

## Feature Selection: Linear Transformations

$$\underline{y}_{new} = M \underline{x}_{old}$$

## Constraint Optimization (insertion)

**Problem:** Given an objective function  $f(x)$  to be optimized and let constraints be given by  $h_k(x) = c_k$ , moving constants to the left,  $\Rightarrow h_k(x) - c_k = g_k(x)$ .  $f(x)$  and  $g_k(x)$  must have continuous first partial derivatives

### A Solution:

Lagrangian Multipliers  $0 = \nabla_x f(x) + \sum \nabla_x \lambda_k g_k(x)$

or starting with the Lagrangian :  $L(x, \lambda) = f(x) + \sum \lambda_k g_k(x)$ .

with  $\nabla_x L(x, \lambda) = 0$ .

# The Covariance Matrix (insertion)

## Definition

Let  $\underline{x} = \{x_1, \dots, x_N\} \in \mathbb{R}^N$  be a real valued random variable (data vectors), with the expectation value of the mean  $E[\underline{x}] = \mu$ .

We define the **covariance matrix**  $\Sigma_x$  of a random variable  $\underline{x}$  as  $\Sigma_x := E[(\underline{x} - \mu)(\underline{x} - \mu)^T]$  with matrix elements  $\Sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)^T]$ .

**Application:** Estimating  $E[\underline{x}]$  and  $E[(\underline{x} - E[\underline{x}])(\underline{x} - E[\underline{x}])^T]$  from data.

We assume  $m$  samples of the random variable  $\underline{x} = \{x_1, \dots, x_N\} \in \mathbb{R}^N$  that is we have a set of  $m$  vectors  $\{\underline{x}_1, \dots, \underline{x}_m\} \in \mathbb{R}^N$  or when put into a data matrix  $X \in \mathbb{R}^{N \times m}$

Maximum Likelihood estimators

for  $\mu$  and  $\Sigma_x$  are:

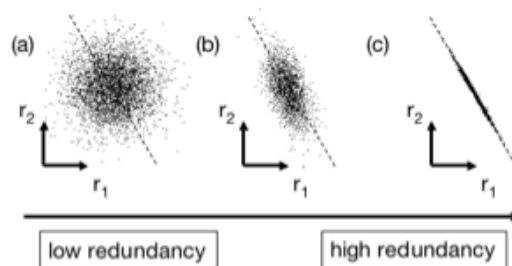
$$\Rightarrow \underline{\mu}_{ML} = \frac{1}{m} \sum_{k=1}^m x_k$$

$$\Rightarrow \Sigma_{ML} = \frac{1}{m} \sum_{k=1}^m (x_k - \underline{\mu}_{ML})(x_k - \underline{\mu}_{ML})^T = \frac{1}{m} X X^T$$

## KLT/PCA Motivation

5

- Find meaningful "directions" in correlated data
- Linear dimensionality reduction
- Visualization of higher dimensional data
- Compression / Noise reduction
- PDF-Estimate



## Karhunen-Loève Transform: 1<sup>st</sup> Derivation <sup>7</sup>

### Problem

Let  $\underline{x} = \{x_1, \dots, x_N\} \in \mathbb{R}^N$  be a feature vector of zero mean, real valued random variables.

We seek the direction  $\underline{a}_1$  of maximum variance:

$\Rightarrow y_1 = \underline{a}_1^T \underline{x}$  for which  $\underline{a}_1$  is such as  $E[y_1^2]$  is maximum with the constraint that  $\underline{a}_1^T \underline{a}_1 = 1$

This is a constrained optimization  $\rightarrow$  use of the Lagrangian:

$$\begin{aligned} L(\underline{a}_1, \lambda_1) &= E[\underline{a}_1^T \underline{x} \underline{x}^T \underline{a}_1] - \lambda_1 (\underline{a}_1^T \underline{a}_1 - 1) \\ &= \underline{a}_1^T \Sigma_{\underline{x}} \underline{a}_1 - \lambda_1 (\underline{a}_1^T \underline{a}_1 - 1) \end{aligned}$$

Lagrange multiplier

## Karhunen-Loève Transform <sup>8</sup>

$$L(\underline{a}_1, \lambda_1) = \underline{a}_1^T \Sigma_{\underline{x}} \underline{a}_1 - \lambda_1 (\underline{a}_1^T \underline{a}_1 - 1)$$

for  $E[y_1^2]$  to be maximum :  $\frac{\partial L(\underline{a}_1, \lambda_1)}{\partial \underline{a}_1} = 0$

$$\Rightarrow \Sigma_{\underline{x}} \underline{a}_1 - \lambda_1 \underline{a}_1 = 0$$

$\Rightarrow \underline{a}_1$  must be *eigenvector* of  $\Sigma_{\underline{x}}$  with *eigenvalue*  $\lambda_1$ .

$$E[y_1^2] = \underline{a}_1^T \Sigma_{\underline{x}} \underline{a}_1 = \lambda_1$$

$\Rightarrow$  for  $E[y_1^2]$  to be maximum,  $\lambda_1$  must be the *largest* eigenvalue.

# Karhunen-Loève Transform

9

Now let's search for a second direction,  $\underline{a}_2$ , such that:

$$y_2 = \underline{a}_2^T \underline{x} \quad \text{such as } E[y_2^2] \text{ is maximum, and}$$
$$\underline{a}_2^T \underline{a}_1 = 0 \quad \text{and} \quad \underline{a}_2^T \underline{a}_2 = 1$$

Similar derivation:  $L(\underline{a}_2, \lambda_2) = \underline{a}_2^T \Sigma_x \underline{a}_2 - \lambda_2 (\underline{a}_2^T \underline{a}_2 - 1)$  with  $\underline{a}_2^T \underline{a}_1 = 0$

=>  $\underline{a}_2$  must be the *eigenvector* of  $\Sigma_x$  associated with the *second largest eigenvalue*  $\lambda_2$ .

We can derive  $N$  orthonormal directions that maximize the variance:  $A = [\underline{a}_1, \underline{a}_2, \dots, \underline{a}_N]$  and  $\underline{y} = A^T \underline{x}$

The resulting matrix  $A$  is known as Principal Component Analysis (PCA) or Karhunen-Loève transform (KLT)  $\underline{y} = A^T \underline{x}$   $\underline{x} = \sum_{i=1}^N y_i \underline{a}_i$

# Karhunen-Loève Transform: 2<sup>nd</sup> Derivation

10

## Problem

Let  $\underline{x} = \{x_1, \dots, x_N\} \in \mathbb{R}^N$  be a feature vector of ~~zero mean~~, real valued random variables.

We seek a transformation  $A$  of  $\underline{x}$  that results in a new set of variables  $\underline{y} = A^T \underline{x}$  (feature vectors) which are uncorrelated ( i.e.  $E[y_i y_j] = 0$  for  $i \neq j$  ).

- Let  $\underline{y} = A^T \underline{x}$ , then by definition of the correlation matrix:

$$R_y \equiv E[\underline{y} \underline{y}^T] = E[A^T \underline{x} \underline{x}^T A] = A^T R_x A$$

- $R_x$  is symmetric  $\Rightarrow$  its eigenvectors are mutually orthogonal

# Karhunen-Loève Transform

11

- i.e. if we choose  $A$  such that its columns  $a_i$  are orthonormal eigenvectors of  $R_x$ , we get:

$$R_y = A^T R_x A = \Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_N \end{bmatrix}$$

- If we further assume  $R_x$  to be positive definite,  
---- > the eigenvalues  $\lambda_i$  will be positive.

The resulting matrix  $A$  is known as

Karhunen-Loève transform (KLT)  $\underline{y} = A^T \underline{x}$   $\underline{x} = \sum_{i=1}^N y_i \underline{a}_i$

# Karhunen-Loève Transform

12

The Karhunen-Loève transform (KLT)

$$\underline{y} = A^T \underline{x} \quad \underline{x} = \sum_{i=1}^N y_i \underline{a}_i$$

For mean-free vectors ( e.g. replace  $\underline{x}$  by  $\underline{x} - E[\underline{x}]$  )

this process diagonalizes the covariance matrix  $\Sigma_y$

## KLT Properties: MSE-Approximation

13

We define a new vector  $\hat{x}$  in  $m$ -dimensional subspace ( $m < N$ ),

using only  $m$  basis vectors:

$$\hat{x} = \sum_{i=1}^m y_i \mathbf{a}_i$$

- Projection of  $\mathbf{x}$  into the subspace spanned by the  $m$  used (orthonormal) eigenvectors.

Now, what is the expected mean square error between  $\mathbf{x}$  and its projection  $\hat{x}$  :

$$E \left[ \|\mathbf{x} - \hat{x}\|^2 \right] = E \left[ \left\| \sum_{i=m+1}^N y_i \mathbf{a}_i \right\|^2 \right] = E \left[ \sum_i \sum_j (y_i \mathbf{a}_i^T)(y_j \mathbf{a}_j) \right] = \dots$$



## KLT Properties: MSE-Approximation

14

$$E \left[ \|\mathbf{x} - \hat{x}\|^2 \right] = \dots = E \left[ \sum_i \sum_j (y_i \mathbf{a}_i^T)(y_j \mathbf{a}_j) \right] = \sum_{i=m+1}^N E \left[ y_i^2 \right] = \sum_{i=m+1}^N \lambda_i$$

The error is minimized if we choose as basis those eigenvectors corresponding to the  $m$  largest eigenvalues of the correlation matrix.

- Amongst all other possible orthogonal transforms KLT is the one leading to minimum MSE

This form of KLT ( as presented here ) is also referred to as **Principal Component Analysis** (PCA).  
The principal components are the eigenvectors ordered (desc.) by their respective eigenvalue magnitudes  $\lambda_i$

## KLT Properties

### Total variance

- Let w.l.o.g.  $E[\mathbf{x}] = 0$  and  $\mathbf{y} = A^T \mathbf{x}$  the KLT (PCA) of  $\mathbf{x}$ .  
From the previous definitions we get:

$$\sigma_{y_i}^2 \equiv E[y_i^2] = \lambda_i$$

- i.e. the eigenvalues of the input covariance matrix are equal to the variances of the transformed coordinates.

- Selecting those features corresponding to  $m$  largest eigenvalues retains the maximal possible total variance (sum of component variances) associated with the original random variables  $x_i$ .

## KLT Properties: Entropy

For a random vector  $\mathbf{y}$  the entropy  $H_y = -E[\ln p_y(\mathbf{y})]$  is a measure for the randomness of the underlying process.

**Example:** for a zero-mean ( $\mu=0$ )  $m$ -dim. Gaussian

$$H_y = -E \left[ \ln \left( (2\pi)^{-\frac{m}{2}} |\Sigma_y|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \underline{\mathbf{y}}^T \Sigma_y^{-1} \underline{\mathbf{y}}\right) \right) \right]$$

$$H_y = \frac{m}{2} \ln(2\pi) + \frac{1}{2} \ln |\Sigma_y| + \frac{1}{2} E[\mathbf{y}^T \Sigma_y^{-1} \mathbf{y}]$$

$$= \frac{m}{2} \ln(2\pi) + \frac{1}{2} \ln \prod_{i=1}^m \lambda_i + \frac{m}{2}$$

$$\begin{aligned} E[\mathbf{y}^T \Sigma_y^{-1} \mathbf{y}] &= E[\text{trace}\{\mathbf{y}^T \Sigma_y^{-1} \mathbf{y}\}] \\ &= E[\text{trace}\{\Sigma_y^{-1} \mathbf{y} \mathbf{y}^T\}] \\ &= E[\text{trace}\{I\}] = m \end{aligned}$$

- Selecting those features corresponding to  $m$  largest eigenvalues maximizes the entropy in the remaining features.
- No wonder: variance and randomness are directly related !

## Computing a PCA:

17

**Problem:** Given mean free data  $X$ , a set of  $n$  feature vectors  $x_i \in \mathbb{R}^m$ . Compute the orthonormal eigenvectors  $a_i$  of the correlation matrix  $R_x$ .

- There are many algorithms that can compute very efficiently eigenvectors of a matrix. However, most of these methods can be very unstable in certain special cases.
- Here we present **SVD**, a method that is in general not the most efficient one. However, the method can be made numerically stable very easily!

## Computing a PCA:

18

### **Singular Value Decomposition:**

an Excursus to Linear Algebra  
( without Proofs )



## Singular Value Decomposition :

19

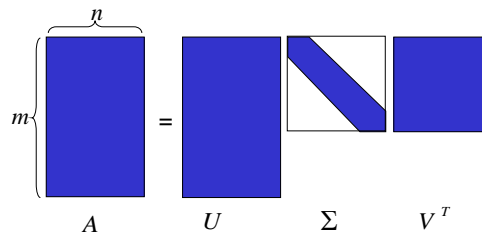
SVD (reduced Version): For matrices  $A \in \mathbb{R}^{m \times n}$  with  $m \geq n$ , there exist matrices

$U \in \mathbb{R}^{m \times n}$  with orthonormal columns ( $U^T U = I$ ),

$V \in \mathbb{R}^{n \times n}$  orthogonal ( $V^T V = I$ ),

$\Sigma \in \mathbb{R}^{n \times n}$  diagonal,

with  $A = U \Sigma V^T$



- The diagonal values of  $\Sigma$  ( $\sigma_1, \sigma_2, \dots, \sigma_n$ ) are called the **singular values**.
- It is accustomed to sort them:  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$

## SVD Applications:

20

SVD is an all-rounder !

Once you have  $U, \Sigma, V$ , you can use it to:

- Solve Linear Systems:  $A \underline{x} = \underline{b}$ 
  - a) If  $A^{-1}$  exists  $\rightarrow$  Compute matrix inverse
  - b) for fewer equations than unknowns
  - c) for more equations than unknowns
  - d) if there is no solution: compute
 
$$x \text{ that } |A \underline{x} - \underline{b}| = \min$$
  - e) compute rank (numerical rank) of a matrix
- .....
- Compute PCA / KLT

## SVD : Matrix inverse $A^{-1}$

21

$$A \underline{x} = \underline{b} :$$

$$A = U \Sigma V^T \quad U, \Sigma, V, \text{ exist for all } A$$

If  $A$  is square  $n \times n$  and not singular, then  $A^{-1}$  exists.

$$\begin{aligned} A^{-1} &= (U \Sigma V^T)^{-1} \\ &= (V^T)^{-1} \Sigma^{-1} U^{-1} \\ &= V \begin{bmatrix} \frac{1}{\sigma_1} & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_n} \end{bmatrix} U^T \end{aligned}$$

Computing  $A^{-1}$  for a singular  $A$  !?

Since  $U, \Sigma, V$  all exist, the only problem can originate if one  $\sigma_i = 0$  or numerically close to zero.

--> singular values indicate if  $A$  is singular or not!!

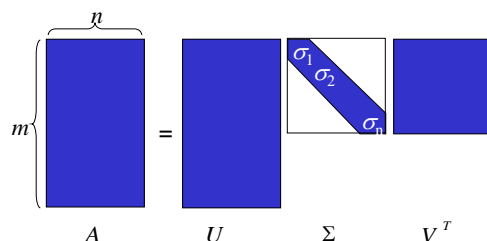
## SVD : Rank of a Matrix

22

- The *rank* of  $A$  is the number of non-zero singular values.
- If there are very small singular values  $\sigma_i$ , then  $A$  is close of being singular.

We can set a threshold  $t$ , and set  $\sigma_i = 0$  if  $\sigma_i \leq t$

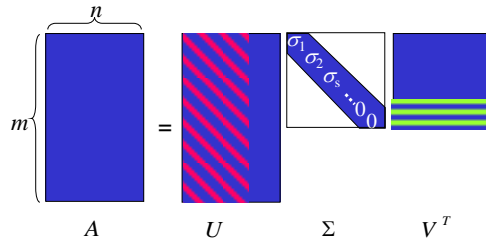
then the  $\text{numeric\_rank}(A) = \# \{ \sigma_i / \sigma_i > t \}$



## SVD : Rank of a Matrix (2)

23

- $\text{numeric\_rank}(A) = \# \{ \sigma_i / \sigma_i > t \}$ ,  
the rank of  $A$  is equal the  $\dim(\text{Img}(A))$



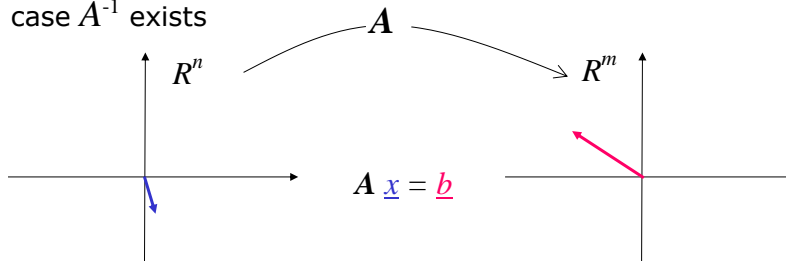
$$n = \dim(\text{Img}(A)) + \dim(\text{Ker}(A))$$

- the columns of  $U$  corresponding to the  $\sigma_i \neq 0$ , span the range of  $A$
- the columns of  $V$  corresponding to the  $\sigma_i = 0$ , span the nullspace of  $A$

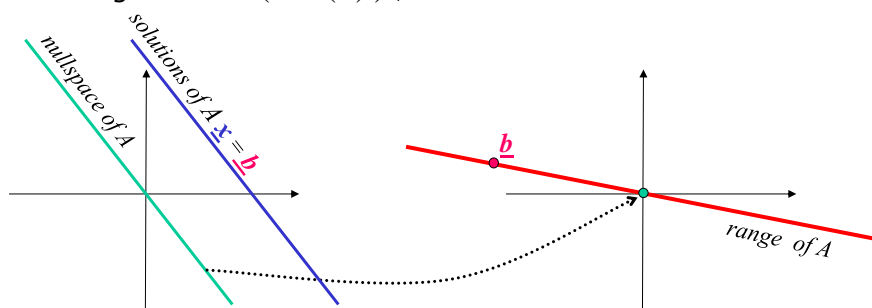
## remember linear mappings $A \underline{x} = \underline{b}$

24

1) case  $A^{-1}$  exists



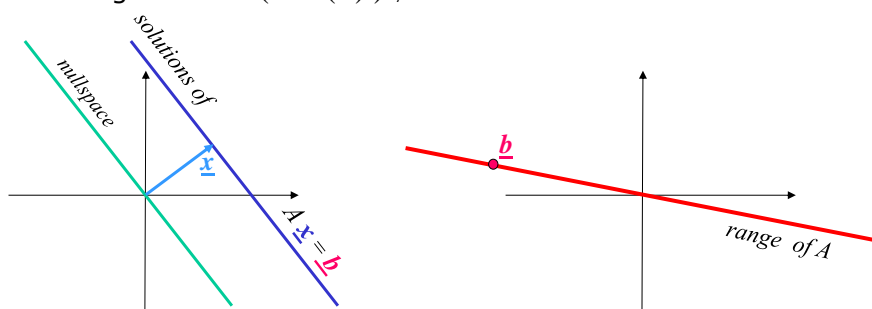
2)  $A$  is singular:  $\dim(\text{Ker}(A)) \neq 0$



## SVD : solving $A \underline{x} = \underline{b}$

25

2)  $A$  is singular:  $\dim(\text{Ker}(A)) \neq 0$



There are an infinite number of different  $\underline{x}$  that solve  $A\underline{x}=\underline{b}$  !!??

Which one should we choose??

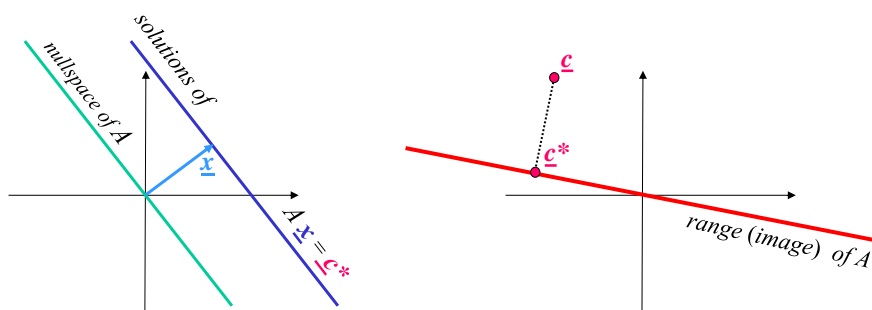
e.g. we can choose the  $\underline{x}$  with  $\|\underline{x}\| = \min$

→ then we have to search in the space *orthogonal to the nullspace*

## SVD : Solving $\|A \underline{x} - \underline{c}\| = \min$

26

3)  $\underline{c}$  is not in the range of  $A$



1) Projecting  $\underline{c}$  into the *range* of  $A$  results in  $\underline{c}^*$

2) From all the solutions of  $A \underline{x} = \underline{c}^*$  we choose the  $\underline{x}$  with  $\|\underline{x}\| = \min$

## SVD : Solving $\|A \underline{x} - \underline{c}\| = \min$

$$A \underline{x} = \underline{c} \quad \text{for any } A \text{ exist } U, \Sigma, V, \text{ with } A = U \Sigma V^T$$

with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$

$$U \Sigma V^T \underline{x} = \underline{c}$$

$$\underline{x} = (U \Sigma V^T)^{-1} \underline{c}$$

$$= (V^T)^{-1} \Sigma^{-1} U^{-1} \underline{c}$$

$$= V \begin{bmatrix} \frac{1}{\sigma_1} & & & \\ & \ddots & & \\ & & \frac{1}{\sigma_n} & \\ & & & \ddots \end{bmatrix} U^T \underline{c}$$

Computing  $A^{-1}$  for a singular  $A$  !?

--> What to do in  $\Sigma^{-1}$  with  $1/0 = \infty$  ????

Some  $\sigma_i = 0$  if  $\sigma_i \leq t$

Remember what we need ---- >


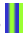
## SVD : Solving: $\|A \underline{x} - \underline{c}\| = \min$

We need to:

- 1) Project  $\underline{c}$  into the *range* of  $A$  to obtain a  $\underline{c}^*$
- 2) From all the solutions of  $A \underline{x} = \underline{c}^*$  we choose the  $\|\underline{x}\| = \min$  that is the  $\underline{x}$  in the space *orthogonal* to the *nullspace*

$$\underline{x} = \begin{bmatrix} \text{blue block} & \text{green block} \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & 0^{-1} \\ & & & & 0^{-1} \end{bmatrix} \begin{bmatrix} \text{red/blue diagonal} \\ \text{blue block} \end{bmatrix} \underline{c}$$

$V \quad \Sigma^{-1} \quad U^T$

- the columns  of  $U$  corresponding to the  $\sigma_i \neq 0$ , span the range of  $A$
- the columns  of  $V$  corresponding to the  $\sigma_i = 0$ , span the nullspace of  $A$

Basically all rows or columns multiplied by  $1/0$  are irrelevant!!

--> so even setting  $1/0 = 0$ , will lead to the correct result.

## SVD at Work:

29

For Linear Systems  $\mathbf{A} \underline{x} = \underline{b}$  :

Case fewer equations than unknowns:

→ fill rows of  $\mathbf{A}$  with zeros so that  $n = m$

Perform SVD on  $\mathbf{A}$  with  $(n \leq m)$ :

→ Compute  $U, \Sigma, V$ , with  $\mathbf{A} = U \Sigma V^T$

→ Compute threshold  $t$  and

→ in  $\Sigma$  set  $\sigma_i = 0$  for all  $\sigma_i \leq t$

→ in  $\Sigma^{-1}$  set  $1/\sigma_i = 0$  for all  $\sigma_i \leq t$

For Linear Systems: compute Pseudoinverse  $\mathbf{A}^+ = V \Sigma^{-1} U^T$   
and compute  $\underline{x} = \mathbf{A}^+ \underline{b}$

## Application: Compute PCA via SVD

30

**Problem:** Given mean free data  $X$ , a set of  $n$  feature vectors  $\underline{x}_i \in \mathbb{R}^m$  compute the orthonormal eigenvectors  $a_i$  of the correlation matrix  $R_x$ .

Now we use SVD

1. Move center of mass to origin:  $x_i' = x_i - \mu$
2. Build data matrix, from mean free data  $X = U \Sigma V^T$
3. The principal axes are eigenvector of the covariance matrix  $C = 1/n X X^T$


$$X X^T = U \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix} U^T$$


## Application: Compute PCA via SVD (2)<sup>31</sup>

with SVD

$$\begin{aligned} XX^T &= U \Sigma V^T (U \Sigma V^T)^T \\ &= U \Sigma V^T (V \Sigma^T U^T) \\ &= U \Sigma \Sigma^T U^T \\ &= U \Sigma^2 U^T \end{aligned}$$

Since  $C = 1/n XX^T$

the eigenvalues compute to  $\lambda_i = 1/n \sigma_i^2$    $\sigma$  from SVD

with  $\lambda_i = \sigma_i^2$    $\sigma^2$  variance of  $E[y_i^2]$

## Example: PCA on Images

32

- Assume we have a set of  $k$  images (of size  $N \times N$ )
- Each image can be seen as  $N^2$ -dimensional point  $\mathbf{p}_i$  (lexicographically ordered); the whole set can be stored as matrix:

$$X = \begin{bmatrix} | & | & & | \\ \mathbf{p}_1 & \mathbf{p}_2 & \cdots & \mathbf{p}_k \\ | & | & & | \end{bmatrix}$$

- Computing PCA the "naïve" way
    - Build correlation matrix  $XX^T$  ( $N^4$  elements)
    - Compute eigenvectors from this matrix:  $O((N^2)^3)$
- Already for small images (e.g.  $N=100$ ) this is far too expensive

# PCA on Images

33

Now we use SVD

1. Move center of mass to origin:  $p'_i = p_i - \mu$

2. Build data matrix, from mean free data

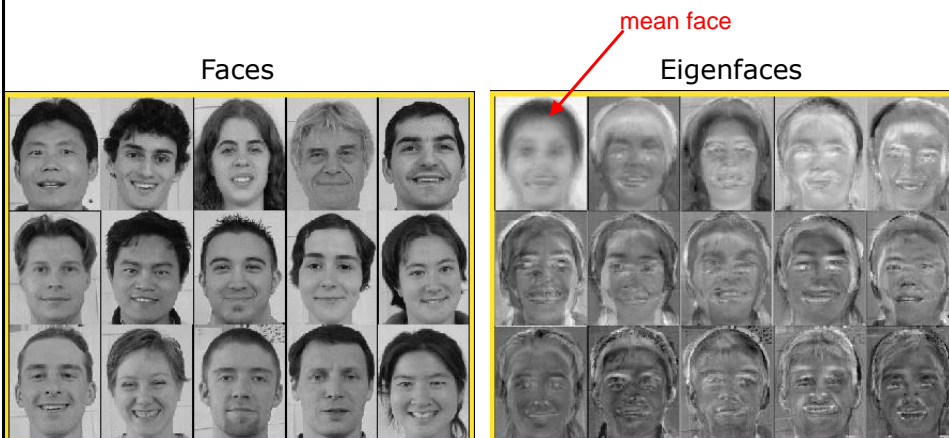
$$X = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{p}'_1 & \mathbf{p}'_2 & \cdots & \mathbf{p}'_n \\ | & | & & | \end{bmatrix}$$

3. The principal axes are eigenvector of

$$XX^T = U \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_d & \end{bmatrix} U^T$$

# PCA on Images

34



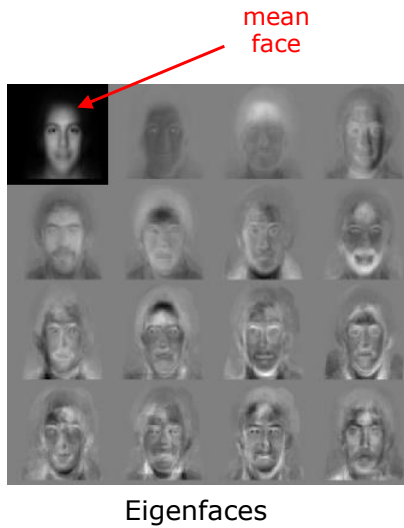
Principal Components can be visualized by adding to the mean vector an eigenvector multiplied by a factor (e.g.  $\lambda$ )



# PCA applied to face images

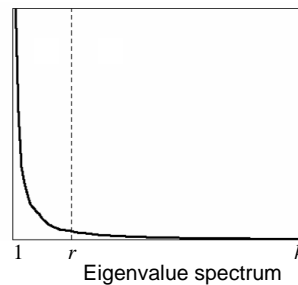
35

Here the faces were normalized in eye distance and eye position.



Choosing subspace dimension  $r$ :

- Look at decay of the eigenvalues as a function of  $r$
- Larger  $r$  means lower expected error in the subspace data approximation

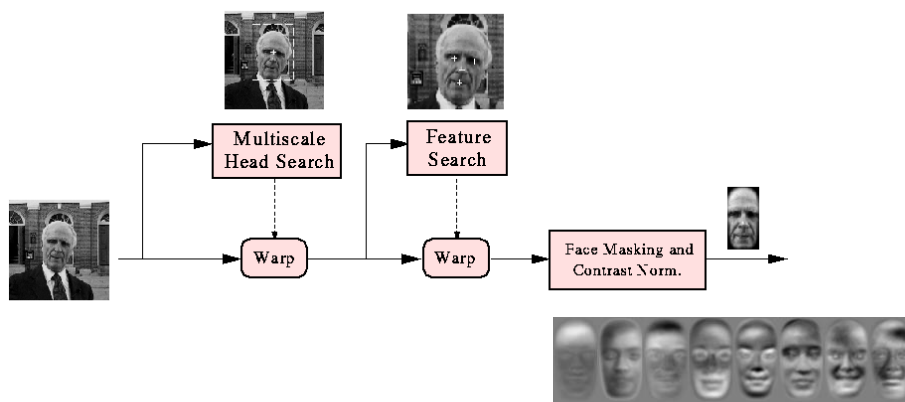


# Eigenfaces for Face Recognition

36

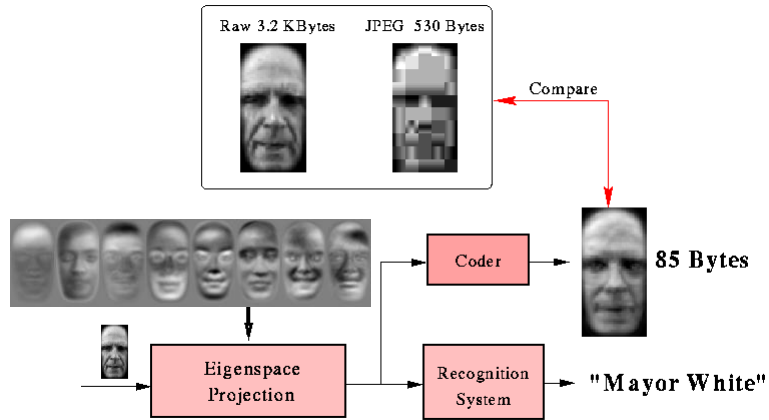
In the 90's the best performing Face Recognition System!

Turk, M. and Pentland, A. (1991).  
Face recognition using eigenfaces.  
In Proceedings of Computer Vision and Pattern Recognition, pages 586--591. IEEE.



# PCA for Face Recognition

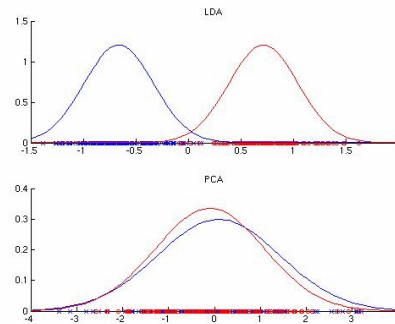
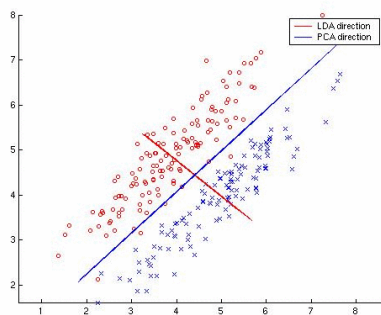
37



# PCA & Discrimination

38

- PCA/KLT do not use any class labels in the construction of the transform.
  - The resulting features may obscure the existence of separate groups.



## PCA Summary

39

- Unsupervised: no assumption about the existence or nature of groupings within the data.
- PCA is similar to learning a Gaussian distribution for the data.
- Optimal basis for compression (if measured via MSE).
- As far as dimensionality reduction is concerned this process is distribution-free, i.e. it's a mathematical method without underlying statistical model.
- Extracted features (PCs) often lack 'intuition'.

## PCA an Neural Networks

40

A three-layer NN with linear hidden units, trained as auto-encoder, develops an internal representation that corresponds to the principal components of the full data set. The transformation  $F_1$  is a linear projection onto a  $k$ -dimensional (Duda, Hart and Stork: chapter 10.13.1).

