

# PCA

Machine Learning

Hamid R Rabiee – Zahra Dehghanian  
Spring 2025



Sharif University  
of Technology

# Dimensionality Reduction: Feature Selection vs. Feature Extraction

## ? Feature **selection**

- ? Select a subset of a given feature set

## ? Feature **extraction**

- ? A linear or non-linear transform on the original feature space

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{bmatrix} x_{i_1} \\ \vdots \\ x_{i_{d'}} \end{bmatrix}$$

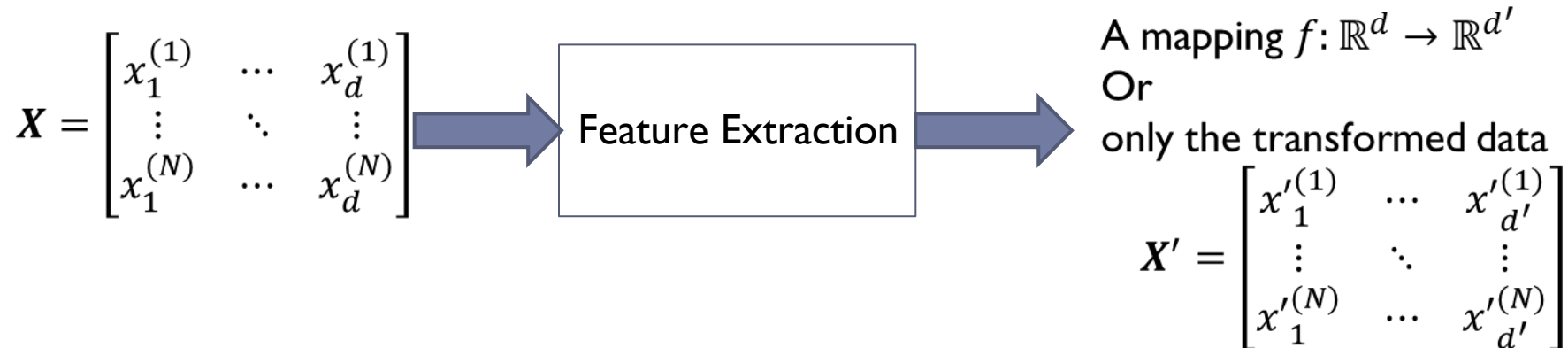
Feature  
Selection  
( $d' < d$ )

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ \vdots \\ y_{d'} \end{bmatrix} = f \left( \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \right)$$

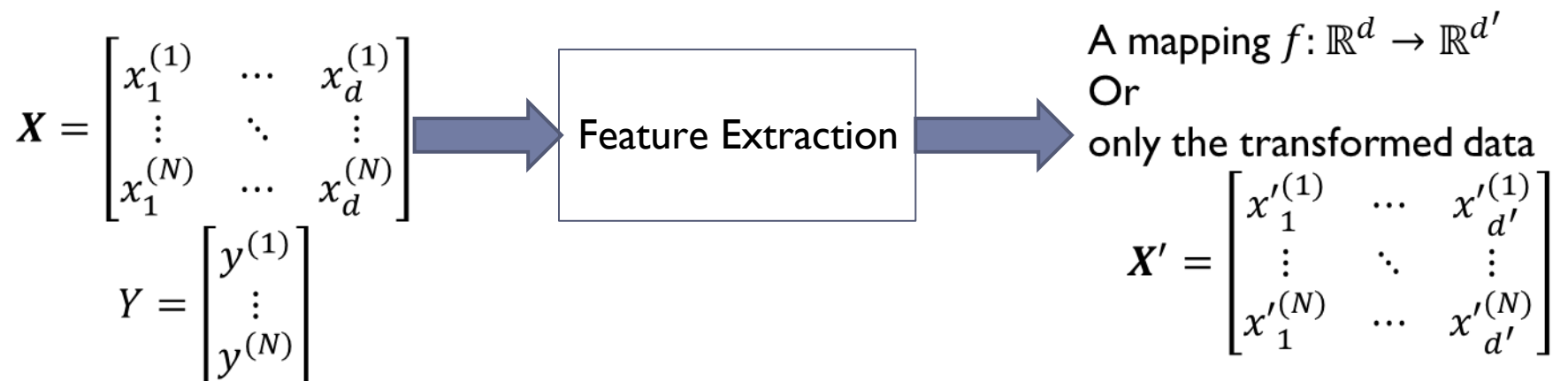
Feature  
Extraction

# Feature Extraction

□ Unsupervised feature extraction:



□ Supervised feature extraction:

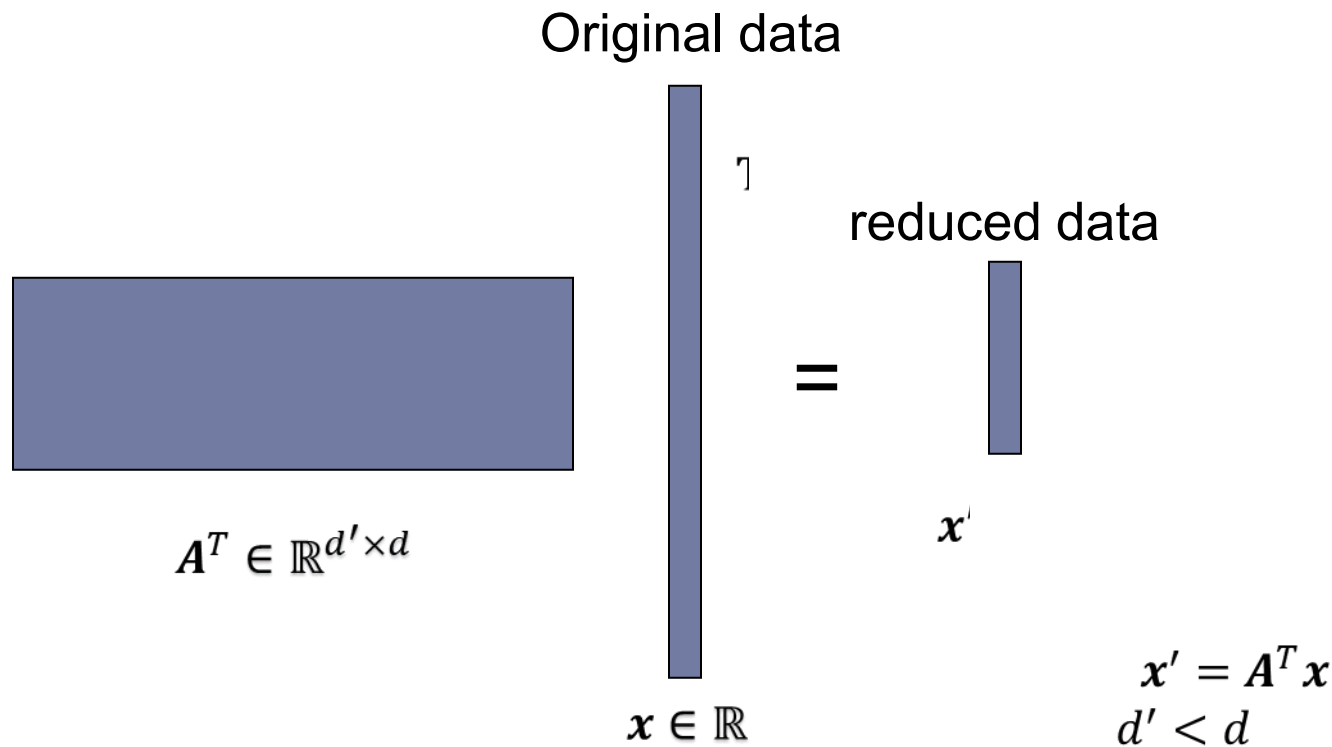


# Unsupervised Feature Reduction

- ❑ Visualization and interpretation: projection of high-dimensional data onto 2D or 3D.
- ❑ Data compression: efficient storage, communication, or and retrieval.
- ❑ Pre-process: to improve accuracy by reducing features
  - ❑ As a preprocessing step to reduce dimensions for supervised learning tasks
  - ❑ Helps avoiding overfitting
- ❑ Noise removal
  - ❑ E.g, “noise” in the images introduced by minor lighting variations, slightly different imaging conditions,

# Linear Transformation

- For linear transformation, we find an explicit mapping  $f(\mathbf{x}) = \mathbf{A}^T \mathbf{x}$  that can transform also new data vectors.



# Linear Transformation

□ Linear transformation are simple mappings

$$\mathbf{x}' = \mathbf{A}^T \mathbf{x} \quad (\mathbf{x}'_j = \mathbf{a}_j^T \mathbf{x}) \quad j = 1, \dots, d$$

$$\mathbf{A} = \begin{bmatrix} \boxed{\begin{matrix} a_{11} \\ \vdots \\ a_{d1} \end{matrix}} & \cdots & \boxed{\begin{matrix} a_{1d} \\ \vdots \\ a_{dd'} \end{matrix}} \end{bmatrix}$$

$\mathbf{a}_1$                        $\mathbf{a}_{d'}$

# Linear Dimensionality Reduction

## □ Unsupervised

- **Principal Component Analysis (PCA)**
- Singular Value Decomposition (SVD)
- Independent Component Analysis (ICA)
- Multi Dimensional Scaling (MDS)
- Canonical Correlation Analysis (CCA)
- ...

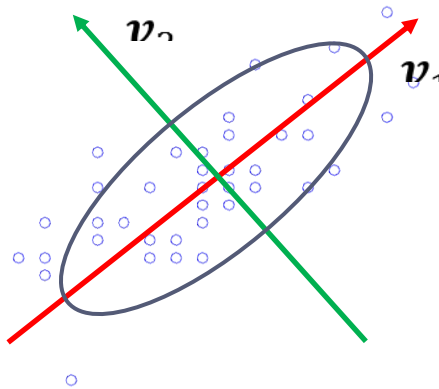
# Principal Component Analysis (PCA)

- Also known as Karhonen-Loeve (KL) transform
- Principal Components (PCs): **orthogonal** vectors that are **ordered** by the fraction of the total information (variation) in the corresponding directions
  - Find the directions at which data approximately lie

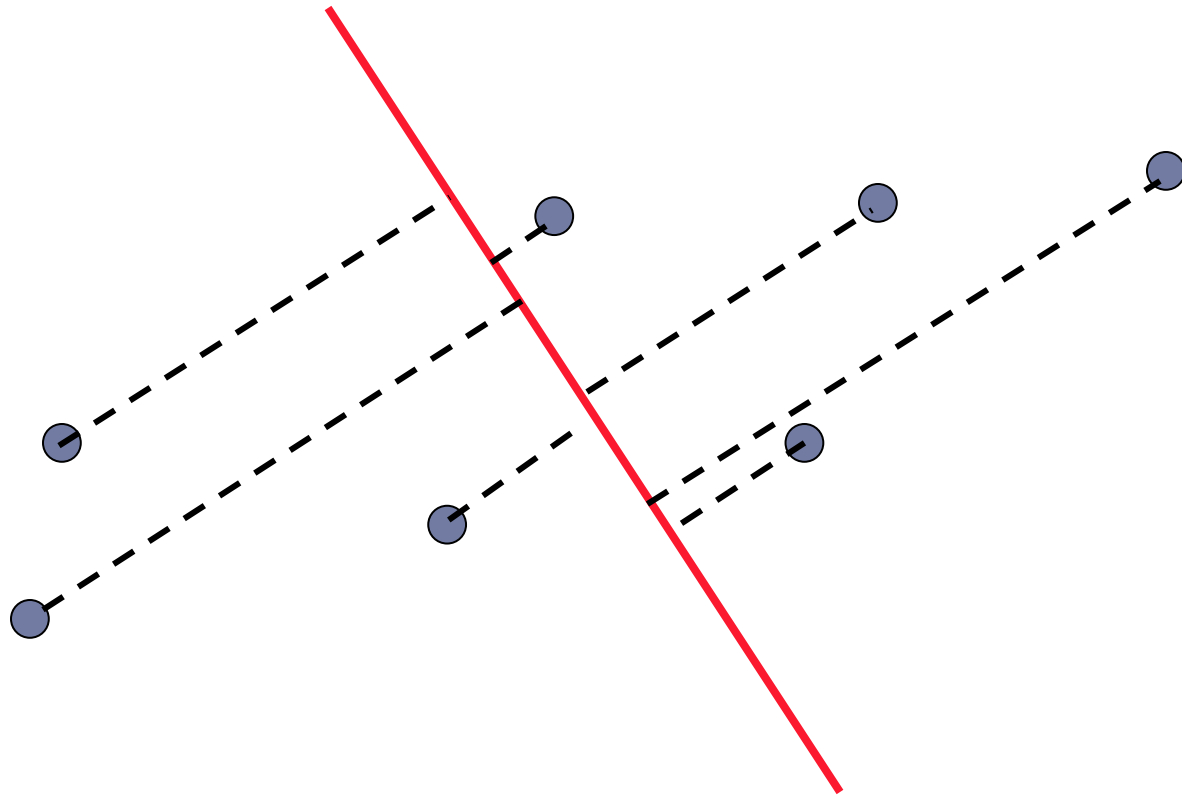


# Principal components

- ▣ If data has a Gaussian distribution  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , the direction of the largest variance can be found by the eigenvector of  $\boldsymbol{\Sigma}$  that corresponds to the largest eigenvalue of  $\boldsymbol{\Sigma}$

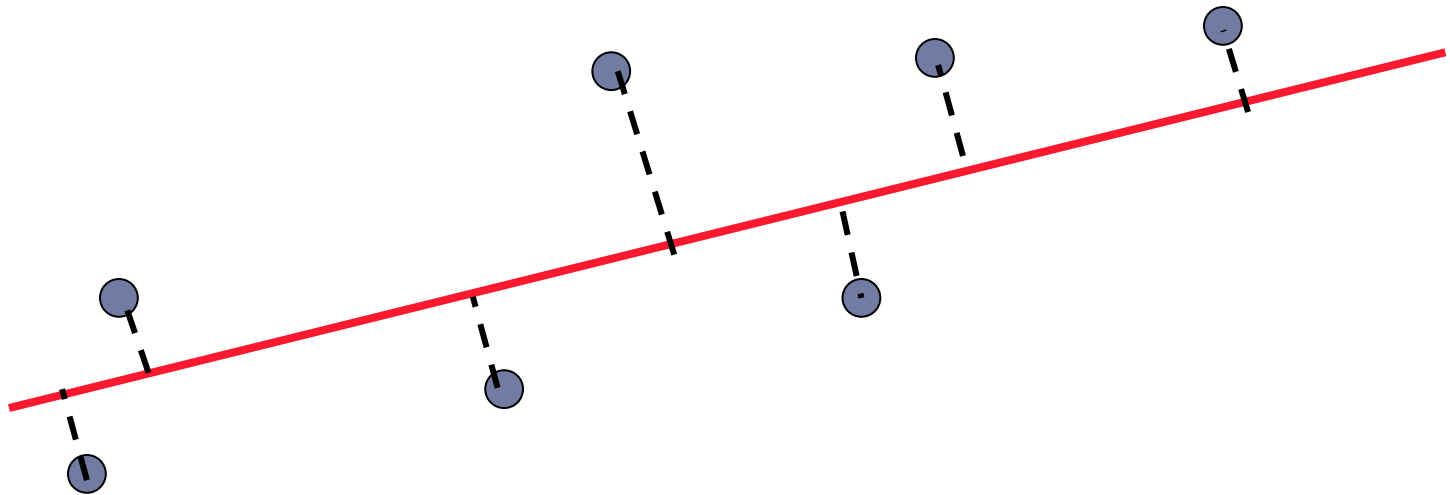


# Example: random direction



# Example: principal component

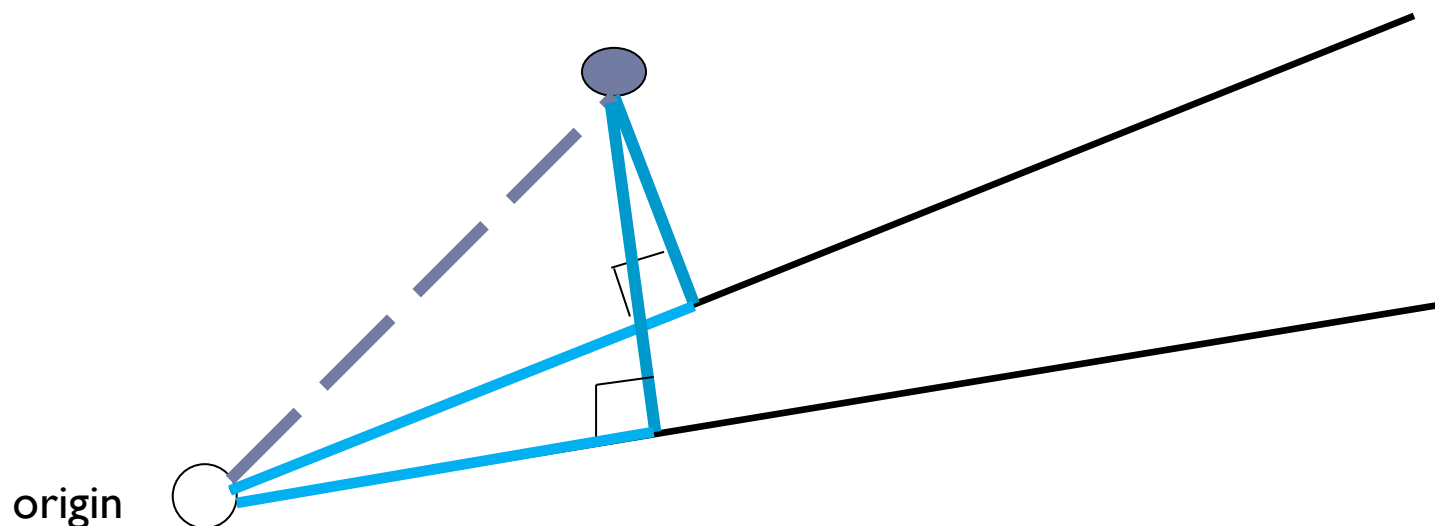
- Find the direction that preserves important aspect of data



# Least Squares Error and Maximum Variance Views Are Equivalent (1-dim Interpretation)

□ When data are mean-removed:

- Minimizing sum of square distances to the line is equivalent to maximizing the sum of squares of the projections on that line (Pythagoras).



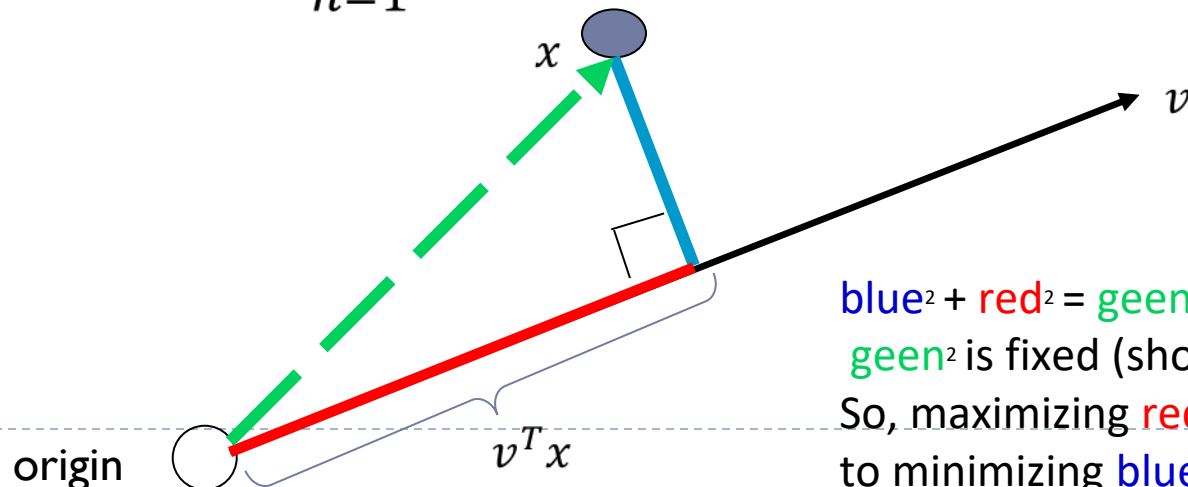
# Two interpretations (for mean centered data)

## Maximum variance subspace

$$\operatorname{argmax}_v \frac{1}{N} \sum_{n=1}^N (v^T x^{(n)})^2 = v^T S v$$

## Minimum reconstruction error

$$\operatorname{argmin}_v \sum_{n=1}^N \|x^{(n)} - (v^T x^{(n)})v\|^2$$



$$\text{blue}^2 + \text{red}^2 = \text{green}^2$$

green<sup>2</sup> is fixed (shows data)

So, maximizing red<sup>2</sup> is equivalent to minimizing blue<sup>2</sup>

# Principal Component Analysis (PCA)

- ❑ Goal: **reducing the dimensionality** of the data while preserving important aspects of the data
- ❑ Two equal views: find directions for which
  - ❑ the **variation** presents in the dataset is as **much** as possible.
  - ❑ the **reconstruction error is minimized**.
- ❑ PCs can be found as the “**best**” **eigenvectors** of the covariance matrix of the data points.

# PCA: Steps

- Input:  $N \times d$  data matrix  $X$  (each row contain a  $d$  dimensional data point)
  - $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$
  - $\tilde{X} \leftarrow$  Mean value of data points is subtracted from rows of  $X$
  - $S = \frac{1}{N} \tilde{X}^T \tilde{X}$  (Covariance matrix)
  - Calculate eigenvalue and eigenvectors of  $S$
  - Pick  $d'$  eigenvectors corresponding to the largest eigenvalues and put them in the columns of  $A = [\mathbf{v}_1, \dots, \mathbf{v}_{d'}]$ 
    - First PC
    - d'-th PC
  - $X' = XA$

# Covariance Matrix



$$\boldsymbol{\mu}_x = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} = \begin{bmatrix} E(x_1) \\ \vdots \\ E(x_d) \end{bmatrix}$$

$$\boldsymbol{\Sigma} = E[(\boldsymbol{x} - \boldsymbol{\mu}_x)(\boldsymbol{x} - \boldsymbol{\mu}_x)^T]$$



# Covariance Matrix



$$\boldsymbol{\mu}_x = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} = \begin{bmatrix} E(x_1) \\ \vdots \\ E(x_d) \end{bmatrix}$$

$$\boldsymbol{\Sigma} = E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T]$$

- ▶ ML estimate of covariance matrix from data points  $\{\mathbf{x}^{(i)}\}_{i=1}^N$ :

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

$$\tilde{\mathbf{X}} = \begin{bmatrix} \tilde{\mathbf{x}}^{(1)} \\ \vdots \\ \tilde{\mathbf{x}}^{(N)} \end{bmatrix} = \begin{bmatrix} \mathbf{x}^{(1)} - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}^{(N)} - \bar{\mathbf{x}} \end{bmatrix}$$

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



# Find 1st principal component

- Find vector  $v_1$  that maximizes sample variance of the projected data:

$$\begin{aligned} & \max_{v_1} \frac{1}{N} \sum_{n=1}^N (v_1^T x^{(n)} - v_1^T \bar{x})^2 \\ &= \frac{1}{N} \sum_{n=1}^N v_1^T (x^{(n)} - \bar{x})(x^{(n)} - \bar{x})^T v_1 \\ &= v_1^T \left( \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \bar{x})(x^{(n)} - \bar{x})^T \right) v_1 = v_1^T S v_1 \\ & \quad \text{s. t. } v_1^T v_1 = 1 \end{aligned}$$

# Find 1st principal component



- Find vector  $v$  that maximizes sample variance of the projected data:

$$\max_v \frac{1}{N} \sum_{n=1}^N (v_1^T x^{(n)} - v_1^T \bar{x})^2 = v_1^T S v_1$$

s. t.  $v_1^T v_1 = 1$

$$L(v_1, \lambda_1) = v_1^T S v_1 + \lambda_1 (1 - v_1^T v_1)$$

$$\begin{aligned} \frac{\partial L}{\partial v_1} = 0 &\Rightarrow 2Sv_1 - 2\lambda_1 v_1 = 0 \\ &\Rightarrow Sv_1 = \lambda_1 v_1 \end{aligned}$$

Eigenvector with maximum eigenvalue maximizes the objective

# PCA Derivation: Relation between Eigenvalues and Variances



$$S\mathbf{v}_j = \lambda_j \mathbf{v}_j$$
$$\Rightarrow \underbrace{\text{var}(\mathbf{v}_j^T \mathbf{x})}_{\downarrow} = \mathbf{v}_j^T S \mathbf{v}_j = \lambda_j \mathbf{v}_j^T \mathbf{v}_j = \lambda_j$$

Variance along j-th eigenvector

- Therefore, eigenvector with maximum eigenvalue maximizes the objective

# Finding second principal component



$$\begin{aligned} \max_{v_2} \quad & v_2^T S v_2 \\ \text{s. t.} \quad & v_2^T v_2 = 1 \\ & v_2^T v_1 = 0 \end{aligned}$$

$$L(v_2, \lambda_2, \alpha) = v_2^T S v_2 + \lambda_2(1 - v_2^T v_2) - \alpha v_2^T v_1$$

# Finding second principal component



$$\begin{aligned} \max_{v_2} \quad & v_2^T S v_2 \\ \text{s. t.} \quad & v_2^T v_2 = 1 \\ & v_2^T v_1 = 0 \end{aligned}$$

$$L(v_2, \lambda_2, \alpha) = v_2^T S v_2 + \lambda_2(1 - v_2^T v_2) - \alpha v_2^T v_1$$

Finding  $\alpha$ :

$$\begin{aligned} \frac{\partial L}{\partial v_2} = 0 &\Rightarrow 2Sv_2 - 2\lambda_2 v_2 - \alpha v_1 = 0 \\ &\Rightarrow 2\mathbf{v}_1^T S v_2 - 2\lambda_2 \mathbf{v}_1^T v_2 - \alpha \mathbf{v}_1^T v_1 = 0 \\ &\Rightarrow 2\lambda_1 \mathbf{v}_1^T v_2 - 2\lambda_2 \times 0 - \alpha = 0 \\ &\Rightarrow \alpha = 0 \end{aligned}$$

# Finding second principal component



$$\begin{aligned} \max_{v_2} \quad & v_2^T S v_2 \\ \text{s. t.} \quad & v_2^T v_2 = 1 \\ & v_2^T v_1 = 0 \end{aligned}$$

$$L(v_2, \lambda_2, \alpha) = v_2^T S v_2 + \lambda_2(1 - v_2^T v_2) - \cancel{\alpha v_2^T v_1}^0$$

Finding  $\lambda_2$ :

$$\begin{aligned} \frac{\partial L}{\partial v_2} = 0 &\Rightarrow 2Sv_2 - 2\lambda_2 v_2 = 0 \\ &\Rightarrow Sv_2 = \lambda_2 v_2 \end{aligned}$$

$v_2$  is the eigenvector corresponding to the second largest eigenvalue

# Find principal components

- For symmetric matrices, there exist eigen-vectors that are orthogonal.
- ▶ Let  $v_1, \dots, v_d$  denote the eigen-vectors of  $S$  such that:

$$\begin{aligned}v_i^T v_j &= 0, & \forall i \neq j \\v_i^T v_i &= 1, & \forall i\end{aligned}$$



# PCA

- Eigenvalues:  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots$ 
  - ▶ The first PC  $v_1$  is the the eigenvector of the sample covariance matrix  $S$  associated with the largest eigenvalue.
  - ▶ The 2nd PC  $v_2$  is the the eigenvector of the sample covariance matrix  $S$  associated with the second largest eigenvalue
  - ▶ And so on ...
- ▶ Find eigenvectors with the top k eigenvalues

# PCA: Steps

- Input:  $N \times d$  data matrix  $\mathbf{X}$  (each row contain a  $d$  dimensional data point)
  - $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$
  - $\tilde{\mathbf{X}} \leftarrow$  Mean value of data points is subtracted from rows of  $\mathbf{X}$
  - $\mathbf{S} = \frac{1}{N} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  (Covariance matrix)
  - Calculate eigenvalue and eigenvectors of  $\mathbf{S}$
  - Pick  $d'$  eigenvectors corresponding to the largest eigenvalues and put them in the columns of  $\mathbf{A} = [\mathbf{v}_1, \dots, \mathbf{v}_{d'}]$ 
    - $\downarrow$  First PC
    - $\downarrow$   $d'$ -th PC
  - $\mathbf{X}' = \mathbf{X}\mathbf{A}$

# Reconstruction



$$\mathbf{x}' = \begin{bmatrix} \mathbf{v}_1^T \mathbf{x} \\ \vdots \\ \mathbf{v}_{d'}^T \mathbf{x} \end{bmatrix}$$

$$\mathbf{A} = [\mathbf{v}_1, \dots, \mathbf{v}_{d'}]$$

$$\begin{aligned} \mathbf{x}' &= \mathbf{A}^T (\mathbf{x} - \bar{\mathbf{x}}) \\ \Rightarrow \hat{\mathbf{x}} &= \bar{\mathbf{x}} + \mathbf{A} \mathbf{x}' = \bar{\mathbf{x}} + \mathbf{A} \mathbf{A}^T (\mathbf{x} - \bar{\mathbf{x}}) \end{aligned}$$

- ▶ Incorporating all eigenvectors in  $\mathbf{A} = [\mathbf{v}_1, \dots, \mathbf{v}_d]$ :  
 $\Rightarrow$  If  $d' = d$  then  $\mathbf{x}$  can be reconstructed exactly from  $\mathbf{x}'$

# PCA on Faces: “Eigenfaces”

## ❓ ORL Database



Some Images

# PCA on Faces: “Eigenfaces”



Average  
face

1<sup>st</sup> PC

6<sup>th</sup> PC



For eigen faces

“gray” = 0,

“white” > 0,

“black” < 0

# PCA on Faces:



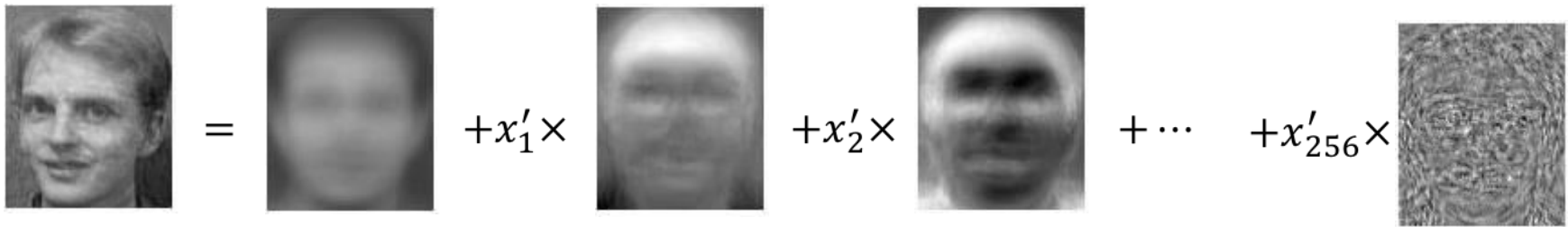
$x$  is a  $112 \times 92 = 10304$  dimensional vector containing intensity of the pixels of this image and  $\tilde{x} = x - \bar{x}$

**Feature vector** =  $[x'_1, x'_2, \dots, x'_{d'}]$

$x'_i$   $\longrightarrow$  The projection of  $x$  on the  $i$ -th PC

$\hat{x}$

$$\hat{x} = \bar{x} + \sum_{i=1}^{d'} (v_i^T \tilde{x}) \times v_i$$



Average  
Face

# PCA on Faces: Reconstructed Face

**$d'=1$**



**$d'=2$**



**$d'=4$**



**$d'=8$**



**$d'=16$**



**$d'=32$**



**$d'=64$**



**$d'=128$**



**$d'=256$**

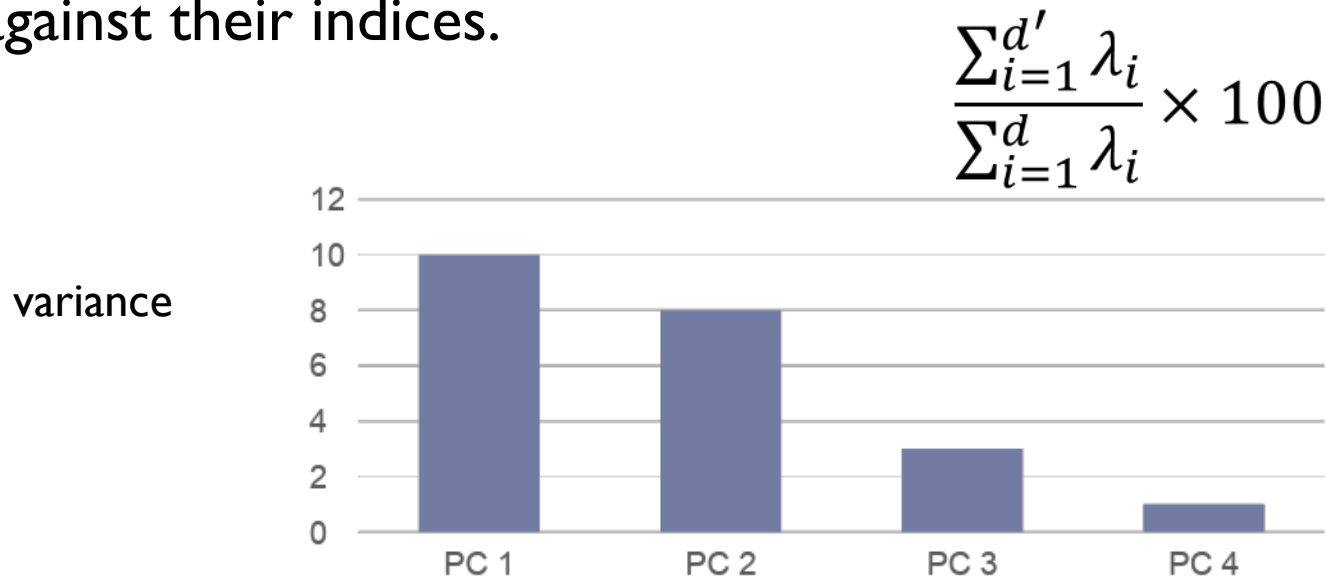


**Original  
Image**



# Dimensionality reduction by PCA

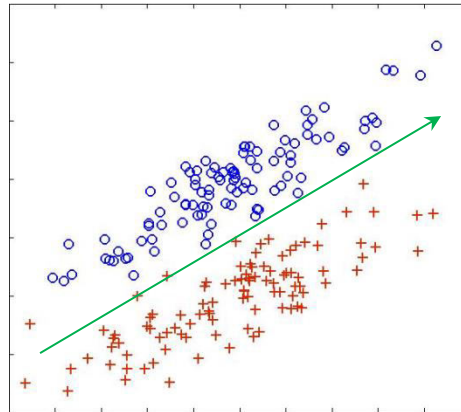
- ❑ Data may lie near a linear subspace of high-dimensional input space
- ❑ Only keep data projections onto principal components with **large** eigenvalues
- ❑ Plot of the eigenvalues (or variances of principal components) against their indices.





# Unsupervised feature extraction drawback

- ❓ PCA drawback: An excellent information packing transform does not necessarily lead to a good class separability.
- ❓ The directions of the maximum variance may be useless for classification purpose



# PCA vs. LDA

- ▶ Although LDA often provide more suitable features for classification tasks, PCA might outperform LDA in some situations:
  - ▶ When there are many unlabeled data while no or small amount of labeled data
    - ▶ when the number of samples per class is small (overfitting problem of LDA)
  - ▶ when the number of the desired features is more than  $C - 1$
  - ▶ when the training data non-uniformly sample the underlying distribution
- ▶ Semi-supervised feature extraction
  - ▶ E.g., PCA+LDA, Regularized LDA, Locally FDA (LFDA)

# PCA: Summary

- ❑ Global optimum is found by eigenvector method
- ❑ No parameter tuning
- ❑ However, it is limited to:
  - ❑ using second order statistics
  - ❑ limited to linear projections

# Resources

- ❑ C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 12.