# Lecture 21

**Do Well!**

20 min

PCA

Statistical learning ← CV
Regularization

Predictors
Training
Architectures
Unsupervised
Stat Learning ⟵
Learning & people

# Lecture Notes VII – Principal Component Analysis

Marina Meilă

mmp@uwaterloo.ca

With Thanks to Pascal Poupart & Gautam Kamath
Cheriton School of Computer Science
University of Waterloo

March 19, 2026

**3.** Representing $X$, $x^i$

$X = U \sqrt{\tilde{\Sigma}} V^T$   $n \times D$   ← SVD($X$)   $u^i_j = u_{ij} \tilde{\sigma}_j$

$\Downarrow$

$x^i = v_1 \cdot u^i_1 + v_2 \cdot u^i_2 + \cdots + v_D u^i_D$  → scalar coefficients

in basis $V$

$x^i \longleftrightarrow \left( u^i_1 \cdots u^i_D \right) = \tilde{x}^i$

$\tilde{\sigma}_1 \geq \tilde{\sigma}_2 \geq \cdots \left( \tilde{\sigma}_d \right) \geq \cdots \tilde{\sigma}_D$

$V_1 \quad V_2 \quad\quad V_d \quad\quad V_d$

**Principal e-values vectors**

$U = [u_{ij}]_{\substack{i=1:n \\ j=1:D}}$

$X = \begin{bmatrix} (x^1)^T \\ \cdots \\ (x^n)^T \end{bmatrix}$  $n \times D$

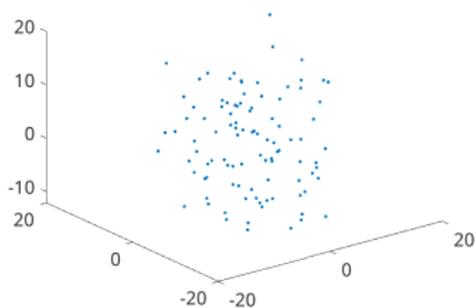**1** $\mathrm{Var}(x) = \frac{1}{n} X^T X =$

$= \frac{1}{n} V \tilde{\Sigma} V^T$  $D \times D$

**2** $Q = X X^T = U \hat{\Sigma} U^T$

$n \times n$

# Example – Gaussian data

$D = 10$



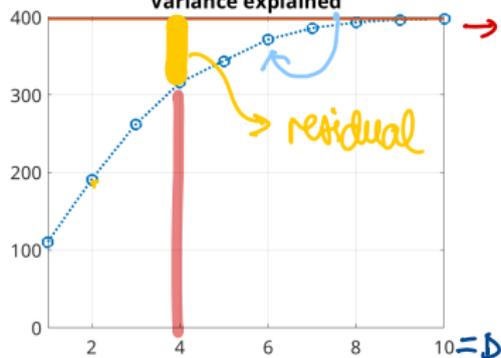**data in dimensions 1:3**



**singular values**

- PCA
- true

$= D$



**data in principal space**

$v_2 =$

$u_2^i$

$u_1^i$

$v_1$



**Variance explained**

$\rightarrow \sum_1^D \tilde{\sigma}_j^2$

residual

$= D$

$\sum_{j=1}^{4} \tilde{\sigma}_j^2$

PCA is "optimal" basis

$$x^i = v_1 \cdot u_1^i + v_2 \cdot u_2^i + \cdots + v_D u_D^i$$

$$\underbrace{\phantom{v_1 \cdot u_1^i + v_2 \cdot u_2^i}}_{d} \quad \underbrace{\phantom{+ v_D u_D^i}}_{\text{residual}}$$

$$\tilde{x}^i = \sum_{j=1}^{d} u_j^i v_j$$

projected on B

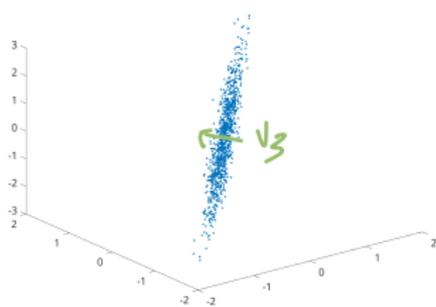$$\text{MSE} \equiv \mathcal{L}_{LS}(B) \equiv \frac{1}{n} \sum_{i=1}^{n} \left( \|x^i\|^2 - \|\tilde{x}^i\|^2 \right)$$

some subspace

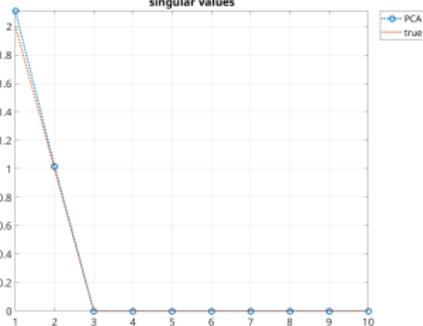$$\mathcal{L}_{LS}(V) = \boxed{\frac{1}{n} \sum_{j=d+1}^{D} \sigma_j^2} = \min_{B} \mathcal{L}_{LS}(B)$$

↑

PCA is optimal
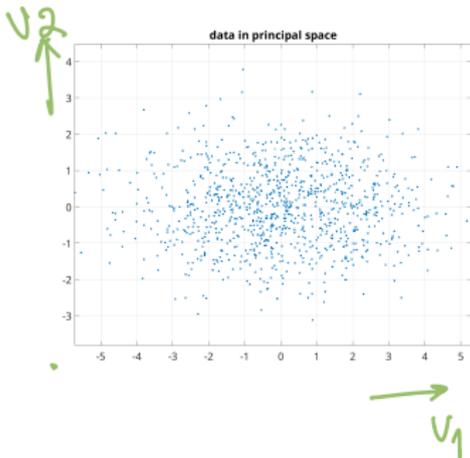
# Example – Gaussian data 2D
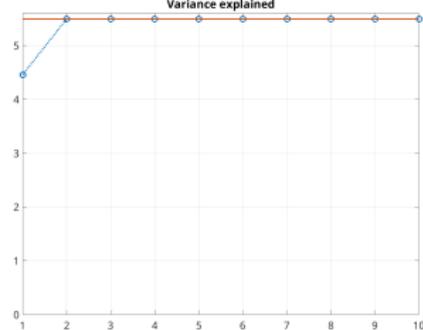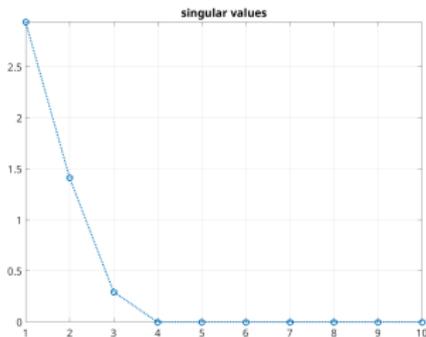
# Example – Brick

d=3 sufficient



data in dimensions 1:3

singular values

data in principal space

Variance explained

pleasure in math

age

# Example – clusters

$d = 3$

data in dimensions 1:3

K = 2 clusters

project on

PCA (K-1)

singular values

data in principal space

$v_2$

$v_1$

Variance explained
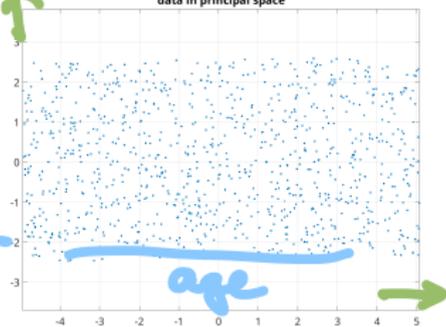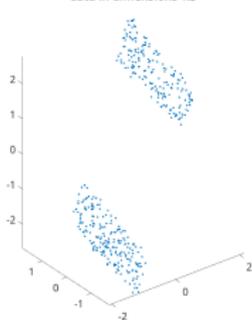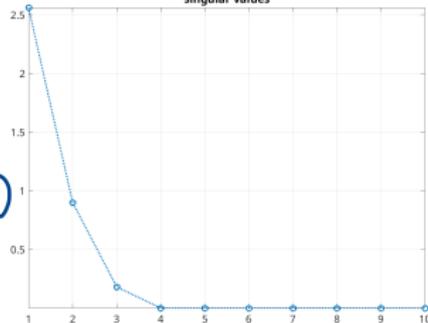
$v_1$

# PCA Summary

*PCAnalysis*

- ▶ Reduces data dimensiont from $D$ to $d$
- ▶ Linear operation (projection)
- ▶ "Optimal" linear method to reduce dimension
- ▶ Can discover if data is low-dimensional
- ▶ For clustering – recommended pre-processing: PCA in $K-1$ dimensions
- ▶ Limitation: fails to discover non-linear low dimensional structure

**• Understand, explore, interpret data**

manifold learning/ non-lin dim red/ embedding

time of day

# Model Selection

- select K for clustering
- select K for K-NN
- compare n.n. architectures
- —''— different predictors $\{DT, nn, K-NN, lin\ regr, \ldots\}$

"Alg". Given $\mathcal{D}$

1. Train $f_1, f_2, \ldots f_M$ predictors on $\mathcal{D}$
2. Select $f^* =$ "best" $\{f_{1:M}\}$

    Model sel

$\Big\langle$ BIC only for Max Likelihood, cheap !! → score

CV any models, any loss → more computation ↑ → score

---

# CV    cross-validation

**Assume** $\mathcal{D}'$ available, $|\mathcal{D}'| = n'$ from same distribution

1. $\ldots$ $f_{1:M}$ trained on $\underline{\mathcal{D}}$ ← training data
2. calculate $Loss(f_m; \mathcal{D}') = $ validation loss
   ↑ validation data
3. $f^* = \underset{m=1:M}{argmin}\ Loss(f_m; \mathcal{D}')$

$n' = ?$ → 1 value

$n \longrightarrow p$ parameters

$n' \leq 1000$

Plenty of data

$n \ggg 1000$

# K-fold CV

1. Partition $\mathcal{D} = \mathcal{D}_1 \cup \ldots \cup \mathcal{D}_K$   disjoint, $|\mathcal{D}_1| \approx \frac{n}{K}$
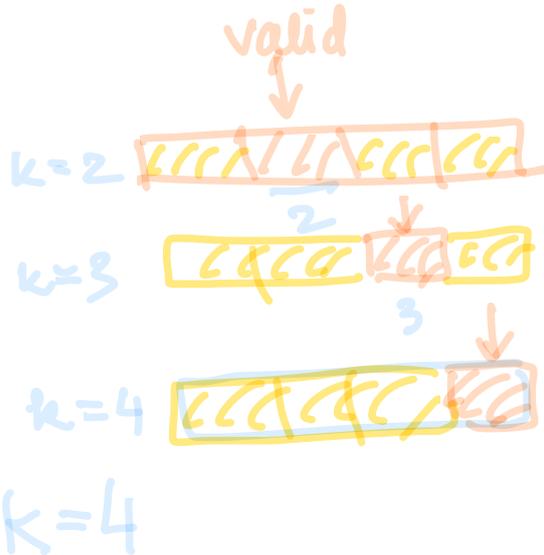
2. for $k = 1:K$

$$\underbrace{\mathcal{D}_{-k} = \mathcal{D} \setminus \overbrace{\mathcal{D}_k}}_{\text{train}}$$

   $\Big[$ train $f_{1:M}^{(k)}$ on $\mathcal{D}_{-k}$

   $\Big[$ compute $L(f_M^{(k)}; \mathcal{D}_k) = L_M^{(k)}$

train   validate

$$\bar{L}_M = \frac{1}{K} \sum_k L_M^{(k)}$$
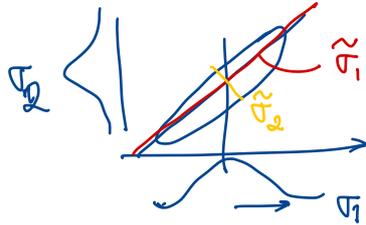
$$m^* = \arg\min_m \bar{L}_M$$

$D = 2$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ & \sigma_2^2 \end{bmatrix}$$

$Var(x_1)$

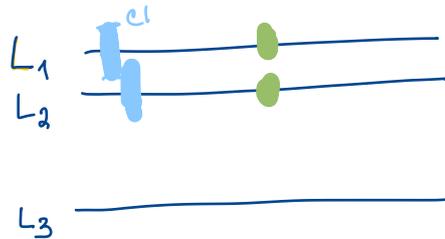$$\sigma_1^2 + \sigma_2^2 = \tilde{\sigma}_1^2 + \tilde{\sigma}_2^2$$



$\sigma_2$    $\tilde{\sigma}_1$    $\tilde{\sigma}_2$    $\sigma_1$

(Questions after class)

$L(f(x))$ r.v. $\longrightarrow$ $\begin{array}{l} Var \\ mean \end{array} = \sigma^2$

$$L = E_{p^*}\left[ L(f(x), y) \right]$$

$$\hat{L} = \frac{1}{n'} \sum L(f(x), y)$$

$$Var\hat{L} = \frac{\sigma^2}{n'}, \quad std = \frac{\sigma}{\sqrt{n'}}$$

$n'$ small

$n'$ large

$L_1$

$L_2$

$c'$

$L_3$