

## PCA algorithms proof

Friday, June 26, 2020 1:47 AM

So the problem is now

$$c^* = \underset{c}{\operatorname{argmin}} (-2x^T D c + c^T c)$$

As our objective function is always increasing then we can solve using derivatives that means

$$\begin{aligned} \nabla_c (-2x^T D c + c^T c) \\ = -2D^T x + 2c \end{aligned}$$

as  $c^*$  is the value where our objective function reaches the minimum then

$$\begin{aligned} -2D^T x + 2c^* &= 0 \\ 2c^* &= 2D^T x \\ c^* &= D^T x \end{aligned}$$

$\Rightarrow f(x) = c = D^T x \rightarrow$  coding function

$g(c) = Dc \rightarrow$  decoding function

$r(x) = g(f(x)) = DD^T x \rightarrow$  reconstruction function

Now we need to know what is the matrix  $D$ .  
Knowing the matrix  $D$  we can map each  $x^{(i)}$  to  $c^{(i)}$ ,  $\mathbb{R}^n \rightarrow \mathbb{R}^p$  with  $f$  and reconstruct the value again with a small error.

As is the same matrix  $D$  for all  $x^{(i)}$  in  $\{x^{(1)}, \dots, x^{(n)}\}$  then we will no longer make the analysis with single  $x^{(i)}$ . Instead let's define

$$X = \begin{pmatrix} x^{(1)T} \\ \vdots \\ x^{(n)T} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

$$X = \begin{pmatrix} x^{(1)} \\ \vdots \\ x^{(m)} \end{pmatrix} \quad n \times 1$$

$$D \in \mathbb{R}^{n \times l}$$

before

$$\underbrace{f(x)}_{n \times l} = \underbrace{D}_{n \times l} \underbrace{x}_{n \times 1}$$

$$\left( \begin{pmatrix} d^{(1)T} \\ \vdots \\ d^{(l)T} \end{pmatrix} \right) x = \begin{pmatrix} d^{(1)T} x \\ \vdots \\ d^{(l)T} x \end{pmatrix}_{l \times 1}$$

$$\Rightarrow d^{(1)T} x \text{ is a dot vector product}$$

$$d^{(1)} \in \mathbb{R}^{l \times 1}$$

$$x \in \mathbb{R}^{n \times 1}$$

Now

$$f(X) = XD$$

$$X \in \mathbb{R}^{m \times n}, D \in \mathbb{R}^{n \times l}$$

$$C = f(X) \in \mathbb{R}^{m \times l}$$

$$\begin{pmatrix} x^{(1)T} \\ \vdots \\ x^{(m)T} \end{pmatrix}_{m \times n} \begin{pmatrix} d^{(1)} & \dots & d^{(l)} \end{pmatrix}_{n \times l} = \begin{pmatrix} x^{(1)T} d^{(1)} & \dots & x^{(1)T} d^{(l)} \\ \vdots & \ddots & \vdots \\ x^{(m)T} d^{(1)} & \dots & x^{(m)T} d^{(l)} \end{pmatrix}_{m \times l}$$

$$= \begin{pmatrix} f(x^{(1)})^T \\ \vdots \\ f(x^{(m)})^T \end{pmatrix} = \begin{pmatrix} x^{(1)T} D \\ \vdots \\ x^{(m)T} D \end{pmatrix}$$

We can see that each row is the corresponding mapping for  $x^{(i)}$   $\forall i \in \{1, \dots, m\}$   
each row is  $c^{(i)}$   $\forall i \in \{1, \dots, m\}$ .

$$f(X) = XD = C \quad \text{where } C \text{ is the reduce matrix of } X$$

Now the decode function for all  $C$  will be

$$g(c^{(i)}) = D c^{(i)} \quad D \in \mathbb{R}^{n \times l} \quad c^{(i)} \in \mathbb{R}^{l \times 1}$$

$$= c^{(i)T} D^T$$

$$\Rightarrow \begin{pmatrix} c^{(1)T} \\ \vdots \\ c^{(m)T} \end{pmatrix}^T = \begin{pmatrix} c^{(1)T} D^T \\ \vdots \\ c^{(m)T} D^T \end{pmatrix} = \begin{pmatrix} g(c^{(1)})^T \\ \vdots \\ g(c^{(m)})^T \end{pmatrix}$$

$$\begin{pmatrix} \vdots \\ c^{(m)T} \end{pmatrix} / 1 \Big) = \begin{pmatrix} \vdots \\ c^{(m)T} D^T \end{pmatrix} / \begin{pmatrix} \vdots \\ g(c^{(m)})^T \end{pmatrix}$$

$$\Rightarrow g(C) = C D^T \quad \Rightarrow r(X) = X D D^T$$

Now that we have our mapping configure to use over all of  $x^{(i)}$  in the collection, to get the value of  $D$  we need to optimize the Frobenius Norm of the difference between  $X$  and  $r(X)$

$$D^* = \underset{D}{\operatorname{argmin}} \|X - X D D^T\|_F \quad D^T D = I_\ell$$

$$\Rightarrow \underset{D}{\operatorname{argmin}} \|X - X D D^T\|_F^2$$

$$\begin{aligned} \|X\|_F^2 &= \operatorname{Tr}(A^T A) \Rightarrow \\ (ABC)^T &= C^T B^T A^T \end{aligned} \quad \Rightarrow \underset{D}{\operatorname{argmin}} \operatorname{Tr} \left( \underbrace{(X - X D D^T)^T (X - X D D^T)}_{(*)} \right) \quad (1)$$

$$(*) \quad (X - X D D^T)^T (X - X D D^T) = X^T X - X^T X D D^T - D D^T X^T X + D D^T X^T X D D^T$$

Replace in (1)

$$\Rightarrow \underset{D}{\operatorname{argmin}} \operatorname{Tr}(\overset{\text{constant}}{X^T X}) - \operatorname{Tr}(X^T X D D^T) - \operatorname{Tr}(D D^T X^T X) + \operatorname{Tr}(D D^T X^T X D D^T)$$

$$\begin{aligned} \operatorname{Tr}(A^T B) &= \operatorname{Tr}(A) \operatorname{Tr}(B) \\ \text{if } AB \text{ and } BA \text{ are possible} &\Rightarrow \operatorname{Tr}(AB) = \operatorname{Tr}(BA) \end{aligned}$$

$$\Leftrightarrow \underset{D}{\operatorname{argmin}} -\operatorname{Tr}(X^T X D D^T) - \operatorname{Tr}(X^T X D D^T) + \operatorname{Tr}(X^T X D D^T D D^T) \quad \text{remember } D^T D = I_\ell$$

$$\Rightarrow \underset{D}{\operatorname{argmin}} -2\operatorname{Tr}(X^T X D D^T) + \operatorname{Tr}(X^T X D D^T)$$

$$\Rightarrow \underset{D}{\operatorname{argmin}} [-\operatorname{Tr}(X^T X D D^T)]$$

one last rearrangement.  
remember that  
 $X \in \mathbb{R}^{m \times n}$   
 $D \in \mathbb{R}^{n \times r}$

$$\Rightarrow \operatorname{argmax}_D \operatorname{Tr}(D^T X^T X D) \quad (2a)$$

Let Analyze (2a)

$$D^T X^T X D$$

$\Rightarrow X^T X$  is a symmetric Matrix

$$(X^T X)^T = X^T X^{TT} = X^T X \quad \checkmark$$

$\Rightarrow$  the eigenvalue decomposition is

$$X^T X = U \Lambda U^T$$

where  $U$  is a orthogonal matrix with each column is the eigenvector  
 $\Leftarrow$  and  $\Lambda$  is a diagonal Matrix with the singular values

So if we create a function

$$u(X^T X) = U^T X^T X U$$

with  $U$  having the properties of  $D$  the maximum value for the

trace correspond to the sum of maximum eigenvalues of the matrix  $X^T X$ , so and value is that

$$\operatorname{Argmax}_D \operatorname{Tr}(D^T X^T X D)$$

proof by induction on  $r$   
number of columns of  $D$ .

$$\text{if } r=1 \Rightarrow D = d \in \mathbb{R}^{n \times 1}$$

$$\operatorname{Argmax}_d \operatorname{Tr}(d^T X^T X d) \quad d^T d = 1$$

$\Rightarrow$  the best value for  $d$  is the eigenvector of  $X^T X$  corresponding to the largest eigenvalue of  $X^T X$

Suppose for  $r \Rightarrow D \in \mathbb{R}^{n \times r}$  (h.i)

$$\operatorname{Argmax} [\operatorname{Tr}(D^T X^T X D)] \quad D^T D = I_r$$

D

it holds that the best value for D are the corresponding eigenvectors for the largest  $l$  eigenvalues of  $X^T X$ .

Proof for  $l+1 \Rightarrow D \in \mathbb{R}^{n \times (l+1)}$

$\text{Argmax}_D \text{Tr}(D^T X^T X D)$  using hypothesis induction

the first  $l$  columns need to be

the first eigenvectors corresponding to the first eigenvalues

$$\begin{aligned} & \begin{pmatrix} d^{(1)T} \\ \vdots \\ d^{(l+1)T} \end{pmatrix}_{(l+1) \times n} \begin{pmatrix} x^{(1)} & \dots & x^{(m)} \end{pmatrix}_{n \times m} \\ & \text{Tr} \left( \begin{pmatrix} d^{(1)T} x^{(1)} & \dots & d^{(1)T} x^{(m)} \\ \vdots & \ddots & \vdots \\ d^{(l+1)T} x^{(1)} & \dots & d^{(l+1)T} x^{(m)} \end{pmatrix} \begin{pmatrix} x^{(1)T} d^{(1)} & \dots & x^{(m)T} d^{(1)} \\ \vdots & \ddots & \vdots \\ x^{(m)T} d^{(1)} & \dots & x^{(m)T} d^{(l+1)} \end{pmatrix} \right) \end{aligned}$$

$$= \underbrace{\left( d^{(1)T} X^T X d^{(1)} + \dots + d^{(l+1)T} X^T X d^{(l+1)} \right)}_l \text{ where } \lambda^{(1)} \geq \lambda^{(2)} \geq \dots \geq \lambda^{(l)} \rightarrow (l, 1)$$

= the first  $l$  is the corresponding Matrix  $D \in \mathbb{R}^{n \times l}$  where each  $d^{(1)} \dots d^{(l)}$  are the eigenvectors corresponding to the first  $l$  largest eigenvalues then

$d^{(l+1)}$  must be the  $(l+1)$  eigenvector corresponding to the  $(l+1)$  eigenvalue of  $X^T X$  to maximize the trace.