

## Dimensionality Reduction

- Many dimensions are often interdependent (correlated);

We can:

- Reduce the dimensionality of problems;
- Transform interdependent coordinates into significant and independent ones;

## Principal Component Analysis

## Principal Component Analysis -- PCA (also called Karhunen-Loeve transformation)

- **PCA** transforms the original input space into a lower dimensional space, by constructing dimensions that are linear combinations of the given features;
- The objective is to consider **independent** dimensions along which data have **largest variance** (i.e., greatest variability);

## Principal Component Analysis -- PCA

- **PCA** involves a linear algebra procedure that transforms a number of possibly correlated variables into a smaller number of uncorrelated variables called **principal components**;
- The first principal component accounts for as much of the variability in the data as possible;
- Each succeeding component (orthogonal to the previous ones) accounts for as much of the remaining variability as possible.

## Principal Component Analysis -- PCA

- So: PCA finds  $n$  linearly transformed components  $s_1, s_2, \dots, s_n$  so that they explain the maximum amount of variance;
- We can define PCA in an intuitive way using a recursive formulation:

## Principal Component Analysis -- PCA

- Suppose data are first centered at the origin (i.e., their mean is  $\mathbf{0}$ );
- We define the direction of the first principal component, say  $w_1$ , as follows

$$w_1 = \arg \max_{\|w\|=1} E[(w^T x)^2]$$

where  $w_1$  is of the same dimensionality  $q$  as the data vector  $x$

- Thus: the first principal component is the projection on the direction along which the variance of the projection is maximized.

## Principal Component Analysis -- PCA

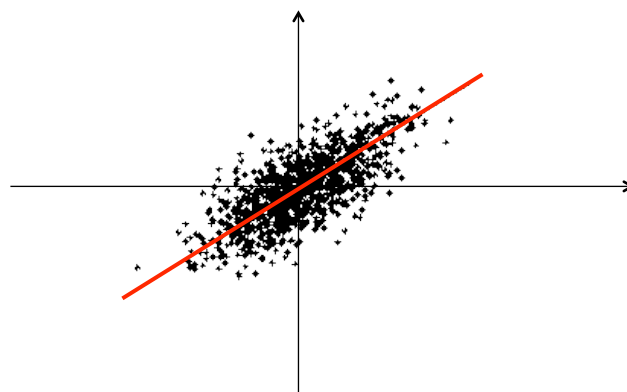
- Having determined the first  $k-1$  principal components, the  $k$ -th principal component is determined as the principal component of the data residual:

$$\mathbf{w}_k = \arg \max_{\|\mathbf{w}\|=1} E\left\{[\mathbf{w}^T (\mathbf{x} - \sum_{i=1}^{k-1} \mathbf{w}_i \mathbf{w}_i^T \mathbf{x})]^2\right\}$$

- The principal components are then given by:

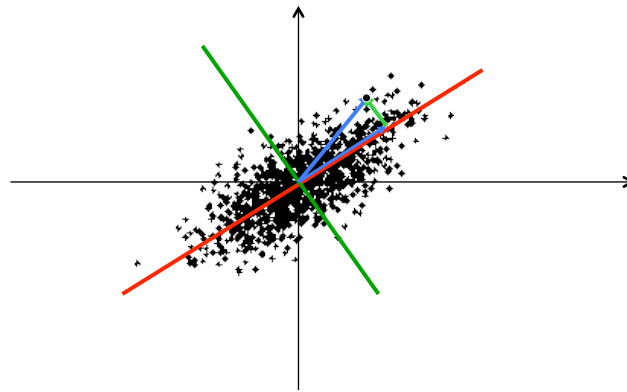
$$s_i = \mathbf{w}_i^T \mathbf{x}$$

### Simple illustration of PCA



First principal component of a two-dimensional data set.

## Simple illustration of PCA



Second principal component of a two-dimensional data set.

## PCA – Geometric interpretation

Basically:

PCA rotates the data (centered at the origin) in such a way that the maximum variability is visible (i.e., aligned with the axes.)

### PCA – How to compute the principal components

Let  $\mathbf{w}$  be the direction of the first principal component, with  $\|\mathbf{w}\| = 1$

$s_i = \mathbf{w}^T \mathbf{x}_i$  is the projection of  $\mathbf{x}_i$  along  $\mathbf{w}$

$$\bar{s} = \frac{1}{N} \sum_{i=1}^N s_i = \frac{1}{N} \sum_{i=1}^N \mathbf{w}^T \mathbf{x}_i$$

Variance of data along  $\mathbf{w}$  :

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (s_i - \bar{s})^2 &= \\ \frac{1}{N} \sum_{i=1}^N \left( \mathbf{w}^T \mathbf{x}_i - \frac{1}{N} \sum_{j=1}^N \mathbf{w}^T \mathbf{x}_j \right)^2 & \end{aligned}$$

### PCA – How to compute the principal components

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (s_i - \bar{s})^2 &= \\ \frac{1}{N} \sum_{i=1}^N \left( \mathbf{w}^T \mathbf{x}_i - \frac{1}{N} \sum_{j=1}^N \mathbf{w}^T \mathbf{x}_j \right)^2 &= \\ \frac{1}{N} \sum_{i=1}^N \left[ \mathbf{w}^T \left( \mathbf{x}_i - \frac{1}{N} \sum_{j=1}^N \mathbf{x}_j \right) \right]^2 &= \\ \frac{1}{N} \sum_{i=1}^N [\mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}})]^2 &= \\ \frac{1}{N} \sum_{i=1}^N [\mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}] &= \\ \mathbf{w}^T \left[ \frac{1}{N} \sum_{i=1}^N [(\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T] \right] \mathbf{w} &= \mathbf{w}^T \Sigma \mathbf{w} \end{aligned}$$

Sample covariance matrix

## PCA – How to compute the principal components

Thus : the variance of data along direction  $\boldsymbol{w}$  can be written as

$$\boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w}$$

Our objective is to find  $\boldsymbol{w}$  such that

$$\boldsymbol{w} = \arg \max_{\boldsymbol{w}} \boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w}$$

with the constraint  $\boldsymbol{w}^T \boldsymbol{w} = 1$

By introducing one Lagrange multiplier  $\lambda$ , we obtain the following unconstrained optimization problem

$$\boldsymbol{w} = \arg \max_{\boldsymbol{w}} [\boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} - \lambda(\boldsymbol{w}^T \boldsymbol{w} - 1)]$$

Setting  $\frac{\partial}{\partial \boldsymbol{w}} = 0$  gives:  $2\boldsymbol{\Sigma} \boldsymbol{w} - 2\lambda \boldsymbol{w} = 0$

That is:  $\boldsymbol{\Sigma} \boldsymbol{w} = \lambda \boldsymbol{w}$

**Our problem is reduced to an eigenvalue problem**

## PCA – How to compute the principal components

Thus : the variance of data along direction  $\boldsymbol{w}$  can be written as

$$\boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w}$$

Our objective is to find  $\boldsymbol{w}$  such that

$$\boldsymbol{w} = \arg \max_{\boldsymbol{w}} \boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w}$$

with the constraint  $\boldsymbol{w}^T \boldsymbol{w} = 1$

By introducing one Lagrange multiplier  $\lambda$ , we obtain the following unconstrained optimization problem

$$\boldsymbol{w} = \arg \max_{\boldsymbol{w}} [\boldsymbol{w}^T \boldsymbol{\Sigma} \boldsymbol{w} - \lambda(\boldsymbol{w}^T \boldsymbol{w} - 1)]$$

Setting  $\frac{\partial}{\partial \boldsymbol{w}} = 0$  gives:  $2\boldsymbol{\Sigma} \boldsymbol{w} - 2\lambda \boldsymbol{w} = 0$

That is:  $\boldsymbol{\Sigma} \boldsymbol{w} = \lambda \boldsymbol{w}$

**The solution  $\boldsymbol{w}$  is the eigenvector of  $\boldsymbol{\Sigma}$  corresponding to the largest eigenvalue  $\lambda$**

## PCA -- Summary

- The computation of the  $w_i$  is accomplished by solving an eigenvalue problem for the sample **covariance** matrix (assuming data have  $0$  mean):

$$\Sigma = E[x x^T]$$

- The eigenvector associated with the **largest eigenvalue** corresponds to the **first** principal component; the eigenvector associated with the **second largest** eigenvalue corresponds to the **second** principal component; and so on...
- Thus: The  $w_i$  are the eigenvectors of  $\Sigma$  that correspond to the  $n$  largest eigenvalues of  $\Sigma$

## PCA -- In practice

- The basic goal of PCA is to reduce the dimensionality of the data. Thus, one usually chooses:

$$n \ll q$$

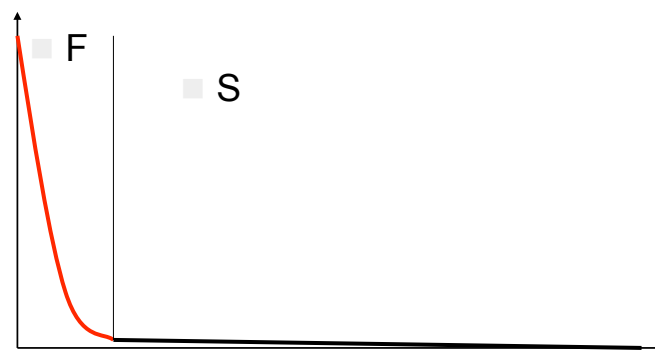
- But how do we select the number of components  $n$  ?



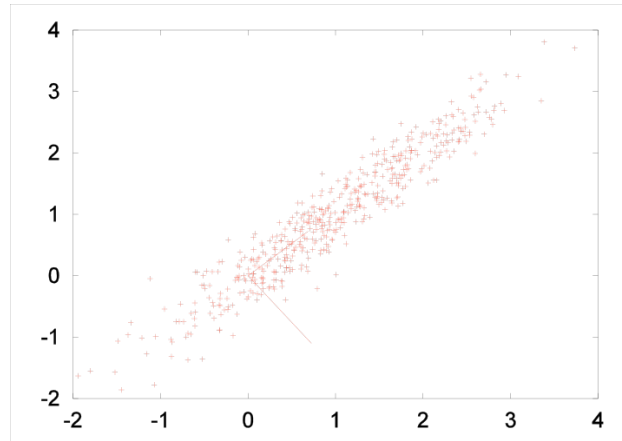
## Determining the number of components

- Plot the eigenvalues – each eigenvalue is related to the amount of variation explained by the corresponding axis (eigenvector);
- If the points on the graph tend to level out (show an “elbow” shape), these eigenvalues are usually close enough to zero that they can be ignored.
- In general: Limit the variance accounted for.

## Critical information lies in low dimensional subspaces



- A typical eigenvalue spectrum and its division into two orthogonal subspaces



$$\lambda_1 = 1.98, \lambda_2 = 0.05$$

### Determining the number of components

$$\mathbf{x}_i \in \mathfrak{R}^q, \quad i = 1, \dots, N$$

$\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q$  :  $q$  eigenvectors (principal component directions)

$\|\mathbf{w}_i\| = 1$  (the  $\mathbf{w}_i$ s are orthonormal vectors)

Representation of  $\mathbf{x}_i$  in eigenvector space :

$$\mathbf{y}_i = (\mathbf{w}_1^T \mathbf{x}_i) \mathbf{w}_1 + (\mathbf{w}_2^T \mathbf{x}_i) \mathbf{w}_2 + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q$$

Suppose we retain the first  $k$  principal components :

$$\mathbf{y}_i^k = (\mathbf{w}_1^T \mathbf{x}_i) \mathbf{w}_1 + (\mathbf{w}_2^T \mathbf{x}_i) \mathbf{w}_2 + \dots + (\mathbf{w}_k^T \mathbf{x}_i) \mathbf{w}_k$$

Then :

$$\mathbf{y}_i - \mathbf{y}_i^k = (\mathbf{w}_{k+1}^T \mathbf{x}_i) \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q$$

### Determining the number of components

$$\begin{aligned}
 (\mathbf{y}_i - \mathbf{y}_i^k)^T (\mathbf{y}_i - \mathbf{y}_i^k) &= \\
 [(\mathbf{w}_{k+1}^T \mathbf{x}_i) \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q]^T [(\mathbf{w}_{k+1}^T \mathbf{x}_i) \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q] &= \\
 \mathbf{w}_{k+1}^T (\mathbf{w}_{k+1}^T \mathbf{x}_i)^2 \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T (\mathbf{w}_q^T \mathbf{x}_i)^2 \mathbf{w}_q &= \\
 \text{(note } \mathbf{w}_i^T \mathbf{w}_j = 0 \ \forall i \neq j \text{ since } \mathbf{w}_i \text{ and } \mathbf{w}_j \text{ are orthogonal vectors)} & \\
 (\mathbf{w}_{k+1}^T \mathbf{x}_i)^2 \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i)^2 \mathbf{w}_q^T \mathbf{w}_q &= \\
 (\mathbf{w}_{k+1}^T \mathbf{x}_i)^2 + \dots + (\mathbf{w}_q^T \mathbf{x}_i)^2 &= \\
 (\mathbf{w}_{k+1}^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{w}_{k+1}) + \dots + (\mathbf{w}_q^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{w}_q) &= \\
 \mathbf{w}_{k+1}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_q &
 \end{aligned}$$

### Determining the number of components

$$\begin{aligned}
 \frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{y}_i^k)^T (\mathbf{y}_i - \mathbf{y}_i^k) &= \text{Mean square error} \\
 \frac{1}{N} \sum_{i=1}^N [(\mathbf{w}_{k+1}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_q)] &= \\
 \mathbf{w}_{k+1}^T \left[ \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right] \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \left[ \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right] \mathbf{w}_q &= \\
 \mathbf{w}_{k+1}^T \Sigma \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \Sigma \mathbf{w}_q &
 \end{aligned}$$

We have:  $\Sigma \mathbf{w}_{k+1} = \lambda_{k+1} \mathbf{w}_{k+1}, \dots, \Sigma \mathbf{w}_q = \lambda_q \mathbf{w}_q$

Thus:

$$\begin{aligned}
 \mathbf{w}_{k+1}^T \Sigma \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \Sigma \mathbf{w}_q &= \\
 \mathbf{w}_{k+1}^T \lambda_{k+1} \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \lambda_q \mathbf{w}_q &= \\
 \lambda_{k+1} + \dots + \lambda_q &
 \end{aligned}$$

**Determining the number of components**

$$\frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{y}_i^k)^T (\mathbf{y}_i - \mathbf{y}_i^k) = \lambda_{k+1} + \dots + \lambda_q$$

⇒ The mean square error of the truncated representation is equal to the sum of the remaining eigenvalues.

**In general: choose  $k$  so that 90-95% of the variance of the data is captured.**