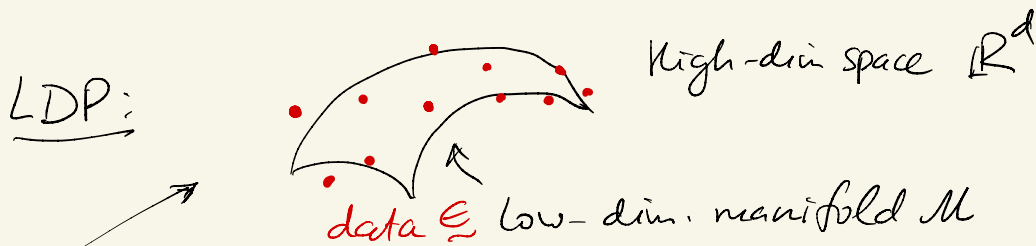


# LECTURE 23

LD Paradigm: High-dimensional data has low-dimensional structure.

This allows us to visualize, think about data, world;  
lift the curse of high dimensionality.

- Ex
- Human decisions are based on  $\sim 5-10$  values
  - Human faces have  $\sim 5-10$  features we recognize.

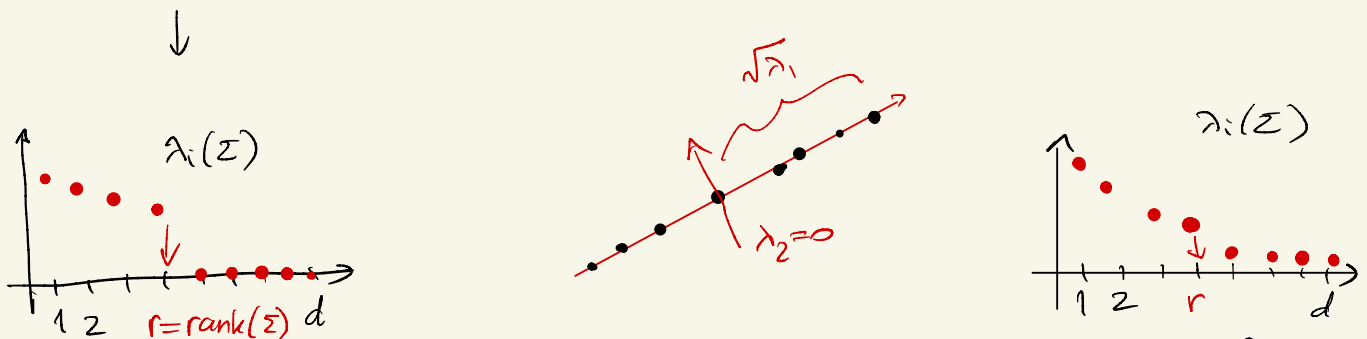


How do we find, parametrize  $M$ ?

How does our brain do it (in vision, thinking)?

How do we know it is true? From the **eigenvalues of  $\Sigma = \text{Cov}(X) = \mathbb{E}XX^T$** :

- If data  $X \in M \leftarrow$  linear subspace of dimension  $r \ll n$   
then  $\Sigma$  has  $\leq r$  nonzero eigenvalues:



- If data  $X$  lies close to  $M$ , the eigs drop at  $r$  (but not exactly to 0)

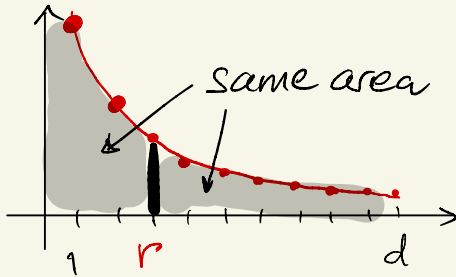
$r =$  "effective rank" of  $\Sigma$

"effective dimension" of the data.

Q Mathematically, how to define the "effective rank" of  $\Sigma$ ,  
 = "effective dimension" of data?

Ans Find  $r$  that splits the area into two equal halves?

Rigorously:



Def The **effective rank** of a PSD matrix  $\Sigma$  is

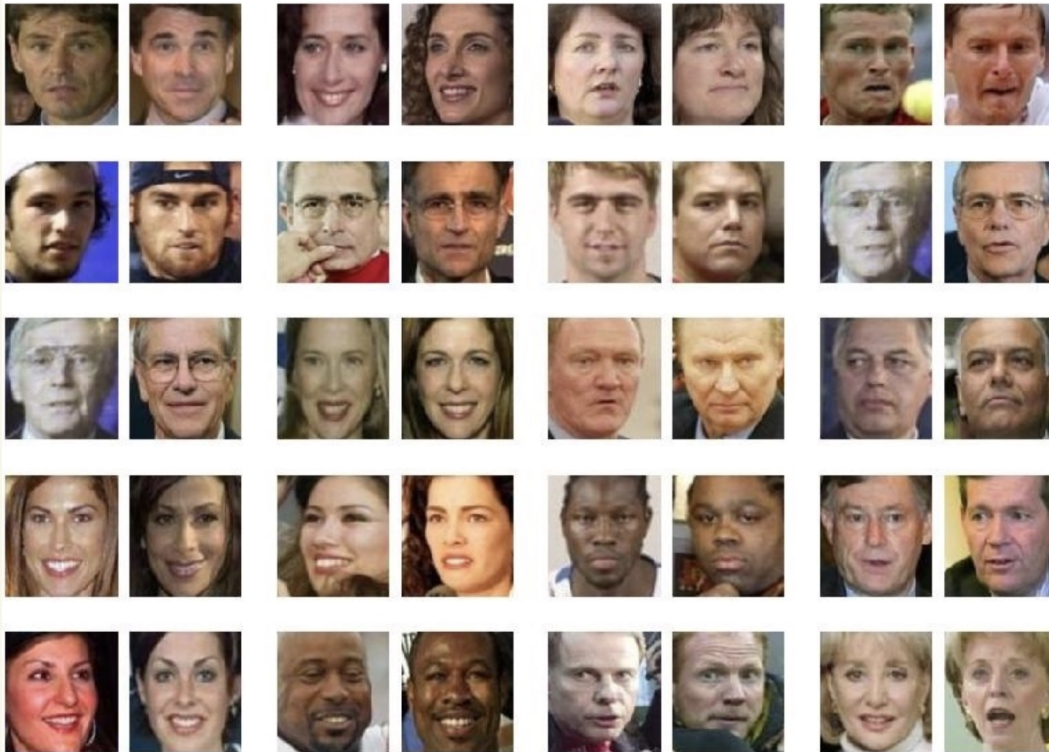
$$r(\Sigma) := \min \left\{ r : \sum_{i \leq r} \lambda_i \geq \sum_{i > r} \lambda_i \right\}$$

↑  
eigs

= effective dimension of  $x$  if  $\Sigma = \text{Cov}(x)$ .

Explains 50% of the variation of the data.

Ex (Eigenfaces) Data  $X \sim \text{Unif}\{\text{images of human faces}\}$



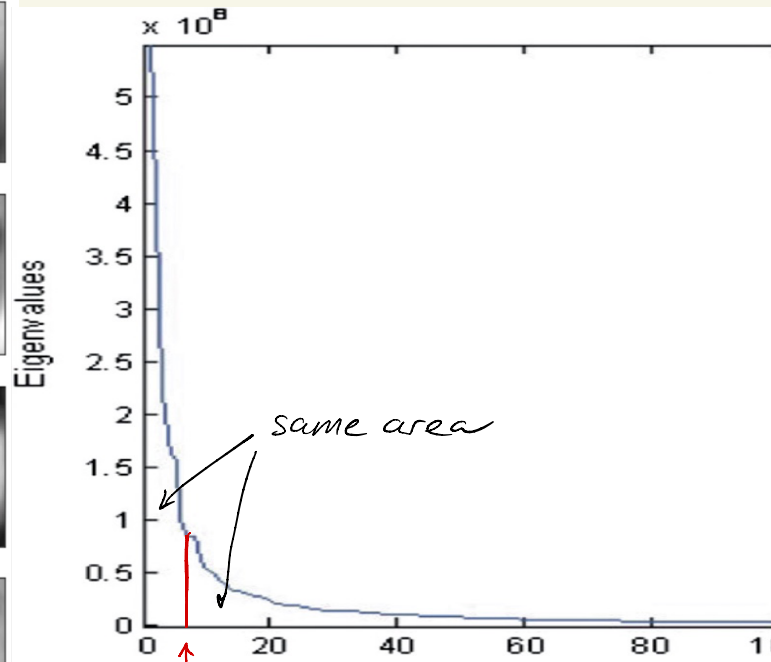
$$X_i \in \mathbb{R}^{50,000 = d}$$

$$n = 2000$$

$$\Sigma_n = \frac{1}{n} \sum_{i=1}^n X_i X_i^T$$

Eigenvectors of  $\Sigma$  = "eigenfaces"

eigenvalues of  $\Sigma$



6 facial features explain 50% of variability

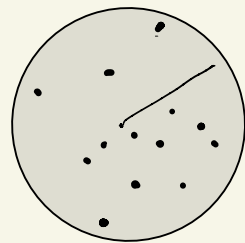
- LD Paradigm  $\Rightarrow$  in all ML results, the **dimension  $d$**  of the data can be replaced by the **effective dimension  $r$**  ( $r \ll d$ ).
- Example: Covariance estimation  $\Rightarrow$  PCA.  
we proved: a sample of size  $n = O(d)$  suffices (lec...)  
We will prove:  $n = O(r)$  suffices: 😊

remove bad outliers prior to applying

Thm Let  $X$  be a <sup>mean zero</sup> random vector in  $\mathbb{R}^d$  such that

$$\|X\|_2^2 \leq 10 \mathbb{E}\|X\|_2^2 \text{ a.s.}$$

If  $n \geq Cr \log d$  then

$$\|\Sigma_n - \Sigma\| \leq 0.1 \|\Sigma\|$$


effective dim.

$$\frac{1}{n} \sum_{i=1}^n X_i X_i^T = \mathbb{E} X X^T$$

sample covariance;  $X_i = \text{iid}$  copies of  $X$

Toward the proof:

$$\text{Lem } \mathbb{E} \|X\|_2^2 \leq 2r \|\Sigma\|$$

$$\mathbb{E} X^T X = \mathbb{E} \text{tr}(X^T X) \stackrel{\text{cyclic property of trace}}{=} \mathbb{E} \text{tr}(X X^T) \stackrel{\text{linearity}}{=} \text{tr} \mathbb{E}[X X^T] = \text{tr}(\Sigma)$$

$$= \sum_{i=1}^d \lambda_i \implies \underbrace{\sum_{i \leq r} \lambda_i}_{\geq} + \sum_{i > r} \lambda_i \leq 2 \sum_{i \leq r} \lambda_i \leq 2r \|\Sigma\|$$

det of  $r$

$\lambda_1 = \|\Sigma\|$

To prove Thm, use matrix Bernstein inequality (last class):

$\forall$  independent mean zero  $d \times d$  symmetric random matrices  $Z_i$  that satisfy  $\|z_i\| \leq 1$  a.s., we have

$$\left\| \sum_{i=1}^n Z_i \right\| \leq C \sigma \sqrt{\log d} + CK \log d \text{ with prob } \geq 0.99$$

where  $\sigma^2 = \left\| \sum_{i=1}^n \mathbb{E} Z_i^2 \right\|$ .



# Proof of Thm

$$\|\Sigma_n - \Sigma\| = \left\| \frac{1}{n} \sum_{i=1}^n \underbrace{(x_i x_i^T - \Sigma)} \right\|$$

$\uparrow$  iid mean 0 random matrices  $\Rightarrow$  use matrix Bernstein inequality:

$$\lesssim \frac{1}{n} (\sigma \sqrt{\log d} + K \log d) \quad \text{if } \|x_i x_i^T - \Sigma\| \leq K \text{ a.s., where } (*)$$

$\uparrow$   
hides an absolute constant factor

$$\sigma^2 = \left\| \sum_{i=1}^n \overbrace{\mathbb{E}(x_i x_i^T - \Sigma)^2}^{\text{all equal}} \right\| = n \|\mathbb{E}(X X^T - \Sigma)^2\|$$

$$\begin{aligned} 0 \leq \mathbb{E}(X X^T - \Sigma)^2 &= \mathbb{E}(X X^T)^2 - \underbrace{\mathbb{E}(X X^T)} \Sigma - \Sigma \underbrace{\mathbb{E}(X X^T)} + \Sigma^2 \\ &= \mathbb{E} X X^T X X^T - \Sigma^2 \leq 20r \|\Sigma\| \cdot \underbrace{\mathbb{E} X X^T}_{\Sigma} = 20r \|\Sigma\| \Sigma \end{aligned}$$

$$\|X\|_2^2 \leq 10 \mathbb{E} \|X\|_2^2 \leq 20r \|\Sigma\|$$

$\uparrow$  assumption       $\uparrow$  lemma

$$0 \leq n \mathbb{E}(X X^T - \Sigma)^2 \leq 20r n \|\Sigma\| \Sigma$$

$$\Rightarrow \sigma^2 = \|\dots\| \leq 20r n \|\Sigma\|^2$$

since  
 $(0 \leq A \leq B \Rightarrow \|A\| \leq \|B\|)$   
HW

$$(*) : \|X X^T - \Sigma\| \leq \|X X^T\| + \|\Sigma\| \leq 21r \|\Sigma\|$$

$\Delta$  inequality  $\uparrow$  HW  
 $\|X\|_2^2 \leq 20r \|\Sigma\|$ , as above

Substitute  $\sigma^2, K$  into  $(*) \Rightarrow$

$$\|\Sigma_n - \Sigma\| \lesssim \frac{1}{n} \left( \sqrt{rn \|\Sigma\|^2 \log d} + r \|\Sigma\| \log d \right) = \underbrace{\left( \sqrt{\frac{r \log d}{n}} + \frac{r \log d}{n} \right)}_{\uparrow} \|\Sigma\|$$

REMARKS 1. Logarithmic oversampling is needed in general

(HW)

2. No distribution assumptions on  $X$ !

$\uparrow$   
0.01 if  $n > Cr \log d$ .