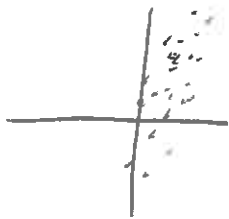


## Lecture #2: 9-6-18

①

• Suppose we have some data  $X_1, \dots, X_n \in \mathbb{R}^D$



• When  $D$  is large, hard to visualize and do statistics (curse of dimensionality).

• We can try to reduce the dimension of the data while retaining important properties of the data.

• Let  $X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times D}$

• Suppose  $U$  is a  $D \times D$  matrix whose columns are orthogonal:  $U = \begin{pmatrix} | & & | \\ u_1 & \dots & u_D \\ | & & | \end{pmatrix}$

$$\text{s.t. } \underbrace{u_i \cdot u_j}_{\text{dot product}} = \underbrace{u_i^T u_j}_{\text{matrix product}} = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$$

Then we can write ~~any~~ any  $x \in \mathbb{R}^D$  as  $x = \sum_{i=1}^D \alpha_i u_i$ . We can even compute  $\alpha_i$  via orthogonality:

$$\Rightarrow x u_j^T = \alpha_j, \quad \forall j = 1, \dots, D$$

$$\Rightarrow \alpha = (\alpha_1, \dots, \alpha_D) \text{ may be computed as } \alpha = U^T x.$$

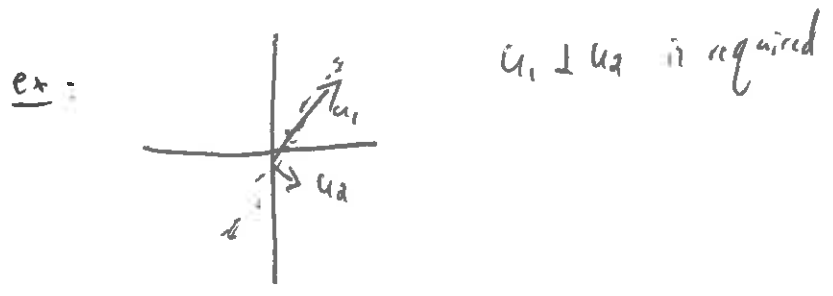
• This is valid for an orthogonal matrix  $U$ , or equivalently, for any orthonormal basis for  $\mathbb{R}^D$  ~~the~~  $u_1, \dots, u_D$ .

**Q**: Which is "best"?

**A**: Depends on how you define "best".

One method - principal component analysis - seeks a basis  $u_1, \dots, u_D$  s.t.

- 1) The projection of  $X \in \mathbb{R}^{n \times D}$  onto  $u_1$  is variance ~~maximizing~~ <sup>maximizing</sup> over all projections onto 1-dimensional spaces.
- 2) The projection of  $X$  onto  $\text{span}\{u_1, u_2\}$  is variance maximizing over all projections onto 2-dimensional spaces.
- etc.



Let's make this rigorous by recalling the definition of variance: for data  $x_1, \dots, x_n \in \mathbb{R}$ , the variance of the data is  $\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$ , where  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$  is the mean.

Let ~~u\*~~  $u_*^T \in \mathbb{R}^{1 \times D}$  be ~~arbitrary~~ the maximizer of  $\frac{1}{n} \sum_{i=1}^n (u_*^T x_i - \mu)^2$ , where  $x_i \in \mathbb{R}^{D \times 1}$  and  $\mu \in \mathbb{R}^{D \times 1}$  is the mean of  $\{x_1, \dots, x_n\}$ .

Suppose, without loss of generality, that  $\mu = 0$ , so

$$u_* = \arg \max_u \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

Let's analyze:  $(u^T x_i)^2 = (u^T x_i) \cdot (u^T x_i)^T$ , since there are just numbers  
 $= u^T x_i \cdot x_i^T \cdot u$ , since  $(AB)^T = B^T A^T$  and  $(A^T)^T = A$

Then,  $u_* = \arg \max_u \frac{1}{n} \sum_{i=1}^n u^T x_i x_i^T u$

$= \arg \max_u u^T \left[ \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right] u \in \mathbb{R}^{D \times D}$

$= \arg \max_u u^T \Sigma u$ , where  $\Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$  is the empirical covariance matrix of  $x_1, \dots, x_n$ .

Remark:  $\Sigma = \frac{1}{n} X^T X$ , where  $X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times D}$

So, how to solve  $\arg \max_u u^T \Sigma u$ ? Derivative! Indeed,

the function  $F: \mathbb{R}^D \rightarrow \mathbb{R}$   
 $u \mapsto u^T \Sigma u$

is smooth, so we can differentiate. However,  $u$  is constrained since  $u^T u = 1$ .

Q = How to constrained optimization? Lagrange multiplier:

$\arg \max_u u^T \Sigma u - \lambda (u^T u - 1)$

to optimize

constraint

- Differentiate and set equal to 0 vector:  $\frac{\partial}{\partial u} [u^T \Sigma u - \lambda (u^T u - 1)] = 0$

$2 \Sigma u - 2 \lambda u = 0$

↙  
exercise

$$\Leftrightarrow \sum u = \lambda u$$

=> 1)  $u$  is an eigenvector of  $\sum$

2)  $\lambda$  is an eigenvalue of  $\sum$

• Conclusion: the potential maximizers are the eigenvectors of  $\sum$ ! Which is

best? Well,  $\sum u = \lambda u$

$$\Rightarrow u^T \sum u = u^T \lambda u$$

$$= \lambda u^T u$$

$$= \lambda$$

So, we should choose  $(u, \lambda)$  the eigenpair corresponding to the largest eigenvalue.

• In general, the principal components of mean-centered data  $x_1, \dots, x_n \in \mathbb{R}^{D \times 1}$ ,  
 i.e.  $\frac{1}{n} \sum_{i=1}^n x_i = 0 \in \mathbb{R}^D$  are the eigenvectors of  $\sum = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$   
 $= \frac{1}{n} X X^T$ , where  $X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times D}$

- Call them  $u_1, \dots, u_D$ . -  $u_1$  points in the direction of maximum variance of  $x_1, \dots, x_n$
- $u_2$  points in the ~~direction~~ " " among those directions orthogonal to  $u_1$
- $u_3$  " " " " among those directions orthogonal to  $u_1$  and  $u_2$
- etc.

The larger the eigenvalue associated to principal component  $u_i$ , the more significant the direction.

- Next time: - geometric considerations
- computational complexity
- numerical examples