

نام و نام خانوادگی: حسنا اویارحسینی	تمرین سوم داده کاوی
شماره دانشجویی: ۹۸۲۳۰۱۰	تاریخ: خرداد - ۱۴۰۲

### سوال (۱)

از روش دسته بندی **partitioning** استفاده میکنیم، یعنی ابتدا فرض میکنیم کل حیوانات یک دسته باشند، سپس این دسته را به نحوی که مجموع فاصله هر داده از مرکز دسته متعلق به آن کمترین شود به دو دسته تقسیم میکنیم. و در نهایت نیز یکی از دو دسته ایجاد شده را به همان نحوه قبلی به دو دسته تقسیم میکنیم. به این صورت حیوانات دسته بندی می شوند و فاصله پارتیشن ها از هم نشان دهنده میزان شباهت می باشد یعنی برای مثال حیوانات متعلق به یک دسته فاصله ۱ دارند، و حیوانات متعلق به دو دسته متفاوت فاصله ۲ یا ۳ میتوانند داشته باشند (با توجه به دسته ها).

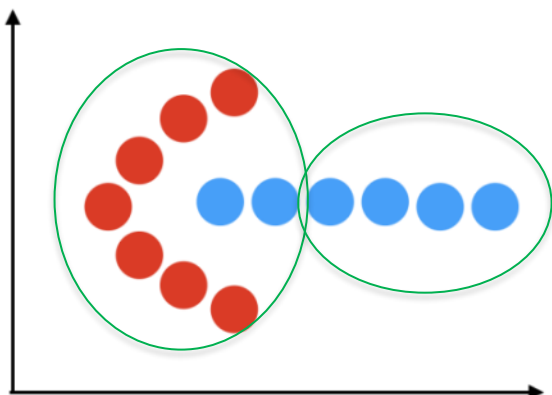
### سوال (۲)

- نرم ۱: میانه. هنگام مینیمم کردن عبارت داده شده با فاصله نرم ۱، میخواهیم مجموع تعدادی قدر مطلق را مینیمم کنیم که این قدر مطلق ها همان تعریف میانه هستند (نقطه ای که مجموع فاصله آن از بقیه نقاط کمتر است)
- نرم ۲: میانگین. ابتدا مشتق این تابع را محاسبه میکنیم و آن را صفر قرار میدهیم تا نقطه مینیمم را پیدا کنیم، خواهیم داشت:

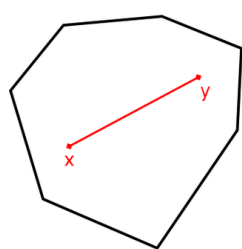
$$c = \frac{\sum d}{n} \rightarrow \text{به ازای یک دسته خواهیم داشت} \rightarrow \sum_d \sum_c 2||d - c|| = 0$$

### سوال ۳)

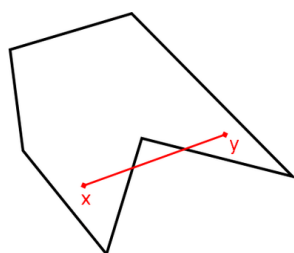
الف) با توجه به اینکه شکل توزیع داده ها  $\text{non-convex}$  می باشد الگوریتم  $\text{k-means}$  نمیتواند به درستی خوشه بندی را انجام دهد زیرا در این الگوریتم هدف آپدیت کردن مراکز خوشه به نحوی است که میانگین فاصله داده ها از مرکز خوشه مینیمم شود و در نتیجه خروجی الگوریتم چیزی همانند شکل زیر خواهد بود که بخشی از داده ها با لیبل آبی را به اشتباه با داده های قرمز در یک خوشه قرار میدهد.



ب) بله، به علت اینکه شکل توزیع داده ها  $\text{non-convex}$  می باشد و الگوریتم  $\text{DBSCAN}$  خوشه ها را بر اساس تراکم نقاط داده تعریف می کند و به آن اجازه می دهد خوشه هایی با مرزهای نامنظم را شناسایی کند. در مقابل،  $\text{k-means}$  خوشه ها را محدب فرض می کند و داده ها را بر اساس فواصل اقلیدسی جدا می کند، که باعث می شود برای خوشه های غیر محدب کارایی کمتری داشته باشد.



Convex region

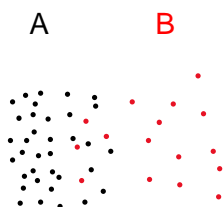


Non-convex region

$\text{DBSCAN}$  به طور خودکار آستانه چگالی مناسب را برای جداسازی خوشه ها از نویز تعیین می کند. برای تعیین مناطق متراکم، تعداد نقاط مجاور را در یک شعاع مشخص در نظر می گیرد. این ماهیت تطبیقی به  $\text{DBSCAN}$  اجازه می دهد تا خوشه هایی با چگالی های مختلف، از جمله خوشه های غیر محدب را کشف کند. در مقابل،  $\text{K-means}$  از کاربر می خواهد که تعداد خوشه ها را از قبل مشخص کند و تأثیر یکسانی از همه نقاط را در نظر می گیرد، که باعث می شود در مدیریت داده های غیر محدب با چگالی های متفاوت انعطاف پذیرتر نباشد.

(ج)

۱. حساسیت به انتخاب پارامتر چگالی: الگوریتم‌های خوشه‌بندی مبتنی بر چگالی نیازمند تعیین پارامترهایی مانند حداقل تعداد نقاط و شعاع همسایگی هستند. انتخاب مقادیر پارامتر مناسب می‌تواند در سناریوهایی با خوشه‌های همپوشانی چالش برانگیز باشد. اگر پارامتر چگالی خیلی زیاد تنظیم شود، منطقه همپوشانی ممکن است به عنوان یک خوشه جداگانه شناسایی نشود. برعکس، تنظیم پارامتر چگالی خیلی کم ممکن است باعث شود الگوریتم خوشه A و ناحیه همپوشانی را ادغام کند و منجر به نتایج خوشه‌بندی نادرست شود.



۲. چگالی متغیر: خوشه‌بندی مبتنی بر چگالی فرض می‌کند که خوشه‌ها چگالی مشابهی دارند. هنگامی که چگالی خوشه‌ها به طور قابل توجهی تغییر می‌کند، تعیین آستانه چگالی مناسب برای الگوریتم چالش برانگیز می‌شود. در چنین مواردی، خوشه‌هایی با چگالی کمتر ممکن است با خوشه‌های نزدیک با چگالی بالاتر ادغام شوند یا به عنوان نویز در نظر گرفته شوند که منجر به نتایج نادرست می‌شود.

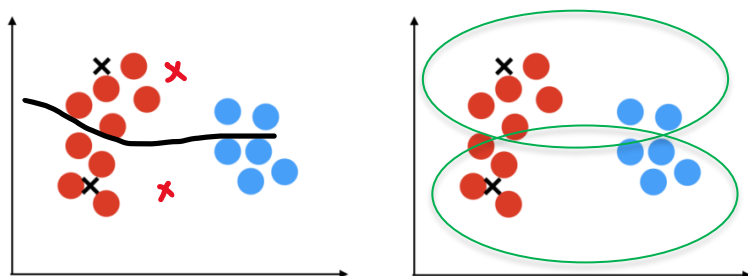
۳. داده‌های با ابعاد بالا: الگوریتم‌های خوشه‌بندی مبتنی بر چگالی در فضاهای با ابعاد بالا به دلیل «تفرین ابعاد» با مشکلاتی روبرو هستند. در داده‌های با ابعاد بالا، مفهوم فاصله کمتر معنادار می‌شود و چگالی نقاط تمایل به یکنواخت‌تر شدن دارد. در نتیجه، الگوریتم ممکن است برای شناسایی خوشه‌های معنی دار مشکل داشته باشد و ممکن است بیشتر نقاط را به عنوان نویز در نظر بگیرد.

۴. شکل‌ها و اندازه‌های خوشه‌ای متفاوت: خوشه‌بندی مبتنی بر چگالی فرض می‌کند که خوشه‌ها نواحی متراکمی هستند که توسط مناطق با چگالی کمتر از هم جدا شده‌اند. با این حال، اگر خوشه‌ها دارای اشکال نامنظم یا اندازه‌های قابل توجهی متفاوت باشند، برای الگوریتم چالش برانگیز است که آستانه‌های چگالی مناسب برای جداسازی خوشه‌ها را به طور موثر تعریف کند. ممکن است منجر به ادغام خوشه‌های کوچکتر به خوشه‌های بزرگتر یا تقسیم نادرست خوشه‌های بزرگ به خوشه‌های کوچکتر شود.

۵. نویز و نقاط پرت: الگوریتم‌های خوشه‌بندی مبتنی بر چگالی معمولاً نویز و نقاط پرت را با برجسب‌گذاری آنها به عنوان نقاط نویز به خوبی کنترل می‌کنند. با این حال، اگر مجموعه داده حاوی مقدار قابل توجهی نویز یا نقاط پرت باشد که به طور متراکم بسته بندی شده اند، ممکن است بر تعیین آستانه چگالی تأثیر بگذارد. الگوریتم ممکن است به اشتباه این نقاط نویز را به عنوان بخشی از یک خوشه در نظر بگیرد یا خوشه های واقعی را شناسایی نکند.

#### سوال (۴)

(الف)



نتیجه الگوریتم به صورت رو به رو خواهد بود که نادرست است:

زیرا مراکز اولیه نزدیک هم انتخاب شده اند و مقدار دهی اولیه نامناسب باعث می شود حتی با پایان رساندن الگوریتم هنوز نتوان به خوشه بندی مناسب دست یافت.

(ب)

#### ۱- استفاده از Medoid به جای Median:

در این روش به جای استفاده از میانگین (میانگین) نقاط داده به عنوان مرکز اولیه، از medoid استفاده می شود. Medoid نقطه داده ای در یک خوشه است که کمترین تفاوت میانگین را با سایر نقاط آن خوشه دارد.

• مزایا:

- Medoids نقاط داده واقعی هستند که اطمینان حاصل می کنند که مرکزهای اولیه نمایندگان معتبر داده ها هستند.

- مدوئیدها در مقایسه با استفاده از میانگین به عنوان مرکز نسبت به نقاط پرت مقاوم تر هستند.

• معایب:

- محاسبه medoidها به محاسبات عدم تشابه زوجی بین تمام نقاط داده نیاز دارد که می تواند از نظر محاسباتی گران باشد.

- شناسایی مناسب‌ترین مدویدها ممکن است چالش‌برانگیز باشد، به‌ویژه زمانی که با داده‌های با ابعاد بالا یا اشکال خوشه‌ای پیچیده سروکار داریم.

۲- انتخاب نقاط اولیه با بیشترین فاصله:

در این روش مرکزهای اولیه به گونه ای انتخاب می شوند که حداکثر فاصله زوجی را از یکدیگر داشته باشند. این تضمین می کند که centroid ها به خوبی در سراسر مجموعه داده توزیع شده اند.

• مزایا:

- انتخاب نقاط اولیه با حداکثر فاصله زوجی به جلوگیری از قرار دادن مرکزها در مجاورت نزدیک کمک می کند.

- احتمال همگرایی به یک راه حل خوب را افزایش می دهد.

• معایب:

- یافتن نقاطی با بیشترین فاصله نیاز به محاسبات فاصله زوجی بین تمام نقاط داده دارد که می تواند از نظر محاسباتی گران باشد.

- اگر مجموعه داده شامل نویز باشد تاثیر نویز بر روی الگوریتم ما در این حالت زیاد میشود و دچار مشکل میشویم.

۳- انتخاب نقاط اولیه بر اساس توزیع داده ها:

در این روش، مرکزهای اولیه بر اساس توزیع داده ها انتخاب می شوند. به عنوان مثال، می توانید مرکزهای اولیه را از مناطق با چگالی بالا یا بر اساس توزیع احتمال متناسب با داده ها انتخاب کنید.

• مزایا:

- این روش توزیع داده ها را در نظر می گیرد و می تواند به گرفتن ساختار یا حالت های اساسی داده ها کمک کند.

• معایب:

- شناسایی روش مبتنی بر توزیع مناسب برای انتخاب نقاط اولیه می تواند چالش برانگیز باشد و ممکن است به دانش حوزه یا تجزیه و تحلیل آماری نیاز داشته باشد.

- اثربخشی این روش به تناسب مدل توزیع انتخاب شده برای داده ها بستگی دارد.

- اگر توزیع داده‌ها نامنظم باشد یا خوشه‌ها اشکال یا اندازه‌های متفاوتی داشته باشند، این رویکرد ممکن است خوب کار نکند.

۴- انتخاب چندگانه مراکز اولیه:

در این روش، مجموعه های متعددی از مراکزهای اولیه به طور تصادفی انتخاب می شوند و الگوریتم k-means چندین بار اجرا می شود. بهترین نتیجه خوشه بندی از اجرای چندگانه انتخاب می شود.

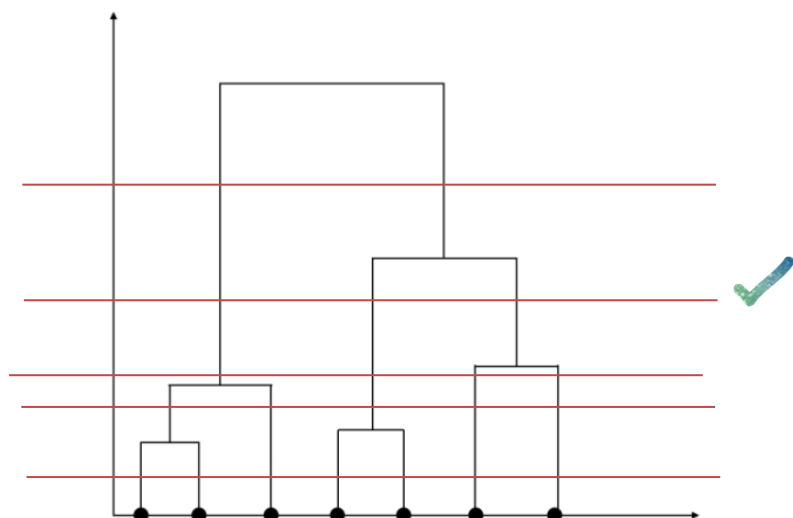
• مزایا:

- انتخاب چندین مجموعه از مراکز به کاهش مشکل گیر کردن در بهینه محلی کمک می کند.
- اجرای الگوریتم با مقداردهی اولیه، شانس یافتن راه حل بهینه را افزایش می دهد.

• معایب:

- این رویکرد هزینه محاسباتی را افزایش می دهد زیرا به اجرای چندین بار الگوریتم k-means نیاز دارد.
- تعیین بهترین نتیجه خوشه بندی در بین چند مقدار اولیه ممکن است به معیارهای ارزیابی اضافی یا قضاوت ذهنی نیاز داشته باشد.

ج) با توجه به دندوگرم داده شده تعداد ۳ خوشه برای خوشه بندی مناسب است پی در k-means نیز  $k = 3$  قرار میدهیم.



سوال ۵)

Standardization:

	Feature1	Feature2
	۱	۱
	۱	۲
	۲	۱

	۱-	۱-
	۱-	۲-
	۲-	۱-
Mean	.	.
variance	$\sqrt{۲}$	$\sqrt{۲}$

Covariance matrix:

$$\frac{\begin{bmatrix} ۱ & ۱ & ۲ & -۱ & -۱ & -۲ \\ ۱ & ۲ & ۱ & -۱ & -۲ & -۱ \end{bmatrix} \times \begin{bmatrix} ۱ & ۱ \\ ۱ & ۲ \\ ۲ & ۱ \\ -۱ & -۱ \\ -۱ & -۲ \\ -۲ & -۱ \end{bmatrix}}{(۶ - ۱)\sqrt{۲}\sqrt{۲}} = \frac{\begin{bmatrix} ۱۲ & ۱۰ \\ ۱۰ & ۱۲ \end{bmatrix}}{۱۰}$$

Eigen value:

$$\det|A - \lambda I| = ۰ \rightarrow \lambda = ۲.۲, ۰.۲$$

Eigen vectors:

$$\lambda = ۲.۲ \rightarrow A.v = \lambda.v \rightarrow x = y \rightarrow v = \begin{bmatrix} ۱ \\ ۱ \end{bmatrix} / \sqrt{۲}$$

Feature vector:

$$\begin{bmatrix} ۱ & ۱ \\ ۱ & ۲ \\ ۲ & ۱ \\ -۱ & -۱ \\ -۱ & -۲ \\ -۲ & -۱ \end{bmatrix} \times \begin{bmatrix} ۱ \\ ۱ \end{bmatrix} = \begin{bmatrix} ۲ \\ ۳ \\ ۳ \\ -۲ \\ -۳ \\ -۳ \end{bmatrix} / \sqrt{۲}$$

جزئیات محاسبات در ادامه آمده است:

1. Standardization:

							mean	Var
Feature 1	1	1	2	-1	-1	-2	0	$\sqrt{2}$
Feature 2	1	2	1	-1	-2	-1	0	$\sqrt{2}$

2. Covariance matrix:

Covariance matrix =  $\begin{bmatrix} 1 & 1 & 2 & -1 & -1 & -2 \\ 1 & 2 & 1 & -1 & -2 & -1 \end{bmatrix}$

3. (A)

4. Eigen values:

5.  $\det |A - \lambda I| = 0 \rightarrow \begin{vmatrix} 1.2 - \lambda & 1 \\ 1 & 1.2 - \lambda \end{vmatrix} = (1.2 - \lambda)^2 - 1$

6.  $0 = 1.44 + \lambda^2 - 2.4\lambda - 1$

7.  $\lambda^2 - 2.4\lambda + 0.44 = 0$

8.  $\Delta = 5.76 - 1.76 = 4$

9.  $\lambda = \frac{2.4 \pm 2}{2} = 2.2, 0.2$

10. Eigen vectors:

11.  $A \cdot v = \lambda \cdot v$

12. 1)  $\lambda = 2.2 \rightarrow \begin{bmatrix} 1.2 & 1 \\ 1 & 1.2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 2.2 \begin{bmatrix} x \\ y \end{bmatrix}$

13.  $\begin{cases} 1.2x + y = 2.2x \rightarrow y = x \\ x + 1.2y = 2.2y \rightarrow y = x \end{cases} \rightarrow (1, 1)$

14. 2)  $\lambda = 0.2 \rightarrow \begin{bmatrix} 1.2 & 1 \\ 1 & 1.2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0.2 \begin{bmatrix} x \\ y \end{bmatrix} \rightarrow 1.2x + y = 0.2x$

15.  $x + y = 0 \rightarrow (1, -1)$

16. Feature vec:  $\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$

سوال ٦

Frequent itemset:

Itemset - 1	freq	support	OK?
Br	5	.72	✓
M	5	.72	✓



D	6	0.185	✓
Be	4	0.57	✓
E	1	0.14	X
C	3	0.43	✓

Itemset - 2	freq	support	OK?
Br, M	3	0.43	✓
Br, D	4	0.57	✓
Br, Be	3	0.43	✓
Br, C	1	0.14	X
M, D	4	0.57	✓
M, Be	2	0.28	X
M, C	3	0.43	✓
D, Be	4	0.57	✓
D, C	3	0.43	✓
Be, C	1	0.14	X

Itemset - 3	freq	support	OK?
Br, M, D	2	0.28	X
Br, D, Be	3	0.43	✓
Br, D, C	1	0.14	X
M, D, Be	2	0.28	X
M, D, C	2	0.28	X
Be, D, C	1	0.14	X
Br, Be, M	1	0.14	X
M, C, Br	1	0.14	X

Rules:

Rule	confidence	OK?
Br $\rightarrow$ Be, D	$3/5 = 0.6$	✓
Be $\rightarrow$ Br, D	$3/4 = 0.75$	✓

$D \rightarrow Br, Be$	$3/5 = 0.6$	✓
$D, Be \rightarrow Br$	$3/4 = 0.75$	✓
$Br, D \rightarrow Be$	$3/4 = 0.75$	✓
$Be, Br \rightarrow D$	$3/3 = 1$	✓
$Br \rightarrow M$	$3/5$	✓
$M \rightarrow Br$	$3/5$	✓
$Br \rightarrow D$	$4/5$	✓
$D \rightarrow Br$	$4/6$	✓
$Br \rightarrow Be$	$3/5$	✓
$Be \rightarrow Br$	$3/4$	✓
$M \rightarrow D$	$4/5$	✓
$D \rightarrow M$	$4/6$	✓
$M \rightarrow C$	$3/5$	✓
$C \rightarrow M$	$3/3$	✓
$D \rightarrow Be$	$4/6$	✓
$Be \rightarrow D$	$4/4$	✓
$Be \rightarrow C$	$3/4$	✓
$C \rightarrow Be$	$3/3$	✓

Final rules:

همه قوانینی که در بالا آمده

سوال (۷)

Frequent itemset:

Itemset - 1	freq	support	OK?
A	۵	۰.۶۲۵	✓
B	۵	۰.۶۲۵	✓
C	۵	۰.۶۲۵	✓
D	۴	۰.۵	✓
E	۲	۰.۲۵	X

Itemset - 2	freq	support	OK?
-------------	------	---------	-----

A, B	3	0.375	✓
A, C	3	0.375	✓
A, D	2	0.25	X
B, C	4	0.5	✓
B, D	2	0.25	X
C, D	2	0.25	X

Itemset - 3	freq	support	OK?
A, B, C	3	0.375	✓

Rules:

Rule	confidence	OK?
$A \rightarrow B, C$	$3/5 = 0.6$	X
$B \rightarrow A, C$	$3/5 = 0.6$	X
$C \rightarrow A, B$	$3/5 = 0.6$	X
$A, B \rightarrow C$	$3/3 = 1$	✓
$A, C \rightarrow B$	$3/3 = 1$	✓
$B, C \rightarrow A$	$3/4 = 0.75$	✓
$A \rightarrow B$	$3/5$	X
$B \rightarrow A$	$3/5$	X
$A \rightarrow C$	$3/5$	X
$C \rightarrow A$	$3/5$	X
$B \rightarrow C$	$4/5$	✓
$C \rightarrow B$	$4/5$	✓

Final rules:

$A, B \rightarrow C$	$A, C \rightarrow B$	$B, C \rightarrow A$
$B \rightarrow C$	$C \rightarrow B$	