





دانشگاه اصفهان

دانشکده ریاضی و آمار

گروه ریاضی کاربردی و علوم کامپیوتر

پروژه کارشناسی رشته علوم کامپیوتر

روش‌های خوشه‌بندی اصلاح شده

استاد راهنما:

دکتر محسن علمبردار میبدی

دکتر نجمه حسینی منجری

دانشجو:

زهرا مشکوتی

بهمن ماه ۱۴۰۳

### چکیده

در این تحقیق هدف بهبود عملکرد الگوریتم  $k$ -means و معرفی نسخه‌های بهبود یافته این الگوریتم برای حل مسائل خوشه‌بندی با دقت و سرعت بالاتر می‌باشد. به این منظور ابتدا به تعریف کلی الگوریتم  $k$ -means معمولی و سپس به تعریف چهار الگوریتم (۱) الگوریتم  $k$ -means اصلاح شده، (۲) الگوریتم  $k$ -means سراسری، (۳) الگوریتم  $k$ -means سریع و (۴) الگوریتم  $k$ -means سریع بهبود یافته می‌پردازیم. در انتها تمام الگوریتم‌ها پیاده‌سازی شده و گزارشی از نتایج حاصل ارائه می‌گردد. واژگان کلیدی: خوشه‌بندی، الگوریتم  $k$ -means، الگوریتم بهبود یافته  $k$ -means.

# فهرست مطالب

ج	۱	مقدمه
۱	۲	تعاریف مقدماتی
۱	۱.۲	الگوریتم k-means . . . . .
۳	۱.۱.۲	مراحل پیاده‌سازی الگوریتم . . . . .
۱	۳	۳ انواع بهبودیافته‌ی الگوریتم k-means
۱	۱.۳	الگوریتم k-means سراسری . . . . .
۲	۱.۱.۳	مراحل پیاده‌سازی الگوریتم . . . . .
۳	۲.۳	الگوریتم k-means سراسری سریع . . . . .
۵	۳.۳	الگوریتم بهبودیافته‌ی k-means سراسری سریع . . . . .
۵	۱.۳.۳	مراحل پیاده‌سازی الگوریتم . . . . .
۷	۴.۳	الگوریتم k-means بهبود یافته . . . . .
۷	۱.۴.۳	الگوریتم یافتن نقطه‌ی اولیه k-th مرکز . . . . .
۸	۲.۴.۳	مراحل پیاده‌سازی الگوریتم . . . . .
۱۰	۵.۳	نتایج پیاده‌سازی و مشاهدات . . . . .



# فصل ۱

## مقدمه

یکی از قضایای اساسی در مبحث آنالیز داده، مبحث خوشه‌بندی است. خوشه‌بندی یا به عبارتی تقسیم داده‌ها و قرار دادن داده‌های مشابه در یک خوشه، برای استخراج الگوها و ساختارهای پنهان موجود در یک مجموعه داده ضروری است. الگوریتم‌های مختلفی برای مسئله‌ی خوشه‌بندی مطرح شده‌اند مانند

- خوشه‌بندی سلسله‌مراتبی

- DBScan

- k-means

- ....

در بین الگوریتم‌های خوشه‌بندی الگوریتم k-means به سبب اینکه در حین سادگی بسیار مؤثر واقع می‌شود مورد توجه است.

این الگوریتم اولین بار توسط Steinhaus Hugo در سال ۱۹۵۷ معرفی شد و سپس توسط James MacQueen در سال ۱۹۶۷ نام‌گذاری و به شهرت رسید و تا به این زمان به یکی از روش‌های اساسی

در مبحث یادگیری بدون نظارت تبدیل شده است.

ایده‌ی اصلی این الگوریتم تقسیم مجموعه داده به  $K$  خوشه بدون داده‌های مشترک است به طوری که مجموع مربعات فاصله‌ی هر نقطه تا مرکز خوشه‌ای که در آن قرار دارد مینیمم باشد. به این معیار مربع فاصله درون خوشه‌ای (WCSS) گفته می‌شود. در نتیجه خوشه‌بندی نقاط متناظر با حل یک مسئله بهینه‌سازی می‌باشد. تابع هدف این مسئله به صورت زیر می‌باشد برای حل این مسئله باید مراکز بهینه همراه با تخصیص هر نقطه به شبیه ترین خوشه به آن را بیابیم.

فرمول‌بندی ریاضی این مسئله به صورت زیر بیان می‌شود:

$$\text{minimize} \quad \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

در این تابع  $k$  تعداد خوشه‌ها،  $C_i$  مجموعه نقاط در خوشه‌ی  $i$ ام،  $\mu_i$  مرکز خوشه‌ی  $i$ ام و  $x$  یک نقطه در مجموعه داده است.

در ادامه در الگوریتم‌های ذکر شده همگی دارای تابع هدف یکسان و برابر با تابع هدف ذکر شده می‌باشند. همچنین معیار شباهت استفاده شده فاصله اقلیدسی می‌باشد و هر نقطه در هر تکرار به خوشه‌ی متناظر با مرکزی که کمترین فاصله از آن را دارد تخصیص داده می‌شود.

امروزه انواع مختلفی از این الگوریتم معرفی شده است مانند Global، fast k-means، modified-kmeans و k-means که انواع بهبود یافته‌ی الگوریتم k-means می‌باشند.

به وضوح قابل مشاهده است که k-means و انواع آن در حوزه‌های مختلف مانند بازاریابی و دسته‌بندی مشتری‌ها با رفتار مشابه، حوزه‌های مربوط به تصویر مانند تشخیص چهره و بخش بندی تصویر

و حوزه‌های پزشکی و ... به طور گسترده در حال استفاده هستند.  
بنابراین به دلیل کاربر گسترده این الگوریتم خوشه‌بندی، بهبود آن هم از نظر دقت و هم سرعت ضروری  
و لازم به نظر می‌رسد.



## فصل ۲

### تعاریف مقدماتی

#### ۱.۲ الگوریتم k-means

روش‌ها و الگوریتم‌های متعددی برای تبدیل اشیاء به گروه‌های هم‌شکل یا مشابه وجود دارد. الگوریتم k-means یکی از ساده‌ترین و محبوب‌ترین الگوریتم‌هایی است که در داده‌کاوی<sup>۱</sup> بخصوص در حوزه یادگیری نظارت نشده<sup>۲</sup> به کار می‌رود. معمولاً در حالت چند متغیره، باید از ویژگی‌های مختلف اشیاء به منظور طبقه‌بندی و خوشه کردن آن‌ها استفاده کرد. به این ترتیب با داده‌های چند بعدی سروکار داریم که معمولاً به هر بعد از آن، ویژگی یا خصوصیت گفته می‌شود. با توجه به این موضوع، استفاده از توابع فاصله مختلف در این جا مطرح می‌شود.

الگوریتم خوشه‌بندی k-means از گروه روش‌های خوشه‌بندی تفکیکی<sup>۳</sup> محسوب می‌شود و درجه پیچیدگی محاسباتی آن برابر با  $O(n^{dk+1})$  است، به شرطی که  $n$  تعداد اشیاء،  $d$  بعد ویژگی‌ها و  $k$  تعداد خوشه‌ها باشد. همچنین پیچیدگی زمانی برای این الگوریتم برابر با  $O(ndki)$  است، که البته منظور از  $i$  تعداد

---

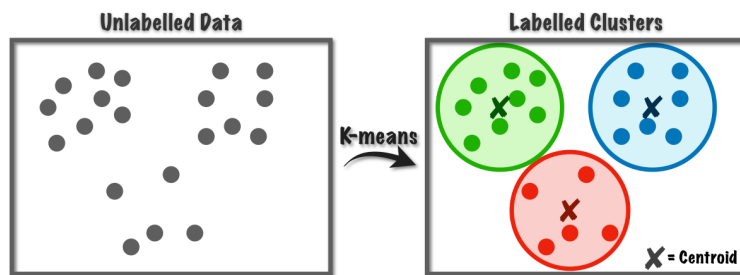
<sup>۱</sup>data mining

<sup>۲</sup>unsupervised learning

<sup>۳</sup>Partitioning Clustering

تکرارهای الگوریتم برای رسیدن به جواب بهینه است.

در خوشه‌بندی k-means از بهینه‌سازی یک تابع هدف استفاده می‌شود. پاسخ‌های حاصل از خوشه‌بندی در این روش، ممکن است به کمک کمینه‌سازی یا بیشینه‌سازی تابع هدف صورت گیرد. به این معنی که اگر ملاک «میزان فاصله» بین اشیاء باشد، تابع هدف براساس کمینه‌سازی خواهد بود پاسخ عملیات خوشه‌بندی، پیدا کردن خوشه‌هایی است که فاصله بین اشیاء هر خوشه کمینه باشد. در مقابل، اگر از تابع مشابهت برای اندازه‌گیری مشابهت اشیاء استفاده شود، تابع هدف را طوری انتخاب می‌کنند که پاسخ خوشه‌بندی مقدار آن را در هر خوشه بیشینه کند.



معمولاً زمانی که هدف کمینه‌سازی باشد، تابع هدف را تابع هزینه<sup>۴</sup> نیز می‌نامند. این کمینه‌سازی به صورت مسئله بهینه‌سازی زیر تعریف می‌شود.

$$\text{minimize } f_k(x) \quad \text{subject to } x = (x^1, \dots, x^k) \in \mathbb{R}^{n \times k},$$

where

$$f_k(x^1, \dots, x^k) = \frac{1}{m} \sum_{j=1}^m \min_{i=1, \dots, k} \|x^j - a^i\|^2.$$

cost function<sup>۵</sup>

در این تابع  $x = (x^1, \dots, x^k) \in \mathbb{R}^{n \times k}$  بیانگر مراکز خوشه‌ها و  $m$  تعداد نقاط می باشد.

## ۱.۱.۲ مراحل پیاده‌سازی الگوریتم

گام مقداردهی اولیه:

در الگوریتم k-means معمولی در ابتدا  $k$  نقطه اولیه به عنوان مراکز خوشه‌ها به صورت تصادفی انتخاب می‌شوند. روش‌های دیگری مانند k-means++ وجود دارند که نقاط اولیه را با معیاری مشخص و غیر تصادفی انتخاب می‌کنند. k-means به شدت به انتخاب نقاط اولیه وابسته است و ممکن است در بهینه‌ی محلی باقی بماند.

گام تکرار:

هر نقطه را به نزدیک‌ترین خوشه‌اش نسبت می‌دهیم. اپدیت مراکز به این صورت که میانگین نقاط درون هر خوشه به عنوان مراکز جدید محاسبه می‌شوند و دوباره هر نقطه به نزدیک‌ترین خوشه‌اش نسبت داده می‌شود.

شرط توقف:

مرحله‌ی قبل تا جایی ادامه پیدا می‌کند که یا به ماکزیمم مقدار تکرار که از قبل معلوم شده برسیم و یا مراکز همگرا شوند. همگرا شدن مراکز به این معناست که دیگر تغییر و بهبود محسوسی در مراکز محاسبه شده دیده نشود.

## فصل ۳

# انواع بهبودیافته‌ی الگوریتم k-means

### ۱.۳ الگوریتم k-means سراسری

همانطور که گفته شد الگوریتم k-means معمولی بسیار به انتخاب نقاط اولیه مراکز وابسته است و انتخاب نقاط اولیه نامناسب می‌تواند در نهایت منجر به گرفتار شدن در بهینه‌ی محلی گردد. بنابراین انواع بهبودیافته‌ای از این الگوریتم معرفی شد که تمرکز بر روی انتخاب مراکز اولیه مناسب باشد. در این الگوریتم‌ها در هر تکرار یک نقطه به عنوان نقطه‌ی اولیه برای  $q$  امین مرکز ( $k \geq q$ ) مسئله با  $k$  خوشه انتخاب می‌شود و به این نوع از انتخاب مراکز، افزایشی<sup>۱</sup> گفته می‌شود.

الگوریتم k-means سراسری یا به اختصار GKM یک الگوریتم از نوع افزایشی می‌باشد. در این الگوریتم در هر تکرار تمام نقاط مجموعه داده به غیر از مراکز فعلی به عنوان نقطه‌ی اولیه احتمالی برای مرکز بعدی در نظر گرفته می‌شود. بنابراین در هر مرحله دارای چندین راه‌حل اولیه هستیم که برای هر کدام

الگوریتم

---

incremental<sup>۱</sup>

k-means معمولی را به کار می‌بریم تا بتوانیم بهترین راه‌حل را میان تمام راه‌حل‌های اولیه بیابیم. این روند تا جایی که هر k خوشه مشخص شود ادامه می‌یابد و همانطور که به نظر می‌رسد بار محاسباتی این الگوریتم بسیار بالاست بنابراین برای مجموعه داده‌های متوسط و بزرگ کاربردی نمی‌باشد.

### ۱.۱.۳ مراحل پیاده‌سازی الگوریتم

• گام اول: مقداردهی اولیه

برای محاسبه‌ی  $x^1$  میانگین تمام نقاط مجموعه داده A به عنوان اولین مرکز در نظر می‌گیریم.

قرار دهید  $q = 1$

• گام دوم: بررسی شرط توقف

قرار دهید  $q = q + 1$  اگر  $q > k$  متوقف شود.

• گام سوم: مراکز  $x^1, \dots, x^{q-1}$  از مرحله‌ی قبل را در نظر بگیرید. سپس هر نقطه‌ی  $a_i$  داخل مجموعه داده که عضو مراکز نباشد را به عنوان مرکز اولیه برای امین q مرکز خوشه در نظر می‌گیریم. به این معنا که در این مرحله دارای چندین راه حل اولیه با q نقطه هستیم:  $(x^1, \dots, x^{k-1}, a_i)$  الگوریتم k-means را بر روی هر کدام از گزینه‌ها اجرا می‌کنیم و بهترین حالت را انتخاب می‌کنیم. بهترین گزینه حالتی است که تابع هدف با در نظر گرفتن آن در مقایسه با بقیه کمترین مقدار را دارا باشد.

سپس مراکز  $(y^1, \dots, y^q)$  بدست آمده را ذخیره می‌کنیم.

• گام چهارم: قرار دهید  $x^i = y^i, i = 1, \dots, q$  و به مرحله‌ی دوم بروید.

## ۲.۳ الگوریتم k-means سراسری سریع

همانطور که اشاره کردیم الگوریتم GKM به دلیل استفاده‌ی متعدد از k-means در هر مرحله دارای زمان پیاده‌سازی بسیار طولانی و بار محاسباتی بالا می‌باشد. بنابراین به دنبال رفع این مشکلات الگوریتم k-means سراسری سریع<sup>۲</sup> معرفی گردید.

این الگوریتم نیز افزایشی می‌باشد و در هر مرحله برای یافتن نقطه‌ی اولیه برای امین  $q$  مرکز  $(k \geq q)$  ( با در نظر گرفتن هر نقطه‌ی  $a_j$  که عضو مراکز نباشد به جای پیاده کردن k-means بر روی هر گزینه یک مقدار با نام  $b_j$  محاسبه می‌شود و بر اساس آن پاسخ را میابیم.

• گام اول: مقداردهی اولیه

محاسبه‌ی  $x^1$ : میانگین تمام نقاط مجموعه داده  $A$  به عنوان اولین مرکز در نظر می‌گیریم.

قرار دهید  $q = 1$

• گام دوم: بررسی شرط توقف

قرار دهید  $q = q + 1$  اگر  $q > k$  متوقف شود.

• گام سوم: مراکز  $x^1, \dots, x^{q-1}$  از مرحله‌ی قبل را در نظر بگیرید. برای انتخاب نقطه‌ی اولیه برای

امین  $q$  مرکز خوشه، برای هر نقطه‌ی  $a_j$  داخل مجموعه داده که عضو مراکز نباشد یک مجموعه

شامل تمام نقاطی که به  $a_j$  نزدیکتر از مرکز خوشه‌ای که در آن قرار دارند هستند، تشکیل می‌دهیم.

که در آن  $d_{q-1}^i$  بیانگر مینیمم فاصله‌ی نقطه‌ی  $i$  از  $q-1$  مرکز خوشه قبلی است.

---

fast global k-means<sup>۲</sup>

$$I = \{i \in \{1, \dots, m\} : \|a^i - a^j\|^2 < d_{q-1}^i\}$$

در این مرحله  $b_j$  را به صورت زیر تعریف می‌کنیم :

$$f_{q-1}(x^1, \dots, x^{q-1}) - f_k(x^1, \dots, x^{q-1}, a^j) = \sum_{i \in I} (d_{q-1}^i - \|a^j - a^i\|^2) = b_j$$

مقدار  $b_j$  نشان می‌دهد با انتخاب  $a_j$  به انواع نقطه‌ی اولیه برای  $q$  امین مرکز به چه اندازه مقدار فعلی تابع هدف کاهش میابد. بنابراین نقطه‌ای که دارای ماکزیمم مقدار  $b$  باشد انتخاب ما برای نقطه‌ی اولیه برای  $q$  امین مرکز خوشه می‌باشد. مراکز بدست آمده ذخیره شده و الگوریتم kmeans روی آنها اجرا می‌شود و پاسخ پیاده‌سازی را به صورت زیر داریم.

$$(y^1, \dots, y^q)$$

• گام چهارم: قرار دهید  $x^i = y^i, i = 1, \dots, q$  و به مرحله‌ی دوم بروید.

بار محاسباتی این الگوریتم بسیار نسبت به GKM پایین‌تر است اما واضح است که دقت این الگوریتم نسبت به آن پایین‌تر می‌باشد زیرا در GKM برای تمام گزینه‌ها k-means اعمال می‌شود و بر اساس مقدار تابع هدف نهایی، جواب مورد نظرم را می‌یابیم ولی در GKM سریع تنها پس از انتخاب مرکز اولیه الگوریتم k-means اعمال شده و به مرحله‌ی بعدی می‌رود. مقدار  $b_j$  مقداری است که در یک تکرار انتظار داریم با انتخاب  $a_j$  به عنوان نقطه‌ی اولیه برای مرکز از تابع هدف کاسته شود اما ممکن است در نهایت این مقدار کاسته شده از تابع هدف نباشد.

پس چگونه می‌توانیم با در نظر گرفتن بار محاسباتی دقت الگوریتم را افزایش دهیم؟

### ۳.۳ الگوریتم بهبودیافته‌ی k-means سراسری سریع

۲

در الگوریتم بهبودیافته‌ی GKM سریع تمام مراحل مانند GKM سریع انجام می‌شوند با این تفاوت که به هنگام انتخاب نقطه‌ی اولیه برای امین  $q$  مرکز ( $k \geq q$ )  $n$  درصد از نقاطی که دارای  $b$  ماکزیمم بودند را انتخاب کرده و روی پاسخ‌های اولیه‌ی بدست آمده الگوریتم k-means اعمال می‌شود.

در الگوریتم GKM برای تمام نقاط به غیر از مراکز این عملیات انجام می‌شود و در GKM سریع تنها برای نقطه با  $b$  ماکزیمم k-means را اعمال می‌کردیم، در نهایت در الگوریتم gkm fast modified برای درصدی از نقاط k-means را اعمال می‌کنیم تا در نهایت پاسخی که مقدار تابع هدف بهتری نسبت به باقی گزینه‌ها داشت را انتخاب کنیم. اینکار سبب بدست آوردن مراکز بهینه‌تر و در نتیجه کاهش احتمال گرفتار شدن در بهینه‌ی محلی می‌شود.

هرچه درصد  $n$  بزرگ‌تری بگیریم تعداد بیشتری از پاسخ‌های اولیه بررسی می‌شود و دقت بالاتر رفته ولی از طرفی زمان اجرا بیشتر خواهد شد.

#### ۱.۳.۳ مراحل پیاده‌سازی الگوریتم

- گام اول : مقداردهی اولیه

برای محاسبه‌ی  $x^1$  میانگین تمام نقاط مجموعه داده  $A$  به عنوان اولین مرکز در نظر می‌گیریم.

$$q = 1 \text{ قرار دهید}$$

---

modified fast global k-means algorithm<sup>۲</sup>



• گام دوم : بررسی شرط توقف

قرار دهید  $q = q + 1$  اگر  $q > k$  متوقف شود.

• گام سوم : مراکز  $x^1, \dots, x^{q-1}$  از مرحله‌ی قبل را در نظر بگیرید. برای انتخاب نقطه‌ی اولیه برای امین  $q$  مرکز خوشه، برای هر نقطه‌ی  $a_j$  داخل مجموعه داده که عضو مراکز نباشد یک مجموعه شامل تمام نقاطی که به  $a_j$  نزدیکتر از مرکز خوشه‌ای که در آن قرار دارند هستند، تشکیل می‌دهیم. که در آن  $d_{q-1}^i$  بیانگر مینیمم فاصله‌ی نقطه‌ی  $i$  از  $q-1$  مرکز خوشه قبلی است.

$$I = \{i \in \{1, \dots, m\} : \|a^i - a^j\|^2 < d_{q-1}^i\}$$

در این مرحله  $b_j$  را به صورت زیر تعریف می‌کنیم:

$$f_{q-1}(x^1, \dots, x^{q-1}) - f_k(x^1, \dots, x^{q-1}, a^j) = \sum_{i \in I} (d_{q-1}^i - \|a^j - a^i\|^2) = b_j$$

برای تمام نقاط به غیر از مراکز این مقدار محاسبه شده و به صورت نزولی مرتب می‌کنیم. سپس  $n < m$  نقطه اول در این لیست که دارای بیشترین مقادیر  $b$  هستند را ذخیره می‌کنیم. بسته به کاربرد می‌توان مقدار  $n$  را درصدی مجموعه داده گرفت (  $m$  تعداد کل نقاط مجموعه داده است).

مراکز بدست آمده ذخیره شده به عنوان نقاط اولیه برای  $q$  امین مرکز در نظر گرفته می‌شوند بنابراین دارای  $n$  راه حل اولیه با  $q$  نقطه هستیم.

$$(x^1, \dots, x^{k-1}, a_i)$$

الگوریتم kmeans را روی هر راه حل اعمال کرده و گزینه‌ای که تابع هدف نسبت به بقیه حالت‌ها

کمتر می‌شود انتخاب می‌شود. پاسخ پیاده‌سازی به صورت زیر در نظر می‌گیریم.

$$(y^1, \dots, y^q)$$

• گام چهارم: قرار دهید  $x^i = y^i$ ,  $i = 1, \dots, q$  و به مرحله‌ی دوم بروید.

انتظار می‌رود دقت این الگوریتم از GKM کمتر ولی از GKM سریع‌تر باشد. همچنین از GKM سریع‌تر باشد که نتایج پیاده‌سازی همین موضوع را تصدیق می‌کنند.

### ۴.۳ الگوریتم k-means بهبود یافته

۴

الگوریتم k-means بهبود یافته نیز یک الگوریتم بهبود یافته‌ی k-means معمولی از نوع افزایشی می‌باشد اما برخلاف الگوریتم‌های معرفی شده در این الگوریتم مقدار k معلوم نیست و روند خوشه‌بندی هنگامی که شرط توقف برقرار شود، متوقف می‌شود. همچنین در این الگوریتم برای یافتن نقطه‌ی اولیه برای k امین مرکز بعدی یک الگوریتم جداگانه به شرح زیر تعریف می‌شود.

### ۱.۴.۳ الگوریتم یافتن نقطه‌ی اولیه k-th مرکز

فرض کنید مراکز  $(x^1, \dots, x^{k-1})$  برای مسئله‌ی با  $k-1$  خوشه معلوم باشد. برای یافتن k امین مرکز در این الگوریتم برای هر نقطه  $a_j$  که عضو مراکز نباشد مجموعه‌ی I که شامل تمام نقاطی است که به این نقطه نسبت به مرکز خوشه فعلیشان نزدیک‌تر هستند را تشکیل می‌دهیم.

modified k-means algorithm<sup>۴</sup>

$$I = \{i \in \{1, \dots, m\} : \|a^i - a^j\|^2 < d_{q-1}^i\}$$

سپس مرکز  $c_j$  که مرکز این مجموعه است را محاسبه می‌کنیم. در این مرحله یک تابع کمکی<sup>۵</sup> به صورت زیر تعریف می‌کنیم:

$$\bar{f}_k(Y) = \frac{1}{m} \sum_{i=1}^m \min\{d_{k-1}^i, \|y - a^i\|^2\}$$

در این تابع  $d_{k-1}^i$  بیانگر مینیمم فاصله‌ی نقطه‌ی  $a_i$  از  $k-1$  مرکز خوشه قبلی است. متناظر با هر  $c_j$  مقدار  $\bar{f}_k(c_j)$  را محاسبه می‌کنیم. هر پارامتر  $c_j$  که مقدار تابع را مینیمم سازد به عنوان نقطه‌ی اولیه برای مرکز بعدی انتخاب می‌شود.

### ۲.۴.۳ مراحل پیاده‌سازی الگوریتم

پارامتر  $\varepsilon$  یک استانه<sup>۶</sup> در شرط توقف می‌باشد. که در این پیاده‌سازی برابر با  $0.001$  است. انتخاب  $\varepsilon$  بزرگ باعث به وجود آمدن خوشه‌های بزرگتر و  $\varepsilon$  کوچکتر باعث به وجود آمدن خوشه‌های مجازی<sup>۷</sup> می‌شود.

• گام اول:

انتخاب استانه مناسب.  $\varepsilon > 0$

محاسبه‌ی  $x^1$ : میانگین تمام نقاط مجموعه داده  $A$  به عنوان اولین مرکز در نظر می‌گیریم.

مقدار  $f^1$  که مقدار متناظر با تابع هدف می‌باشد را محاسبه می‌کنیم.

<sup>۵</sup>auxiliary cluster function

<sup>۶</sup>tolerance

<sup>۷</sup>artificial clusters

قرار می‌دهیم  $k = 1$

• گام دوم : محاسبه‌ی مرکز خوشه بعدی

قرار دهید :  $k+1 = k$

فرض کنید  $(x^1, \dots, x^{k-1})$  مراکز خوشه‌ها برای مسئله‌ی با  $k-1$  خوشه باشد .

با استفاده از الگوریتمی که در ادامه به آن می‌پردازیم، نقطه شروع  $\bar{y}$  را برای  $k$  امین مرکز می‌یابیم.

• گام سوم: اصلاح مراکز

نقاط  $(x^1, \dots, x^{k-1}, \bar{y})$  به عنوان مراکز اولیه در الگوریتم kmeans در نظر می‌گیریم و مسئله‌ی

k-partition را حل می‌کنیم.

فرض کنید  $(y^1, \dots, y^k)$  مراکز به دست آمده‌ی مسئله‌ی k-partition باشند. مقدار تابع هدف

$f^k$  را محاسبه می‌کنیم.

• گام چهارم : شرط توقف، اگر

$$\frac{f^{k-1} - f^k}{f^1} < \epsilon$$

متوقف شود و در غیرای صورت قرار دهید:

$$x^i = y^i, \quad i = 1, \dots, k$$

و به گام ۲ برو.

## ۵.۳ نتایج پیاده‌سازی و مشاهدات

در این جداول مقادیر  $f$ ،  $\alpha$  و  $t$  برای پیاده‌سازی الگوریتم‌های معرفی شده بر روی مجموعه داده‌های مختلف نمایش داده شده‌اند.

به ترتیب مقدار  $f$  بیانگر مقدار تابع هدف،  $\alpha$  بیانگر تعداد نرم‌های محاسبه شده و  $t$  بیانگر زمان پیاده‌سازی می‌باشد.

توجه کنید در پیاده‌سازی الگوریتم k-means بهبود یافته مقدار  $k$  توسط کاربر معلوم نمی‌شود بنابراین ترتیب نتایج بدست آمده از الگوی باقی پیاده‌سازی‌ها پیروی نمی‌کند.

در سه جدول اول نتایج نمایش داده شده بر روی دو مجموعه داده به نام‌های breast و concrete می‌باشد که به ترتیب دارای ۶۸۳ و ۱۰۳۰ تعداد داده و هر دو دارای ۹ بعد می‌باشند.

در سه جدول دوم نتایج نمایش داده شده بر روی دو مجموعه داده به نام‌های TSPLIB ۱۰۶۰ و TSPLIB ۳۰۳۸ می‌باشد که به ترتیب دارای ۱۰۶۰ و ۳۰۳۶ تعداد داده و هر دو دارای ۲ بعد می‌باشند.

تصاویری که به عنوان نمونه مشاهده می‌کنید نتیجه‌ی خوشه‌بندی مجموعه داده ی TSPLIB ۱۰۶۰ با الگوریتم‌های تعریف شده می‌باشد .

$$k = 5$$

Table :۳'۱ Breast and Concrete data sets

No.	$f_{opt}$	k-means	gkm	fast gkm	modified fast gkm	modif
		$f$	$f$	$f$	$f$	$f$
۲	◦	۱۹۳۲۳,۱۷۴	۱۹۳۲۳,۱۷۴	$۱,۹۳E + ۰۴$	۱۹۳۲۳,۲۰۵	$۱,۹۳E + ۰۴$
۵	◦	۱۳۷۷۰,۵۵۷	۱۳۷۰۶,۳۸۶	$۱,۳۸E + ۰۴$	۱۳۷۰۶,۷۹۶	$۱,۴۲E + ۰۴$
۱۰	◦	۱۱۷۱۴,۵۰۲	۱۰۲۰۲,۲۶۰	$۱,۲۴E + ۰۴$	۱۰۲۵۲,۲۴۸	$۱,۰۶E + ۰۴$
۱۵	◦	۱۰۳۸۱,۶۱۴	۸۶۹۹,۷۹۹	$۹,۹۳E + ۰۳$	۹۰۱۳,۶۱۳	$۹,۰۰E + ۰۳$
۲۰	◦	۸۱۸۶,۲۹۱	۷۶۵۶,۵۴۰	$۸,۴۷E + ۰۳$	۸۰۰۴,۹۲۱	$۸,۱۹E + ۰۳$
۲۳	-	-	-	-	-	$۸,۰۰E + ۰۳$
۲	◦	۳۱۴۷۳۱۵۷,۳۷۲	۳۱۴۷۳۱۵۷,۳۷۲	$۳,۲۲E + ۰۷$	۳۱۴۷۳۱۵۷,۳۷۲	$۳,۱۵E + ۰۷$
۵	◦	۲۱۱۰۷۴۴۳,۴۴۹	۱۹۵۳۵۷۳۱,۹۸۷	$۲,۰۱E + ۰۷$	۱۹۵۳۵۷۳۱,۹۸۶	$۲,۰۰E + ۰۷$
۱۰	◦	۱۳۱۳۲۱۹۱,۳۶۶	۱۲۲۸۳۹۷۸,۹۱۵	$۱,۴۳E + ۰۷$	۱۲۳۰۳۷۳۷,۵۳۷	$۱,۲۴E + ۰۷$
۱۵	◦	۹۹۴۸۴۴۴,۷۱۵	۹۲۲۱۸۵۸,۱۴۹	$۹,۶۴E + ۰۶$	۹۲۸۴۸۰۴,۰۰۲	$۹,۵۲E + ۰۶$
۲۰	◦	۸۳۶۷۰۷۶,۵۴۳	۷۶۲۰۲۸۵,۹۲۴	$۸,۶۱E + ۰۶$	۷۶۸۳۶۸۸,۶۹۸	$۸,۰۸E + ۰۶$
۲۶	-	-	-	-	-	$۷,۰۶E + ۰۶$

Table :۳۲ Breast and Concrete data sets

No.	$f_{opt}$	k-means	gkm	fast gkm	modified fast gkm	modified k-me
		$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$
۲	◦	۶۸۳۰	۵۰۱۸۰۰۱	۲۰۷۵۷۳	۱۱۳۶۴۵۳	۱۴۰۸۳۴۶
۵	◦	۱۱۹۵۲۵	۵۶۴۰۰۷۷۴	۳۵۱۶۹۰۷	۱۸۳۵۴۰۰۱	۱۱۲۸۹۳۰۷
۱۰	◦	۶۸۳۰۰	۲۱۲۴۰۴۸۰۴	۱۸۰۵۶۰۲۷	۶۳۷۷۱۴۳۰	۴۶۶۱۷۴۸۲
۱۵	◦	۱۷۴۱۶۵	۴۰۲۱۱۶۹۳۳	۴۴۳۳۸۰۲۸	۱۲۵۲۹۲۲۹۹	۱۰۵۳۸۸۲۶۶
۲۰	◦	۲۱۸۵۶۰	۶۳۶۲۰۶۳۰۴	۸۲۶۸۷۳۱۹	۲۰۴۴۵۹۶۲۷	۱۸۷۴۰۱۵۴۰
۲۳	-	-	-	-	-	۲۴۸۰۴۱۰۱۲
۲	◦	۱۴۴۲۰	۲۳۱۸۶۳۳۰	۲۰۸۷۰۷	۵۱۷۳۳۰۷	۳۲۰۴۳۳۰
۵	◦	۸۲۴۰۰	۱۹۶۰۹۱۴۰۰	۵۷۵۱۲۴۰	۴۱۶۲۳۷۶۱	۲۵۵۶۴۶۰۰
۱۰	◦	۳۷۰۸۰۰	۶۶۳۶۰۸۴۰۰	۳۴۳۰۶۴۰۰	۱۷۶۵۰۹۷۰۷	۱۰۵۸۱۳۹۶۰
۱۵	◦	۲۰۰۸۵۰	۱۱۹۹۱۳۷۳۳۰	۹۱۲۴۰۸۸۸	۳۵۰۵۱۲۲۵۳	۲۳۹۲۵۶۶۴۰
۲۰	◦	۳۲۹۶۰۰	۱۷۹۲۵۵۶۳۸۰	$1/76E + 0.8$	۵۷۱۸۵۴۵۷۱	۴۲۵۶۲۳۸۱۰
۲۶	-	-	-	-	-	۷۱۹۶۷۰۲۷۰

Table :۳۳ Breast and Concrete data sets

No.	$f_{opt}$	k-means $t$	gkm $t$	fast gkm $t$	modified fast gkm $t$	modified k-means $t$
۲	۰	۰/۰۴	۰/۴۳	۵/۴۹	۶/۲۷	۸/۱۷
۵	۰	۰/۰۴	۲/۱۰	۳۸/۷۲	۴۰/۵۵	۱۵/۰۱
۱۰	۰	۰/۰۴	۵/۳۰	۱۲۸/۸۹	۱۷۷/۵۴	۲۲/۲۱
۱۵	۰	۰/۰۴	۷/۰۱	۲۵۴/۸۲	۲۴۸/۷۷	۳۳/۹۲
۲۰	۰	۰/۰۴	۱۰/۵۷	۴۱۷/۰۳	۴۱۰/۴۶	۴۰/۷۴
۲۳	-	-	-	-	-	۵۶۶/۷۳
۲	۰	۰/۰۳	۰/۹۶	۹/۰۰	۷/۷۲	۱۹/۴۴
۵	۰	۰/۰۴	۴/۷۶	۷۷/۰۱	۷۰/۲۳	۳۷/۳۳
۱۰	۰	۰/۰۴	۹/۹۲	۲۷۶/۰۰	۲۵۱/۸۶	۵۱/۹۱
۱۵	۰	۰/۰۴	۱۵/۰۴	۵۵۵/۱۴	۴۷۸/۶۴	۶۷/۴۳
۲۰	۰	۰/۰۴	۲۲/۳۷	۹۳۶/۱۰	۹۱۰/۹۸	۸۱/۸۹
۲۶	-	-	-	-	-	۱۵۷۹/۵۸



Table :۳'۴ TSPLIB۱۰۶۰ and TSPLIB۳۰۳۸ data sets

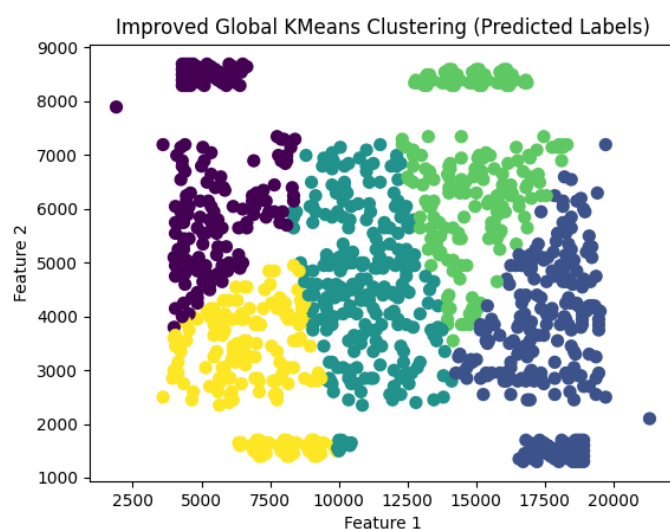
No.	$f_{opt}$	k-means	gkm	fast gkm	modified fast gkm
		$f$	$f$	$f$	$f$
۲	◦	۹۸۳۱۹۴۹۸۸۴/۲۰۱	$۹,۸۳۱۹۴۹۸۸e + ۰۹$	$۹,۸۳e + ۰۹$	$۹,۸۳۱۹۴۹۸۸e + ۰۹$
۵	◦	۳۸۰۰۹۶۸۳۲۹/۷۱	$۳,۷۹۱۵۲۹۸۳e + ۰۹$	$۳,۸۰e + ۰۹$	$۳,۸۰۰۳۳۷۷۲e + ۰۹$
۱۰	◦	۱۷۶۱۷۸۳۰۳۴/۱۲۴	$۱,۷۵۴۸۵۸۰۶e + ۰۹$	$۱,۷۸e + ۰۹$	$۱,۷۵۴۸۵۸۰۶e + ۰۹$
۱۵	◦	۱۲۵۵۱۹۰۲۴۳/۲۷۸	$۱,۱۲۱۶۲۰۱۲e + ۰۹$	$۱,۱۴e + ۰۹$	$۱,۱۲۵۸۳۱۶۷e + ۰۹$
۲۰	◦	۸۸۰۹۰۶۱۴۱/۹۶۶	$۷,۹۲۰۹۶۹۸۵e + ۰۸$	$۸,۸۵e + ۰۸$	$۸,۰۳۱۱۴۵۶۱e + ۰۸$
۲	◦	۳۱۶۸۸۲۶۲۹۸/۴۲۳	$۳,۱۶۸۸۰۴۶۸e + ۰۹$	$۳,۱۷E + ۰۹$	$۳,۱۶۸۸۱۳۲۱e + ۰۹$
۵	◦	۱۲۰۱۱۲۳۱۲۶/۳۶۱	$۱,۱۹۸۲۰۲۶۹e + ۰۹$	$۱,۲۰E + ۰۹$	$۱,۱۹۸۲۰۸۶۰e + ۰۹$
۱۰	◦	۵۷۴۲۶۸۷۰۳/۳۳۷	$۵,۶۰۴۱۱۰۷۸e + ۰۸$	$۵,۸۳E + ۰۸$	$۵,۶۰۴۱۱۰۷۸e + ۰۸$
۱۱	-	-	-	-	-
۱۵	◦	۳۵۶۱۰۴۵۴۸/۰۵۷	$۳,۵۶۱۳۲۹۸۸e + ۰۸$	$۳,۶۴E + ۰۸$	$۳,۵۶۳۰۲۵۰۶e + ۰۸$
۲۰	◦	۲۷۹۸۰۳۲۹۳/۹۰۱	$۲,۶۷۰۲۱۸۳۵e + ۰۸$	$۲,۸۳E + ۰۸$	$۲,۶۷۳۷۸۲۲۷e + ۰۸$

Table :۳'۵ TSPLIB۱۰۶۰ and TSPLIB۳۰۳۸ data sets

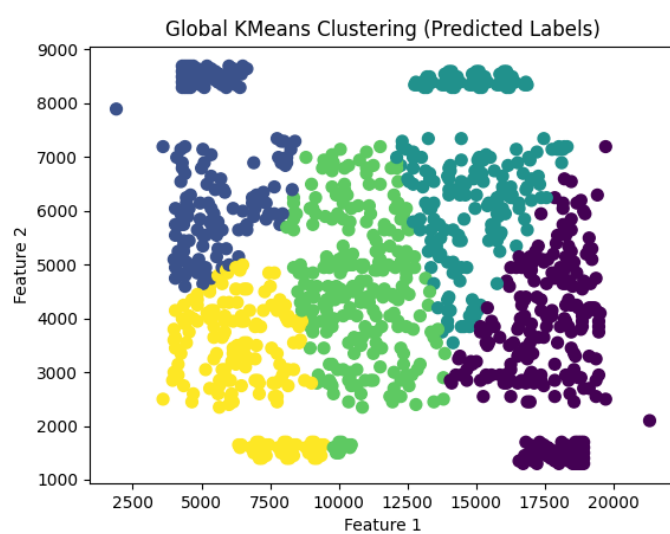
No.	$f_{opt}$	k-means	gkm	fast gkm	modified fast gkm	modified k-
		$\alpha$	$\alpha$	$\alpha$	$\alpha$	$\alpha$
۲	◦	۱۴۸۴۰	۱۳۱۴۷۱۸۰	۳۹۲۵۳۶	۲۶۸۶۳۷۶	۳۳۸۶۷۰۰
۵	◦	۵۳۰۰۰	۲۵۱۶۲۴۹۲۰	۸۵۲۹۸۹۱	۴۷۸۸۸۳۵۸	۲۷۱۴۹۷۸۰
۱۰	◦	۲۳۳۲۰۰	۹۲۶۹۳۵۰۲۰	۴۳۸۹۱۸۴۳	۲۲۳۷۴۰۷۳۴	۱۱۱۸۱۷۲۸۰
۱۵	◦	۴۲۹۳۰۰	۱۶۹۴۸۶۸۹۸۰	۱۰۷۲۷۰۱۵۹	۴۵۹۳۲۴۶۹۸	–
۲۰	◦	۴۸۷۶۰۰	۲۶۱۷۵۸۰۹۶۰	۱۹۷۶۸۱۶۶۲	۷۸۳۴۸۳۹۹۰	–
۲	◦	۷۲۹۱۲	۱۶۰۹۵۰۲۰۲	۲۷۶۸۵۴۴	۲۹۱۳۸۳۸۴	۲۷۷۳۹۹۷۸
۵	◦	۳۳۴۱۸۰	۱۷۷۳۴۳۵۵۳۸	۶۲۷۳۸۱۲۳	۳۵۳۷۰۸۰۰۹	۲۲۲۰۸۶۹۱۴
۱۰	◦	۱۳۰۶۳۴۰	۷۸۸۳۴۸۵۴۴۲	۳۴۶۱۸۱۵۱۶	۱۸۴۷۰۳۸۲۱۵	۹۱۶۱۴۸۳۹۴
۱۱	–	–	–	–	–	۱۱۱۰۳۰۰۸۹
۱۵	◦	۱۴۱۲۶۷۰	۱۸۳۱۹۱۵۸۲۲۸	۸۶۴۳۳۷۰۶۱	۴۲۷۶۲۹۴۷۵۸	–
۲۰	◦	۱۵۱۹۰۰۰	۳۱۱۸۸۱۷۱۷۹۸	۱۶۱۵۱۶۶۲۸۷	۷۹۶۵۸۴۰۸۹۰	–

Table :۳'۶ TSPLIB۱۰۶۰ and TSPLIB۳۰۳۸ data sets

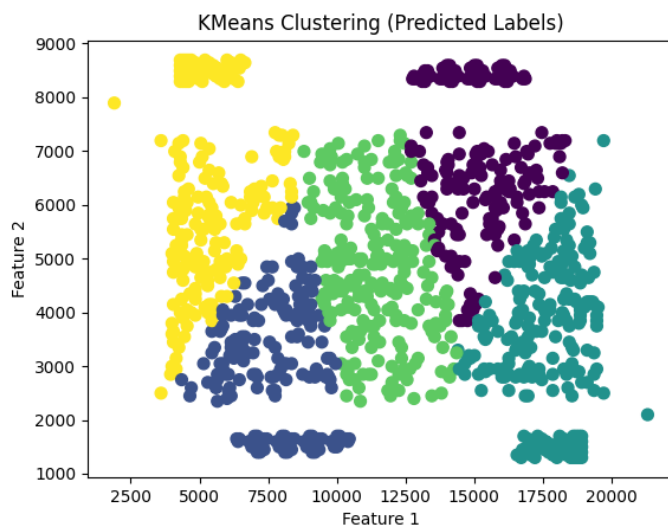
No.	$f_{opt}$	k-means	gkm	fast gkm	modified fast gkm	modified k-means
		$t$	$t$	$t$	$t$	$t$
۲	۰	۰/۰۳	۰/۷۱	۱۲/۷۶	۱۰/۶۳	۱۵/۳۷
۵	۰	۰/۰۴	۴/۳۹	۹۶/۳۹	۷۷/۲۳	۲۶/۰۰
۱۰	۰	۰/۰۴	۹/۹۹	۳۱۲/۴۹	۲۳۹/۱۸	۲۵۸/۴۵
۱۵	۰	۰/۰۴	۱۵/۲۰	۵۹۸/۱۶	۴۸۴/۴۱	—
۲۰	۰	۰/۰۴	۲۰/۷۰	۱۰۰۱/۷۱	۷۴۲/۵۳	—
۲	۰	۰/۰۴	۳/۱۳	۹۶/۳۲	۹۳/۵۴	۱۲۲/۹۷
۵	۰	۰/۰۴	۱۶/۲۴	۷۶۶/۶۴	۷۸۵/۳۴	۲۳۳/۹۹
۱۰	۰	۰/۰۴	۴۳/۵۶	۲۵۹۳/۲۴	۲۵۸۲/۴۹	۳۳۶/۴۰
۱۱	—	—	—	—	—	۳۲۱۱/۹۵
۱۵	۰	۰/۰۴	۸۲/۷۱	۵۶۵۸/۳۶	۴۱۰۱/۵۲	—
۲۰	۰	۰/۰۴	۱۱۸/۱۴	۸۲۵۰/۰۵	۶۵۵۸/۲۸	—



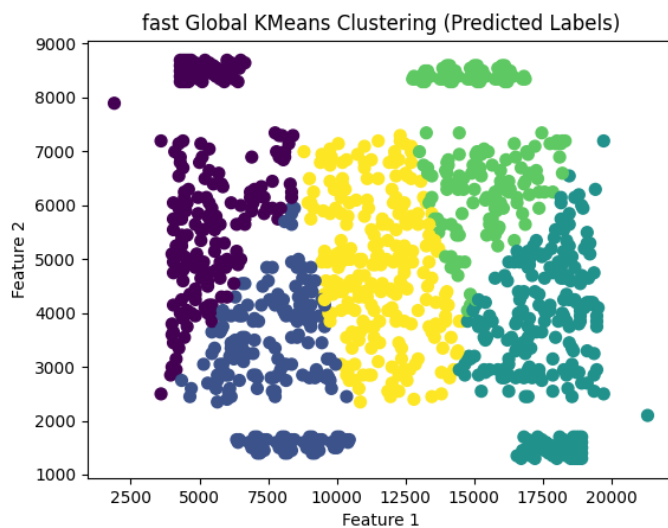
$$k = 5$$



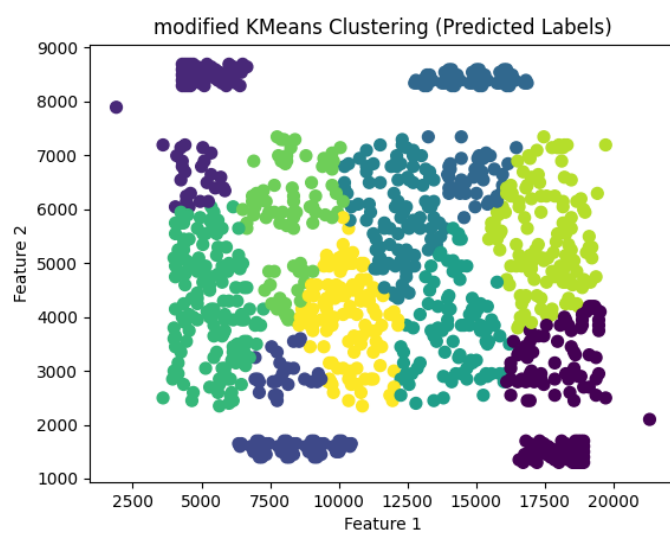
$$k = 5$$



$$k = 5$$



$$k = ۱۰$$



## مراجع

- [1] Bagirov, A., Hoseini-Monjezi, N., Taheri, S. (2023). A novel optimization approach towards improving separability of clusters. *Computer & Operations Research*, 152, 106135.
- [2] Ye, J., Zhao, Z., Liu, H. (2007, June). Adaptive distance metric learning for clustering. In *2007 IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1-7). IEEE.
- [3] Xing, E., Jordan, M., Russell, S. J., Ng, A. (2002). Distance metric learning with application to clustering with side-information. *Advances in neural information processing systems*, 15.
- [4] Bagirov, A. M. (2008). Modified global k-means algorithm for minimum sum-of-squares clustering problems. *Pattern Recognition*, 41(10), 3192-3199.

**Abstract.**

In this project ....

**Key words:** Clustering, k-means algorithm, modified k-means algorithm





University of Isfahan

Facility of Mathematics and Statistics

Departemant of Applied Mathematic and Computer Science

Bachelor Project

## **Modified Clustring Methods**

Supervisor:

Dr. Mohsen Alambardar Meybodi

Dr. Najmeh Hoseini Monjezi

By:

Zahra Meshkati

Feburay 2025