

SVD no PCA

PEDRO PATRÍCIO

2-1

Recorde que temos n "indivíduos", x_1, \dots, x_n .

Cada um com m "inputs" ou atributos.

Assumimos que os dados estão centrados, i.e., c/ média 0.

Podemos escrever o indivíduo i como

$$x_i = \begin{bmatrix} x_i^{(1)} \\ \vdots \\ x_i^{(m)} \end{bmatrix} \text{ onde } x_i^{(j)} \text{ denota}$$

o input do atributo j em x_i .

O atributo j toma, assim, os valores

$$x^{(j)} = (x_1^{(j)}, x_2^{(j)}, \dots, x_n^{(j)})$$

Se considerarmos a matriz B cujas colunas são os vetores x_i , então $x^{(j)}$ será a linha j de B .

A matriz de covariância é

$$S = \frac{1}{n} B B^T = \begin{bmatrix} \text{Cov}(x^{(1)}, x^{(1)}) & \text{Cov}(x^{(1)}, x^{(2)}) & \dots \\ \text{Cov}(x^{(2)}, x^{(1)}) & \text{Cov}(x^{(2)}, x^{(2)}) & \dots \\ \vdots & \vdots & \ddots \\ \text{Cov}(x^{(m)}, x^{(1)}) & \text{Cov}(x^{(m)}, x^{(2)}) & \dots \end{bmatrix}$$

cujas entradas (i, j) iguais $\text{Cov}(x^{(i)}, x^{(j)})$
(= $\text{Cov}(x^{(j)}, x^{(i)})$)

Recorde ainda : 1) $M M^T$ é SDP, i.e., é simétrica e

$$\sigma(M M^T) \subseteq \mathbb{R}_0^+$$

(todos os valores próprios são reais e não-negativos)

2) a soma dos valores próprios de uma matriz

igual a soma dos elementos diagonais

(*) - contando multiplicidade

3) Se $M_{n \times n}$ tem valores próprios $\lambda_1, \dots, \lambda_n$, [2.2]
 então $\alpha M_{n \times n}$ tem valores próprios $\alpha \lambda_1, \dots, \alpha \lambda_n$
 Se v é vector prop. de M assoc. a λ então
 v é vector prop. de αM assoc. a $\alpha \lambda$.
 Ou seja, os valores próprios são afectados por α ,
 mas mantêm os mesmos vectores próprios.
 (repare que v vect. próprio entre αv também
 é vector prop., $\forall \alpha \neq 0$)

4) MM^T e M^TM têm os mesmos valores próprios
 não nulos, contando com a multiplicidade.
 Mais, se v é vector prop. de M^TM assoc.
 a λ então Mv é " " " " MM^T assoc. a λ .
 Se u é vect. prop. de MM^T assoc. a λ então
 M^Tu é " " " " M^TM " " λ .

Este facto (4) pode ser muito relevante na capacidade
 de poder implementar o PCA.

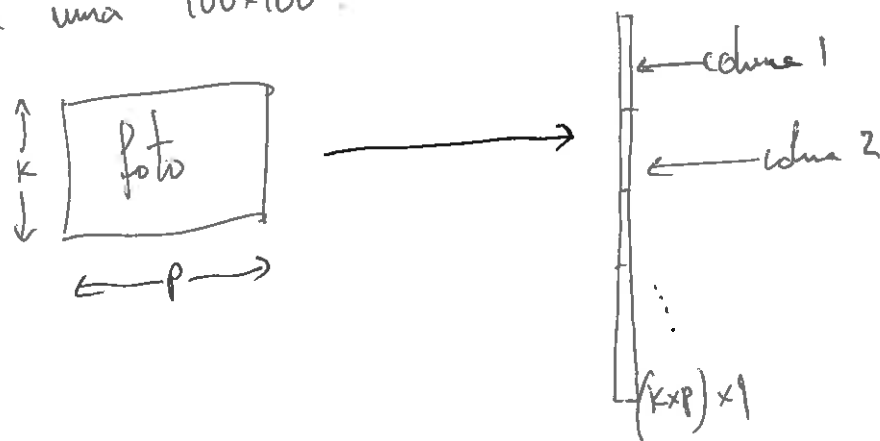
No caso de termos n indivíduos, cada um com m
 atributos, obtemos a matriz de covariância.

$$S = \frac{1}{n} BB^T, \text{ de ordem } m \times m.$$

No caso de termos menos atributos do que indivíduos, é
 preferível usar o PCA de acordo c/ os valores e vectores
 próprios de S .

Suponha agora que tem 400 fotografias ($n=400$)
 cada uma 100×100 .

2.3



obtemos assim 400 fotografias, sendo cada uma um
 vetor com 10^4 atributos. A matriz de covariância
 será $S = \frac{1}{400} BB^T$ onde B é uma matriz $10^4 \times 400$

e portanto S será do tipo $10^4 \times 10^4$.

Por forma a aplicar o PCA, teremos que calcular os maiores
 valores próprios de S e respectivos vetores próprios!

Ao invés de se utilizar BB^T (a menos por um produto
 por um escalar - como vimos sabemos como alterar o espectro)

fazemos o estudo de B^TB . No exemplo, é uma

matriz 400×400 , com os mesmos valores próprios n -úlos
 de BB^T , ... e cujos vetores próprios se relacionam.

Repare que se v é vect. prop. de BB^T , $BB^T v = \lambda v$,

então $(B^TB)(Bv) = \lambda(Bv)$, i.e., Bv é vect. prop. de B^TB ;

reciprocamente, se μ é vtp. de B^TB , $B^TB \mu = \lambda \mu$, então

$(BB^T)(B\mu) = \lambda(B\mu)$ e $B\mu$ é vect. prop. de BB^T .

Considere a matriz $X_{n \times m}$, com característica

2.4

$$\text{rank}(X) = r = \text{rank}(X^T X) = \text{rank}(X X^T)$$

Seja $X^T X$ SDP então $\sigma(X^T X) \subseteq \mathbb{R}_0^+$

Sejam $\lambda_1, \lambda_2, \dots, \lambda_n$ os valores próprios n -últimos de $X^T X$ contando as multiplicidades.

Sejam v_1, \dots, v_m vectores próprios associados aos valores próprios $\lambda_1, \lambda_2, \dots, \lambda_m$, $\lambda_{m+1} = 0, \lambda_{m+2} = 0, \dots, \lambda_n = 0$.

Como $X^T X$ é simétrica, os m vectores próprios são ortogonais 2 a 2. Sem perda de generalidade, sup. v_1, \dots, v_m são ortogonais 2 a 2 e de norma 1. Então seja

$$v_i \perp v_j \text{ se } i \neq j \text{ e } \|v_i\| = 1, \text{ satisfazendo}$$
$$(i.e., v_i^T v_j = 0)$$

$$X^T X v_i = \lambda_i v_i, \quad i = 1, \dots, r$$

$$X^T X v_i = 0, \quad i = r+1, \dots, n \text{ i.e., } v_i \in \text{Ker}(X^T X) \\ \text{Ker}(X)$$

$$\text{Sejam } \sigma_i = \sqrt{\lambda_i}, \quad \mu_i = \frac{1}{\sigma_i} X v_i, \quad i = 1, \dots, r.$$

(Recorde que $X^T X$ é SDP, logo os valores próprios n -últimos são positivos)

Para $i = 1, \dots, r$, temos

$$\|\mu_i\| = \frac{1}{\sigma_i} \|X v_i\| = \frac{1}{\sigma_i} \sqrt{(X v_i)^T (X v_i)} = \frac{1}{\sigma_i} \sqrt{v_i^T X^T X v_i}$$

$$= \frac{1}{\sigma_i} \sqrt{v_i^T \lambda_i v_i^T} = \frac{\sqrt{\lambda_i}}{\sigma_i} \sqrt{v_i^T v_i}$$

2.5

$$= \|v_i\| = 1$$

On seja, $\|u_i\| = 1$, $i=1, \dots, k$

Também temos $u_i \perp u_j$, se $i \neq j$, $i, j=1, \dots, k$

$$\text{De facto, } \langle u_i, u_j \rangle = u_i^T u_j = \frac{1}{\sigma_i} (X v_i)^T \frac{1}{\sigma_j} (X v_j)$$

$$= \frac{1}{\sigma_i \sigma_j} v_i^T \underbrace{X^T X v_j}_{\lambda_j v_j} = \frac{\lambda_j}{\sigma_i \sigma_j} v_i^T v_j = 0$$

$$\text{pois } v_i^T v_j = \langle v_i, v_j \rangle = 0 \text{ por } v_i \perp v_j$$

Temos, portanto, $X v_i = \sigma_i u_i$

$$u_i \perp u_j, \text{ se } i \neq j$$

$$\|u_i\| = 1 \text{ com } i, j=1, \dots, k$$

On seja

$$\begin{cases} X v_1 = \sigma_1 u_1 \\ X v_2 = \sigma_2 u_2 \\ \vdots \\ X v_k = \sigma_k u_k \end{cases} \Leftrightarrow X \begin{bmatrix} v_1 & v_2 & \dots & v_k \end{bmatrix} = \begin{bmatrix} u_1 & u_2 & \dots & u_k \end{bmatrix} \begin{bmatrix} \sigma_1 & \sigma_2 & 0 & \dots \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \sigma_k \end{bmatrix}$$

Recorde que $X v_{k+1} = X v_{k+2} = \dots = X v_m = 0$

$$\text{e portanto } X \begin{bmatrix} v_1 & v_2 & \dots & v_k & v_{k+1} & \dots & v_m \end{bmatrix} = \begin{bmatrix} u_1 & u_2 & \dots & u_k & 0 & \dots & 0 \end{bmatrix} \underbrace{\begin{bmatrix} \sigma_1 & \sigma_2 & \dots & \sigma_k \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}}_{\Sigma}$$

$$\text{onde a matriz } \Sigma = \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ & & \sigma_k & 0 \\ 0 & & & 0 \end{bmatrix}_{n \times m}$$

Note que $\mu_1, \dots, \mu_n \in \mathbb{R}^n$, $n \leq n$
são 2 a 2 ortogonais e de norma 1.

(2.6)

Sabendo que $\mu_i = \frac{1}{\sigma_i} X v_i$, $i=1, \dots, n$, onde μ_i
é vector prop. de $X^T X$ assoc. ao valor próprio $\lambda_i = \sigma_i^2$,

$$(X X^T) \mu_i = X X^T \frac{1}{\sigma_i} X v_i = \frac{1}{\sigma_i} X \underbrace{(X^T X \mu_i)}_{= \sigma_i^2 \mu_i}$$

$$= \sigma_i X v_i = \sigma_i^2 \mu_i = \lambda_i \mu_i$$

On seja, se v_i é vect. prop. $X^T X$ assoc. λ_i então
então μ_i é vect. prop. $X X^T$ assoc. $\lambda_i \neq 0$ (ao mesmo
valor próprio)

onde $i=1, \dots, n$.

Seja $X X^T$ simétrica, então é possível completar
o conjunto $\{\mu_1, \dots, \mu_n\}$ com $n-n$ vectores próprios de
 $X X^T$ assoc. $\lambda = 0$, μ_{n+1}, \dots, μ_n , por forma a que
 $\mu_1, \dots, \mu_n, \mu_{n+1}, \dots, \mu_n$ sejam ortogonais 2 a 2
e cada um com norma 1.

Obtemos assim

$$X \underbrace{\begin{bmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{bmatrix}}_{V_{n \times n}} = \underbrace{\begin{bmatrix} | & & | \\ \mu_1 & \dots & \mu_n \\ | & & | \end{bmatrix}}_{U_{n \times n}} \underbrace{\begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n & & 0 \\ & & & \ddots & \\ 0 & & & & 0 \end{bmatrix}}_{\Sigma_{n \times n}}$$

2.7

$$U^{-1} = V^T$$

$$XV = U\Sigma \Rightarrow X = U\Sigma V^{-1} = U\Sigma V^{-1}$$

$\sigma_1 = 1, \sigma_2$ são os valores singulares,
singulares.

regulares. onde λ_i é valor próprio não nulo de XX^T
 $\sigma_i = \sqrt{\lambda_i}$ onde λ_i " " " " " " " $X^T X$
 i.e. λ_i " " " " " " " $X^T X$

Aplicação no PCA

As principais

Suponha que B é a matriz dos dados indivíduo/inputs da forma usual. Sup. que os dados estão centrados p/ média. (como as linhas indicam os vetores aleatórios, cada linha tem média 0). A matriz de covariância é $S = \frac{1}{n} B B^T$.

Seja $Y = \frac{1}{\sqrt{n}} B^T$; cada coluna terá média 0.

$$Y^T Y = \left(\frac{1}{\sqrt{n}} B^T \right)^T \frac{1}{\sqrt{n}} B^T = \frac{1}{n} B B^T = S$$

Recorde que as Componentes principais de B são os
vetores próprios de S assoc. a valores próprios e/ certa "magnitude".

Aplicando SVD a Y , $Y = U\Sigma V^T$,

2.8

as colunas de V são os vetores próprios de Y^TY

Se $X = B^T$ então $X^TX = nY^TY$, e X^TX e Y^TY têm os mesmos vetores próprios.

Assim as colunas de V são vect. prop. de X^TX

Ou seja, basta aplicar SVD a $X = B^T$.

Obtem-se assim $X = U\Sigma V^T$ onde as colunas de V são vetores próprios de X^TX .

$$B = X^T = (U\Sigma V^T)^T = (U\Sigma V^T)^T = V\Sigma U^T$$

Ou seja, as colunas de V são os componentes principais de B . Aplicar SVD a B ou a $B^T = X$ dependerá da dimensão da matriz.

Projeções!

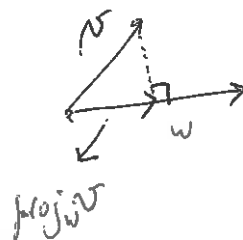
12.9

Qual a importância de termos vectores ortogonais 2 a 2?

Sejam $w, v \in \mathbb{R}^n$. A projecção de v ao longo de w ,

$\text{proj}_w v$ é um vector de \mathbb{R}^n definido por ~~$\text{proj}_w v$~~

$$\text{proj}_w v = \frac{\langle v, w \rangle}{\|w\|^2} w.$$



Dado um espaço vectorial (de dimensão finita) $V \subseteq \mathbb{R}^n$, seja $B = (u_1, \dots, u_k)$ uma base ortogonal de V .

Seja $w \in \mathbb{R}^n$. Define-se projecção ortogonal de w em V como o único $v \in V$ t.q. $v = \arg \min_{v \in V} \|w - v\|$

O facto de B ser uma base ortogonal permite calcular a projecção (ortogonal) de forma fácil:

$$\text{proj}_V w = \text{proj}_{u_1} w + \text{proj}_{u_2} w + \dots + \text{proj}_{u_k} w$$

Podemos, assim, calcular a projecção do novo "indivíduo" nas componentes principais, já que são vectores próprios ortogonais 2 a 2 e cada um de norma 1.

Sejam v_1, \dots, v_k Componentes principais, i.e.,
 vectores próprios de BB^T assoc. aos maiores k
 valores próprios da matriz de covariância.

Seja ϕ um novo input, que assumimos centrado
 na média.

Para calcular a projecção de ϕ em $\langle v_1, \dots, v_k \rangle$,
 espaço gerado por v_1, \dots, v_k , como v_i são ortogonais 2a2
 e de norma 1, $\text{proj}_{\langle v_1, \dots, v_k \rangle} \phi = \sum_{i=1}^k \text{proj}_{v_i} \phi$
 $= \sum_{i=1}^k \langle \phi, v_i \rangle v_i = \sum_{i=1}^k (\phi^T v_i) v_i$

Seja c o vector dos coeficientes da projecção,
 $c = \begin{bmatrix} \phi^T v_1 \\ \phi^T v_2 \\ \vdots \\ \phi^T v_k \end{bmatrix}$ Então $\text{proj}_{\langle v_1, \dots, v_k \rangle} \phi = \begin{bmatrix} v_1 & \dots & v_k \end{bmatrix} \begin{bmatrix} \phi^T v_1 \\ \vdots \\ \phi^T v_k \end{bmatrix}$
 $= V_k c$

Este processo permite "identificar" novos indivíduos. Para tal
 basta projectar o novo indivíduo ϕ (depois de centrado) e
 procurar qual o indivíduo conhecido cuja projecção minimiza
 a distância.

On seja :

Sejam v_1, \dots, v_K as K componentes principais mais representativas (i.e., as K maiores valores próprios)

Seja $V_K = [v_1 \dots v_K]$ (pode ser obtida por SVD)

Seja ϕ um novo indivíduo (centrado na média!)

Sejam $x_1, \dots, x_n \in \mathbb{R}^m$ os elem^{tos} da amostra inicial, centrados
(i.e., $B = [x_1 \dots x_n]$)

Se denotarmos $CS(V_K)$ o espaço gerado pelas colunas de V_K , então

$$\text{proj}_{CS(V_K)} x_i = \langle x_i, v_1 \rangle v_1 + \langle x_i, v_2 \rangle v_2 + \dots + \langle x_i, v_K \rangle v_K$$

Seja c_i o vector dos coeficientes da projecção

$$c_i = \begin{bmatrix} \langle x_i, v_1 \rangle \\ \langle x_i, v_2 \rangle \\ \vdots \\ \langle x_i, v_K \rangle \end{bmatrix} = \begin{bmatrix} x_i^T v_1 \\ x_i^T v_2 \\ \vdots \\ x_i^T v_K \end{bmatrix}_{K \times 1} \in \mathbb{R}^K, \text{ e } K \ll m$$

Façamos o mesmo para ϕ :

$$\text{proj}_{CS(V_K)} \phi = \sum_{i=1}^K \langle \phi, v_i \rangle v_i \quad \text{e} \quad \text{seja}$$

$$c \text{ o vector dos coeficientes, } c = \begin{bmatrix} \phi^T v_1 \\ \vdots \\ \phi^T v_K \end{bmatrix}$$

Identificamos ϕ com x_j se

$$c_j = \arg \min d(c_i, c)$$

e $d(c_j, c) < L$
onde L é uma certa tolerância

Podemos usar a distância euclidiana

(2.12)

$$d(v, w) = \|v - w\|_2 = \sqrt{\langle v - w, v - w \rangle}$$

No entanto, há relatos em como a distância de Mahalanobis tem um melhor comportamento:

$$d_H\left(\begin{bmatrix} c_1 \\ \vdots \\ c_k \end{bmatrix}, \begin{bmatrix} \hat{c}_1^i \\ \hat{c}_2^i \\ \vdots \\ \hat{c}_k^i \end{bmatrix}\right) = \sum_{j=1}^k \frac{1}{\lambda_j} (c_j - \hat{c}_j^i)^2$$

Intuitivamente, dá-se mais importância às entradas correspondentes aos maiores valores próprios, quando pretendemos minimizar a distância.

• Dada $X_{n \times m}$, prova-se como

2.x

$$\text{Ker}(X^T X) = \text{Ker}(X) \quad \text{onde } \text{K}(A) \text{ denota o núcleo de } A : \text{Ker}(A) = \{v : Av = 0\}$$

$$\text{Se } v \in \text{Ker}(X) \text{ então } Xv = 0 \Rightarrow X^T X v = 0 \Rightarrow v \in \text{Ker}(X^T X)$$

Reciprocamente se $v \in \text{Ker}(X^T X)$ então

$$X^T X v = 0 \Rightarrow v^T X^T X v = 0 \quad (\text{multiplicando à esquerda por } v^T)$$

$$\Rightarrow (Xv)^T Xv = 0$$

$$\Rightarrow \langle Xv, Xv \rangle = 0 \Rightarrow \|Xv\| = 0$$

$$\Rightarrow Xv = 0 \Rightarrow v \in \text{Ker}(X)$$

• O ponto anterior permite mostrar que $\text{rank}(X) = \text{rank}(X^T X)$.

Recorde-se que se X é $n \times m$ então

$$m = \text{rank}(X) + \dim \text{Ker}(X)$$

Ora $X^T X$ é $m \times m$, o que leva a

$$m = \text{rank}(X^T X) + \dim \underbrace{\text{Ker}(X^T X)}_{= \text{Ker}(X)} = \text{rank}(X^T X) + \dim \text{Ker}(X)$$

$$\text{logo } \text{rank}(X) = \text{rank}(X^T X)$$