

RWTH Aachen University
Chair of Computer Science VI
Prof. Dr.-Ing. Hermann Ney

Seminar Data Mining WS 2003/2004

Preprocessing and Visualization

Jonathan Diehl

January 19, 2004

Preprocessing and Visualization

- Introduction
- Theoretical Fundamentals
- Visualization
- Preprocessing

Literature

- **D. Hand, H. Manila, P. Smyth: Principles of Data Mining. MIT Press, Cambridge, MA, 2001, Chapters 2 and 3**
- **J. Han, M. Kamber: Data Mining: Concepts and Techniques. Academic Press, San Diego, CA, 2001, Chapter 3**
- **S. Roweis, L.Saul: Locally Linear Embedding Homepage.
<http://www.cs.toronto.edu/~roweis/lle/>**

Overview

Data Preprocessing:

- data manipulation prior to mining
- improves quality or speed of actual mining

Data Exploration:

- combined human and computer analysis
- utilizes human's natural abilities

Data Visualization:

- methods to display data to the human
- data is plotted graphically

Preprocessing and Visualization

- Introduction and Motivation
- Theoretical Fundamentals
- Visualization
- Preprocessing

Data Representation

- data consists of N samples (measures, etc.) and M attributes (variables, dimensions, etc.)
- for each sample n and each attribute m there exists one data value $x_{n,m}$

Data Representation

Data Set:

$$\mathbf{X} := \left(\mathbf{X}(1), \dots, \mathbf{X}(N) \right)^T \quad \text{with } \mathbf{X}(n) = (x_{n,1}, \dots, x_{n,M})$$

Data Matrix:

$$\mathbf{X} := \begin{pmatrix} x_{1,1} & \dots & x_{1,M} \\ \vdots & & \vdots \\ x_{N,1} & \dots & x_{N,M} \end{pmatrix}$$

with $x_{n,m}$ value of n -th sample and m -th attribute

Measures of Location

Sample Mean:

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N X(n) \quad \text{with } X \text{ set of } N \text{ data values}$$

Median:

50% of values below median

Mode:

data point which occurs most often

Measures of Dispersion

Variance:

$$\hat{\sigma}^2 = \frac{1}{N} \cdot \sum_{n=1}^N (X(n) - \mu)^2 \quad \text{with } X \text{ set of } N \text{ data values and mean } \mu$$

Quartiles:

25%/75% of values below quartile

Range:

difference between lowest and highest data value

Measures of Correlation

Covariance (of attributes A and B):

$$\text{Cov}(A, B) = \frac{1}{N} \cdot \sum_{n=1}^N (A(n) - \mu_A) \cdot (B(n) - \mu_B)$$

Covariance Matrix:

$$V = \frac{1}{N} \cdot X^T X \quad \text{with } X \text{ zero-centered data matrix}$$

Correlation Coefficient:

$$\rho_{A,B} = \frac{\text{Cov}(A, B)}{\sigma_A \cdot \sigma_B}$$

Modelling Data

Linear Regression (for two-dimensional data set):

$$Y = a + b \cdot X \quad \text{with } a, b \text{ coefficients and } X, Y \text{ attributes}$$

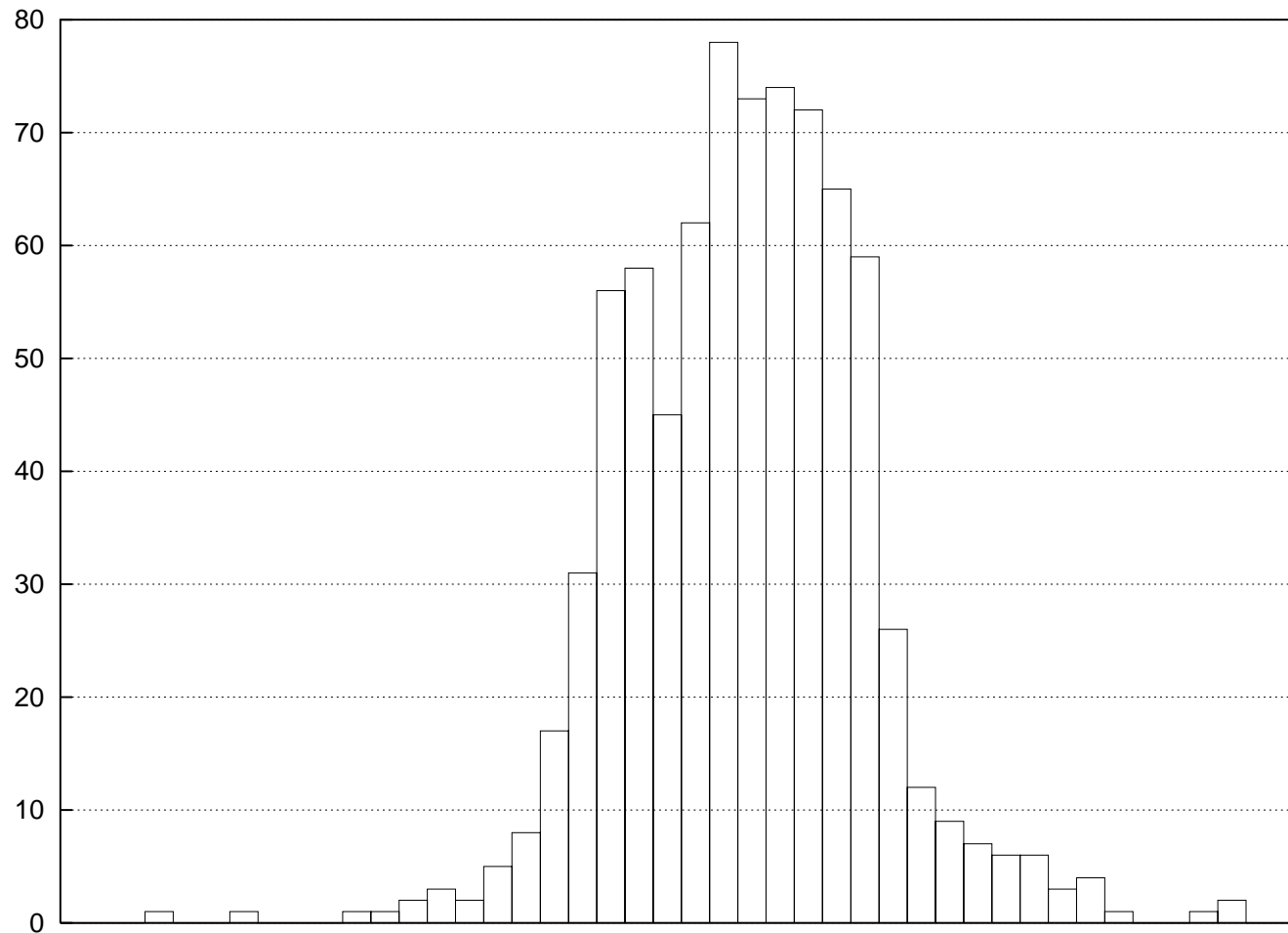
Multiple Linear Regression:

$$Y = a + \sum_{n=1}^N b_n \cdot X_n$$

Preprocessing and Visualization

- Introduction and Motivation
- Theoretical Fundamentals
- Visualization
- Preprocessing

Histogram



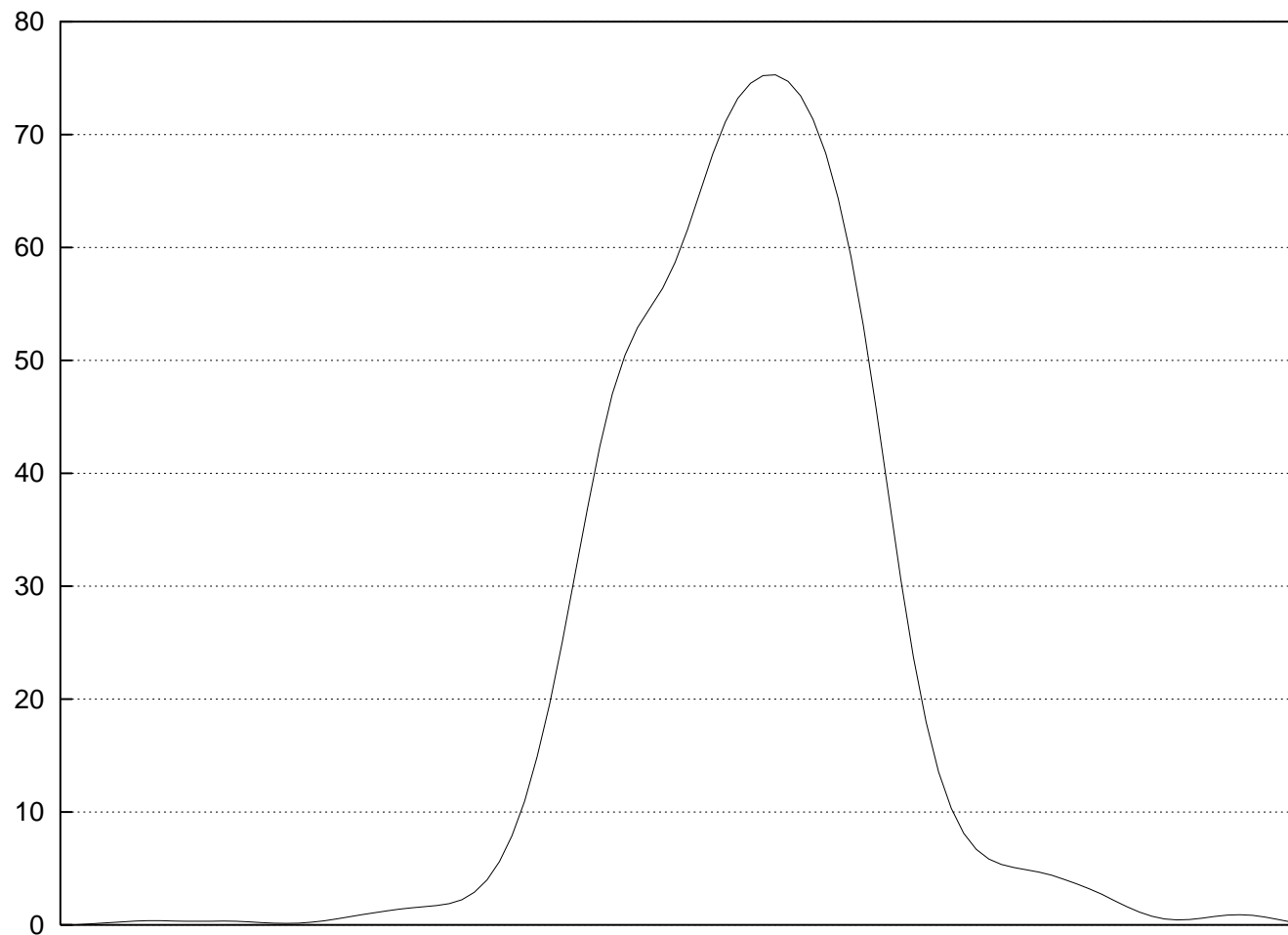
Kernel Method

Estimated density of kernel function K with bandwidth h :

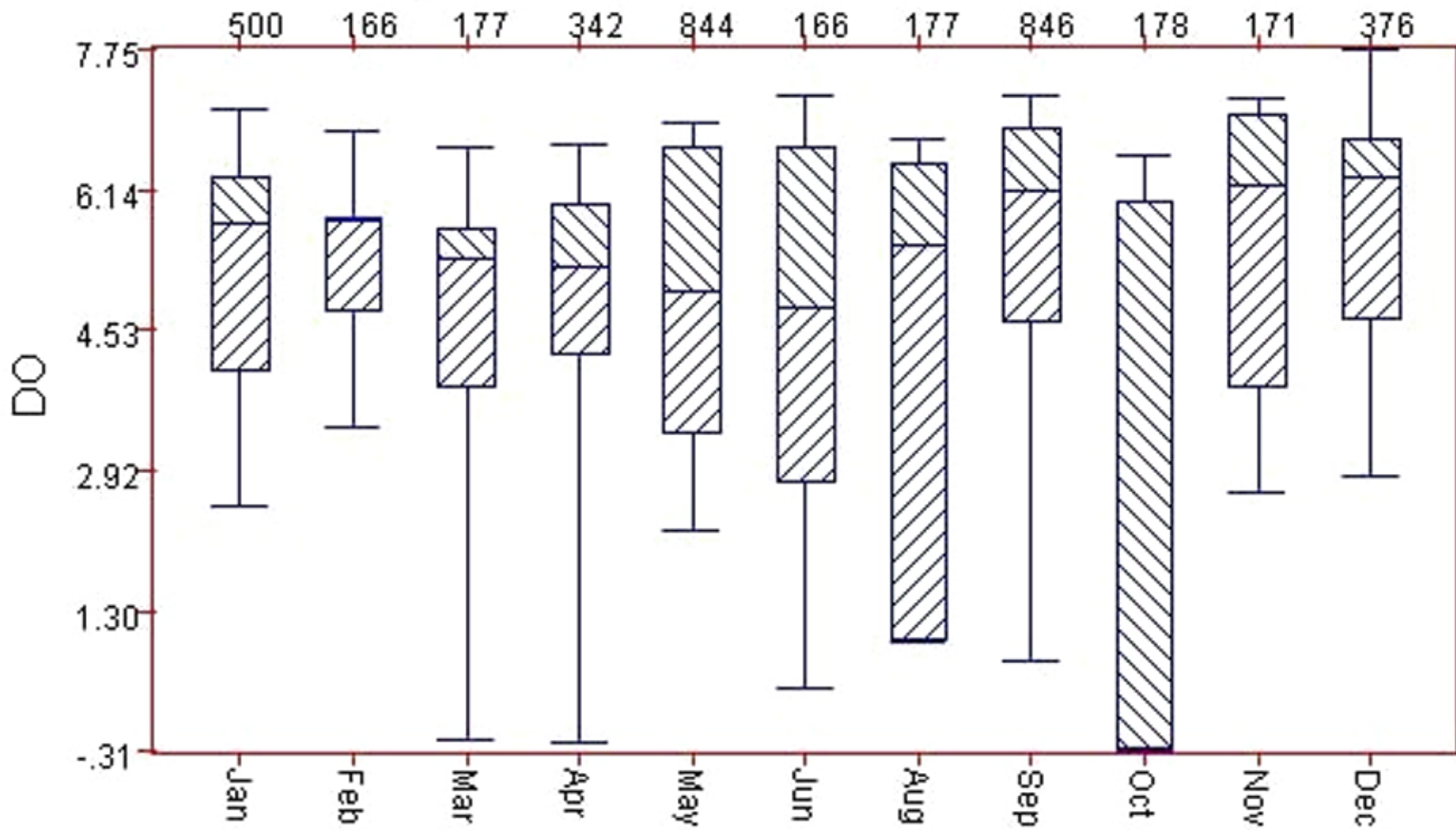
$$f(x) = \frac{1}{N} \sum_{n=1}^N K \left(\frac{x - X(n)}{h} \right) \quad \text{for given data set } X \text{ of } N \text{ values}$$

where $\int K(t)dt = 1$

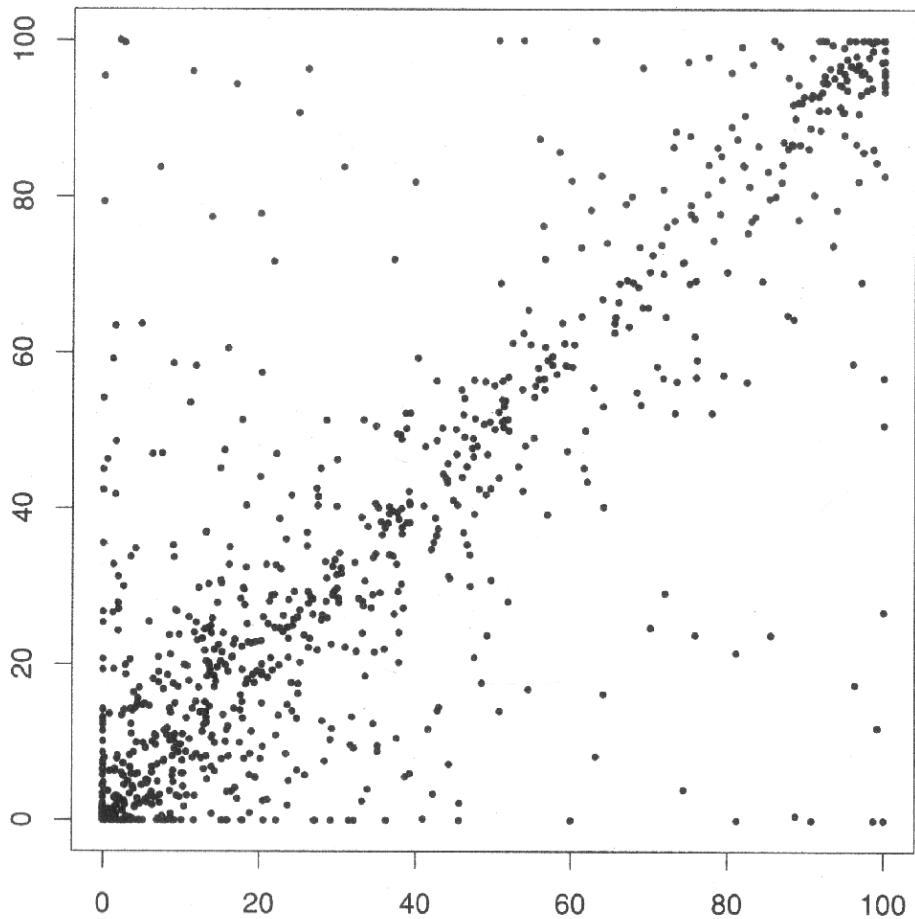
Kernel Method



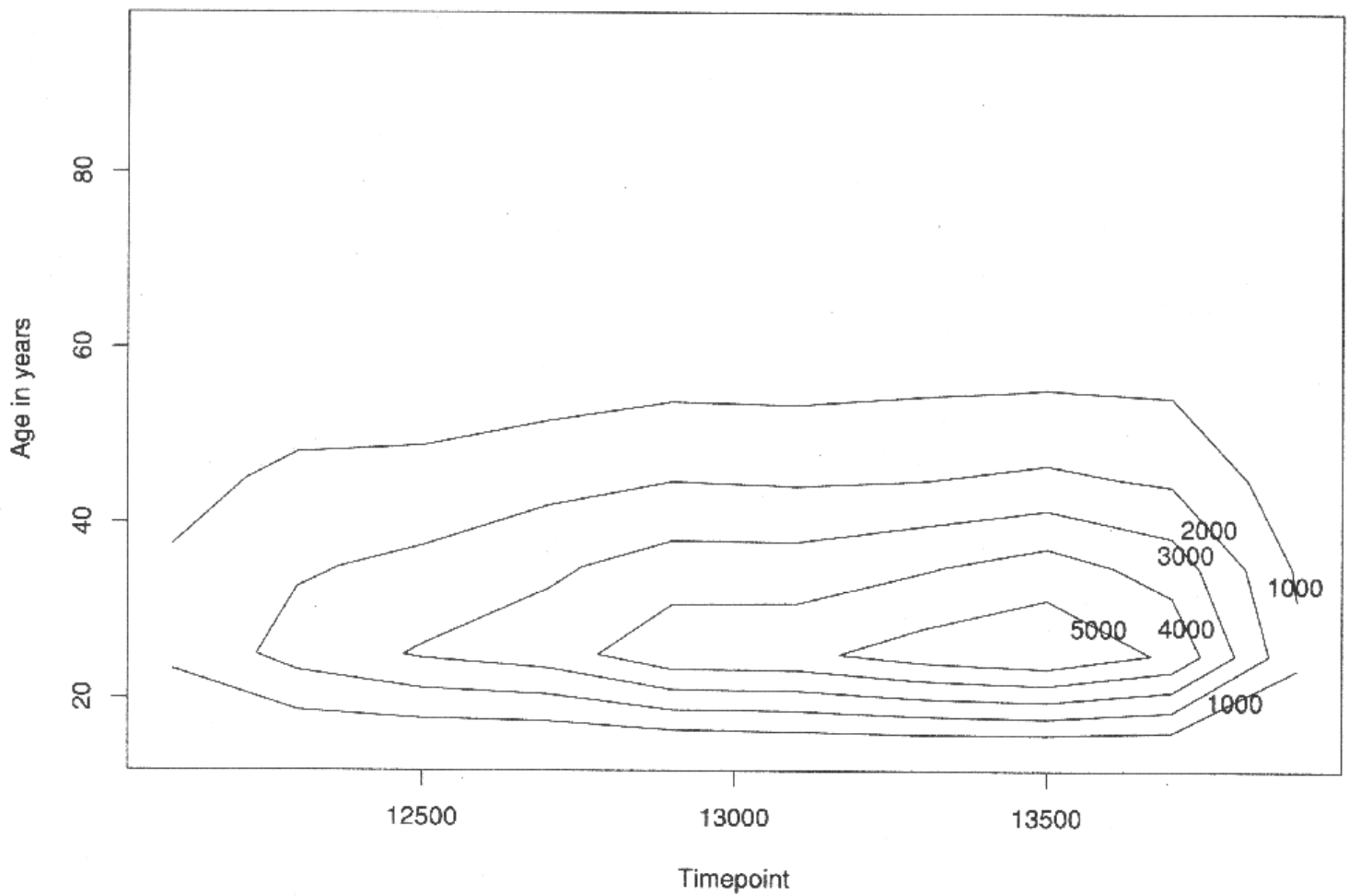
Boxplot



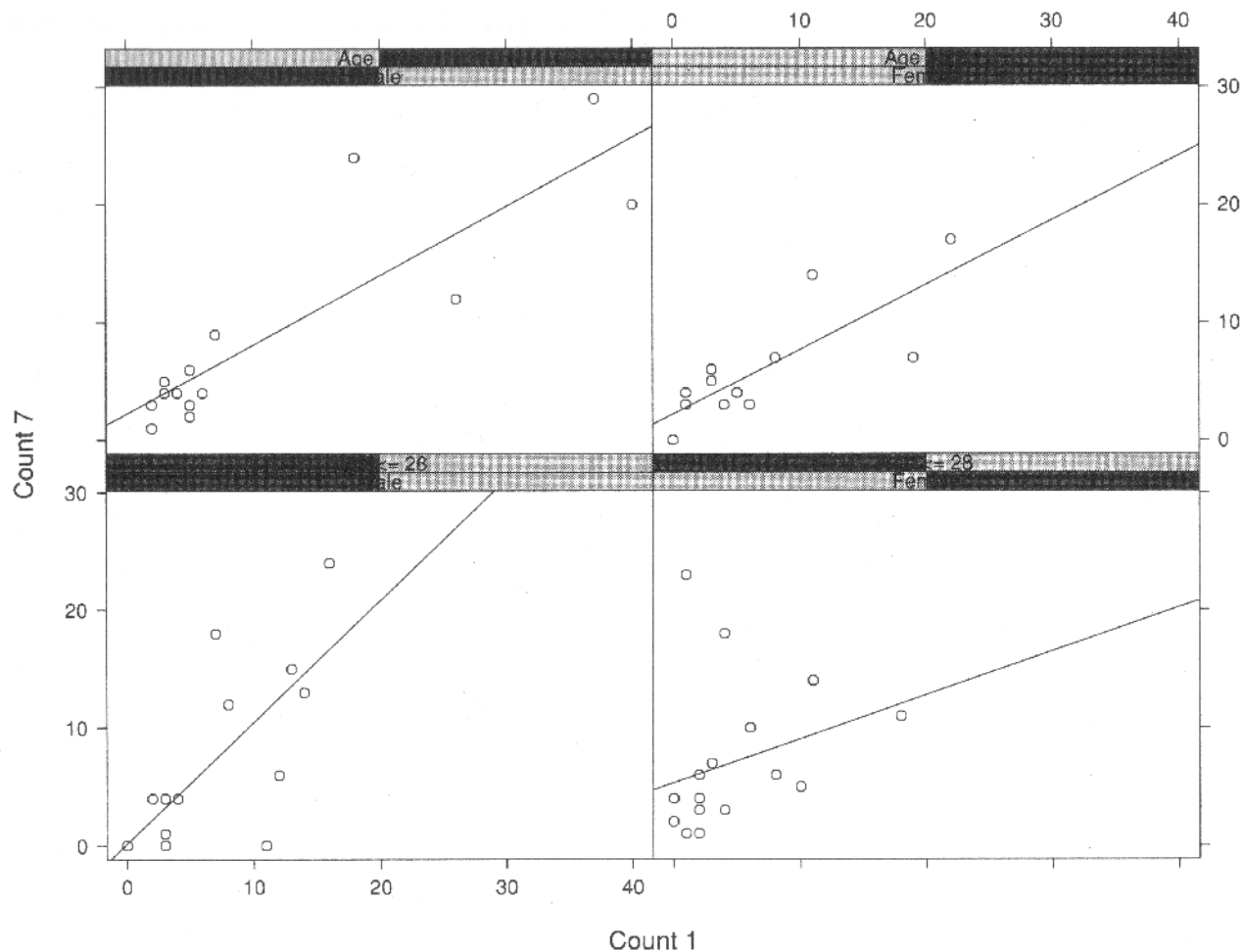
Scatterplot



Contour Plot



Trellis Plot



Preprocessing and Visualization

- Introduction and Motivation
- Theoretical Fundamentals
- Visualization
- Preprocessing

Data Cleaning

Problems:

- incomplete data (missing values)
- erroneous (noisy) data
- inconsistent data (→ data integration)

Missing Values

Solutions:

- ignore samples → lose important information
- determine most likely value:
 1. make use of statistical measures (mean/median, class mean/median)
 2. construct model of attribute relations (regression) and calculate value
 3. construct decision tree and derive value
- ...

Noisy Data: Regression

1. construct model of entire data set:

- linear Regression

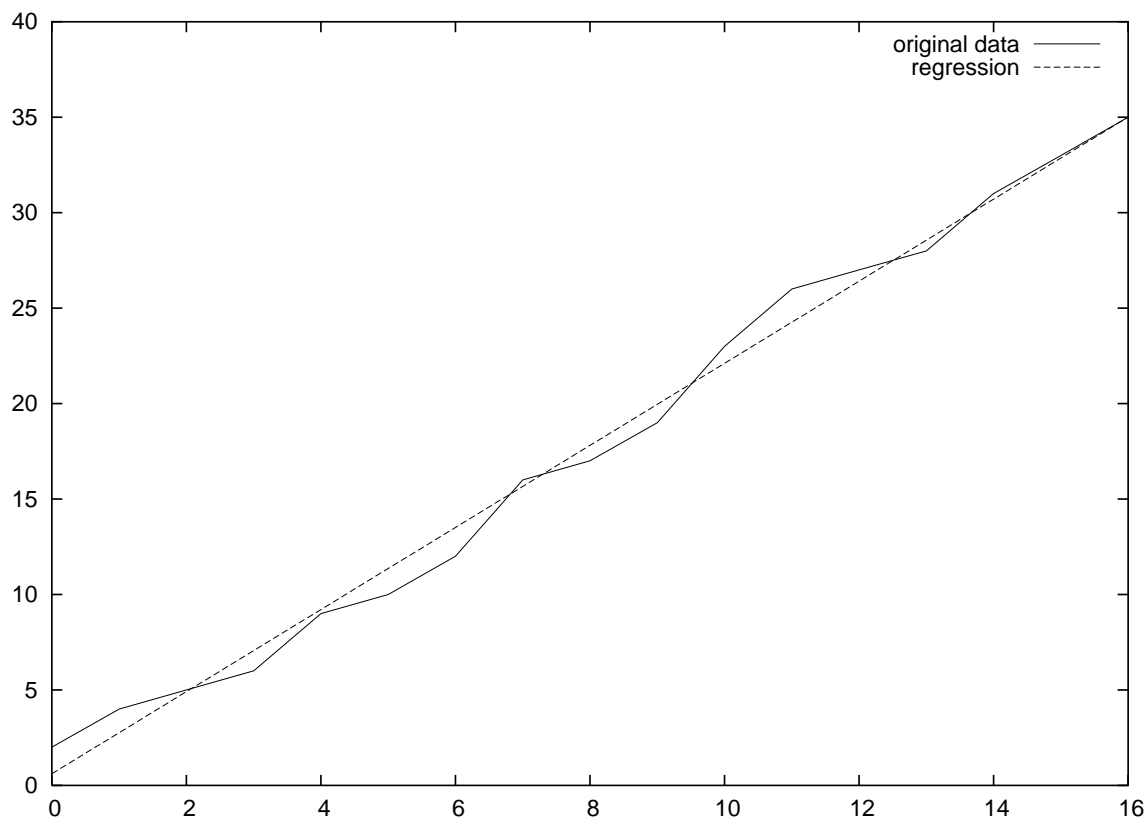
- multiple Regression

(regression over weighted linear combination of attributes)

2. calculate data points from regression equation

⇒ global smoothing, strong data reduction

Noisy Data: Regression



Data Integration

Problems:

- entity identification problem (→ metadata)
- data redundancy (→ correlation analysis)
- value conflicts (→ data transformation)

Data Transformation

Problem:

- data has to meet certain criteria before further processing
- e.g. attribute values must be weighted equally for many analysis methods

Normalization

Fit attribute A of data X into predefined range (e.g. $[0.0, 1.0]$)

- min-max normalization:

$$f(i) = \frac{A(i) - \min_A}{\max_A - \min_A} \cdot (\text{newmax}_A - \text{newmin}_A) + \text{newmin}_A$$

- z-score normalization:

$$g(i) = \frac{A(i) - \mu_A}{\sigma_A}$$

Attribute Construction

Construct new (redundant) attributes summarizing stored information

- sums, products of samples/data classes

⇒ Fast online access to summarized data

Data Reduction

Problem:

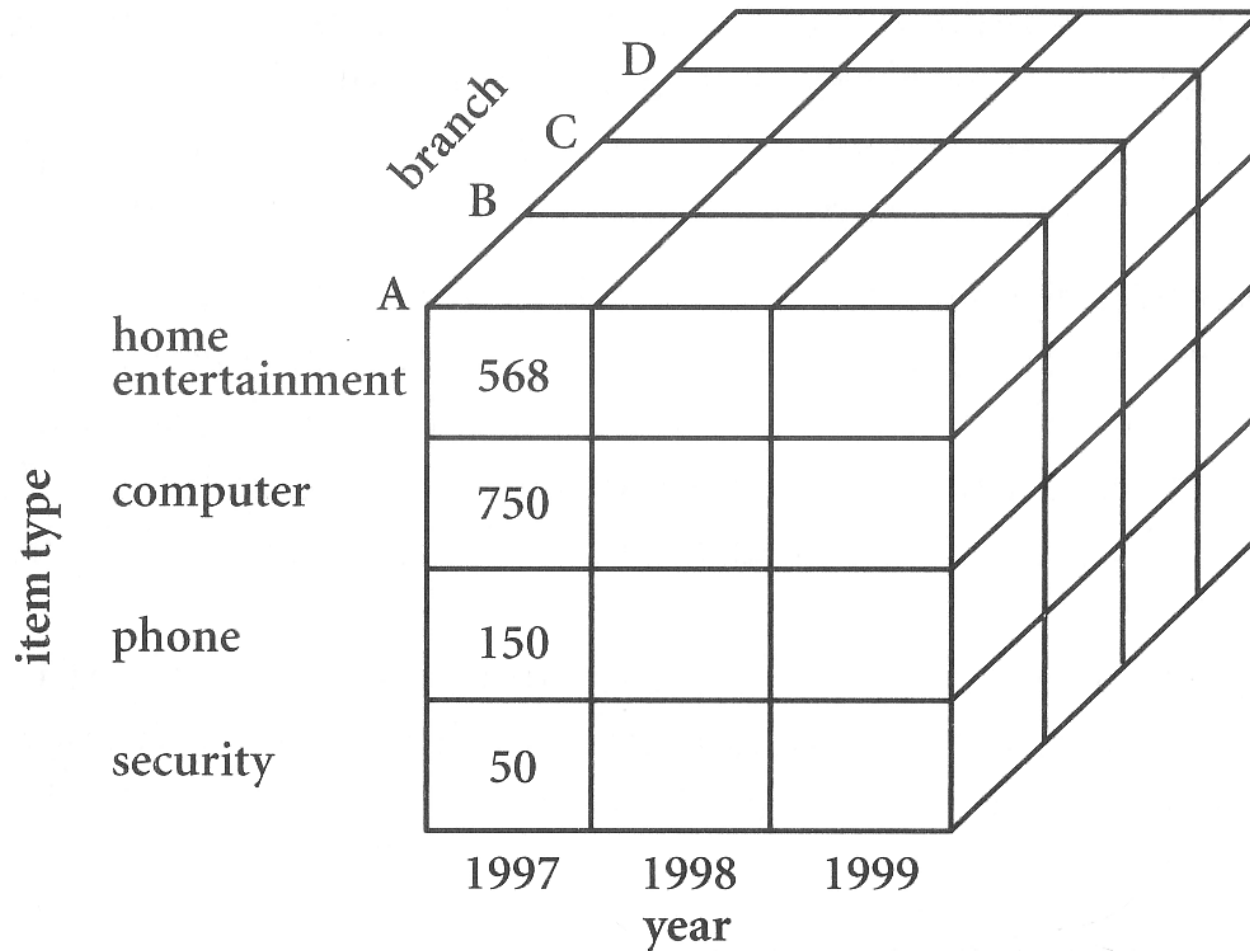
- huge databases with many attributes

⇒ very slow mining process

Data Cube Aggregation

- multidimensional arrangement of aggregated information
 - dimensions (axes) represent attributes
 - multiple levels of abstraction possible (concept hierarchy)
- ⇒ Efficient organization of summarized data

Data Cube Aggregation



Attribute Subset Selection

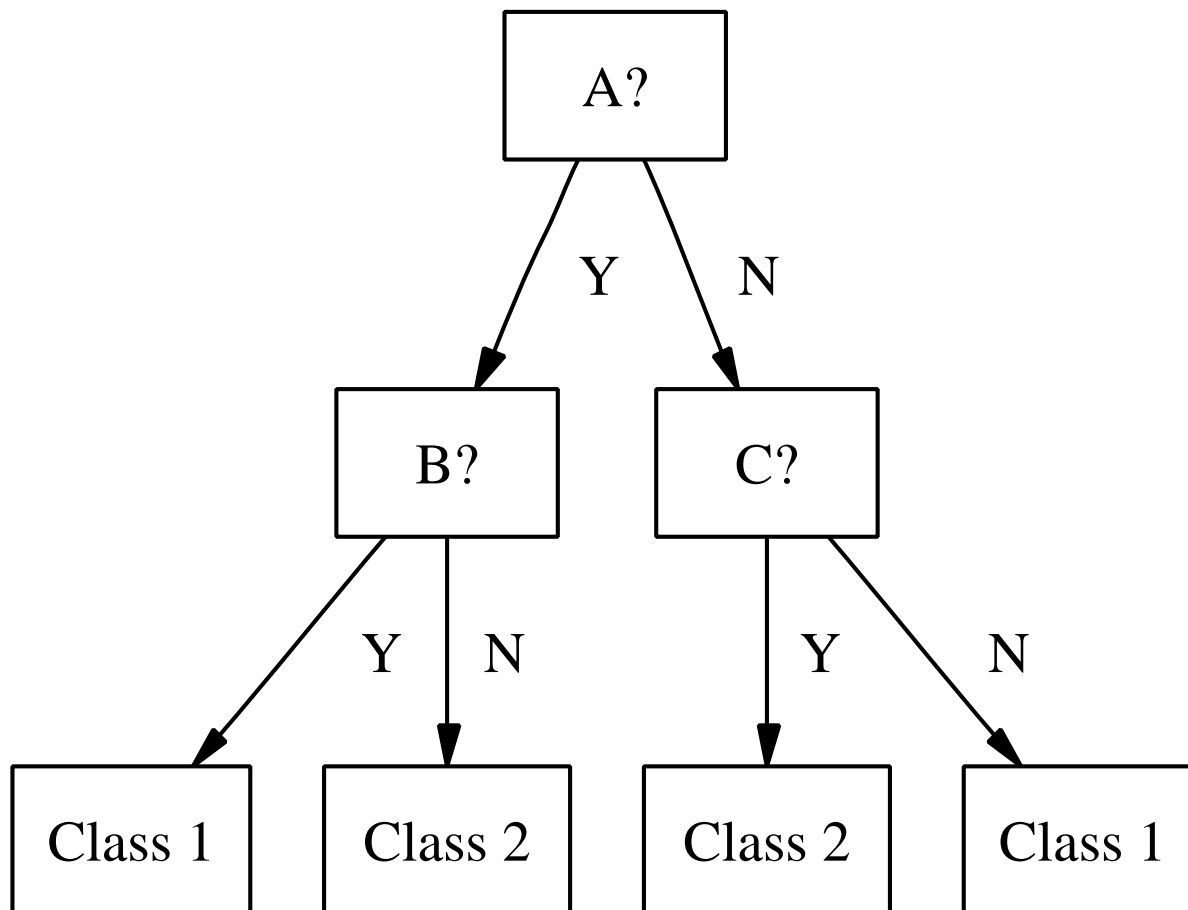
Determine significance of attributes:

- correlation coefficient
- information gain

Make selection:

- stepwise selection/elimination
- decision tree induction

Attribute Subset Selection



Principal Component Analysis

Aim: Maximize variability

- **project data onto principal components
(linear combination of attributes)**
- **normalize principal components**
- **zero-centered data set (subtract mean from all values)**

Principal Component Analysis

Variability of projection vector a ($S := X^T X$ scatter matrix):

$$\begin{aligned} s_a^2 &= (Xa)^T \cdot (Xa) \\ &= a^T X^T X a \\ &= a^T S a \end{aligned}$$

Impose normalization constraint and maximize:

$$\begin{aligned} u &= a^T S a - \lambda(a^T a - 1) \quad \rightarrow \max \\ \frac{\delta u}{\delta a} &= \underbrace{2Sa - 2\lambda a}_{\text{eigenvalue form}} = 0 \end{aligned}$$

⇒ eigenvector a with largest eigenvalue is first principal component

Principal Component Analysis

- data is projected onto the first K eigenvectors
- small values of K sufficient because late principal components are insignificant (eigenvalue approaches zero)
- for $K = 2$ principal component analysis can be used to project data into a plane for visualization

Multidimensional Scaling

Aim: Preserve Distances

- **fit multidimensional data into plane**
- **minimize squared distance error**
- **zero-centered data set (subtract mean from all values)**

Multidimensional Scaling

Given $B = XX^T$ solve (euclidian distance)

$$\begin{aligned}d_{ij}^2 &= \|x_i - x_j\|^2 \\ &= x_i^T x_i - 2x_i^T x_j + x_j^T x_j \\ &= b_{ii} + b_{jj} - 2b_{ij}\end{aligned}$$

then minimize

$$\sum_{i=1}^N \sum_{j=1}^N (\delta_{ij} - d_{ij})^2$$

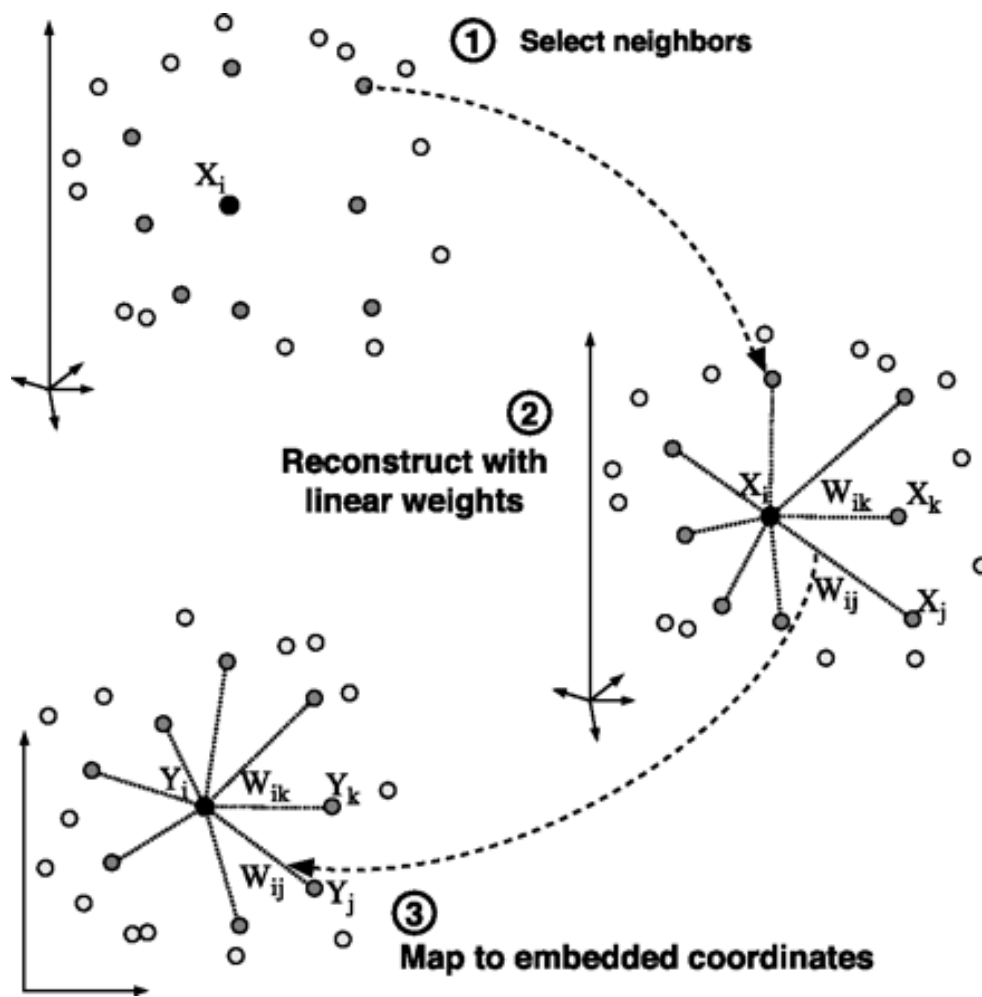
with δ_{ij} multidimensional distance between data point i and j

Locally Linear Embedding

Aim: Preserve Neighborhoods

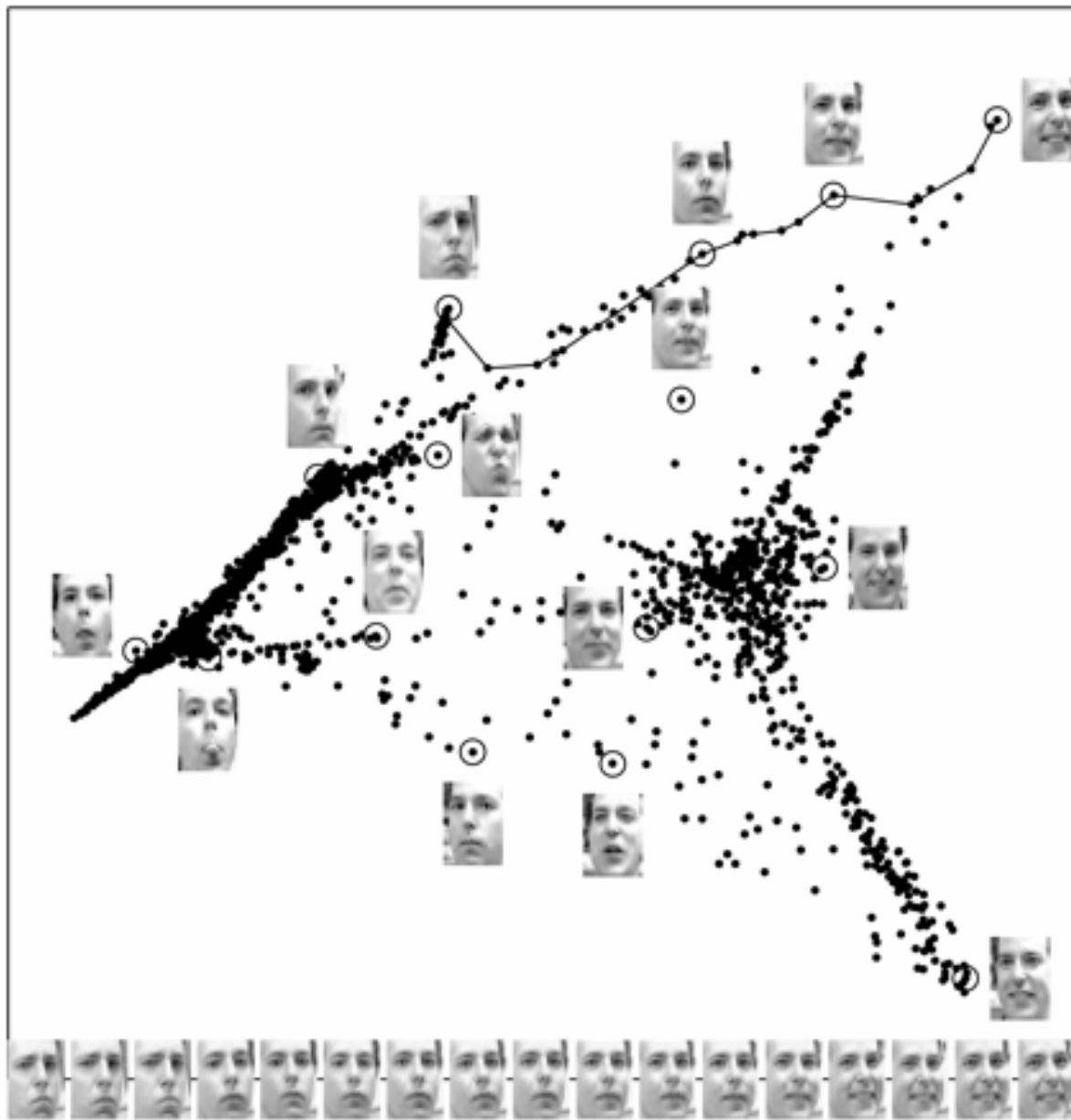
- **reduce data to lower dimensional embedding**
- **find locally linear patches of the data**
- **reconstruct data points from its neighbors**
- **invariant to rotations, rescalings and translations**

Locally Linear Embedding



Locally Linear Embedding

1. Find neighbors
2. Reconstruct data point as linear combination of neighbors
⇒ Weight matrix
3. Calculate low-dimensional representation of weight matrix using eigenvector analysis



Summary and Discussion