

Метод головних компонент з точки зору методів оптимізації

Н. Фордуй, О. Галганов

Задано набір $\{x_1, \dots, x_m\}$ із m точок в \mathbb{R}^n . Необхідно для заданого $k < n$ знайти k -вимірну гіперплощину, яка буде найближчою до цих точок в сенсі евклідової норми: тобто, мінімальною має бути різниця між цими точками та їх проекціями на шукану площину.

Ця задача називається **методом головних компонент** (англ. principal component analysis, PCA) і широко застосовується в статистиці та машинному навчанні.

Постановка задачі

Розглянемо задані $\{x_1, \dots, x_m\}$ — m векторів з \mathbb{R}^n . Відомо, що k -вимірна гіперплощина H_k в \mathbb{R}^n задається k ортонормованими векторами $\{u_1, \dots, u_k\}$, які можна доповнити до базису (далі — ОНБ), та вектором зсуву відносно нуля b :

$$H_k = \{x = b + c_1 u_1 + \dots + c_k u_k : c_1, \dots, c_k \in \mathbb{R}\}.$$

Нехай $\{u_i\}_{i=1}^n$ — деякий ОНБ в \mathbb{R}^n . Тоді $\forall i = 1, \dots, n : x_i = b + \sum_{j=1}^n c_{i,j} u_j$, де $c_{i,j} = (x_i, u_j)$ (це — розклад Фур'є в \mathbb{R}^n).

Покажемо, що після застосування деякого перетворення заданих векторів можна вважати $b = 0$.

Допоміжна задача

Для заданих m точок $\{x_1, \dots, x_m\}$ в \mathbb{R}^n знайти точку, яка знаходиться найближче до них в сенсі евклідової норми.

Розв'язання допоміжної задачі

Задача має вигляд $F(x) = \sum_{k=1}^m \|x - x_k\|^2 \rightarrow \min, x \in \mathbb{R}^n$.

$F'_x(x^*) = 2 \sum_{k=1}^m (x^* - x_k) = 0 \Rightarrow x^* = \frac{1}{m} \sum_{k=1}^m x_k$ — стаціонарна точка.

Оскільки $F(x)$ — опукла функція, то $x^* = \frac{1}{m} \sum_{k=1}^m x_k$ — розв'язок.

Постановка задачі

Таким чином, якщо замінити x_i на $y_i = x_i - \frac{1}{m} \sum_{i=1}^m x_i$, то можна вважати $b = 0$, оскільки тепер найближчою точкою \mathbb{R}^n буде точка 0, тому й шукана гіперплощина теж має проходити через 0.

Отже, $\forall i = 1, \dots, n : y_i = \sum_{j=1}^n c_{i,j} u_j$, де $c_{i,j} = (y_i, u_j)$. Проекції на гіпер-

площину y_i будемо шукати у вигляді $\hat{y}_i = \sum_{j=1}^k c_{i,j} u_j$, $k < n$. Введемо

вектори похибок $\varepsilon_i = y_i - \hat{y}_i = \sum_{j=k+1}^n c_{i,j} u_j$, які зберемо у матрицю:

$$E = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_m) = (u_{k+1}, u_{k+2}, \dots, u_n) \cdot \begin{pmatrix} c_{1,k+1} & c_{2,k+1} & \dots & c_{m,k+1} \\ c_{1,k+2} & c_{2,k+2} & \dots & c_{m,k+2} \\ \dots & \dots & \dots & \dots \\ c_{1,n} & c_{2,n} & \dots & c_{m,n} \end{pmatrix}$$

Коротко: $E = UC$, причому $U^T U = I$, бо ці вектори ортонормовані, а $C = Y^T U$.

Задача

Знайти такі ортонормовані вектори $\{u_i\}_{i=1}^n$, що $\|E\|^2 \rightarrow \min$, де $\|E\| = \sqrt{\sum_{i,j=1}^n e_{ij}^2}$ — норма Фробеніуса матриці похибок E .

Введемо позначення $Y = (y_1, \dots, y_m)$, $F = YY^T$.

$\|E\|^2 = \text{Tr}(E^T E) = \text{Tr}(C^T U^T U C) = \text{Tr}(C^T C) = \text{Tr}(U^T Y Y^T U) = \text{Tr}(U^T F U)$. Оскільки $U = \sum_{j=k+1}^n (0, \dots, 0, u_j, 0, \dots, 0) = \sum_{j=k+1}^n U_j$, за лінійністю Tr (слід матриці, сума діагональних елементів) маємо $\text{Tr}(U^T F U) = \sum_{j=k+1}^n \text{Tr}(U_j^T F U_j)$. В кожній матриці U_j лише один стовпець не рівний нулю, тому $\sum_{j=k+1}^n \text{Tr}(U_j^T F U_j) = \sum_{j=k+1}^n (F u_j, u_j)$.

Таким чином, отримуємо задачу умовної оптимізації.

Розв'язання задачі

$$\begin{cases} F(u_{k+1}, \dots, u_n) = \sum_{j=k+1}^n (Fu_j, u_j) \rightarrow \min \\ \|u_j\|^2 = 1, \quad j = k+1, \dots, n \\ \{u_{k+1}, \dots, u_n\} \text{ — лінійно незалежні} \end{cases}$$

Ця задача є регулярною, бо градієнти функцій, що задають обмеження, є лінійно незалежними за умовою.

Маємо функцію Лагранжа:

$$\mathcal{L}(u_{k+1}, \dots, u_n, \lambda_{k+1}, \dots, \lambda_n) = \sum_{j=k+1}^n \left((Fu_j, u_j) + \lambda_j (\|u_j\|^2 - 1) \right)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial u_j} = 2Fu_j + 2\lambda_j u_j = 0 \\ (j = k+1, \dots, n) \\ \|u_j\|^2 = 1, \quad j = k+1, \dots, n \\ \{u_{k+1}, \dots, u_n\} \text{ — лінійно незалежні} \end{cases}$$

Оскільки $2Fu_j + 2\lambda_j u_j = 0 \Leftrightarrow Fu_j = -\lambda_j u_j$, то розв'язком системи будуть u_j — власні вектори F одиничної норми. Але з умови мінімізації цільової функції та лінійної незалежності $\{u_{k+1}, \dots, u_n\}$ ці власні вектори мають відповідати найменшим власним числам $\mu_j = -\lambda_j$.

$F^T = (YY^T)^T = YY^T$, $F \geq 0$, бо $\forall x \in \mathbb{R}^n : (Fx, x) = (YY^T x, x) = (Y^T x, Y^T x) \geq 0$. Таким чином, всі $\mu_j = -\lambda_j \geq 0$.

Оскільки цільова функція є нескінченно зростаючою, то отримані u_j є розв'язками задачі. Зрозуміло, що інші складові ОНБ $\{u_i\}_{i=1}^n$ можна покласти рівними іншим власним векторам матриці F , розташувавши всі отримані у порядку спадання відповідних власних значень, причому векторам u_1, \dots, u_k відповідатимуть k найбільших власних чисел (нагадаємо, що за цими векторами розкладалися наближення \hat{y}_i).

Залишилося згадати, що $y_i = x_i - \frac{1}{m} \sum_{i=1}^m x_i$. Позначимо $X = (x_1, \dots, x_m)$,

тоді $Y = X - \left(\frac{1}{m} \sum_{i=1}^m x_i \right) \cdot \underbrace{(1, 1, \dots, 1)}_m$. Тепер обчислимо матрицю YY^T .

Введемо позначення $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$.

$$\begin{aligned} YY^T &= (X - \bar{x} \cdot (1, 1, \dots, 1)) \cdot (X^T - (1, 1, \dots, 1)^T \cdot \bar{x}^T) = \\ &= XX^T - \bar{x} \cdot (1, 1, \dots, 1) \cdot X^T - X \cdot (1, 1, \dots, 1)^T \cdot \bar{x}^T + \bar{x} \cdot (1, 1, \dots, 1) \cdot \\ &\quad (1, 1, \dots, 1)^T \cdot \bar{x}^T = XX^T - m \cdot \bar{x} \cdot \bar{x}^T - m \cdot \bar{x} \cdot \bar{x}^T + m \cdot \bar{x} \cdot \bar{x}^T = \\ &= XX^T - m \cdot \bar{x} \cdot \bar{x}^T = \\ &= XX^T - \frac{1}{m} \left(\sum_{i=1}^m x_i \right) \cdot \left(\sum_{i=1}^m x_i^T \right) \end{aligned}$$

Таким чином, знайдені вектори u_j є власними векторами матриці $XX^T - m \cdot \bar{x} \cdot \bar{x}^T$.

Шуканою k -вимірною гіперплощиною є $\bar{x} + L(u_1, \dots, u_k)$.

Тут $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$, а $L(u_1, \dots, u_k)$ — лінійна оболонка власних векторів матриці $XX^T = m \cdot \bar{x} \cdot \bar{x}^T$ ($X = (x_1, \dots, x_m)$), які відповідають k найбільшим власним числам. Зазначимо, що при цьому похибка (сума квадратів евклідових норм векторів ε_i) рівна $\sum_{j=k+1}^n \mu_j$ — сумі $n - k$ найменших власних чисел цієї матриці.

Для практичного застосування є корисною формула для обчислення проєкцій $y_i = x_i - \bar{x}$ на $L(u_1, \dots, u_k)$: $\text{pr}(y_i) = (u_1, \dots, u_k)^T \cdot y_i$, або в матричному вигляді: $\text{pr}(Y) = (u_1, \dots, u_k)^T \cdot Y$

Додаток: диференціювання за векторним аргументом

У розв'язанні задачі було використано похідні скалярної функції з декількома векторними аргументами. Під частинною похідною за векторним аргументом вважається вектор з частинних похідних цієї функції за координатами вектора. Доречно навести виведення використаної формули для похідної від квадратичної форми.

$F(x) = (Ax, x)$, $x \in \mathbb{R}^n$, A — дійсна симетрична $n \times n$ матриця.

$$F(x+h) - F(x) = (Ax + Ah, x+h) - (Ax, x) = (Ax, x) + (Ax, h) + (Ah, x) + (Ah, h) - (Ax, x) = [A = A^T] = (2Ax, h) + (Ah, h)$$

Отже, лінійна частина приросту рівна $2Ax$, звідки $F'(x) = 2Ax$. Зокрема, похідна квадрата норми $\|x\|^2$ рівна $2x$.