

ارائه‌ی مدلی برای کمینه‌سازی مجموع متوسط فواصل درون‌خوشه‌یی در طبقه‌بندی مشتریان

ابوالفضل کاظمی^{*} (استادیار)

زهرا ضیایی کوچصفهانی (کارشناس ارشد)

دانشکده‌ی مهندسی صنایع و مکانیک، دانشگاه آزاد اسلامی واحد قزوین

در سال‌های اخیر اهمیت داده‌ها به عنوان منابع دارای پتانسیل اطلاعاتی بسیار بالا بهنحو گستردگی مورد توجه قرار گرفته شده است. داده‌کاوی با استخراج و کشف سریع و دقیق اطلاعات با ارزش و پنهان از پایگاه داده‌ها به منظور تصمیم‌گیری و پشتیبانی تصمیم از جمله اموری است که هر کشور، سازمان و شرکتی به منظور توسعه علمی، فنی و اقتصادی خود به آن نیاز دارد. با توجه به ضرورت استفاده از فنون داده‌کاوی -- خصوصاً خوشه‌بندی -- در این نوشتار یک مدل ریاضی براساس رویکرد خوشه‌بندی ارائه‌یی شود که در خوشه‌بندی مشتریان شرکت صنعتی پارس خزر کاربرد دارد. مسئله‌ی خوشه‌بندی به صورت مدل ریاضی با هدف کمینه‌سازی مجموع متوسط فواصل درون‌خوشه‌یی در طبقه‌بندی مشتریان فرمول‌بندی می‌شود که در تمامی موارد آزمایش شده، با بهبودبخشی شدید در فواصل درون‌خوشه‌یی همراه است. عملکرد این شیوه در یک مسئله‌ی واقعی آزموده شده و تحلیل نتایج حاکی از کارایی محاسبات این شیوه است.

abkaazemi@gmail.com
shziaeem@yahoo.com

واژگان کلیدی: داده‌کاوی، خوشه‌بندی، روش‌های خوشه‌بندی سلسه‌یی و غیرسلسله‌یی، مدیریت ارتباط با مشتری، الگوریتم k-means

۱. مقدمه

طرح‌های خوشه‌بندی مؤثر و پرمعنا. الگوریتم‌هایی دقیق و اکتشافی برای این مدل‌ها مطرح شده است.

تحقیقات درخصوص مباحث ریاضی داده‌کاوی و خوشه‌بندی، در سال ۱۹۷۱ به دو طریق انجام شد: (الف) با کمینه‌سازی مجموع مجدد رهای میان‌خوشه‌یی؛ (ب) با کمینه‌سازی بینگاههای k (k-means) ارائه شد؛ این الگوریتم شیوه‌یی فراینده برای خوشه‌بندی انجامیده است. با توجه به رقابت فشرده در بازار و گزینه‌های متنوع محصولات و خدمات که پیش روی مشتریان قرار دارد، تحلیل رفتار مشتریان عامل بسیار مهمی برای بقای بینگاههای تلقی می‌شود. با توجه به موارد یادشده، یک از مؤثث‌ترین ابزارهای بررسی رفتار مشتریان استفاده از تکنیک‌های داده‌کاوی است.

در سال‌های اخیر، شرکت‌ها بر شناخت نیازها و توقعات مشتریان و نیز برگروه‌بندی مشتریان فعلی و احتمالی -- با هدف بهبودبخشی بازدهی راهبردهای بازاریابی و افزایش سهم بازار خود -- متمرکز شده‌اند. فرآونی مجموعه‌های بزرگ داده و نیاز به استخراج دانش نهفته در آن‌ها، به ایجاد الگوریتم‌هایی برای تشخیص الگوهای ناشناخته موجود در مجموعه داده‌ها انجامیده است. آنالیز خوشه‌بندی یک تکنیک داده‌کاوی است که با هدف شناسایی گروه‌های مشابه از نقطه نظر معیارهای معین تشابه، توسعه یافته است. رویکردهای زیادی نسبت به مسئله‌ی خوشه‌بندی وجود دارد، شامل بهینه‌سازی براساس روش‌های در برگیرنده مدل‌های ریاضی برای تشکیل

در این نوشتار از داده‌های حقیقی حاصل از شرکت صنعتی پارس خزر در آزمایشات محاسباتی استفاده شده است. این شرکت در بهمن ماه ۱۳۴۷ در شهر رشت، با نام شرکت صنعتی پارس توشیبا و به منظور ساخت انواع لوازم خانگی بر قی شروع به فعالیت کرد. در اردیبهشت ماه سال ۱۳۵۲، نوع شرکت به شهامت خاص

امروزه رشد سریع فناوری اطلاعات در بخش‌های مختلف کسب و کار، مرایای فراوانی برای صاحبان بینگاههای اقتصادی به ارungan آورده است؛ حوزه‌ی بازاریابی نیز به عنوان زیربخش اصلی سازمان از این قاعده مستثنی نیست. فناوری اطلاعات به بروز تغییراتی در روش‌های بازاریابی سازمان‌ها و ایجاد حجم رسیعی از داده‌های مشتریان انجامیده است. با رقابت فشرده در بازار و گزینه‌های متنوع محصولات و خدمات که پیش روی مشتریان قرار دارد، تحلیل رفتار مشتریان عامل بسیار مهمی برای بقای بینگاههای تلقی می‌شود. با توجه به موارد یادشده، یک از مؤثث‌ترین ابزارهای بررسی رفتار مشتریان استفاده از تکنیک‌های داده‌کاوی است.

* نویسنده مسئول
تاریخ: دریافت ۱۳۸۹/۰۶/۱۵، اصلاحیه ۱۳۹۰/۰۱/۱۱، پذیرش ۱۳۹۰/۰۱/۱۱.

محاسبه‌ی پیچیدگی بین الگوریتم‌های خوش‌بندی k-means و k-medoids برای توزیع نرمال و یکنواخت داده‌های نقطه‌ی انجام شود.^[۷]

روش‌های خوش‌بندی سلسله‌ی (HC)^۳، به صورت یک سری عملیات تقسیمی پیش می‌رود که شروع شان با یک خوشی واحد در برگزینه‌ی تمامی موارد (مثال‌ها) و پایان‌شان با رسیدن به معیار از پیش تعیین شده خالقه‌ی بخشی مشخص می‌شود. بنابراین، روشن‌های HC برای ایجاد تقسیمی بر دسته داده کاربرد دارند، و حل HC را می‌توان به عنوان یک ورودی برای روشن خوش‌بندی غیرسلسله‌ی (NHC)^۴ و به منظور بهبود بخشی حل خوش‌ی حاصله به کار برد.^[۸]

روشن خوش‌بندی غیرسلسله‌ی (NHC)، موسوم به خوش‌بندی تقسیمی (تفکیکی)، به حالتی اطلاق می‌شوند که در آن تعداد تقسیمات معلوم است. در این روشن معمولاً داده ابتدا به k خوش‌تقسیم می‌شود، و پس از آن برای تمامی انتقالات ممکن نقاط داده‌ی مابین خوش‌های تشکیل شده پیش می‌رود تا این که معیار توقف تحقق یابد. بنا بر اظهار محققین، در این شیوه هر خوش‌های توسط مرکز خوش (k) یا یک مورد واقع در مرکز خوش (میانه k) نمایش داد. شاید پرمصرف‌ترین الگوریتم خوش‌بندی الگوریتم k-means باشد که در عمل خوب کار می‌کند، هرچند پیچیدگی زمان در بدترین حالت به طور تصاعدی بالا می‌رود.^[۹] در بررسی‌های اخیر یک الگوریتم جدید خوش‌بندی سلسله‌ی برای داده‌های فازی (FHCA) مورد بررسی قرار گرفته که در آن، خوش‌بندی فازی اعداد با استفاده از روشن‌های سلسله‌ی مرتبی انجام می‌گیرد. در همه‌ی روشن‌های مشابه قبلی، روشن (FCM)^۵ بسط یافته (خوش‌بندی فازی)، داده‌های فازی را پوشش می‌دهد. وجه تمایز این روشن، کاربرد الگوریتم خوش‌بندی سلسله‌ی بر در خوش‌بندی داده‌های فازی است. مهم‌ترین ویژگی این الگوریتم مقاومت بالای آن در مقابل داده‌های مغفوش است.^[۱۰]

در بررسی روشن‌های عدد صحیح با هدف کمیه‌سازی مجموع فواصل میان هر نقطه‌ی داده و مرکز خوش‌بندی منطبق با آن، یا به عبارتی بررسی اجرای الگوریتم k-means در مدل بهینه‌سازی، محققین روشن‌های خوش‌بندی را با مرکز بر فرمولاسیون‌های ریاضی به سه دسته تقسیم کردند: روشن‌های مبتنی بر مقیاس (متريک)، روشن‌های مبتنی بر مدل، و روشن‌های مبتنی بر تفکیک.^[۱۱]

از سوی دیگر، در برخی از الگوریتم‌های NHC -- مثل روشن اکتشافی جهانی k-means -- خوش‌هایها به صورت صعودی تشکیل می‌شوند. طراحان این الگوریتم بر این باورند که حل دسته‌ی k بهینه را می‌توان از حل بهینه‌ی مسئله دسته‌ی (k-1) با استفاده از یک جست‌وجوی محلی روی تمام مجموعه داده‌ها به دست آورد، طوری که الگوریتم k-means مکرراً انجام شود تا بهترین مرکز دسته‌ی k ممکن‌های حاصل شود. اگرچه ممکن است این الگوریتم به راه حل‌های امیدبخشی بینجامد، متضمن رسیدن به بهترین حل کلی نیست.^[۱۲]

در سری دیگری از تحقیقات انجام شده، معیار کمیه‌سازی مجموع مجزوهرهای میان‌گروه‌ها در نظر گرفته شده است. همچنین با بهره‌گیری و بسط برخی نتایج، یک فرمول برنامه‌نویسی خطی عدد صحیح برای مسئله‌ی کاربردی k-medoids ممکن است این الگوریتم پویا و مؤثر برای موردی عرضه می‌شود که در آن بتوان اجزاء را نقاطی روی خط حقیقی در نظر گرفت؛ این الگوریتم کارآمدتر از الگوریتم‌های فعلی برنامه‌نویسی عدد صحیح است. به علاوه در این الگوریتم، از پیش مشخص شدن تعداد خوش‌بندی‌ها ضرورت ندارد، بلکه راه حلی برای هر تعداد خوش‌بندی ارائه می‌شود. زمانی که اجزاء نقاطی از یک فضای اقلیدسی p بعدی باشند، فرمول برنامه‌نویسی خطی عدد صحیح برای مسئله ارائه می‌شود. این فرمول اگرچه تحت شرایط معینی معتبر است، این مزیت را

تبديل شد و در سال ۱۳۶۱ پس از واگذاری سهام متعلق به توشیبا ای ژاپن، نام آن به شرکت صنعتی پارس خزر تغییر یافت.

این شرکت شبکه‌ی گسترده، مرکب از حدود ۳۵۰ تعمیرگاه مجاز سرویس و خدمات پس از فروش در اقصی نقاط ایران دارد. در پی کسب فرصت‌هایی برای بازاریابی را طرح موقعيت جغرافیایی طبقه‌ی بندی کند. قفسیری از خوش‌بندی با استفاده از الگوریتم خوش‌بندی داده‌های شرکت صنعتی پارس خزر در بخش ۴ ارائه خواهد شد. در این نوشتار، در بخش ۲ باسas روشن‌های خوش‌بندی، تحقیقات انجام شده مرور می‌شود و در بخش ۳ نیز مدل پیشنهادی ارائه می‌شود. سپس با طرح مثالی نمادین در بخش ۴، حل مدل را بر روی آن شرح داده و نتایج را با مدل ارائه شده قابل مقایسه می‌کنیم. نهایتاً در بخش ۵ نتایج مطالعه‌ی حاضر، نیز بحث و بررسی آرا برای تحقیقات آتی در زمینه‌ی بهینه‌سازی و خوش‌بندی ارائه می‌شود.

۲. مرور مقالات و تجارب کسب شده

شاید بتوان گفت نخستین گزارش درمورد داده‌کاوی تحت عنوان «شبیه‌سازی فعالیت داده‌کاوی» در سال ۱۹۸۳ ارائه شده است. هم‌زمان، پژوهش‌گران و متخصصین علوم رایانه، آمار، هوش مصنوعی، یادگیری ماشین و... نیز به پژوهش در این زمینه و زمینه‌های مرتبط با آن پرداختند.^[۱]

در سال ۲۰۰۲ نیز در مقاله‌ی «طبقه‌بندی مشتریان با استفاده از داده‌کاوی» از بعضی روشن‌های داده‌کاوی برای شناسایی مشتریان (مشتریان یک بانک و یک کلوب کتاب) استفاده شد.^[۲] برای توصیف این روشن در مجموعه‌ی بانک، از روشن خوش‌بندی به منظور یافتن ساختارها و الگوهای داخل سیستم استفاده شده است. داده‌کاوی (DM)^۳ بخشی لاینک از مطالعات مدیریت ارتباط با مشتری (CRM)^۴ بوده، با این فرضیه که شرکت‌ها می‌توانند در صورت شناخت خصوصیات و علاقه مشتریان خود، با آن‌ها رابطه‌ی موقوفیت‌آمیز برقرار سازند.^[۲]

برخی از محققین به مرور داده‌کاوی، مصارف آن در صنعت و تکنیک‌های مصرفی تحت این عنوان پرداخته‌اند و رابطه و تعامل میان کاربردهای داده‌کاوی و مدیریت ارتباط با مشتری (CRM) را از جوانب مختلف شرح داده‌اند.^[۲] عده‌ی نیز نقش بهینه‌سازی در داده‌کاوی را مورد بررسی قرار داده، و عنوان کردنده که داده‌کاوی و بهینه‌سازی می‌توانند برای توسعه‌ی روابط توین CRM -- روابطی مثل بیشینه‌سازی میزان طول عمر مشتری، آنالیز مشتری و تعاملات مشتریان -- به هم کم کنند.^[۵] با توجه به اهمیت روشن‌های داده‌کاوی در مدیریت ارتباط با مشتری، مرور ادبیات و دسته‌بندی جامع باعث ایجاد برنامه‌ی جامع و کامل شده است. در این مورد ابتدا مرور ادبیات علمی درخصوص کاربرد روشن‌های داده‌کاوی در CRM مورد بررسی قرار گرفته و سپس یک سری اطلاعات در دوره زمانی ۲۰۰۶ تا ۲۰۰۰ شامل ۲۴ مقاله تهیه شده که طی آن ۹۰٪ نویسنده ارتباط مستقیم بین روشن‌های داده‌کاوی و CRM را شناسایی کرده‌اند.^[۶]

«خوش‌بندی» یکی از مهم‌ترین موارد تحقیق در حوزه‌ی داده‌کاوی، و به معنی ایجاد گروه‌هایی است که اجزاء داخل هر گروه بیشترین تشابه، و اجزاء در گروه‌های مجزا همچگونه شبههایی نداشته باشند. خوش‌بندی یک روشن یادگیری بدون نظرت است که مهم‌ترین مزیت آن -- ناشی از ساختار و ویژگی‌هاییش -- امکان گروه‌بندی مستقیم داده‌های بزرگ با کمترین یا بدون اطلاع قبلی است. الگوریتم خوش‌بندی در حوزه‌های زیادی کاربرد دارد و همین نکته باعث شده تا تحقیقاتی درمورد

مراحل انجام کار در این نوشتار بدین ترتیب است که ابتدا فرمول پیشنهادی ارائه شده و سپس برای صحنه‌گذاری، مراحل مقایسه‌ی خروجی آن با خروجی فرمول مشابه استقاده شده در تحلیل خوش‌بندی در نرم‌افزار (R2007b) ۷.۰.۵ MATLAB انجام می‌شود. نتایج این مقایسه در جداول ۱ تا ۴ و نمودارهای ۱ تا ۴

دارد که علاوه بر کاستن از تعداد قبود در حدی قبل توجه، مسئله را به تفکیک‌بندی مجموعه با تعدادی تقسیمات یا مجموعه‌های از پیش مشخص شده کاهش می‌دهد. پژوهش‌گرانی که این مسئله را بررسی کرده‌اند، اظهار داشته‌اند که حل مسئله‌های این چنینی نیازمند زمان محاسباتی مطلوبی است.^[۱۲] در بخش بعد مدل ریاضی پیشنهادی برای مسئله‌ی خوش‌بندی ارائه خواهد شد.

جدول ۱. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $M=2$.

کمترین فاصله		تعداد نمونه
ب	الف	
۰,۵	۲	۳
۰,۶۶	۲	۴
۱	۴	۵
۱,۸۶	۱۱,۱۴	۶
۱,۸۷	۱۱	۷
۲,۷۷	۲۱,۷۷	۸
۴,۴۵	۴۵,۰۹	۹
۴,۱۰	۴۲	۱۰

جدول ۲. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $M=3$.

کمترین فاصله		تعداد نمونه
ب	الف	
۰,۳۷	۱,۶۲	۴
۰,۶۶	۴	۵
۱,۰۴	۹,۱۱	۶
۱,۱۰	۹,۷۰	۷
۱,۶۱	۱۸,۳۳	۸
۲,۹۲	۴۵	۹
۲,۷۳	۴۲	۱۰

جدول ۳. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $M=4$.

کمترین فاصله		تعداد نمونه
ب	الف	
۰,۴۵	۳,۸۰	۵
۰,۷۸	۹,۱۱	۶
۰,۸۳	۹,۷۰	۷
۱,۰۷	۱۶	۸
۱,۹۵	۳۹,۱۶	۹
۲,۰۵	۴۲	۱۰

جدول ۴. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $M=5$.

کمترین فاصله		تعداد نمونه
ب	الف	
۰,۶۱	۹,۱۱	۶
۰,۶۶	۹,۷۰	۷
۰,۸۶	۱۶	۸
۱,۰۶	۳۸,۷۲	۹
۱,۶۴	۴۲	۱۰

۳. مدل ریاضی پیشنهادی برای مسئله‌ی خوش‌بندی

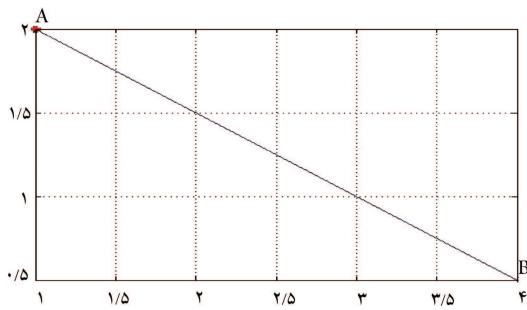
تحلیل‌های خوش‌بندی زیادی (نظیر کمینه‌سازی مجموع مجذورهای میان‌خوش‌بندی و بررسی اجزای متعلق به فضای اقلیدسی P بعدی) توسط رائو انجام شده است. اولین مورد در نظر گرفته شده در بررسی اجزای متعلق به فضای اقلیدسی P بعدی این است که وقتی اجزاء به یک فضای اقلیدسی P بعدی تعلق داشته باشند مسئله پیچیده‌تر می‌شود.

معیار مصرفی در این مدل، کمینه‌سازی مجموع متوسط مجذور فواصل میان‌خوش‌بندی است. به طور عام، فرمول یک مسئله‌ی برنامه‌نویسی $\min_{x \in P}$ و $\max_{x \in P}$ غیرخطی کسری است که حل آن در حالت عمومی خیلی سخت است.

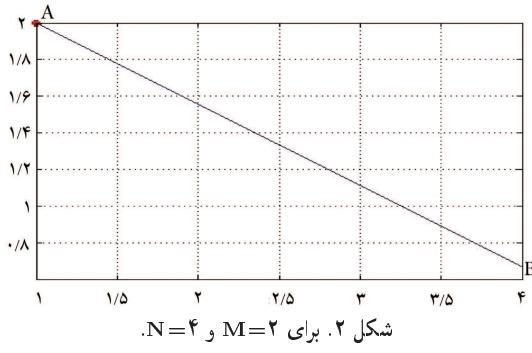
در مدل پیشنهادی دیگر توسط ادورزو کوالی، تعداد خوش‌ها برابر ۲ است، و برای تحلیل خوش‌بندی آن اجزاء به دو خوش با بیشترین فشردگی تقسیم می‌شود، و این فرایند به صورت متوالی تکرار می‌شود. این امر حاکی از آن است که در مرحله از فرایند $M=2$ است. مسئله‌ی برنامه‌نویسی $\min_{x \in P}$ و $\max_{x \in P}$ غیرخطی کسری فاقد قیدهای ضمیمی است، و یک مسئله‌ی برنامه‌نویسی boolean است.

مدل‌های دیگر، کمینه‌سازی کل فاصله‌ی میان‌خوش‌بندی و کمینه‌سازی بیشترین فاصله‌ی میان‌خوش‌بندی است که هر دو مدل مسئله‌ی برنامه‌نویسی خطی عدد صحیح است، اما تعداد قبود سریعاً برحسب N و M افزایش می‌یابد و لذا این فرمول از حيث محاسباتی فقط برای مقادیر کوچک N و M مفید است. باید خاطرنشان ساخت که اگر این مسئله را به صورت یک مسئله‌ی برنامه‌نویسی خطی حل کنیم، بدون آن که متغیرها به اعداد صحیحی محدود شوند، آنگاه یک راه حل $x_{ik} = 1/M$ برای k است. لذا هیچ وقت آنقدر خوش‌شانس نیستیم که به یک راه حل عدد صحیح مثل حل برنامه‌نویسی خطی، دست یابیم. البته وقتی تعداد خوش‌ها M برابر ۲ باشد، این مشکل را می‌توان به صورت مؤثی حل کرد.

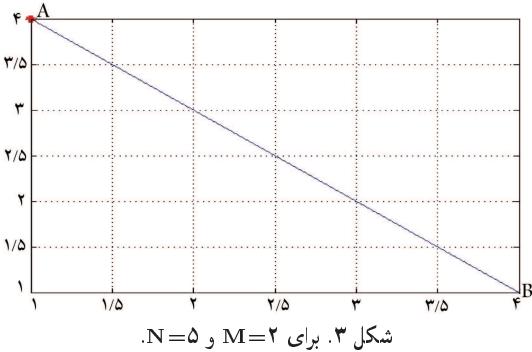
در تمامی تحقیقات گذشته فرمول‌های مختلفی درخصوص فواصل درون‌خوش‌بندی ارائه شده است. با توجه به توضیحات ارائه شده، در این مقاله مدل پیشنهادی «کمینه‌سازی مجموع متوسط فواصل درون‌خوش‌بندی در طبقه‌بندی مشتریان» با مدل مشابه «کمینه‌سازی مجموع متوسط مجذور فواصل میان‌خوش‌بندی» ارائه شده توسعه را تو مورد بررسی قرار گرفته است. در مدل رائو معيار مصرفی، کمینه‌سازی مجموع مجذور فواصل میان‌خوش‌بندی است که حل آن در حالت عمومی خیلی سخت است. البته اگر تعداد اجزاء در هر خوش‌های قبل مشخص باشد می‌توان تابع حقیقی آن را تصحیح کرد، با این حال هنوز یک مسئله‌ی برنامه‌نویسی $\min_{x \in P}$ و $\max_{x \in P}$ غیرخطی با مجموعه شرط‌های جدید است. در این حالت دست‌کم دو روش ممکن برای حل این مسئله وجود دارد؛ روش اول تلقی آن به منزله‌ی یک مسئله‌ی برنامه‌نویسی غیرخطی boolean است، و روش دیگر خطی‌کردن تابع حقیقی است. در خطی‌کردن تابع حقیقی تعداد قبود بیشتر می‌شود، و مسئله‌ی حاصله را باید با یکی از شیوه‌های شناخته شده برای برنامه‌نویسی عدد صحیح خطی حل کنیم.



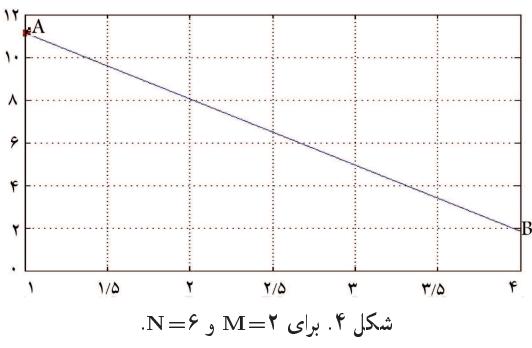
شکل ۱. برای $N=3$ و $M=2$.



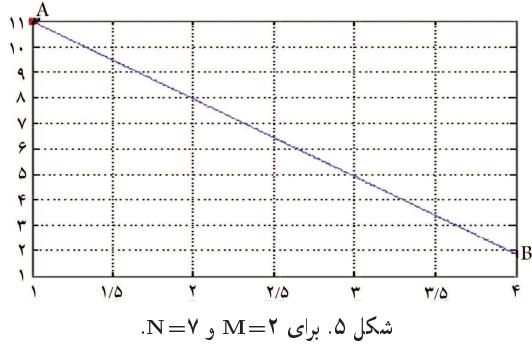
شکل ۲. برای $N=4$ و $M=2$.



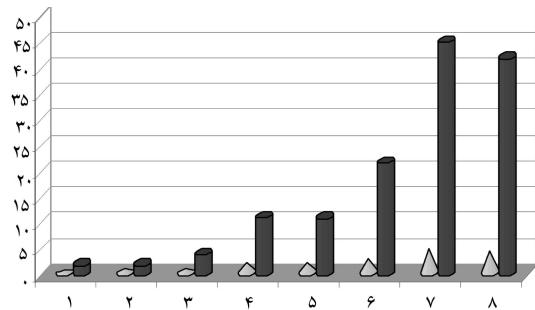
شکل ۳. برای $N=5$ و $M=2$.



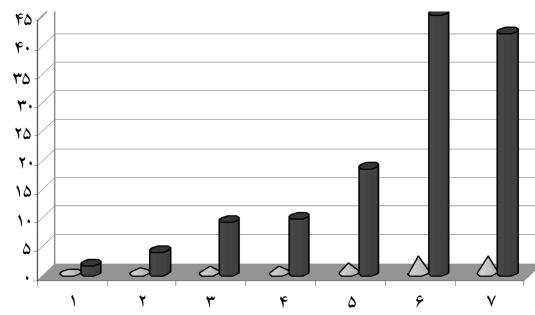
شکل ۴. برای $N=6$ و $M=2$.



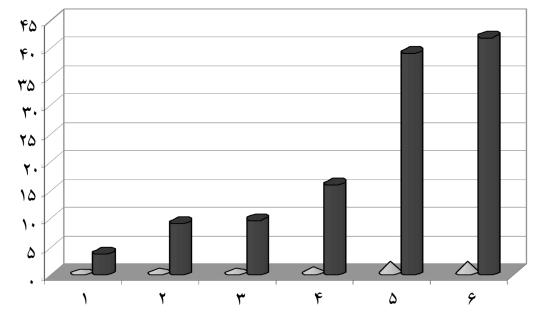
شکل ۵. برای $N=7$ و $M=2$.



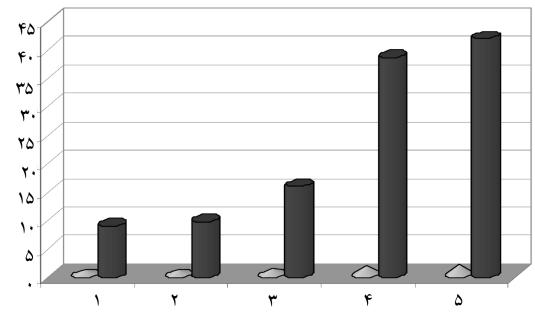
نمودار ۱. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $N=2$ و $M=2$.



نمودار ۲. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $N=3$ و $M=3$.



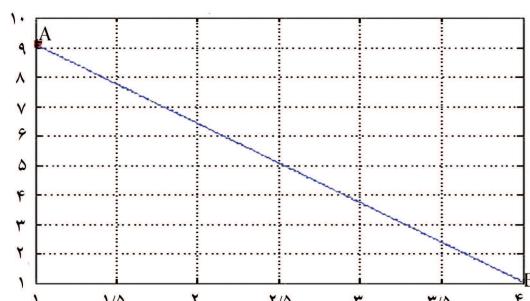
نمودار ۳. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $N=4$ و $M=4$.



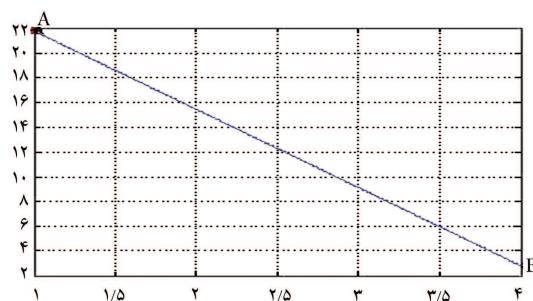
نمودار ۴. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله برای $N=5$ و $M=5$.

و شکل‌های ۱ تا ۲۶ آمده است. در نمودارهای ۱ تا ۴ شکل مثلث نشان‌دهنده نتایج فرمول پیشنهادی، و شکل استوانه نشان‌دهنده نتایج فرمول مشابه را تو است. ضمناً منبع داده‌های آزمون، ماتریس‌های تصادفی از عناصر 0_1 بوده‌اند که سطر و ستون آن براساس N انتخاب شده است.

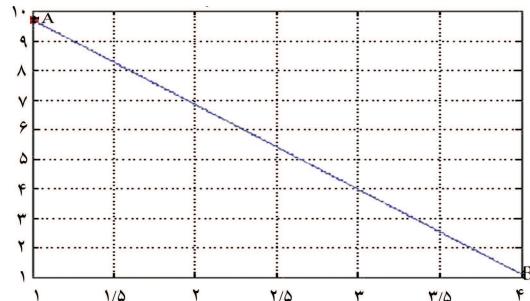
فرمول پیشنهادی که در بقیه‌ی مراحل با علامت اختصاری (ب) نشان داده خواهد شد طبق فرمول ۱ است. در اینجا تحلیل، به «تحلیل خوشبندی مبتنی بر



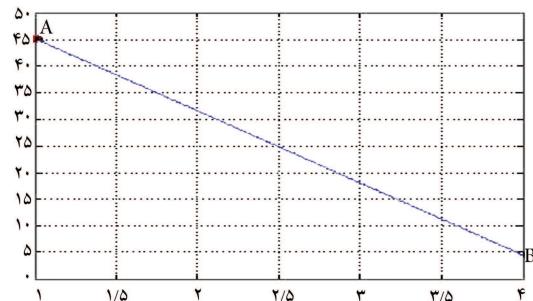
شکل ۱۱. برای $N=6$ و $M=3$



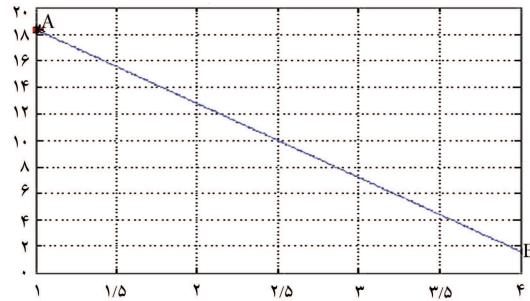
شکل ۱۲. برای $N=8$ و $M=2$



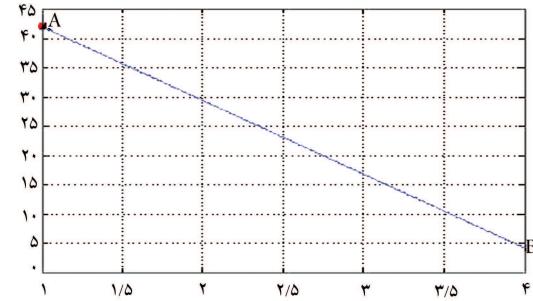
شکل ۱۳. برای $N=7$ و $M=3$



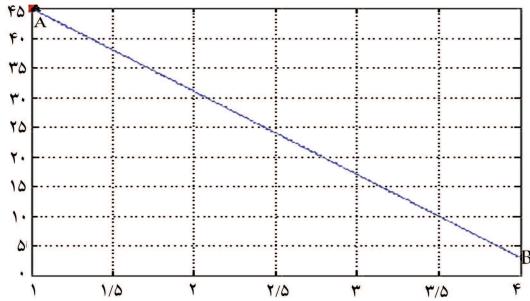
شکل ۱۴. برای $N=9$ و $M=2$



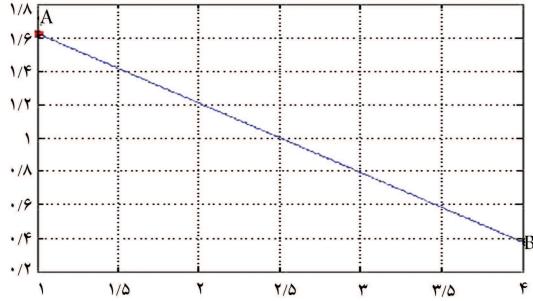
شکل ۱۵. برای $N=8$ و $M=3$



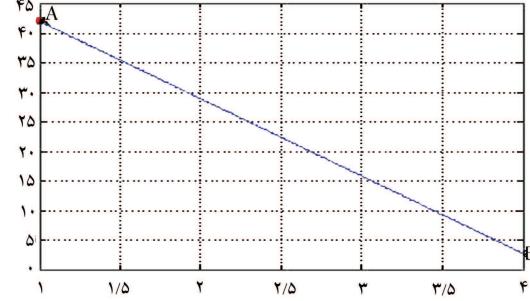
شکل ۱۶. برای $N=10$ و $M=2$



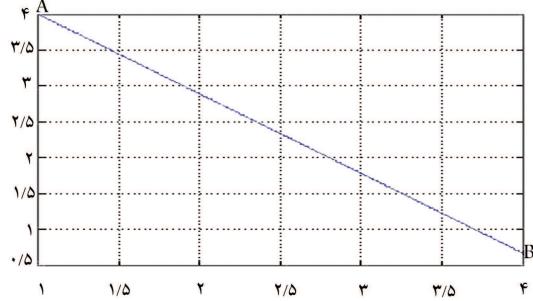
شکل ۱۷. برای $N=9$ و $M=3$



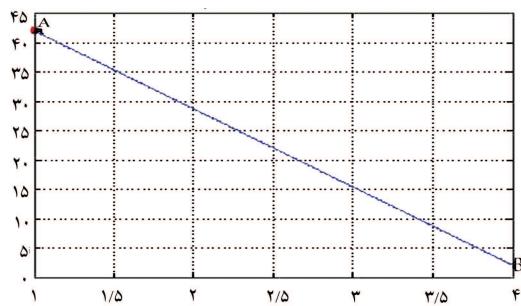
شکل ۱۸. برای $N=4$ و $M=3$



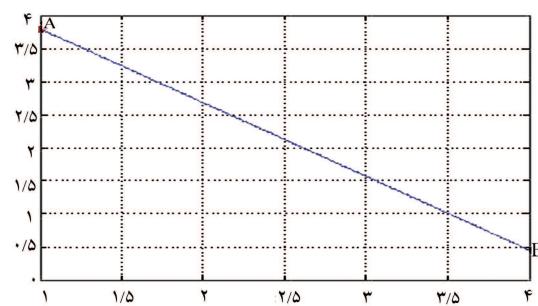
شکل ۱۹. برای $N=10$ و $M=3$



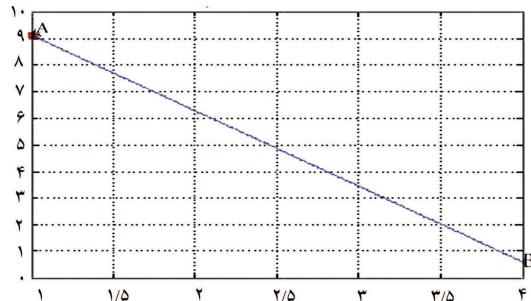
شکل ۲۰. برای $N=5$ و $M=3$



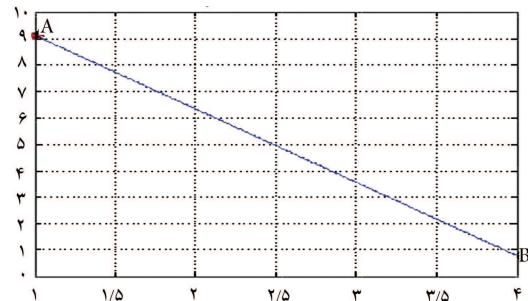
شکل ۲۱. برای $N=10$ و $M=4$



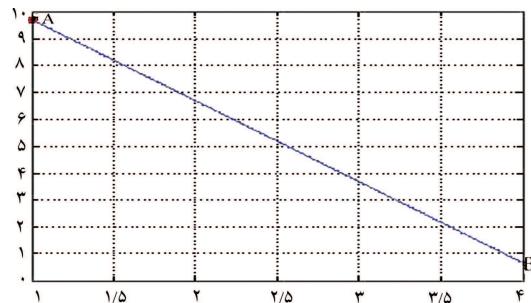
شکل ۱۶. برای $N=5$ و $M=4$



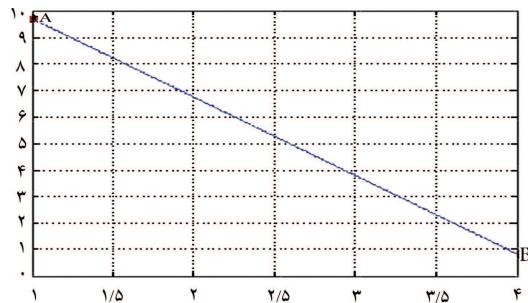
شکل ۲۲. برای $N=6$ و $M=5$



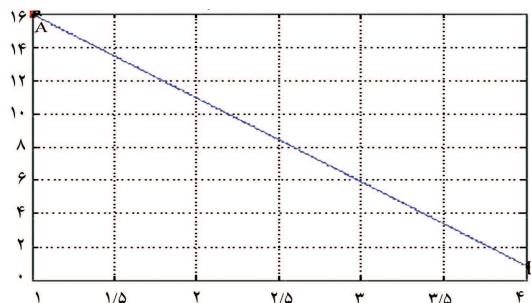
شکل ۱۷. برای $N=6$ و $M=4$



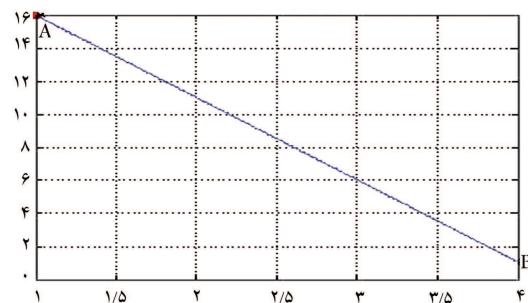
شکل ۲۳. برای $N=7$ و $M=5$



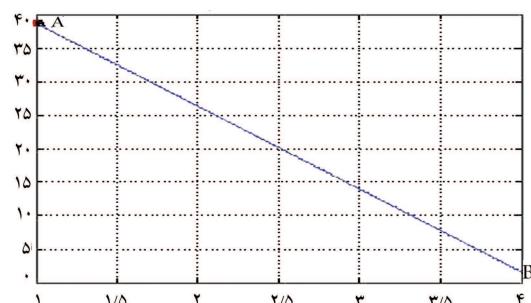
شکل ۱۸. برای $N=7$ و $M=4$



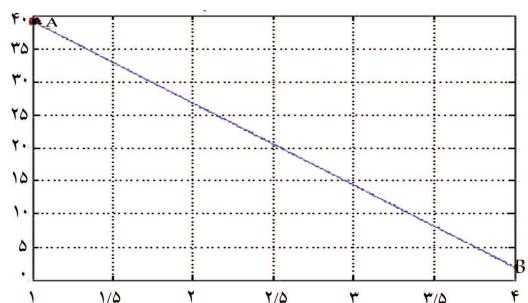
شکل ۲۴. برای $N=8$ و $M=5$



شکل ۱۹. برای $N=8$ و $M=4$



شکل ۲۵. برای $N=9$ و $M=5$



شکل ۲۰. برای $N=9$ و $M=4$

با کم ترین فاصله‌ی درون خوشبی است — نشان داده شده است. ضمناً چنان‌که در جدول ۱ و نمودار ۱ نشان داده شده فواصل درون خوشبی مدل ارائه شده قلی، به مراتب پیشتر از مدل پیشنهادی است.

مرحله‌ی دوم. در جدول ۲ داده‌های ورودی از $N=4$ الی $N=10$ در سه خوشه‌ی متفاوت قرار داده شده است. چنان که در جدول ۲ و نمودار ۲ نشان داده شده باز هم کمترین فاصله‌ی درونخوشه‌ی مدل پیشنهادی به مرتبه بهتر از مدل ارائه شده‌ی قبلی است؛ همچنین در مدل ارائه شده‌ی قبلی کمترین فاصله‌ی درونخوشه‌ی براساس مجموع متوسط مجدد فواصل نیز دارای اعداد بزرگ‌تری است. باز تغییرات در مدل پیشنهادی کمتر از مدل ارائه شده‌ی قبلی است.

در شکل های ۹ تا ۱۵، مدل پیشنهادی با مدل ارائه شده قبایی برای خوشه بمندی داده ها در سه گروه مقایسه شده، که در همه مدل های ارائه شده روند بهبود مدل پیشنهادی بهوضوح دیده شده است.

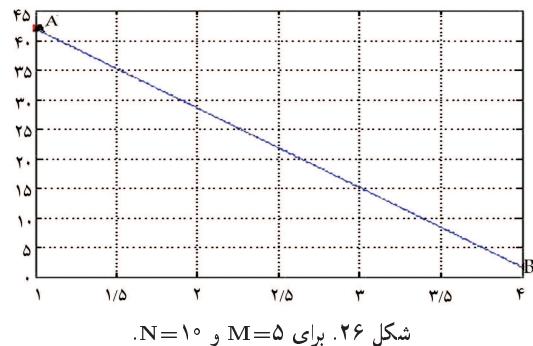
مرحله‌ی سوم. در جدول ۳ داده‌های ورودی از $N=5$ الی $N=10$ در چهار خوشه‌ی متفاوت قرار داده شده است. چنان که در جدول ۳ و نمودار ۳ مشاهده می‌شود کمترین فواصل درون خوشه‌ی بسیار سلسیوس مدل پیشنهادی بهتر از مدل ارائه‌شده‌ی قبلی بوده که روند بهبود به طور کامل در همه‌ی شکل‌های ۱۶ تا ۲۱ تا به صورت روند نزولی بهبود نشان داده شده است. در جدول ۳ نشان داده شده که با افزایش مقدار خوشه‌ها (M) فواصل درون خوشه‌ی تغییرات کمتری دارد؛ این نکته گویای آن است که خوشبندی با مدل پیشنهادی مناسب‌تر و بهتر از خوشبندی از طرق مدل، ارائه‌شده‌ی قبلی، است.

مرحله‌ی چهارم. در جدول ۴ نمودار ۴ داده‌های ورودی از $N=6$ تا $N=10$ در پنج خوشة‌ی متقاوت قرار داده شده است. همانند سه مرحله‌ی قبلی (طبق شکل‌های ۲۲ تا ۲۶) روند بهبود و نیز مقایسه‌ی کمترین فاصله‌ی درون خوشه‌ی بی مدل پیشنهادی با مدل ارائه‌شده‌ی قبلی نشان داده شده، که همچنان شاهد همان روند بهبود هستیم.

۴. ارزیابی مدل پیشنهادی روی یک مجموعه داده حقیقی

در این بخش عملکرد و صحبت مدل ریاضی و مدل خوشبندی پیشنهادی، بر روی یک مجموعه داده حقیقی شرکت صنعتی پارس خزر آزمایش می‌شود. در بازار فرازهای قابلیت^۱ امروز، پارس خزر زمانی قادر به بقایت که نیاز مشتری و تغییرات بازار را به درستی شناسایی کند و منابع در دسترس خود را برای خلق فاکتورهای حیاتی موقوفیت و مزیت های رقابتی خود به کار گیرد. لذا شناسایی ویژگی های مشتریان و انجام خوشبندی مناسب، طوری که مشتریان با پیشترین شbahت در یک خوشبندی قرار بگیرند، از اهمیت به سرزایی برخوردار است. از طرفی با توجه به ۳۵۰ مرکز خدمات پس از فروش شرکت صنعتی پارس خزر برای ارائه خدمات مناسب تر به مشتریان، از مدل پیشنهادی برای خوشبندی مناسب مراکز استفاده شده، و نتیجه آن با مدل ارائه شده قبلي براساس فرمول مشابه «کمینه سازی مجموع متوسط مجدول فواید میان خوشبندی» (رابطه^۲ ۴) مقایسه شده است. مراحل انجام شده عبارت است از:

۱. با توجه به پراکندگی مشتریان شرکت در سراسر ایران، ابتدا نقشه‌ی جغرافیای ایران به چهار ناحیه‌ی «شمال شرقی، شمال غربی، جنوب شرقی و جنوب غربی» تقسیک شد.



شکل ۲۶. برای $M=5$ و $N=10$

فاصله» محدود می شود که در آن معیار سنجش فاصله بین موجودیت های مختلف فراهم است. اگرچه معیار فاصله بین اجزای مختلف شناخته شده، اما هنوز هم معیار تقسیم بهینه به کار بردن نظر بستگی دارد و همیشه، N بیان گر تعداد اجزاء و M بیان گر تعداد خوشه هاست.

$$\min \sum_{k=1}^M \left[\left(\sum_{i=1}^{N-k} \sum_{j=i+1}^N |d_{ij}| x_{ik} x_{jk} \right) / \sum_{i=1}^N x_{ik} \right] / M \quad (1)$$

$$\sum_{k=1}^M x_{ik} = 1$$

$$i = 1, 2, \dots, N \quad (2)$$

این فرمول پیشنهادی با فرمول مشابه «کمینه‌سازی مجموع متوسط مجذور فواصل میان خوشبی» طبق فرمول شماره ۴ — که در بقیه مراحل با علامت اختصاری (الف) نشان داده خواهد شد — مقایسه شده است: در مدل مذکور نیز از نرم‌افزار MATLAB استفاده شده است:

$$\min \sum_{k=1}^M \left[\left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N d_{ij} x_{ik} x_{jk} \right) / \sum_{i=1}^N x_{ik} \right] \quad (4)$$

$$\sum_{k=1}^M x_{ik} = 1$$

$$i = 1, 2, \dots, N \quad (5)$$

$$x_{ik} \geq 0 \quad k = 1, 2, \dots, M, \quad i = 1, 2, \dots, N \quad (8)$$

در این رابطه‌ها d_{ij} فاصله‌ی بین جزء i و j است. برحسب این که آیا جزء k به گروه K نسبت داده شده یا خیر، x_{ik} برابر ۱ یا صفر خواهد بود.

چنان‌که در جداول و نمودارهای ۱ تا ۴ مشاهده می‌شود هدف خوشبندی داده‌های شمشهاری مختلف در گروههای مقاوت از $M=2$ تا $M=5$ خوبه است.

۱.۳. مراحل انجام کار

مراحل انجام کار و نتایج حاصل از آن به شرح زیر است:

مرحله‌ی اول. در جدول ۱ داده‌های ورودی از $N=3$ الی $N=10$ در دو خوشة قرار داده شده است و نشان داده شده که کمترین فاصله‌ی درون خوشه‌ی در مدل پیشنهادی در مقایسه با مدل ارائه شده قبلی (طبق فرمول (۴)) با توجه به افزایش داده‌ها دارای روند بهبود مناسب‌تری بوده و خوشة‌های ایجاد شده از لحاظ کمترین فاصله نیز دارای وضعیت بهتری هستند. همچنین خروجی داده‌های اعمال شده و روند بهبود آن‌ها در شکل‌های ۱ تا ۸ -- که مایع ابعاد خوشه‌های

جدول ۵. مقایسه‌ی نتایج فرمول الف و ب براساس کمترین فاصله.

حداقل فاصله		تعداد خوشه	تعداد نمونه
الف	ب		
۱۰۶/۳۰	۴۶۷۵/۸۵	۳	۳۰
۷۹/۷۲	۴۶۷۵/۸۵	۴	۳۰

۲. سپس تعیین شد که استان مورد نظر در کدام چهار ناحیه‌ی تعیین شده قرار دارد.

۳. براساس چهار ناحیه‌ی تعیین شده، اعداد ۱ تا ۴ به آن استان اختصاص داده شد.

۴. با توجه به محل استقرار استان در ناحیه‌ی تعیین شده، اعداد مختصات جغرافیایی همان ناحیه به استان مربوطه اختصاص یافته است.

با توجه به توضیحات داده شده، استان‌هایی که در چهار ناحیه قرار گرفته‌اند عبارت‌اند از:

۳. فارس (۳۴۴)

۴. کهکیلویه و بویراحمد (۴۴۴)

۵. بوشهر (۵۴۴)

اعداد مختصاتی که به هر استان داده شده براساس ترتیب استقرار در چهار ناحیه‌ی انتخاب شده بوده که بعد از تعیین محل جغرافیایی هریک از استان‌ها، داده‌های تعیین شده وارد نرم‌افزار Matlab ۷.۵.۰ (R2007b) شد، وکلیه‌ی داده‌ها در سه و چهار خوشه قرار گرفت (جدول ۵). چنان‌که در جدول ۵ مشاهده می‌شود، خروجی مدل پیشنهادی برای تعداد نمونه‌های ۳۰ تایی و در خوشه‌های ۳ و ۴ تایی به مرتبه بهتر از مدل ارائه شده قابلی است و روند بهبود آن نیز به‌وضوح مشخص است.

۵. نتیجه‌گیری

در این نوشتار نشان داده شد که برخی مسائل تحلیل خوشه‌بندی مبتنی بر فاصله، مسائل برنامه‌نویسی ریاضی‌اند. مدل ریاضی ارائه شده در این تحقیق مبتنی بر مدل «کمینه‌سازی مجموع متوسط فواصل درون‌خوشه‌یی» است که مزیت اصلی آن مطلوب بودن نتایج محاسباتی است. از مدل ارائه شده به صورت موردی برای داده‌های مشتریان شرکت صنعتی پارس‌خزر در اقصی نقاط ایران استفاده شد، و براین اساس خوشه‌های پیشنهادی تشکیل شد. صحبت مدل پیشنهادی با نتایج حاصل از مدل ارائه شده قابلی «کمینه‌سازی مجموع متوسط فواصل مجازی درون‌خوشه‌یی» بر روی مثالی نمایین و مجموعه داده‌ی پیشنهادی مقتضیه و مشاهده شد که خروجی مدل پیشنهادی نسبت به مدل ارائه شده قابلی بهتر است.

در چند دهه‌ی اخیر، مدل‌های ریاضی کاربرد موفقیت‌آمیزی در تحلیل خوشه‌بندی داشته‌اند. اگرچه نتایج زیادی حاصل شده، اما همچنان برای ادغام کامل تحلیل خوشه‌بندی در مدل‌های ریاضی باید با جدیت کوشید و در این راه به اصول بدیهی، بهویه خوشه‌بندی نیازمندیم. بنابراین موضوعات پیشنهادی برای تحقیقات بعدی، ابداع الگوریتم‌های دقیق جدید عمدتاً برای خوشه‌بندی سلسه‌ی تقسیمی، خوشه‌بندی دنباله‌یی و جمعی است. انجام مقایسه‌ی تجربی روش‌های ابداعی نیز خیلی نادر است، همچنین این روش‌ها هیچ اطلاعاتی درمورد ساختار خوشه‌ها در اختیار ما قرار نمی‌دهند، مثلاً کدام داده‌های خوشه به مرکز یا متوسط آن نزدیک ترند یا داده‌هایی که به نوعی به دو یا چند خوشه شباهت دارند چگونه بیان می‌شوند. برای رفع این مشکلات بهتر است از رویکرد فازی در مدل‌های تحلیل خوشه‌بندی استفاده شود. همچنین پیشنهاد می‌شود مدل پیشنهادی در این نوشتار به صورت تک‌بیانی با روش BRICH -- که الگوریتمی جدید برای خوشه‌بندی مجموعه داده‌های خیلی بزرگ است و مستلزم خوشه‌بندی بزرگ را از طریق مرکز بر بخش‌های متراکم و خلقی خلاصه‌بی فشرده ایجاد می‌کند -- نیز در بررسی‌های آتی استفاده شود.

• ناحیه‌ی شمال شرقی با اعداد مختصات (۱):

۱. خراسان رضوی (۱۱۱)

۲. خراسان جنوبی (۲۱۱)

۳. خراسان شمالی (۳۱۱)

۴. سمنان (۴۱۱)

۵. یزد (۵۱۱)

۶. گلستان (۶۱۱)

• ناحیه‌ی شمال غربی با اعداد مختصات (۲):

۱. آذربایجان شرقی (۱۲۲)

۲. اردبیل (۲۲۲)

۳. آذربایجان غربی (۳۲۲)

۴. گیلان (۴۲۲)

۵. زنجان (۵۲۲)

۶. کردستان (۶۲۲)

۷. قزوین (۷۲۲)

۸. مازندران (۸۲۲)

۹. همدان (۹۲۲)

۱۰. کرمانشاه (۱۰۲۲)

۱۱. تهران (۱۱۲۲)

۱۲. قم (۱۲۲)

۱۳. مرکزی (۱۳۲۲)

۱۴. لرستان (۱۴۲۲)

۱۵. ایلام (۱۵۲۲)

۱۶. اصفهان (۱۶۲۲)

• ناحیه‌ی جنوب شرقی با اعداد مختصات (۳):

۱. سیستان و بلوچستان (۱۳۳)

۲. کرمان (۲۳۳)

۳. هرمزگان (۳۳۳)

• ناحیه‌ی جنوب غربی با اعداد مختصات (۴):

۱. خوزستان (۱۴۴)

۲. چهارمحال و بختیاری (۲۴۴)

پانوشت ها

1. Data Mininy
2. customer relationship management
3. hierarchical clustering
4. nonhierarchical clustering
5. Fuzzy Clusterinity Method
6. hyper competitive

منابع (References)

1. Saeedi, A. "Concept and application of data mining in high education", **18** (In Persain)(2005).
2. Ahola, J. and Rinta-Runsala, E. "Data mining case studies in customer profiling", Research Report TTE1-2001-29, VTT Information Tecknology, pp. 1-24 (2002).
3. Nemati, H.R. and Barko, C.D. "Enhancing enterprise decisions through organizational data mining", *J. of Computer Information Systems*, **42**(4) pp. 21-28 (2002).
4. Rygielski, C., Wang, J.C. and Yen, D. "Data mining techniques for customer relationship management", *Technology in Society*, **24**, pp. 483-502 (2002).
5. Padmanabhan, B. and Tuzhilin, A. "On the use of optimization for data mining: Theoretical interactions and CRM opportunities Management", *Science*, **49**, pp. 1327-1343 (2003).
6. Ngai, E.W.T., Xiu, L. and Chau, D.C.K. "Application of data mining techniques in customer relationship management: A literature review and classification", *Expert Systems with Applications*, **36**, pp. 2592-2602 (2009).
7. Velmurugan, T. and Santhanam, T. "Computational complexity between k-means and k-medoids clustering algorithms for normal and uniform distributions of data points", *J. of Computer Science*, **6**(3), pp. 363-368 (2010).
8. Sharma, S., *Applied Multivariate Techniques*, John Wiley & Sons, USA (1996).
9. Zhang, T., Ramakrishnan, R. and Livny, M. "BIRCH: A new data clustering algorithm and its applications", *Data Mining and Knowledge Discovery*, **1**, pp. 141-182 (1997).
10. GhasemiGol, M., Sadoghi Yazdi, H. and Monsefi, R. "A new hierarchical clustering algorithm on fuzzy data (FHCA)", *Int. J. of Computer and Electrical Eng.*, **2**(1), pp. 1793-8163 (2010).
11. Bradley, P.S., Fayyad, U.M. and Angasarian, O.L. "Mathematical programming for data mining: Formulations and challenges", *J. on Computing*, **11**, pp. 217-238 (1999).
12. Likas, A., Vlassis, N. and Verbeek, J.J. "The global k-means clustering algorithm", *Pattern Recognition*, **36**, pp. 451-461 (2003).
13. Rao, M.R. "Cluster analysis and mathematical programming", *J. of the American Statistical Association*, **66**, pp. 622-626 (2009).