



Automatic Clustering of Mixed Data Using Genetic Algorithm

M.Yaghini* & M.Vard

*Masoud Yaghini, Assistance professor of School of Railway Engineering - Iran University of Science and Technology
Mahdi Vard, MSc, School of Railway Engineering - Iran University of Science and Technology*

Keywords

Data mining
Clustering
Mixed data
Genetic algorithm
Davies-Bouldin index

ABSTRACT

In the real world clustering problems, it is often encountered to perform cluster analysis on data sets with mixed numeric and categorical values. However, most existing clustering algorithms are only efficient for the numeric data rather than the mixed data set. In addition, traditional methods, for example, the K-means algorithm, usually ask the user to provide the number of clusters. In this paper, we propose a new method to cluster mixed data and automatically evolve the number of clusters as well as clustering of data set. In the proposed method, Davies-Bouldin Index is used as fitness function and we use the genetic algorithm to optimize fitness function. Also, we use a more accurate distance measure for calculating the distance between categorical values. The performance of this algorithm has been studied on real world and simulated data sets. Comparisons with other clustering algorithms illustrate the effectiveness of this approach.

© 2012 IUST Publication, IJIEPM. Vol. 23, No. 2, All Rights Reserved

*
Corresponding author: Masoud Yaghini
Email: yaghini@iust.ac.ir

۲. موری بر ادبیات موضوع

۱-۲. روش‌های خوشبندی مبتنی بر الگوریتم ژنتیک

کریشنا و مورتی [۷] یک روش خوشبندی مبتنی بر الگوریتم ژنتیک (GKA) ارائه کردند که می‌توان برای حل مساله خوشه‌بندی با تعداد خوشه‌های مشخص از آن بهره گرفت.

لو و همکاران [۸] الگوریتم سریع k-means مبتنی بر ژنتیک (FGKA)^۱ را با الهام از الگوریتم GKA ابداع نمودند که این روش در بسیاری از جنبه‌ها نسبت به GKA بهبود داده شده است. از جمله اینکه روش FGKA بسیار سریع‌تر از GKA عمل می‌کند.

لو و همکاران [۹] الگوریتم افزاینده k-means ژنتیک^۲ (IGKA) را طراحی کردند که توسعه‌ای بر الگوریتم خوشبندی قابلی آنها (روش FGKA) بود. روش IGKA این ویژگی برجسته را از روش FGKA به ارت برده است که همواره به جواب بهینه کلی همگرا می‌باشد.

ماولیک و همکاران [۱۰] یک روش خوشبندی مبتنی بر الگوریتم ژنتیک پیشنهاد کردند که در آن از قابلیت جستجوی الگوریتم ژنتیک به منظور تعیین K مرکز خوشه در فضای R^N بهره گرفته می‌شود. در این روش مقدار K از پیش معلوم در نظر گرفته شده است و یکی از مقادیر ورودی مساله می‌باشد.

در روش دیگری که توسط لین و همکاران [۱۱] ارائه گردید، مراکز خوشه‌ها مستقیماً از میان نقاط مجموعه داده انتخاب می‌شود. اتخاذ این رویکرد سبب سرعت بیشتر در محاسبه تابع برازنده و در نتیجه سریع تر شدن کل الگوریتم می‌گردد.

بندیوپد و ماولیک [۱۲] روشی را ارائه دادند که در آن همزمان با خوشه‌بندی داده‌ها، مقدار مناسب برای تعداد خوشه‌ها نیز تعیین می‌شود. در این روش جهت اعتبار سنجی خوشه‌های حاصله از شاخص Davies–Bouldin استفاده شده است.

یک روش خوشه‌بندی ترکیبی مبتنی بر الگوریتم ژنتیک تحت عنوان روش خوشه‌بندی HGA توسط لیو و همکاران [۱۳] ارائه شد. این الگوریتم با بهره گیری از لیست ممنوعه و معیار انتظار، بین تنوع در جمعیت و سرعت همگرایی، هماهنگی ایجاد می‌کند.

چیانگ و همکاران [۱۴]، روش k-modes را که با ایجاد تغییراتی در روش k-means به خوشه‌بندی داده‌های دسته‌ای می‌پردازد، توسعه دادند. به این ترتیب که با بهره گیری از الگوریتم ژنتیک، شاخصی برای سنجش عدم تشابه با نام مقیاس فاصله ژنتیک (GDM) طراحی نمودند.

روی داده‌های مختلط را دارا باشد، بسیار مورد توجه و حائز اهمیت خواهد بود و چنین روشی می‌تواند در حل بسیاری از مسائل دنیای واقعی و تحلیل پایگاه‌های داده در سازمانهای مختلف مورد استفاده قرار گیرد.

برای طراحی روشی که بتواند علاوه بر داده‌های عددی، بر روی داده‌های دسته‌ای نیز قابلیت کاربرد داشته باشد، نخستین موضوعی که باید به آن پرداخته شود طراحی معیار سنجش فاصله میان داده‌های دسته‌ای است. در این زمینه روش‌های مختلفی ارائه شده است که روش هوآنگ [۳] با توجه به سادگی و کاربرد آسان، از روش‌های مورد توجه است. اما روش هوآنگ علیرغم سادگی، از دقت خوبی برخوردار نبوده و دارای نواقصی است که سبب می‌شود عملکرد الگوریتم خوشبندی با این معیار فاصله از کیفیت خوبی برخوردار نباشد. لذا در این مقاله سعی نموده‌ایم با بهره گیری از یک معیار فاصله دقیق‌تر، الگوریتمی را طراحی کنیم که از دقت بالاتری نسبت به دیگر روش‌های خوشبندی داده‌های مختلط برخوردار باشد.

یکی دیگر از بخش‌های تشکیل‌دهنده یک روش خوشبندی، انتخاب رویکردی جهت جستجوی فضای جواب و بهینه‌سازی تابع هدف است. یکی از تکنیک‌های جستجوی فضای جواب و بهینه‌سازی، روش‌های فرا ابتکاری هستند که از آن جمله می‌توان به الگوریتم ژنتیک [۴]، الگوریتم کلونی مورچگان [۵] و روش جستجوی من نوع [۶] اشاره نمود.

از میان روش‌های فرا ابتکاری، الگوریتم ژنتیک با توجه به قدرت و قابلیت بالایی که در جستجوی فضای جواب دارد، از اهمیت ویژه‌ای برخوردار است و لذا بسیاری از روش‌های خوشبندی نیز از این روش به عنوان ابزار مورد استفاده جهت بهینه‌سازی تابع هدف استفاده نموده‌اند. لذا با توجه به قابلیت و کارایی بالای الگوریتم ژنتیک، ما نیز در طراحی روش خوشبندی خود از این الگوریتم جهت جستجوی فضای جواب بهره برده‌ایم.

یکی از ویژگیهای روش خوشبندی ارائه شده در این مقاله، محاسبه تعداد بهینه خوشبندی است؛ به این معنی که برخلاف بسیاری از روش‌های خوشبندی که تعداد خوشه‌ها را به عنوان یک مقدار ورودی از کاربر دریافت می‌کنند، روش بیشنهادی ما می‌تواند به موازات گروه‌بندی داده‌ها و تشکیل خوشه‌های بهینه، مقدار بهینه برای تعداد خوشه‌ها را نیز محاسبه و ارائه دهد.

در ادامه ابتدا موری بر روش‌های خوشبندی مبتنی بر الگوریتم ژنتیک و نیز روش‌های خوشبندی داده‌های مختلط صورت گرفته است. سپس به تشریح روش مورد استفاده برای محاسبه فاصله بین مقادیر داده‌های دسته‌ای خواهیم پرداخت و پس از آن مراحل الگوریتم ژنتیک مورد استفاده در فرآیند خوشبندی معرفی خواهد شد.

^۱ Fast Genetic K-means Algorithm

^۲ Incremental Genetic K-means Algorithm

الگوریتم‌های خوشبندی مبتنی بر افزای داده‌ها و بر روی مجموعه‌های داده مختلط قابل استفاده است. نحوه عملکرد این تابع هزینه به این شکل است که تشابه بین دو عنصر از مجموعه داده‌ها را به صورت مجموع دو مقدار فاصله، یکی برای مشخصه‌های عددی و دیگری برای مشخصه‌های دسته‌ای، محاسبه می‌نماید. از آنجا که تابع هزینه هوآنگ قابلیت کاربرد با الگوریتم‌های مبتنی بر افزای دارد، هزینه‌های محاسباتی آن در حد مناسب و قابل قبولی است.

پس از آن، هوآنگ و همکاران [۲۱] روش خوشبندی k-prototypes را برای اجرا بر روی داده‌های مختلط ارائه کردند. در این روش وزن هر یک از مشخصه‌ها به صورت خودکار بر اساس افزای کنونی داده‌ها محاسبه می‌شود.

لو و همکاران [۲۲] روشی را پیشنهاد کردند که در آن مشخصه‌های عددی و دسته‌ای به صورت جداگانه خوشبندی می‌شوند و از تکیک جمع آوری شواهد برای ترکیب نتایج خوشبندی‌ها و حصول خوشبندی نهایی استفاده می‌شود.

هی و همکاران [۲۳]، روش قبلی خود را که به خوشبندی داده‌های دسته‌ای می‌پرداخت و الگوریتم فشرده^۱ نام داشت توسعه دادند و روشی را پیشنهاد کردند که قابلیت کاربرد بر روی داده‌های مختلط را دارد.

احمد و دی [۱۸] تابع هزینه ای را ارائه دادند که در آن سعی شده با بهبود و رفع نقاطی تابع هوآنگ، فرآیند خوشبندی با کیفیت بیشتری صورت گرفته و جوابهای بهتری حاصل آید. نخستین تفاوت تابع هزینه ارائه شده توسط احمد با تابع هوآنگ در اینست که تابع هوآنگ برای محاسبه فاصله بین داده‌های دسته‌ای، براساس تطابق یا عدم تطابق مقدار مشخصه موردنظر در دو شیء مورد بررسی، یکی از مقادیر صفر یا یک را (صفر برای تطابق و یک برای عدم تطابق) به عنوان مقدار فاصله تخصیص می‌دهد. اما احمد و دی یک تابع فاصله پیوسته را تعریف نمودند که با توجه به مقادیر سایر مشخصه‌ها برای دو شیء مورد بررسی، مقادیر بین صفر و یک را محاسبه و آنرا به عنوان فاصله دو شیء در نظر می‌گیرد.

۲. معرفی روش خوشبندی پیشنهادی

روش خوشبندی پیشنهادی در این مقاله از نوع روش‌های افزای داده‌ها و مبتنی بر الگوریتم ژنتیک است که در آن از الگوریتم ژنتیک جهت جستجوی جواب بهینه استفاده می‌شود. همچنین این روش قادر به خوشبندی داده‌های مختلط می‌باشد. به طور کلی مهم‌ترین مشخصات روش پیشنهادی عبارتند از:

^۴ Squeezed Algorithm

الگوریتم k-means ژنتیک وزن دار^۱ (GWKMA) که تلفیقی از الگوریتم ژنتیک و الگوریتم k-means وزن دار است، توسط ژبانگ و همکاران [۱۵] پیشنهاد گردید.

HGACLUS مدلی ترکیبی جهت خوشبندی بر مبنای الگوریتم ژنتیک است که توسط پن و ژو [۱۶] ارائه گردید. در این مدل از روش Simulated Annealing برای یافتن مراکز خوشبندی استفاده شده است.

کاتاری و همکاران [۱۷]، روشی را برای خوشبندی داده‌ها ارائه نمودند که در آن از الگوریتم ژنتیک بهبود یافته^۲ استفاده شده و عملگرهای باز ترکیب و جهش به شکل کاراتری تعریف شده اند.

۲-۲. روش‌های خوشبندی داده‌های مختلط

امروزه با پایگاههای داده بسیار بزرگی مواجه هستیم که مشتمل بر مشخصه‌های مختلط هستند. در حالیکه اغلب روش‌های خوشبندی که تا کنون ارائه شده اند تنها بر روی داده‌های عددی و یا دسته‌ای عملکرد خوبی داشته و قابلیت کاربرد بر روی داده‌هایی که دارای مشخصه‌هایی از هر دو نوع عددی و دسته‌ای هستند را ندارند.

برای غلبه بر این مشکل، برخی از استراتژی‌های به کار گرفته شده به شرح زیر می‌باشد [۱۸]:

۱) مشخصه‌های دسته‌ای به مقادیر عدد صحیح تبدیل شده و سپس مقیاسهای موجود برای اندازه گیری فواصل داده‌های عددی برای محاسبه تشابه بین هر جفت از داده‌ها به کار گرفته می‌شود. در این روش، تخصیص مقادیر عددی صحیح به مقادیر دسته‌ای نظیر رنگ کاری بسیار مشکل است.

۲) رویکرد دیگر به این صورت است که مقادیر مشخصه‌های عددی را گسسته سازی می‌کنند و از این طریق آنها را به صورت مقادیر دسته‌ای در می‌آورند و سپس از الگوریتم خوشبندی داده‌های دسته‌ای استفاده می‌نمایند. اشکال این نوع رویکرد اینست که فرآیند گسسته‌سازی مقادیر عددی با از دست دادن اطلاعات توأم است.

لی و بیزوگار [۱۹] یک الگوریتم خوشبندی ترکیب کننده مبتنی بر تشابهات^۳ را ارائه کردند که بر اساس شاخص تشابه گodal [۲۰] عمل می‌کرد. این الگوریتم بر روی مشخصه‌های عددی و دسته‌ای به خوبی عمل می‌کند، اما اشکال آن اینست که از نظر محاسباتی هزینه زیادی دارد.

هوانگ [۳] تابع هزینه ای پیشنهاد کرد که مشخصه‌های عددی و دسته‌ای را به صورت مجزا در نظر می‌گیرد. این تابع هزینه در

¹ Genetic Weighted K-means Algorithm

² Improved Genetic Algorithm (IGA)

³ Similarity Based Agglomerative Clustering (SBAC)

۲-۳. تعیین فاصله میان دو داده

فرض کنید که D_1 و D_2 دو داده از مجموعه داده های مختلط باشند که مجموعاً دارای m مشخصه می باشند که m_r مشخصه اول عددی و m_c مشخصه بعدی دسته‌ای هستند و فاصله بین D_1 و D_2 برابر است با:

$$Dist(D_1, D_2) = \sum_{t=1}^{m_r} (w_t(X_t - Y_t))^2 + \sum_{t=1}^{m_c} (\delta(X_t, Y_t))^2 \quad (4)$$

۳-۳. محاسبه مراکز خوشة ها برای داده های مختلط

تعریف اصلاح شده مرکز خوشه که در اینجا ارائه شده است، با نحوه تعریف مراکز خوشه ها در خوشه بندی فازی شباهت های دارد. اما در این مقاله از این تعریف برای مراکز خوشه ها در حالت خوشه بندی با مز مشخص^۱ استفاده شده است.

در روش پیشنهادی برای تعریف مراکز خوشه ها، مقدار مرکزی به ازای مشخصه های عددی همچنان با مقدار میانگین نمایش داده می شود. اما برای مشخصه های دسته های از نحوه نمایش متفاوتی استفاده شده است. از آنجا که در روش پیشنهادی فاصله بین دو مقدار دسته های براساس توزیع کلی آنها در سراسر مجموعه داده ها تعریف می شود، این مقدار فاصله به ازای زوجهای مختلف از مقادیر، متفاوت خواهد بود.

بنابراین اگر به عنوان مثال فاصله مقدار r تا مقدار s کمتر از فاصله r تا r باشد، یعنی $\delta(r, s) < \delta(r, t)$ آنگاه انتظار می‌رود که در یک خوش بندی مناسب از داده‌ها، تعداد رخداد‌های هم‌زمان r و s از تعداد رخداد‌های هم‌زمان r و t بیشتر باشد. با در نظر گرفتن این مطالب، مقدار مرکزی a این مشخصه دسته‌های برای خوش C به شکل زیر محاسبه می‌گردد:

$$1/N_c \left\langle \left(N_{1,1,c}, N_{1,2,c}, \dots, N_{1,p_1,c} \right), \dots, \left(N_{m,1,c}, N_{m,2,c}, \dots, N_{m,p_m,c} \right) \right\rangle \quad (\textcircled{d})$$

در رابطه فوق، N_c تعداد داده های موجود در خوشه C را نشان می دهد، $N_{i,k,c}$ نمایانگر تعداد داده هایی در خوشه C است که مشخصه i ام آنها دارای k امین مقدار ممکن باشد، با این فرض که مشخصه i ام دارای p_i مقدار مختلف باشد.

در نتیجه مرکز خوشه، توزیع نسبی هر یک از مقادیر دسته ای را در خوشه مورد نظر نشان می دهد.

- (۱) قابلیت کار بر روی داده‌های مختلط
 - (۲) تعیین مقدار بهینه برای تعداد خوشšeها
 - (۳) استفاده از روش دقیق‌تری برای تعیین فاصله بین داده‌های دسته‌ای
 - (۴) استفاده از شاخص Davies-Bouldin به عنوان تابع برآزندگی نوآوری مقاله حاضر، ارائه الگوریتمی است که مجموعه ویژگیهای فوق را به صورت توان دارا می‌پاشد.

۱-۳. روش مورد استفاده برای محاسبه فاصله بین دو مقدار از یک متغیر دسته‌ای

در روش خوشبندی پیشنهادی در این مقاله، از تابع فاصله ارائه شده توسط احمد و دی [۱۸] جهت محاسبه فاصله میان نقاط و نیز محاسبه مراکز خوشبها استفاده شده است و جهت پیاده‌سازی الگوریتم ژنتیک، نحوه نمایش جوابها در کروموزومها و نیز عملگرهای الگوریتم ژنتیک مناسب با این تابع فاصله تعریف شده است.

- تعریف ۱ -

فاصله بین دو مقدار y , x از متغیر A_i نسبت به متغیر A_j و یک زیرمجموعه خاص \mathcal{W} به صورت زیر تعریف می شود:

$$\delta_w^i(x, y) = P_i(w|x) + P_i(\sim w|y) \quad (1)$$

- ۲ - تعریف

فاصله بین دو مقدار x و y از مشخصه A_i نسبت به مشخصه A_j به صورت زیر تعریف می شود:

$$\delta^{ij}(x, y) = P_i(\omega|x) + P_i(\sim\omega|y) - 1 \quad (5)$$

که در رابطه فوق، ω ، آن زیر مجموعه‌ای از مقادیر A_i است
که به ازای آن مقدار عبارت $P_i(\omega|x) + P_i(\sim\omega|y)$ مانگزیم گردید.

تعريف - ٣

برای یک مجموعه از داده‌ها که از m مشخصه، شامل مشخصه‌های عددی و دسته‌ای، تشکیل شده است، و مشخصه‌های عددی را در آن به صورت گسسته در آورده‌ایم، فاصله بین دو مقدار y از یک مشخصه دسته‌ای نسبت به یکدیگر برابر خواهد بود با:

$$\delta(x, y) = (1/m - 1) \sum_{j=1 \dots m, j \neq i} \delta^{ij}(x, y) \quad (\dagger)$$

تخصیص داده می‌شود تا مشخص شود که ژن مربوطه خالی است و مرکز خوش ای در آن قرار نگرفته است.

۳-۵-۲. مقدار دهی اولیه جمعیت

به ازای هر کروموزوم i در جمعیت (P , ..., $P_{i=1}$) و P برابر با اندازه جمعیت است، یک مقدار تصادفی k_i در بازه تعریف شده تولید می‌شود. سپس k_i نقطه به صورت تصادفی از میان داده‌ها انتخاب می‌شود و به صورت تصادفی در میان ژنهای کروموزوم قرار داده می‌شود. در نهایت به ژنهای خالی کروموزوم مقدار (۱) تخصیص داده می‌شود.

۳-۵-۳. عملگرهای تغییر

در روش پیشنهادی از عملگر باز ترکیب تک نقطه‌ای استفاده شده است. مقدار p_c در آزمایش‌های مختلف بین ۰/۵ تا ۰/۷ قرار داده شد و نتایج حاصله مقایسه شد و در نهایت ترخ بازترکیب برابر ۰/۵ انتخاب گردید. در مورد عملگر جهش نیز یک تعریف جدید برای اعمال عملگر جهش بر روی مقدار دسته‌ای ابداع شد. همانطور که در بخش قبلی اشاره شد، نمایش مختصات مراکز خوش اها در مورد مشخصه‌های دسته‌ای به صورت نسبت فراوانی هر یک از مقدار مشخصه مذبور در میان نقاط موجود در خوش مورد نظر است. برای اجرای عملگر جهش، در صورت انتخاب مشخصه دسته‌ای A_i (احتمال انتخاب هر یک از مشخصه‌ها برای اجرای عملگر جهش برای با مقدار p_m است) از ردیف مربوط به مشخصه i در ماتریس نمایش مختصات مرکز خوشها، دو نسبت a_{ij} و a_{ik} انتخاب می‌شود و مقدار آنها با هم عوض می‌شود.

۳-۵-۴. تابع برازنده‌گی

یک دیگر از مواردی که در طراحی الگوریتم‌های خوشبندی باید مدنظر قرار گیرد، انتخاب مقیاس اعتبار مناسب جهت انتخاب به عنوان تابع برازنده‌گی است. شاخص‌های اعتبار مختلفی نظیر شاخص Dunn، شاخص XB (Xie-Beni)، شاخص BM و شاخص DB در این زمینه ارائه شده‌اند. شاخص DB که به صورت تابعی از نسبت مجموع پراکندگی نقاط در داخل خوش به جدایی بین خوشها تعریف می‌شود، در مقایسه با سایر شاخصهایی که در بالا به آنها اشاره شد نتایج دقیق‌تری را به دست می‌دهد. در روش ارائه شده در این مقاله از شاخص DB به عنوان تابع برازنده‌گی استفاده شده است. مقدار کوچکتر این شاخص نشاندهنده خوش بندی داده‌ها به نحوی بهتر خواهد بود. با توجه به اینکه فرآیند الگوریتم ژنتیک به دنبال بیشینه سازی

۴-۳. فاصله بین یک داده و مرکز خوش ای متناظرش

فاصله بین یک داده و مرکز خوش ای متناظرش برابر با مجموع فواصل مقادیر عددی و دسته ای می‌باشد. در مورد مشخصه‌های عددی، فاصله اقلیدسی میان مقدار مشخصه عددی و میانگین مقادیر آن مشخصه در خوش مورد نظر مورد استفاده قرار می‌گیرد. اما در مورد مشخصه‌های دسته‌ای، تمامی مقادیر ممکن آن مشخصه، همانطور که در بخش قبلی مشاهده شد، سهمی نسبی را در تعریف مرکز خوش دارا می‌باشد. به ازای مشخصه دسته ای A_i ، اگر مقدار مشخصه برای داده مورد نظر برابر با v باشد، فاصله بین این داده و مرکز خوش به صورت تابع وزن داری از مقادیر $\delta(r, v)$ محاسبه می‌شود که در آن، v تمامی مقادیر ممکن مشخصه A_i را اختیار می‌کند.

از آنجا که مرکز خوش دارای نمایشی به صورت نسبتی از تک تک مقادیر ممکن مشخصه‌های دسته ای می‌باشد، به هر یک از مقادیر فاصله $\delta(r, v)$ یک ضریب وزنی که نمایانگر نسبت حضور مقدار v در خوش است، تخصیص داده می‌شود. فرض کنید $a_{i,k}$ نمایانگر k -امین مقدار ممکن برای مشخصه دسته ای A_i باشد. همچنین فرض کنید تعداد مقادیر متمایز برای مشخصه A_i برابر با p_i باشد. با این فرضیات، فاصله به صورت رابطه زیر تعریف می‌گردد:

$$\Omega(X, C) = \left(N_{i,1,c} / N_c \right) * \delta(X, A_{i,1}) + \dots + \left(N_{i,p_i,c} / N_c \right) * \delta(X, A_{i,p_i}) \quad (6)$$

در نهایت فاصله کل میان یک داده و یک مرکز خوش برای مجموعه داده‌ای شامل داده‌های مختلط به صورت زیر تعریف خواهد شد:

$$D(d_i, C_j) = \sum_{t=1}^{m_r} (d_{it}^r - C_{jt}^r)^2 + \sum_{t=1}^{m_c} (\Omega(d_{it}^c, C_{jt}^c))^2 \quad (7)$$

۳-۵-۵. اجزاء الگوریتم ژنتیک در روش پیشنهادی

۳-۵-۶. نحوه نمایش رشته ها

در روش پیشنهادی، کروموزوم‌ها از اعداد حقیقی تشکیل شده‌اند و مقادیر و مختصات مربوط به مراکز خوشها را در خود جای داده اند. طول کروموزوم‌ها ثابت و برابر مقدار k_{max} است. مقدار k یعنی تعداد خوشها به صورت تصادفی از بازه $[k_{min}, k_{max}]$ انتخاب می‌شود که مقادیر w و k_{max} کمینه و مаксیمه مساله بوده که می‌بایست توسط کاربر معین شوند. پس از مشخص شدن مقدار k تعداد k ژن، مراکز خوشها را در خود جای می‌دهند و به مابقی ژنهای یک عدد خاص (در روش پیشنهادی مقدار ۱) نشانه دهند.

^۱ UCI استخراج شده اند. مجموعه داده‌های مورد استفاده به صورت از پیش طبقه بندی شده هستند و کلاس متناظر با هر داده از پیش معلوم است. در نتیجه برای سنجش میزان دقت الگوریتم، از میزان انطباق نحوه خوش بندی داده‌ها با کلاسهای واقعی آنها استفاده گردیده است. در ادامه نتایج اجرای الگوریتم بر روی هر یک از این دو مجموعه داده استاندارد آورده شده است.

جدول ۱. مقادیر پارامترهای ورودی الگوریتم پیشنهادی

مقدار پارامتر	پارامتر ورودی
2	حداقل تعداد خوش‌ها (k_{min})
15	حداکثر تعداد خوش‌ها (k_{max})
25	حداکثر تعداد نسلها (max-gen)
0.5	نرخ باز ترکیب (φ_c)
0.01	نرخ جهش (p_m)
40	اندازه جمعیت

۲-۱. داده‌های بیماران قلبی^۲

این داده‌ها اطلاعات مربوط به تعدادی از بیماران قلبی را شامل می‌شود و در کلینیک کلولند تولید شده است. پایگاه داده اصلی شامل ۷۶ مشخصه است، اما مقالات تحقیقی مختلف برای آزمایش الگوریتم خود از یک زیر مجموعه از این پایگاه داده که شامل ۱۴ مشخصه است استفاده نموده اند. مجموعه داده‌های بیماران قلبی یک مجموعه داده مختلط است که شامل ۹ مشخصه دسته‌ای و ۵ مشخصه عددی می‌باشد. این مجموعه داده شامل ۳۰۳ نمونه است که در دو کلاس طبقه بندی شده اند و مجموعاً ۱۶۴ نمونه متعلق به کلاس نرمال (عدم بیماری) و ۱۳۹ مورد متعلق به کلاس بیمار می‌باشند.

جدول ۲ نتایج حاصل از خوش بندی این مجموعه داده را به وسیله روش خوش بندی پیشنهادی در این مقاله نشان می‌دهد. همچنین نتایج حاصل از پنج روش خوش بندی داده‌های مختلط که برای آزمایش نتایج خود از مجموعه داده بیماران قلبی استفاده کرده اند، یعنی روش‌های SBAC^[۱۹] روش ECOWEB^[۲۴]، روش COBWEB^[۲۵] و روش هوآنگ^[۳] آورده شده است. مقادیر دقت بدست آمده برای این الگوریتمها، دقت بیشتر و برتری روش خوش بندی ارائه شده را نسبت به سایر روش‌ها نشان می‌دهد.

تابع هدف است، معکوس این شاخص به عنوان مقدار تابع برازنگی تعريف شده است. روابط مربوط به محاسبه شاخص DB در ادامه ارائه می‌گردد.

$$S_i = \frac{1}{|C_i|} \sum_{x \in C_i} \|x - z_i\| \quad (8)$$

$$R_i = \underset{j, j \neq i}{Max} \left\{ \frac{S_i + S_j}{d_{ij}} \right\} \quad (9)$$

$$d_{ij} = d(C_i, C_j) = \|z_i - z_j\| \quad (10)$$

$$DB_r = \frac{1}{k_r} \sum_{i=1}^{k_r} R_i \quad (11)$$

$$Fitness(Ch_r) = \frac{1}{DB_r} \quad (12)$$

که در روابط فوق C_i نمایانگر خوش نام، $|C_i|$ بیانگر تعداد نقاط موجود در خوش نام، z_i نشانده‌نده مرکز خوش نام و DB_r نمایانگر مقدار شاخص DB برای کروموزوم r است.

۴. آزمایش الگوریتم پیشنهادی

در این بخش نتایج اجرای الگوریتم خوش بندی پیشنهادی بر روی مجموعه داده‌های استاندارد و داده‌های شبیه سازی شده آورده شده است. در ادامه دقت عملکرد روش پیشنهادی نسبت به روش‌های پیشین مقایسه گردیده است.

۱-۱. تنظیم پارامترهای ورودی الگوریتم

با استفاده از مسایل شبیه سازی شده، مقادیر پارامترهای ورودی الگوریتم شامل تعداد جمعیت، نرخ باز ترکیب، نرخ جهش، حداکثر تعداد نسلها و نیز k_{min} تعیین گردید که مقادیر این پارامترها در جدول ۱ آورده شده است. همچنین در الگوریتم پیشنهادی، انتخاب والدها بر اساس روش متناسب با برازنگی و انتخاب بازماندها بر اساس سن کروموزوم‌ها صورت می‌گیرد.

۲-۱. نتایج الگوریتم بر روی داده‌های استاندارد

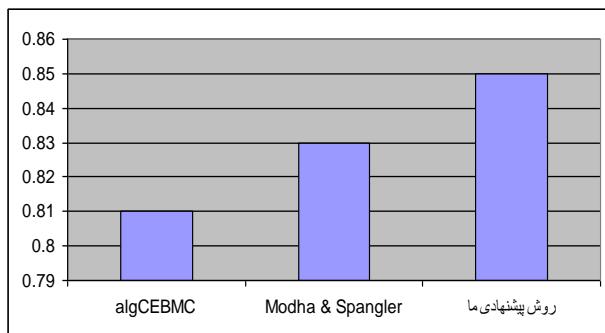
در این قسمت نتایج اجرای الگوریتم خوش بندی پیشنهادی بر روی مجموعه داده‌های استاندارد آورده شده است. به منظور مقایسه نحوه عملکرد الگوریتم ارائه شده با روش‌های قبلی، از دو مجموعه داده استاندارد استفاده شده است که بسیاری از مقالات معتبر برای آزمایش الگوریتم خود از آن استفاده کرده اند و لذا امکان مقایسه نتایج وجود خواهد داشت. این داده‌ها از مخزن داده

¹ <http://archive.ics.uci.edu>

² Heart Disease Data

جدول ۳. مقایسه نتایج روش‌های مختلف بر روی مجموعه داده کارتهای اعتباری

نام الگوریتم	تعداد داده‌هایی که در خوشبندی انتظار قرار گرفته‌اند	دقت
روش پیشنهادی	587	0.85
روش مودها و اسپنگلر	572	0.83
algCEBMC	559	0.81



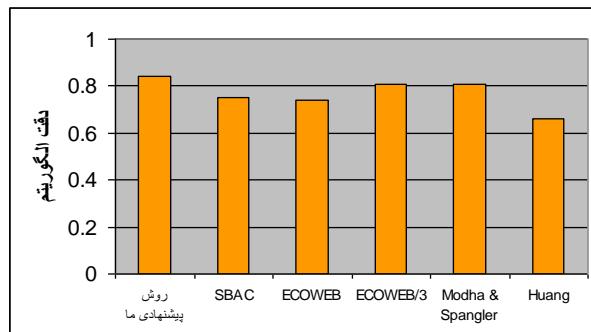
شکل ۲. مقایسه میزان دقت روش خوشبندی پیشنهادی نسبت به دیگر روش‌ها بر روی داده‌های موسسه اعتباری

۴-۳. نتایج الگوریتم پیشنهادی بر روی داده‌های شبیه-سازی شده

در این بخش جهت سنجش نتایج الگوریتم پیشنهادی، مجموعه داده‌هایی با ابعاد ۲۰۰۰۰، ۲۰۰۰ و ۱۲۰۰۰۰ داده شبیه سازی شده و فرآیند خوشبندی در مورد آنها انجام گرفت. همچنین با تغییر نحوه تعریف فاصله بین مقادیر دسته‌ای در روش خوشبندی هوآنگ یعنی روش k -prototypes شکل بهبود یافته این الگوریتم نیز برای خوشبندی این مجموعه از داده‌ها مورد استفاده قرار گرفته و نتایج حاصله استخراج گردید. به منظور اینکه مقایسه بین الگوریتم پیشنهادی و روش k -prototypes بهبود یافته با دقت بیشتری صورت گیرد، از دو شاخص اعتبار سنگی یعنی شاخص مجموع مربعات خطاهای (SSE) و شاخص DB استفاده گردید و مقادیر هر دوی این شاخصها برای جواب حاصل از هر یک از دو روش خوشبندی محاسبه شد. شکلهای ۳ تا ۸ نسبت مقدار شاخص DB و شاخص SSE را برای روش پیشنهادی به مقدار این دو شاخص برای روش k -prototypes قابل مشاهده است، اینست که علیرغم اینکه در روش پیشنهادی از شاخص DB به عنوانتابع برآزندگی استفاده شده و الگوریتم پیشنهادی به دنبال کمینه کردن این تابع هزینه است،

جدول ۲. مقایسه نتایج روش‌های مختلف بر روی مجموعه داده بیماران قلبی

نام الگوریتم	تعداد داده‌هایی که در خوشبندی انتظار قرار گرفته‌اند	دقت
روش پیشنهادی	255	0.84
SBAC	228	0.75
ECOWEB	224	0.74
COBWEB/3	245	0.81
روش مودها و اسپنگلر	244	0.81
روش هوآنگ	200	0.66



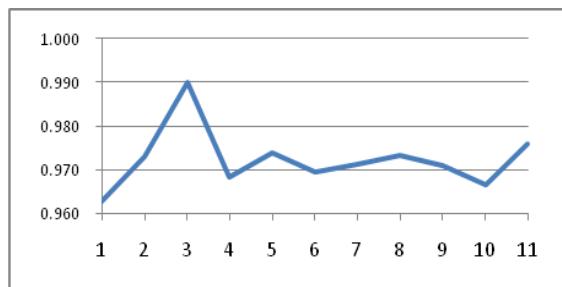
شکل ۱. مقایسه میزان دقت روش خوشبندی پیشنهادی نسبت به دیگر روش‌ها بر روی داده‌های بیماران قلبی

۴-۲-۲. داده‌های کارتهای اعتباری^۱

این مجموعه داده، اطلاعات مربوط به یک موسسه مالی و اعتباری در استرالیا را شامل می‌شود. این داده‌ها یک مجموعه داده مختلط هستند که دارای هشت مشخصه دسته‌ای و شش مشخصه عددی می‌باشند. این مجموعه داده دارای ۶۹۰ نمونه است که به دو کلاس تقسیم می‌شوند: کلاس منفی شامل ۳۸۳ نمونه و کلاس مثبت شامل ۳۰۷ نمونه.

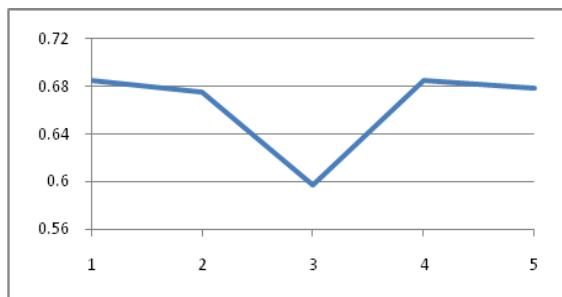
نتایج حاصل از خوشبندی این مجموعه داده نیز توسط روش پیشنهادی و دو روش خوشبندی داده‌های مختلط که برای آزمایش نتایج خود از مجموعه داده‌های کارتهای اعتباری استفاده کرده‌اند، در جدول ۳ آورده شده است. مقادیر این جدول نیز بار دیگر برتری روش پیشنهادی را نسبت به روش ارائه شده توسط مودها و اسپنگلر [۲۶] تأیید می‌کند.

^۱ Australian Credit Approval

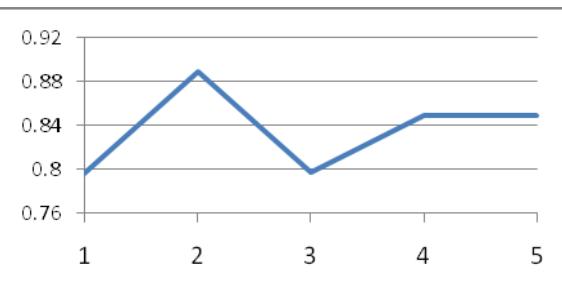


شکل ۶. نسبت مقدار شاخص SSE برای روش پیشنهادی به مقدار این شاخص برای روش k-prototypes بهبود یافته؛ مجموعه داده های ۲۰۰۰۰ تایی

شکلهای ۵ و ۶ بیانگر برتری نتایج حاصل از روش پیشنهادی ۲۰۰۰۰ k-prototypes در مورد مجموعه داده های ۲۰۰۰۰ تایی است؛ به نحویکه مقدار شاخص DB برای روش پیشنهادی بین ۰.۹٪ تا ۱.۸٪ و مقدار شاخص SSE برای روش پیشنهادی بین ۱٪ تا ۴٪ نسبت به شاخصهای متناظر برای روش k-prototypes کمتر است.



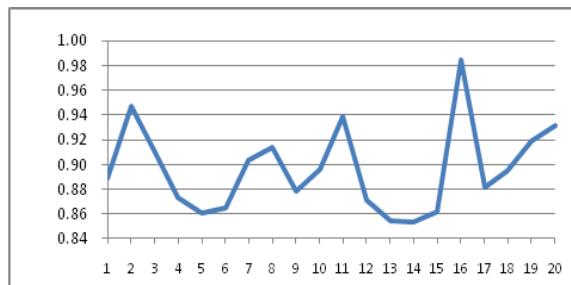
شکل ۷. نسبت مقدار شاخص DB برای روش پیشنهادی به مقدار این شاخص برای روش k-prototypes بهبود یافته؛ مجموعه داده های ۱۲۰۰۰ تایی



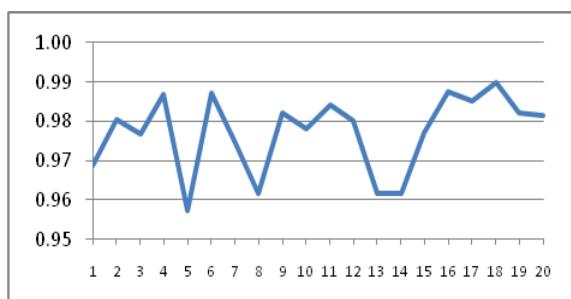
شکل ۸: نسبت مقدار شاخص SSE برای روش پیشنهادی به مقدار این شاخص برای روش k-prototypes بهبود یافته؛ مجموعه داده های ۱۲۰۰۰ تایی

شکلهای ۷ و ۸ نیز بار دیگر برتری روش پیشنهادی را در مورد هر دو معیار DB و SSE تایید می‌نماید. به نحویکه مقدار شاخص

اما به ازای تمامی مجموعه داده های شبیه سازی شده، مقدار شاخص SSE نیز برای الگوریتم پیشنهادی مقدار کمتری را به خود اختصاص داده است.

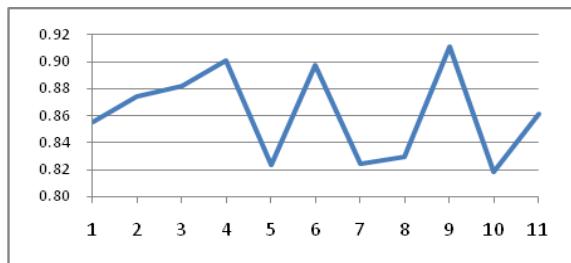


شکل ۳. نسبت مقدار شاخص DB برای روش پیشنهادی به مقدار این شاخص برای روش k-prototypes بهبود یافته؛ مجموعه داده های ۲۰۰۰ تایی



شکل ۴. نسبت مقدار شاخص SSE برای روش پیشنهادی به مقدار این شاخص برای روش k-prototypes بهبود یافته؛ مجموعه داده های ۲۰۰۰ تایی

شکلهای ۳ و ۴ نشان می‌دهد که در مورد مجموعه داده های ۲۰۰۰ تایی، مقدار شاخص DB برای روش پیشنهادی بین ۰.۹۵٪ تا ۰.۹۹٪ نسبت به مقدار این شاخص برای روش k-prototypes کمتر است. شاخص SSE روش پیشنهادی نیز بین ۰.۹٪ تا ۱٪ نسبت به شاخص SSE روش k-prototypes کمتر است.



شکل ۵. نسبت مقدار شاخص DB برای روش پیشنهادی به مقدار این شاخص برای روش k-prototypes بهبود یافته؛ مجموعه داده های ۲۰۰۰ تایی

- [6] Glover, F., Laguna, M., *Tabu search*, Kluwer Academic Publishers, Boston, 1997.
- [7] Krishna, K., Murty, M.N., "Genetic K-Means Algorithm", IEEE Transaction On Systems, Man, And cybernetics—Part B: CYBERNETICS, Vol. 29, No. 3, 1999, pp. 433-439.
- [8] Lu, Y., Lu, S., Fotouhi, F., "FGKA: A Fast Genetic K-means Clustering Algorithm", SAC'04 Nicosia, Cyprus., ACM, 2004.
- [9] Lu, Y., Lu, S., Fotouhi, F., Deng, Y., Susan, D., Brown, J., "an Incremental Genetic K-Means Algorithm and its Application in Gene Expression Data Analysis", BMC Bioinformatics, Vol. 5, 2004, pp. 172-181.
- [10] Maulik, U., Bandyopadhyay, S., "Genetic Algorithm-Based Clustering Technique", Pattern Recognition, Vol. 33, No. 9, 2000, pp. 1455-1465.
- [11] Lin, H.J., Yang, F.W., Kao, Y.T., "An Efficient GA Based Clustering Technique", Tamkang Journal of Science and Engineering, Vol. 8, No. 2, 2005, pp. 113- 122.
- [12] Bandyopadhyay, S., Maulik, U., "Genetic Clustering for Automatic Evolution of Clusters and Application to Image Classification", Pattern Recognition, Vol. 35, No. 6, 2002, pp. 1197- 1208.
- [13] Liu, Y., Kefe, S., Liz, X., "A Hybrid Genetic Based Clustering Algorithm", Proceedings of the Third International Conference on Machine Learning and Cybernetics, Shanghai, 2004.
- [14] Chiang, S., Chu, S.C., Hsin, Y.C., Wang, M.H., "Genetic Distance Measure for K-Modes Algorithm", International Journal of Innovative Computing, Information and Control, Vol. 2, No. 1, 2006, pp. 33- 40.
- [15] Wu, F.X., Kusalik, A.J., Zhang, W.J., "Genetic Weighted K-means for Large-Scale Clustering Problems", University of Saskatchewan, CANADA, 2003.
- [16] Pan, H., Zhu, J., "Genetic Algorithms Applied to Multi-Class Clustering for Gene Expression Data", Genomics Proteomics Bioinformatics, Vol. 1, No. 4, 2003, pp. 279-287.
- [17] Katari, V., Satapathy, S.C., Murthy, J., Reddy, P., "Hybridized Improved Genetic Algorithm with

DB برای روش پیشنهادی بین ۳۱٪ تا ۴۰٪ و مقدار شاخص SSE برای روش پیشنهادی بین ۱۱٪ تا ۲۰٪ نسبت به شاخصهای متناظر برای روش k-prototypes کمتر است.

۵ نتیجه‌گیری

در این مقاله یک روش خوشبندی داده‌های مختلط مبتنی بر الگوریتم ژنتیک ارائه شده است. در روش پیشنهادی، برخلاف اکثر روش‌های خوشبندی داده‌های مختلط که از معیار فاصله صفر و یک برای اندازه‌گیری فاصله بین داده‌های دسته‌ای بهره می‌برند، از تعریف دقیق‌تری جهت سنجش فاصله بین داده‌های دسته‌ای و نیز محاسبه مراکز خوشبندی استفاده شده است و سپس اجزاء الگوریتم ژنتیک (عملگرهای بازترکیب و جهش) مناسب با ساختار جدید نمایش مراکز خوشبندی برای هر یک از داده‌های عددی و دسته‌ای تعریف گردیده است.

از دیگر مزایای روش پیشنهادی اینست که نیازی به تعیین تعداد خوشبندی به عنوان ورودی الگوریتم نداشته و قادر است با بهره گیری از قابلیت جستجوی الگوریتم ژنتیک در فضای جواب، ضمن خوشبندی داده‌ها، مقدار بهینه تعداد خوشبندی را نیز محاسبه نماید که این ویژگی در مورد بسیاری از مسائل دنیای واقعی که در آنها با مجموعه داده‌های بسیار بزرگ با تعداد خوشبندی نامعین سر و کار داریم، بسیار حائز اهمیت است. آزمایش الگوریتم پیشنهادی توسط داده‌های استاندارد و نیز داده‌های شبیه‌سازی شده، نشان از برتری این روش و دقت بالاتر آن نسبت به سایر روش‌های خوشبندی داده‌های مختلط دارد.

مراجع

- [1] Han, J., Kamber, M., *Data Mining Concepts And Techniques*, Elsevier, 2006.
- [2] Witten, I.H., Frank, E., *Data Mining-Practical Machine Learning Tools And Techniques*, Elsevier, 2005.
- [3] Huang, Z., "Clustering Large Data Sets with Mixed Numeric & Categorical Values", in: Proceedings of the First Pacific-Asia Conference on Knowledge Discovery and Data Mining, World Scientific, Singapore, 1997.
- [4] Goldberg, D.E., *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley, 1989.
- [5] Dorigo, M., Stützle, T., *Ant Colony Optimization*, Cambridge: MIT Press, 2004.

Variable Length Chromosome for Image Clustering", International Journal of Computer Science and Network Security, Vol. 7, No. 11, 2007, pp. 121-131.

- [18] Ahmad, A., Dey, L., "A k -Mean Clustering Algorithm for Mixed Numeric and Categorical Data", Data & Knowledge Engineering, Vol. 63, No. 2, 2007, pp. 503- 527.
- [19] Li, C., Biswas, G., "Unsupervised Learning with Mixed Numeric and Nominal Data", IEEE Transactions on Knowledge and Data Engineering Vol. 14, No. 4, 2002, pp. 673- 690.
- [20] Goodall, D.W., "A New Similarity Index Based on Probability", Biometric 22, 1966, pp. 882- 907.
- [21] Huang, J.Z., Ng, M.K., Rong, H., Li, Z., "Automated Variable Weighting in k -Mean Type Clustering", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 27, No. 5, 2005, pp. 657- 668.
- [22] Luo, H., Kong, F., Li, Y., "Clustering Mixed Data Based on Evidence Accumulation", Lecture Notes on Artificial Intelligence 4093, 2006.
- [23] He, Z., Xu, X., Deng, S., "Scalable Algorithms for Clustering Large Datasets with Mixed Type Attributes", International Journal of Intelligent Systems, Vol 20, No. 10, 2005, pp. 1077- 1089.
- [24] Reich, Y., Fenves, S.J., "The Formation and use of Abstract Concepts in Design", Morgan Kaufmann Series In Machine Learning, 1991, pp. 323- 353.
- [25] McKusick, K., Thompson, K., "COBWEB/3: A Portable Implementation", Technical Report FIA-90-6-18-2, NASA Ames Research Center, 1990.
- [26] Modha, D.S., Spangler, W.S., "Feature Weighting in k -Mean Clustering", Machine Learning, Vol. 52, No. 3, pp. 217- 237.
- [27] Zengyou, H., Xiaofe, X., Shengchun, D., "Clustering Mixed Numeric and Categorical Data: A Cluster Ensemble Approach", Department Of Computer Science And Engineering, Harbin Institute Of Technology, Harbin, China, 2001.

