

# Foundations of Machine Learning

## AI2000 and AI5000

FoML-29  
PCA

Dr. Konda Reddy Mopuri  
Department of AI, IIT Hyderabad  
July-Nov 2025

# So far in FoML

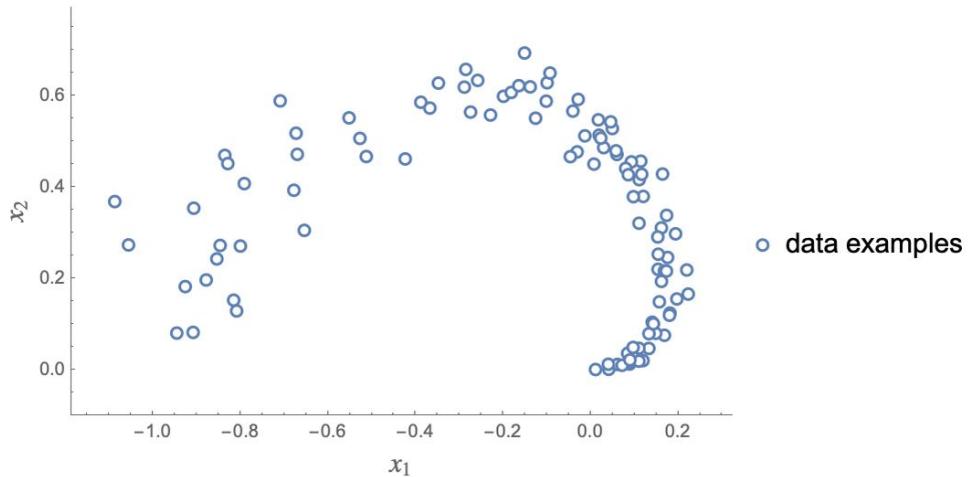
- Intro to ML and Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Supervised learning
  - a. Linear Regression with basis functions
  - b. Bias-Variance Decomposition
  - c. Decision Theory - three broad classification strategies
  - d. Neural Networks
- Unsupervised learning
  - a. K-Means, Hierarchical, and GMM for clustering



# For today

- PCA - Principal Component Analysis (Pearson, 1901) & (Hotelling, 1933)

# Manifold coordinates as Latent variables



$$\{x_1, x_2\} = \{t \cos(3 t), t \sin(3 t)\}$$



# Example - facial image data

- Possible degrees of freedom
  - Skull size
  - Skin color
  - Eye color
  - Facial attributes
  - Horizontal orientation
  - Vertical orientation
  - Mood (e.g., happy)
  - etc.



# Example - facial image data

- Possible degrees of freedom
  - Skull size
  - Skin color
  - Eye color
  - Facial attributes
  - Horizontal orientation
  - Vertical orientation
  - Mood (e.g., happy)
  - etc.

Latent subspace will be a nonlinear transformation of image data



# PCA

- Linear latent subspaces

# What is PCA?

- Tool to summarize a large set of variables (dimensions) with a smaller set



ભારતીય નોંકેટિક વિજ્ઞાન સંસ્કૃત પ્રૌદ્યોગિકી  
ભારતીય પ્રૌદ્યોગિકી સંસ્થાન હૈદરાબાદ  
Indian Institute of Technology Hyderabad

# What is PCA?

- Tool to summarize a large set of variables (dimensions) with a smaller set
  - Of ‘representative’ variables that explain the ‘variability’ in the original set



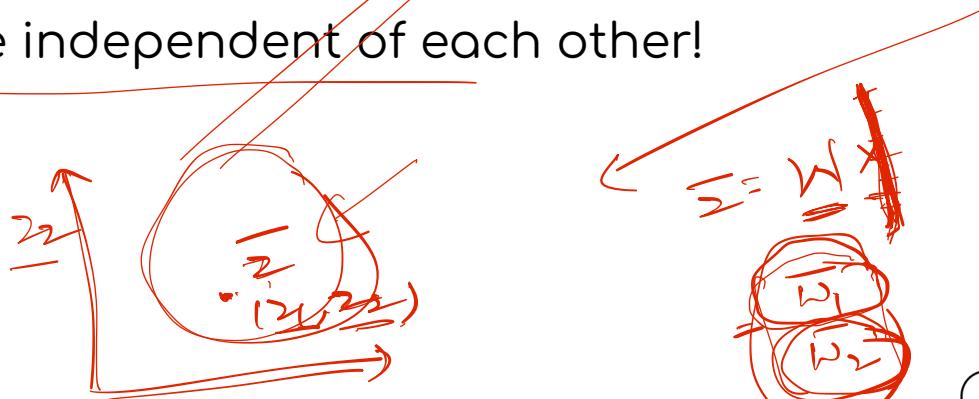
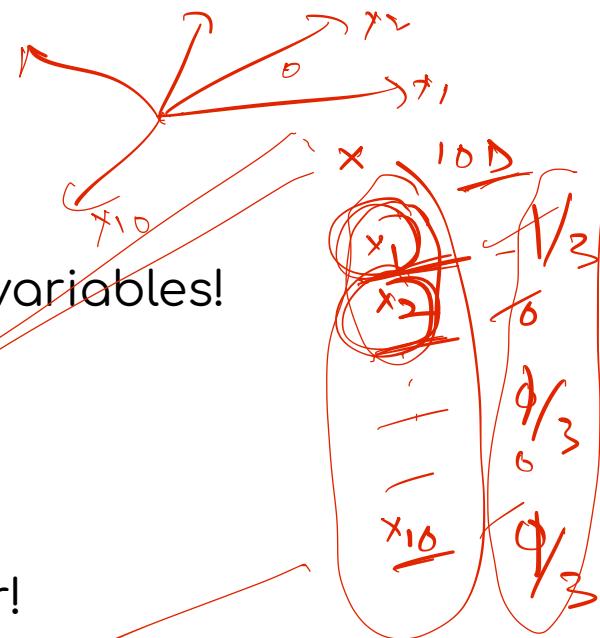
# What is PCA?

- Tool to summarize a large set of variables (dimensions) with a smaller set
  - Of ‘representative’ variables that explain the ‘variability’ in the original set



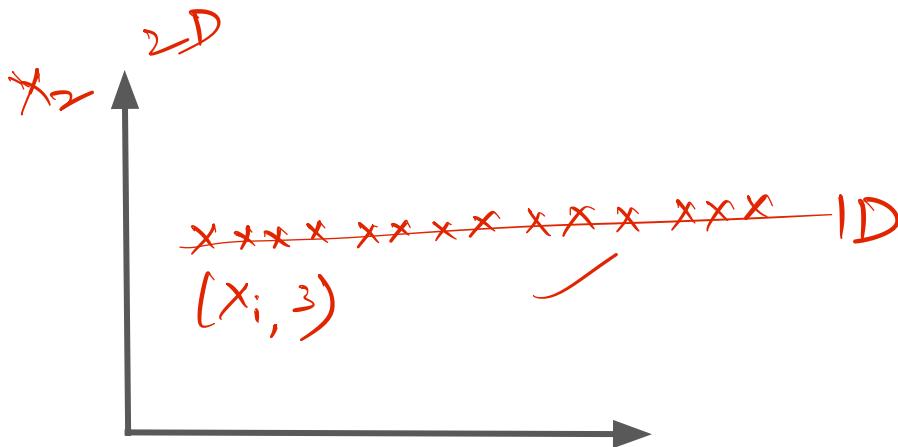
# What is PCA?

- Smaller set need not be a subset of original variables!
  - Rather, combinations of original variables
- They may not mean the same as originals
  - lost interpretability!
- New variables are independent of each other!



# What is PCA?

- Gives the directions along which the data are highly 'variable'
  - Projects linearly such that the variance in the projected space is maximal



1-dim  
 $(x_1)$   $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$



# PCA

- Data  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$



$$\mathbf{x}_i \in \mathbb{R}^D$$



- Aim: project data onto  $M$  dimensional space ( $M < D$ ) maximizing the variance of the projected data



# PCA

- Mean

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

- Covariance

$$\underline{\mathbf{S}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Symmetric  
↓

$\lambda_i \geq 0$  ↗  
↑  
lin ind Eig Vecs

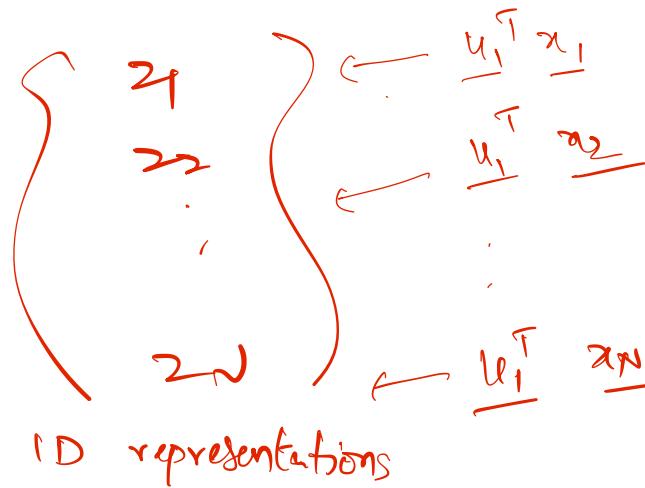
PSD



# 1D representation using PCA

- Project data onto a direction where most of the variance is preserved
- Projecting onto  $u_1$  gives a scalar  $\rightarrow$  1D representation

$$\mathbf{z}_i = \mathbf{u}_1^T \mathbf{x}_i$$

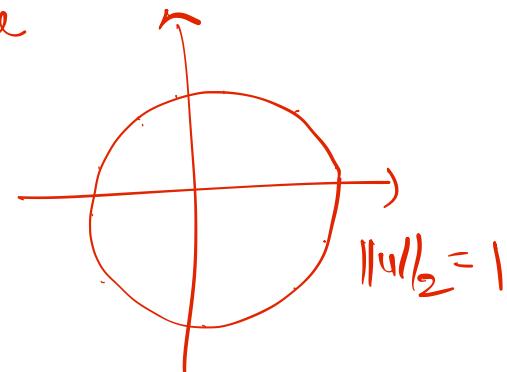


# 1D representation using PCA

- Direction of  $u_1$  is important → consider unit vector in that direction

$$\|\mathbf{u}_1\|_2 = 1$$


↓  
But not its magnitude  
w/o this constraint,  
we can't optimize



# 1D representation using PCA

- Consider the variance in the new subspace

$$\begin{aligned}\text{Var}[\mathbf{z}] &= \frac{1}{N} \sum_{i=1}^N (\underline{z}_i - \bar{z})^2 = \frac{1}{N} \sum_{i=1}^N (\underline{u}_i^\top \underline{x}_i - \underline{u}_i^\top \bar{\underline{x}})^2 \\ &= \frac{1}{N} \sum_{i=1}^N [\underline{u}_i^\top (\underline{x}_i - \bar{\underline{x}})]^2 = \frac{1}{N} \sum_{i=1}^N \underline{u}_i^\top (\underline{x}_i - \bar{\underline{x}}) (\underline{x}_i - \bar{\underline{x}})^\top \underline{u}_i \\ &= \frac{1}{N} \underline{u}_i^\top \sum_{i=1}^N (\underline{x}_i - \bar{\underline{x}}) (\underline{x}_i - \bar{\underline{x}})^\top \underline{u}_i = \underline{u}_i^\top S \underline{u}_i\end{aligned}$$

$S$



# 1D representation using PCA

1st principal component  
PCI  
↓  
 $u_1$

- Let's find the direction ( $u_1$ ) that maximizes the variance in  $z_i$

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \text{ such that } \mathbf{u}_1^T \mathbf{u}_1 = 1$$

using Lagrange multiplier

$$\begin{aligned} \text{argmax}_{\mathbf{u}_1, \lambda} \quad & \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda(1 - \mathbf{u}_1^T \mathbf{u}_1) = 0 \\ \frac{\partial}{\partial \mathbf{u}_1} \quad & \downarrow = 0 \\ & \mathbf{R} \mathbf{S} \mathbf{u}_1 - 2\lambda \mathbf{u}_1 = 0 \end{aligned}$$

$$\begin{aligned} \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) &= 0 \\ f(\mathbf{x}) + \lambda g(\mathbf{x}) &= 0 \end{aligned}$$

$$\mathbf{S} \mathbf{u}_1 = \lambda \mathbf{u}_1$$

$\mathbf{u}_1$  is eig. vector  
 $\lambda$  is corresponding eig. value of  $\mathbf{S}$



# PCA via maximum variance

$$\left( \underline{x_i} - \underline{z_i} \underline{u_1} \right) \frac{\downarrow}{d\underline{x_i}}$$

- Repeat the procedure for the next M-1 orthogonal vectors
  - Maximize the variance by projecting onto a direction orthogonal to the found ones
  - These are the next M-1 eigenvectors of the covariance matrix (S)

$$\underline{u_1^T S u_1} = \underline{u_1^T} \lambda \underline{u_1} = \lambda \underline{u_1^T} \underline{u_1} = \lambda$$

The variance captured by a PC is its eigen value ( $\lambda_i$ )  
(of the projected data)



# PCA - Eigen decomposition

- For the symmetric PSD matrix  $\mathbf{S} = \mathbf{U}\Lambda\mathbf{U}^T$
- Eigenvectors are orthonormal (contained in U)
- Eigenvalues are non-negative (contained in  $\Lambda$ )



# PCA - Eigen decomposition

- Variance is  $\text{Tr}(S)$   
of the projected data

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$$

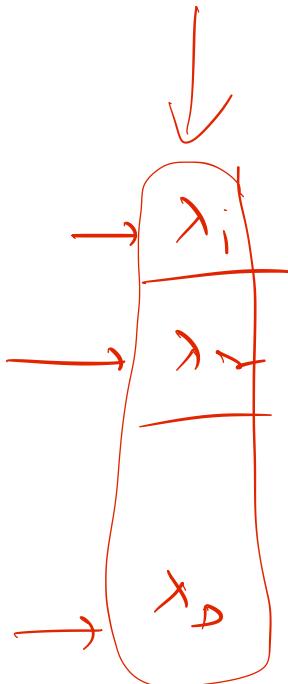
$$\begin{aligned}\text{Tr}[\text{cov}[z]] &= \text{variance of} \\ &\text{data along each of the new dims} \\ &= \lambda_1 + \lambda_2 + \dots + \lambda_D\end{aligned}$$

$$\underline{\mathbf{S}} = \underline{\mathbf{U}} \underline{\boldsymbol{\Lambda}} \underline{\mathbf{U}}^T$$

$$\rightarrow \left( u_1^T \underline{S} u_1 \right)$$

$$\rightarrow \left( u_2^T \underline{S} u_2 \right)$$

$$\left( u_D^T \underline{S} u_D \right)$$



# PCA - some notes

# PCA - Scaling the features

- Sensitive to the scales of the features/variables
  - Features with greater range dominate the process of finding the PCs
- Perform standardization to prevent this

mean subtraction and variance normalization



# PCA - Proportion of the Variance Explained

- How much of the information is lost by projecting onto PCs?

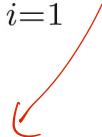


# PCA - Proportion of the Variance Explained

- Total variance

$p = \text{no. of dimensions in the original space}$

$$\sum_{j=1}^p \text{Var}(X_j) = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n x_{ij}^2$$



assume mean subtracted



# PCA - Proportion of the Variance Explained

- Total variance
- Variance explained by the ' $m^{\text{th}}$  PC' 

$$\sum_{j=1}^p \text{Var}(X_j) = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n x_{ij}^2$$

$$\frac{1}{n} \sum_{i=1}^n z_{im}^2$$

 linear transformation  
doesn't disturb the mean  
centering



# PCA - Proportion of the Variance Explained

- Total variance
- Variance explained by the 'm<sup>th</sup>' PC

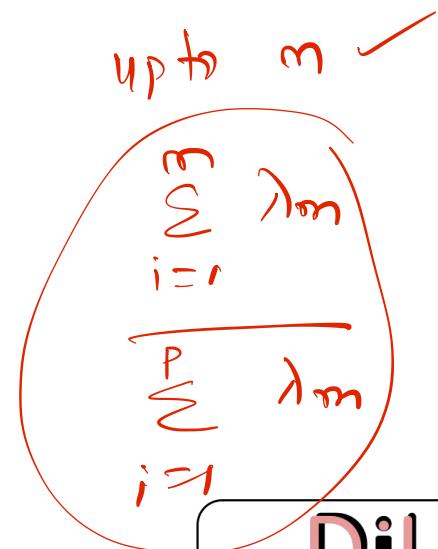
$$\sum_{j=1}^p \text{Var}(X_j) = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n x_{ij}^2$$

$$\frac{1}{n} \sum_{i=1}^n z_{im}^2$$

PVE of the 'm<sup>th</sup>' PC =



$$\frac{\lambda_m}{\sum_{i=1}^D \lambda_i}$$



# PCA - How many PCs to consider?

- For an  $n \times p$  data matrix
  - $\min(n-1, p)$  PCs are possible
  - Why?

mean centering makes  
the rows dependent,

hence  $n-1$

$$X = \underbrace{\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}}_{n \times D}$$

$\frac{500}{n}$  - # samples  
 $D$  - dims

$n > D$  generally

$n < D$  rarely



# PCA - How many PCs to consider?

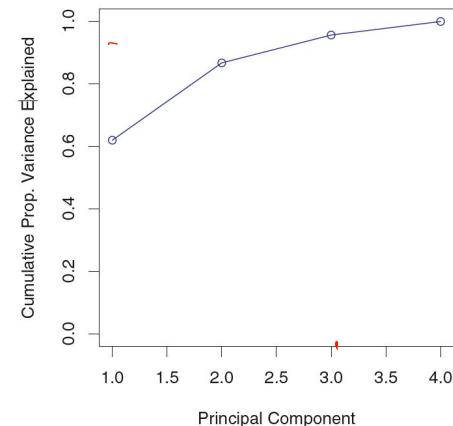
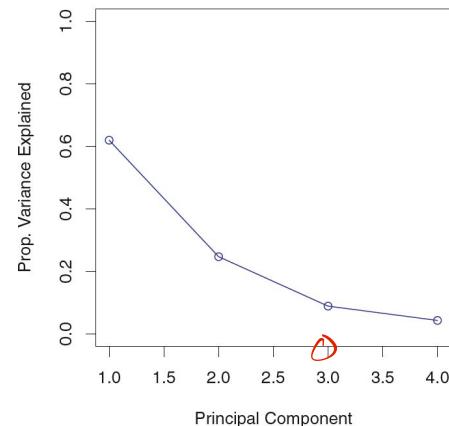
- For an  $n \times p$  data matrix
  - $\min(n-1, p)$  PCs are possible
- Not all of them may be interesting

(leads to acceptable  
loss of information)



# PCA - How many PCs to consider?

- Generally, we want the ‘smallest’ number of them → good understanding of the data
- → scree plot & elbow



# PCA

(not feature selection,  
rather feature extraction)

- Doesn't discard the redundant variables
  - Finds new variables (linear combinations of the ' $p$ ' variables) that summarize the data well
  - The 'best' variables (among the all possible linear combinations)
  - Resulting new features are uncorrelated (covariance matrix will be diagonal)

' $p$ ' and ' $D$ ' both indicate the original dimension of the data



# Applications of PCA

- Dimensionality reduction → tackles curse of dimensionality
- Less compute requirement
- Less prone to overfitting
- Useful preprocessing



भारतीय नॉर्केटिक विज्ञान संस्था  
भारतीय प्रौद्योगिकी संस्थान हैदराबाद  
Indian Institute of Technology Hyderabad



# Next

- PCA continued

on the constraint set (green level set)

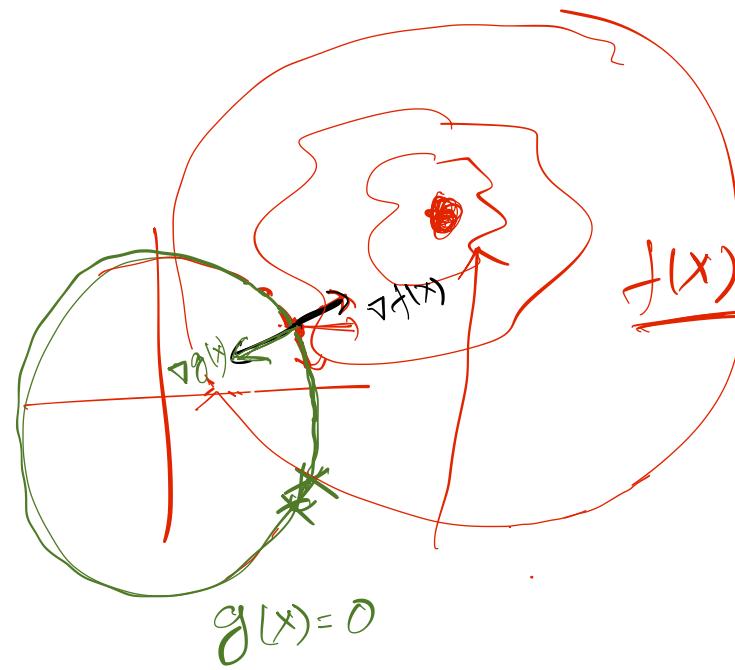
$\nabla g(x)$  is  $\perp$  to the level set surface

$$g(x+\epsilon) = g(x) + \epsilon^T \nabla g(x)$$

$g(x+\epsilon) = g(x)$  on the level set

$$\Rightarrow \epsilon^T \nabla g(x) = 0 \quad \nabla g(x) \perp \epsilon$$

$\epsilon$  is the movement along the levelset



$\nabla f(x)$  will also be  $\perp$  to the level set at  $x^*$ . otherwise we can increase  $f(x)$  by moving along the  $\nabla f(x)$



भारतीय नॉर्केटिक विज्ञान संस्था पूर्वभाग

भारतीय प्रौद्योगिकी संस्थान हैदराबाद

Indian Institute of Technology Hyderabad

DIL

Data-driven Intelligence  
& Learning Lab