

Section 1. Some important Linear Algebra background for PCA

• Some useful conclusion

For a matrix $A \in \mathbb{R}^{m \times n}$

- ① $A^T A$ and AA^T is symmetric matrix
- ② $A^T A$ and AA^T can be diagonalizable and get an orthonormal eigenvector.
- ③ $\text{rank}(A) = \text{rank}(A^T) = \text{rank}(A^T A) = \text{rank}(A^T A)$
- ④ $A^T A$ is positive semi-definite. If all the column of A is independent, then $A^T A$ is positive definite.
- ⑤ $A^T A$ and AA^T have the same non-zero eigenvalues. The number of non-zero eigenvalues is equal to $\text{rank}(A)$.

• Spectral theorem

For a symmetric matrix, $A = U D U^T$, where $U U^T = U^T U = I$

• SVD decomposition

- ① For a matrix $A_{m \times n}$, $A = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T$, where $U U^T = U^T U = I$
 $V^T V = V V^T = I$
 Σ is diagonal. $\Sigma = (\sigma_1, \sigma_2, \dots, \sigma_r, 0_{n-r})$
 $\sigma_1 > \sigma_2 > \dots > \sigma_r$
- ② Question: How to compute SVD?
 - 1) Right Singular Vector: Find eigenvectors of $A^T A$
 This give V matrix
 - 2) Singular Value: Find eigenvalues of $A^T A$. $A^T A$ has $\text{rank}(A)$ eigenvalues and $A^T A / AA^T$ has the same non-zero eigenvalue. (Conclusion ⑤)
 $\sigma_j = \sqrt{\lambda_j} \quad 1 \leq j \leq r$
 $= 0 \quad r+1 \leq j \leq n$
 - 3) Left Singular Vector: $U_j = \frac{1}{\sigma_j} A V_j \quad 1 \leq j \leq r$

• Truncated SVD

Ignore some of the small singular vector.

$$\tilde{A} = \sigma_1 U_1 V_1^T + \sigma_2 U_2 V_2^T \quad (\text{Only use the Top 2 largest } \sigma)$$

• Matrix norm and Eckart-Young Theorem

① Frobenius norm of A is $\|A\|_F^2 = \sum_{i,j} a_{ij}^2 = \text{trace}(A^T A) = \sum_{i=1}^r \sigma_i^2$

② ℓ_2 norm of a matrix $\|A\|_2 = \max_{1 \leq j \leq n} \sigma_{\max} = \sigma_1$

③ Eckart-Young Theorem:

$A^{m \times n}$ is a rank r matrix, $B^{m \times n}$ is a rank k matrix where $k \leq r$
 Define $\hat{A}_k = \sum_{i=1}^k \sigma_i u_i v_i^T$, then $\|A - \hat{A}_k\|_F^2 = \sum_{i=k+1}^r \sigma_i^2$
 $\|A - \hat{A}_k\|_2 = \sigma_{k+1}$

④ Question: How to decide K in low rank approximation?

$$\text{where } \hat{A}_k = \sum_{i=1}^k \sigma_i u_i v_i^T$$

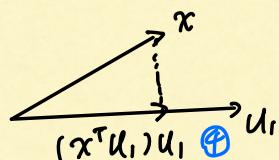
$$\text{Answer: } \frac{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_k^2}{\sigma_{k+1}^2 + \sigma_{k+2}^2 + \dots + \sigma_n^2} \rightarrow \text{threshold (0.95 / 0.90 \dots)}$$

Section 2. PCA

• 2 - prospectives

• Maximum variance of projection

Project onto u_1 : Find a direction u_1 such that the variance of the projection of the data onto u_1 is maximized.



$$\begin{aligned} & \underset{u_1}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \left[((x - \bar{x})^T u_1) u_1 \right]^2, \quad u_1^T u_1 = 1 \\ & \underset{u_1}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n [(x^T u_1) u_1]^2 \quad \textcircled{1} \\ & \Rightarrow \underset{u_1}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n (u_1^T x)^2 \quad \textcircled{2} \\ & \Rightarrow \underset{u_1}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n (u_1^T x)(u_1^T x)^T \\ & \Rightarrow \underset{u_1}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n u_1^T x x^T u_1 \\ & \Rightarrow \underset{u_1}{\operatorname{argmax}} u_1^T \left(\frac{1}{n} \sum_{i=1}^n x x^T \right) u_1 \end{aligned}$$

Note that $\frac{1}{n} \sum_{i=1}^n x x^T$ is a covariance matrix.

$$\text{define } S = \frac{1}{n} \sum_{i=1}^n x x^T$$

$$\text{Then: } \underset{u_1}{\operatorname{argmax}} u_1^T S u_1$$

The optimization problem above can be written as

$$\left\{ \begin{array}{l} \arg \max_{U_i} U_i^T S U_i \\ U_i^T U_i = 1 \end{array} \right.$$

Apply Lagrange multiplier:

$$L(U_i, \lambda) = U_i^T S U_i - \lambda(U_i^T U_i - 1) = 0$$

$$\frac{d L(U_i, \lambda)}{d U_i} = 2U_i^T S - 2\lambda U_i = 0 \quad \textcircled{3}$$

$$\Rightarrow \underbrace{U_i^T S}_{= \lambda U_i}$$

where λ is the eigenvalue of S and U_i is the eigenvector of S

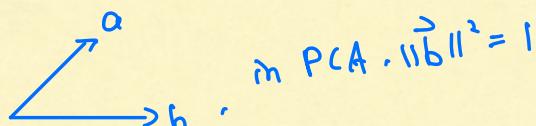
Remark

- ① For convenience - assume x has been centered such that $\bar{x} = 0$
- ② $(X^T U)$ is a coefficient (constant value), therefore $x^T U = U^T X$
- ③ Apply the formula of matrix derivation:

$$\left\{ \begin{array}{l} \frac{\partial X^T A X}{\partial X} = 2AX \quad (A \text{ is symmetric}) \\ \frac{\partial X^T X}{\partial X} = 2X \end{array} \right.$$

- ④ Vector projection formula:

$$\text{Proj}_{\vec{b}} \vec{a} = \frac{\vec{a} \cdot \vec{b}}{\|\vec{b}\|^2} \vec{b}$$



$$\text{in PCA}, \|\vec{b}\|^2 = 1$$

Summary:

According to the procedure above, PCA can be generalized as 4 steps.

- ① centralize data.
- ② Form covariance matrix $S = X X^T$
- ③ Find eigenvectors and eigenvalues of S $U_i^T S = \lambda U_i$
- ④ The larger the eigenvalue \Leftrightarrow the more significant direction

Minimum reconstruction error

Choose k PCs (principal components)

$$\tilde{V} = \begin{pmatrix} v_1 & v_2 & \cdots & v_k \end{pmatrix} \quad k \ll d$$

Project data point x_i onto the subspace spanned by the first k PCs.

$$\begin{aligned}\tilde{x}_i &= (x_i^T v_1) v_1 + (x_i^T v_2) v_2 + \cdots + (x_i^T v_k) v_k \\ &= (x_i^T v_1) v_1 + (x_i^T v_2) v_2 + \cdots + (x_i^T v_k) v_k \\ &= v_1^T v_1 x_i + v_2^T v_2 x_i + \cdots + v_k^T v_k x_i \quad (\text{Remark ②}) \\ &= \tilde{V}^T V x_i\end{aligned}$$

The PCA reconstruction error (Pearson 1991) is defined as

$$E = \frac{1}{n} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$$

$$\text{Given that } x_i = \sum_{j=1}^d v_j^T v_j x_i, \quad \tilde{x}_i = \sum_{j=1}^k v_j^T v_j x_i$$

$$\begin{aligned}\text{Then, } E &= \frac{1}{n} \sum_{i=1}^N \left\| \sum_{j=1}^d v_j^T v_j x_i - \sum_{j=1}^k v_j^T v_j x_i \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^N \left\| \left(\sum_{j=k+1}^d v_j^T v_j \right) x_i \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^N x_i^T \left(\sum_{j=k+1}^d v_j^T v_j \right) \left(\sum_{q=k+1}^d v_q^T v_q \right)^T x_i \\ &= \frac{1}{n} \sum_{i=1}^N \sum_{j=k+1}^d (v_j^T x_i) (x_i^T v_j) \\ &= \frac{1}{n} \sum_{i=1}^N \sum_{j=k+1}^d v_j^T (x_i x_i^T) v_j \\ &= \sum_{j=k+1}^d v_j^T \frac{1}{n} \sum_{i=1}^N (x_i x_i^T) v_j\end{aligned}$$

$$E = \sum_{j=k+1}^d v_j^T S v_j$$

According to our intuition, we want to minimize E .

Note that:

$$\underbrace{\sum_{j=1}^K v_j^\top S v_j}_{\text{First perspective: maximize this term}} + \underbrace{\sum_{j=k+1}^d v_j^\top S v_j}_{\text{second perspective: minimize this term}} = \sum_{j=1}^d v_j^\top S v_j$$

First perspective: maximize this term
Second perspective: minimize this term

Section 3 Kernel PCA

We have n points x_1, x_2, \dots, x_n each lying in \mathbb{R}^d

- standard PCA doesn't yield good features on highly non-linear datasets.
- In detail, define ϕ is some non-linear transformation
 $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m, m > d$
- We imitate PCA procedures in above highly non-linear datasets
 - Assume that $\phi(x_i)$ has been centralized $\frac{1}{N} \sum_{i=1}^N \phi(x_i) = 0$
 - Form covariance matrix $S = \frac{1}{N} \sum_{i=1}^N \phi(x_i) \phi(x_i)^\top \quad \textcircled{1}$
 - Find eigs $V_k S = \lambda_k V_k, k = 1, 2, \dots, M \quad \textcircled{2}$

Combined \textcircled{1}, \textcircled{2}:

$$\frac{1}{N} \sum_{i=1}^N \phi(x_i) (\phi(x_i)^\top V_k) = \lambda_k V_k \quad \textcircled{3}$$

Note that each v_k is a linear combination of $\phi(x_i)$.

$$v_k = \sum_{j=1}^n a_{kj} \phi(x_j) \quad \textcircled{4}$$

Combined \textcircled{3}, \textcircled{4}:

$$\frac{1}{N} \sum_{i=1}^N \phi(x_i) \phi(x_i)^\top \sum_{j=1}^n a_{kj} \phi(x_j) = \lambda_k \sum_{i=1}^N a_{ki} \phi(x_i) \quad \textcircled{5}$$

Define $K(x_i, x_j) = \phi(x_i)^\top \phi(x_j) \quad \textcircled{6}$

Multiply $\phi(x_i)^\top$ both sides in \textcircled{5}:

$$\frac{1}{N} \sum_{i=1}^N \phi(x_i)^\top \phi(x_i) \phi(x_i)^\top \sum_{j=1}^n a_{kj} \phi(x_j) = \lambda_k \sum_{i=1}^N \phi(x_i)^\top a_{ki} \phi(x_i) \quad \textcircled{7}$$

Rewrite ⑦ using ⑥:

$$\frac{1}{N} \sum_{i=1}^N k(x_i, x_i) = \sum_{j=1}^N a_{kj} k(x_i, x_j) = \lambda_k \sum_{i=1}^N a_{ki} k(x_i, x_i) \quad ⑧$$

Define $K_{ij} = k(x_i, x_j)$

Rewrite ⑧:

$$\frac{1}{N} \sum_{i=1}^N k_{ei} \sum_{j=1}^N a_{kj} K_{ij} = \lambda_k \sum_{i=1}^N a_{ki} k_{ei} \quad ⑨$$

Left side of ⑨:

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N k_{ei} \sum_{j=1}^N a_{kj} K_{ij} &= \frac{1}{N} \left(\sum_{i=1}^N k_{ei} \right) \left(\sum_{j=1}^N K_{ij} a_{kj} \right) \\ &= \frac{1}{N} \underline{k} \cdot \underline{K} \underline{a}_k \\ &= \frac{1}{N} K^2 a_k \end{aligned}$$

i and l cannot be anything.
by using all the values of
i/l, we will get the vector
 $k \cdot a_k$

Right side of ⑨:

$$\begin{aligned} \lambda_k \sum_{i=1}^N a_{ki} k_{ei} &= \lambda_k \left(\sum_{i=1}^N \underline{k}_{ei} \underline{a}_{ki} \right) \\ &= \lambda_k \underline{K} \underline{a}_k \end{aligned}$$

Hence, rewrite ⑨:

$$\frac{1}{N} K^2 a_k = \lambda_k K a_k \quad ⑩$$

For non-zero eigenvalues,

$$\begin{aligned} \frac{1}{N} K a_k &= \lambda_k a_k \\ K a_k &= N \lambda_k a_k \end{aligned} \quad ⑪$$

Converted into eigens problem!

- Question: Project $\phi(x)$ onto V_k ?

A: $\phi(x) = \underline{(\phi(x)^T V_k)} V_k$

$$\begin{aligned} \phi(x)^T V_k &= \phi(x)^T \sum_{j=1}^N a_{kj} \phi(x_j) \\ &= \sum_{j=1}^N a_{kj} \phi(x)^T \phi(x_j) \\ &= \sum_{j=1}^N a_{kj} k(x, x_j) \end{aligned}$$

• Example:

Given a non-linear mapping: $\phi \left(\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \right) = \begin{pmatrix} u_1^2 \\ u_1 u_2 \\ u_2 u_1 \\ u_2^2 \end{pmatrix}$

compute $\phi^T(u) \cdot \phi(v)$, given $u = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$, $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$

Solution:

$$\begin{aligned} & \begin{pmatrix} u_1^2 & u_1 u_2 & u_2 u_1 & u_2^2 \end{pmatrix} \begin{pmatrix} v_1^2 \\ v_1 v_2 \\ v_2 v_1 \\ v_2^2 \end{pmatrix} \\ &= u_1^2 v_1^2 + 2 u_1 u_2 v_1 v_2 + u_2^2 v_2^2 \\ &= (u_1 v_1 + u_2 v_2)^2 \\ &= (u^T v)^2 \quad \leftarrow \text{no need to form 4-dim vectors to compute } \phi(u)^T \phi(v) \end{aligned}$$

Note: We don't need to know ϕ is explicitly. All we care is being able to compute the kernel.

Question: What is the kernel ($k(x_i, x_j)$)?

Answer: Give some kernel following:

Polynomial kernel: $k(x_1, x_2) = (1 + x_1^T x_2)^m$

Gaussian kernel (RBF): $k(x_1, x_2) = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$

...

Question: What are the conditions to be a kernel?

Answer: Mercer's condition: $k(x, x')$ is valid kernel function if and only if the kernel matrix is always symmetric positive semi-definite for any given $\{x_1, x_2, \dots, x_n\}$

Example: Check $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ is a "qualified" kernel function.

Solution:

Theorem: if K is positive semi-definite, its quadratic form must ≥ 0

Hence, we just need to check if $y^T K y \geq 0$

$$y^T K y = y^T \phi(x_i)^T \phi(x_j) y$$

$$= \sum_{i,j} \phi(x_i)^T \phi(x_j) y_i y_j$$

$$= \left[\sum_i y_i \phi(x_i)^T \right] \sum_j y_j \phi(x_j)$$

$$= \Phi^T y^T y \Phi, \text{ where } \Phi = \begin{pmatrix} \phi(x_1) & \cdots & \phi(x_n) \end{pmatrix}$$

$$= (y \Phi)^T (y \Phi) \geq 0$$