضرورتاً ویژگی‌های مهمتر نیستند و وزن هر ویژگی، تنها بازه هر ویژگی را به نحوی تغییر می‌دهد که خوشه‌بندی بهتری صورت بگیرد. لذا، سپس با استفاده از یک مدل ریاضی جدید، کسری از ویژگی‌های وزندار شده داده‌های هر خوشه انتخاب می‌شود به نحوی که کمترین تغییر در خوشه‌ها حاصل شود. تعداد ویژگی‌های منتخب هر خوشه در روش پیشنهادی، برخلاف روش‌های مشابهی چون k-means تنک و fuzzy c-means تنک بصورت صریح تعیین می‌شود. در ضمن، آزمایش‌های تجربی روی چهار مجموعه داده واقعی نشان می‌دهد که روش پیشنهادی، از دقت بیشتری نسبت به روش‌های LIPCA, LLE و روش K-means تنک برخوردار است.

واژه‌های کلیدی: کاهش ویژگی، خوشه بندی، K-means وزن دار

بطورکلی دو روش کاهش ویژگی وجود دارد: استخراج ویژگی و انتخاب ویژگی. در روش انتخاب ویژگی، کسری از ویژگی‌ها به عنوان ویژگی‌های موثر انتخاب می‌شوند و در روش استخراج ویژگی، تعداد اندکی ویژگی موثر تولید می‌شود که هر یک، تابعی از یک یا چند ویژگی از ویژگی‌های اولیه است.

از جمله روش‌های استخراج ویژگی می‌توان به روش PCA (Principle Component Analysis) [۱] اشاره کرد. در این روش، بهترین ویژگی استخراج شده، آن ویژگی است که بیشترین واریانس را داشته باشد و ویژگی با واریانس نزدیک به صفر، ویژگی بی اهمیت یا نویز محسوب می‌شود. واریانس داده، متناسب است با مجموع مربع اختلاف داده‌ها تا میانگین داده‌ها. در صورت وجود داده نویزی، مربع اختلاف داده نویزی تا میانگین داده‌ها، اثر زیادی بر روی واریانس داده‌ها و عملکرد PCA خواهد گذاشت. برای کاهش اثر داده‌های نویزی، در روش LIPCA [۲] از قدرمطلق اختلاف داده تا میانگین داده‌ها بجای مربع اختلاف داده تا میانگین داده‌ها استفاده می‌شود.

روش استخراج ویژگی Locally Linear Embedding (LLE) [۳] شامل دو مرحله است. در مرحله نخست، هر داده بصورت ترکیب خطی سایر داده‌ها نوشته می‌شود و ضرایب این ترکیب خطی ذخیره می‌گردد. سپس بردارهایی با بعد کمتر از داده‌های اولیه، به نحوی تولید می‌شوند که همچون داده‌های اولیه، هر کدام، یک ترکیب خطی با همان ضرایب ذخیره شده از سایر بردارهای با بعد کمتر باشد.

هدف از کاهش ویژگی، کاهش حجم داده‌ها، حذف ویژگی‌های نامرتب و کشف ویژگی‌های موثر در تحلیل داده‌هاست. ویژگی‌های نامرتب در داده‌ها حتی می‌تواند نتیجه تحلیل داده‌ها را منحرف سازد. به آن دسته از ویژگی‌ها، ویژگی‌های موثر گفته می‌شود که بتواند به خوبی ساختار داده‌ها مثل خوشه‌های موجود در داده‌ها را حفظ کند. لذا، حفظ ساختار داده‌ها یا خوشه‌های داده باید در

زمان کاهش ویژگی مدنظر باشد. در روش‌های انتخاب ویژگی k-means تنک [۴] و fuzzy c-means تنک [۵]، کسری از ویژگی‌ها به نحوی انتخاب می‌شوند که فاصله درون خوشه‌های داده‌ها نسبت به فاصله داده‌ها از کل داده‌ها کمتر شود و امکان خوشه‌بندی بهتر میسر گردد. متأسفانه، امکان تعیین صریح تعداد ویژگی‌های منتخب در این دو روش وجود ندارد. در این تحقیق، با توجه به موفقیت روش‌های خوشه‌بندی مبتنی بر k-means [۴, ۵]، یک روش انتخاب ویژگی بر پایه روش خوشه‌بندی k-means وزن دار [۵] ارائه می‌شود. به بیان دقیقتر، در روش پیشنهادی، نخست با استفاده از روش k-means وزن دار، به ویژگی‌های داده‌ها وزن داده می‌شود. ویژگی‌های وزن‌تر ضرورتاً ویژگی‌های مهم‌تر نیستند و وزن هر ویژگی، تنها بازه هر ویژگی یا شکل خوشه‌ها را به نحوی تغییر می‌دهد که خوشه‌بندی بهتری صورت بگیرد. لذا، سپس با استفاده از یک مدل ریاضی جدید، کسری از ویژگی‌های وزن‌دار شده هر خوشه انتخاب می‌شود به نحوی که کمترین تغییر در خوشه‌ها حاصل شود. تعداد ویژگی‌های منتخب هر خوشه در روش پیشنهادی، برخلاف روش‌های k-means تنک و fuzzy c-means تنک بصورت صریح تعیین می‌شود. درضمن، آزمایش‌های تجربی روی چهار مجموعه داده واقعی نشان می‌دهد که روش پیشنهادی، از دقت بیشتری نسبت به روش‌های LIPCA, LLE و روش K-means تنک برخوردار است.

در ادامه، در بخش ۲ به بیان پیش نیازهای تحقیق پرداخته می‌شود. در فصل ۳، روش پیشنهادی ارائه می‌شود. در فصل ۴ با استفاده از چهار مجموعه داده واقعی، روش پیشنهادی با سه روش کاهش ویژگی دیگر بصورت تجربی مقایسه می‌شود. در فصل ۵ به جمع بندی نتایج و کارهای آتی می‌پردازیم.

۲- پیش نیازها

۲-۱- روش خوشه بندی k-means

فرض کنید می‌خواهیم n داده $\{x_1, x_2, \dots, x_n\}$ را به c خوشه تقسیم کنیم که $x_i \in R^m$ و m بعد هر داده است. مدل خوشه بندی k-means به صورت زیر است:

(۱)

$$\min_{u,z} \sum_{k=1}^c \sum_{i=1}^n u_{ik} \|x_i - z_k\|_2^2$$

subject to $\begin{cases} \sum_{k=1}^c u_{ik} = 1, & i = 1, 2, \dots, n; \\ u_{ik} \in \{0, 1\}, & i = 1, 2, \dots, n; k = 1, 2, \dots, c. \end{cases}$

که z_k مرکز خوشه k -ام است و u_{ik} تعلق داده i -ام به خوشه k -ام را مشخص می‌کند و

$$\|x_i - z_k\|_2^2 = \sum_{j=1}^m (x_{ij} - z_{kj})^2$$

هدف این مدل، پیدا کردن مراکز خوشه‌ها یعنی z_k ها و u_{ik} ها به نحوی است که فاصله داده‌های هر خوشه تا مرکز آن خوشه، حداقل شود. اولین قید این مدل مشخص می‌کند که داده i -ام به یک خوشه تعلق دارد. مساله خوشه‌بندی، یک مساله NP-Complete است. الگوریتم تکراری خوشه بندی k-means برای به دست آوردن مقدار بهینه محلی مدل k-means بصورت الگوریتم ۱ است.

الگوریتم ۱:

۱. مراکز اولیه خوشه‌ها (θ) را به صورت تصادفی در نظر بگیر.
۲. مراکز خوشه‌ها (θ) را ثابت فرض کرده و مقدار بهینه مدل k-means را نسبت به درجات عضویت (u) بدست آور.
۳. درجات عضویت (u) را ثابت فرض کرده و مقدار بهینه مدل k-means را نسبت به مراکز خوشه (θ) بدست آور.
۴. در صورت عدم حصول شرط همگرایی برو به ۲.

۲-۲ روش خوشه بندی k-means وزندار (wk-means)

مدل ریاضی خوشه بندی k-means وزن دار به صورت زیر است:

(۲)

$$\min_{u,z,w} \sum_{i=1}^n \sum_{k=1}^c u_{ik} \sum_{j=1}^m w_{kj}^\beta (x_{ij} - z_{kj})^2$$

Subject to

$$\begin{cases} \sum_{k=1}^c u_{ik} = 1, & i = 1, 2, \dots, n; \\ \sum_{j=1}^m w_{kj} = 1, & k = 1, 2, \dots, c; \\ u_{ik} \in \{0, 1\}, & i = 1, 2, \dots, n; k = 1, 2, \dots, c; \\ w_{kj} \geq 0, & k = 1, 2, \dots, c; j = 1, 2, \dots, m. \end{cases}$$

در این مدل، برای ویژگی j -ام خوشه k -ام، وزن متفاوتی در نظر گرفته شده است که با w_{kj} نشان داده شده است و مقدار بهینه آن در طی فرآیند خوشه‌بندی، تعیین می‌شود. در این روش، به ویژگی‌هایی از یک خوشه، وزن بیشتری داده می‌شود که بازه کوچکتری در آن خوشه داشته باشند. لذا، در این مدل، سعی می‌شود با تغییر بازه ویژگی‌های هر خوشه یا تغییر شکل خوشه‌ها از طریق وزندار کردن ویژگی‌های هر خوشه، خوشه‌بندی بهتری انجام شود. پارامتر β نحوه اختصاص وزن به ویژگی‌ها را کنترل می‌کند. اگر $\beta = 1$ باشد، تنها یک ویژگی از هر خوشه، وزنی غیرصفر خواهد داشت. الگوریتم ۲ روش خوشه‌بندی wk-means را نشان می‌دهد.

الگوریتم ۲

۱. مراکز اولیه خوشه‌ها (z) و وزن‌ها (w) را بصورت تصادفی در نظر بگیر.
۲. مراکز خوشه‌ها (z) و وزن‌ها (w) را ثابت فرض کرده و مقدار بهینه مدل wk-means را نسبت به درجات عضویت (u) با استفاده از رابطه زیر بدست آور:

$$u_{ik} = \begin{cases} 1 & \text{if } \forall l: \sum_{j=1}^m w_{kj}^\beta (x_{ij} - z_{kj})^2 \leq \sum_{j=1}^m w_{lj}^\beta (x_{ij} - z_{lj})^2 \\ 0 & \text{otherwise.} \end{cases}$$

۳. مراکز خوشه‌ها (Z) و درجات عضویت (u) را نسبت به وزن‌ها (w) با استفاده از رابطه زیر بدست آور:

$$w_{kj} = \begin{cases} \frac{1}{m_i} & \text{if } \sum_{i=1}^n u_{ik}(x_{ij} - z_{kj})^2 = 0 \text{ and } m_i = \left\| \left\{ t: \sum_{i=1}^n u_{ik}(x_{it} - z_{kt})^2 = 0 \right\} \right\|; \\ 0 & \text{if } \sum_{i=1}^n u_{ik}(x_{ij} - z_{kj})^2 \neq 0 \text{ and } \sum_{i=1}^n u_{ik}(x_{it} - z_{kt})^2 = 0 \text{ for some } t; \\ \frac{1}{\sum_{t=1}^m \left[\frac{\sum_{i=1}^n u_{ik}(x_{ij} - z_{kj})^2}{\sum_{i=1}^n u_{ik}(x_{it} - z_{kt})^2} \right]^{\frac{1}{\beta-1}}} & \text{if } \sum_{i=1}^n u_{ik}(x_{it} - z_{kt})^2 = 0 \text{ for each } t. \end{cases}$$

۴. وزن‌ها (w) و درجات عضویت (u) را ثابت فرض کرده و مقدار بهینه مدل wk-means را نسبت به مراکز خوشه‌ها (Z) با استفاده از رابطه زیر بدست آور:

$$\min_{u, z, \tilde{w}} F = \sum_{i=1}^n \sum_{k=1}^c u_{ik} \sum_{j=1}^m \tilde{w}_{kj} w_{kj}^\beta (x_{ij} - z_{kj})^2$$

$$\text{subject to } \begin{cases} \sum_{j=1}^m \tilde{w}_{kj} = S, \quad k = 1, 2, \dots, c; \\ \sum_{k=1}^c u_{ik} = 1, \quad i = 1, 2, \dots, n; \\ \tilde{w}_{kj} \in \{0, 1\}, \quad k = 1, 2, \dots, c; j = 1, 2, \dots, m; \\ u_{ik} \in \{0, 1\}, \quad i = 1, 2, \dots, n; k = 1, 2, \dots, c. \end{cases}$$

$$z_{kj} = \frac{\sum_{i=1}^n u_{ik} x_{ij}}{\sum_{i=1}^n u_{ik}}$$

۵. در صورت عدم حصول شرط همگرایی برو به ۲.

۳- روش پیشنهادی (استخراج ویژگی مبتنی بر روش خوشه بندی k-means وزندار)

روش کاهش ویژگی پیشنهادی، شامل دو مرحله است. در مرحله اول، با استفاده از روش wk-means، خوشه‌بندی صورت می‌گیرد. همانطور که قبلاً گفته شد در روش wk-means، سعی می‌شود با وزن‌دار کردن ویژگی‌های داده‌های هر خوشه یا تغییر شکل خوشه‌ها، خوشه‌بندی بهتری صورت گیرد. در مرحله دوم، یعنی پس از تعیین مقدار بهینه وزن ویژگی‌های هر خوشه (w_{kj} ها)، با استفاده از مدل زیر، S ویژگی هر خوشه به نحوی انتخاب می‌شود که کمترین تغییر در خوشه بندی یا کمترین تغییر در تابع هدف مدل wk-means حاصل شود:

که w_{kj} وزن ویژگی j -ام از خوشه k -ام است که قبلاً با استفاده از مدل wk-means به دست آمد. \tilde{w}_{kj} متغیری برای انتخاب ویژگی است. اگر $\tilde{w}_{kj} = 1$ باشد یعنی j -امین ویژگی برای k -امین خوشه حفظ می‌شود وگرنه حذف می‌شود. برای حل مدل (۳)، الگوریتم زیر پیشنهاد می‌شود:

الگوریتم ۳

۱. مراکز اولیه خوشه‌ها (Z) و مقدار متغیرهای انتخاب ویژگی (\tilde{w}) را بصورت تصادفی در نظر بگیر.
۲. مراکز خوشه‌ها (Z) و مقدار متغیرهای انتخاب ویژگی (\tilde{w}) را ثابت فرض کرده و مقدار بهینه مدل (۳) را نسبت

به درجات عضویت (u) با استفاده از رابطه زیر بدست آور:

$$u_{ik} = \begin{cases} 1 & \text{if } \forall l: g_{ik} \leq g_{il} \\ 0 & \text{otherwise.} \end{cases}$$

که $g_{il} = \sum_{j=1}^m \tilde{w}_{lj} w_{lj}^{\beta} (x_{ij} - z_{lj})^2$

۳. مراکز خوشه‌ها (Z) و درجات عضویت (u) را ثابت فرض کرده و مقدار بهینه مدل (۳) را نسبت به متغیرهای انتخاب ویژگی (\tilde{w}) با استفاده از رابطه زیر بدست آور:

$$\tilde{w}_{kj} = \begin{cases} 1 & f_{kj} \leq S \min\{f_{kl}\}; \\ 0 & \text{otherwise.} \end{cases}$$

که $f_{kl} = \sum_{i=1}^n u_i w_{ik}^{\beta} (x_{il} - z_{kl})^2$ و منظور از $S \min$ کمترین مقدار می باشد.

۴. مقدار متغیرهای انتخاب ویژگی (\tilde{w}) و درجات عضویت (u) را ثابت فرض کرده و مقدار بهینه مدل (۳) را نسبت به مراکز خوشه‌ها (Z) با استفاده از رابطه زیر بدست آور:

$$z_{kj} = \frac{\sum_{i=1}^n u_{ik} x_{ij}}{\sum_{i=1}^n u_{ik}}$$

۵. در صورت عدم حصول شرط همگرایی برو به ۲.

شرط همگرایی هر یک از الگوریتم‌های ۱ تا ۳، عدم تغییر مراکز خوشه‌ها در دو تکرار متوالی است. در ادامه، درستی هر یک از مراحل الگوریتم ۳ اثبات می‌شود.

قضیه ۱- اگر مراکز خوشه‌ها (Z) و مقدار متغیرهای انتخاب ویژگی (\tilde{w}) را ثابت فرض کنیم مقدار بهینه مدل (۲۶) نسبت به درجات عضویت (u) با استفاده از رابطه زیر بدست می‌آید:

$$u_{ik} = \begin{cases} 1 & \text{if } \forall l: g_{ik} \leq g_{il} \\ 0 & \text{otherwise.} \end{cases}$$

که $g_{il} = \sum_{j=1}^m \tilde{w}_{lj} w_{lj}^{\beta} (x_{ij} - z_{lj})^2$

اثبات.

اگر مراکز خوشه‌ها (Z) و مقدار متغیرهای انتخاب ویژگی (\tilde{w}) را ثابت فرض کنیم مدل (۳) بصورت زیر در می‌آید:

$$\min_u \sum_{i=1}^n \sum_{k=1}^c u_{ik} \sum_{j=1}^m \tilde{w}_{kj} w_{kj}^{\beta} (x_{ij} - z_{kj})^2$$

subject to $\begin{cases} \sum_{k=1}^c u_{ik} = 1, & i = 1, 2, \dots, n; \\ u_{ik} \in \{0, 1\}, & i = 1, 2, \dots, n; k = 1, 2, \dots, c. \end{cases}$

از آنجا که متغیرهای مدل فوق فقط u_{ik} ها هستند و در هر عبارت از مدل فوق، فقط یک متغیر وجود دارد مدل فوق را نیز می‌توان بصورت جمع چند مدل کوچکتر و بصورت زیر نوشت:

$$\sum_{i=1}^n \left(\min_u \sum_{k=1}^c u_{ik} \sum_{j=1}^m \tilde{w}_{kj} w_{kj}^{\beta} (x_{ij} - z_{kj})^2 \right)$$

subject to $\begin{cases} \sum_{k=1}^c u_{ik} = 1; \\ u_{ik} \in \{0, 1\}, & k = 1, 2, \dots, c. \end{cases}$

i -امین مدل از این مدل‌های کوچکتر را در نظر بگیرید:

$$\min_u \sum_{k=1}^c u_{ik} \sum_{j=1}^m \tilde{w}_{kj} w_{kj}^{\beta} (x_{ij} - z_{kj})^2$$

subject to $\begin{cases} \sum_{k=1}^c u_{ik} = 1; \\ u_{ik} \in \{0, 1\}, & k = 1, 2, \dots, c. \end{cases}$

با توجه به قیود این مدل، فقط یکی از u_{ik} برابر با یک و سایر u_{ik} برابر با صفر می‌شوند. بنابراین برای حداقل شدن تابع هدف مدل، $u_{ik} = 1$ می‌شود اگر ضریب آن در تابع هدف، یعنی $g_{ik} = \sum_{j=1}^m \tilde{w}_{kj} w_{kj}^{\beta} (x_{ij} - z_{kj})^2$ حداقل باشد.

پایان اثبات.

قضیه ۲- اگر مراکز خوشه‌ها (Z) و درجات عضویت (u) را ثابت فرض کنیم مقدار بهینه مدل (۳) نسبت به متغیرهای انتخاب ویژگی (\tilde{w}) با استفاده از رابطه زیر بدست می‌آید:

$$\tilde{w}_{kj} = \begin{cases} 1 & f_{kj} \leq S \min\{f_{kl}\}; \\ 0 & \text{otherwise.} \end{cases}$$

که $f_{kl} = \sum_{i=1}^n u_i w_{ik}^{\beta} (x_{il} - z_{kl})^2$ و منظور از $S \min$ کمترین مقدار می‌باشد.

اثبات.

اگر مراکز خوشه‌ها (Z) و درجات عضویت (u) را ثابت فرض کنیم مدل (۳) بصورت زیر در می‌آید:

برای بدست آوردن مقدار بهینه متغیر Z کفایت مشتق تابع هدف مدل (۳) نسبت به Z را برابر مقدار صفر قرار دهیم:

$$\frac{\partial F}{\partial z_{kj}} = 0 \quad \rightarrow \sum_{i=1}^n -2u_{ik}w_{kj}^{\beta}\tilde{w}_{kj}(x_{ij} - z_{kj}) = 0;$$

$$\rightarrow 2w_{kj}^{\beta}\tilde{w}_{kj}\sum_{i=1}^n u_{ik}x_{ij} = 2w_{kj}^{\beta}\tilde{w}_{kj}z_{kj}\sum_{i=1}^n u_{ik}; \rightarrow z_{kj} = \frac{\sum_{i=1}^n u_{ik}x_{ij}}{\sum_{i=1}^n u_{ik}}.$$

پایان اثبات.

۴- آزمایش‌ها و ارزیابی

در این تحقیق از مجموعه داده‌های واقعی ذیل از مخزن UCI برای مقایسه تجربی روش کاهش ویژگی پیشنهادی با روش‌های کاهش ویژگی L1PCA, LLE و روش K-means تنگ استفاده می‌گردد:

- ۱) Wine: این مجموعه داده دارای ۳ خوشه با مجموعاً ۱۷۸ داده و ۱۳ ویژگی است.
- ۲) Vertebral2: این مجموعه داده دارای ۲ خوشه با مجموعاً ۳۱۰ داده و ۶ ویژگی است.
- ۳) Vertebral3: این مجموعه داده دارای ۳ خوشه با مجموعاً ۳۱۰ داده و ۶ ویژگی است.
- ۴) Parkinsons: این مجموعه داده دارای ۲ خوشه با مجموعاً ۱۹۵ داده و ۲۲ ویژگی است.

در هر آزمایش، یکی از روش‌های کاهش ویژگی بر روی یک مجموعه داده اعمال شده و سپس مجموعه داده کاهش بعد یافته، با استفاده از روش wk-means خوشه‌بندی می‌شود. روش‌های خوشه‌بندی و کاهش ابعاد مبتنی بر wk-means به شدت به مقادیر اولیه مرکز خوشه‌ها، که به صورت تصادفی انتخاب می‌شوند، وابسته است. لذا هر یک از آزمایش‌ها ۲۰ بار تکرار شده و میانگین و انحراف معیار نتایج گزارش می‌شود. در ضمن، مقدار پارامتر β در مدل wk-means و فاز اول روش پیشنهادی برابر با ۶ در نظر گرفته شده است.

$$\min_{\tilde{w}} \sum_{k=1}^c \sum_{j=1}^m \tilde{w}_{kj} \sum_{i=1}^n u_{ik} w_{kj}^{\beta} (x_{ij} - z_{kj})^2$$

$$\text{subject to } \begin{cases} \sum_{j=1}^m \tilde{w}_{kj} = S, \quad k = 1, 2, \dots, c; \\ \tilde{w}_{kj} \in \{0, 1\}, \quad k = 1, 2, \dots, c; j = 1, 2, \dots, m; \end{cases}$$

از آنجا که متغیرهای مدل فوق فقط \tilde{w}_{kj} ها هستند و در هر عبارت از مدل فوق، فقط یک متغیر وجود دارد، مدل فوق را نیز می‌توان بصورت جمع چند مدل کوچکتر و بصورت زیر نوشت:

$$\sum_{k=1}^c \left(\min_{\tilde{w}} \sum_{j=1}^m \tilde{w}_{kj} \sum_{i=1}^n u_{ik} w_{kj}^{\beta} (x_{ij} - z_{kj})^2 \right)$$

$$\text{subject to } \begin{cases} \sum_{j=1}^m \tilde{w}_{kj} = S; \\ \tilde{w}_{kj} \in \{0, 1\}, \quad j = 1, 2, \dots, m; \end{cases}$$

k-امین مدل از این مدل‌های کوچکتر را در نظر بگیرد:

$$\min_{\tilde{w}} \sum_{j=1}^m \tilde{w}_{kj} \sum_{i=1}^n u_{ik} w_{kj}^{\beta} (x_{ij} - z_{kj})^2$$

$$\text{subject to } \begin{cases} \sum_{j=1}^m \tilde{w}_{kj} = S; \\ \tilde{w}_{kj} \in \{0, 1\}, \quad j = 1, 2, \dots, m; \end{cases}$$

با توجه به قیود این مدل، فقط S تا از \tilde{w}_{kj} برابر با یک و سایر \tilde{w}_{kj} برابر با صفر می‌شوند. بنابراین برای حداقل شدن تابع هدف مدل، $\tilde{w}_{kj} = 1$ می‌شود. اگر ضریب آن در تابع هدف، یعنی $f_{kj} = \sum_{i=1}^n u_{ik} w_{kj}^{\beta} (x_{ij} - z_{kj})^2$ کمتر از S -امین کوچکترین ضریب \tilde{w}_{lj} باشد.

پایان اثبات.

قضیه ۳- اگر مقدار متغیرهای انتخاب ویژگی (\tilde{w}) و درجات عضویت (u) را ثابت فرض کنیم مقدار بهینه مدل (۲۶) نسبت به مراکز خوشه‌ها (Z) با استفاده از رابطه زیر به دست می‌آید:

$$z_{kj} = \frac{\sum_{i=1}^n u_{ik} x_{ij}}{\sum_{i=1}^n u_{ik}}$$

اثبات.

هنگامی که متغیرهای u و \tilde{w} را ثابت در نظر می‌گیریم

جدول ۱: مقایسه دقت خوشه بندی به روش wk-means با استفاده از داده‌های کاهش ابعاد یافته به روش‌های LLE, L1PCA و K-means تنگ و روش پیشنهادی

Vertebral2	Vertebral3	Parkinsons	Wine	الگوریتم کاهش ابعاد	تعداد ویژگی پس از کاهش ویژگی
۰,۶۲۶۰±۰,۰۳۰	۰,۴۹۳۵±۰,۰۰۸	۰,۷۱۵۳±۰,۰۵۰	۰,۶۵۱۰±۰,۰۵۵	روش پیشنهادی	D = 2
۰,۵۳۵۵±۰,۰۶۰	۰,۶۲۹۴±۰,۰۱۰	۰,۶۱۵۳±۰,۰۱۰	۰,۵۶۱۸±۰,۰۱۰	Sparse K-means[6]	
۰,۶۲۰۲±۰,۰۲۲	۰,۵۴۱۶±۰,۰۱۳	۰,۶۲۱۳±۰,۰۱۰	۰,۵۴۳۸±۰,۰۱۰	LLE[3]	
۰,۵۷۸۱±۰,۰۱۲	۰,۴۹۴۵±۰,۰۲۴	۰,۴۹۱۲±۰,۰۱۲	۰,۶۳۲۱±۰,۰۱۴	L1PCA[۲]	
۰,۶۱۰۶±۰,۰۱۶	۰,۳۳۵±۰,۰۳۰	۰,۷۴۷۶±۰,۰۲۳	۰,۶۰۸۹±۰,۰۸۶	روش پیشنهادی	D = 3
۰,۶۷۱۰±۰,۰۲۰	۰,۳۹۵۷±۰,۰۲۰	۰,۷۱۳۳±۰,۰۲۰	۰,۴۹۶۵±۰,۰۱۰	Sparse K-means[6]	
۰,۶۸۸۳±۰,۰۴۲	۰,۳۹۶۰±۰,۰۱۰	۰,۷۲۲۸±۰,۰۱۰	۰,۴۰۴۵±۰,۰۱۱	LLE[3]	
۰,۶۹۰۱±۰,۰۲۹	۰,۳۵۲۱±۰,۰۱۹	۰,۷۴۶۸±۰,۰۱۳	۰,۸۱۱۲±۰,۰۲۲	L1PCA[۲]	
۰,۶۳۳۴±۰,۰۱۱	۰,۵۵۸۰±۰,۰۱۵	۰,۷۶۷۶±۰,۰۲۰	۰,۶۹۱۵±۰,۰۵۳	روش پیشنهادی	D = 4
۰,۶۸۷۰±۰,۰۳۰	۰,۵۴۳۰±۰,۰۳۰	۰,۷۱۳۳±۰,۰۲۰	۰,۶۷۴۰±۰,۰۴۲	Sparse K-means[6]	
۰,۶۶۲۵±۰,۰۲۰	۰,۵۰۶۵±۰,۰۱۰	۰,۷۵۳۵±۰,۰۱۴	۰,۶۸۱۷±۰,۰۴۳	LLE[3]	
۰,۶۳۹۶±۰,۰۲۲	۰,۵۳۰۹±۰,۰۱۲	۰,۷۳۹۸±۰,۰۱۱	۰,۶۷۸۴±۰,۰۱۹	L1PCA[۲]	
۰,۶۴۶۴±۰,۰۱۱	۰,۵۱۸۳±۰,۰۱۸	۰,۷۹۱۷±۰,۰۳۴	۰,۶۳۱۴±۰,۰۹۳	روش پیشنهادی	D = 5
۰,۶۹۹۰±۰,۰۳۰	۰,۵۵۷۷±۰,۰۱۰	۰,۷۲۳۳±۰,۰۱۰	۰,۷۵۴۲±۰,۰۲۰	Sparse K-means[6]	
۰,۶۶۴۵±۰,۰۲۰	۰,۵۰۶۵±۰,۰۱۰	۰,۷۶۳۵±۰,۰۲۰	۰,۶۶۰۴±۰,۰۱۰	LLE[3]	
۰,۶۳۹۸±۰,۰۱۹	۰,۵۴۹۷±۰,۰۱۲	۰,۷۵۶۲±۰,۰۱۱	۰,۶۵۹۷±۰,۰۲۱	L1PCA[۲]	
-	-	۰,۷۱۸۴±۰,۰۳۳	۰,۶۶۷۹±۰,۰۵۴	روش پیشنهادی	D = 6
-	-	۰,۶۳۹۰±۰,۰۱۸	۰,۶۵۱۶±۰,۰۲۶	Sparse K-means[6]	
-	-	۰,۶۹۰۲±۰,۰۳۷	۰,۴۹۲۱±۰,۰۳۵	LLE[3]	
-	-	۰,۵۶۲۷±۰,۰۳۲	۰,۵۶۲۰±۰,۰۵۱	L1PCA[۲]	

جدول ۲: مقایسه نتایج خوشه بندی به روش wk-means بدون کاهش بعد داده‌ها.

Vertebral2	Vertebral3	Parkinsons	Wine
۰,۶۵۳۷±۰,۰۱۰	۰,۴۸۵۹±۰,۰۳۰	۰,۶۶۱۵±۰,۰۱۰	۰,۸۹۹۸±۰,۰۳۰

کاهش ویژگی و نوع مجموعه داده دارد. به عنوان مثال، میانگین دقت خوشه‌بندی مجموعه داده ۱۳ بعدی Wine پس از کاهش ویژگی به ۳ ویژگی با استفاده از روش L1PCA، ۱۲،۸۱٪ شده است، حال آنکه میانگین دقت

با توجه به جدول ۱ می‌توان به نکات زیر اشاره کرد: همانطور که در جدول ۱ مشخص است، دقت خوشه‌بندی هر مجموعه داده پس از کاهش ویژگی، بستگی به تعداد ابعاد پس از کاهش ویژگی یا میزان کاهش ویژگی، روش

خوشه‌بندی مجموعه داده ۶ بعدی Vertebral3 پس از کاهش ویژگی به ۳ ویژگی با استفاده از روش L1PCA، ۳۵،۲۱٪ شده است. در جدول ۱ بهترین روش کاهش ویژگی به ازای هر مجموعه داده که بیشترین دقت خوشه‌بندی را پس از کاهش ویژگی به دست آورده است به صورت پررنگ تر مشخص شده است.

در حالت کلی کاهش بیش از حد ابعاد در مجموعه داده‌های مختلف می‌تواند دقت خوشه‌بندی را به شدت کاهش دهد. برخی نتایج در ابعاد ۲ و ۳ تقریباً غیر قابل اتکا می‌باشد، لیکن با کاهش ابعاد داده به میزان مناسب در هر مجموعه داده، می‌توان امیدوار بود که روش پیشنهادی علاوه بر تامین اهداف روش‌های کاهش ابعاد (از جمله کاهش فضای ذخیره سازی، فشردن سازی و افزایش سرعت پردازش بر روی داده‌ها) دقت مناسبی را در حل مسائل داده کاوی ارائه خواهد داد.

مقایسه جدول ۱ و ۲ نشان می‌دهد که در بسیاری موارد، انتخاب ویژگی مناسب در مجموعه‌های داده می‌تواند نتایج بهتری را در مسائل داده کاوی ارائه دهد. به این معنا که وجود بعضی ویژگی‌ها در برخی از مجموعه داده‌ها می‌تواند نتایج تحلیل داده‌ها را منحرف کند و دقت خوشه‌بندی را کاهش دهد.

برای مقایسه کارایی چند خوشه‌بند بر روی بیش از دو مجموعه داده، استفاده از میانگین آنها نادرست است [۷]. لذا از روش رتبه‌بندی خوشه‌بندها برای هر مجموعه داده به‌طور مستقل از سایر مجموعه داده‌ها استفاده می‌گردد [۸]. روش پیشنهادی در مقایسه با روش L1PCA از ۱۸ حالت آزمایش شده ۱۳ بار رتبه بهتری به دست آورده است. لازم به ذکر است خوشه‌بندی داده‌های کاهش ابعاد یافته به روش پیشنهادی در ابعاد بالاتر عملکرد بهتری از خود نشان می‌دهد. روش پیشنهادی در ۲۲ حالت از ۱۳ حالت آزمایش نسبت به روش LLE عملکرد بهتری داشته است. همانطور که در جدول ۱ مشاهده می‌شود، روش پیشنهادی در اکثر آزمایش‌ها در ابعاد ۵ و ۶ از روش LLE پیشی

گرفته است. اگر بخواهیم روش پیشنهادی را با روش K-means ^{تُنک}، که یک روش مبتنی بر K-means می‌باشد، مقایسه کنیم، در جدول ۱ برتری روش پیشنهادی مشاهده می‌شود. چرا که روش پیشنهادی در ۱۲ مورد از ۱۸ مورد آزمایش این مقاله، عملکرد بهتری از خود نشان داده است. بنابراین، می‌توان گفت که دقت خوشه‌بندی پس از کاهش ویژگی با روش پیشنهادی بهتر از سایر روش‌های کاهش ویژگی مورد مقایسه در این پایان نامه است.

آزمایش‌های مذکور با استفاده از پردازنده ۶۴ بیتی Core i5 2.5Ghz به همراه حافظه اصلی 6GB و با استفاده از سیستم عامل Windows 7 Ultimate آزمایش شده است. جدول ۳ سرعت اجرای الگوریتم‌های کاهش ابعاد مورد آزمایش را نشان می‌دهد. براین اساس می‌توان گفت که سرعت اجرای الگوریتم پیشنهادی بهتر از روش LLE و بدتر از روش L1PCA و k-means ^{تُنک} است.

۵- نتیجه گیری

در این مقاله، با توجه به موفقیت روش‌های خوشه‌بندی مبتنی بر k-means، یک روش کاهش ویژگی بر پایه k-means وزن دار ارائه شد. به بیان دقیقتر، در روش پیشنهادی، نخست با استفاده از روش k-means وزن دار، به ویژگی‌های داده‌ها وزن داده شد. ویژگی‌های وزن‌تر ضرورتاً ویژگی‌های مهمتر نیستند و وزن هر ویژگی، تنها بازه هر ویژگی را به نحوی تغییر می‌دهد که خوشه‌بندی بهتری صورت بگیرد. لذا، سپس با استفاده از یک مدل ریاضی جدید، کسری از ویژگی‌های وزن‌دار شده داده‌های هر خوشه انتخاب شد به نحوی که کمترین تغییر در خوشه‌ها حاصل شود. تعداد ویژگی‌های منتخب هر خوشه در روش پیشنهادی، برخلاف روش‌های مشابهی چون k-means ^{تُنک} و fuzzy c-means ^{تُنک} بصورت صریح تعیین می‌شود. درضمن، آزمایش‌های تجربی روی چهار مجموعه‌داده واقعی نشان می‌دهد که روش پیشنهادی، از دقت بیشتری نسبت به روش‌های L1PCA، LLE و روش

K-means تنگ برخوردار است.

روش پیشنهادی در ۲۲ حالت از ۱۳ حالت آزمایش نسبت به روش LLE عملکرد بهتری داشته است. روش پیشنهادی در ۱۲ مورد از ۱۸ مورد آزمایش این مقاله، عملکرد بهتری نسبت به روش k-means تنگ داشته است. بنابراین، می‌توان گفت که دقت خوشه‌بندی پس از کاهش ویژگی با روش پیشنهادی بهتر از سایر روش‌های کاهش ویژگی مورد مقایسه در این مقاله است. در ضمن، سرعت اجرای الگوریتم پیشنهادی بهتر از روش LLE و بدتر از روش L1PCA و k-means تنگ است.

برای مقایسه روش کاهش ویژگی پیشنهادی با سایر روش‌های کاهش ویژگی، در هر آزمایش، یکی از روش‌های کاهش ویژگی بر روی یک مجموعه داده واقعی اعمال شده و سپس مجموعه داده کاهش بعد یافته، با استفاده از روش wk-means خوشه بندی شدند. آزمایش‌های تجربی روی چهار مجموعه داده واقعی نشان داد که روش پیشنهادی در مقایسه با روش L1PCA از ۱۸ حالت آزمایش شده ۱۳ بار رتبه بهتری بدست آورده است.

جدول ۳: مقایسه سرعت اجرای الگوریتم‌های کاهش ابعاد LLE, L1PCA و K-means تنگ و روش پیشنهادی (به ثانیه).

Vertebral2	Vertebral3	Parkinsons	Wine	الگوریتم کاهش ابعاد	اندازه ابعاد آزمایش
۰,۰۴۰	۰,۰۴۸	۰,۰۷۷	۰,۰۵۸	روش پیشنهادی	D = 2
۰,۰۲۷	۰,۰۳۸	۰,۰۳۶	۰,۰۲۵	Sparse K-means[6]	
۰,۰۹۴	۰,۱۱۱	۰,۱۶۸	۰,۰۸۰	LLE[3]	
۰,۰۱۷	۰,۰۲۸	۰,۰۲۵	۰,۰۱۳	L1PCA[۲]	
۰,۰۴۰	۰,۰۴۵	۰,۰۷۷	۰,۰۵۵	روش پیشنهادی	D = 3
۰,۰۲۸	۰,۰۳۱	۰,۰۳۶	۰,۰۲۱	Sparse K-means[6]	
۰,۱۱۱	۰,۱۱۱	۰,۱۵۳	۰,۰۸۰	LLE[3]	
۰,۰۱۷	۰,۰۲۶	۰,۰۲۵	۰,۰۱۷	L1PCA[۲]	
۰,۰۴۰	۰,۰۴۸	۰,۰۹۴	۰,۰۶۱	روش پیشنهادی	D = 4
۰,۰۳۳	۰,۰۳۶	۰,۰۳۶	۰,۰۲۵	Sparse K-means[6]	
۰,۱۱۱	۰,۱۰۸	۰,۱۷۰	۰,۰۷۵	LLE[3]	
۰,۰۱۹	۰,۰۱۸	۰,۰۲۱	۰,۰۲۳	L1PCA[۲]	
۰,۰۴۰	۰,۰۴۹	۰,۰۹۷	۰,۰۵۷	روش پیشنهادی	D = 5
۰,۰۳۵	۰,۰۳۹	۰,۰۴۵	۰,۰۲۶	Sparse K-means[6]	
۰,۱۱۶	۰,۰۹۶	۰,۱۷۲	۰,۰۷۰	LLE[3]	
۰,۰۲۹	۰,۰۳۱	۰,۰۳۲	۰,۰۱۹	L1PCA[۲]	
-	-	۰,۰۹۴	۰,۰۶۴	روش پیشنهادی	D = 6
-	-	۰,۰۴۵	۰,۰۲۷	Sparse K-means[6]	
-	-	۰,۱۷۲	۰,۰۷۳	LLE[3]	
-	-	۰,۰۳۹	۰,۰۲۱	L1PCA[۲]	

1. Ding, C. and T. Li. *Adaptive dimension reduction using discriminant analysis and k-means clustering*. in *Proceedings of the 24th International Conference on Machine Learning*. 2007. Corvallis.
2. Markopoulos, P.P., et al., *Efficient L1-Norm Principal-Component Analysis via Bit Flipping*. *IEEE Transactions on Signal Processing*, 2017. **65**(16): p. 4252-4264.
3. Rovis, S.T. and L.K. Saul, *Nonlinear Dimensionality Reduction by Locally Linear Embedding*. *Science*, 2000. **290**: p. 2323-2326.
4. Qiu, X., et al., *A sparse fuzzy c-means algorithm based on sparse clustering frame work*. *Jornal of Elsevier on Neurocomputing*, 2015. **157**: p. 290-295.
5. Huang, X., Y. Ye, and H. Zhang, *Extensions of Kmeans-Type Algorithms: A New Clustering Framework by Integrating Intracluster Compactness and Intercluster Separation*. *IEEE Transactions on Neural Networks and Learning Systems*, 2014. **25**(8): p. 1433 - 1446.
6. Tibshirani, R. and D.M. Witten, *A framework for feature selection in clustering*. *NIH-PA Author Manuscript*, 2010. **105**(490): p. 713–726.
7. Demsar, J., *Statistical comparison of classifier over multiple data sets*. *Journal of Machine Learning*, 2006. **7**: p. 1-30.
8. García, S., A. Fernández, and J. Luengo, *Advanced nonparametric tests for multiple comparisons in the design of experiments in computational intelligence and data mining: Experimental analysis of power*. *Information Sciences*, 2010. **180**: p. 2044-2064.