



Machine Learning
– winter term 2016/17 –

Chapter 06: Dimensionality Reduction

(and Anomaly Detection)

Prof. Adrian Ulges
Masters "Computer Science"
DCSM Department
University of Applied Sciences RheinMain

1

Outline



1. Dimensionality Reduction: Motivation
2. Feature Selection
3. Principal Component Analysis
4. PCA Example: "Eigenfaces"
5. Anomaly Detection

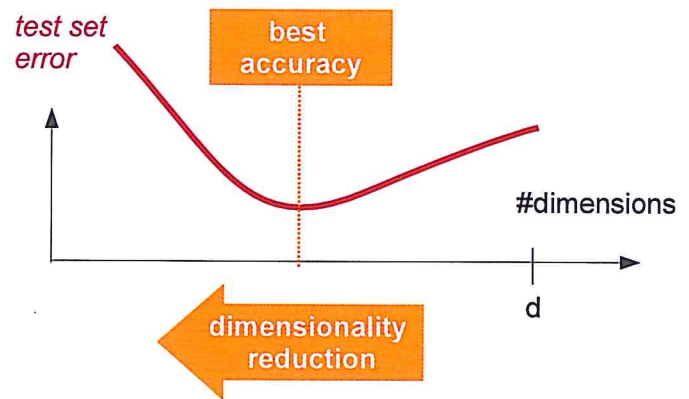
2

Dimensionality Reduction: Motivation



The Challenge

- ▶ **Reminder:** In machine learning, we are usually given **d-dimensional** samples $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- ▶ The number of features d can be high!
- ▶ **Goal: Reduce d** while preserving (or even improving) discriminativity



Why?

- ▶ **Efficiency:** faster training, faster application, less storage
- ▶ Avoid **overfitting:** overcoming the curse of dimensionality
- ▶ Better **interpretability** (and maybe visualization) of data

3

Dimensionality Reduction: Overview

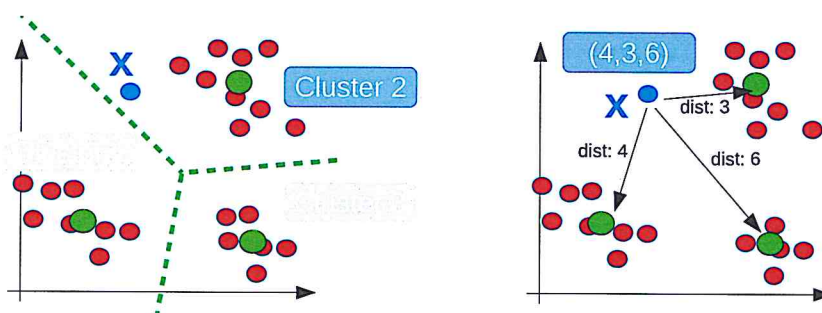


Approaches

- ▶ **feature selection**
- ▶ **feature derivation**, typically, by applying transformations

Example: K-Means for Feature Derivation

- ▶ **Approach:** map samples \mathbf{x} to clusters (*vector quantization*)
- ▶ **Alternative 1:** store only the cluster number (*1 dimension!*)
- ▶ **Alternative 2:** store the distances to all centers (*K dimensions*) (see `sklearn > KMeans > transform()`)



4



1. Dimensionality Reduction: Motivation
2. Feature Selection
3. Principal Component Analysis
4. PCA Example: "Eigenfaces"
5. Anomaly Detection

5

Feature Selection: Strategies



- ▶ **Goal:** Find a **subset of features** $\mathcal{F} \subseteq \{1, \dots, d\}$ for which we can solve our machine learning problem 'best' (*for example: minimizing classification error*)
- ▶ This is a **search problem** (*brute-force effort: $O(2^d)$*)
- ▶ There are **three common approaches: wrappers, filters, and embedded methods**

1. Wrappers

- ▶ Wrappers use an explicit evaluation of feature subsets (*training and validating classifiers*)
- ▶ Search can be done in a **greedy** fashion (adding or removing the 'best' features to the feature set), or by **backtracking**
- ▶ Benefit: It takes the **underlying classifier** into account!
- ▶ Usually the most reliable way, but very slow

6

Feature Selection: Strategies

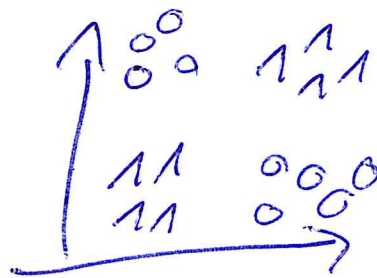


2. Filters

- ▶ Assess feature quality by a **proxy measure**
 - ▶ example: *mutual information* between feature X and class labels C

$$I(X, C) = \sum_{x \in X} \sum_{c \in C} p(x, c) \cdot \log_2 \left(\frac{P(x, c)}{P(x) \cdot P(c)} \right)$$

- ▶ **Search strategy:** rank features by their quality, pick the K top ones (K determined via cross-validation)
- ▶ Filters are cheaper than wrappers, but not as accurate



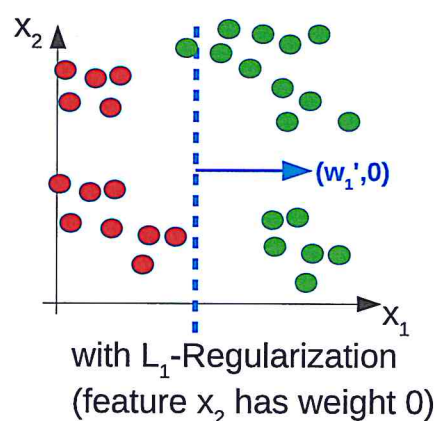
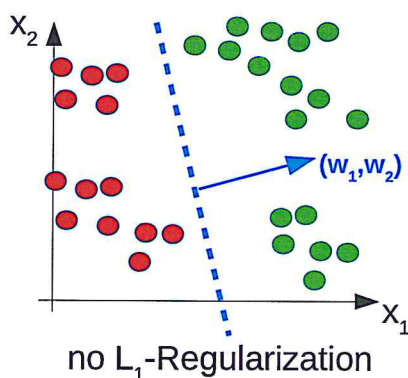
7

Feature Selection: Strategies



3. Embedded Methods

- ▶ ... treat feature selection as a part of **model construction** (i.e., we find classifier and features in the *same process*)
- ▶ Optimization is driven towards models with **few features**
- ▶ **Example:** logistic regression with **regularization**
 - ▶ the classifier penalizes feature weights, shrinking them to zero
 - ▶ features with weight zero are “filtered”!



8



1. Dimensionality Reduction: Motivation
2. Feature Selection
3. Principal Component Analysis
4. PCA Example: "Eigenfaces"
5. Anomaly Detection

9

Principal Component Analysis (PCA)



Idea (unsupervised!)

- ▶ The dataset exhibits a large stretch along some directions. These are the *important* directions.
- ▶ Along other directions, there is only little variation. These directions are merely *noise* and can be discarded.
- ▶ PCA is about finding the *important* directions (or **principal components**) of the data

Principal Components: Formalization

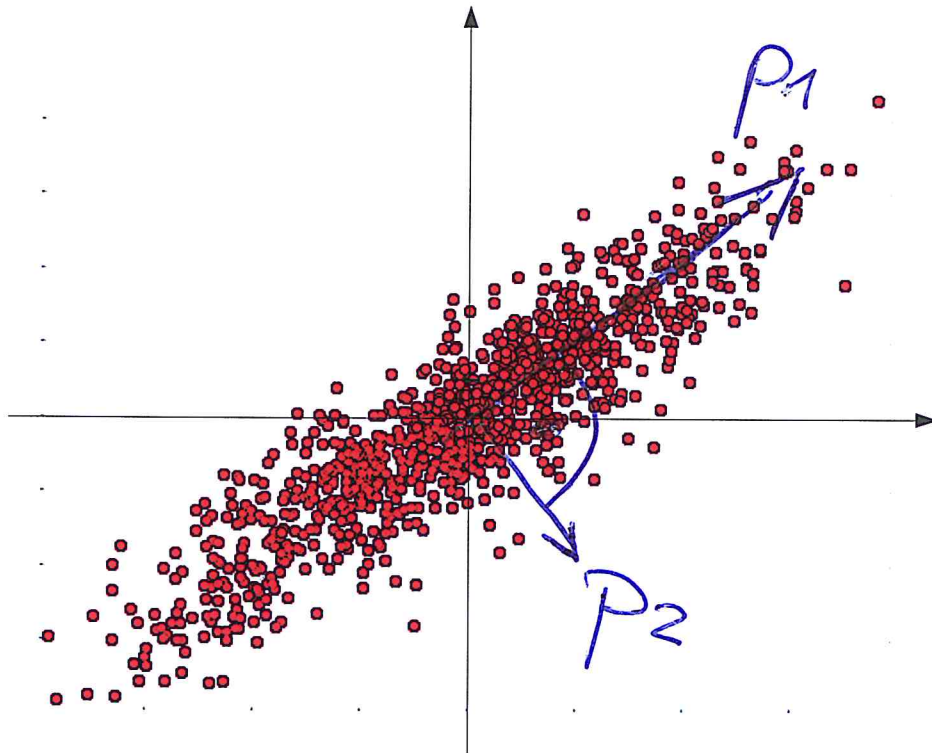
- ▶ In the following, we assume our data points to be centered around the origin (otherwise, we simply shift the data beforehand)
- ▶ What is the *most important direction*? What is the *second most important direction*, etc.?

10

PCA: Illustration



What are the Principal Components here?



11

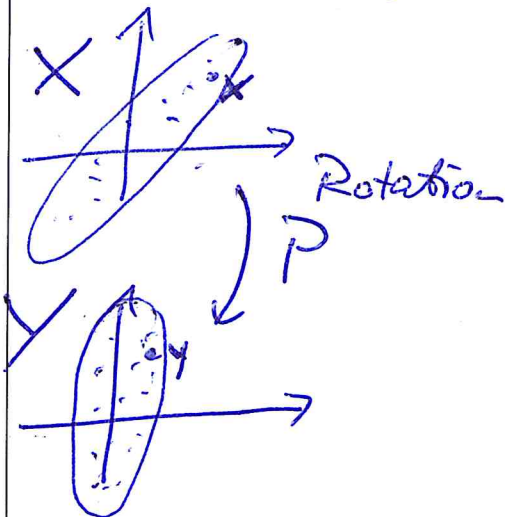
Principal Components: Derivation



Obviously, the principal components seem to be related to the data's **covariance matrix**. Let's find out how¹

We assume: $\bar{X} = \vec{0}$

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} x_{11} \dots x_{1d} \\ \vdots \\ x_{n1} \dots x_{nd} \end{pmatrix}$$



Goal: Find rotation matrix $P \in \mathbb{R}^{d \times d}$

$$Y^T = P \cdot X^T$$

such that $\text{cov}(Y) = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_d \end{pmatrix}$

¹cmp. Marsland, page 228f

12

Principal Components: Derivation

$$\text{cov}(Y) = \frac{1}{n} Y^T Y$$

$$Y_{ij} = \frac{1}{n} \sum x_{sj} \cdot x_{sj}^*$$

$$= \frac{1}{n} (P \cdot X^T) (P \cdot X^T)^T$$

$$(A \cdot B)^T = B^T \cdot A^T$$

$$= \frac{1}{n} P \cdot X^T X \cdot P^T$$

$$= P \cdot \underbrace{\frac{1}{n} X^T X}_{\text{cov}(X)} \cdot P^T$$

$$\text{cov}(Y) = P \cdot \text{cov}(X) \cdot P^T \quad / P^T$$

$$P^T \cdot \text{cov}(Y) = P^T P \cdot \text{cov}(X) \cdot P^T$$

$$\begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_d \end{pmatrix} \cdot \begin{pmatrix} \\ \\ \end{pmatrix} \cdot \begin{pmatrix} \\ \\ \end{pmatrix}$$

13

Principal Components: Derivation

*

Let p_j be a column of P^T .
(= a row of P)

$$\lambda_j \cdot p_j = \text{cov}(X) \cdot p_j$$

$$\lambda \cdot x = A \cdot x$$

The rows p_1, \dots, p_d are the eigenvectors of the covariance matrix $\text{cov}(X)$.

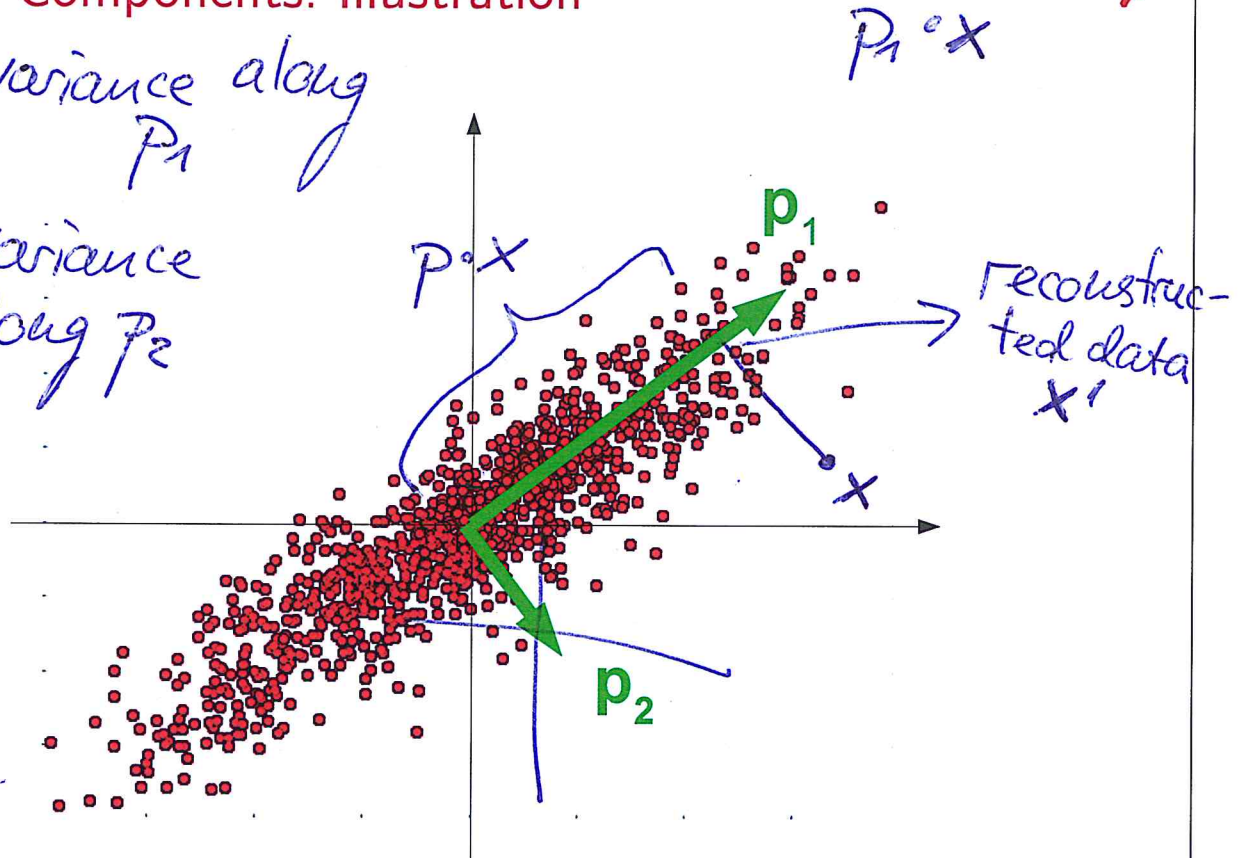
14

Principal Components: Illustration



$\lambda_1 = \text{variance along } P_1$

$\lambda_2 = \text{variance along } P_2$



15

PCA: Training



```
1 function PCA_TRAIN(x1, ..., xn, K) // given: samples x1, ..., xn ∈ ℝd
2   μ := 1/n ∑i xi
3   for i = 1, ..., n:
4     xi := xi - μ // shift the samples to mean zero
5   stack the samples into an n × d data matrix X
6   Σ := 1/n (XT · X) // covariance matrix
7   compute Σ's eigenvalues λ1, ..., λd and eigenvectors p1, ..., pd
8   (sorted in descending order of the eigenvalues)
9   stack p1, ..., pK as rows into a K × d-Matrix PK
10  return PK, μ
11
```

Remarks

- ▶ **Training:** Compute the top K eigenvectors of the data's covariance matrix
- ▶ **Application:** reduces the samples from d to K dimensions

```
1 function PCA_APPLY(x) // given: a new sample x
2   return PK · (x - μ) // dimension of return value: K
3
```

16

PCA: Remarks



- ▶ Heuristic for choosing K : Preserve α (e.g., 90%) of the eigenvector's total energy

$$K := \min_{K \in \{1, \dots, d\}} \text{ such that } \sum_{k=1}^K \lambda_k \geq \alpha \cdot \sum_{k=1}^d \lambda_k$$

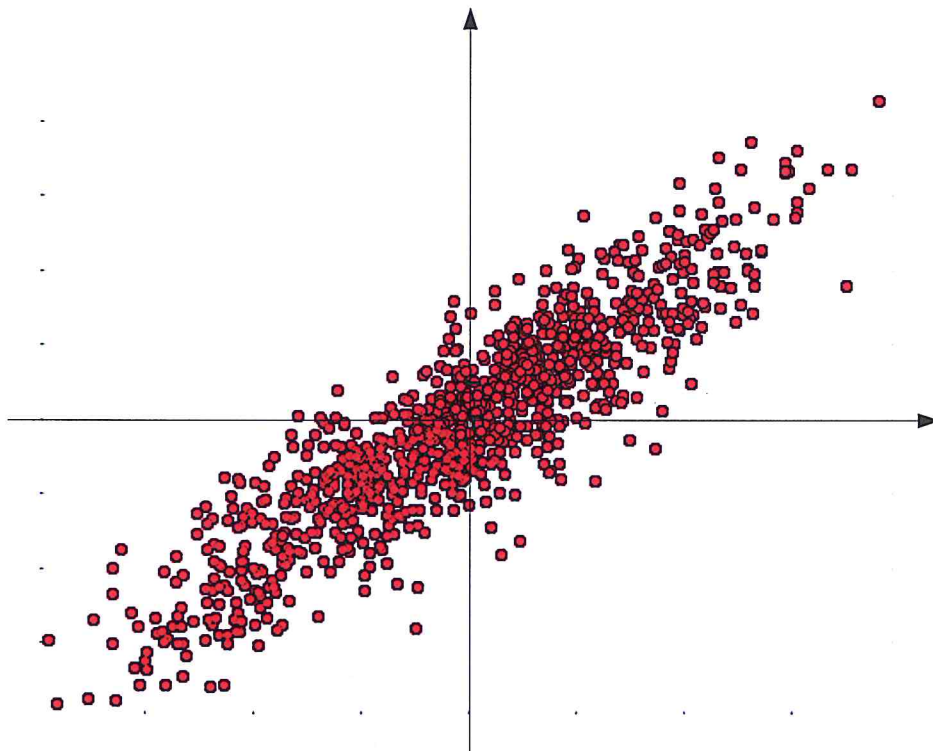
- ▶ Given a reduced K -dimensional feature $\mathbf{y} = (y_1, \dots, y_K)$, we can reconstruct \mathbf{x} (to some extent) from \mathbf{x}' :

$$\mathbf{x}' = \mu + y_1 \cdot \mathbf{p}_1 + y_2 \cdot \mathbf{p}_2 + \dots + y_K \cdot \mathbf{p}_K$$

- ▶ We call $\|\mathbf{x} - \mathbf{x}'\|$ the **reconstruction error**.

17

PCA: Illustration



18

From PCA to Whitening



We learned above: Given a feature vector \mathbf{x} (and the $d \times d$ matrix P with Σ 's eigenvectors as rows), by applying the transformation $P \cdot \mathbf{x}$, the **covariance matrix** of the data becomes

$$\begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & 0 & \lambda_d \end{pmatrix}$$

Each value λ_i is a variance within one dimension. We turn these variances into 1, simply by dividing by the standard deviation $\sqrt{\lambda_i}$. This leads to the **whitening** transform (see *Chapter 'Features'*):

$$\text{Diag}\left(\frac{1}{\sqrt{\lambda_1}}, \frac{1}{\sqrt{\lambda_2}}, \dots, \frac{1}{\sqrt{\lambda_d}}\right) \cdot P \cdot \mathbf{x}.$$

Applying this transformation turns the data's covariance matrix into the identity I and decorrelates the features.

19

Outline



1. Dimensionality Reduction: Motivation
2. Feature Selection
3. Principal Component Analysis
4. PCA Example: "Eigenfaces"
5. Anomaly Detection

20

Example: Applying PCA to Images



- ▶ We can apply PCA to **images**, simply by scaling images to a standard resolution (say, $N \times M$) and stacking all pixel values into sample vectors of dimension $d = N \times M$
- ▶ Note: The principal components are $(N \cdot M)$ -dimensional, too! (i.e., we can *visualize them as images*)
- ▶ Example: PCA for **face recognition** → “eigenfaces” (apply PCA to lots of face images)



21

Example: Applying PCA to Images



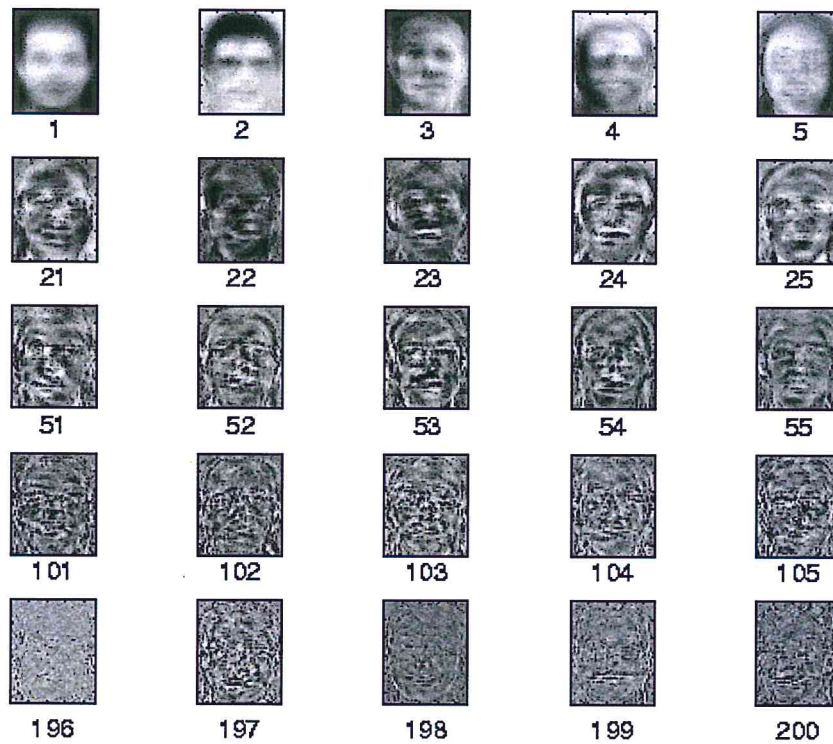
Mean face (top left) and the first 7 principal components (=eigenfaces)

22

Example: Applying PCA to Images



Some eigenfaces from a different face dataset

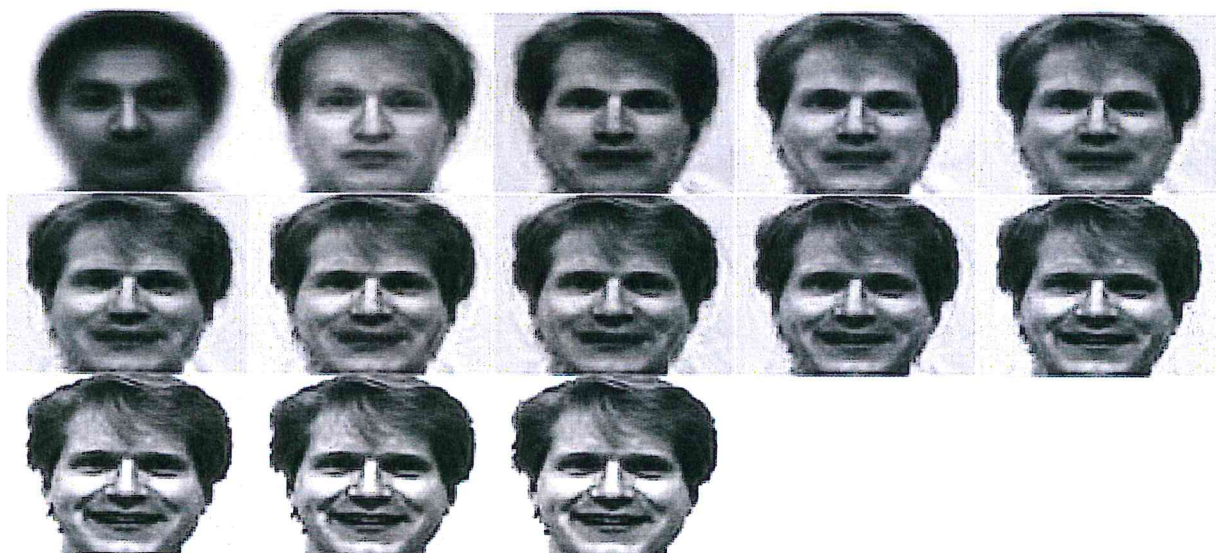


23

Example: Applying PCA to Images



Reconstruction of a face using 0, 8, 16, ... eigenfaces

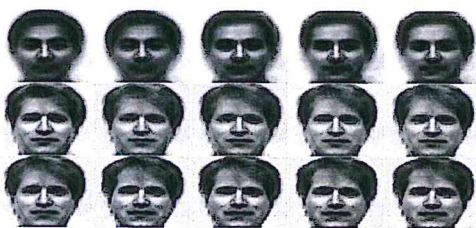


24

Example: Applying PCA to Images



- ▶ Which of these images does not show a human face?
- ▶ Approach: Compare the **original image** with its **PCA reconstruction**



faces are well reconstructed



non-faces are poorly reconstructed

25

Eigenfaces: Discussion



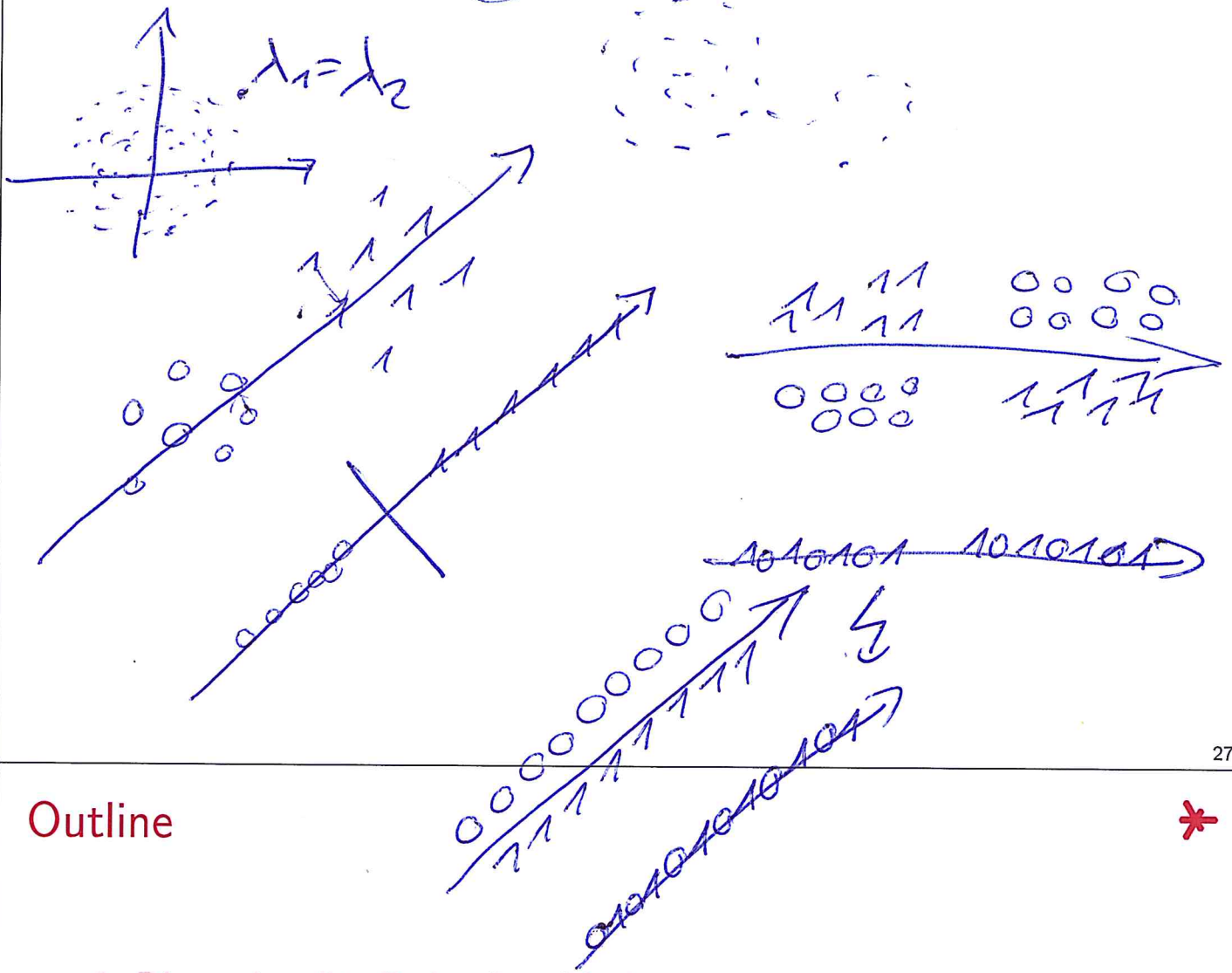
The Eigenfaces approach for **face detection** is **simple and powerful**, but it has some shortcomings:

- ▶ The results depend strongly on *illumination*, *shadows*, and *local changes*, e.g. glasses or beards
 - ▶ include those effects in the training set
 - ▶ learn smaller components (eye, mouth, ...)
 - ▶ (illumination): normalize all images
- ▶ Only faces of *fixed size* are detected
 - ▶ create scaled versions of the input image before searching
- ▶ Already small *rotations* of the head change the result
 - ▶ include rotated faces into the training set.
 - ▶ try to make the image upright before the applying PCA

Since eigenfaces [8], there have been at least two generations of more elaborate face detection methods [9, 7].

26

PCA: Discussion



27

Outline



1. Dimensionality Reduction: Motivation
2. Feature Selection
3. Principal Component Analysis
4. PCA Example: "Eigenfaces"
5. Anomaly Detection