

# Machine Learning for Data Science (CS4786)

## Lecture 2

Dimensionality Reduction  
&  
Principal Component Analysis

Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2015sp/>

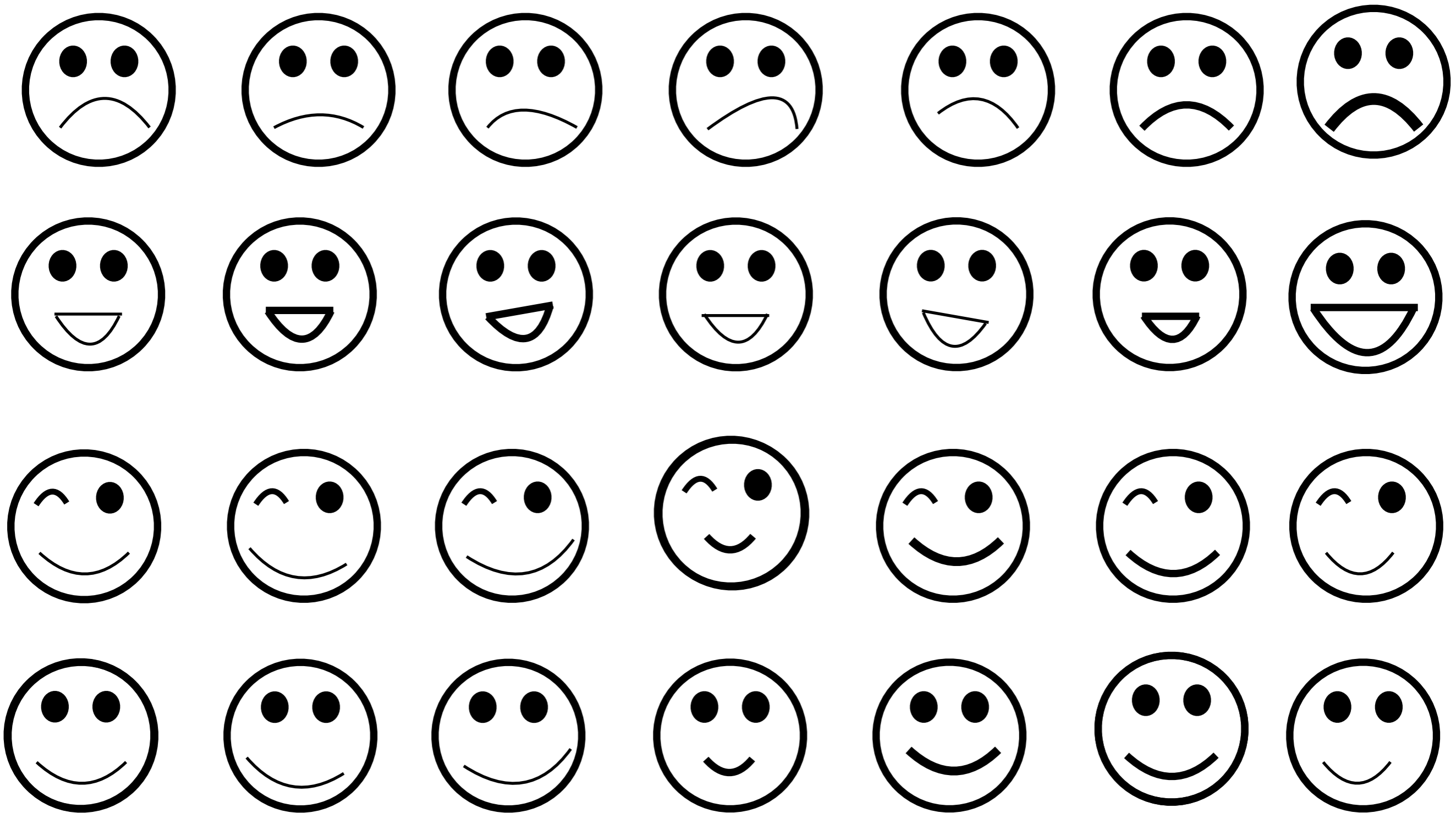
# ANNOUNCEMENTS

- Diagnostic assignment due on 29th (Thursday) beginning of class
- Course webpage is the official source of all class related information
- Please make sure to add both Prof. Lee and Prof. Sridharan on all emails

# REPRESENTING DATA AS FEATURE VECTORS

- How do we represent data?
- Each data-point often represented as vector referred to as feature vector
- Eg. text document represented by vector in which each coordinate represents a word and value represents number of times the word occurred in the document
- Eg. Image represented as a vector where each coordinate represents a pixel and value represents the grayscale value of that pixel

# EXAMPLE: IMAGES



# DIMENSIONALITY REDUCTION

- You are provided with  $n$  data points each in  $\mathbb{R}^d$
- Goal: Compress data into  $n$  points in  $\mathbb{R}^K$  where  $K \ll d$ 
  - Retain as much information about the original data set
  - Retain desired properties of the original data set

# WHY DIMENSIONALITY REDUCTION?

- For computational ease
  - As input to supervised learning algorithm
  - Before clustering to remove redundant information and noise
- Data compression & Noise reduction
- Data visualization

# DIMENSIONALITY REDUCTION

Given feature vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ , compress the data points into low dimensional representation  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$  where  $K \ll d$

# DIMENSIONALITY REDUCTION

Desired properties:

- 1 Original data can be (approximately) reconstructed
- 2 Preserve distances between data points
- 3 “Relevant” information is preserved
- 4 Noise is reduced



# DIM REDUCTION: LINEAR TRANSFORMATION

- Pick a low dimensional subspace
- Project linearly to this subspace
- Subspace retains as much information

# DIM REDUCTION: LINEAR TRANSFORMATION

$$X = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1.1 & 2 & 3 & 4 \\ 3 & 2 & 3 & 4 \\ -1 & 2 & 3 & 4 \\ -0.2 & 2 & 3 & 4 \\ -2 & 2 & 3 & 4 \\ 1.4 & 2 & 3 & 4 \\ 1.4 & 2 & 3 & 4 \\ -0.1 & 2 & 3 & 4 \\ 0.5 & 2 & 3 & 4 \end{bmatrix}$$



# ORTHONORMAL PROJECTIONS

- (Centered) Data-points as linear combination of some orthonormal basis, i.e.

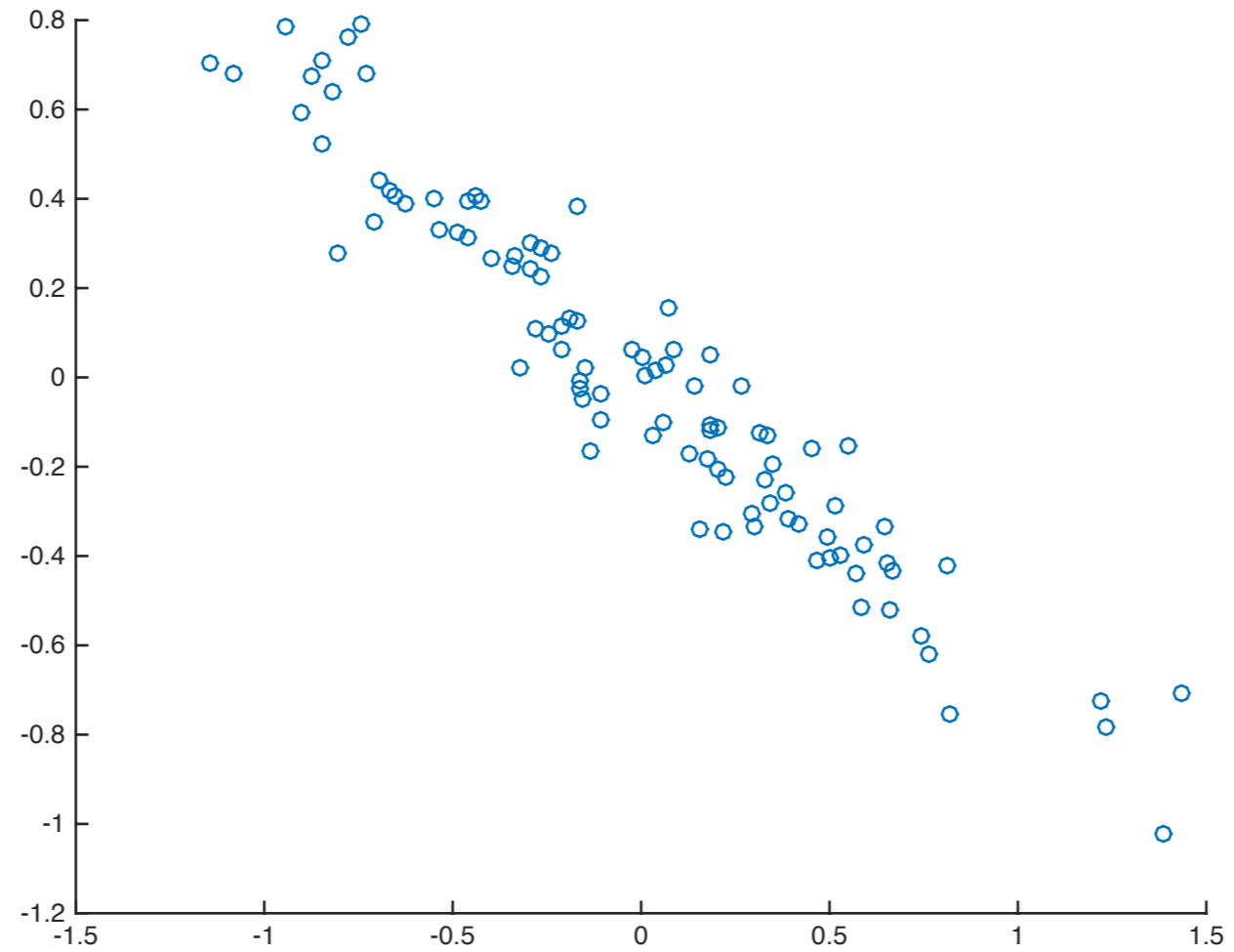
$$\mathbf{x}_t = \boldsymbol{\mu} + \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j$$

where  $\mathbf{w}_1, \dots, \mathbf{w}_d \in \mathbb{R}^d$  are the orthonormal basis and  $\boldsymbol{\mu} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t$ .

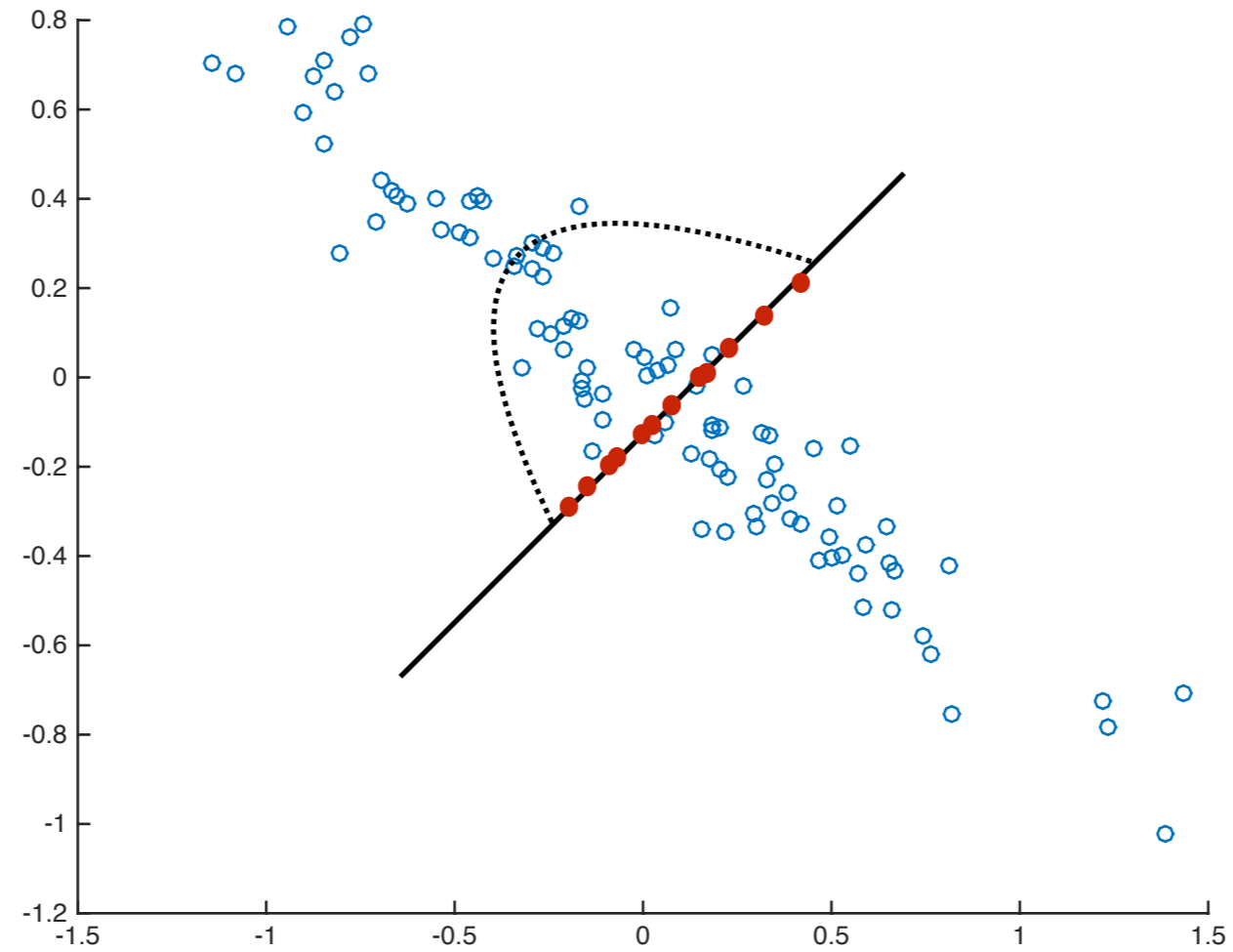
- Represent data as linear combination of just  $K$  orthonormal basis,

$$\hat{\mathbf{x}}_t = \boldsymbol{\mu} + \sum_{j=1}^K \mathbf{y}_t[j] \mathbf{w}_j$$

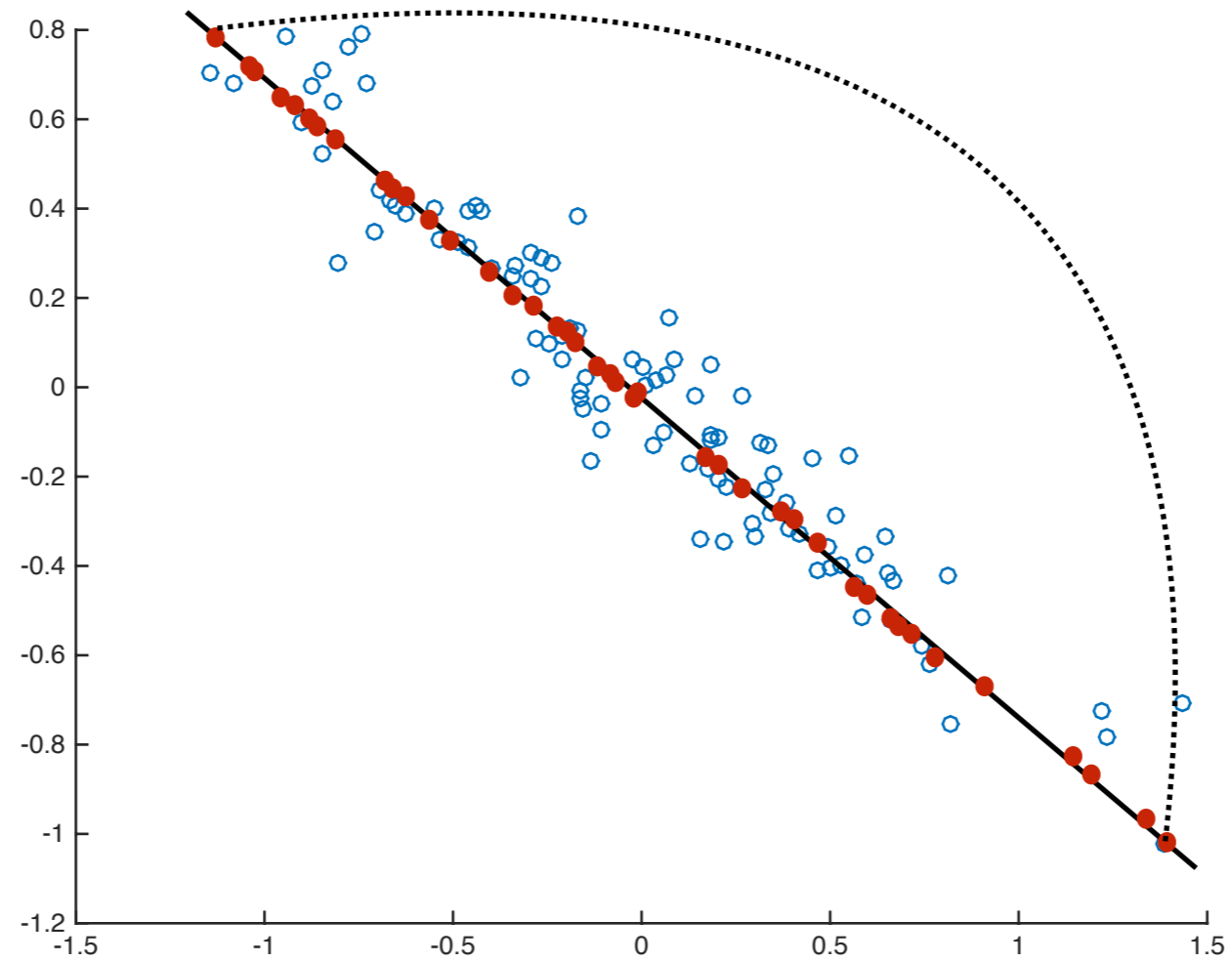
# PCA: VARIANCE MAXIMIZATION



# PCA: VARIANCE MAXIMIZATION



# PCA: VARIANCE MAXIMIZATION



# PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most



# PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2$$

# PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2\end{aligned}$$

# PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}\end{aligned}$$

# PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left( \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w} \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}\end{aligned}$$

$\boldsymbol{\Sigma}$  is the covariance matrix

# PCA: VARIANCE MAXIMIZATION

- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \Sigma \mathbf{w} \quad (1)$$

To solve the above maximization problem, we use Lagrange multipliers. Specifically there exists  $\lambda$  such that solution  $\mathbf{w}_1$  is:

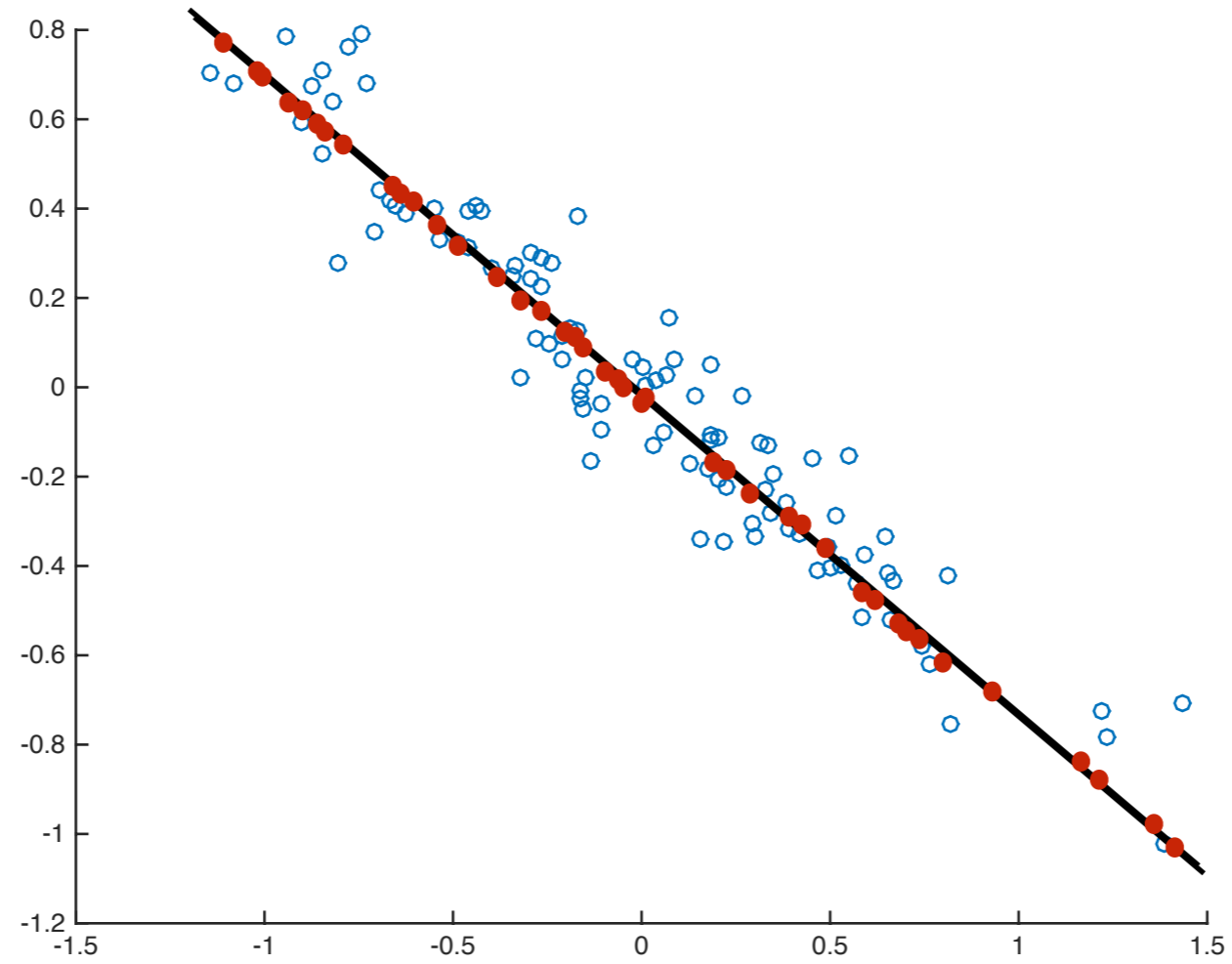
$$\mathbf{w}_1 = \arg \max_{\mathbf{w}} \mathbf{w}^\top \Sigma \mathbf{w} - \lambda \|\mathbf{w}\|_2^2$$

Taking derivative and equality to 0 we find that  $\Sigma \mathbf{w} = \lambda \mathbf{w}$  (ie. eigenvector). Plugging this back into Eq. 1,

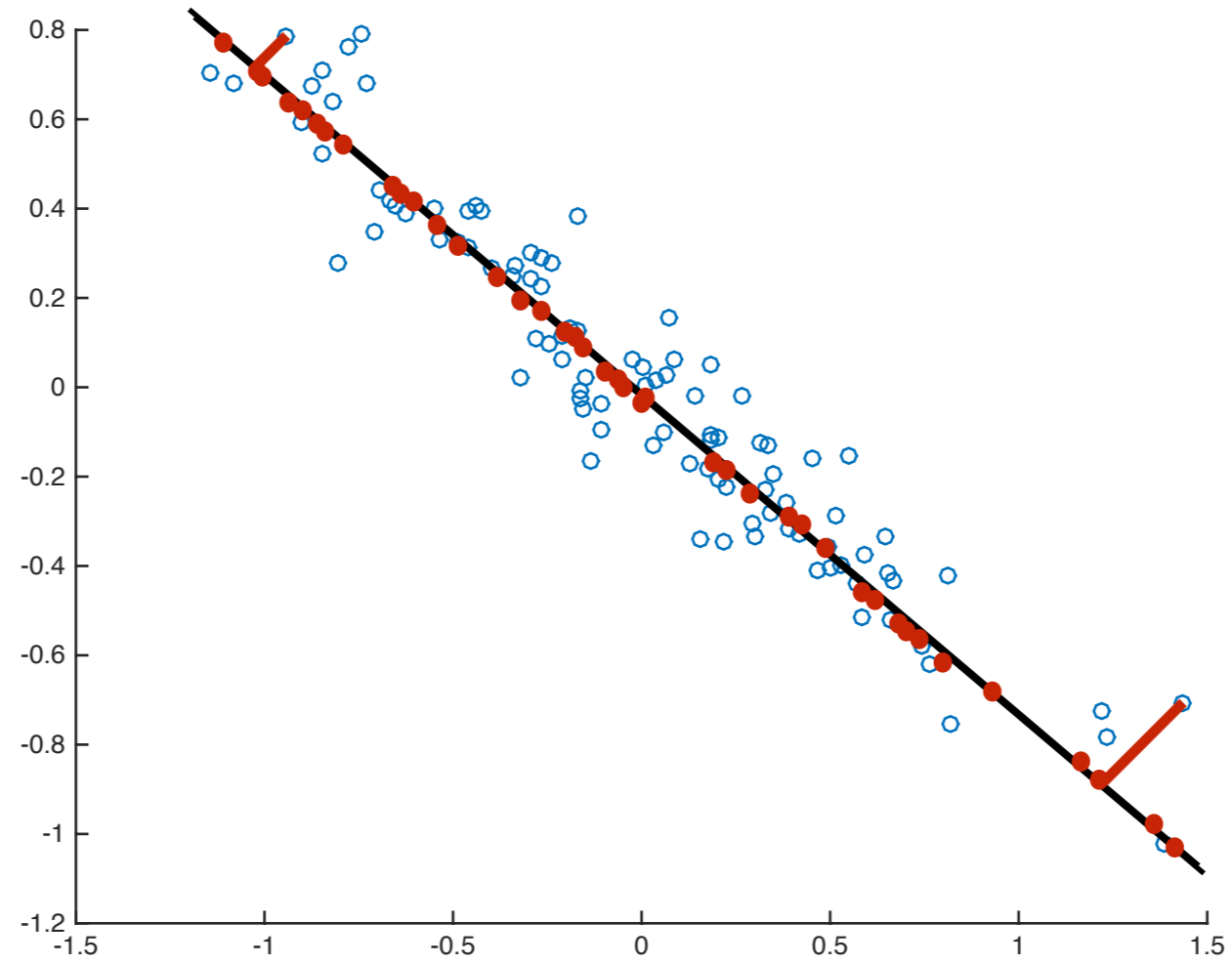
$$\frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \Sigma \mathbf{w} = \mathbf{w}^\top (\lambda \mathbf{w}) = \lambda$$

Hence to maximize variance we pick direction with largest eigenvalue

# PCA: MINIMIZING RECONSTRUCTION ERROR



# PCA: MINIMIZING RECONSTRUCTION ERROR



# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2$$



# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \mu - \mathbf{x}_t \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \mu - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \mu \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \mathbf{x}_t \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \boldsymbol{\mu} \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{note that } \mathbf{y}_t[j] = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}))\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \mathbf{x}_t \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \boldsymbol{\mu} \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{note that } \mathbf{y}_t[j] = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d (\mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \mathbf{w}_j \right\|_2^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \mathbf{x}_t \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \boldsymbol{\mu} \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{note that } \mathbf{y}_t[j] = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d (\mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \mathbf{w}_j \right\|_2^2 \\ &= \sum_{t=1}^n \sum_{j=k+1}^d \left( \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\begin{aligned}\sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \mathbf{x}_t \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=1}^k \mathbf{y}_t[j] \mathbf{w}_j + \boldsymbol{\mu} - \sum_{j=1}^d \mathbf{y}_t[j] \mathbf{w}_j - \boldsymbol{\mu} \right\|_2^2 \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d \mathbf{y}_t[j] \mathbf{w}_j \right\|_2^2 \quad (\text{note that } \mathbf{y}_t[j] = \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \\ &= \sum_{t=1}^n \left\| \sum_{j=k+1}^d (\mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu})) \mathbf{w}_j \right\|_2^2 \\ &= \sum_{t=1}^n \sum_{j=k+1}^d \left( \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2 = \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j\end{aligned}$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j$$



# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j = \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j = \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

Minimize w.r.t.  $\mathbf{w}$ 's that are orthonormal,

$$\operatorname{argmin}_{\forall j, \|\mathbf{w}_j\|_2=1} \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j = \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

Minimize w.r.t.  $\mathbf{w}$ 's that are orthonormal,

$$\operatorname{argmin}_{\forall j, \|\mathbf{w}_j\|_2=1} \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

Using Lagrangian multipliers, there exists  $\lambda_{k+1}, \dots, \lambda_d$  such that solution to above is given by:

$$\operatorname{minimize} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j + \sum_{j=k+1}^d \lambda_j \|\mathbf{w}_j\|_2^2$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Goal: find the basis that minimizes reconstruction error,

$$\frac{1}{n} \sum_{t=1}^n \|\hat{\mathbf{x}}_t - \mathbf{x}_t\|_2^2 = \frac{1}{n} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}_j = \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

Minimize w.r.t.  $\mathbf{w}$ 's that are orthonormal,

$$\operatorname{argmin}_{\forall j, \|\mathbf{w}_j\|_2=1} \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j$$

Using Lagrangian multipliers, there exists  $\lambda_{k+1}, \dots, \lambda_d$  such that solution to above is given by:

$$\operatorname{minimize} \sum_{t=1}^n \sum_{j=k+1}^d \mathbf{w}_j^\top \boldsymbol{\Sigma} \mathbf{w}_j + \sum_{j=k+1}^d \lambda_j \|\mathbf{w}_j\|_2^2$$

Setting derivate to 0,  $\boldsymbol{\Sigma} \mathbf{w}_j = \lambda_j \mathbf{w}_j$ . That is  $\mathbf{w}_j$ 's are eigenvectors and  $\lambda_j$ 's are eigenvalues.

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Solution :  $\mathbf{w}_j$ 's are eigenvectors and  $\lambda_j$ 's are corresponding eigenvalues

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Solution :  $\mathbf{w}_j$ 's are eigenvectors and  $\lambda_j$ 's are corresponding eigenvalues
- Further, reconstruction error can be written as:

$$\operatorname{argmin}_{\mathbf{w}: \|\mathbf{w}_j\|_2=1} \sum_{j=k+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j = \sum_{j=k+1}^d \lambda_j \mathbf{w}_j^\top \mathbf{w}_j = \sum_{j=k+1}^d \lambda_j$$

# PCA: MINIMIZING RECONSTRUCTION ERROR

- Solution :  $\mathbf{w}_j$ 's are eigenvectors and  $\lambda_j$ 's are corresponding eigenvalues
- Further, reconstruction error can be written as:

$$\operatorname{argmin}_{\mathbf{w}: \|\mathbf{w}_j\|_2=1} \sum_{j=k+1}^d \mathbf{w}_j^\top \Sigma \mathbf{w}_j = \sum_{j=k+1}^d \lambda_j \mathbf{w}_j^\top \mathbf{w}_j = \sum_{j=k+1}^d \lambda_j$$

- Clearly to minimize reconstruction error, we need to minimize  $\sum_{j=k+1}^d \lambda_j$ . In other words we discard the  $d - k$  directions that have the smallest eigenvalue

# PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components



# PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top  $K$  principal components are the eigenvectors with  $K$  largest eigenvalues

# PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top  $K$  principal components are the eigenvectors with  $K$  largest eigenvalues
- $\text{Projection} = \text{Data} \times \text{Top } K \text{ Eigenvectors}$

# PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top  $K$  principal components are the eigenvectors with  $K$  largest eigenvalues
- Projection = Data  $\times$  Top  $K$  eigenvectors
- Reconstruction = Projection  $\times$  Transpose of top  $K$  eigenvectors

# PRINCIPAL COMPONENT ANALYSIS

- Eigenvectors of the covariance matrix are the principal components
- Top  $K$  principal components are the eigenvectors with  $K$  largest eigenvalues
- $\text{Projection} = \text{Data} \times \text{Top } K \text{ eigenvectors}$
- $\text{Reconstruction} = \text{Projection} \times \text{Transpose of top } K \text{ eigenvectors}$
- Independently discovered by Pearson in 1901 and Hotelling in 1933.

# PRINCIPAL COMPONENT ANALYSIS: DEMO