

## به کارگیری الگوریتم k-means وزنی هوشمند با متریک مینکوفسکی به منظور خوشه‌بندی مشتریان

(مورد مطالعاتی مشتریان یک شرکت ارائه دهنده خدمات ارتباطی)

نرگس احمدزاده گلی

شرکت مخابرات ایران، منطقه تهران،

سید علی ملک جعفریان

شرکت مخابرات ایران، منطقه تهران،

حسن متقی محمودی

شرکت مخابرات ایران، منطقه تهران،

امیرحسین شیخ انصاری

شرکت مخابرات ایران، منطقه تهران،

### چکیده

خوشه‌بندی k-means وزنی هوشمند مینکوفسکی تعمیمی از خوشه‌بندی k-میانگین است که در آن تعداد خوشه‌ها و مراکز مربوطه را می‌توان مشخص کرد درحالی‌که تابع زیان مینکوفسکی به‌عنوان معیار عدم تشابه در نظر گرفته می‌شود. بنابراین برخلاف الگوریتم‌های k-میانگین ساده، انتخاب مراکز در هر خوشه تصادفی نمی‌باشد، در نتیجه نیازی به تکرار چندین مرتبه الگوریتم به منظور حذف اثرات تصادفی مراکز اولیه نیست. انتخاب معیار عدم تشابه مینکوفسکی به جای اقلیدسی یا سایر معیارهای عدم تشابه، به پژوهشگر در بهینه‌سازی بهتر مراکز بویژه زمانی که میان داده‌ها، داده پرت وجود داشته باشد، کمک می‌کند و در نتیجه دقت خوشه‌بندی افزایش می‌یابد. در این پژوهش، کارکرد الگوریتم یادشده بر برخی پایگاه داده‌های شبیه‌سازی‌شده که خوشه‌های آنها، از پیش تعیین شده است، بررسی می‌شود. به منظور سنجش دقت الگوریتم در خوشه‌بندی داده‌ها از معیارهای اطلاعات نرمال شده (NVI) و همچنین درصد داده‌هایی که به درستی خوشه‌بندی شده‌اند، استفاده می‌شود. در نهایت عملکرد الگوریتم بر پایگاه داده‌های مشترکین یک شرکت ارائه دهنده خدمات ارتباطی به منظور شناسایی مشتریان خاص و ارائه بسته‌های تشویقی در راستای توسعه سرویس‌ها و همچنین نگهداشت مشتریان ارزشمند سنجیده می‌شود.

واژگان کلیدی: معیار عدم تشابه، k-میانگین هوشمند وزنی، خوشه‌بندی، تابع زیان مینکوفسکی

## ۱- مقدمه

### ۱-۱- الگوریتم k-میانگین

الگوریتم های k-میانگین یکی از محبوبترین و مورد استفادهترین الگوریتمهای خوشه‌بندی می‌باشند و برای اولین بار توسط مک کوئین (Macqueen, ۱۹۶۷) معرفی شد. این الگوریتم ها، از الگوریتمهای سلسله‌مراتبی بسیار کارآمدتر می‌باشند و بدین‌منظور طراحی شده‌اند تا داده‌های عددی را خوشه‌بندی کنند به‌طوری که هر خوشه دارای یک مرکز به نام میانگین است. در این الگوریتم فرض بر آن است که تعداد خوشه‌ها یعنی  $K$  ثابت باشد، همچنین در آن یک تابع خطا وجود دارد. نتایج الگوریتم به معیار فاصله‌ای که در آن به کار برده می‌شود، می‌تواند بستگی داشته‌باشد. الگوریتم مرسوم k-میانگین به‌صورت زیر تعریف می‌شود:

فرض کنید  $X$  یک مجموعه داده با  $n$  مشاهده باشد، به‌عبارتی،  $X = (X_1, \dots, X_n)'$ ، همچنین هر یک از مشاهدات دارای  $m$  ویژگی باشد، یعنی،  $X_i = (X_{i1}, \dots, X_{im})$  k-میانگین الگوریتمی است که پایگاه داده  $X$  را به  $K$  خوشه مجزای  $S = \{S_1, \dots, S_K\}$  از مشاهدات مشابه تقسیم می‌کند به‌طوری که در هر خوشه عدم‌تشابه بین هر مشاهده و مرکز آن (که به طور تصادفی از مجموعه داده‌ها انتخاب می‌شود) مینیمم می‌شود. تابع زیان بین مشاهده  $X_i$  در خوشه  $k$  و مرکز مربوط به آن یعنی  $C_k$  به‌صورت زیر است:

$$J_{k-\text{میانگین}} = \sum_{k=1}^K \sum_{X_i \in S_k} L(X_i, C_k), \quad (1-1)$$

که در آن،  $C_k = (C_{k1}, \dots, C_{km})$  برای  $k = 1, \dots, K$  و  $L$  فاصله اقلیدسی بین  $X_i \in S_k$  و  $C_k$  می‌باشد. هر مشاهده به نزدیکترین مرکز  $C_k$  یعنی میانگین آن تخصیص داده می‌شود و تابع زیان میانگین  $J_{k-\text{میانگین}}$  محاسبه می‌شود. مراکز خوشه‌ها به میانگین  $S_k$  به‌روزرسانی می‌شود تا زمانی که مراکز جدید در مقایسه با مراکز مرحله قبل تغییر نکنند. برای یک مقدار اولیه  $K$ ، با قرار دادن داده‌های به‌جای مانده به نزدیکترین خوشه‌ها و سپس تکرار تغییر اعضای آن خوشه با توجه به تابع زیان، تا زمانی که مقدار تابع زیان دیگر به‌طور معنادار تغییری نکند و اعضای خوشه‌ها ثابت بمانند، ادامه می‌یابد. رابطه (۱-۱) را می‌توان به‌صورت زیر بازنویسی کرد.

$$G_{k-\text{میانگین}}(H, C) = \sum_{k=1}^K \sum_{i=1}^n h_{ik} L(X_i, C_k), \quad (1-2)$$

که در آن  $C = (C_1, \dots, C_K)'$  و  $H$  یک ماتریس  $n \times K$  است به‌طوری که برای هر  $i = 1, \dots, n$  و  $h_{ik} \in \{0, 1\}$  و  $\sum_{k=1}^K h_{ik} = 1$  می‌باشد. اجرای الگوریتم به پیدا کردن تعداد خوشه‌ها و مراکز اولیه بستگی دارد. بنابراین الگوریتم می‌بایست به‌دفعات بسیار تکرار شود تا تأثیر مرکز اولیه کمتر شود. فرآیند تا زمانی ادامه می‌یابد که هیچ مشاهده‌ای خارج از خوشه‌ها باقی نماند.

در یک الگوریتم خوشه‌بندی، ویژگی‌هایی که بار اطلاعات کمی دربر دارند یا از اهمیت برخوردار نیستند، ممکن است باعث ایجاد سردرگمی و طبقه‌بندی اشتباه شوند. زیرا، در مقایسه با سایر ویژگی‌ها ممکن است در حمل اطلاعات، به آنها درجه اهمیت یکسانی اختصاص یابد. چنین

مجموعه داده‌هایی که دارای ویژگی‌های زائد می‌باشند، زیاد مشاهده می‌شوند. در الگوریتم  $k$ -میانگین، می‌توان ویژگی وزن را به‌منظور کاهش تأثیر ویژگی‌های اضافه یا خطا، اضافه نمود. هانگ و همکاران (Haung et. Al., ۲۰۰۸)، الگوریتم خوشه‌بندی  $k$ -میانگین وزنی را که به هر یک از ویژگی‌ها، وزن‌های متفاوتی اختصاص می‌دهد، معرفی نمودند.

یک تابع زیان با توجه به ساختار پایگاه داده انتخاب می‌شود. احمدزاده و همکاران (Ahmadzadehgoli, ۲۰۱۹) الگوریتم خوشه‌بندی  $k$ -میانگین لاینکس بر پایه تابع زیان نامتقارن لاینکس به جای اقلیدسی را ارائه کردند. آنها نشان دادند که زمانی که بیش برآوردی یا کم برآوردی مراکز خوشه‌ها دارای اهمیت باشد، می‌توان به یاری آن مشاهداتی را به خوشه‌ای خاص هدایت کرد. آمورین و کامیزرکزاک (Komisarczuk, Amorim, ۲۰۱۲) نیز در الگوریتم  $k$ -میانگین وزنی به جای تابع زیان اقلیدسی، از تابع مینکوفسکی به‌عنوان معیار عدم‌تشابه استفاده کردند که به آن پرداخته می‌شود.

مشخص نبودن تعداد خوشه‌ها از پیش و همچنین مراکز اولیه، یکی از ضعف‌های بسیاری از الگوریتم‌های خوشه‌بندی است. الگوریتم  $k$ -میانگین هوشمند که توسط میرکین (Mirkin, ۲۰۰۵) معرفی شد، به‌منظور تعیین تعداد خوشه‌ها و مراکز دقیق آنها مفید است. در این روش مشاهدات خوشه‌بندی نشده، در دورترین نقطه از گرانیگاه اولیه قرار می‌گیرد و دورترین نقطه به عنوان یک مرکز آزمایشی در نظر گرفته می‌شود. آنگاه خوشه توسط همه مشاهداتی که به مرکز آزمایشی نسبت به مرکز اولیه نزدیک‌تر است، پر می‌شود. پس از آنکه همه مشاهدات خوشه‌بندی شدند، خوشه‌های کوچک با استفاده از یک مقدار آستانه‌ای از پیش تعیین‌شده، حذف می‌شوند. این روش به علت سادگی، حتی برای افرادی که پیش‌زمینه‌ای از آمار و علوم کامپیوتر ندارند، نیز خوشایند است و نیازی به تکمیل چندین بار الگوریتم برای یافتن بهترین تعداد خوشه را ندارد و درواقع یک الگوریتم قطعی است که تنها به یکبار اجرا نیاز دارد.

## ۲-۱- الگوریتم خوشه‌بندی $k$ -میانگین وزنی مینکوفسکی

هانگ و همکاران (۲۰۰۸)، الگوریتم خوشه‌بندی  $k$ -میانگین وزنی را که به هر یک از ویژگی‌ها، وزن‌های متفاوتی اختصاص می‌دهد، معرفی نمودند. این الگوریتم تابع هدف زیر را مینیمم می‌سازد:

$$G_{\beta, \gamma}(H, C, W) = \sum_{k=1}^K \sum_{i=1}^n \sum_{j=1}^m h_{ik} w_j^{\beta} (X_{ij} - C_{kj})^{\gamma} \quad (1-2)$$

که در آن،  $W_j$  وزن ویژگی‌ها، مقداری نامنفی می‌باشد و  $\sum_{j=1}^m w_j = 1$  همچنین  $\beta$  نشانگر تأثیر وزن و  $C = (C_1, \dots, C_K)'$  و  $H$  یک ماتریس  $n \times K$  به‌طوری‌که برای هر  $i = 1, \dots, n$  و  $h_{ik} \in \{0, 1\}$  و  $\sum_{k=1}^K h_{ik} = 1$  می‌باشد. هدف مینیمم نمودن تابع (۱-۲) بین  $X_i$  و مراکز خوشه مربوطه در هر خوشه با در نظر گرفتن قید  $\sum_{j=1}^m w_j = 1$  می‌باشد. هر یک از وزن‌ها با استفاده از رابطه زیر به‌روزرسانی می‌شود،

$$w_j^* = \begin{cases} \cdot & \text{if } E_j = \cdot \\ \frac{1}{\sum_{u \in M} (\frac{E_j}{E_u})^{\beta-1}} & \text{if } E_j \neq \cdot \end{cases} \quad (1-3)$$

که در آن

$$E_j = \sum_{k=1}^K \sum_{i=1}^n h_{ik} |X_{ij} - C_{kj}|^\beta.$$

الگوریتم k-میانگین وزنی به صورت زیر می باشد:

۱- تعداد خوشه ها یعنی  $K$ ، مشخص می شود. مراکز اولیه خوشه ها و وزن های ویژگی ها، به صورت تصادفی انتخاب می شود به گونه ای که  $\sum_{j=1}^m w_j = 1$ .

وزن های تصادفی ابتدایی می تواند به این صورت در نظر گرفته شوند،  $w_j = \frac{1}{m}$ .

۲- هر شی به نزدیکترین مرکز با حداقل نمودن (۱-۲) اختصاص می یابد.

۳- تمام مراکز به میانگین خوشه مربوطه به روزرسانی می شوند. فرآیند ادامه می یابد و اگر تغییری در مراکز خوشه های مرحله ۲ ایجاد نشود، متوقف و سپس  $S = \{S_1, \dots, S_K\}$  خوشه های نهایی خواهد بود.

۴- با استفاده از رابطه (۱-۳)، وزن ها با در نظر گرفتن شرط  $\sum_{j=1}^m w_j = 1$  به روزرسانی می شوند.

آمورین و کامیزرکزاک (۲۰۱۲) این الگوریتم را به k-میانگین وزنی مینکوفسکی تعمیم دادند که در آن فاصله مینکوفسکی به عنوان معیار عدم تشابه به جای فاصله اقلیدسی به صورت زیر به کار می رود.

$$d_{Min}(X, Y) = \left[ \sum_{j=1}^m (X_j - Y_j)^p \right]^{\frac{1}{p}}, \quad p \geq 1,$$

که در آن  $p$ ، مرتبه فاصله مینکوفسکی می باشد. اگر وال و همکاران (Agrawal, ۱۹۹۳) و یی و فالوتساس (Faloutsos, Yi, ۲۰۰۰) و لی (Li, ۲۰۰۰) برخی از کاربردهای این فاصله را بررسی نمودند. اگر  $p$  برابر با ۱، ۲ و  $\infty$  باشد آنگاه به ترتیب فاصله منهتن، اقلیدسی و فاصله ماکسیمم نتیجه می شود.

حال در این الگوریتم تابع هدف زیر (با توجه به  $\sum_{j=1}^m w_j = 1$ ) مینیمم می شود.

$$J_{\beta, \beta}(H, C, W) = \sum_{k=1}^K \sum_{i=1}^n \sum_{j=1}^m h_{ik} w_j^\beta |X_{ij} - C_{kj}|^\beta \quad (1-4)$$

اگر قرار دهیم

$$E_j = \sum_{k=1}^K \sum_{i=1}^n h_{ik} |X_{ij} - C_{kj}|^\beta.$$

با استفاده از  $E_j$ ، می‌توان (۴-۱) را بصورت زیر بازنویسی نمود،

$$J_{\beta,\beta}(H, C, W) = \sum_{j=1}^m w_j^\beta E_j$$

حال به منظور حداقل نمودن آن باتوجه به قید  $\sum_{j=1}^m w_j = 1$  می‌بایست با مشتق گرفتن از تابع لاگرانژ زیر نسبت به  $w_j$  و قرار دادن آن با صفر

$$w_j = \left(\frac{\lambda}{\beta E_j}\right)^{\frac{1}{\beta-1}}$$

از آنجا که مشتق مرتبه دوم برای  $\beta \geq 1$  نامنفی می باشد لذا نقطه بدست آمده، مینیمم می باشد. حال

با جمع بستن این عبارت بر روی  $j$ ،

$$1 = \sum_{j=1}^m \left(\frac{\lambda}{\beta E_j}\right)^{\frac{1}{\beta-1}}$$

$$w_j^* = \frac{1}{\sum_{u \in M} \left(\frac{E_j}{E_u}\right)^{\frac{1}{\beta-1}}} \quad (1-5)$$

که مقدار بهینه وزن ویژگی‌ها در الگوریتم است.

الگوریتم k-میانگین وزنی مینکوفسکی در برخی مراحل با الگوریتم k-میانگین وزنی تفاوت دارد که بصورت زیر می‌باشند:

- ۱- در مرحله به‌روزرسانی خوشه‌ها، برای یک مجموعه داده شده از مراکز و یک مجموعه از وزن‌ها به ازاء هر ویژگی، خوشه‌ها با استفاده از قانون مینیمم سازی فاصله‌ها، با فاصله مینکوفسکی با توان  $\beta$ ، به‌روزرسانی می‌شوند.
  - ۲- در مرحله به‌روزرسانی مراکز، برای یک مجموعه داده شده از خوشه‌ها و یک مجموعه از وزن‌ها، برای هر ویژگی و خوشه، مرکز هر خوشه به مرکز مینکوفسکی آن باتوجه به مینیمم نمودن معادله (۳-۲) به‌روزرسانی می‌شود.
  - ۳- در مرحله به‌روزرسانی وزن‌ها، برای یک مجموعه داده شده از خوشه‌ها، با توجه به مراکزشان، وزن‌ها با توجه به رابطه (۵-۱) به شرط آنکه جمع آنها برابر با ۱ باشد، به‌روزرسانی می‌شوند.
- برای آنکه الگوریتم بالا کاربردی باشد، به یک فرآیند امکان پذیر محاسباتی به‌منظور بدست آوردن مرکز مینکوفسکی ویژگی‌ها برای هر  $\beta \geq 1$ ، نیاز است. همچنین مقدار مینیمم برای  $\beta = 1$  که معادل با فاصله بلوک شهری است، برابر با میانه  $X_i$  ها می‌باشد. لذا در اینجا جهت مینیمم سازی کافی است  $\beta > 1$  در نظر گرفته شود. مقادیر  $X_i$  به‌صورت صعودی  $X_1 \leq X_2 \leq \dots \leq X_{N_k}$  مرتب می‌شوند که  $N_k$  تعداد اشیاء در خوشه  $k$  می‌باشد. در ابتدا ثابت

می‌شود که حداقل مقدار یعنی  $c$ ، در فاصله  $[X_i, X_{N_k}]$  قرار دارد. اگر حداقل مقدار در خارج از این محدوده باشد، برای مثال  $X_{N_k} > c$ ، آنگاه از آنجایی که برای هر  $i = 1, 2, \dots, N_k$ ،  $|X_i - X_{N_k}| < |X_i - c|$  می‌شود، لذا  $d(X_{N_k}) < d(c)$ ، که این یک تناقض بوده، پس  $c$  در فاصله  $[X_i, X_{N_k}]$  است. بنابراین،

$$d(c) = \sum_{i \in I^+} (c - X_i)^\beta + \sum_{i \in I^-} (X_i - c)^\beta$$

که در آن  $I^+$  نشانگر اندیس‌هایی از  $N_k, \dots, 2, 1$  است که به ازای آن  $c > X_i$  بوده و  $I^-$  اندیس‌هایی است که  $c \leq X_i$  می‌شود. مشتق اول  $d'(c)$  به صورت زیر می‌باشد.

$$d'(c) = \beta \left( \sum_{i \in I^+} (c - X_i)^{\beta-1} - \sum_{i \in I^-} (X_i - c)^{\beta-1} \right)$$

و مشتق دوم نیز به صورت زیر است.

$$d''(c) = \beta(\beta - 1) \left( \sum_{i \in I^+} (c - X_i)^{\beta-2} + \sum_{i \in I^-} (X_i - c)^{\beta-2} \right)$$

برای  $\beta > 1$  معادله بالا مقدار مثبت دارد که محدب بودن  $d(c)$  را نشان می‌دهد. محدب بودن این رابطه برای  $\beta > 1$ ، به ویژگی‌های مفید دیگری نیز منجر می‌شود. فرض کنید  $d(X_i^*)$  مینیمم  $d(X_i)$  برای هر  $i = 1, 2, \dots, N_k$  باشد و  $X_i'$  مقداری از  $X_i$  باشد که در میان تمام مقادیری که از  $X_i^*$  کوچک‌ترند به  $X_i^*$  نزدیک‌تر باشد. به همین ترتیب، اگر  $X_i''$  مقداری از  $X_i$  باشد که در میان تمام مقادیری که از  $X_i^*$  بزرگ‌ترند، به  $X_i^*$  نزدیک‌تر باشد، آنگاه مینیمم  $d(c)$  در فاصله  $[X_i', X_i'']$  قرار می‌گیرد. به کمک این ویژگی می‌توان الگوریتم زیر را به منظور یافتن مرکز مینکوفسکی مجموعه  $\{X_i\}$  از مقادیر  $X_1 \leq X_2 \leq \dots \leq X_{N_k}$  برای  $\beta > 1$  به کار برد. مرکز مینکوفسکی برای  $\beta > 1$  به شرح زیر است:

۱- مقداردهی اولیه: قرار دهید  $c = X_i^*$ ، که مینیمم  $d(c)$  بر مجموعه  $\{X_i\}$  است. همچنین یک مقدار مثبت  $\lambda$  که تقریباً می‌تواند ۱۰٪ از فاصله  $X_{N_k} - X_1$  را اختیار کند، تعریف می‌شود.

۲- به‌روزرسانی  $c_1$ :  $c - \lambda d'(c)$  محاسبه می‌شود. اگر حاصل آن در فاصله  $[X_i', X_i'']$ ، قرار گیرد، آنگاه  $c_1$  در نظر گرفته می‌شود. در غیر این صورت،  $\lambda$  حدوداً ۱۰ درصد کاهش می‌یابد و این مرحله دوباره تکرار می‌شود.

۳- شرط توقف: آزمایش می‌شود که آیا  $c$  و  $c_1$  همزمان در آستانه تعریف‌شده قرار دارند یا خیر؟ اگر در محدوده تعریف‌شده باشند، فرآیند متوقف و  $c_1$  به عنوان مقدار بهینه  $c$  در نظر گرفته می‌شود. در غیر این صورت به مرحله بعد ارجاع داده می‌شود.

۴- به روزرسانی  $c$ : اگر  $d(c_1) \leq d(c)$ ، آنگاه  $c_1 = c$  و  $d(c_1) = d(c)$  را قرار داده و به مرحله ۲ باز می‌گردد. در غیر این صورت  $\lambda$  تقریباً ۱۰٪ کاهش یافته و سپس به مرحله ۲ بازگشت داده می‌شود.

حال الگوریتم خوشه‌بندی  $k$ -میانگین وزنی را می‌توان به نسخه هوشمند تعمیم داد. در این الگوریتم که توسط میرکین (Mirkin) معرفی شده، مراکز اولیه تصادفی نمی‌باشد و خوشه‌های اولیه توسط الگوریتم  $k$ -میانگین هوشمند که پیشتر توضیح داده شد، بوجود می‌آید و در ادامه با الگوریتم  $k$ -میانگین وزنی ادغام می‌شود. الگوریتم فوق به صورت زیر می‌باشد:

۱- هریک از وزن‌های اولیه ویژگی‌ها، بصورت مساوی مقداردهی اولیه می‌شوند، به طوریکه جمع آنها برابر با ۱ شود. در این الگوریتم یک پارامتر خارجی به کار می‌رود تا حداقل اندازه خوشه‌ها تعریف شود و مقدار آن به طور معمول برابر با ۲ می‌باشد.

۲- اعضای که دورترین فاصله را تا مرکز ثقل کل مجموعه داده‌ها دارند، به عنوان مراکز آزمایشی خوشه‌ها قرار می‌گیرند. خوشه  $S$  که از اعضای که به مرکز آزمایشی نسبت به مرکز ثقل، نزدیک‌ترند، ساخته می‌شود. همچنین وزن‌ها در محاسبات به کمک معیار فاصله زیر دخالت داده می‌شوند،

$$\sum_{j=1}^m w_j^\beta (X_{ij} - c_{kj})^2$$

مرکز آزمایشی به مرکز ثقل خوشه  $S$  یعنی میانگین آن، به روزرسانی می‌شود. با توجه به مرکز مربوطه وزن‌ها، با در نظر گرفتن رابطه (۵-۱) به روزرسانی می‌شوند. اگر مرکز جدید با مرکز قبلی متفاوت باشد این فرآیند از ابتدای مرحله ۲ تکرار می‌شود. در غیر این صورت متوقف شده و خوشه  $S$  از پایگاه داده‌ها حذف می‌شود.

۳- گام ۲ تا زمانی که تمام مشاهدات خوشه‌بندی شوند، ادامه می‌یابد.

۴- خوشه‌هایی که تعداد اعضای آن‌ها کمتر از مقدار آستانه‌ای تعریف شده باشند، حذف می‌شوند.

۵- حال الگوریتم  $k$ -میانگین وزنی با مراکز و وزن‌های ایجاد شده، اجرا می‌شود.

اگر در این الگوریتم، به جای فاصله اقلیدسی، فاصله مینکوفسکی به کار گرفته شود، به الگوریتم حاصل  $k$ -میانگین وزنی مینکوفسکی هوشمند گفته می‌شود. با روشی مشابه آنچه قبلاً گفته شد، یک ثابت مثبت بسیار کوچک به  $E_j$  و  $E_u$  اضافه می‌شود تا دقت افزایش یابد. سایر قسمت‌های الگوریتم مشابه الگوریتم  $k$ -میانگین وزنی مینکوفسکی می‌باشد.

الگوریتم  $k$ -میانگین مینکوفسکی وزنی هوشمند نسبت به سایر الگوریتم‌های  $k$ -میانگین، ویژگی‌های مثبتی دارد از جمله آنکه از آنجا که مراکز اولیه تصادفی نمی‌باشند، لذا یک بار اجرای الگوریتم کفایت می‌کند و نیازی به تکرار چندین مرتبه آن به منظور حذف اثرات مراکز تصادفی نمی‌باشد.

پیش از اجرای الگوریتم بر پایگاه داده ها، نرمال سازی اطلاعات به نتایج بهتر منجر می شود. جین و دویز (Jin, ۱۹۸۸) اظهار داشتند که در بسیاری از مطالعات تحلیل خوشه‌ای، داده‌های خام به‌طور مستقیم مورد استفاده قرار نمی‌گیرند و آماده‌سازی داده‌ها برای تحلیل خوشه‌ای مستلزم برخی تبدیلات مانند استانداردسازی یا نرمال سازی می‌باشد. به کمک استانداردسازی داده‌ها می‌توان اندازه داده‌ها را کم نمود. پس از استانداردسازی می‌توان از داده‌ها یک اندازه مکانی را کسر نمود و حاصل را به یک اندازه مقیاسی برای هر متغیر تقسیم کرد. اگر فرض کنید  $\underline{X}^* = (X_1^*, \dots, X_n^*)$  دلالت بر داده‌های استاندارد شده داشته‌باشد. آن گاه ماتریس داده‌های استاندارد شده یک ماتریس  $n \times m$  به صورت زیر می‌باشد:

$$(X_1^*, \dots, X_n^*) = \begin{bmatrix} X_{11}^* & \dots & X_{1m}^* \\ \vdots & \ddots & \vdots \\ X_{n1}^* & \dots & X_{nm}^* \end{bmatrix}$$

که در آن

$$X_{ij}^* = \frac{X_{ij} - \mu_j}{\sigma_j} \quad (۱-۶)$$

و  $X_{ij}^*$  مقدار استاندارد شده،  $\mu_j$  اندازه مکانی و  $\sigma_j$  اندازه مقیاسی می‌باشد. می‌توان روش‌های استانداردسازی مختلفی را با انتخاب متفاوت  $\mu_j$  و  $\sigma_j$  در معادله (۱-۶) به کار گرفت. در برخی از روش‌های استانداردسازی معروف  $\mu_j$  را میانگین یا میانه و  $\sigma_j$  را انحراف معیار یا دامنه تغییرات در نظر می‌گیرند.

### ۳- روش تحقیق

در این پژوهش انگیزه اصلی، استفاده از الگوریتم k-میانگین وزنی هوشمند بر پایه تابع زیان مینکوفسکی به منظور دسته بندی مشتریان یک شرکت ارائه دهنده خدمات ارتباطی است. در این بخش ۲، چند پایگاه داده‌های شبیه سازی شده از چند توزیع آماری تولید می‌شود و به کمک الگوریتم k-میانگین وزنی هوشمند با معیار عدم تشابه مینکوفسکی در نرم افزار MATLAB نسخه ۲۰۲۳ خوشه بندی می‌گردد. به کارگیری مشاهدات برچسب دار سبب می‌شود تا دقت الگوریتم به کمک برخی معیارهای بیرونی سنجش اعتبار خوشه بندی (که در بخش بعدی عنوان می‌شود)، ارزیابی گردد. در ادامه نتایج الگوریتم بر پایگاه داده واقعی با بیش از ۳۲۴۰۰۰ مشتری (هر یک با ۷ ویژگی) مورد بررسی قرار می‌گیرد و نتایج جمع بندی می‌شود.

#### ۳-۱- پایگاه داده های شبیه سازی شده

۵ پایگاه داده از توزیع های لاگ نرمال، نرمال، گاما و پواسن هر کدام دارای ۱۰۰ متغیر وابسته و  $m$  ویژگی، تولید می شود به گونه ای که در دو گروه، خوشه بندی شوند. فرآیند تولید داده ها در ادامه نشان داده می شود.





در پایگاه داده‌های تولید شده، به ۵۰ متغیر نخست، برچسب ۱ و به ۵۰ متغیر دوم، برچسب ۲، تعلق می‌گیرد. حال الگوریتم k-میانگین هوشمند وزنی مینکوفسکی را بر پایگاه داده‌ها اجرا کرده و با معیار بیرونی تغییر اطلاعات نرمال شده (NVI<sup>۱</sup>) (Duda، ۲۰۰۳) دقت آن سنجیده می‌شود. اگر NVI در فاصله [۰،۱] قرار گیرد، بدان معناست که کارکرد الگوریتم مناسب است. هر چقدر مقدار آن کوچک‌تر باشد خوشه‌ها همگن‌ترند.

جدول ۱: معرفی توزیع‌ها و روابط بین متغیرهای تصادفی و مستقل تولید شده از این توزیع‌ها

روابط میان متغیرهای تصادفی مستقل	$f(x)$
ضرب $n$ متغیر تصادفی مستقل و هم توزیع LN دارای توزیع LN است.	$x > 0, \sigma > 0, \mu \in R, \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{(\log(x)-\mu)^2}{2\sigma^2}\right)$
جمع $n$ متغیر تصادفی مستقل و هم توزیع نرمال، گاما و پواسن به ترتیب دارای همان توزیع‌ها هستند.	$x, \mu \in R, \sigma > 0, f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$
گاما	$x \geq 0, \alpha, \beta > 0, \frac{\beta^\alpha x^{\alpha-1} e^{-x\beta}}{\Gamma(\alpha)}$
پواسن	$x = 0, 1, \dots, \frac{e^{-\lambda} \lambda^x}{x!}$

<sup>۱</sup> Normalized variation information



جدول ۲: فرآیند تولید چند پایگاه داده از ۱۰۰ متغیر وابسته هر کدام با  $m$  ویژگی که در دو گروه خوشه‌بندی می‌شوند

چند متغیره با ۳ ویژگی	پایگاه داده	$Y = (Y_1, \dots, Y_{100})'$
		$Y_i = (X_{i1} + X_{i2}, X_{i2} + X_{i3}, X_{i3} + X_{i1}) = (Y_{ij}) \text{ for } j = 1, 2, 3$
	نرمال	$Y_{ij} \stackrel{iid}{\sim} \text{Normal}(\cdot, \sigma^2) \text{ for } i = 1, \dots, 50, X_{ij} \stackrel{iid}{\sim} \text{Normal}(\cdot, \sigma^2)$ $X_{ij} \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2), Y_{ij} \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2) \text{ for } i = 51, \dots, 100$
	گاما	$Y_{ij} \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta) \text{ for } i = 1, \dots, 50, X_{ij} \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$ $X_{ij} \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta), Y_{ij} \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta) \text{ for } i = 51, \dots, 100$
تک متغیره	لاگ نرمال	$Y_i = (X_{i1}X_{i2}, X_{i2}X_{i3}, X_{i3}X_{i1}) = (Y_{ij}) \text{ for } j = 1, 2, 3$ $Y_{ij} \stackrel{iid}{\sim} \text{LN}(\mu, \sigma^2) \text{ for } i = 1, \dots, 50, X_{ij} \stackrel{iid}{\sim} \text{LN}(\mu, \sigma^2)$ $Y_{ij} \stackrel{iid}{\sim} \text{LN}(\mu, \sigma^2) \text{ for } i = 51, \dots, 100, X_{ij} \stackrel{iid}{\sim} \text{LN}(\mu, \sigma^2)$
	پواسن	$Y_i = (X_{i1} + X_{i2}, X_{i2} + X_{i3}, X_{i3} + X_{i1}) = (Y_{ij})$ $Y_{ij} \stackrel{iid}{\sim} \text{Poisson}(\lambda) \text{ for } i = 1, \dots, 50, X_{ij} \stackrel{iid}{\sim} \text{Poisson}(\lambda)$ $X_{ij} \stackrel{iid}{\sim} \text{Poisson}(\lambda), Y_{ij} \stackrel{iid}{\sim} \text{Poisson}(\lambda) \text{ for } i = 51, \dots, 100$

## ۲-۳- پایگاه داده های مشتریان یک شرکت ارائه دهنده خدمات ارتباطی

شرکت مورد مطالعه دارای ۳۲۴۶۱۹ مشترک است که از سرویس های ارتباطی بهره مند هستند و در ۲۳ محدوده جغرافیایی در کشور ایران واقع شده اند. به دلیل محدودیت های امنیتی اسامی این ۲۳ محدوده از A1 تا A23 مشخص شده است. هر مشترک دارای ۷ ویژگی کلیدی استفاده از سرویس ها است که با B1 تا B7 مشخص شده است. تمام ۷ ویژگی مورد اشاره (شامل میزان استفاده از سرویس ها، سرعت سرویس مورد استفاده، میزان شارژ سرویس، درآمد حاصل از مصرف سرویس ها و صورتحساب و بدهکاری و بستنکاری) برای تمام مشترکین مشخص است و جامعه فاقد هر گونه داده گم شده می باشد. مراحل تجزیه و تحلیل اطلاعات به شرح زیر است:

### ۱- نرمال سازی داده ها و کاهش بُعد داده ها به کمک تکنیک PCA



۲- خوشه بندی داده ها بر اساس ویژگی های انتخاب شده به کمک الگوریتم  $k$ -میانگین وزنی هوشمند مینکوفسکی با برنامه نویسی در محیط MATLAB نسخه ۲۰۲۳

۳- انتخاب نهایی برچسب ها بر اساس های معیارهای ارزیابی

۴-توصیف خوشه ها به کمک مراکز بهینه هر خوشه و معیارهای تشابه و عدم تشابه و نمودارها

۵-برنامه ریزی به منظور توسعه سرویس ها و نگهداشت مشتریان بر اساس نتایج خوشه بندی

پس از نرمال سازی داده ها کاهش بُعد داده ها با روش PCA انجام گرفت. با توجه به حجم بالای داده ها، کاهش بُعد سبب افزایش دقت الگوریتم و سرعت پردازش اطلاعات می شود. به منظور انتخاب تعداد خوشه ها به ترتیب داده ها را در ۲، ۳، ۴ و ۵ دسته خوشه بندی نموده و با استفاده از ۱۰ معیار ارزیابی خوشه ها که بیشترین همبستگی درون خوشه ای و بیشترین فاصله بین خوشه ها را در نظر می گیرد، تعداد ۳ خوشه تصمیم گیری می شود.

#### ۴-یافته ها

با اجرای الگوریتم بر پایگاه داده های شبیه سازی شده بخش ۳-۱، و تفکیک داده ها در ۲ خوشه، و مقایسه نتایج با برچسب های اولیه، متوسط درصد مشاهداتی که به درستی خوشه بندی شده اند (MA) و معیار تغییر اطلاعات نرمال شده (NVI) در جدول ۳ اندازه گیری شده است.

جدول ۳: نتایج الگوریتم  $k$ -میانگین وزنی مینکوفسکی هوشمند، در پایگاه داده های شبیه سازی شده

پایگاه داده	AM	NVI
نرمال	۱۰۰.۰	۰.۰۰
لاگ نرمال	۱۰۰.۰	۰.۰۰
گاما	۹۸.۰	۰.۲۱
پواسن	۹۴.۰	۰.۴۲

نتایج اجرای الگوریتم بر ۴ پایگاه داده شبیه سازی شده گویای دقت بالای الگوریتم در تفکیک سازی خوشه ها است. مقادیر کوچک NVI نشان دهنده خوشه بندی دقیق داده هاست.



در ادامه مشترکین شرکت ارائه دهنده خدمات ارتباطی مورد بررسی قرار می گیرد. به منظور تصمیم گیری جهت انتخاب تعداد خوشه بهینه برای خوشه بندی ۳۲۴۶۱۹ مشترک شرکت هر کدام با ۷ ویژگی، به ترتیب در ۲، ۳، ۴ و ۵ خوشه تقسیم می شوند و بر اساس نتایج خوشه ها ۱۰ معیار ارزیابی در جدول ۴ سنجیده می شود.

جدول ۴: نتایج معیارهای ارزیابی به منظور انتخاب بهینه تعداد خوشه ها

validity index	>	<	۲ خوشه	۳ خوشه	۴ خوشه	۵ خوشه
Dunn index	*		۰.۰۶۴۹۰	۰.۱۱۱۹۱	۰.۰۰۰۰۶	۰.۰۰۰۰۷
Generalized Dunn-۳۳ criterion	*		۰.۵۶۶۰۳	۰.۹۵۳۴۸	۰.۶۵۴۲۲	۰.۰۰۳۲۵
Score Function criterion	*		۰.۰۰۰۰۰	۰.۲۹۰۷۵	۰.۰۰۰۰۰	۰.۰۰۰۰۰
Silhouette index.	*		۰.۹۷۸۳۱	۰.۹۹۶۰۳	۰.۲۰۳۸۹	۰.۱۲۵۹۹
C-criterion		*	۰.۴	۰.۱	۰.۸	۰.۸
COP INDEX		*	۰.۴	۱.۳	۹.۶	۷.۴
CS index		*	۳.۴	۲.۶	۱۰.۵	۹.۳
Davies-Bouldin criterion		*	۱.۱	۰.۷	۲.۷	۳.۵
Davies-Bouldin criterion		*	۱.۱	۰.۷	۲.۷	۳.۵
Xie-Beni criterion.		*	۰.۱	۱.۰	۶.۴	۴.۷

در جدول ۴، مقادیر بالاتر ۴ شاخص اول و مقادیر پایین تر ۶ شاخص دوم نشان دهنده نتایج بهتر است. در هر ردیف بهترین نتایج پررنگ شده است. مطابق نتایج انتخاب ۳ خوشه بهترین معیارهای ارزیابی را نتیجه می دهد. لذا مشتریان شرکت مورد مطالعه در ۳ خوشه اصلی دسته بندی خواهند شد.

جدول ۵- توصیف خوشه ها بر اساس ۷ ویژگی مشترکین

خوشه ها	متوسط سرعت سرویس	میزان استفاده از سرویس	میزان خرید بسته	متوسط درآمد	متوسط صورتحساب	بدهکاری	بستانکاری
خوشه ۱	۱۲,۹۸۲	۱۱۹,۴۴۰	۵۸۸,۴۳۵	۳,۹۶۶,۰۲۲	۵,۸۰۸,۲۶۳	۲۰۶,۹۲۵	۴۴,۶۹۸
خوشه ۲	۱۸,۳۴۹	۱۶۳,۹۷۷	۳,۱۵۶,۵۱۷	۹,۷۴۷,۲۵۸	۹,۳۴۰,۱۵۹	۳۰۴,۹۹۰	۵۹,۳۲۸
خوشه ۳	۲۹,۶۴۹	۳۹۵,۴۲۱	۱۹,۳۰۵,۵۸۹	۳۱,۴۱۷,۶۹۸	۱۵,۲۴۶,۰۲۰	۵۰۶,۷۸۷	۹۰,۲۳۲



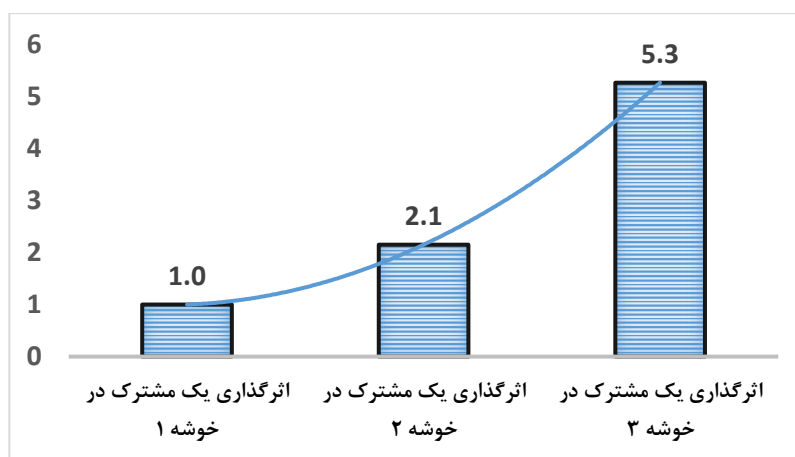
در ادامه در جدول ۶، ۲۳ ناحیه زیر مجموعه شرکت مورد مطالعه در سه خوشه دسته بندی می شوند و به طور کلی خوشه اول با ۵۴ درصد جمعیت، خوشه دوم با ۳۳ درصد و خوشه سوم با ۱۳ درصد از مشترکین شکل می گیرد.

جدول ۶- فراوانی و درصدهای فراوانی خوشه ها به تفکیک ۲۳ ناحیه

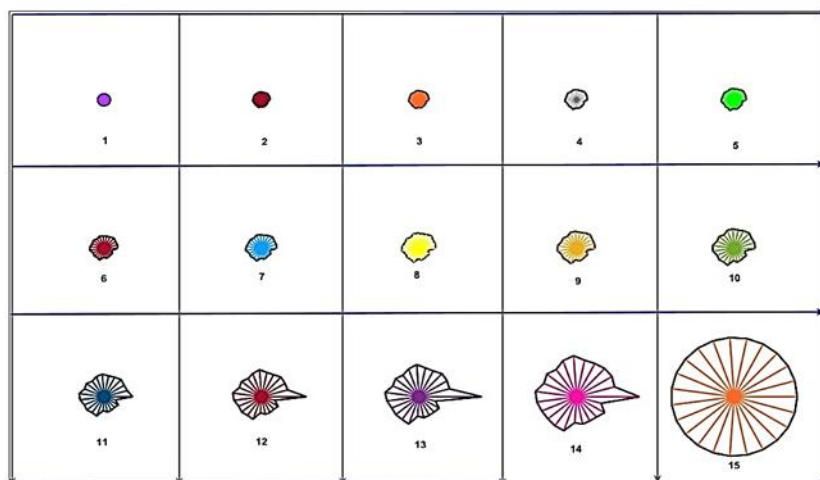
ناحیه	تعداد مشترک خوشه ۱	تعداد مشترک خوشه ۲	تعداد مشترک خوشه ۳	تعداد کل مشترکین	درصد مشترک خوشه ۱	درصد مشترک خوشه ۲	درصد مشترک خوشه ۳
A1	۳,۶۰۱	۱,۹۹۷	۷۵۵	۶,۳۵۳	۵۷	۳۱	۱۲
A2	۴,۴۸۴	۱,۶۴۷	۶۸۶	۶,۸۱۷	۶۶	۲۴	۱۰
A3	۷,۲۹۳	۴,۲۳۷	۲,۰۷۶	۱۳,۶۰۶	۵۴	۳۱	۱۵
A4	۱,۹۵۹	۱,۲۰۷	۴۷۹	۳,۶۴۵	۵۴	۳۳	۱۳
A5	۹,۸۷۸	۵,۹۷۵	۳,۲۱۲	۱۹,۰۶۵	۵۲	۳۱	۱۷
A6	۱۳,۶۳۲	۱۰,۲۶۹	۲,۹۰۰	۲۶,۸۰۱	۵۱	۳۸	۱۱
A7	۶,۰۳۷	۴,۰۸۶	۲,۲۲۰	۱۲,۳۴۳	۴۹	۳۳	۱۸
A8	۱۱,۱۱۳	۶,۸۳۷	۲,۶۹۴	۲۰,۶۴۴	۵۴	۳۳	۱۳
A9	۴,۶۱۳	۲,۶۳۰	۹۹۷	۸,۲۴۰	۵۶	۳۲	۱۲
A10	۱۱,۷۸۱	۷,۹۰۴	۳,۰۱۶	۲۲,۷۰۱	۵۲	۳۵	۱۳
A11	۱۱,۲۴۱	۹,۰۶۰	۴,۸۳۶	۲۵,۱۳۷	۴۵	۳۶	۱۹
A12	۸,۹۵۷	۶,۳۸۷	۲,۹۳۰	۱۸,۲۷۴	۴۹	۳۵	۱۶
A13	۶,۹۰۸	۵,۲۳۰	۱,۷۶۹	۱۳,۹۰۷	۵۰	۳۸	۱۳
A14	۲,۷۲۱	۱,۸۱۱	۶۶۳	۵,۱۹۵	۵۲	۳۵	۱۳
A15	۲۰,۸۷۷	۱۱,۹۱۵	۳,۳۱۸	۳۶,۱۱۰	۵۸	۳۳	۹
A16	۷,۴۳۹	۳,۹۲۰	۱,۳۶۱	۱۲,۷۲۰	۵۸	۳۱	۱۱
A17	۵,۸۱۰	۳,۰۰۶	۹۲۴	۹,۷۴۰	۶۰	۳۱	۹
A18	۱۰,۶۹۰	۵,۹۵۱	۲,۱۰۴	۱۸,۷۴۵	۵۷	۳۲	۱۱
A19	۲,۹۰۷	۱,۲۳۷	۶۰۴	۴,۷۴۸	۶۱	۲۶	۱۳
A20	۸,۰۸۲	۴,۴۷۵	۱,۷۵۳	۱۴,۳۱۰	۵۶	۳۱	۱۲
A21	۸,۸۶۸	۴,۵۳۳	۱,۶۵۲	۱۵,۰۵۳	۵۹	۳۰	۱۱
A22	۳,۹۶۵	۲,۷۳۰	۷۷۴	۷,۴۶۹	۵۳	۳۷	۱۰
A23	۱,۶۸۵	۱,۰۲۷	۲۸۴	۲,۹۹۶	۵۶	۳۴	۹

مجموع	۱۷۴,۵۴۱	۱۰۸,۰۷۱	۴۲,۰۰۷	۳۲۴,۶۱۹	۵۴	۳۳	۱۳
-------	---------	---------	--------	---------	----	----	----

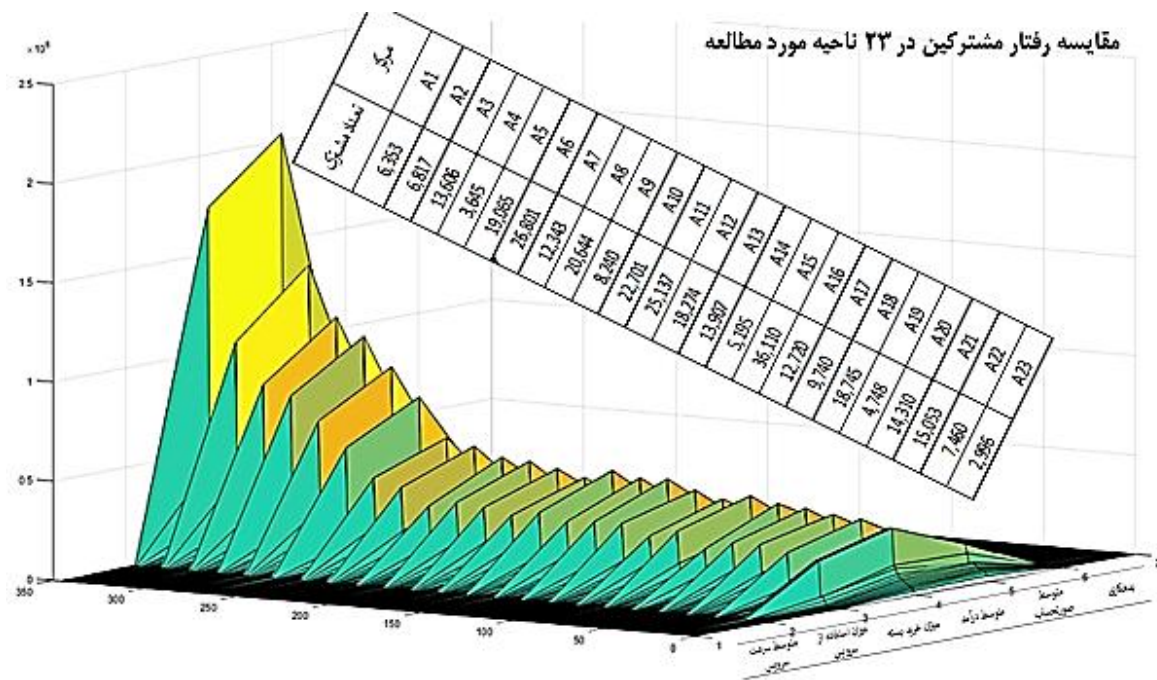
با محاسبه سرانه درآمدزایی هر یک از مشترکین در هریک از خوشه ها و تخصیص ضریب ۱ به سرانه درآمد خوشه ۱ و محاسبه نسبت سرانه سایر خوشه ها به خوشه ۱، مشخص می شود که رشد درآمد در خوشه های دوم و سوم نسبت به خوشه اول دارای رشد نمایی است به طوریکه یک مشترک در خوشه سوم حدود ۵.۳ برابر یک مشترک در خوشه اول و یک مشترک در خوشه دوم حدود ۲.۱ برابر یک مشترک در خوشه اول برای شرکت درآمد زایی دارد.



شکل ۱- نسبت تأثیر درآمدی یک مشترک در هر خوشه در مقایسه با خوشه ۱

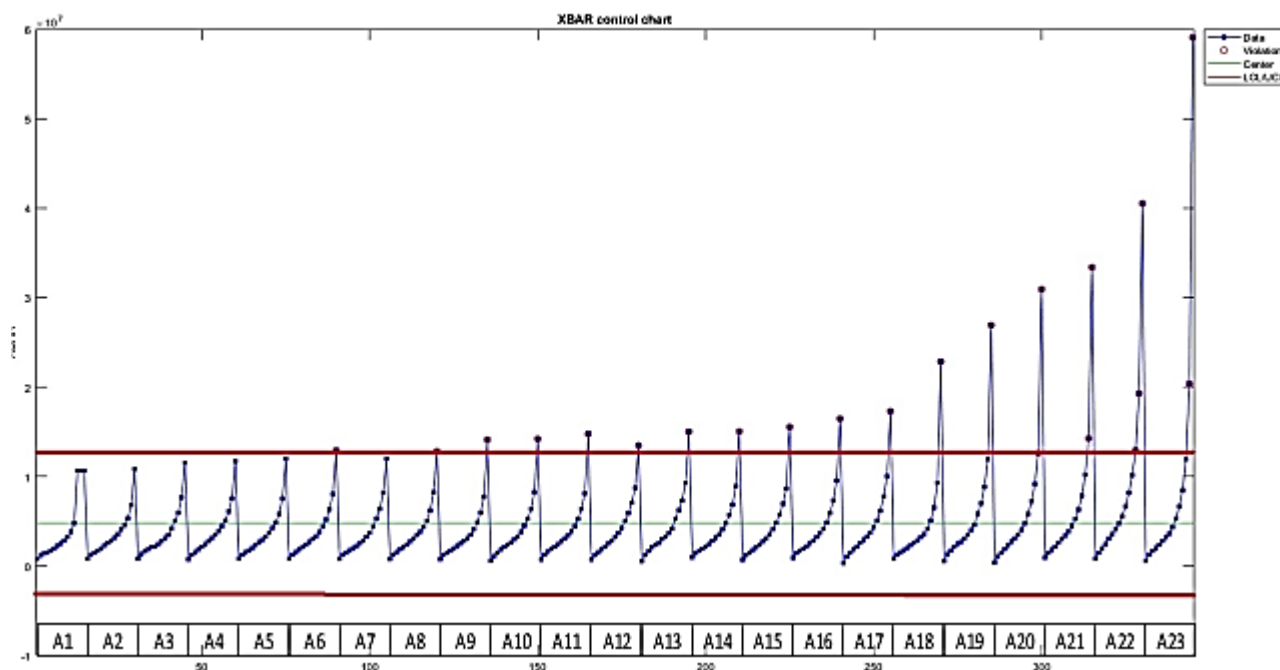


شکل ۲- سیر تکاملی هر مشترک در خوشه ۱ تا ۳ (هر خوشه به ۵ کران تقسیم شده است)



شکل ۳- نمودار سه بعدی مقایسه رفتار مشترکین در ۲۳ ناحیه مورد مطالعه

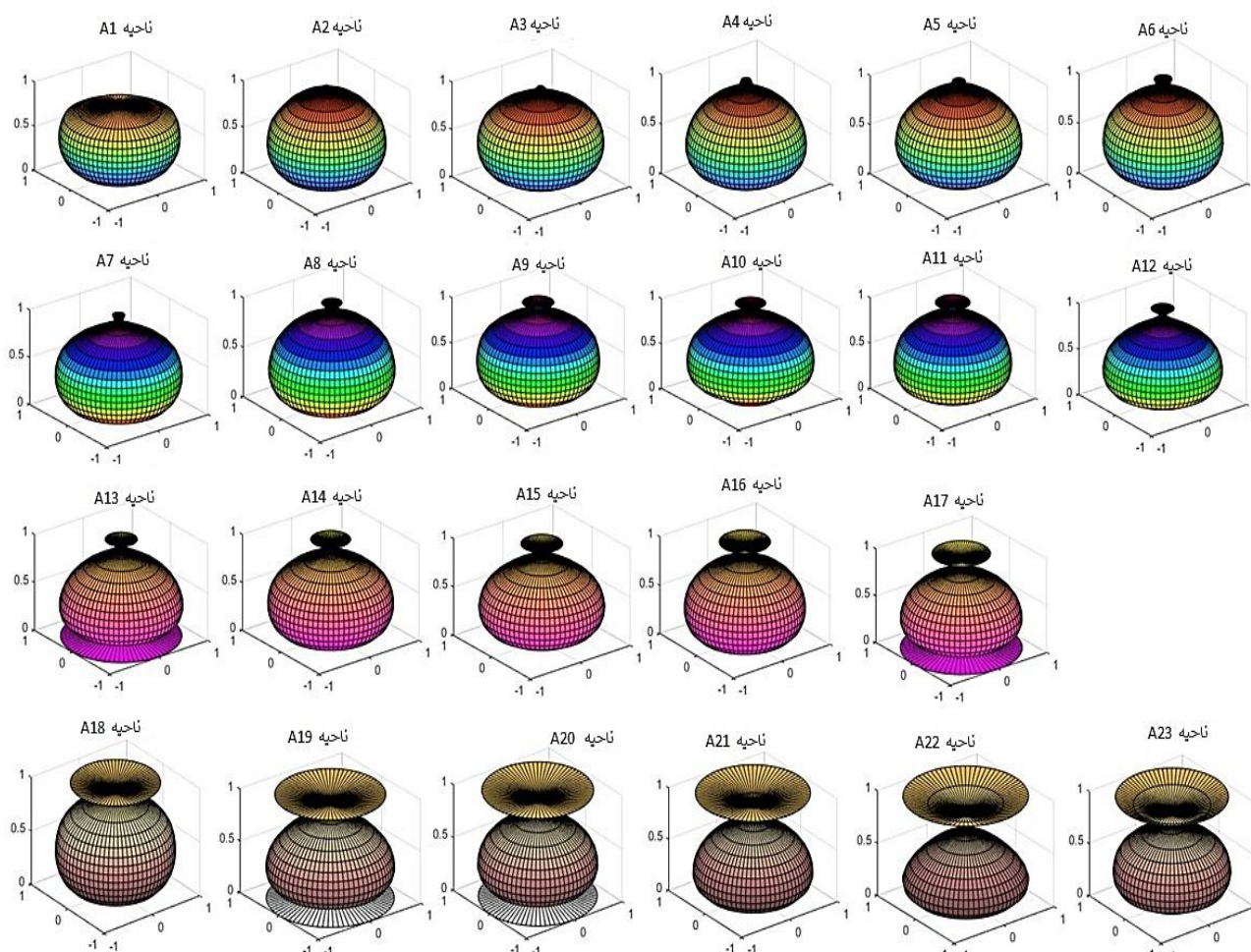




شکل ۴- شناسایی ناحیه هایی که مشترکین خاص در آن قرار گرفته اند

در شکل ۳ رشد ویژگی های درآمدی و مصرف خدمات تمام مشترکین ۲۳ ناحیه مورد مطالعه به تفکیک ۳ خوشه اصلی و هر خوشه ۵ کران بندی، نشان داده شده است. شکل ۴ با در نظر گرفتن ۱۵ کران بندی برای مشترکین هر یک از ناحیه های ۲۳ گانه، ناحیه هایی که بالاتر از خط بالایی قرار گرفته است، دارای مشترکین بسیار خاص و با درآمد زایی بسیار زیاد است که می بایست به منظور خدمات بهتر و توسعه سرویس ها در اولویت قرار گیرند.

در ادامه در شکل ۵، مشترکین هر یک از ناحیه ها را به کمک تلفیقی از توابع مثلثاتی، در مقیاس واحد [۱ و -۱] بصورت اشکال هندسی با طیف های رنگی ترسیم شده اند. به کمک تصاویر بدست آمده می توان نواحی که مشترکین آن دارای رفتار مشابه (از منظر درآمد زایی برای شرکت و میزان مصرف خدمات) هستند را بر اساس تشابه اشکال هندسی بدست آمده، تفکیک کرد و برای هر ناحیه A<sup>۱</sup> تا A<sup>۲۳</sup> متناسب با شکل بدست آمده (که متأثر از بضاعت مشترکین، فرهنگ و... است)، برنامه ریزی مربوطه را انجام داد. صفحه های زیر اشکال هندسی در سه ناحیه A<sup>۱۷</sup>، A<sup>۱۹</sup> و A<sup>۲۰</sup> نشان دهنده مشترکینی است که به تازگی سرویس را خریده اند یا استفاده آنها از سرویس بسیار کم است.



شکل ۵- تبدیل مشترکین هر ناحیه به اشکال هندسی

شکل های ۱ تا ۵ به رده بندی مشتریان هر ناحیه و شناسایی ناحیه هایی که رفتار مشترکین در آن مشابه است کمک می کند. با مقایسه اشکال بدست آمده، ناحیه های (A۱-A۸) در رده پایین تر، ناحیه های (A۹-A۱۷) در رده متوسط و (A۱۸-A۲۳) در قوی ترین رده قرار دارند لذا می توان ناحیه های جغرافیایی رده های بالاتر را در اولویت برنامه های توسعه قرار داد. همچنین در هر ناحیه مشتریان خوشه سوم با بیشترین درآمدزایی برای شرکت می بایست با ارائه بسته های تشویقی و سرویس های خاص در نگهداشت و وفاداری آنها تأثیر گذاشت. مشتریان خوشه دوم با تدوین برنامه های خاص به خوشه سوم سوق داده شوند تا از تنزل آن به خوشه دوم جلوگیری شوند. بسیاری از مشترکین خوشه اول به تازگی سرویس را خریداری کرده اند و مجدد در شش ماه آینده می بایست



مورد بررسی قرار گیرند تا به خوشه های دوم و سوم هدایت شوند. و در مواردی که سرویس در خوشه اول بلااستفاده مانده پس از ۶ ماه می بایست مورد بررسی قرار گیرد.

## ۵-نتایج

در این پژوهش تابع زیان مینکوفسکی با  $p=1$  به عنوان معیار عدم تشابه به جای سایر معیارهای رایج از قبیل مربع اقلیدسی در خوشه بندی  $k$ -میانگین هوشمند وزنی، به کار برده شده است. با در نظر گرفتن  $p=1$ ، مراکز در هر خوشه به جای میانگین میانه مشاهدات است لذا داده های پرت تأثیر کمتری بر مرکز داده ها دارند. در اینجا با به کارگیری الگوریتم  $k$ -میانگین هوشمند وزنی مینکوفسکی، کارکرد آن بر شماری از پایگاه داده های شبیه سازی شده که از پیش برچسب گذاری شده اند، به یاری برخی معیارهای درونی و بیرونی همچون  $NVI$ ،  $AM$ ، سنجیده شده است. در ادامه به منظور استفاده از کاربرد الگوریتم در تحلیل مشتریان یک شرکت ارائه دهنده خدمات مخابراتی، مشتریان در ۳ خوشه اصلی دسته بندی می شوند و بهینه بودن نتایج به کمک ۱۰ معیار ارزیابی سنجیده می شود. از آنجاکه انتخاب مراکز اولیه به تصادف نیست در نتیجه سرعت پردازش اطلاعات (به دلیل نیاز نداشتن به تکرار الگوریتم به منظور حذف اثرات تصادفی مراکز تصادفی) افزایش می یابد و تنها یک باز اجرای الگوریتم کافی است. پس از تفسیر هر یک از خوشه ها به کمک شاخص های آماری، نمودارها و اشکال هندسی مدیران سازمان برای توسعه خدمات و نگهداشت مشتریان ارزشمند، بسته های تشویقی و خاص برای هر یک از مشتریان در هر یک از خوشه ها در نظر گرفته اند همچنین با برنامه ریزی به منظور تقویت مشتریان خوشه ۱ جهت ترغیب آنها به استفاده از خدمات شرکت، از ریزش آنها جلوگیری بعمل می آید.



## ۶-منابع

۱. Ahmadzadehgoli, N., Mohammadpour, A., Behzadi M.H. (۲۰۱۷). LINEX k-means: Clustering by an Asymmetric Dissimilarity Measure, Journal of Statistical Theory and Applications, volume ۱۷, issue ۱, pages ۲۹-۳۸.
۲. Davies D., and Bouldin D. (۱۹۷۹). A Cluster Separation Measure. IEEE Transactions on Pattern Analysis and Machine Intelligence, ۱(۲): ۲۲۴-۲۲۷.
۳. Duda, R. O., Hart, P. E., and Stork, D. G. (۲۰۰۱). Pattern Classification. John Wiley and Sons, Inc. New York, ۲nd Edition.
۴. Harris T. J. (۱۹۹۲). Optimal Controllers for Nonsymmetric and Nonquadratic Loss Functions, Technometrics ۳۴, ۲۹۸-۳۰۶.
۵. Hartigan, J. (۱۹۷۵). Clustering Algorithms. Toronto: JohnWiley & Sons.
۶. Jiang, D., Tang, C., and Zhang, A. (۲۰۰۴). Cluster analysis for Gene Expression Data: A survey. IEEE Transactions on Knowledge and Data Engineering, ۱۶(۱۱):۱۳۷۰-۱۳۸۶.
۷. Kummamuru K., Krishnapuram R. and Agrawal R. (۲۰۰۵). On Learning Asymmetric Dissimilarity Measures, ICDM '۰۵Proceedings of the Fifth IEEE International Conference on Data Mining, ۶۹۷-۷۰۰.



۸. Macqueen, J. (۱۹۶۷). Some methods for classification and analysis of multivariate observations. In Proceedings of the ۵th Berkeley symposium on mathematical statistics and probability, volume ۱, pages ۲۸۱–۲۹۷. Berkeley, CA: University of California Press .
- ۹.
۱۰. Mirkin, B. (۲۰۰۵). Clustering for Data Mining: A Data Recovery Approach. Computer Science and Data Analysis Series. Boca Raton, FL: Chapman & Hall/CRC.
۱۱. Modha D. S. and Spangler W. S. (۲۰۰۳). Feature Weighting in k-means Clustering, Machine Learning ۵۲, ۲۱۷–۲۳۷
۱۲. Parsian A. and Kirmani S. N. U. A. (۲۰۰۲). Estimation under LINEX Loss Function, Handbook of Applied Econometrics and Statistical Inference ۱۶۵, ۵۳–۷۶
۱۳. Reichart R. and Rappoport A. (۲۰۰۹). The NVI Clustering Evaluation Measure. In Proceedings of the Thirteenth Conference on Computational Natural Language Learning, ۱۶۵–۱۷۳
۱۴. Seidpishe, M., Mohammadpour, A. (۲۰۱۷). Hierarchical Clustering of Heavy-Tailed Data Using a New Similarity Measure .
۱۵. Varian H. R., (۱۹۷۵). A Bayesian Approach to Real Estate Assessment. Studies in Bayesian Econometrics and Statistics in Honor of Leonard J. Savage (Eds S. E. Fienberg and A. Zellner), ۱۹۵–۲۰۸