



## تمرین دوم

## پاسخ مسئله ۱.

الف) در ابتدا دقت کنید که برای محاسبه بردار میانگین، داریم:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

پس در اینجا خواهیم داشت:

$$\begin{aligned} \bar{x} &= \frac{1}{4} \sum_{i=1}^4 x_i \\ &= \frac{1}{4} \left( \begin{bmatrix} 4 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ 5 \end{bmatrix} + \begin{bmatrix} 3 \\ 4 \end{bmatrix} \right) \\ &= \frac{1}{4} \begin{bmatrix} 10 \\ 10 \end{bmatrix} = \begin{bmatrix} \frac{5}{2} \\ \frac{5}{2} \end{bmatrix} \end{aligned}$$

در ادامه دقت کنید که برای به دست آوردن ماتریس کوواریانس با استفاده از بردار میانگین، داریم:

$$\text{Cov}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

بنابراین در اینجا می توان نوشت:

$$\begin{aligned} \text{Cov}(X) &= \frac{1}{3} \sum_{i=1}^4 (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \frac{1}{3} \left( \begin{bmatrix} \frac{1}{2} \\ -\frac{3}{2} \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ -\frac{3}{2} \end{bmatrix}^T + \begin{bmatrix} -\frac{3}{2} \\ -\frac{5}{2} \end{bmatrix} \begin{bmatrix} -\frac{3}{2} \\ -\frac{5}{2} \end{bmatrix}^T + \begin{bmatrix} -\frac{1}{2} \\ \frac{5}{2} \end{bmatrix} \begin{bmatrix} -\frac{1}{2} \\ \frac{5}{2} \end{bmatrix}^T + \begin{bmatrix} \frac{1}{2} \\ \frac{3}{2} \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ \frac{3}{2} \end{bmatrix}^T \right) \\ &= \frac{1}{3} \begin{bmatrix} 5 & 1 \\ 1 & 17 \end{bmatrix} = \begin{bmatrix} \frac{5}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} \end{bmatrix} \end{aligned}$$

حال با حل معادله مشخصه، مقادیر ویژه ماتریس  $\text{Cov}(X)$  را می یابیم:

$$\begin{aligned}\det(\text{Cov}(X) - \lambda I) = 0 &\implies \det\left(\begin{bmatrix} \frac{5}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} \end{bmatrix} - \lambda I\right) = 0 \\ &\implies \det\left(\begin{bmatrix} \frac{5}{3} - \lambda & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} - \lambda \end{bmatrix}\right) = 0 \\ &\implies \left(\frac{5}{3} - \lambda\right)\left(\frac{17}{3} - \lambda\right) - \frac{1}{9} = 0 \\ &\implies (5 - 3\lambda)(17 - 3\lambda) - 1 = 0 \\ &\implies \lambda_1 = \frac{11 + \sqrt{37}}{3}, \quad \lambda_2 = \frac{11 - \sqrt{37}}{3}\end{aligned}$$

از آنجا که هدف انتقال داده ها به فضای یک بعدی است، در ابتدا مقدار ویژه بزرگ تر (یا همان  $\lambda_1$ ) را در نظر می گیریم و بردار ویژه متناظر به آن را پیدا می کنیم:

$$\text{Cov}(X)v = \lambda_1 v, \quad \|v\|_2 = 1$$

از تساوی اول نتیجه می شود:

$$\begin{aligned}\begin{bmatrix} \frac{5}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} \end{bmatrix} v = \left(\frac{11 + \sqrt{37}}{3}\right) v &\implies \begin{bmatrix} \frac{5}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \frac{11 + \sqrt{37}}{3} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \\ &\implies \begin{cases} 5v_1 + v_2 = (11 + \sqrt{37})v_1 \\ v_1 + 17v_2 = (11 + \sqrt{37})v_2 \end{cases} \\ &\implies v_2 = (6 + \sqrt{37})v_1\end{aligned}$$

و از ترکیب تساوی اخیر و تساوی دوم به دست می آید:

$$\begin{aligned}v_1^2 + v_2^2 = 1 &\implies v_1^2 + (6 + \sqrt{37})^2 v_1^2 = 1 \\ &\implies v_1^2 \left(1 + (6 + \sqrt{37})^2\right) = 1 \\ &\implies \begin{cases} v_1 = \frac{1}{\sqrt{1 + (6 + \sqrt{37})^2}} \\ v_2 = (6 + \sqrt{37})v_1 = \frac{6 + \sqrt{37}}{\sqrt{1 + (6 + \sqrt{37})^2}} \end{cases} \\ &\implies v = \begin{bmatrix} \frac{1}{\sqrt{1 + (6 + \sqrt{37})^2}} \\ \frac{6 + \sqrt{37}}{\sqrt{1 + (6 + \sqrt{37})^2}} \end{bmatrix}\end{aligned}$$

بنابراین نقاط انتقال داده شده به شکل زیر خواهند بود:

$$x'_i = x_i^T v$$

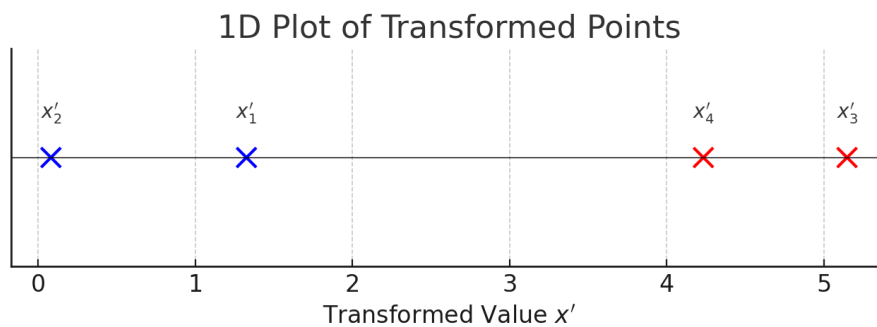
$$\Rightarrow x'_1 = \frac{10 + \sqrt{37}}{\sqrt{1 + (6 + \sqrt{37})^2}}$$

$$x'_2 = \frac{1}{\sqrt{1 + (6 + \sqrt{37})^2}}$$

$$x'_3 = \frac{32 + 5\sqrt{37}}{\sqrt{1 + (6 + \sqrt{37})^2}}$$

$$x'_4 = \frac{27 + 4\sqrt{37}}{\sqrt{1 + (6 + \sqrt{37})^2}}$$

حال نقاط به دست آمده را به تفکیک کلاس رسم می کنیم:



همانطور که مشاهده می شود، در این حالت کلاس ها به سادگی و به شکل خطی تمایز پذیرند. حال مقدار ویژه کوچک تر (یا همان  $\lambda_2$ ) را در نظر می گیریم و بردار ویژه متناظر به آن را می یابیم:

$$\text{Cov}(X)v = \lambda_2 v, \quad \|v\|_2 = 1$$

کاملاً مشابه حالت قبل، می توان نوشت:

$$\begin{aligned} \begin{bmatrix} \frac{5}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} \end{bmatrix} v &= \left( \frac{11 - \sqrt{37}}{3} \right) v \Rightarrow \begin{bmatrix} \frac{5}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{17}{3} \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \frac{11 - \sqrt{37}}{3} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \\ &\Rightarrow \begin{cases} 5v_1 + v_2 = (11 - \sqrt{37}) v_1 \\ v_1 + 17v_2 = (11 - \sqrt{37}) v_2 \end{cases} \\ &\Rightarrow v_2 = (6 - \sqrt{37}) v_1 \end{aligned}$$

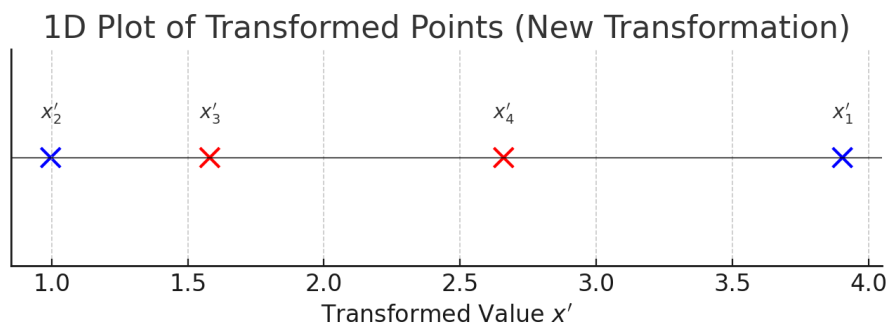
$$\begin{aligned}
v_1^2 + v_2^2 = 1 &\Rightarrow v_1^2 + (6 - \sqrt{37})^2 v_1^2 = 1 \\
&\Rightarrow v_1^2 \left( 1 + (6 - \sqrt{37})^2 \right) = 1 \\
&\Rightarrow \begin{cases} v_1 = \frac{1}{\sqrt{1 + (6 - \sqrt{37})^2}} \\ v_2 = (6 - \sqrt{37}) v_1 = \frac{6 - \sqrt{37}}{\sqrt{1 + (6 - \sqrt{37})^2}} \end{cases} \\
&\Rightarrow v = \begin{bmatrix} \frac{1}{\sqrt{1 + (6 - \sqrt{37})^2}} \\ \frac{6 - \sqrt{37}}{\sqrt{1 + (6 - \sqrt{37})^2}} \end{bmatrix}
\end{aligned}$$

بنابراین در این حالت، نقاط انتقال داده شده به شکل زیر خواهند بود:

$$x'_i = x_i^T v$$

$$\begin{aligned}
\Rightarrow x'_1 &= \frac{10 - \sqrt{37}}{\sqrt{1 + (6 - \sqrt{37})^2}} \\
x'_2 &= \frac{1}{\sqrt{1 + (6 - \sqrt{37})^2}} \\
x'_3 &= \frac{32 - 5\sqrt{37}}{\sqrt{1 + (6 - \sqrt{37})^2}} \\
x'_4 &= \frac{27 - 4\sqrt{37}}{\sqrt{1 + (6 - \sqrt{37})^2}}
\end{aligned}$$

حال نقاط به دست آمده را به تفکیک کلاس رسم می کنیم:



در این حالت، دیگر نقاط انتقال یافته به شکل خطی جدایی پذیر نیستند و تمایز کلاس ها نسبت به حالت قبل دشوار تر می باشد. به عبارت دیگر با انتخاب مقدار ویژه کوچک تر، تمایزپذیری کلاس ها کاهش یافت.

ب) نشان می دهیم که اگر بردار دلخواه  $x' \in F$  (در اینجا منظور از  $F$ ، فضای ویژگی ها یا همان Space Feature است) به همه داده ها اضافه شود، مولفه اساسی اول تغییری نخواهد کرد.

در ابتدا دقت کنید که با افزودن  $x'$  به همه داده ها، بردار میانگین جدید (که آن را با  $\bar{x}'$  نمایش می دهیم) برابر  $\bar{x} + x'$  می شود؛ چرا که می توان نوشت:

$$\begin{aligned}\bar{x}' &= \frac{1}{n} \sum_{i=1}^n x_i + x' \\ &= \frac{1}{n} \left( \sum_{i=1}^n x_i + nx' \right) \\ &= \frac{1}{n} (n\bar{x} + nx') \\ &= \bar{x} + x'\end{aligned}$$

در نتیجه اگر ماتریس کوواریانس داده های جدید را با  $\text{Cov}'(X')$  نمایش دهیم، خواهیم داشت:

$$\begin{aligned}\text{Cov}'(X') &= \frac{1}{n-1} \sum_{i=1}^n (x_i + x' - \bar{x}')(x_i + x' - \bar{x}')^T \\ &= \frac{1}{n-1} \sum_{i=1}^n (x_i + x' - \bar{x} - x')(x_i + x' - \bar{x} - x')^T \\ &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \text{Cov}(X)\end{aligned}$$

بنابراین از آنجا که  $\text{Cov}'(X')$  برابر  $\text{Cov}(X)$  می باشد، بردار های ویژه این دو ماتریس یکسان بوده و بردار ویژه متناظر به بزرگ ترین مقدار ویژه (مولفه اساسی اول برای دو مجموعه داده) نیز در آن ها یکسان خواهند بود. پس با افزودن  $x'$  به مجموعه داده اولیه، مولفه اساسی اول تغییر نمی کند و ادعای ما صحیح است.

## پاسخ مسئله‌ی ۲.

الف) در ابتدا دقت کنید که طبق تعریف، بایاس و واریانس مدل  $f$  که آن‌ها را به ترتیب با  $b(f)$  و  $\text{Var}(f)$  نمایش می‌دهیم، از روابط زیر به دست می‌آیند:

$$b(f) = \mathbb{E}[f - y], \quad \text{Var}(f) = \mathbb{E}[(f - \bar{f})^2]$$

به طوری که در عبارات بالا،  $y = y(x)$  مقدار واقعی برچسب برای داده  $x$  و  $\bar{f} = \mathbb{E}_X[f]$  مدل مورد انتظار (یا Expected Model) می‌باشد. حال دقت کنید که طبق فرض مسئله، داریم:

$$b(f_1) = b(f_2) = \dots = b(f_M) = b, \quad \text{Var}(f_1) = \text{Var}(f_2) = \dots = \text{Var}(f_M) = \sigma^2$$

با توجه به خاصیت خطی بودن امید ریاضی، می‌توان نوشت: (فارغ از مستقل بودن یا نبودن  $f_i$  ها)

$$\begin{aligned} b(f_{\text{ensemble}}) &= b\left(\frac{1}{M} \sum_{i=1}^M f_i\right) \\ &= \mathbb{E}\left[\frac{1}{M} \sum_{i=1}^M f_i - y\right] \\ &= \mathbb{E}\left[\frac{1}{M} \sum_{i=1}^M f_i - \frac{M}{M} y\right] \\ &= \mathbb{E}\left[\frac{1}{M} \sum_{i=1}^M (f_i - y)\right] \\ &= \frac{1}{M} \sum_{i=1}^M \mathbb{E}[f_i - y] \\ &= \frac{1}{M} \sum_{i=1}^M b = b \end{aligned}$$

در نتیجه بایاس  $f_{\text{ensemble}}$  با بایاس مدل‌های ضعیف برابر می‌باشد و با افزایش تعداد این مدل‌ها نیز مقدار  $b(f_{\text{ensemble}}) = b$  تغییری نخواهد کرد.

از طرف دیگر، با توجه به مستقل بودن مدل‌های ضعیف از یکدیگر، خواهیم داشت:

$$\begin{aligned} \text{Var}(f_{\text{ensemble}}) &= \text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i\right) \\ &= \frac{1}{M^2} \sum_{i=1}^M \text{Var}(f_i) \\ &= \frac{1}{M^2} (M\sigma^2) = \frac{\sigma^2}{M} \end{aligned}$$

بنابراین با افزایش مقدار  $M$  (تعداد مدل‌های ضعیف)،  $\text{Var}(f_{\text{ensemble}}) = \frac{\sigma^2}{M}$  کاهش خواهد یافت.

ب) در ابتدا دقت کنید که بایاس  $f_{\text{ensemble}}$  هیچگونه وابستگی به استقلال یا عدم استقلال مدل های ضعیف ندارد، چرا که خاصیت خطی بودن امید ریاضی، فارغ از مستقل بودن یا نبودن متغیر های تصادفی جمع شونده برقرار است. در نتیجه بایاس در این قسمت برابر بایاس در قسمت قبل می باشد:

$$b(f_{\text{ensemble}}) = b(f_i) = b$$

در نتیجه بایاس مدلی نهایی کاملاً مستقل از تعداد مدل های ضعیف ( $M$ ) و وابستگی بین آنها ( $\rho$ ) می باشد و با تغییر این پارامتر ها،  $b(f_{\text{ensemble}}) = b$  تغییری نخواهد کرد. اما برای واریانس  $f_{\text{ensemble}}$  در این حالت خواهیم داشت:

$$\begin{aligned} \text{Var}(f_{\text{ensemble}}) &= \text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i\right) \\ &= \frac{1}{M^2} \text{Var}\left(\sum_{i=1}^M f_i\right) \\ &= \frac{1}{M^2} \left( \sum_{i=1}^M \text{Var}(f_i) + 2 \sum_{i < j} \text{Cov}(f_i, f_j) \right) \\ &= \frac{1}{M^2} \left( \sum_{i=1}^M \sigma^2 + 2 \sum_{i < j} \rho \sigma^2 \right) \\ &= \frac{1}{M^2} (M \sigma^2) + \frac{2}{M^2} \binom{M}{2} \rho \sigma^2 \\ &= \frac{\sigma^2}{M} + \left(1 - \frac{1}{M}\right) \rho \sigma^2 = \frac{\sigma^2 - \rho \sigma^2}{M} + \rho \sigma^2 \end{aligned}$$

بنابراین با افزایش مقدار  $\rho$  (وابستگی بین مدل های ضعیف) با فرض ثابت بودن  $M$ ،  $\text{Var}(f_{\text{ensemble}})$  افزایش می یابد و با افزایش مقدار  $M$  (تعداد مدل های ضعیف) با فرض ثابت بودن  $\rho$ ، این کمیت کاهش خواهد یافت (دقت کنید که  $\sigma^2 - \rho \sigma^2 \geq 0$ ). البته توجه کنید که اگر  $\rho = 1$  باشد،  $\text{Var}(f_{\text{ensemble}})$  مستقل از  $M$  خواهد شد.

پ)

- خیر، یادگیرنده های ضعیف در AdaBoost نیاز به مشتق پذیر بودن ندارند. AdaBoost با تغییر وزن داده های نادرست طبقه بندی شده در هر تکرار کار می کند و به بهینه سازی مبتنی بر گرادیان متکی نیست. بنابراین، مدل های غیردیفرانسیبل، مانند درخت های تصمیم، می توانند به عنوان یادگیرنده های ضعیف مورد استفاده قرار گیرند.

- boosting به طور کلی از نظر محاسباتی گران تر از bagging است. در boosting، هر مدل جدید بر روی نسخه ای اصلاح شده ای از داده ها که در آن مثال های نادرست طبقه بندی شده وزن بیشتری دارند، آموزش می بیند، و این فرآیند به صورت ترتیبی انجام می شود. این وابستگی بین تکرارها پیچیدگی محاسباتی را افزایش می دهد. در مقابل bagging مدل ها را به صورت مستقل آموزش می دهد که امکان پردازش موازی را فراهم می کند.

### پاسخ مسئله‌ی ۳.

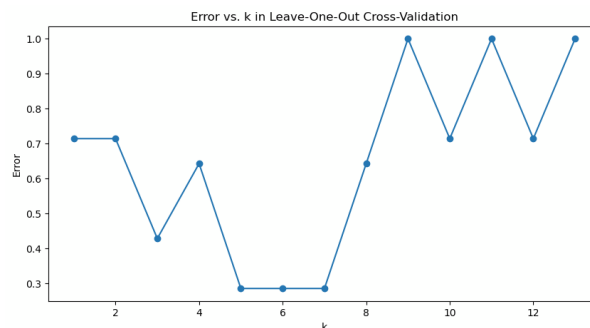
الف) به وضوح اگر فرض کنیم که هر داده همسایه خودش هم می باشد، با انتخاب  $k = 1$  خطا روی مجموعه داده آموزش برابر صفر خواهد شد؛ چرا که نزدیک ترین همسایه هر داده خود آن داده است و در نتیجه مقدار حقیقی برچسب هر داده توسط دسته بند 1-NN برای آن داده پیش بینی خواهد شد که به کمترین میزان خطا (خطای صفر) منجر می شود.

ب) اگر  $k$  خیلی کوچک باشد، دسته‌بند به نویزها حساس شده و واریانس آن زیاد می شود و در نتیجه ممکن است دچار بیش‌برازش (overfitting) شود. به طور مثال برای  $k = 1$ ، دسته‌بند تنها نزدیک‌ترین نقطه را در نظر می‌گیرد که این انتخاب بسیار نسبت به داده های نویزی و outlier ها حساس می باشد.

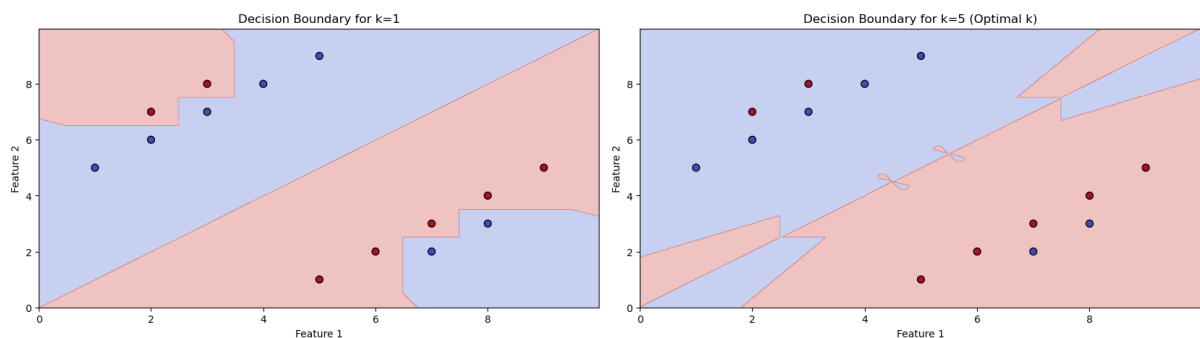
از طرفی اگر  $k$  خیلی بزرگ باشد، به طور مثال  $k = n$  (تعداد کل نقاط آموزشی)، همه‌ی داده‌ها برای طبقه‌بندی هر نقطه جدید استفاده می‌شوند و دسته‌بند به سمت کم‌برازش شدن (underfitting) می‌رود، زیرا تفاوت بین کلاس‌ها را نادیده می‌گیرد و به هر داده آزمون برچسب بیشینه بین داده های آموزشی را نسبت می دهد. بنابراین انتخاب مقدار بهینه برای  $k$  نیازمند تعادلی است که هم الگوهای محلی و هم کلی را در داده‌ها در نظر بگیرد.

پ و ت) در ابتدا فرض می کنیم که برچسب داده هایی که با علامت (-) مشخص شده اند، برابر صفر و برچسب داده هایی که با علامت (+) مشخص شده اند، برابر یک می باشد.

برای حل این قسمت از یک python notebook استفاده شده است که می توانید آن را در [این لینک](#) مشاهده کنید. بر اساس این برنامه، بهترین مقدار  $k$  با توجه به روش Leave One Out Cross Validation، برابر ۵ است که خطای آن برابر ۰/۲۸۶ می باشد. نمودار های خطا بر حسب مقدار  $k$  و مرز تصمیم گیری برای  $k = 1, 5$  در ادامه آورده شده اند.



شکل ۱: نمودار خطای validation بر حسب مقدار  $k$



شکل ۳: مرز تصمیم گیری برای  $k = 1$

شکل ۲: مرز تصمیم گیری برای  $k = 5$



## پاسخ مسئله‌ی ۴.

الف) جهت پیدا کردن مراکز خوشه‌ها، به گونه‌ای که تابع هزینه کمینه شود، از این تابع نسبت به هر کدام از آنها گرادیان می‌گیریم و تمامی این گرادیان‌ها را برابر صفر قرار می‌دهیم (دقت کنید که تابع هزینه محدب است و در نتیجه کمینه سراسری آن در نقطه‌ای رخ می‌دهد که گرادیان صفر شود). برای این منظور، در ابتدا تابع  $L$  را به شکل زیر بازنویسی می‌کنیم:

$$L = \sum_{j=1}^k \sum_{x_i \in S_j} \|x_i - \mu_j\|^2 = \sum_{j=1}^k \sum_{x_i \in S_j} (x_i - \mu_j)^T (x_i - \mu_j) = \sum_{j=1}^k \sum_{x_i \in S_j} x_i^T x_i - 2\mu_j^T x_i + \mu_j^T \mu_j$$

حال گرادیان تابع هزینه را نسبت به مرکز  $l$  امین دسته (برای  $1 \leq l \leq k$ ) محاسبه می‌کنیم و آن را برابر صفر قرار می‌دهیم:

$$\begin{aligned} \frac{\partial L}{\partial \mu_l} &= \frac{\partial \left( \sum_{j=1}^k \sum_{x_i \in S_j} x_i^T x_i - 2\mu_j^T x_i + \mu_j^T \mu_j \right)}{\partial \mu_l} \\ &= \frac{\partial \left( \sum_{x_i \in S_l} x_i^T x_i - 2\mu_l^T x_i + \mu_l^T \mu_l \right)}{\partial \mu_l} \\ &= \sum_{x_i \in S_l} \frac{\partial (x_i^T x_i - 2\mu_l^T x_i + \mu_l^T \mu_l)}{\partial \mu_l} \\ &= 2 \sum_{x_i \in S_l} \mu_l - x_i \\ &= 2|S_l|\mu_l - 2 \sum_{x_i \in S_l} x_i \\ \frac{\partial L}{\partial \mu_l} = 0 &\implies 2|S_l|\mu_l - 2 \sum_{x_i \in S_l} x_i = 0 \implies \mu_l = \frac{\sum_{x_i \in S_l} x_i}{|S_l|} \end{aligned}$$

و در نتیجه مقدار بهینه برای  $\mu_l$ ، میانگین داده‌های موجود در دسته  $l$  ام می‌باشد و حکم درست است.

ب) الگوریتم k-means نسبت به مقداردهی اولیه مراکز خوشه‌ها حساس می‌باشد. به عنوان مثال، فرض کنید داده‌هایی که می‌خواهیم آن‌ها را به  $k = 2$  خوشه تقسیم کنیم، به شکل زیر باشند:

$$x_1 = -2, \quad x_2 = -1, \quad x_3 = 1, \quad x_4 = 2$$

اگر مراکز خوشه‌ها را به صورت  $\mu_1 = 0$  و  $\mu_2 = 5$  مقداردهی اولیه کنیم، الگوریتم k-means پس از یک مرحله همگرا می‌شود و خوشه‌های زیر را به دست می‌آورد:

$$S_1 = \{x_1, x_2, x_3, x_4\}, \quad S_2 = \emptyset$$

اما اگر مراکز خوشه‌ها را به شکل  $\mu_1 = -1/5$  و  $\mu_2 = 1/5$  مقداردهی اولیه کنیم، الگوریتم k-means مجدداً پس از یک مرحله همگرا شده و خوشه‌های زیر را به دست می‌آورد:

$$S_1 = \{x_1, x_2\}, \quad S_2 = \{x_3, x_4\}$$

بنابراین مقداردهی اولیه متفاوت، می تواند خروجی نهایی k-means را به کل تغییر دهد و اگر این مقداردهی نامناسب باشد، می تواند منجر به همگرایی k-means به یک کمینه محلی شود. برای رفع این مشکل، راهکارهایی مانند الگوریتم **k-means++** معرفی شده اند که احتمال همگرایی به یک خوشه بندی بهتر را افزایش می دهند.

الگوریتم k-means حتما به یک پاسخ (نه لزوما پاسخ بهینه) همگرا خواهد شد.

برای اثبات این موضوع، در ابتدا به این نکته دقت کنید که در کل حداکثر  $k^n$  حالت مختلف برای خوشه بندی داده ها وجود دارد ( $n$  برابر تعداد داده ها است)، و در نتیجه  $k^n$  حالت برای مقادیر مراکز خوشه ها و مقدار تابع هزینه وجود خواهد داشت. از طرفی دقت کنید که در هر مرحله از الگوریتم k-means، در صورتی که الگوریتم همگرا و متوقف نشود، مقدار تابع هزینه حتما کاهش می یابد و مقدار جدیدی برای تابع هزینه (نسبت به تمامی مراحل قبل) حاصل می شود. در نتیجه پس از حداکثر  $k^n$  مرحله، الگوریتم متوقف خواهد شد (چرا که اگر پس از این تعداد مرحله متوقف نشود، حداقل  $k^n + 1$  حالت مختلف برای مقدار تابع هزینه مشاهده می شود که تناقض است).

برای نشان دادن اینکه مقدار تابع هزینه در هر مرحله از الگوریتم اکیدا کاهش می یابد، دقت کنید که هر مرحله شامل دو بخش زیر می باشد:

۱. انتقال هر داده به خوشه ای که کمترین فاصله را از مرکزش دارد.

به وضوح اگر حداقل یک داده در این مرحله به خوشه دیگری منتقل شود (که برای متوقف نشدن الگوریتم باید این اتفاق بیفتد)، مقدار تابع  $L$  اکیدا کاهش پیدا می کند.

۲. بروزرسانی مرکز خوشه ها به میانگین داده های موجود در آنها.

در این مرحله نیز، مطابق آنچه که در قسمت الف نشان داده شد، مقدار تابع هزینه کم می شود؛ چرا که اگر مرکز هر خوشه را میانگین آن در نظر بگیریم، کمترین میزان تابع هزینه را برای خوشه بندی ایجاد شده خواهیم داشت.

در نتیجه مقدار تابع  $L$  در هر مرحله اکیدا کاهش پیدا می کند و حکم درست می باشد.

پ) در الگوریتم k-means، انتخاب همان خوشه قبلی برای  $X_j$  (در صورت فاصله مساوی از چند مرکز خوشه) باعث پایداری بیشتر می شود. این تصمیم کمک می کند که الگوریتم از تغییرات غیرضروری جلوگیری کند و نقاط به صورت تصادفی بین خوشه ها جابه جا نشوند. این پایداری باعث می شود الگوریتم سریع تر به همگرایی برسد و تعداد تکرارها کاهش یابد.

اگر این اصل رعایت نشود و  $X_j$  به طور تصادفی به یکی از خوشه های دیگر اختصاص یابد، ممکن است الگوریتم در تکرارهای بعدی با مشکل نوسان مواجه شود؛ یعنی نقاط بین خوشه ها جابه جا شوند بدون اینکه همگرایی رخ دهد. این جابه جایی ها باعث افزایش تعداد تکرارها می شود و ممکن است الگوریتم در چرخه بی پایانی از جابه جایی ها گرفتار شود، به خصوص زمانی که چند نقطه به یک اندازه از چند مرکز خوشه فاصله داشته باشند. در نتیجه، الگوریتم نمی تواند به طور مؤثری واریانس درون خوشه ای را کاهش دهد و حتی ممکن است هیچ گاه همگرا نشود.

## پاسخ مسئله‌ی ۵.

در ابتدا مقدار تابع  $f_u(x)$  را برحسب  $u$  محاسبه می‌کنیم. از آنجا که نرم ۲ یک تابع محدب است، برای پیدا کردن کمینه‌کننده سراسری تابع  $f_u(x)$ ، کافی است گرادیان آن را نسبت به  $a$  محاسبه کرده و آن را برابر صفر قرار دهیم:

$$\begin{aligned} f_u(x) &= \arg \min_{a \in \mathbb{R}} \|x - au\|^2 \\ &= \arg \min_{a \in \mathbb{R}} (x - au)^T (x - au) \\ &= \arg \min_{a \in \mathbb{R}} x^T x - 2au^T x + a^2 u^T u \end{aligned}$$

$$\frac{\partial f_u(x)}{\partial a} = -2u^T x + 2au^T u$$

$$\begin{aligned} \frac{\partial f_u(x)}{\partial a} = 0 &\implies -2u^T x + 2au^T u = 0 \implies a = \frac{u^T x}{u^T u} \\ &\implies f_u(x) = au = \frac{u^T x u}{u^T u} \end{aligned}$$

در ادامه تابع MSE تعریف شده را، به شکل زیر بازنویسی می‌کنیم:

$$\begin{aligned} \text{MSE} &= \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - f_u(x^{(i)})\|^2 \\ &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - f_u(x^{(i)}))^T (x^{(i)} - f_u(x^{(i)})) \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - 2f_u(x^{(i)})^T x^{(i)} + f_u(x^{(i)})^T f_u(x^{(i)}) \end{aligned}$$

با توجه به فرض  $u^T u = 1$ ، خواهیم داشت  $f_u(x^{(i)}) = u^T x^{(i)} u$  و در نتیجه می‌توان نوشت: (دقت کنید که با توجه به خاصیت جابجایی برای ضرب داخلی،  $u^T x^{(i)} = x^{(i)T} u$  برقرار می‌باشد. همچنین چون  $u^T x^{(i)}$  یک عدد اسکالر است، می‌توان در ضرب با ماتریس و بردار آن را جابجا نمود)

$$\begin{aligned} \text{MSE} &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - 2u^T x^{(i)T} u x^{(i)} + u^T x^{(i)T} u u^T x^{(i)} u \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - 2u^T (u^T x^{(i)}) u x^{(i)} + u^T (u^T x^{(i)}) (u^T x^{(i)}) u \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - 2(u^T x^{(i)})^2 + (u^T x^{(i)})^2 u^T u \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - (u^T x^{(i)})^2 \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - (u^T x^{(i)})^2 \end{aligned}$$

دقت کنید که اگر تعریف کنیم

$$\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$$

می توان نوشت:

$$\begin{aligned} \text{MSE} &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - (u^T x^{(i)})^2 \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - u^T \left( \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right) u \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} - u^T \Sigma u \end{aligned}$$

حال توجه کنید که با توجه به فرض صورت سوال مبنی بر اینکه داده ها در هر بعد دارای میانگین صفر و واریانس یک هستند، از یک طرف برای ماتریس کوواریانس داده ها داریم:

$$\begin{aligned} \text{Cov}(X) &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \\ &= \Sigma \end{aligned}$$

و از طرف دیگر، با توجه به تساوی  $\text{trace}(A) = \text{trace}(A^T)$ ، می توان نوشت:

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m x^{(i)T} x^{(i)} &= \frac{1}{m} \text{trace} \left( \sum_{i=1}^m x^{(i)T} x^{(i)} \right) \\ &= \frac{1}{m} \text{trace} \left( \sum_{i=1}^m x^{(i)} x^{(i)T} \right) \\ &= \frac{1}{m} \text{trace}(\text{Cov}(X)) \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^d \text{Var}(x_j^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^d 1 \\ &= \frac{md}{m} \\ &= d \end{aligned}$$

به طوری که در تساوی های بالا،  $d$  همان بعد فضای ویژگی (Feature Space) می باشد.

بنابراین خواهیم داشت:

$$\text{MSE} = d - u^T \Sigma u$$

$$\implies \arg \min_{u: u^T u = 1} \text{MSE} = \arg \min_{u: u^T u = 1} d - u^T \Sigma u = \arg \max_{u: u^T u = 1} u^T \Sigma u$$

می دانیم که پاسخ مسئله بهینه سازی  $\arg \max_{u: u^T u = 1} u^T \Sigma u$  که معادل مسئله بهینه سازی اولیه است، مولفه اساسی اول داده ها می باشد (با توجه به آنچه که در کلاس ها و اسلاید های درسی بیان شده است) و در نتیجه حکم برقرار می باشد.