

PCA

We use notations from the following reference:

Goodfellow-et-al-2016,

Deep Learning,

Jan Goodfellow and Yoshua Bengio and Aaron Courville,
MIT press,

<http://www.deeplearningbook.org>

year=2016

it's a convention:
an hypothesis; in PCA

$D \in \mathbb{R}^{n \times \ell}$ decode matrix

$c \in \mathbb{R}^{\ell}$ code of a line

$X \in \mathbb{R}^{m \times n}$ data

$$\textcircled{H} D^T D = I_{\ell}$$

matrix dimensions
provide valuable
information for
checking
calculations

Reconstruction function: $\pi: \mathbb{R}^n \rightarrow \mathbb{R}^n$
 $x \mapsto DD^T x$

$$\text{Our aim is to minimize } \sum_{j=1}^n \sum_{i=1}^m \left(x_{ij}^{(ci)} - (DD^T x^{(ci)})_j \right)^2 = (*)$$

! Be careful: $(DD^T x^{(ci)})_j \neq DD^T x_{ij}^{(ci)}$

$$(*) = \sum_j \sum_i x_{ij}^2 - 2 \sum_j \sum_i x_{ij} (DD^T x_{ij}^T)_j + \sum_j \sum_i (DD^T x_{ij}^T)_j^2$$

So we are looking for:

$$D^* = \underset{D}{\operatorname{argmin}} \quad -2 \sum_j \sum_i x_{ij} (DD^T x_{ij}^T)_j + \sum_j \sum_i (DD^T x_{ij}^T)_j^2$$

s.t. $D^T D = I_{\ell}$

always explicit matrix products

$$\textcircled{1} \sum_j^n \sum_i^m X_{ij} (DD^T X_{i,:}^T)_j = \sum_j^n \sum_i^m X_{ij} \sum_{k,l} D_{jk} D_{lk} X_{il}$$

put the X's in front because we'll differentiate w.r.t. later

$$= \sum_j^n \sum_i^m \sum_{k,l} X_{ij} X_{il} D_{jk} D_{lk}$$

$$= \sum_j^n (X^T X D D^T)_{jj}$$

$$= \text{Tr}(X^T X D D^T)$$

use definition of matrix trace!

organize indexes to prepare the product:
 $j \rightarrow i \rightarrow l \rightarrow k \rightarrow j$
 $(X^T)_{ji} X_{il} D_{lk} (D^T)_{kj}$

② Exercise: I let you do the same with

$$\sum_j^n \sum_i^m (DD^T X_{i,:}^T)_j^2$$

We obtain $\text{Tr}(X^T X \underbrace{DD^T DD^T}_{\textcircled{H} = I_l}) = \text{Tr}(X^T X DD^T)$

$$\begin{aligned} \rightarrow \mathcal{D}^* &= \underset{\substack{D \\ \text{s.t. } D^T D = I_l}}{\text{argmin}} - \text{Tr}(X^T X DD^T) \\ &= \underset{D^T D = I_l}{\text{argmax}} \text{Tr}(X^T X DD^T) \end{aligned}$$

We write the Lagrangian $\Delta \in \mathbb{R}^{l \times l}$

$$\mathcal{L}(D, \Delta) = \text{Tr}(X^T X DD^T) - \text{Tr}(\Delta (DD^T - I_l))$$

It is easy to see that in the formula $D^T D = I_l$, there are l^2 constraints by only $\frac{l(l+1)}{2}$ useful constraints because the formula is entirely symmetric.

Therefore, we could decide that Λ is upper triangular, but it seems more interesting to choose Λ symmetric: $\Lambda^T = \Lambda$

It's a choice.

Now let's compute $\mathcal{D}_\Lambda^*(D)(H)$

$$\begin{aligned}\mathcal{L}(D+H, \Lambda) &= \text{Tr}(X^T X (D+H)(D+H)^T) - \text{Tr}(\Lambda (D+H)^T (D+H) - I_2) \\ &= \mathcal{L}(D, \Lambda) + \text{Tr}(X^T X H D^T) + \text{Tr}(X^T X D H^T) \\ &\quad - \text{Tr}(\Lambda H^T D) - \text{Tr}(\Lambda D^T H) + o(\|H\|) \\ &= \mathcal{L}(D, \Lambda) + \text{Tr}(D^T X^T X H) + \text{Tr}(H D^T X^T X) \\ &\quad - \text{Tr}(D^T H \Lambda^T) - \text{Tr}(\Lambda D^T H) + o(\|H\|) \\ &= \mathcal{L}(D, \Lambda) + \text{Tr}(D^T X^T X H) + \text{Tr}(D^T X^T X H) \\ &\quad - \text{Tr}(\Lambda^T D^T H) - \text{Tr}(\Lambda D^T H) + o(\|H\|)\end{aligned}$$

(KKT)

$$\begin{aligned}\text{So } D^* X^T X &= \Lambda^T D^* \\ X^T X D^* &= D^* \Lambda \\ \Lambda &= D^{*T} X^T X D^*\end{aligned}$$

- Let λ be an eigenvalue of $X^T X$ and v an eigenvector for this eigenvalue

$$D^{*T} X^T X v = \Lambda^T D^{*T} v \quad \text{so} \quad \lambda D^{*T} v = \underbrace{\Lambda^T}_{\Lambda} D^{*T} v$$

$$\text{So } v \in \text{Ker } D^{*T} \text{ or } \lambda \in \text{Sp}(\Lambda)$$

$$\text{Conversely } \tilde{\lambda}, \tilde{v} \text{ eigenvalue/vector of } \Lambda \Rightarrow \tilde{v} \in \text{Ker } D^{*T} \text{ or } \tilde{\lambda} \in \text{Sp}(X^T X)$$

In the meantime, we see that we are solving

$$D^* = \underset{\Lambda \in \text{Sp}(X^T X)}{\text{argmax}} \text{Tr}(\Lambda) \quad \text{with } \text{Sp}(\Lambda) \subset \text{Sp}(X^T X)$$

$$\begin{aligned}\Lambda D^T &= D^T X^T X \\ D^T D &= I_2\end{aligned}$$

But the trace is the sum of the eigenvalues

Therefore, one best Λ is the diagonal matrix of the first l biggest eigenvalues of $X^T X$

not
necessarily
unique

achieves

$$\text{Tr}(\Lambda) = \max_{D^T D = I_l} \text{Tr}(X^T X D D^T)$$

Let's write the eigenvalue equation for $X^T X$:

$$\text{Diag}(\lambda_1, \dots, \lambda_n) = Q X^T X Q^T$$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ \uparrow

To remind $Q \in \mathbb{R}^{n \times n}$

$Q^T Q = I_n$] orthogonal matrix

Q orthogonal matrix

$\in \mathbb{R}^{n \times n}$ such

that $Q_{i,:}^T$ is

a unit vector corresponding to λ_i

Finally we just have to write D^* as the first l rows of Q .