

بسمه تعالیٰ



دانشگاه صنعتی اصفهان

پاسخنامه تکلیف دوم

مبانی ریاضی علوم داده

پاییز ۱۴۰۴

(الف)

روش PCA یک روش خطی برای کاهش بُعد است که بر مبنای بیشینه‌سازی واریانس داده‌ها عمل می‌کند. فرض کنید داده‌ها به صورت ماتریس $X \in \mathbb{R}^{d \times n}$ (هر ستون یک نمونه) داده شده‌اند و داده‌ها مرکز شده‌اند. در PCA به دنبال زیرفضای خطی k -بعدی هستیم که تصویر داده‌ها روی آن، بیشترین واریانس ممکن را حفظ کند (یا به‌طور معادل، خطای بازسازی کمینه شود).

اگر تجزیه مقادیر تکین داده‌ها به صورت

$$X = U\Sigma V^T$$

باشد، آنگاه نمایش k -بعدی داده‌ها با نگهداشتن k مقدار تکین اول به صورت

$$X_k = U_k \Sigma_k V_k^T$$

به دست می‌آید. ستون‌های U_k جهت‌های اصلی داده هستند و این روش از نظر بهینگی در نرم فربنیوس و طیفی تضمین شده است. از نقاط قوت PCA می‌توان به تفسیرپذیری بالا و حذف نویز خطی اشاره کرد؛ اما این روش نیازمند محاسبه SVD بوده و به خود داده‌ها وابسته است.

(ب)

در روش Random Projection (و در قالب کلی‌تر، Compressive Sensing) کاهش بُعد با استفاده از یک تصویرسازی تصادفی انجام می‌شود. ایده اصلی آن است که داده‌های با بُعد بالا را با ضرب در یک ماتریس تصادفی به فضای کم‌بعد منتقل کنیم، بدون آن‌که فاصله‌ها (یا ساختار هندسی کلی داده‌ها) به‌طور معنادار تخریب شوند.

به طور دقیق، یک ماتریس تصادفی

$$R \in \mathbb{R}^{k \times d}$$

(مثلاً با درایه‌های گاووسی یا رادماخر) انتخاب می‌شود و نگاشت به صورت

$$y = Rx$$

تعریف می‌گردد. طبق لم جانسون-لیندن‌اشترووس، اگر

$$k = O\left(\frac{\log n}{\varepsilon^2}\right),$$

آنگاه با احتمال بالا فاصله‌های زوج نقاط تا خطای نسبی ε حفظ می‌شوند.

در فشرده‌سازی یا Compressive Sensing نیز ایده مشابه است، با این تفاوت که فرض می‌شود داده (یا سیگنال) در یک پایه‌ی مناسب تُنگ است و می‌توان آن را از تعداد کمی اندازه‌گیری تصادفی با حل یک

مسئله بهینه‌سازی بازیابی کرد. نکته‌ی کلیدی این است که در Random Projection طول بردار تصویرشده حفظ نمی‌شود، اما روابط نسبی بین نقاط (مانند فاصله‌ها و زوایا) تقریباً ثابت می‌مانند.

(پ)

مقایسه‌ی دو روش به صورت زیر است:

- PCA یک روش **داده‌محور** است و جهت‌های تصویرسازی مستقیماً از ساختار داده استخراج می‌شوند، در حالی که Random Projection کاملاً **داده‌ناوابسته** و تصادفی است.
- PCA بهینه‌ترین تقریب کمرنگ را (در معنای نرم فربنیوس و طیفی) ارائه می‌دهد، اما محاسبات آن پرهزینه‌تر است؛ در مقابل Random Projection بسیار سریع و مقیاس‌پذیر است.
- در PCA تفسیرپذیری جهت‌ها بالا است، ولی در Random Projection جهت‌ها فاقد معنای تفسیری مستقیم‌اند.
- PCA برای فشرده‌سازی و نویززدایی مناسب‌تر است، در حالی که Com- Random Projection برای حفظ فاصله‌ها، محاسبات سریع و کار با داده‌های بسیار بزرگ کاربرد بیشتری دارند.

(الف)

اگر $A \in \mathbb{R}^{d \times n}$ با مقادیر تکین $r = \text{rank}(A)$ و $\sigma_1 \geq \dots \geq \sigma_r > 0$ باشد، داریم:

$$\|A\|_F^2 = \sum_{i=1}^r \sigma_i^2, \quad \|A\|_2 = \sigma_1.$$

چون $\sigma_i \leq \sigma_1$ برای همه i ، نتیجه می‌شود:

$$\|A\|_F^2 = \sum_{i=1}^r \sigma_i^2 \leq \sum_{i=1}^r \sigma_1^2 = r\sigma_1^2 = r\|A\|_2^2 \Rightarrow \|A\|_F \leq \sqrt{r}\|A\|_2.$$

(ب)

از مرتب بودن مقادیر تکین داریم $\sigma_1 \geq \dots \geq \sigma_k$ ، بنابراین:

$$\sum_{i=1}^k \sigma_i^2 \geq k\sigma_k^2.$$

از طرف دیگر $\sum_{i=1}^k \sigma_i^2 \leq \sum_{i=1}^r \sigma_i^2 = \|A\|_F^2$ پس:

$$k\sigma_k^2 \leq \|A\|_F^2 \Rightarrow \sigma_k \leq \frac{\|A\|_F}{\sqrt{k}}.$$

(پ)

تقریب کم‌رنگ بهینه با برش SVD داده می‌شود:

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T, \quad \text{rank}(A_k) = k.$$

طبق قضیه اکهارت-یانگ:

$$\|A - A_k\|_2 = \sigma_{k+1}.$$

همچنین از نامساوی بخش (ب) نتیجه می‌شود $B := A_k \leq \|A\|_F / \sqrt{k}$ داریم:

$$\|A - B\|_2 \leq \frac{\|A\|_F}{\sqrt{k}}.$$

(الف)

اگر A مربعی و وارون پذیر باشد، با نوشتن $A = U\Sigma V^T$ داریم:

$$A^{-1} = V\Sigma^{-1}U^T.$$

از طرفی شبیه وارون مور-پنروز در این حالت برابر است با:

$$A^* = V\Sigma^{-1}U^T,$$

پس:

$$A^* = A^{-1}.$$

(ب)

برای گزاره ها:

(۱)

$$\|v_1^T A^*\|_2 = \min_{\|x\|=1} \frac{1}{\|Ax\|_2}.$$

از آنجا که $A^* = \sum_{i=1}^r \frac{1}{\sigma_i} v_i u_i^T$ داریم:

$$v_1^T A^* = \frac{1}{\sigma_1} u_1^T \Rightarrow \|v_1^T A^*\|_2 = \frac{1}{\sigma_1}.$$

و نیز $\max_{\|x\|=1} \|Ax\|_2 = \|A\|_2 = \sigma_1$ پس

$$\min_{\|x\|=1} \frac{1}{\|Ax\|_2} = \frac{1}{\sigma_1}.$$

بنابراین گزاره درست است.

$$(درست؛ از شرایط مور-پنروز) \quad A^* A A^* = A^* \quad (۲)$$

$$(درست؛ از شرایط مور-پنروز) \quad A = A A^* A \quad (۳)$$

$$(نادرست در حالت کلی) \quad A A^* = I \quad (۴)$$

مثال:

$$A = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \Rightarrow A^* = \begin{bmatrix} 1 & 0 \end{bmatrix}, \quad A A^* = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \neq I_2.$$

$$(درست) \quad A^* = (A^T A)^* A^T \quad (۵)$$

(الف)

با فرض $\sum_{i=1}^r c_i^2 = 1$ و قید $\sigma_1 \geq \dots \geq \sigma_r > 0$ داریم:

$$\sum_{i=1}^r c_i^2 \sigma_i^2 \leq \sigma_1^2 \sum_{i=1}^r c_i^2 = \sigma_1^2.$$

برابری وقتی رخ می‌دهد که $c_1 = \pm 1$ و $c_{i \neq 1} = 0$. بنابراین:

$$\max = \sigma_1^2.$$

(ب)

اگنون قیود $\sum c_i = 0$ و $\sum c_i^2 = 1$ برقرار است. به دلیل قید جمع صفر، لااقل دو مؤلفه با علامت مخالف لازم است. بهترین حالت آن است که این دو مؤلفه روی بزرگ‌ترین σ ها قرار گیرند:

$$c_1 = \frac{1}{\sqrt{2}}, \quad c_2 = -\frac{1}{\sqrt{2}}, \quad c_3 = \dots = c_r = 0,$$

که در نتیجه:

$$\sum_{i=1}^r c_i^2 \sigma_i^2 = \frac{1}{2}\sigma_1^2 + \frac{1}{2}\sigma_2^2 \quad \Rightarrow \quad \max = \frac{\sigma_1^2 + \sigma_2^2}{2}.$$

(الف)

محاسبه دقیق Ax زمان $O(nd)$ دارد. طبق نتیجه سؤال ۲، برای هر k می‌توان تقریب رتبه- k از A (مثلاً برش SVD) را به صورت A_k در نظر گرفت به طوری که:

$$\|A - A_k\|_2 \leq \frac{\|A\|_F}{\sqrt{k}}.$$

اکتون خروجی را

$$y := A_k x$$

می‌گیریم. آنگاه:

$$\|y - Ax\|_2 = \|(A_k - A)x\|_2 \leq \|A - A_k\|_2 \|x\|_2 \leq \frac{\|A\|_F}{\sqrt{k}} \|x\|_2.$$

برای اینکه $\|y - Ax\|_2 \leq \varepsilon \|A\|_F \|x\|_2$ برقرار باشد کافی است:

$$\frac{1}{\sqrt{k}} \leq \varepsilon \quad \Rightarrow \quad k \geq \frac{1}{\varepsilon^2}.$$

پس انتخاب حروفهای:

$$k = \left\lceil \frac{1}{\varepsilon^2} \right\rceil$$

کفایت می‌کند.

(ب)

اگر $A_k = U_k \Sigma_k V_k^T$ باشد:

$$y = A_k x = U_k \Sigma_k (V_k^T x).$$

محاسبه $V_k^T x$ زمان $O(dk)$ و سپس ضرب در U_k زمان $O(nk)$ دارد، لذا کل زمان:

$$O((d+n)k).$$

(الف)

برای

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

داریم:

$$A^T A = \begin{bmatrix} 4 & 0 \\ 0 & 16 \end{bmatrix}.$$

$$v^{(0)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \text{ با پاور مت روی } A^T A$$

$$v^{(t+1)} = \frac{A^T A v^{(t)}}{\|A^T A v^{(t)}\|}.$$

گام ۱:

$$w^{(1)} = \begin{bmatrix} 4 \\ 16 \end{bmatrix}, \quad v^{(1)} = \frac{1}{\sqrt{4^2 + 16^2}} \begin{bmatrix} 4 \\ 16 \end{bmatrix} = \begin{bmatrix} 0.2425356 \\ 0.9701425 \end{bmatrix}.$$

گام ۲:

$$w^{(2)} = \begin{bmatrix} 0.9701424 \\ 15.52228 \end{bmatrix}, \quad v^{(2)} \approx \begin{bmatrix} 0.0623783 \\ 0.9980526 \end{bmatrix}.$$

گام ۳:

$$w^{(3)} = \begin{bmatrix} 0.2495132 \\ 15.96884 \end{bmatrix}, \quad v^{(3)} \approx \begin{bmatrix} 0.0156231 \\ 0.9998780 \end{bmatrix}.$$

پس:

$$v_1 \approx \begin{bmatrix} 0.0156 \\ 0.9999 \end{bmatrix}.$$

(ب)

چون $A^T A = \text{diag}(4, 16)$, مقادیر تکین:

$$\sigma_1 = 4, \quad \sigma_2 = 2.$$

بردار تکین راست اول متناظر با بزرگ‌ترین مقدار ویژه 16 برابر است با:

$$v_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

و

$$u_1 = \frac{Av_1}{\sigma_1} = \frac{1}{4} \begin{bmatrix} 2 \\ 2 \\ -2 \\ -2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}.$$

تقریب پاورمتد پس از ۳ تکرار با v_1 دقیق هم‌راستا و بسیار نزدیک است.

(پ)

در تفسیر توصیه‌گر، v_1 الگوی غالب در فضای آیتم‌ها (rstaurant) و u_1 میزان هم‌راستایی کاربران با آن الگو را نشان می‌دهد. اختلاف/نسبت σ_1 و σ_2 شدت «تک‌بعدی بودن» ساختار غالب را مشخص می‌کند.

(الف)
اگر

$$A = \begin{bmatrix} B & 0 \\ 0 & C \end{bmatrix},$$

آنگاه:

$$A^T A = \begin{bmatrix} B^T B & 0 \\ 0 & C^T C \end{bmatrix}.$$

بنابراین مقادیر ویژه $A^T A$ اجتماع مقادیر ویژه $B^T B$ و $C^T C$ است. چون مقادیر تکین ریشه‌ی مربع مقادیر ویژه‌اند، نتیجه می‌شود مقادیر تکین A دقیقاً اجتماع مقادیر تکین B و C (با احتساب تکرار) خواهد بود.

(ب)

برای بردارهای تکین راست نیز: اگر v_B بردار ویژه $B^T B$ باشد، آنگاه

$$\begin{bmatrix} v_B \\ 0 \end{bmatrix}$$

بردار ویژه $A^T A$ است و مشابه آن برای C . با تعریف $u = \frac{Av}{\sigma}$ ، بردارهای تکین چپ نیز به همین روش از بردارهای تکین بلوک‌ها ساخته می‌شوند.

(الف)

فرض کنید $X \in \mathbb{R}^{d \times n}$ ماتریس مختصات باشد (ستون i برابر x_i) و ماتریس فاصله‌ها $D = [d_{ij}]$ با $d_{ij} = \|x_i - x_j\|$ داده شده باشد. همچنین فرض کنید مرکز هندسی در مبدأ است، یعنی:

$$X\mathbf{1} = 0.$$

اگر $G = X^T X$ باشد، آنگاه:

$$d_{ij}^2 = \|x_i\|^2 + \|x_j\|^2 - 2G_{ij}.$$

با «دو بار مرکزسازی» به رابطه استاندارد می‌رسیم:

$$G = -\frac{1}{2} JD^{(2)} J, \quad J = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T, \quad D^{(2)} = [d_{ij}^2].$$

(ب)

الگوریتم بازسازی مختصات (MDS) کلاسیک():

- ابتدا $D^{(2)} = [d_{ij}^2]$ را تشکیل دهید.

- سپس $J = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$ را بسازید.

- ماتریس گرام را حساب کنید: $.G = -\frac{1}{2} JD^{(2)} J$

- تجزیه ویژه انجام دهید: $.G = V\Lambda V^T$

- با انتخاب d مقدار ویژه مثبت اول:

$$X = \Lambda_d^{1/2} V_d^T.$$

(در صورت وجود نویز، مقادیر ویژه منفی کوچک معمولاً صفر در نظر گرفته می‌شوند.)