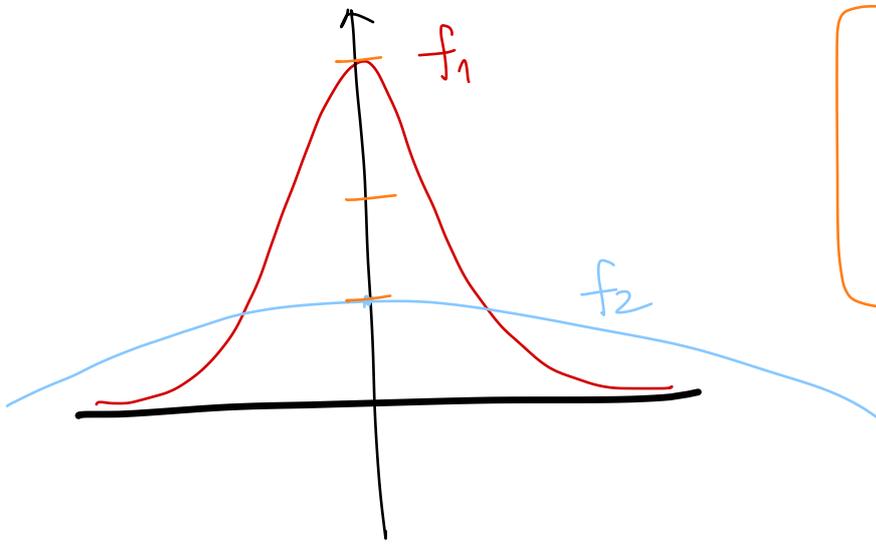# Lecture 21

Do Well!

20 min

PCA
CV

$f_1$

$f_2$

$\sigma_1 = 1$    $f_1(0) = \frac{1}{\sqrt{2\pi}}$

$\sigma_2 = 3$    $f_3(0) = \frac{1}{3}\frac{1}{\sqrt{2\pi}}$

on Q3 Pb 4

# Predictors

Bias - Var

# Training

# Neural Architectures

# Unsupervised

Statistical
CV
Regularisation

ML & People
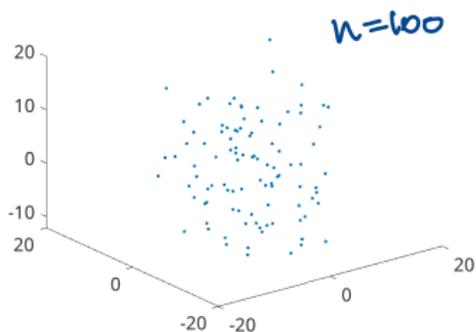
# Lecture Notes VII – Principal Component Analysis

Marina Meilă

`mmp@uwaterloo.ca`

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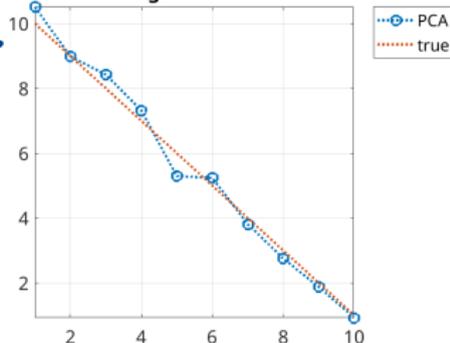Cheriton School of Computer Science
University of Waterloo

March 19, 2026

Eigendecompositions of Variance

**1.** $Var(X) \equiv \Sigma = \frac{1}{n} X^T X = \frac{1}{n} V \tilde{\tilde{S}} V^T$

**3.** $X = U \tilde{S} V^T$ ← SVD

$u_j^i = u_{ij} \tilde{\sigma}_j$

**2.** $G = X X^T = U \tilde{\tilde{S}} U^T$

$\Rightarrow \tilde{X}^i = u_1^i V_1 + \cdots u_d^i V_d$

↑ new basis vectors

optimal approximation

$\tilde{X} \in \mathbb{R}^d$

$X^i = \begin{bmatrix} x_1^i \\ \vdots \\ x_D^i \end{bmatrix} \rightarrow \tilde{X}^i = \begin{bmatrix} u_1^i \\ \vdots \\ u_d^i \end{bmatrix}$

$\tilde{\sigma}_1 \geq \tilde{\sigma}_2 \geq \cdots \tilde{\sigma}_d \geq \cdots$

Principal $d$ e-values

e-vectrs

# Example – Gaussian data

**data in dimensions 1:3** — $n=600$

**singular values** — PCA, true

$\sigma_j = 11-j$
$D = 10$

**data in principal space** — $v_2$, $v_1$, $u_2^i$, $u_1^i$

**Variance explained** — $\sum_{j=d+1}^{D}\sigma^2$

choosing $d$

Var explained

$d$

$D$

# Example – Gaussian data

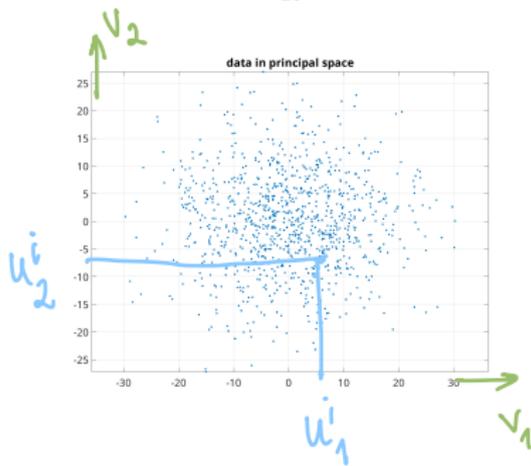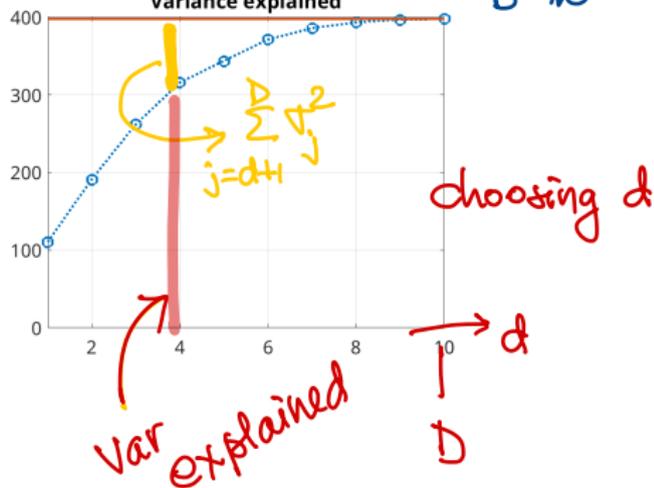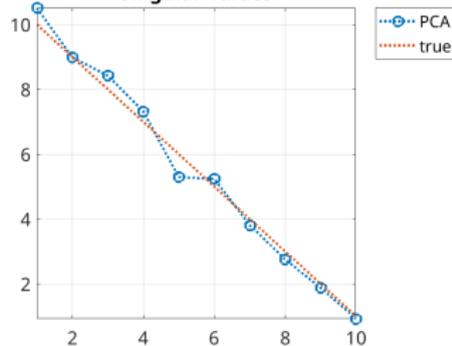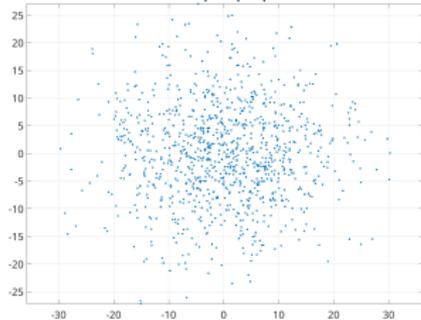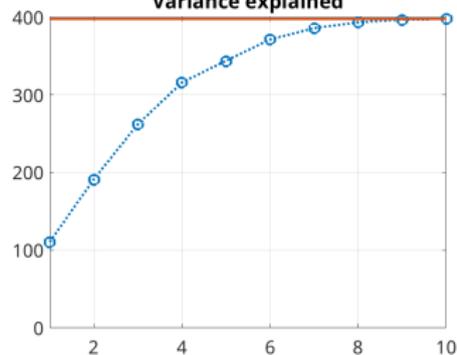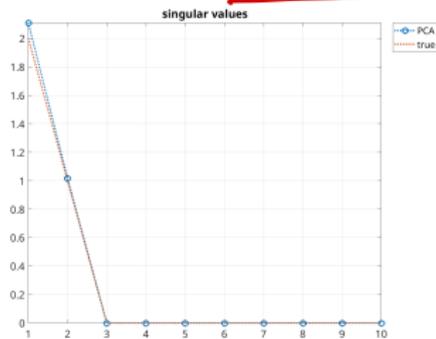# Example – Gaussian data 2D
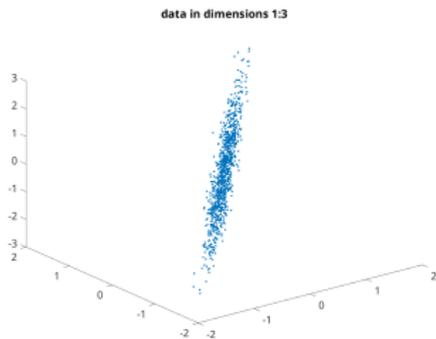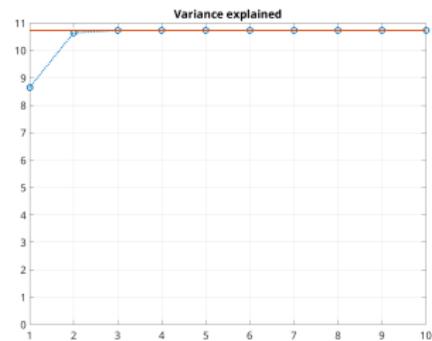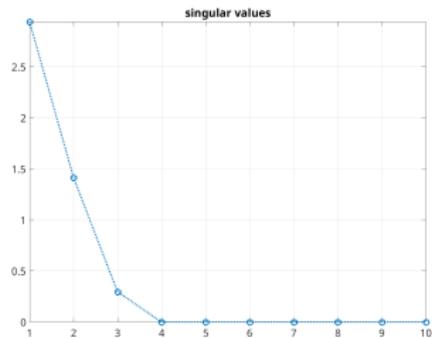
PC Analysis



data in dimensions 1:3

singular values

data in principal space

Variance explained

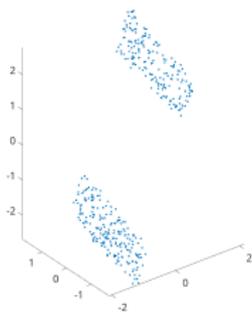# Example – Brick

# Example – clusters



**data in dimensions 1:3**

$k=2$

**singular values**

**data in principal space**

$d=1$

**Variance explained**

In general
K clusters
PCA (K-1)

$d = K - 1$

⇓

clustering
easier
after PCA

⇔

in high
dim

# PCA Summary

- ▶ Reduces data dimensiont from $D$ to $d$
- ▶ Linear operation (projection) → *manifold learning, embedding, non-linear dim reduction*
- ▶ "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

- Low Rank approx
- Interpretation

PCA: $d = 1$

$\leftarrow V_1$

$d = 1$
$D = 2$

Torus

$D = 2$
$d = 2$

# Model Selection

- Selecting $k$ for clustering
- —"— $k$ in $k$-nearest n.
- comparing 2 neural networks
- comparing different types of predictors $\{DT, k\text{-nn}, \text{logistic r}, \dots\}$

Given $\mathcal{D}$, Loss $L_{01}, L_{LS}, \dots$

1. Train predictors $f_1, f_2, \dots f_M$
2. Do <u>Model Selection</u> (= compare $f_{1:M}$)

BIC — only for ML estimation, cheap → score
CV — any predictors, expensive → score $L^{valid}$ = validation loss

$$f^* = \underset{m=1:M}{\arg\max} \{ \text{score}(f_m) \\ \min \}$$

# Cross Validation

Assume $\mathcal{D}'$ available $|\mathcal{D}'| = n'$ (from same distribution)

1. for $m=1:M$ train on $\mathcal{D}$ $f_m$
2. —"— calculate $L(f_m; \mathcal{D}') = L^v_m$
3. choose $f^* = \arg\min L^v_{1:M}$

$n' = $ how large? $\sim 1000$

$n$

If $n \gg 1000$    K-fold CV

1. Partition $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \ldots \cup \mathcal{D}_K$    (disjoint)

$$|\mathcal{D}_k| \approx \frac{n}{K}$$

$\mathcal{D}_{-k} = \mathcal{D} \setminus \mathcal{D}_k$

train    validation

2. for $k = 1:K$

  train $f_{1:m}$ on $\mathcal{D}_{-k} \Rightarrow f_m^{(k)}$

  calculate loss    $L_m^{(k)} = Loss\left(f_m^{(k)}; \mathcal{D}_k\right)$

· $\overline{L}^v(f_m) = \frac{1}{K} \sum_{k=1}^{K} L_m^{(k)}$

· $m^* = \arg\min_{m=1:M} \overline{L}_m^v$

· retrain $f_{m^*}$ on $\mathcal{D}$

$K \nearrow$    $|\mathcal{D}_{-k}| = \frac{k-1}{K} n$    more data for training

more computation

very small $n \leftsquigarrow 100 \rightarrow K=n$    Leave-One-Out CV