

به نام خدا

تمرین سری دوم
درس یادگیری ماشین
دکتر علی شریفی زارچی

فرزان رحمانی
۴۰۳۲۱۰۷۲۵

سوال اول
از این الگوریتم استفاده میکنیم:

Algorithm 1 Sample Covariance Matrix

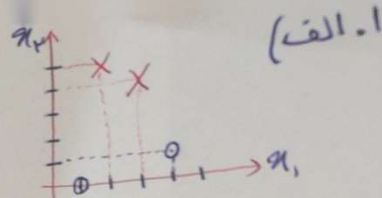
- 1: **Input:** $X \in \mathbb{R}^{N \times d}$ (data matrix with N data points and d dimensions)
 - 2: Compute the mean of each feature: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$
 - 3: Subtract the mean from each data point (center the data): $\tilde{X} \leftarrow X - \bar{x}^T$
 - 4: Compute the covariance matrix: $\Sigma = \frac{1}{N-1} \tilde{X}^T \tilde{X}$
 - 5: Compute the eigenvalues and eigenvectors of Σ : $[\lambda_1, \lambda_2, \dots, \lambda_d], [v_1, v_2, \dots, v_d] = \text{eig}(\Sigma)$
 - 6: Select the top k eigenvectors corresponding to the largest eigenvalues: $A \leftarrow [v_1, v_2, \dots, v_k]$
 - 7: Transform the data into the new subspace: $X' \leftarrow X \cdot A$
 - 8: **Output:** $X' \in \mathbb{R}^{N \times k}$ (transformed data with reduced dimensions)
-

ادامه جواب در صفحه بعد آمده است:

$$X = \begin{bmatrix} 4 & 1 \\ 1 & 0 \\ 2 & 5 \\ 3 & 4 \end{bmatrix}$$

ابتدا کلاً مجموعه داده را به صورت ماتریس می نویسیم.

حال برای محاسبه بردار میانگین داریم:



الف)

$$\mu = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{4} \sum_{i=1}^4 X_i = \frac{1}{4} ([4, 1] + [1, 0] + [2, 5] + [3, 4]) = \frac{1}{4} [10, 10] = [2.5, 2.5]$$

$$\mu = [2.5, 2.5] \xrightarrow{\text{حال با استفاده از میانگین داده ها را Center می کنیم}} \tilde{X} = X - \mu = \begin{bmatrix} 1.5 & -1.5 \\ -1.5 & -2.5 \\ -0.5 & 2.5 \\ 0.5 & 1.5 \end{bmatrix}$$

$$\Sigma = \frac{1}{n} \tilde{X}^T \tilde{X} = \frac{1}{4} \tilde{X}^T \tilde{X}$$

حال ماتریس Covariance را با استفاده از $n=4$ محاسبه می کنیم.

چون فرقی که در کلاً مجموعه داده را داریم.

$$\Sigma = \frac{1}{4} \begin{bmatrix} 1.5 & -1.5 & -0.5 & 0.5 \\ -1.5 & -2.5 & 2.5 & 1.5 \end{bmatrix} \begin{bmatrix} 1.5 & -1.5 \\ -1.5 & -2.5 \\ -0.5 & 2.5 \\ 0.5 & 1.5 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 5 & 1 \\ 1 & 17 \end{bmatrix} = \begin{bmatrix} 1.25 & 0.25 \\ 0.25 & 4.25 \end{bmatrix}$$

$$\det(\Sigma - \lambda I) = 0$$

$$\det \begin{bmatrix} 1.25 - \lambda & 0.25 \\ 0.25 & 4.25 - \lambda \end{bmatrix} = \begin{vmatrix} 1.25 - \lambda & 0.25 \\ 0.25 & 4.25 - \lambda \end{vmatrix} = (1.25 - \lambda)(4.25 - \lambda) - (0.25)^2 = 0$$

$$\lambda^2 - 5.5\lambda + 5.3125 - 0.0625 = 0$$

$$\lambda^2 - 5.5\lambda + 5.25 = 0$$

مقادیر ویژه $\lambda_1 \approx 4.2707$
 $\lambda_2 \approx 1.2293$

با کمک sklearn می توانیم معادله را حل کردیم.

$$\lambda_1 = 4.2707 \xrightarrow{\text{sklearn}} V_1 \approx [0.0825 \quad -0.9944]$$

$$\lambda_2 = 1.2293 \xrightarrow{\text{sklearn}} V_2 \approx [-0.9944 \quad -0.0825]$$

بردار
تای

ویژه

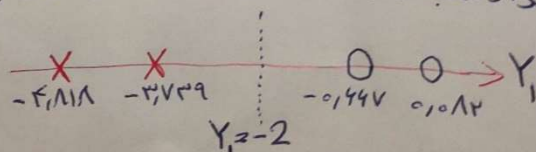
$$PC1 \rightarrow \lambda_1, V_1$$

$$PC2 \rightarrow \lambda_2, V_2$$

$$Y_1 = X \cdot V_1^T$$

$$Y_1 = \begin{bmatrix} 4 & 1 \\ 1 & 0 \\ 2 & 5 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 0.0825 \\ -0.9944 \end{bmatrix} = \begin{bmatrix} -0.9447 \\ 0.082 \\ -4.118 \\ -3.739 \end{bmatrix}$$

ابتدا داده ها را با $PC1$ project می کنیم.



همان طور که می بینیم کلاس ها W_1 و W_2 تمایز پذیری

هستند و در فضای Y_1 (همان $PC1$) می توانیم با خطی

نظیر $Y_1 = -2$ این دو کلاس را جدا کنیم (linearly separable).

حال داده‌ها را با $PC2$ (V_2) Project می‌کنیم.

$$Y_2 = X V_2^T$$

$$Y_2 = \begin{bmatrix} 4 & 1 \\ 1 & 0 \\ 2 & 5 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} -0,9944 \\ -0,0825 \end{bmatrix} = \begin{bmatrix} -4,069 \\ -0,997 \\ -2,404 \\ -3,320 \end{bmatrix}$$

همان طور که می‌بینیم کلاس‌های W_1 و W_2 در فضای Y_2 تمایز پذیری نسبت به و در فضای $PC2$ نمی‌توانیم یک خط این دو کلاس را از هم جدا کنیم. (not linearly separable)

ب) مؤلفه اساسی اول به وسیله V_1 تعریف می‌شود. برای اینکه مؤلفه اساسی اول ($PC1$) پس از جمع کردن یک بردار نظیر C تغییر نکند باید بردار جمع شونده C بر بردار V_1 عمود باشد.

ضرب داخلی آنها باید صفر شود. $C^T V_1 = 0$ (orthogonal)

$$C = [C_1 \ C_2]$$

$$V_1 = [0,0825 \ -0,9944]^T \rightarrow 0,0825 C_1 - 0,9944 C_2 = 0$$

یک جواب برای C می‌تواند $C_1 = 0,9944$ و $C_2 = 0,0825$ یا هر ضریب scalar از این بردار باشد.

جواب‌های دیگر $C = k [0,9944, 0,0825]$ (ضریب scalar)

بعد از جمع C قبل از جمع: $[x_1 \ x_2] [V_1 \ V_2]^T = x_1 V_1 + x_2 V_2$

بعد از جمع: $[x_1 + C_1 \ x_2 + C_2] [V_1 \ V_2]^T = x_1 V_1 + C_1 V_1 + x_2 V_2 + C_2 V_2$

معادله $0 = C_1 V_1 + C_2 V_2$ (چون $C^T V_1 = 0$)

برای اینکه first principal component بدون تغییر بماند، برداری که اضافه می‌کنیم باید متعامد به مؤلفه اصلی اول باشد. دلیل آن این است که وقتی نقاط داده را بر روی $PC1$ project می‌کنیم، projection یک عملیات خطی است. افزودن بردار متعامد به $PC1$ ، نمایش نقاط داده را به $PC1$ تغییر نمی‌دهد، زیرا حاصل ضرب نقطه ای $PC1$ و بردار متعامد صفر خواهد بود. افزودن این بردار به تمام نقاط داده، جهت اولین جزء اصلی را تغییر نمی‌دهد.

سوال دوم

الف) مدل‌های ضعیف مستقل

- M مدل ضعیف داریم: $(f_1(X), f_2(X), \dots, f_M(X))$
- ensemble model

$$f_{\text{ensemble}}(X) = \frac{1}{M} \sum_{i=1}^M f_i(X)$$

- هر مدل ضعیف دارای سوگیری (b) و واریانس (v) است.
- مدل‌های ضعیف مستقل از یکدیگر هستند.

هدف: بایاس و واریانس $f_{\text{ensemble}}(X)$ را محاسبه کنیم.

محاسبه بایاس

بایاس یک مدل به عنوان تفاوت بین پیش‌بینی مورد انتظار مدل و خروجی واقعی (Y) تعریف می‌شود:

$$\text{Bias}(f_{\text{ensemble}}(X)) = E[f_{\text{ensemble}}(X)] - Y$$

بیاپید پیش‌بینی مورد انتظار $f_{\text{ensemble}}(X)$ را محاسبه کنیم:

$$E[f_{\text{ensemble}}(X)] = E\left[\frac{1}{M} \sum_{i=1}^M f_i(X)\right]$$

از آنجایی که امید ریاضی یک عملگر خطی است:

$$E[f_{\text{ensemble}}(X)] = \frac{1}{M} \sum_{i=1}^M E[f_i(X)]$$

برای هر مدل ضعیف ($f_i(X)$)، پیش‌بینی مورد انتظار این است:

$$E[f_i(X)] = Y + b$$

بدین ترتیب:

$$E[f_{\text{ensemble}}(X)] = \frac{1}{M} \sum_{i=1}^M (Y + b) = \frac{1}{M} \cdot M \cdot (Y + b) = Y + b$$

بنابراین، بایاس ensemble model این است:

$$\text{Bias}(f_{\text{ensemble}}(X)) = (Y + b) - Y = b$$

نتیجه‌گیری: بایاس ($f_{\text{ensemble}}(X)$) برابر با بایاس مدل‌های ضعیف (b) است.

محاسبه واریانس

واریانس ($f_{\text{ensemble}}(X)$) به صورت زیر تعریف می‌شود:

$$\text{Var}(f_{\text{ensemble}}(X)) = E[(f_{\text{ensemble}}(X) - E[f_{\text{ensemble}}(X)])^2]$$

جایگزین کردن ($f_{\text{ensemble}}(X)$):

$$\text{Var}(f_{\text{ensemble}}(X)) = \text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i(X)\right)$$

استفاده از ویژگی واریانس برای متغیرهای مستقل:

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{1}{M^2} \sum_{i=1}^M \text{Var}(f_i(X))$$

با توجه به اینکه واریانس هر مدل ضعیف (v) است:

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{1}{M^2} \sum_{i=1}^M v = \frac{Mv}{M^2} = \frac{v}{M}$$

نتیجه گیری:

- بایاس ensemble model (b) است.
- واریانس ensemble model ($\frac{v}{M}$) است.

اثر افزایش (M):

- با افزایش (M) بایاس بدون تغییر باقی می ماند.
- با افزایش (M) ، واریانس کاهش می یابد زیرا ($\frac{v}{M}$) کوچکتر می شود.

ب) مدل های ضعیف همبسته

- M مدل ضعیف داریم: $(f_1(X), f_2(X), \dots, f_M(X))$
- ensemble model:

$$f_{\text{ensemble}}(X) = \frac{1}{M} \sum_{i=1}^M f_i(X)$$

- هر مدل ضعیف دارای سوگیری (b) و واریانس (v) است.
 - مدل های ضعیف یک همبستگی (ρ) بین هر جفت ($f_i(X)$) و ($f_j(X)$) که در آن ($i \neq j$) است.
- می خواهیم بایاس و واریانس ensemble model را محاسبه کنیم و تأثیر همبستگی (ρ) را روی این معیارها ببینیم.

مرحله ۱: محاسبه بایاس

بایاس ensemble model دقیقاً مانند قسمت (الف) است زیرا سوگیری به همبستگی بین مدل ها بستگی ندارد.

برای هر مدل ضعیف:

$$E[f_i(X)] = Y + b$$

برای ensemble model:

$$E[f_{\text{ensemble}}(X)] = E\left[\frac{1}{M} \sum_{i=1}^M f_i(X)\right]$$

استفاده از خطی بودن امید ریاضی:

$$E[f_{\text{ensemble}}(X)] = \frac{1}{M} \sum_{i=1}^M E[f_i(X)] = \frac{1}{M} \sum_{i=1}^M (Y + b) = Y + b$$

بایاس میشود:

$$\text{Bias}(f_{\text{ensemble}}(X)) = E[f_{\text{ensemble}}(X)] - Y = (Y + b) - Y = b$$

مرحله ۲: محاسبه واریانس با همبستگی

برای محاسبه واریانس $(f_{\text{ensemble}}(X))$ با تعریف آن شروع می کنیم:

$$\text{Var}(f_{\text{ensemble}}(X)) = \text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i(X)\right)$$

فرمول واریانس را گسترش میدهیم:

برای هر متغیر تصادفی (X) و (Y) :

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \cdot \text{Cov}(X, Y)$$

با استفاده از این، می توانیم واریانس مجموع را گسترش دهیم:

$$\text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i(X)\right) = \frac{1}{M^2} \text{Var}\left(\sum_{i=1}^M f_i(X)\right)$$

$$\text{Var}\left(\sum_{i=1}^M f_i(X)\right) = \sum_{i=1}^M \text{Var}(f_i(X)) + \sum_{i \neq j} \text{Cov}(f_i(X), f_j(X))$$

مقادیر واریانس و کوواریانس را جایگزین میکنیم.

$$\begin{aligned} - \quad & \text{Var}(f_i(X)) = v \quad \text{برای هر } (i) \\ - \quad & \text{Cov}(f_i(X), f_j(X)) = \rho \cdot \sqrt{\text{Var}(X) \cdot \text{Var}(Y)} = \rho \cdot \sqrt{v \cdot v} = \rho v \quad \text{برای } (i \neq j) \end{aligned}$$

اکنون هر قسمت را محاسبه می کنیم:

۱. مجموع واریانس ها:

$$\sum_{i=1}^M \text{Var}(f_i(X)) = \sum_{i=1}^M v = Mv$$

۲. مجموع کوواریانس ها:

$$\begin{aligned} - \quad & (M(M-1)) \text{ جفت های متمایز } ((i, j)) \text{ وجود دارد که در آن } (i \neq j) \\ - \quad & \sum_{i \neq j} \text{Cov}(f_i(X), f_j(X)) = \rho v \cdot M(M-1) \end{aligned}$$

مجموع واریانس را محاسبه میکنیم.

$$\text{Var}\left(\sum_{i=1}^M f_i(X)\right) = Mv + \rho v \cdot M(M-1)$$

اکنون، بر (M^2) تقسیم کنید تا واریانس $f_{\text{ensemble}}(X)$ را بدست آوریم:

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{Mv + \rho v \cdot M(M-1)}{M^2}$$

عبارت را ساده میکنیم:

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{v}{M} + \rho v \cdot \frac{M(M-1)}{M^2}$$

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{v}{M} + \rho v \cdot \frac{M-1}{M}$$

نتیجه گیری

- بایس ensemble model (b) است.
- واریانس مدل مجموعه عبارت است از:

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{v}{M} + \rho v \cdot \frac{M-1}{M}$$

تجزیه و تحلیل: اثر افزایش (M) و همبستگی (ρ)

۱. اثر افزایش (M) :

- اولین عبارت $(\frac{v}{M})$ با افزایش (M) کاهش می یابد.
- ترم دوم $(\rho v \cdot \frac{M-1}{M})$ با افزایش (M) به (ρ v) نزدیک می شود.
- زمانی که $(M \rightarrow \infty)$:

$$\text{Var}(f_{\text{ensemble}}(X)) \rightarrow \rho v$$

- بنابراین، اثر کاهش واریانس با افزایش تعداد مدل های ضعیف (M) در صورت وجود همبستگی (ρ) بالا کاهش می یابد.

۲. اثر همبستگی (ρ) :

- اگر (ρ = 0) (مدل های مستقل)، واریانس $(\frac{v}{M})$ است که با افزایش (M) کاهش می یابد.
- اگر (ρ > 0)، واریانس به دلیل عبارت اضافی $(\rho v \cdot \frac{M-1}{M})$ افزایش می یابد.
- همبستگی بالاتر (ρ) منجر به کاهش کمتری در واریانس با افزایش (M) می شود.

خلاصه:

- برای مدل های مستقل (ρ = 0)، واریانس به صورت $\frac{v}{M}$ کاهش می یابد.
- برای مدل های همبسته (ρ > 0)، کاهش واریانس کمتر موثر است زیرا عبارت $(\rho v \cdot \frac{M-1}{M})$ قابل توجه است.
- با افزایش (ρ)، مزیت افزایش (M) کاهش می یابد و واریانس به (ρ v) برای (M) بزرگ نزدیک می شود.

(ب)

- نه، یادگیرنده های ضعیف در AdaBoost نیازی به مشتق پذیر بودن ندارند. AdaBoost یک الگوریتم یادگیری گروهی است که چندین دسته بند ضعیف (یا «یادگیرندگان ضعیف») را برای ایجاد یک طبقه بندی قوی ترکیب می کند. یادگیرندگان ضعیف می توانند به هر شکلی باشند، به شرطی که بهتر از حدس زدن تصادفی عمل کنند. مشتق پذیری یک الزام برای AdaBoost نیست، زیرا الگوریتم به جای بهینه سازی یک تابع ضرر مشتق پذیر، بر تنظیم مکرر وزن نمونه های آموزشی تمرکز دارد. در واقع، AdaBoost با تنظیم وزن نمونه های طبقه بندی اشتباه کار می کند و به بهینه سازی مبتنی بر گرادینان متکی نیست (برخلاف روش های gradient boosting)، بنابراین نیازی به مشتق پذیری نیست.

- Boosting معمولاً از نظر محاسباتی گران‌تر از bagging است. دلایل:
 - آموزش متوالی: Boosting یادگیرندگان ضعیف را به طور متوالی آموزش می‌دهد و هر مدل جدید بر روی خطاهای قبلی تمرکز می‌کند. این فرآیند iterative است و نمی‌توان آن را به طور موثر موازی کرد. این فرآیند مکرر متوالی می‌تواند از نظر محاسباتی گران باشد، به ویژه با افزایش تعداد یادگیرندگان ضعیف.
 - وزن دهی مجدد داده‌ها: Boosting وزن نمونه‌ها را در هر تکرار تنظیم می‌کند و در هر مرحله محاسبات اضافی را اضافه می‌کند.
 - در مقابل، Bagging (به عنوان مثال، Random Forest) چندین مدل را به طور مستقل بر روی bootstrap samples های مختلف آموزش می‌دهد، که امکان آموزش موازی را فراهم می‌کند. در واقع، یادگیرندگان ضعیف در Bagging به طور مستقل آموزش می‌بینند و پیش‌بینی نهایی با جمع‌آوری خروجی‌های تک‌تک یادگیرندگان (مثلاً از طریق رأی اکثریت یا میانگین‌گیری) به دست می‌آید. که این کار که می‌تواند موازی هم باشد از نظر محاسباتی کارآمدتر است.
 - هزینه محاسباتی در boosting با تعداد یادگیرندگان ضعیف scale میشود، زیرا هر یادگیرنده جدید باید در وزن‌های به روز شده نمونه‌های آموزشی آموزش ببیند. از سوی دیگر، Bagging می‌تواند از محاسبات موازی برای آموزش همزمان یادگیرندگان ضعیف استفاده کند، که باعث می‌شود هزینه محاسباتی نسبتاً کمتری داشته باشد.

سوال سوم

در این سوال فرض میکنیم هر نقطه همسایه خودش هم هست. (در بحث‌های کوئرا بیان شد.)

الف) $k=1$ و مقدار خطا $= 0$ است.

تعیین مقدار k که خطا را به حداقل می‌رساند:

- برای $k=1$: طبقه‌بندی‌کننده KNN تنها نزدیک‌ترین همسایه را برای طبقه‌بندی هر نقطه در نظر می‌گیرد. از آنجایی که هر نقطه نزدیک‌ترین همسایه خود است، خطای طبقه‌بندی ۰ خواهد بود زیرا هر نقطه بر اساس خودش به درستی طبقه‌بندی می‌شود. (فاصله $L2$ هر نقطه از خودش ۰ است و ما هم فقط نزدیک‌ترین همسایه را که خود نقطه هست در نظر می‌گیریم.)
- برای $k>1$: با افزایش k ، طبقه‌بندی‌کننده همسایگان بیشتری را در نظر می‌گیرد. اگر نقاط همسایه متعلق به کلاس‌های مختلف باشند، ممکن است منجر به طبقه‌بندی اشتباه و افزایش شود.

در حالت کلی، برای یافتن مقدار دقیق k که خطا را به حداقل می‌رساند، باید مقادیر k مختلف را روی مجموعه داده cross validation آزمایش کنیم. در این سوال، زمانی که k از ۱ بالاتر می‌رود، شروع به مشاهده خطاهای طبقه‌بندی می‌کنیم.

برای مجموعه داده‌های کوچکی مانند این سوال، $k=1$ معمولاً خطا را به حداقل می‌رساند زیرا:

- هر نقطه همسایه خود است، بنابراین به درستی طبقه‌بندی شده است.
- مقدار خطا (نسبت دسته‌بندی‌های نادرست به کل دسته‌بندی‌ها) در $k=1$ برابر با $0/14=0$ است (۱۴ نقطه داده داریم).

ب) مقادیر k بسیار کوچک:

- طبقه‌بندی‌کننده به نویز و outliers های موجود در دیتاست بسیار حساس می‌شود زیرا فقط نزدیک‌ترین همسایه را در نظر می‌گیرد. که این نکته منجر به poor generalization to new data می‌شود.
- ممکن است منجر به overfitting و high variance شود زیرا مرز تصمیم می‌تواند بسیار پیچیده شود و با نقاط داده individual سازگار شود و نویزها را نیز در برگیرد.

مقادیر k بزرگ:

- طبقه‌بندی کننده همسایه‌های زیادی را در نظر می‌گیرد، احتمالاً شامل نقاطی از کلاس‌های مختلف است که می‌تواند منجر به طبقه‌بندی اشتباه شود.
- به عنوان مثال، اگر k نزدیک به تعداد کل نقاط باشد، تصمیم دسته بند به سمت کلاس اکثریت (majority class) در کل مجموعه داده متمایل می‌شود و الگوهای محلی را نادیده می‌گیرد و منجر به underfitting و low variance می‌شود.
- دسته بند بسیار ساده میشود و نمی‌تواند الگوهای اساسی در داده‌ها را ثبت کند و منجر به عملکرد ضعیف می‌شود.

پ) $k=5$ مقدار خطای دسته بند را کمینه میکند. مقدار این خطا $0.2857142857142857 = 4/14$ می‌شود.

Optimal k : 5, Error: 0.2857142857142857

Leave One Out Cross-Validation (LOOCV) روشی است که در آن هر نمونه یک بار به عنوان نمونه آزمایشی استفاده می‌شود در حالی که بقیه مجموعه داده به عنوان مجموعه آموزشی عمل می‌کند.

Leave-one-out cross validation is K-fold cross validation taken to its logical extreme, with K equal to N, the number of data points in the set. That means that N separate times, the function approximator is trained on all the data except for one point and a prediction is made for that point.

(<https://www.cs.cmu.edu/~schneide/tut5/node42.html>)

k بهینه را می‌توان با استفاده از کتابخانه ای مانند scikit-learn در پایتون پیدا کرد. این کار شامل fit کردن مدل KNN برای مقادیر مختلف k و بررسی اینکه کدام k خطای طبقه بندی را به حداقل می‌رساند.

این کد به شما k بهینه را با استفاده از LOOCV می‌دهد:

```
from sklearn.model_selection import LeaveOneOut
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
import numpy as np

# Example input data (x, y) and labels (Replace with actual data if available)
X = np.array([[5, 1], [6, 2], [7, 2], [7, 3], [8, 3], [8, 4], [9, 5],
              [1, 5], [2, 6], [2, 7], [3, 7], [3, 8], [4, 8], [5, 9]]) #
Coordinates from the image
y = np.array([1, 1, 0, 1, 0, 1, 1,
              0, 0, 1, 0, 1, 0, 0]) # Labels (1 for positive, 0 for negative)

loo = LeaveOneOut()
best_k = 1
best_score = 0

# Trying out different k values
for k in range(1, len(X)):
    # knn = KNeighborsClassifier(n_neighbors=k)
    knn = KNeighborsClassifier(n_neighbors=k, metric='euclidean')
    scores = []
    # Perform LOOCV
    for train_index, test_index in loo.split(X):
        X_train, X_test = X[train_index], X[test_index]
        y_train, y_test = y[train_index], y[test_index]
        knn.fit(X_train, y_train)
```

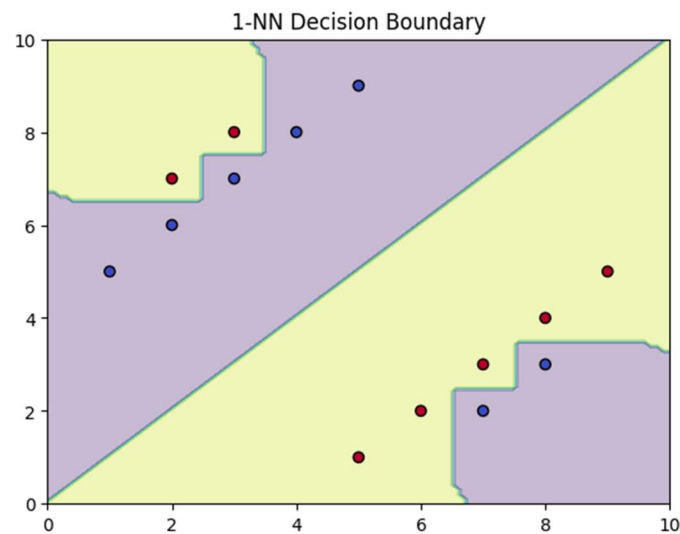
```

y_pred = knn.predict(X_test)
scores.append(accuracy_score(y_test, y_pred))
# Calculate mean accuracy for current k
mean_score = np.mean(scores)
if mean_score > best_score:
    best_score = mean_score
    best_k = k
print(f"Optimal k: {best_k}, Error: {1 - best_score}")

```

ت) برای دسته بند ۱-NN:

- مرز تصمیم با نقاطی که از نزدیکترین نمونه های مثبت و منفی فاصله دارند، تعریف می شود.
- برای این مجموعه داده، مرزهای خطی تکه تکه ای بین نقاط کلاس های مختلف تشکیل می دهد، و در فضای ویژگی های دوی بعدی، شکلی زیر را ایجاد می کند.
- این مرز بر اساس تراکم و مکان نقاط در هر کلاس تغییر خواهد کرد.



این کد به شما مرز تصمیم گیری 1-NN را می دهد:

```

import matplotlib.pyplot as plt
from sklearn.neighbors import KNeighborsClassifier
import numpy as np
# Fitting 1-NN classifier
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X, y)
# Plotting decision boundary
x_min, x_max = 0, 10
y_min, y_max = 0, 10
xx, yy = np.meshgrid(np.linspace(x_min, x_max, 100), np.linspace(y_min, y_max, 100))
Z = knn.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
plt.contourf(xx, yy, Z, alpha=0.3)

```

```
plt.scatter(X[:, 0], X[:, 1], c=y, edgecolor='k', cmap='coolwarm')
plt.title("1-NN Decision Boundary")
plt.show()
```

سوال چهارم

الف) به روز رسانی مراکز خوشه ها برای به حداقل رساندن تابع هزینه

تابع هزینه برای خوشه بندی k-means به صورت زیر تعریف می شود:

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

هدف:

نشان دهیم که برای به حداقل رساندن تابع هزینه، مقدار بهینه (μ_j) میانگین تمام نقاط در (S_j) است.

اثبات:

۱. برچسب های (y_j) را ثابت در نظر میگیریم:

در این مرحله فرض می کنیم که تخصیص نقاط به خوشه ها (یعنی y_j) ثابت است.

۲. بازنویسی تابع هزینه برای یک خوشه واحد (j) :

$$L_j = \sum_{x_i \in S_j} |x_i - \mu_j|^2$$

۳. گرفتن مشتق با توجه به (μ_j) :

برای کمینه کردن (L_j) ، باید مشتق (L_j) را نسبت به (μ_j) پیدا کرده و آن را صفر کنیم.

$$\frac{\partial L_j}{\partial \mu_j} = \frac{\partial}{\partial \mu_j} \left(\sum_{x_i \in S_j} |x_i - \mu_j|^2 \right)$$

از آنجایی که $|x_i - \mu_j|^2 = (x_i - \mu_j) \cdot (x_i - \mu_j)$ ، این را گسترش داده و مشتق میگیریم:

$$\frac{\partial L_j}{\partial \mu_j} = \frac{\partial}{\partial \mu_j} \left(\sum_{x_i \in S_j} (x_i - \mu_j) \cdot (x_i - \mu_j) \right) = \sum_{x_i \in S_j} 2 \cdot -1 \cdot (x_i - \mu_j) = \sum_{x_i \in S_j} 2 (\mu_j - x_i)$$

۴. برابر قرار دادن مشتق با صفر:

$$\sum_{x_i \in S_j} 2 (\mu_j - x_i) = 0$$

ساده کردن:

$$\sum_{x_i \in S_j} (\mu_j - x_i) = 0$$

$$\sum_{x_i \in S_j} (\mu_j) - \sum_{x_i \in S_j} (x_i) = 0$$

$$\sum_{x_i \in S_j} (\mu_j) = \sum_{x_i \in S_j} (x_i)$$

$$|S_j| \mu_j = \sum_{x_i \in S_j} x_i$$

$$\mu_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i$$

نتیجه گیری:

مرکز خوشه (μ_j) که تابع هزینه را به حداقل می رساند میانگین نقاط موجود در خوشه (S_j) است.

(ب) حساسیت به مقداردهی اولیه:

بله، k-means به مقداردهی اولیه حساس است:

- مقداردهی اولیه متفاوت مراکز خوشه می تواند منجر به خوشه های نهایی متفاوت شود.
- مقداردهی اولیه ضعیف ممکن است باعث شود که الگوریتم به کمینه های محلی همگرا شود و در k-means لزوماً به global minimum نمیرسیم. در نتیجه ممکن است به suboptimal clustering solutions بدی برسیم.
- مقداردهی اولیه ضعیف می تواند منجر به خوشه های خالی یا راه حل های غیربهبوده شود.
- برای کاهش این امر، روش هایی مانند ++K-means برای انتخاب مراکز خوشه اولیه بهتر استفاده می شود. همچنین میتوانیم چند بار الگوریتم را اجرا کنیم.

همگرایی:

بله، k-means همگرا می شود، اما نه لزوماً به global minimum:

- الگوریتم همگرایی را تضمین می کند زیرا تابع هزینه L در هر مرحله (reassignment and center update) افزایش نمی یابد. (non-increasing at each step)
- تابع هزینه bounded below by 0 است و در هر step کاهش می یابد. تنها تعداد محدودی (possibly many) از تخصیص نقاط به خوشه ها وجود دارد بنابراین، الگوریتم در نهایت به یک حالت پایدار میرسد.
- با این حال، k-means ممکن است بسته به محل قرارگیری اولیه مراکز خوشه، به local minimum همگرا شود.
- الگوریتم زمانی متوقف می شود که هیچ تغییری در مراکز خوشه وجود نداشته باشد یا زمانی که برچسب نقاط تغییر نمی کند.

(پ) سناریو:

- فرض کنید یک نقطه x_j از چندین مرکز خوشه در مرحله به روز رسانی برچسب، فاصله یکسانی دارد.
- اگر یکی از گزینه ها این است که x_j را در همان خوشه ای که در تکرار قبلی بوده باقی بماند، بهتر است این گزینه را انتخاب کنید.

دلیل:

۱. تثبیت همگرایی (Stabilizing Convergence):

- اگر x_j را در خوشه قبلی خود در فاصله مساوی نگه داریم، با اجتناب از تغییرات غیرضروری در انتساب های خوشه، به تثبیت الگوریتم کمک می کند.
- جابجایی مکرر x_j بین خوشه های هم فاصله ممکن است باعث نوسانات در مراکز خوشه ای شود و همگرایی را به تاخیر بیاندازد یا حتی از آن جلوگیری کند.
- نقاط ممکن است به جابجایی بین خوشه ها ادامه دهند، حتی زمانی که مراکز تغییر نمیکنند.
- ۲. معیارهای توقف الگوریتم (Stopping Criteria):
- الگوریتم k-means معمولاً زمانی متوقف می شود که مراکز خوشه دیگر پس از یک تکرار تغییر نکنند یا تغییر خیلی ناچیزی داشته باشند.
- اگر نقاطی مانند x_j مکرر بین خوشه ها جابه جا می شوند، این شرط ممکن است برآورده نشود و باعث شود الگوریتم برای تکرارهای بیشتر اجرا شود.
- این روش کمک می کند تا از حلقه های بی نهایت جلوگیری شود و الگوریتم را قطعی تر می کند.

مشکل بدون این اصل:

مسئله نوسان:

- نقطه x_j ممکن است بین دو یا چند خوشه در تکرارهای متوالی در نوسان پیوسته باشد و از همگرا شدن الگوریتم جلوگیری کند. در واقع، نقاط می توانند بین خوشه ها جابجا شوند حتی زمانی که مراکز خوشه ها ثابت هستند. که ممکن است منجر به ایجاد چرخه ای که هرگز ثابت نمی شود بشود.
 - این مسئله می تواند منجر به افزایش زمان محاسباتی شود و ممکن است قبل از رسیدن به حالت پایدار به تکرارهای بیشتری نیاز داشته باشد.
- با انتخاب این روش برای نگه داشتن نقاط با فاصله مساوی در خوشه قبلی، احتمال این مسائل را کاهش می دهیم، بنابراین از همگرایی سریعتر و پایدارتر الگوریتم k-means اطمینان حاصل می کنیم.

سوال پنجم

الف) برای حل این مسئله، می خواهیم نشان دهیم که بردار u که میانگین مربعات خطا (MSE) را بین نقاط و تصویر آنها به حداقل می رساند، در واقع اولین مؤلفه اصلی (PC1) است.

داریم:

- مجموعه ای از m نقاط داده نرمال شده: $(x^{(1)}, x^{(2)}, \dots, x^{(m)} \in R^n)$
- داده ها نرمال شده اند، یعنی میانگین هر بعد صفر و واریانس هر بعد یک است.
- تصویر x بر روی بردار u :

$$f_u(x) = \arg \min_{v \in V} |x - v|^2 \text{ where } V = \{au: a \in R\}$$

- هدف این است که:

$$\arg \min_{u: u^T u = 1} \frac{1}{m} \sum_{i=1}^m |x^{(i)} - f_u(x^{(i)})|^2$$

مرحله ۱: بیان خطای تصویر کردن

برای بیان $f_u(x)$ ، توجه داشته باشید که تصویر x بر روی u به صورت زیر داده می شود:

$$f_u(x) = (x^T u)u$$

که $x^T u$ ضرب داخلی است و $(x^T u)u$ برداری در جهت بردار u است.

بنابراین، بردار خطا برای یک نقطه $(x - f_u(x))$ به صورت زیر است:

$$x - f_u(x) = x - (x^T u)u$$

پس مربع خطای یک نقطه به صورت زیر است:

$$|x - (x^T u)u|^2$$

مرحله ۲: خطای مربع را گسترش می دهیم.

بباید این عبارت را گسترش دهیم:

$$|x - (x^T u)u|^2 = (x - (x^T u)u)^T (x - (x^T u)u)$$

گسترش ضرب نقطه ای:

$$|x - (x^T u)u|^2 = x^T x - 2(x^T u)(u^T x) + (x^T u)^2 (u^T u)$$

از آنجایی که $(u^T u = 1)$

$$|x - (x^T u)u|^2 = x^T x - 2(x^T u)^2 + (x^T u)^2 = x^T x - (x^T u)^2$$

مرحله ۳: مجموع خطا را محاسبه میکنیم.

حالا بباید این خطا را روی تمام نقاط داده جمع کنیم و میانگین مربعات خطا را به حداقل برسانیم:

$$\text{Total Error} = \text{MSE} = \frac{1}{m} \sum_{i=1}^m \left(x^{(i)T} x^{(i)} - (x^{(i)T} u)^2 \right)$$

از آنجایی که داده ها نرمال شده اند، $(x^{(i)T} x^{(i)} = |x^{(i)}|^2)$ برابر واریانس است (چون میانگین صفر است) که برابر با ۱ است:

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m (x^{(i)T} x^{(i)}) - \frac{1}{m} \sum_{i=1}^m ((x^{(i)T} u)^2) = \text{Var}(X) - \frac{1}{m} \sum_{i=1}^m ((x^{(i)T} u)^2)$$

این معادله ساده می شود:

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

مرحله ۴: به حداکثر رساندن $\frac{1}{m} \sum_{i=1}^m (x^{(i)\top} u)^2$

کمینه کردن خطای کل معادل به حداکثر رساندن عبارت دوم در معادله بالاست (چرا که ضریب آن منفی است عبارت اول ثابت است) است:

$$\arg \max_{u: u^\top u=1} \frac{1}{m} \sum_{i=1}^m (x^{(i)\top} u)^2 = \arg \max_{u: u^\top u=1} \frac{1}{m} \sum_{i=1}^m (u^\top x^{(i)}) (x^{(i)\top} u) = \arg \max_{u: u^\top u=1} \frac{1}{m} \sum_{i=1}^m u^\top (x^{(i)} x^{(i)\top}) u$$

مرحله ۵: بازنویسی با ماتریس کوواریانس

Covariance Matrix یا همان ماتریس کوواریانس (Σ) را به صورت زیر تعریف میکنیم: (فرض میکنیم کل داده های را داریم اگر فقط نمونه ها را داشتیم باید در مخرج از $m-1$ استفاده میکردیم که البته تفاوت زیادی در محاسبه ایجاد نمیکند.)

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

تابع هدف تبدیل می شود به:

$$\arg \max_{u: u^\top u=1} u^\top \Sigma u$$

مرحله ۶: ضرایب لاگرانژ (Lagrange Multipliers)

برای حل این مشکل بهینه سازی مشروط شده، از Lagrange Multipliers استفاده می کنیم. لاگرانژ را به صورت زیر تعریف میکنیم:

$$\mathcal{L}(u, \lambda) = f(u) + \lambda g(u) = u^\top \Sigma u - \lambda (u^\top u - 1)$$

گرفتن مشتق با توجه به u و صفر قرار دادن آن:

$$\frac{\partial \mathcal{L}}{\partial u} = 2\Sigma u - 2\lambda u + 0 = 0$$

این عبارت ساده می شود:

$$\Sigma u - \lambda u = 0$$

$$\Sigma u = \lambda u$$

این دقیقاً همان معادله مقدار ویژه برای ماتریس کوواریانس (Σ) است! جواب u باید بردار ویژه Σ باشد و λ مقدار ویژه مربوطه است. از آنجایی که می خواهیم $u^\top \Sigma u$ را به حداکثر برسانیم، باید بردار ویژه مربوط به بزرگترین مقدار ویژه را انتخاب کنیم که دقیقاً تعریف اولین مؤلفه اساسی (PC1) است.

مرحله ۷: نتیجه گیری

بردار u که explained variance را به حداکثر می رساند، بردار ویژه مربوط به بزرگترین مقدار ویژه (Σ) است که طبق تعریف همان اولین مؤلفه اصلی (PC1) است.

بدین ترتیب نشان دادیم:

$$\arg \min_{u: u^T u = 1} \frac{1}{m} \sum_{i=1}^m |x^{(i)} - f_u(x^{(i)})|^2 = \text{PC1}$$

که این اثبات را کامل می کند.

همچنین این اثبات نشان می دهد که:

۱. تصویری (projected) که خطای بازسازی (reconstruction error) را به حداقل می رساند.

۲. جهتی که واریانس را به حداکثر می رساند.

۳. اولین مولفه اصلی (PC1).

همه معادل هستند!

پایان