

Intelligent Data Analysis

Christian Borgelt

Intelligent Data Analysis and Graphical Models Research Unit
European Center for Soft Computing
c/ Gonzalo Gutiérrez Quirós s/n, 33600 Mieres, Spain

`christian.borgelt@softcomputing.es`
`http://www.softcomputing.es/`
`http://www.borgelt.net/`

Intelligent Data Analysis

- **Introduction**
- **Data and Knowledge**
 - Characteristics and Differences of Data and Knowledge
 - Quality Criteria for Knowledge
 - Example: Tycho Brahe and Johannes Kepler
- **Knowledge Discovery and Data Mining**
 - How to Find Knowledge?
 - The Knowledge Discovery Process (KDD Process)
 - Data Analysis / Data Mining Tasks
 - Data Analysis / Data Mining Methods
- **Summary**

Introduction

- Today every enterprise uses electronic information processing systems.
 - Production and distribution planning
 - Stock and supply management
 - Customer and personnel management
- Usually these systems are coupled with a database system (e.g. databases of customers, suppliers, parts etc.).
- Every possible individual piece of information can be retrieved.
- However: **Data alone are not enough.**
 - In a database one may “not see the wood for the trees”.
 - General patterns, structures, regularities go undetected.
 - Often such patterns can be exploited to increase turnover (e.g. joint sales in a supermarket).

Data

Examples of Data

- “Columbus discovered America in 1492.”
- “Mr Jones owns a Volkswagen Golf.”

Characteristics of Data

- refer to single instances
(single objects, persons, events, points in time etc.)
- describe individual properties
- are often available in huge amounts (databases, archives)
- are usually easy to collect or to obtain
(e.g. cash registers with scanners in supermarkets, Internet)
- do not allow us to make predictions

Knowledge

Examples of Knowledge

- “All masses attract each other.”
- “Every day at 5 pm there runs a train from Magdeburg to Berlin.”

Characteristic of Knowledge

- refers to *classes* of instances
(*sets* of objects, persons, points in time etc.)
- describes general patterns, structure, laws, principles etc.
- consists of as few statements as possible (this is an objective!)
- is usually difficult to find or to obtain
(e.g. natural laws, education)
- allows us to make predictions

Criteria to Assess Knowledge

- Not all statements are equally important, equally substantial, equally useful.
⇒ Knowledge must be assessed.

Assessment Criteria

- Correctness (probability, success in tests)
- Generality (range of validity, conditions of validity)
- Usefulness (relevance, predictive power)
- Comprehensibility (simplicity, clarity, parsimony)
- Novelty (previously unknown, unexpected)

Priority

- Science: correctness, generality, simplicity
- Economy: usefulness, comprehensibility, novelty

Tycho Brahe (1546–1601)

Who was Tycho Brahe?

- Danish nobleman and astronomer
- In 1582 he built an observatory on the island of Ven (32 km NE of Copenhagen).
- He determined the positions of the sun, the moon and the planets (accuracy: one angle minute, without a telescope!).
- He recorded the motions of the celestial bodies for several years.

Brahe's Problem

- He could not summarize the data he had collected in a uniform and consistent scheme.
- The planetary system he developed (the so-called Tychonic system) did not stand the test of time.

Johannes Kepler (1571–1630)

Who was Johannes Kepler?

- German astronomer and assistant of Tycho Brahe
- He advocated the Copernican planetary system.
- He tried all his life to find the laws that govern the motion of the planets.
- He started from the data that Tycho Brahe had collected.

Kepler's Laws

1. Each planet moves around the sun in an ellipse, with the sun at one focus.
2. The radius vector from the sun to the planet sweeps out equal areas in equal intervals of time.
3. The squares of the periods of any two planets are proportional to the cubes of the semi-major axes of their respective orbits: $T \sim a^{\frac{3}{2}}$.

How to find Knowledge?

We do not know any universal method to discover knowledge.

Problems

- Today huge amounts of data are available in databases.

*We are drowning in information,
but starving for knowledge.*

John Naisbett

- Manual methods of analysis have long ceased to be feasible.
- Simple aids (e.g. displaying data in charts) are too limited.

Attempts to Solve the Problems

- Intelligent Data Analysis
- Knowledge Discovery in Databases
- Data Mining

Knowledge Discovery and Data Mining

Knowledge Discovery and Data Mining

As a response to the challenge raised by the growing volume of data a new research area has emerged, which is usually characterized by one of the following phrases:

- **Knowledge Discovery in Databases (KDD)**

Usual characterization:

KDD is the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data. [Fayyad et al. 1996]

- **Data Mining**

- Data mining is that step of the knowledge discovery process in which data analysis methods are applied to find interesting patterns.
- It can be characterized by a set of types of tasks that have to be solved.
- It uses methods from a variety of research areas.
(statistics, databases, machine learning, artificial intelligence, soft computing etc.)

The Knowledge Discovery Process (KDD Process)

Preliminary Steps

- estimation of potential benefit
- definition of goals, feasibility study

Main Steps

- check data availability, data selection, if necessary, data collection
- preprocessing (60-80% of total overhead)
 - unification and transformation of data formats
 - data cleaning (error correction, outlier detection, imputation of missing values)
 - reduction / focusing (sample drawing, feature selection, prototype generation)
- **Data Mining** (using a variety of methods)
- visualization (also in parallel to preprocessing, data mining, and interpretation)
- interpretation, evaluation, and test of results
- deployment and documentation

Data Analysis / Data Mining Tasks

- **Classification**

Is this customer credit-worthy?

- **Segmentation, Clustering**

What groups of customers do I have?

- **Concept Description**

Which properties characterize fault-prone vehicles?

- **Prediction, Trend Analysis**

What will the exchange rate of the dollar be tomorrow?

- **Dependence/Association Analysis**

Which products are frequently bought together?

- **Deviation Analysis**

Are there seasonal or regional variations in turnover?

Data Analysis / Data Mining Methods 1

- **Classical Statistics**

(charts, parameter estimation, hypothesis testing, model selection, regression)

tasks: classification, prediction, trend analysis

- **Bayes Classifiers**

(probabilistic classification, naive and full Bayes classifiers)

tasks: classification, prediction

- **Decision and Regression Trees**

(top down induction, attribute selection measures, pruning)

tasks: classification, prediction

- **k-nearest Neighbor/Case-based Reasoning**

(lazy learning, similarity measures, data structures for fast search)

tasks: classification, prediction

Data Analysis / Data Mining Methods 2

- **Artificial Neural Networks**

(multilayer perceptrons, radial basis function networks, learning vector quantization)

tasks: classification, prediction, clustering

- **Cluster Analysis**

(k -means and fuzzy clustering, hierarchical agglomerative clustering)

tasks: segmentation, clustering

- **Association Rule Induction**

(frequent item set mining, rule generation)

tasks: association analysis

- **Inductive Logic Programming**

(rule generation, version space, search strategies, declarative bias)

tasks: classification, association analysis, concept description

Statistics

Statistics

- **Descriptive Statistics**

- Tabular and Graphical Representations
- Characteristic Measures
- Principal Component Analysis

- **Inductive Statistics**

- Parameter Estimation
(point and interval estimation, finding estimators)
- Hypothesis Testing
(parameter test, goodness-of-fit test, dependence test)
- Model Selection
(information criteria, minimum description length)

- **Summary**

Statistics: Introduction

Statistics is the art to collect, to display, to analyze, and to interpret data in order to gain new knowledge.

[Sachs 1999]

[...] statistics, that is, the mathematical treatment of reality, [...]

Hannah Arendt

There are lies, damned lies, and statistics.

Benjamin Disraeli

Statistics, n. Exactly 76.4% of all statistics (including this one) are invented on the spot. However, in 83% of cases it is inappropriate to admit it.

The Devil's IT Dictionary

Basic Notions

- **Object, Case**

Data describe objects, cases, persons etc.

- **(Random) Sample**

The objects or cases described by a data set is called a *sample*, their number the *sample size*.

- **Attribute**

Objects and cases are described by *attributes*, patients in a hospital, for example, by age, sex, blood pressure etc.

- **(Attribute) Value**

Attributes have different possible *values*.

The age of a patient, for example, is a non-negative number.

- **Sample Value**

The value an attribute has for an object in the sample is called *sample value*.

Scale Types / Attribute Types

Scale Type	Possible Operations	Examples
nominal (categorical, qualitative)	equality	sex blood group
ordinal (rank scale, comparative)	equality greater/less than	exam grade wind strength
metric (interval scale, quantitative)	equality greater/less than difference maybe ratio	length weight time temperature

- Nominal scales are sometimes divided into *dichotomic* (two values) and *polytomic* (more than two values).
- Metric scales may or may not allow us to form a ratio: weight and length do, temperature does not. time as duration does, time as calendar time does not.

Descriptive Statistics

Tabular Representations: Frequency Table

- Given data set: $x = (3, 4, 3, 2, 5, 3, 1, 2, 4, 3, 3, 4, 4, 1, 5, 2, 2, 3, 5, 3, 2, 4, 3, 2, 3)$

a_k	h_k	r_k	$\sum_{i=1}^k h_i$	$\sum_{i=1}^k r_i$
1	2	$\frac{2}{25} = 0.08$	2	$\frac{2}{25} = 0.08$
2	6	$\frac{6}{25} = 0.24$	8	$\frac{8}{25} = 0.32$
3	9	$\frac{9}{25} = 0.36$	17	$\frac{17}{25} = 0.68$
4	5	$\frac{5}{25} = 0.20$	22	$\frac{22}{25} = 0.88$
5	3	$\frac{3}{25} = 0.12$	25	$\frac{25}{25} = 1.00$

- Absolute Frequency** h_k (frequency of an attribute value a_k in the sample).
- Relative Frequency** $r_k = \frac{h_k}{n}$, where n is the sample size (here $n = 25$).
- Cumulated Absolute/Relative Frequency** $\sum_{i=1}^k h_i$ and $\sum_{i=1}^k r_i$.

Tabular Representations: Contingency Tables

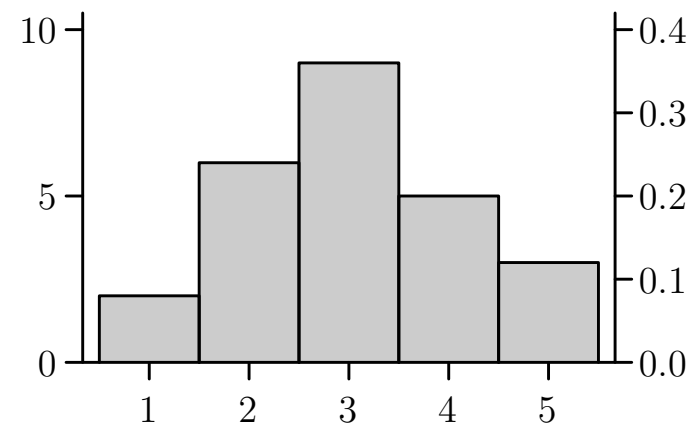
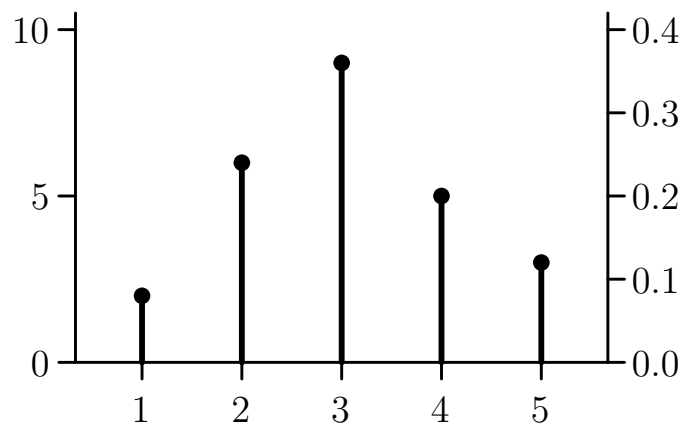
- Frequency tables for two or more attributes are called **contingency tables**.
- They contain the absolute or relative frequency of **value combinations**.

	a_1	a_2	a_3	a_4	Σ
b_1	8	3	5	2	18
b_2	2	6	1	3	12
b_3	4	1	2	7	14
Σ	14	10	8	12	44

- A contingency table may also contain the **marginal frequencies**, i.e., the frequencies of the values of individual attributes.
- Contingency tables for a higher number of dimensions (> 4) may be difficult to read.

Graphical Representations: Pole and Bar Chart

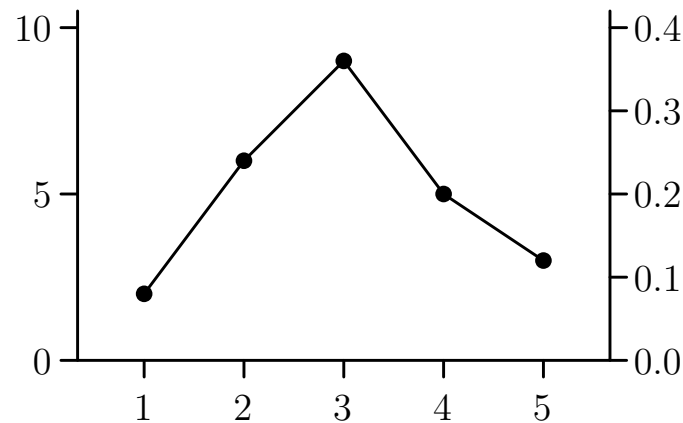
- Numbers, which may be, for example, the frequencies of attribute values are represented by the lengths of poles (left) or bars (right).



- Bar charts are the most frequently used and most comprehensible way of displaying absolute frequencies.
- A wrong impression can result if the vertical scale does not start at 0 (for frequencies or other absolute numbers).

Frequency Polygon and Line Chart

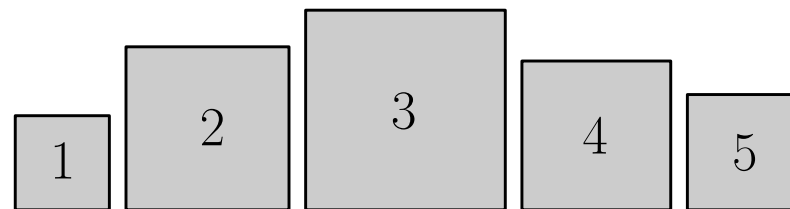
- Frequency polygon: the ends of the poles of a pole chart are connected by lines. (Numbers are still represented by lengths.)



- If the attribute values on the horizontal axis are not ordered, connecting the ends of the poles does not make sense.
- Line charts are frequently used to display time series.

Area and Volume Charts

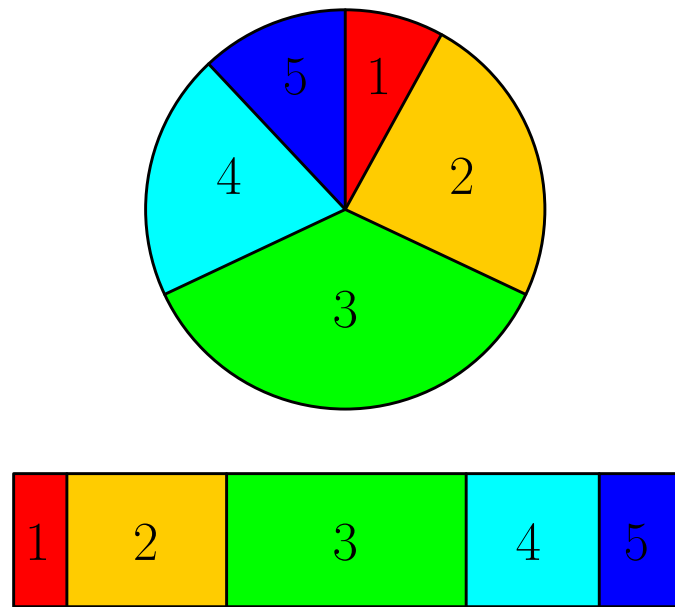
- Numbers may also be represented by other geometric quantities than lengths, like areas or volumes.
- Area and volume charts are usually less comprehensible than bar charts, because humans have more difficulties to compare areas and especially volumes than lengths. (exception: the represented numbers are areas or volumes)



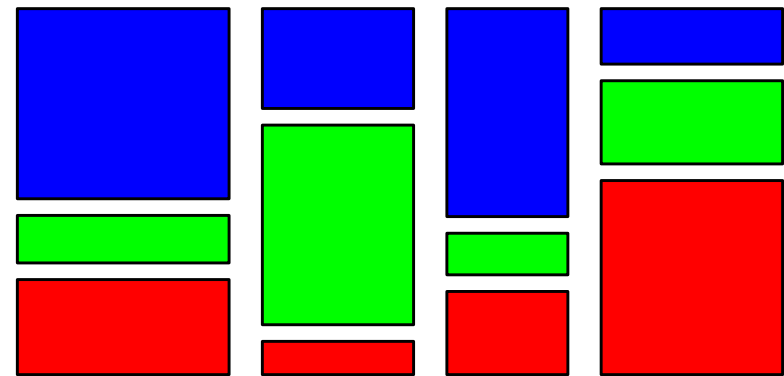
- Sometimes the height of a two- or three-dimensional object is used to represent a number. The diagram then conveys a misleading impression.

Pie and Stripe Charts

- Relative numbers may be represented by angles or sections of a stripe.



Mosaic Chart



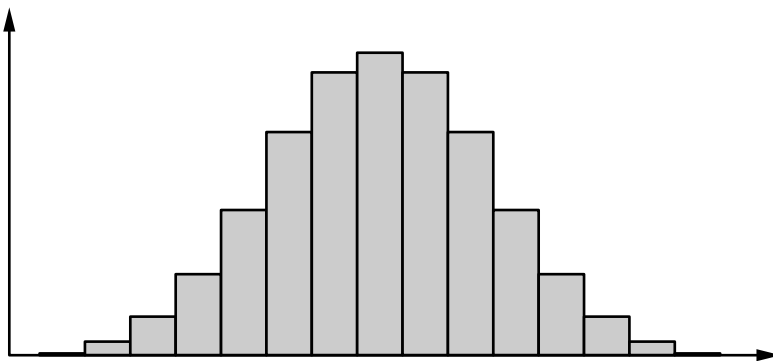
- Mosaic charts can be used to display contingency tables.
- More than two attributes are possible, but then separation distances and color must support the visualization to keep it comprehensible.

Histograms

- Intuitively: **Histograms are frequency bar charts for metric data.**
- However: Since there are so many different values, **values have to be grouped** in order to arrive a proper representation.

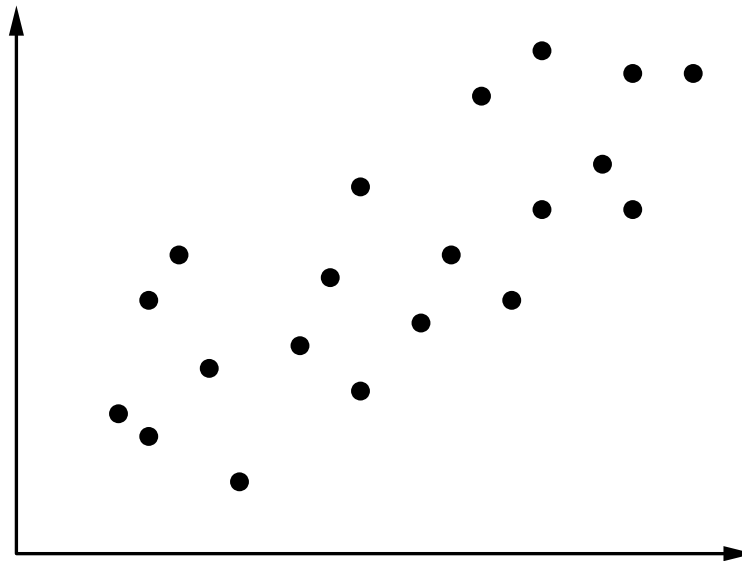
Most common approach: form equally sized intervals (so-called **bins**) and count the frequency of sample values inside each interval.

- **Attention:** Depending on the size and the position of the bins the histogram may look considerably different.
- In sketches often only a rough outline of a histogram is drawn:



Scatter Plots

- Scatter plots are used to display two-dimensional metric data sets.
- Sample values are the coordinates of a point.
(Numbers are represented by lengths.)



- Scatter plots provide a simple means for checking for dependency.

How to Lie with Statistics

pictures not available in online version

Often the vertical axis of a pole or bar chart does not start at zero, but at some higher value.

In such a case the conveyed impression of the ratio of the depicted values is completely wrong.

This effect is used to brag about increases in turnover, speed etc.

Sources of these diagrams and those on the following transparencies:

D. Huff: How to Lie with Statistics.

W. Krämer: So lügt man mit Statistik.

How to Lie with Statistics

pictures not available in online version

- Depending on the position of the zero line of a pole, bar, or line chart completely different impressions can be conveyed.

How to Lie with Statistics

pictures not available in online version

- Poles and bars are frequently replaced by (sketches of) objects in order to make the diagram more aesthetically appealing.
- However, objects are perceived as 2- or even 3-dimensional and thus convey a completely different impression of the numerical ratios.

How to Lie with Statistics

pictures not available in online version

How to Lie with Statistics

pictures not available in online version

- In the left diagram the areas of the barrels represent the numerical value. However, since the barrels are drawn 3-dimensional, a wrong impression of the numerical ratios is conveyed.
- The right diagram is particularly striking: an area measure is represented by the *side length* of a rectangle representing the apartment.

Descriptive Statistics: Characteristic Measures

Descriptive Statistics: Characteristic Measures

Idea: Describe a given sample by few characteristic measures and thus summarize the data.

- **Localization Measures**

Localization measures describe, usually by a single number, where the data points of a sample are located in the domain of an attribute.

- **Dispersion Measures**

Dispersion measures describe how much the data points vary around a localization parameter and thus indicate how well this parameter captures the localization of the data.

- **Shape Measures**

Shape measures describe the shape of the distribution of the data points relative to a reference distribution. The most common reference distribution is the normal distribution (Gaussian).

Localization Measures: Mode and Median

- **Mode** x^*

The mode is the attribute value that is most frequent in the sample.

It need not be unique, because several values can have the same frequency.

It is the most general measure, because it is applicable for all scale types.

- **Median** \tilde{x}

The median minimizes the sum of absolute differences:

$$\sum_{i=1}^n |x_i - \tilde{x}| = \min. \quad \text{and thus it is} \quad \sum_{i=1}^n \text{sgn}(x_i - \tilde{x}) = 0$$

If $x = (x_{(1)}, \dots, x_{(n)})$ is a sorted data set, the median is defined as

$$\tilde{x} = \begin{cases} x_{(\frac{n+1}{2})}, & \text{if } n \text{ is odd,} \\ \frac{1}{2} \left(x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)} \right), & \text{if } n \text{ is even.} \end{cases}$$

The median is applicable to ordinal and metric attributes.

Localization Measures: Arithmetic Mean

- **Arithmetic Mean** \bar{x}

The arithmetic mean minimizes the sum of squared differences:

$$\sum_{i=1}^n (x_i - \bar{x})^2 = \min. \quad \text{and thus it is} \quad \sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - n\bar{x} = 0$$

The arithmetic mean is defined as

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

The arithmetic mean is only applicable to metric attributes.

- Even though the arithmetic mean is the most common localization measure, the **median** is preferable if
 - there are few sample cases,
 - the distribution is asymmetric, and/or
 - one expects that outliers are present.

How to Lie with Statistics

pictures not available in online version

Dispersion Measures: Range and Interquantile Range

A man with his head in the freezer and feet in the oven
is *on the average* quite comfortable.

old statistics joke

- **Range R**

The range of a data set is the difference between the maximum and the minimum value.

$$R = x_{\max} - x_{\min} = \max_{i=1}^n x_i - \min_{i=1}^n x_i$$

- **Interquantile Range**

The p -quantile of a data set is a value such that a fraction of p of all sample values are smaller than this value. (The median is the $\frac{1}{2}$ -quantile.)

The p -interquantile range, $0 < p < \frac{1}{2}$, is the difference between the $(1 - p)$ -quantile and the p -quantile.

The most common is the *interquartile range* ($p = \frac{1}{4}$)

Dispersion Measures: Average Absolute Deviation

- **Average Absolute Deviation**

The average absolute deviation is the average of the absolute deviations of the sample values from the median or the arithmetic mean.

- Average Absolute Deviation from the **Median**

$$d_{\tilde{x}} = \frac{1}{n} \sum_{i=1}^n |x_i - \tilde{x}|$$

- Average Absolute Deviation from the **Arithmetic Mean**

$$d_{\bar{x}} = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

- It is always $d_{\tilde{x}} \leq d_{\bar{x}}$, since the median minimizes the sum of absolute deviations (see the definition of the median).

Dispersion Measures: Variance and Standard Deviation

- **Variance s^2**

It would be natural to define the variance as the average squared deviation:

$$v^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

However, inductive statistics suggests that it is better defined as

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

- **Standard Deviation s**

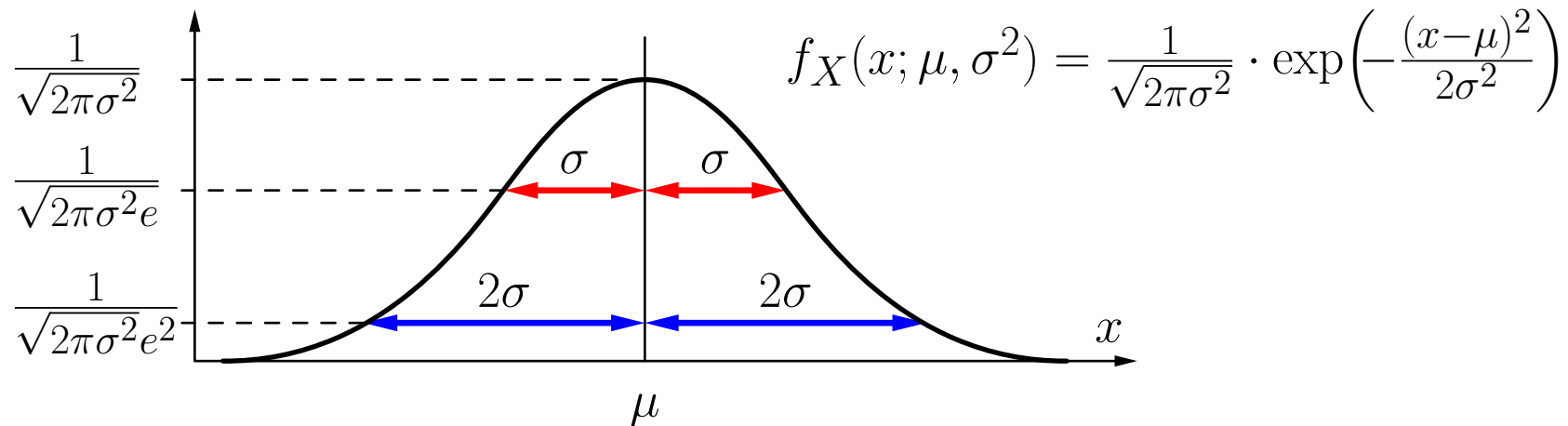
The standard deviation is the square root of the variance, i.e.,

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}.$$

Dispersion Measures: Variance and Standard Deviation

- **Special Case: Normal/Gaussian Distribution**

The variance/standard deviation provides information about the height of the mode and the width of the curve.



- μ : expected value, estimated by mean value \bar{x}
 σ^2 : variance, estimated by (empirical) variance s^2
 σ : standard deviation, estimated by (empirical) standard deviation s
(Details about parameter estimation are studied later.)

Dispersion Measures: Variance and Standard Deviation

Note that it is often more convenient to compute the variance using the formula that results from the following transformation:

$$\begin{aligned}s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\&= \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2 \right) \\&= \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - 2n\bar{x}^2 + n\bar{x}^2 \right) = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \\&= \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right)\end{aligned}$$

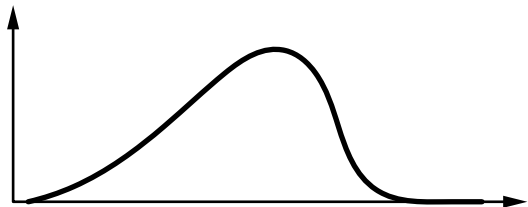
- Advantage: The sums $\sum_{i=1}^n x_i$ and $\sum_{i=1}^n x_i^2$ can both be computed in the same traversal of the data and from them both mean and variance are computable.

Shape Measures: Skewness

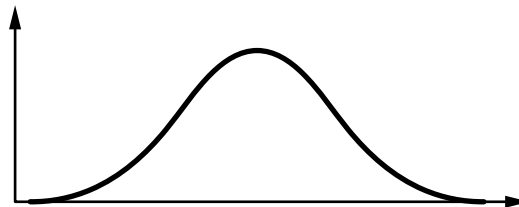
- The **skewness** α_3 (or **skew** for short) measures whether, and if, how much, a distribution differs from a symmetric distribution.
- It is computed from the 3rd moment about the mean, which explains the index 3.

$$\alpha_3 = \frac{1}{n \cdot v^3} \sum_{i=1}^n (x_i - \bar{x})^3 = \frac{1}{n} \sum_{i=1}^n z_i^3$$

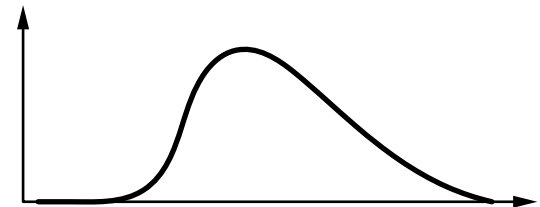
$$\text{where } z_i = \frac{x_i - \bar{x}}{v} \text{ and } v^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$



$\alpha_3 < 0$: right steep



$\alpha_3 = 0$: symmetric



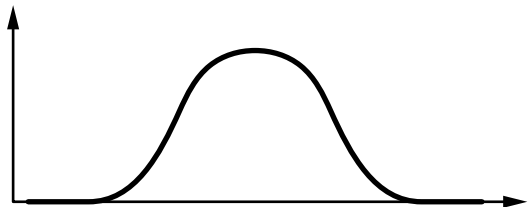
$\alpha_3 > 0$: left steep

Shape Measures: Kurtosis

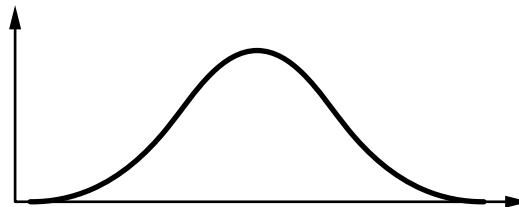
- The **kurtosis** or **excess** α_4 measures how much a distribution is arched, usually compared to a Gaussian distribution.
- It is computed from the 4th moment about the mean, which explains the index 4.

$$\alpha_4 = \frac{1}{n \cdot v^4} \sum_{i=1}^n (x_i - \bar{x})^4 = \frac{1}{n} \sum_{i=1}^n z_i^4$$

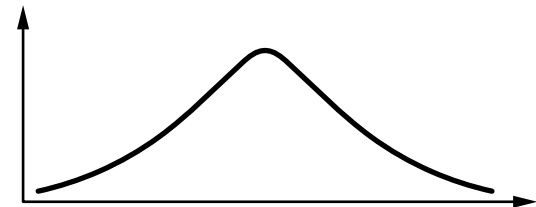
$$\text{where } z_i = \frac{x_i - \bar{x}}{v} \text{ and } v^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$



$\alpha_4 < 3$: leptokurtic



$\alpha_4 = 3$: Gaussian



$\alpha_4 > 3$: platikurtic

Moments of Data Sets

- The k -th **moment** of a dataset is defined as

$$m'_k = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

The first moment is the **mean** $m'_1 = \bar{x}$ of the data set.

Using the moments of a data set the **variance** s^2 can also be written as

$$s^2 = \frac{1}{n-1} \left(m'_2 - \frac{1}{n} m_1'^2 \right) \quad \text{and also} \quad v^2 = \frac{1}{n} m'_2 - \frac{1}{n^2} m_1'^2.$$

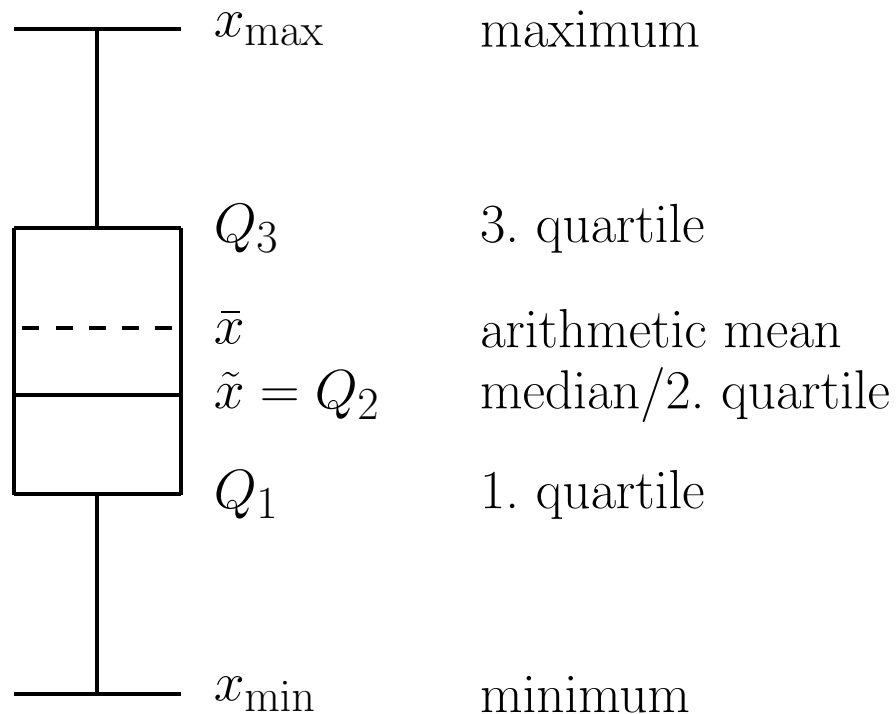
- The k -th **moment about the mean** is defined as

$$m_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^k.$$

It is $m_1 = 0$ and $m_2 = v^2$ (i.e., the **average squared deviation**).

The **skewness** is $\alpha_3 = \frac{m_3}{m_2^{3/2}}$ and the **kurtosis** is $\alpha_4 = \frac{m_4}{m_2^2}$.

Visualizing Characteristic Measures: Box Plots



A box plot is a common way to combine some important characteristic measures into a single graphic.

Often the central box is drawn laced $\langle \rangle$ w.r.t. the arithmetic mean in order to emphasize its location.

Box plots are often used to get a quick impression of the distribution of the data by showing them side by side for several attributes.

Multidimensional Characteristic Measures

General Idea: Transfer the characteristic measures to vectors.

- **Arithmetic Mean**

The arithmetic mean for multi-dimensional data is the vector mean of the data points. For two dimensions it is

$$\overline{(x, y)} = \frac{1}{n} \sum_{i=1}^n (x_i, y_i) = (\bar{x}, \bar{y})$$

For the arithmetic mean the transition to several dimensions only combines the arithmetic means of the individual dimensions into one vector.

- Other measures are transferred in a similar way.

However, sometimes the transfer leads to new quantities, as for the variance.

Excursion: Vector Products

For the variance, the square of the difference to the mean has to be generalized.

Inner Product
Scalar Product

$$\vec{v}^\top \vec{v} \quad \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{pmatrix}$$

$$(v_1, v_2, \dots, v_m) \quad \sum_{i=1}^m v_i^2$$

Outer Product
Matrix Product

$$\vec{v} \vec{v}^\top \quad (v_1, v_2, \dots, v_m)$$

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{pmatrix} \quad \begin{pmatrix} v_1^2 & v_1 v_2 & \cdots & v_1 v_m \\ v_1 v_2 & v_2^2 & \cdots & v_2 v_m \\ \vdots & \vdots & \ddots & \vdots \\ v_1 v_m & v_2 v_m & \cdots & v_m^2 \end{pmatrix}$$

- In principle both vector products may be used for a generalization.
- The second, however, yields more information about the distribution:
 - a measure of the (linear) dependence of the attributes,
 - a description of the direction dependence of the dispersion.

Covariance Matrix

- **Covariance Matrix**

Compute variance formula with vectors (square: outer product $\vec{v}\vec{v}^\top$):

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n \left(\begin{pmatrix} x_i \\ y_i \end{pmatrix} - \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} \right) \left(\begin{pmatrix} x_i \\ y_i \end{pmatrix} - \begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix} \right)^\top = \begin{pmatrix} s_x^2 & s_{xy} \\ s_{xy} & s_y^2 \end{pmatrix}$$

where

$$s_x^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \quad (\text{variance of } x)$$

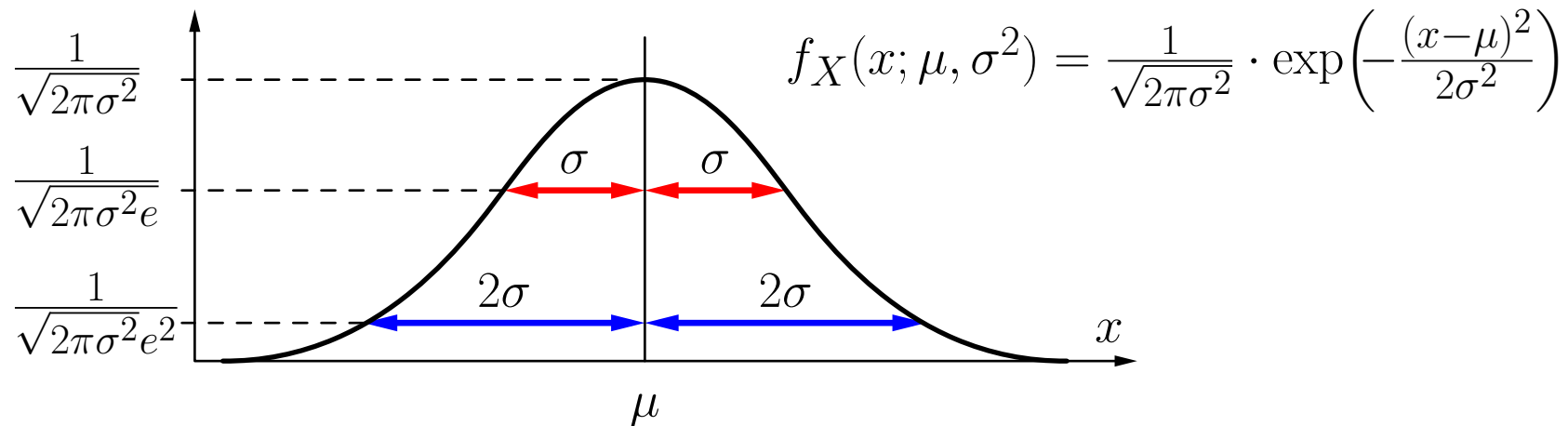
$$s_y^2 = \frac{1}{n-1} \left(\sum_{i=1}^n y_i^2 - n\bar{y}^2 \right) \quad (\text{variance of } y)$$

$$s_{xy} = \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \right) \quad (\text{covariance of } x \text{ and } y)$$

Reminder: Variance and Standard Deviation

- **Special Case: Normal/Gaussian Distribution**

The variance/standard deviation provides information about the height of the mode and the width of the curve.



- μ : expected value, estimated by mean value \bar{x} ,
 σ^2 : variance, estimated by (empirical) variance s^2 ,
 σ : standard deviation, estimated by (empirical) standard deviation s .
Important: standard deviation has same unit as expected value.

Multivariate Normal Distribution

- A **univariate normal distribution** has the density function

$$f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

μ : expected value, estimated by mean value \bar{x} ,
 σ^2 : variance, estimated by (empirical) variance s^2 ,
 σ : standard deviation, estimated by (empirical) standard deviation s .

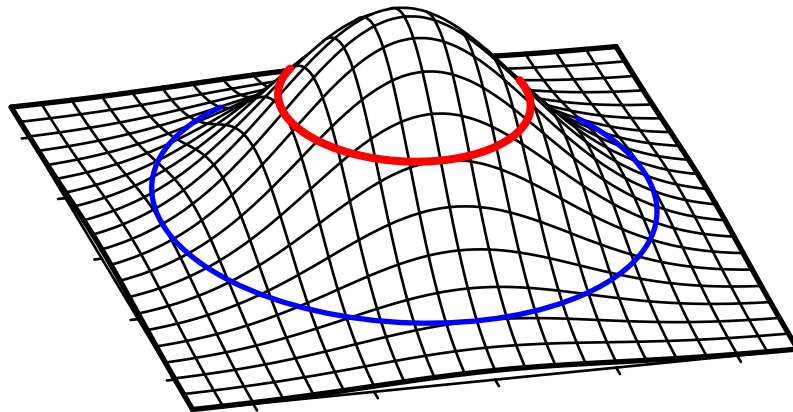
- A **multivariate normal distribution** has the density function

$$f_{\vec{X}}(\vec{x}; \vec{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \cdot \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^\top \Sigma^{-1}(\vec{x} - \vec{\mu})\right)$$

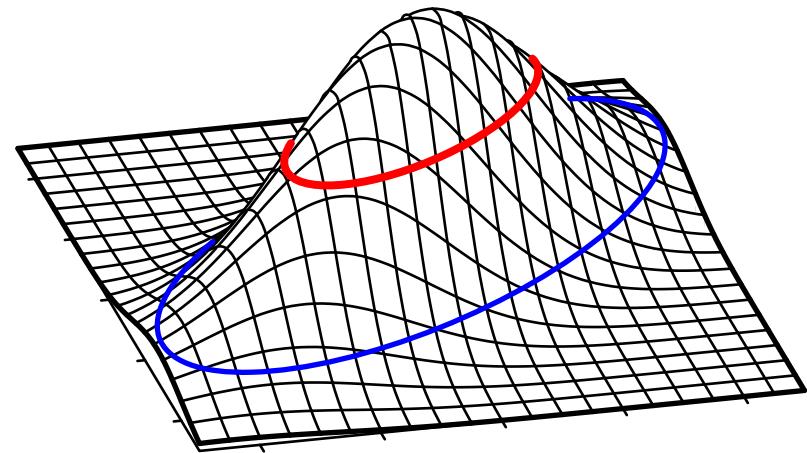
m : size of the vector \vec{x} (it is m -dimensional),
 $\vec{\mu}$: mean value vector, estimated by (empirical) mean value vector $\bar{\vec{x}}$,
 Σ : covariance matrix, estimated by (empirical) covariance matrix \mathbf{S} ,
 $|\Sigma|$: determinant of the covariance matrix Σ .

Interpretation of a Covariance Matrix

- The variance/standard deviation relates the spread of the distribution to the spread of a **standard normal distribution** ($\sigma^2 = \sigma = 1$).
- The covariance matrix relates the spread of the distribution to the spread of a **multivariate standard normal distribution** ($\Sigma = \mathbf{1}$).
- Example: bivariate normal distribution



standard



general

- **Question:** Is there a multivariate analog of standard deviation?

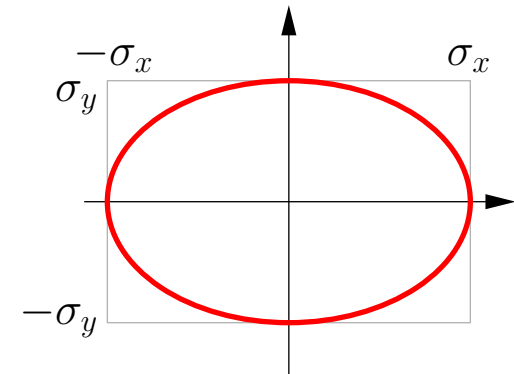
Interpretation of a Covariance Matrix

Question: Is there a multivariate analog of standard deviation?

First insight:

If all covariances vanish,
the contour lines are axes-parallel ellipses.
The upper ellipse is inscribed into the
rectangle $[-\sigma_x, \sigma_x] \times [-\sigma_y, \sigma_y]$.

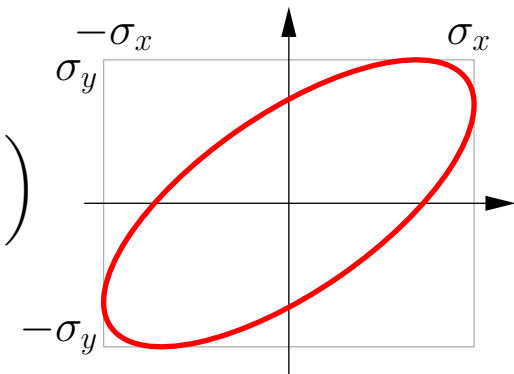
$$\Sigma = \begin{pmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{pmatrix}$$



Second insight:

If the covariances do not vanish,
the contour lines are rotated ellipses.
Still the upper ellipse is inscribed into the
rectangle $[-\sigma_x, \sigma_x] \times [-\sigma_y, \sigma_y]$.

$$\Sigma = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$$



Consequence: A covariance matrix describes a scaling and a rotation.

Cholesky Decomposition

- Intuitively: **Compute an analog of standard deviation.**
- Let \mathbf{S} be a symmetric, positive definite matrix (e.g. a covariance matrix). Cholesky decomposition serves the purpose to compute a “square root” of \mathbf{S} .
 - symmetric: $\forall 1 \leq i, j \leq m : s_{ij} = s_{ji}$
 - positive definite: for all m -dimensional vectors $\vec{v} \neq \vec{0}$ it is $\vec{v}^\top \mathbf{S} \vec{v} > 0$
- Formally: Compute a left/lower triangular matrix \mathbf{L} such that $\mathbf{L}\mathbf{L}^\top = \mathbf{S}$. (\mathbf{L}^\top is the transpose of the matrix \mathbf{L} .)

$$l_{ii} = \left(s_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 \right)^{\frac{1}{2}}$$
$$l_{ji} = \frac{1}{l_{ii}} \left(s_{ij} - \sum_{k=1}^{i-1} l_{ik} l_{jk} \right), \quad j = i + 1, i + 2, \dots, m.$$

Cholesky Decomposition

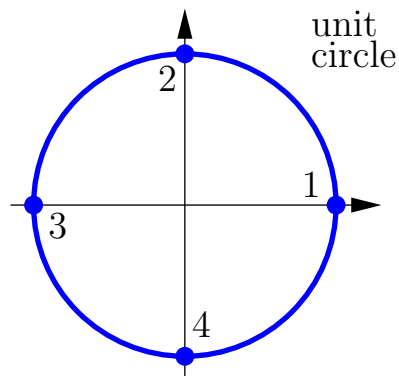
Special Case: Two Dimensions

- Covariance matrix

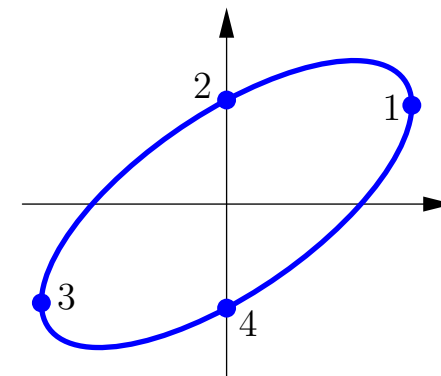
$$\Sigma = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$$

- Cholesky decomposition

$$\mathbf{L} = \begin{pmatrix} \sigma_x & 0 \\ \frac{\sigma_{xy}}{\sigma_x} & \frac{1}{\sigma_x} \sqrt{\sigma_x^2 \sigma_y^2 - \sigma_{xy}^2} \end{pmatrix}$$



→
mapping with \mathbf{L}
 $\vec{v}' = \mathbf{L}\vec{v}$



Eigenvalue Decomposition

- Also yields an **analog of standard deviation**.
- Computationally more expensive than Cholesky decomposition.
- Let \mathbf{S} be a symmetric, positive definite matrix (e.g. a covariance matrix).
 - \mathbf{S} can be written as

$$\mathbf{S} = \mathbf{R} \operatorname{diag}(\lambda_1, \dots, \lambda_m) \mathbf{R}^{-1},$$

where the λ_j , $j = 1, \dots, m$, are the eigenvalues of \mathbf{S}
and the columns of \mathbf{R} are the (normalized) eigenvectors of \mathbf{S} .

- The eigenvalues λ_j , $j = 1, \dots, m$, of \mathbf{S} are all positive
and the eigenvectors of \mathbf{S} are orthonormal ($\rightarrow \mathbf{R}^{-1} = \mathbf{R}^\top$).
- Due to the above, \mathbf{S} can be written as $\mathbf{S} = \mathbf{T} \mathbf{T}^\top$, where

$$\mathbf{T} = \mathbf{R} \operatorname{diag} \left(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m} \right)$$

Eigenvalue Decomposition

Special Case: Two Dimensions

- Covariance matrix

$$\Sigma = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$$

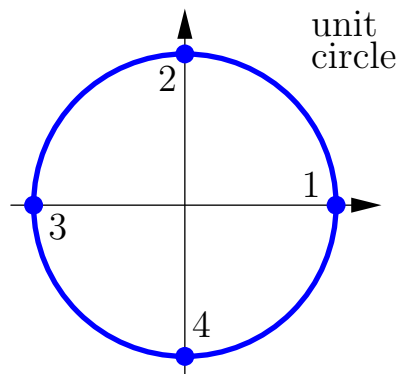
- Eigenvalue decomposition

$$\mathbf{T} = \begin{pmatrix} c & -s \\ s & c \end{pmatrix} \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix},$$

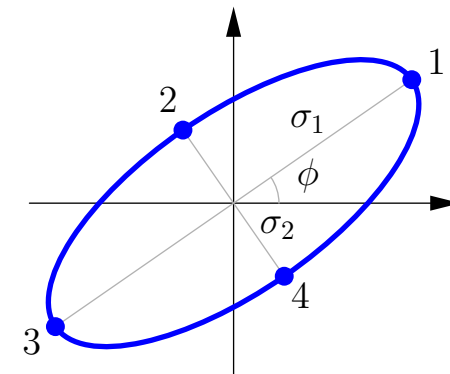
$$s = \sin \phi, c = \cos \phi, \phi = \frac{1}{2} \arctan \frac{2\sigma_{xy}}{\sigma_x^2 - \sigma_y^2},$$

$$\sigma_1 = \sqrt{c^2 \sigma_x^2 + s^2 \sigma_y^2 + 2sc\sigma_{xy}},$$

$$\sigma_2 = \sqrt{s^2 \sigma_x^2 + c^2 \sigma_y^2 - 2sc\sigma_{xy}}.$$



mapping with \mathbf{T}
 $\vec{v}' = \mathbf{T}\vec{v}$



Eigenvalue Decomposition

Eigenvalue decomposition enables us to write a covariance matrix Σ as

$$\Sigma = \mathbf{T}\mathbf{T}^\top \quad \text{with} \quad \mathbf{T} = \mathbf{R} \operatorname{diag} \left(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m} \right).$$

As a consequence we can write its inverse Σ^{-1} as

$$\Sigma^{-1} = \mathbf{U}^\top \mathbf{U} \quad \text{with} \quad \mathbf{U} = \operatorname{diag} \left(\lambda_1^{-\frac{1}{2}}, \dots, \lambda_m^{-\frac{1}{2}} \right) \mathbf{R}^\top.$$

\mathbf{U} describes the inverse mapping of \mathbf{T} , i.e., rotates the ellipse so that its axes coincide with the coordinate axes and then scales the axes to unit length. Hence:

$$(\vec{x} - \vec{y})^\top \Sigma^{-1} (\vec{x} - \vec{y}) = (\vec{x} - \vec{y})^\top \mathbf{U}^\top \mathbf{U} (\vec{x} - \vec{y}) = (\vec{x}' - \vec{y}')^\top (\vec{x}' - \vec{y}'),$$

where $\vec{x}' = \mathbf{U}\vec{x}$ and $\vec{y}' = \mathbf{U}\vec{y}$.

Result: $(\vec{x} - \vec{y})^\top \Sigma^{-1} (\vec{x} - \vec{y})$ is equivalent to the squared **Euclidean distance** in the properly scaled eigensystem of the covariance matrix Σ .

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^\top \Sigma^{-1} (\vec{x} - \vec{y})} \quad \text{is called } \mathbf{Mahalanobis \ distance}.$$

Eigenvalue Decomposition

Eigenvector decomposition also shows that the determinant of the covariance matrix $\mathbf{\Sigma}$ provides a measure of the (hyper-)volume of the (hyper-)ellipsoid. It is

$$|\mathbf{\Sigma}| = |\mathbf{R}| |\text{diag}(\lambda_1, \dots, \lambda_m)| |\mathbf{R}^\top| = |\text{diag}(\lambda_1, \dots, \lambda_m)| = \prod_{i=1}^m \lambda_i,$$

since $|\mathbf{R}| = |\mathbf{R}^\top| = 1$ as \mathbf{R} is orthogonal with unit length columns, and thus

$$\sqrt{|\mathbf{\Sigma}|} = \prod_{i=1}^m \sqrt{\lambda_i},$$

which is proportional to the (hyper-)volume of the (hyper-)ellipsoid.

To be precise, the volume of the m -dimensional (hyper-)ellipsoid a (hyper-)sphere with radius r is mapped to with a covariance matrix $\mathbf{\Sigma}$ is

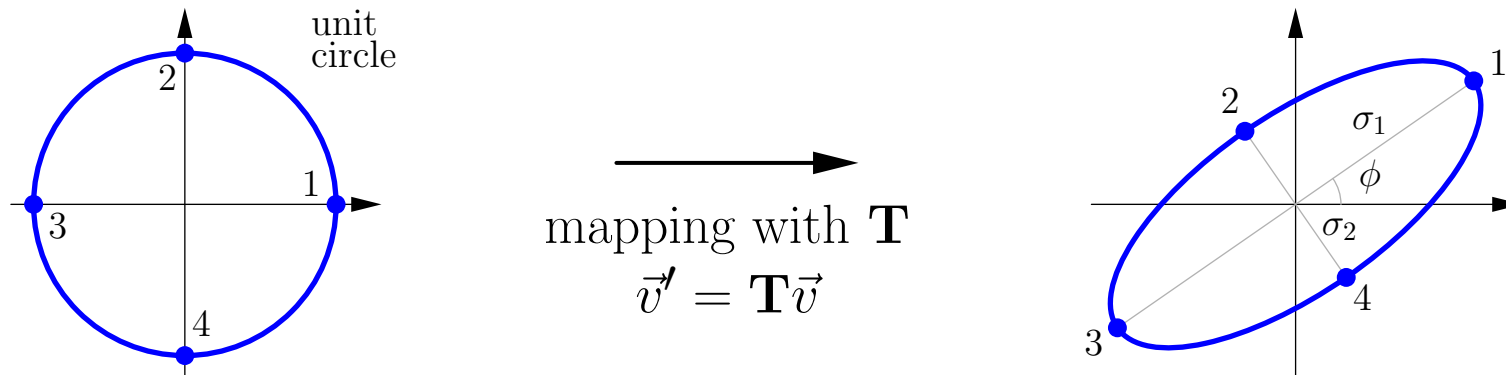
$$V_m(r) = \frac{\pi^{\frac{m}{2}} r^m}{\Gamma\left(\frac{m}{2} + 1\right)} \sqrt{|\mathbf{\Sigma}|}, \quad \text{where} \quad \begin{aligned} \Gamma(x) &= \int_0^\infty e^{-t} t^{x-1} dt, \quad x > 0, \\ \Gamma(x+1) &= x \cdot \Gamma(x), \quad \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(1) = 1. \end{aligned}$$

Eigenvalue Decomposition

Special Case: Two Dimensions

- Covariance matrix and its eigenvalue decomposition:

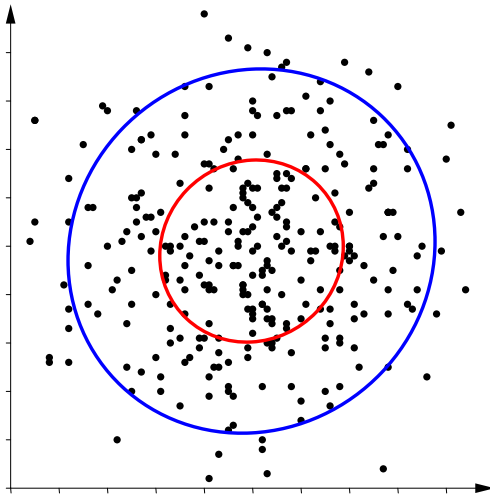
$$\Sigma = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix} \quad \text{and} \quad \mathbf{T} = \begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix} \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}.$$



- The area of the ellipse, to which the unit circle (area π) is mapped, is

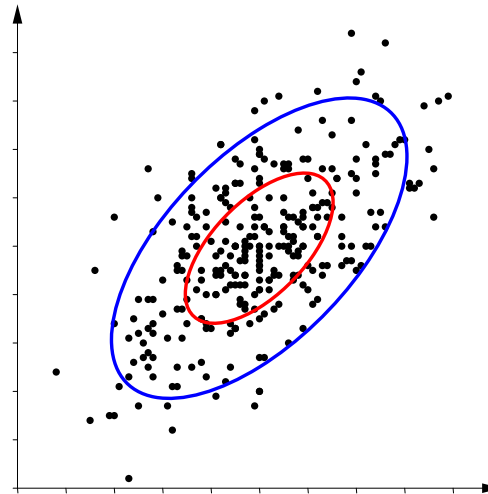
$$A = \pi \sigma_1 \sigma_2 = \pi \sqrt{|\Sigma|}.$$

Covariance Matrices of Example Data Sets



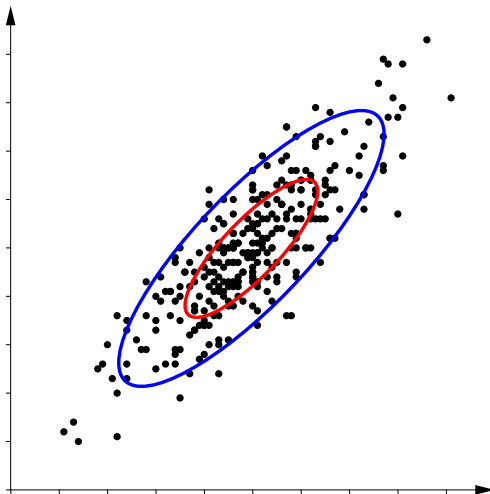
$$\Sigma = \begin{pmatrix} 3.59 & 0.19 \\ 0.19 & 3.54 \end{pmatrix}$$

$$\mathbf{L} = \begin{pmatrix} 1.90 & 0 \\ 0.10 & 1.88 \end{pmatrix}$$



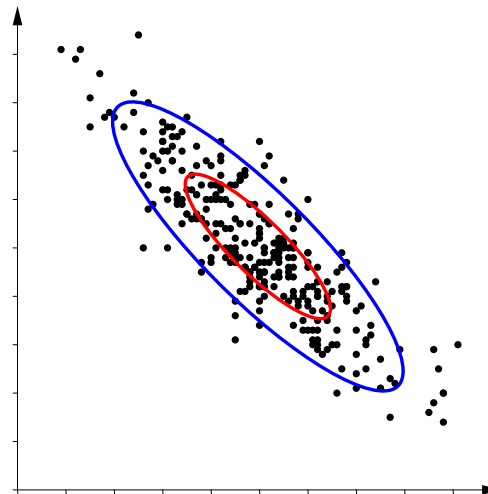
$$\Sigma = \begin{pmatrix} 2.33 & 1.44 \\ 1.44 & 2.41 \end{pmatrix}$$

$$\mathbf{L} = \begin{pmatrix} 1.52 & 0 \\ 0.95 & 1.22 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 1.88 & 1.62 \\ 1.62 & 2.03 \end{pmatrix}$$

$$\mathbf{L} = \begin{pmatrix} 1.37 & 0 \\ 1.18 & 0.80 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} 2.25 & -1.93 \\ -1.93 & 2.23 \end{pmatrix}$$

$$\mathbf{L} = \begin{pmatrix} 1.50 & 0 \\ -1.29 & 0.76 \end{pmatrix}$$

Covariance Matrix: Summary

- A covariance matrix provides information about the **height of the mode** and about the **spread/dispersion** of a multivariate normal distribution (or of a set of data points that are roughly normally distributed).
- A multivariate **analog of standard deviation** can be computed with Cholesky decomposition and eigenvalue decomposition. The resulting matrix describes the distribution's shape and orientation.
- The shape and the orientation of a two-dimensional normal distribution can be visualized as an **ellipse** (curve of equal probability density; similar to a **contour line** — line of equal height — on a map.)
- The shape and the orientation of a three-dimensional normal distribution can be visualized as an **ellipsoid** (surface of equal probability density).
- The (square root of the) **determinant** of a covariance matrix describes the spread of a multivariate normal distribution with a single value. It is a measure of the area or (hyper-)volume of the (hyper-)ellipsoid.

Correlation and Principal Component Analysis

Correlation Coefficient

- The covariance is a measure of the strength of **linear dependence** of the two quantities.
- However, its value depends on the variances of the individual dimensions.
⇒ Normalize to unit variance in the individual dimensions.

- **Correlation Coefficient**

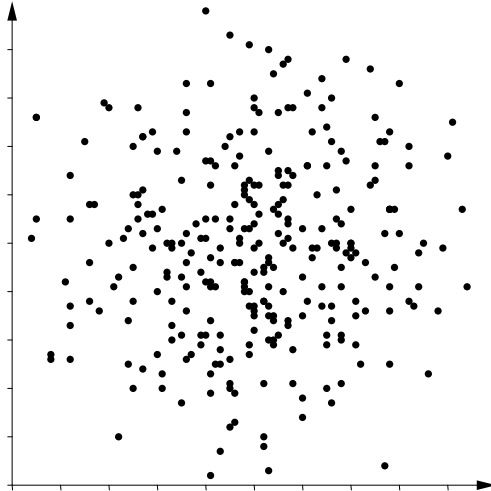
(more precisely: Pearson's Product Moment Correlation Coefficient)

$$r = \frac{s_{xy}}{s_x s_y}, \quad r \in [-1, +1].$$

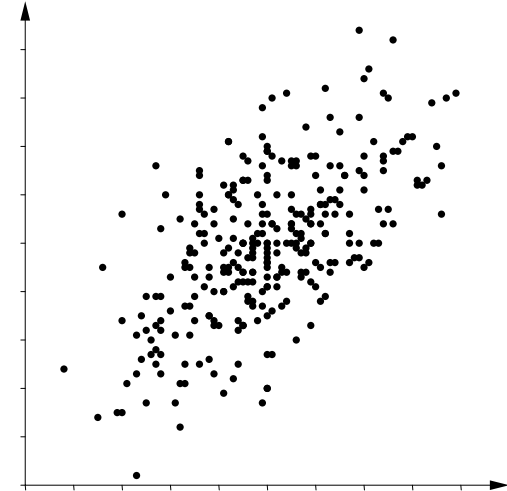
- r measures the strength of linear dependence:
 $r = -1$: the data points lie perfectly on a descending straight line.
 $r = +1$: the data points lie perfectly on an ascending straight line.
- $r = 0$: there is no **linear** dependence between the two attributes
(but there may be a non-linear dependence!).

Correlation Coefficients of Example Data Sets

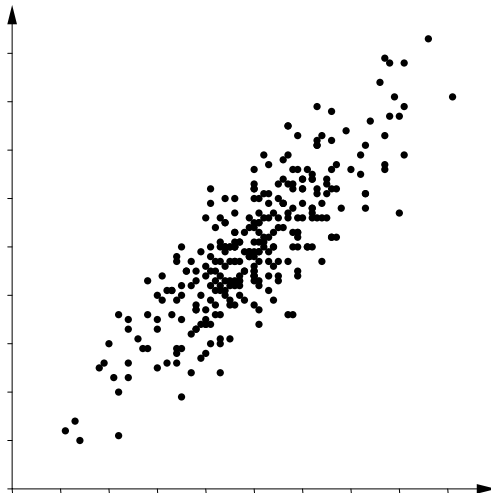
no
correlation
($r \approx 0.05$)



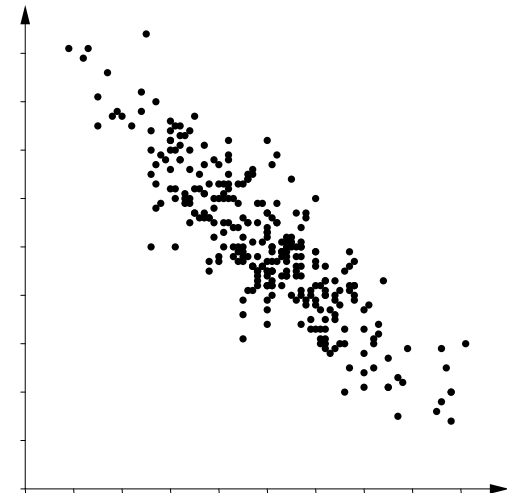
weak
positive
correlation
($r \approx 0.61$)



strong
positive
correlation
($r \approx 0.83$)



strong
negative
correlation
($r \approx -0.86$)



Correlation Matrix

- **Normalize Data**

Transform data to mean value 0 and variance/standard deviation 1:

$$\forall i; 1 \leq i \leq n : \quad x'_i = \frac{x_i - \bar{x}}{s_x}, \quad y'_i = \frac{y_i - \bar{y}}{s_y}.$$

- **Compute Covariance Matrix of Normalized Data**

Sum outer products of transformed data vectors:

$$\Sigma' = \frac{1}{n-1} \sum_{i=1}^n \begin{pmatrix} x'_i \\ y'_i \end{pmatrix} \begin{pmatrix} x'_i \\ y'_i \end{pmatrix}^\top = \begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix}$$

Subtraction of mean vector is not necessary (because it is $(0, 0)^\top$).

Diagonal elements are always 1 (because of unit variance in each dimension).

- Normalizing the data and then computing the covariances or computing the covariances and then normalizing them has the same effect.

Correlation Matrix: Interpretation

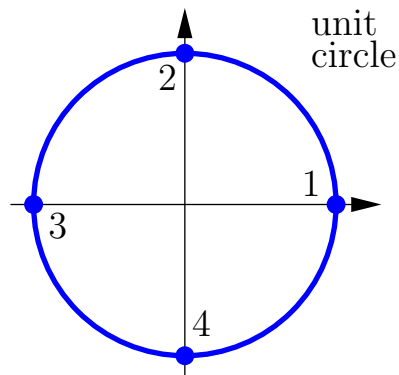
Special Case: Two Dimensions

- Correlation matrix

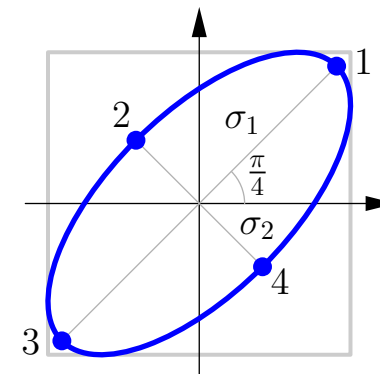
$$\Sigma' = \begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix}, \quad \begin{array}{l} \text{eigenvalues: } \sigma_1^2, \sigma_2^2 \\ \text{correlation: } r = \frac{\sigma_1^2 - \sigma_2^2}{\sigma_1^2 + \sigma_2^2} \end{array}$$

- Eigenvalue decomposition

$$\mathbf{T} = \begin{pmatrix} c & -s \\ s & c \end{pmatrix} \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}, \quad \begin{array}{l} s = \sin \frac{\pi}{4} = \frac{1}{\sqrt{2}}, \quad \sigma_1 = \sqrt{1+r}, \\ c = \cos \frac{\pi}{4} = \frac{1}{\sqrt{2}}, \quad \sigma_2 = \sqrt{1-r}. \end{array}$$



mapping with \mathbf{T}
 $\vec{v}' = \mathbf{T}\vec{v}$



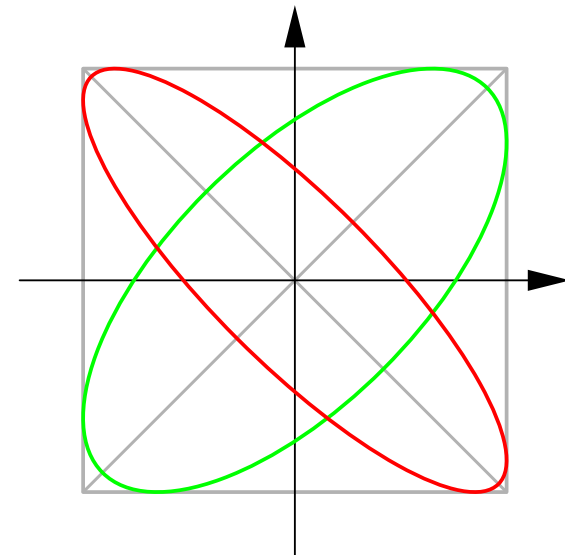
Correlation Matrix: Interpretation

- For two dimensions the eigenvectors of a correlation matrix are always

$$\vec{v}_1 = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right) \quad \text{and} \quad \vec{v}_2 = \left(-\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right)$$

(or their opposites $-\vec{v}_1$ or $-\vec{v}_2$ or exchanged).

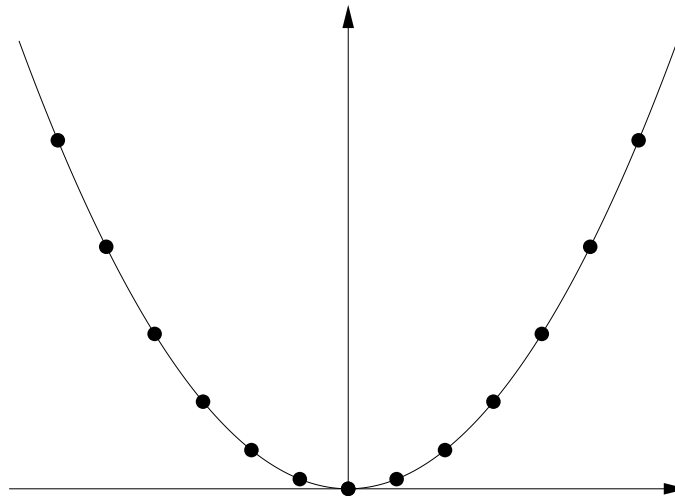
The reason is that the normalization transforms the data points in such a way, that the ellipse, the unit circle is mapped to by the “square root” of the covariance matrix of the normalized data, is always inscribed into the square $[-1, 1] \times [-1, 1]$. Hence the ellipse’s major axes are the square’s diagonals.



- The situation is analogous in m -dimensional spaces: the eigenvectors are always m of the 2^{m-1} diagonals of the m -dimensional unit (hyper-)cube around the origin.

Correlation and Stochastic (In)Dependence

- Note: stochastic independence $\Rightarrow r = 0$,
but: $r = 0 \not\Rightarrow$ stochastic independence.
- Example: Suppose the data points lie symmetrically on a parabola.



- The correlation coefficient of this data set is $r = 0$,
because there is **no linear** dependence between the two attributes.
However, there is a perfect **quadratic** dependence,
and thus the two attributes are **not** stochastically independent.

Regression Line

- Since the covariance/correlation measures linear dependence, it is not surprising that it can be used to define a **regression line**:

$$(y - \bar{y}) = \frac{s_{xy}}{s_x^2}(x - \bar{x}) \quad \text{or} \quad y = \frac{s_{xy}}{s_x^2}(x - \bar{x}) + \bar{y}.$$

- The regression line can be seen as a conditional arithmetic mean: there is one arithmetic mean for the y -dimensions for each x -value.
- This interpretation is supported by the fact that the regression line minimizes the sum of squared differences in y -direction.
(Reminder: the arithmetic mean minimizes the sum of squared differences.)
- More information on **regression** and the **method of least squares** in the corresponding chapter.

Principal Component Analysis

- Correlations between the attributes of a data set can be used to **reduce the number of dimensions**:
 - Of two strongly correlated features only one needs to be considered.
 - The other can be reconstructed approximately from the regression line.
 - However, the feature selection can be difficult.
- Better approach: **Principal Component Analysis** (PCA)
 - Find the direction in the data space that has the highest variance.
 - Find the direction in the data space that has the highest variance among those perpendicular to the first.
 - Find the direction in the data space that has the highest variance among those perpendicular to the first and second and so on.
 - Use first directions to describe the data.

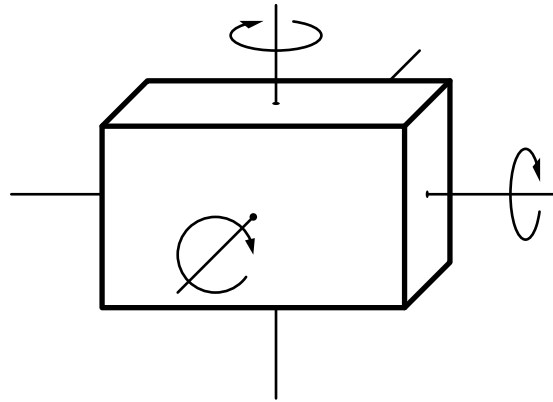
Principal Component Analysis: Physical Analog

- The rotation of a body around an axis through its center of gravity can be described by a so-called **inertia tensor**, which is a 3×3 -matrix

$$\Theta = \begin{pmatrix} \Theta_{xx} & \Theta_{xy} & \Theta_{xz} \\ \Theta_{xy} & \Theta_{yy} & \Theta_{yz} \\ \Theta_{xz} & \Theta_{yz} & \Theta_{zz} \end{pmatrix}.$$

- The diagonal elements of this tensor are called the **moments of inertia**. They describe the “resistance” of the body against being rotated.
- The off-diagonal elements are the so-called **deviation moments**. They describe forces vertical to the rotation axis.
- All bodies possess three perpendicular axes through their center of gravity, around which they can be rotated without forces perpendicular to the rotation axis. These axes are called **principal axes of inertia**.
There are bodies that possess more than 3 such axes (example: a homogeneous sphere), but all bodies have at least three such axes.

Principal Component Analysis: Physical Analog



The principal axes
of inertia of a box.

- The deviation moments cause “rattling” in the bearings of the rotation axis, which cause the bearings to wear out quickly.
- A car mechanic who balances a wheel carries out, in a way, a principal axes transformation. However, instead of changing the orientation of the axes, he/she adds small weights to minimize the deviation moments.
- A statistician who does a principal component analysis, finds, in a way, the axes through a weight distribution with unit weights at each data point, around which it can be rotated most easily.

Principal Component Analysis: Formal Approach

- Normalize all attributes to arithmetic mean 0 and standard deviation 1:

$$x' = \frac{x - \bar{x}}{s_x}$$

- Compute the **correlation matrix** Σ
(i.e., the covariance matrix of the normalized data)
- Carry out a **principal axes transformation** of the correlation matrix, that is, find a matrix \mathbf{R} , such that $\mathbf{R}^\top \Sigma \mathbf{R}$ is a diagonal matrix.
- Formal procedure:
 - Find the **eigenvalues** and **eigenvectors** of the correlation matrix, i.e., find the values λ_i and vectors \vec{v}_i , such that $\Sigma \vec{v}_i = \lambda_i \vec{v}_i$.
 - The eigenvectors indicate the desired directions.
 - The eigenvalues are the variances in these directions.

Principal Component Analysis: Formal Approach

- Select dimensions using the **percentage of explained variance**.
 - The eigenvalues λ_i are the variances σ_i^2 in the principal dimensions.
 - It can be shown that the sum of the eigenvalues of an $m \times m$ correlation matrix is m . Therefore it is plausible to define $\frac{\lambda_i}{m}$ as the share the i -th principal axis has in the total variance.
 - Sort the λ_i descendingly and find the smallest value k , such that

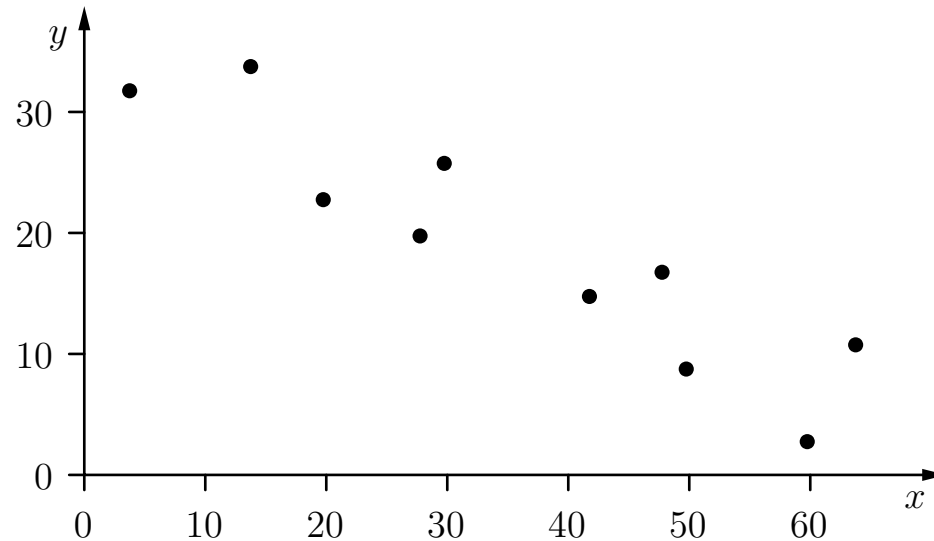
$$\sum_{i=1}^k \frac{\lambda_i}{m} \geq \alpha,$$

where α is a user-defined parameter (e.g. $\alpha = 0.9$).

- Select the corresponding k directions (given by the eigenvectors).
- Transform the data to the new data space by multiplying the data points with a matrix, the rows of which are the eigenvectors of the selected dimensions.

Principal Component Analysis: Example

x	5	15	21	29	31	43	49	51	61	65
y	33	35	24	21	27	16	18	10	4	12



- Strongly correlated features \Rightarrow Reduction to one dimension possible.

Principal Component Analysis: Example

Normalize to arithmetic mean 0 and standard deviation 1:

$$\bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i = \frac{370}{10} = 37,$$

$$\bar{y} = \frac{1}{10} \sum_{i=1}^{10} y_i = \frac{200}{10} = 20,$$

$$s_x^2 = \frac{1}{9} \left(\sum_{i=1}^{10} x_i^2 - 10\bar{x}^2 \right) = \frac{17290 - 13690}{9} = 400 \Rightarrow s_x = 20,$$

$$s_y^2 = \frac{1}{9} \left(\sum_{i=1}^{10} y_i^2 - 10\bar{y}^2 \right) = \frac{4900 - 4000}{9} = 100 \Rightarrow s_y = 10.$$

x'	-1.6	-1.1	-0.8	-0.4	-0.3	0.3	0.6	0.7	1.2	1.4
y'	1.3	1.5	0.4	0.1	0.7	-0.4	-0.2	-1.0	-1.6	-0.8

Principal Component Analysis: Example

- Compute the correlation matrix (covariance matrix of normalized data).

$$\Sigma = \frac{1}{9} \begin{pmatrix} 9 & -8.28 \\ -8.28 & 9 \end{pmatrix} = \begin{pmatrix} 1 & -\frac{23}{25} \\ -\frac{23}{25} & 1 \end{pmatrix}.$$

- Find the eigenvalues and eigenvectors, i.e., the values λ_i and vectors \vec{v}_i , $i = 1, 2$, such that

$$\Sigma \vec{v}_i = \lambda_i \vec{v}_i \quad \text{or} \quad (\Sigma - \lambda_i \mathbf{1}) \vec{v}_i = \vec{0}.$$

where $\mathbf{1}$ is the unit matrix.

- Here: Find the eigenvalues as the roots of the characteristic polynomial.

$$c(\lambda) = |\Sigma - \lambda \mathbf{1}| = (1 - \lambda)^2 - \frac{529}{625}.$$

For more than 3 dimensions, this method is numerically unstable and should be replaced by some other method (Jacobi-Transformation, Householder Transformation to tridiagonal form followed by the QR algorithm etc.).

Principal Component Analysis: Example

- The roots of the characteristic polynomial $c(\lambda) = (1 - \lambda)^2 - \frac{529}{625}$ are

$$\lambda_{1/2} = 1 \pm \sqrt{\frac{529}{625}} = 1 \pm \frac{23}{25}, \quad \text{i.e.} \quad \lambda_1 = \frac{48}{25} \quad \text{and} \quad \lambda_2 = \frac{2}{25}$$

- The corresponding eigenvectors are determined by solving for $i = 1, 2$ the (underdetermined) linear equation system

$$(\mathbf{\Sigma} - \lambda_i \mathbf{1}) \vec{v}_i = \vec{0}$$

- The resulting eigenvectors (normalized to length 1) are

$$\vec{v}_1 = \left(\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right) \quad \text{and} \quad \vec{v}_2 = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right),$$

(Note that for two dimensions always these two vectors result.

Reminder: directions of the eigenvectors of a correlation matrix.)

Principal Component Analysis: Example

- Therefore the transformation matrix for the principal axes transformation is

$$\mathbf{R} = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}, \quad \text{for which it is} \quad \mathbf{R}^\top \mathbf{\Sigma} \mathbf{R} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$

- However, instead of \mathbf{R}^\top we use $\sqrt{2}\mathbf{R}^\top$ to transform the data:

$$\begin{pmatrix} x'' \\ y'' \end{pmatrix} = \sqrt{2} \cdot \mathbf{R}^\top \cdot \begin{pmatrix} x' \\ y' \end{pmatrix}.$$

Resulting data set:

x''	-2.9	-2.6	-1.2	-0.5	-1.0	0.7	0.8	1.7	2.8	2.2
y''	-0.3	0.4	-0.4	-0.3	0.4	-0.1	0.4	-0.3	-0.4	0.6

- y'' is discarded ($s_{y''}^2 = 2\lambda_2 = \frac{4}{25}$) and only x'' is kept ($s_{x''}^2 = 2\lambda_1 = \frac{96}{25}$).

Inductive Statistics

Inductive Statistics: Main Tasks

- **Parameter Estimation**

Given an assumption about the type of distribution of the underlying random variable the parameter(s) of the distribution function is estimated.

- **Hypothesis Testing**

A hypothesis about the data generating process is tested by means of the data.

- *Parameter Test*

Test whether a parameter can have certain values.

- *Goodness-of-Fit Test*

Test whether a distribution assumption fits the data.

- *Dependence Test*

Test whether two attributes are dependent.

- **Model Selection**

Among different models that can be used to explain the data the best fitting is selected, taking the complexity of the model into account.

Inductive Statistics: Random Samples

- In inductive statistics probability theory is applied to make inferences about the process that generated the data. This presupposes that the sample is the result of a random experiment, a so-called **random sample**.
- The random variable yielding the sample value x_i is denoted X_i . x_i is called a **instantiation** of the random variable X_i .
- A random sample $x = (x_1, \dots, x_n)$ is an instantiation of the **random vector** $X = (X_1, \dots, X_n)$.
- A random sample is called **independent** if the random variables X_1, \dots, X_n are (stochastically) independent, i.e. if

$$\forall c_1, \dots, c_n \in \mathbb{R} : \quad P \left(\bigwedge_{i=1}^n X_i \leq c_i \right) = \prod_{i=1}^n P(X_i \leq c_i).$$

- An independent random sample is called **simple** if the random variables X_1, \dots, X_n have the same distribution function.

Inductive Statistics: Parameter Estimation

Parameter Estimation

Given:

- A data set and
- a family of parameterized distributions functions of the same type, e.g.
 - the family of binomial distributions $b_X(x; p, n)$ with the parameters p , $0 \leq p \leq 1$, and $n \in \mathbb{N}$, where n is the sample size,
 - the family of normal distributions $N_X(x; \mu, \sigma^2)$ with the parameters μ (expected value) and σ^2 (variance).

Assumption:

- The process that generated the data can be described well by an element of the given family of distribution functions.

Desired:

- The element of the given family of distribution functions (determined by its parameters) that is the best model for the data.

Parameter Estimation

- Methods that yield an estimate for a parameter are called **estimators**.
- Estimators are **statistics**, i.e. functions of the values in a sample.

As a consequence they are functions of (instantiations of) random variables and thus (instantiations of) random variables themselves.

Therefore we can use all of probability theory to analyze estimators.

- There are two types of parameter estimation:
 - **Point Estimators**
Point estimators determine the best value of a parameter w.r.t. the data and certain quality criteria.
 - **Interval Estimators**
Interval estimators yield a region, a so-called **confidence interval**, in which the true value of the parameter lies with high certainty.

Inductive Statistics:

Point Estimation

Point Estimation

Not all statistics, that is, not all functions of the sample values are reasonable and useful estimator. Desirable properties are:

- **Consistency**

With growing data volume the estimated value should get closer and closer to the true value, at least with higher and higher probability.

Formally: If T is an estimator for the parameter θ , it should be

$$\forall \varepsilon > 0 : \lim_{n \rightarrow \infty} P(|T - \theta| < \varepsilon) = 1,$$

where n is the sample size.

- **Unbiasedness**

An estimator should not tend to over- or underestimate the parameter.

Rather it should yield, on average, the correct value.

Formally this means

$$E(T) = \theta.$$

Point Estimation

- **Efficiency**

The estimation should be as precise as possible, that is, the deviation from the true value should be as small as possible. Formally: If T and U are two estimators for the same parameter θ , then T is called *more efficient* than U if

$$D^2(T) < D^2(U).$$

- **Sufficiency**

An estimator should exploit all information about the parameter contained in the data. More precisely: two samples that yield the same estimate should have the same probability (otherwise there is unused information).

Formally: an estimator T for a parameter θ is called sufficient iff for all samples $x = (x_1, \dots, x_n)$ with $T(x) = t$ the expression

$$\frac{f_{X_1}(x_1; \theta) \cdots f_{X_n}(x_n; \theta)}{f_T(t; \theta)}$$

is independent of θ .

Point Estimation: Example

Given: a family of **uniform distributions** on the interval $[0, \theta]$, i.e.

$$f_X(x; \theta) = \begin{cases} \frac{1}{\theta}, & \text{if } 0 \leq x \leq \theta, \\ 0, & \text{otherwise.} \end{cases}$$

Desired: an estimate for the unknown parameter θ .

- We will now consider two estimators for the parameter θ and compare their properties.
 - $T = \max\{X_1, \dots, X_n\}$
 - $U = \frac{n+1}{n} \max\{X_1, \dots, X_n\}$
- **General approach:**
 - Find the probability density function of the estimator.
 - Check the desirable properties by exploiting this density function.

Point Estimation: Example

To analyze the estimator $T = \max\{X_1, \dots, X_n\}$, we compute its density function:

$$\begin{aligned} f_T(t; \theta) &= \frac{d}{dt} F_T(t; \theta) = \frac{d}{dt} P(T \leq t) \\ &= \frac{d}{dt} P(\max\{X_1, \dots, X_n\} \leq t) \\ &= \frac{d}{dt} P\left(\bigwedge_{i=1}^n X_i \leq t\right) = \frac{d}{dt} \prod_{i=1}^n P(X_i \leq t) \\ &= \frac{d}{dt} (F_X(t; \theta))^n = n \cdot (F_X(t; \theta))^{n-1} f_X(t, \theta) \end{aligned}$$

where

$$F_X(x; \theta) = \int_{-\infty}^x f_X(x; \theta) dx = \begin{cases} 0, & \text{if } x \leq 0, \\ \frac{x}{\theta}, & \text{if } 0 \leq x \leq \theta, \\ 1, & \text{if } x \geq \theta. \end{cases}$$

Therefore it is

$$f_T(t; \theta) = \frac{n \cdot t^{n-1}}{\theta^n} \quad \text{for } 0 \leq t \leq \theta, \quad \text{and } 0 \text{ otherwise.}$$

Point Estimation: Example

- The estimator $T = \max\{X_1, \dots, X_n\}$ is **consistent**:

$$\begin{aligned}\lim_{n \rightarrow \infty} P(|T - \theta| < \epsilon) &= \lim_{n \rightarrow \infty} P(T > \theta - \epsilon) \\ &= \lim_{n \rightarrow \infty} \int_{\theta - \epsilon}^{\theta} \frac{n \cdot t^{n-1}}{\theta^n} dt = \lim_{n \rightarrow \infty} \left[\frac{t^n}{\theta^n} \right]_{\theta - \epsilon}^{\theta} \\ &= \lim_{n \rightarrow \infty} \left(\frac{\theta^n}{\theta^n} - \frac{(\theta - \epsilon)^n}{\theta^n} \right) \\ &= \lim_{n \rightarrow \infty} \left(1 - \left(\frac{\theta - \epsilon}{\theta} \right)^n \right) = 1\end{aligned}$$

- It is **not unbiased**:

$$\begin{aligned}E(T) &= \int_{-\infty}^{\infty} t \cdot f_T(t; \theta) dt = \int_0^{\theta} t \cdot \frac{n \cdot t^{n-1}}{\theta^n} dt \\ &= \left[\frac{n \cdot t^{n+1}}{(n+1)\theta^n} \right]_0^{\theta} = \frac{n}{n+1} \theta < \theta \quad \text{for } n < \infty.\end{aligned}$$

Point Estimation: Example

- The estimator $U = \frac{n+1}{n} \max\{X_1, \dots, X_n\}$ has the density function

$$f_U(u; \theta) = \frac{n^{n+1}}{(n+1)^n} \frac{u^{n-1}}{\theta^n} \quad \text{for } 0 \leq u \leq \frac{n+1}{n}\theta, \text{ and } 0 \text{ otherwise.}$$

- The estimator U is **consistent** (without formal proof).
- It is **unbiased**:

$$\begin{aligned} E(U) &= \int_{-\infty}^{\infty} u \cdot f_U(u; \theta) \, du \\ &= \int_0^{\frac{n+1}{n}\theta} u \cdot \frac{n^{n+1}}{(n+1)^n} \frac{u^{n-1}}{\theta^n} \, du \\ &= \frac{n^{n+1}}{(n+1)^n \theta^n} \left[\frac{u^{n+1}}{n+1} \right]_0^{\frac{n+1}{n}\theta} \\ &= \frac{n^{n+1}}{(n+1)^n \theta^n} \cdot \frac{1}{n+1} \left(\frac{n+1}{n} \theta \right)^{n+1} = \theta \end{aligned}$$

Point Estimation: Example

Given: a family of **normal distributions** $N_X(x; \mu, \sigma^2)$

$$f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Desired: estimates for the unknown parameters μ and σ^2 .

- The median and the arithmetic mean of the sample are both consistent and unbiased estimators for the parameter μ .
The median is less efficient than the arithmetic mean.
- The function $V^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is a consistent, but **biased** estimator for the parameter σ^2 (it tends to underestimate the variance).
The function $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$, however, is a consistent and **unbiased** estimator for σ^2 (this explains the definition of the empirical variance).

Point Estimation: Example

Given: a family of **polynomial distributions**

$$f_{X_1, \dots, X_k}(x_1, \dots, x_k; \theta_1, \dots, \theta_k, n) = \frac{n!}{\prod_{i=1}^k x_i!} \prod_{i=1}^k \theta_i^{x_i},$$

(n is the sample size, the x_i are the frequencies of the different values a_i , $i = 1, \dots, k$, and the θ_i are the probabilities with which the values a_i occur.)

Desired: estimates for the unknown parameters $\theta_1, \dots, \theta_k$

- The relative frequencies $R_i = \frac{X_i}{n}$ of the different values a_i , $i = 1, \dots, k$, are
 - consistent,
 - unbiased,
 - most efficient, and
 - sufficient estimators for the θ_i .

Inductive Statistics:

Finding Point Estimators

How Can We Find Estimators?

- Up to now we analyzed given estimators, now we consider the question how to find them.
- There are three main approaches to find estimators:
 - **Method of Moments**
Derive an estimator for a parameter from the moments of a distribution and its generator function.
(We do not consider this method here.)
 - **Maximum Likelihood Estimation**
Choose the (set of) parameter value(s) that makes the sample most likely.
 - **Maximum A-posteriori Estimation**
Choose a prior distribution on the range of parameter values, apply Bayes' rule to compute the posterior probability from the sample, and choose the (set of) parameter value(s) that maximizes this probability.

Maximum Likelihood Estimation

- General idea: **Choose the (set of) parameter value(s) that makes the sample most likely.**
- If the parameter value(s) were known, it would be possible to compute the probability of the sample. With unknown parameter value(s), however, it is still possible to state this probability as a function of the parameter(s).
- Formally this can be described as choosing the value θ that maximizes

$$L(D; \theta) = f(D | \theta),$$

where D are the sample data and L is called the **Likelihood Function**.

- Technically the estimator is determined by
 - setting up the likelihood function,
 - forming its partial derivative(s) w.r.t. the parameter(s), and
 - setting these derivatives equal to zero (necessary condition for a maximum).

Brief Excursion: Function Optimization

Task: Find values $\vec{x} = (x_1, \dots, x_m)$ such that $f(\vec{x}) = f(x_1, \dots, x_m)$ is optimal.

Often feasible approach:

- A necessary condition for a (local) optimum (maximum or minimum) is that the partial derivatives w.r.t. the parameters vanish (Pierre Fermat).
- Therefore: (Try to) solve the equation system that results from setting all partial derivatives w.r.t. the parameters equal to zero.

Example task: Minimize $f(x, y) = x^2 + y^2 + xy - 4x - 5y$.

Solution procedure:

1. Take the partial derivatives of the objective function and set them to zero:

$$\frac{\partial f}{\partial x} = 2x + y - 4 = 0, \quad \frac{\partial f}{\partial y} = 2y + x - 5 = 0.$$

2. Solve the resulting (here: linear) equation system: $x = 1, \quad y = 2$.

Maximum Likelihood Estimation: Example

Given: a family of **normal distributions** $N_X(x; \mu, \sigma^2)$

$$f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Desired: estimators for the unknown parameters μ and σ^2 .

The **Likelihood Function**, which describes the probability of the data, is

$$L(x_1, \dots, x_n; \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right).$$

To simplify the technical task of forming the partial derivatives, we consider the natural logarithm of the likelihood function, i.e.

$$\ln L(x_1, \dots, x_n; \mu, \sigma^2) = -n \ln\left(\sqrt{2\pi\sigma^2}\right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.$$

Maximum Likelihood Estimation: Example

- Estimator for the **expected value** μ :

$$\frac{\partial}{\partial \mu} \ln L(x_1, \dots, x_n; \mu, \sigma^2) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \stackrel{!}{=} 0$$

$$\Rightarrow \sum_{i=1}^n (x_i - \mu) = \left(\sum_{i=1}^n x_i \right) - n\mu \stackrel{!}{=} 0 \quad \Rightarrow \quad \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Estimator for the **variance** σ^2 :

$$\frac{\partial}{\partial \sigma^2} \ln L(x_1, \dots, x_n; \mu, \sigma^2) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 \stackrel{!}{=} 0$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n^2} \left(\sum_{i=1}^n x_i \right)^2 \quad (\text{biased!})$$

Maximum A-posteriori Estimation: Motivation

Consider the following three situations:

- A drunkard claims to be able to predict the side on which a thrown coin will land (head or tails). On ten trials he always states the correct side beforehand.
- A tea lover claims that she is able to taste whether the tea or the milk was poured into the cup first. On ten trials she always identifies the correct order.
- An expert of classical music claims to be able to recognize from a single sheet of music whether the composer was Mozart or somebody else. On ten trials he is indeed correct every time.

Maximum likelihood estimation treats all situations alike, because formally the samples are the same. However, this is implausible:

- We do not believe the drunkard at all, despite the sample data.
- We highly doubt the tea drinker, but tend to consider the data as evidence.
- We tend to believe the music expert easily.

Maximum A-posteriori Estimation

- Background knowledge about the plausible values can be incorporated by
 - using a **prior distribution** on the domain of the parameter and
 - adapting this distribution with **Bayes' rule** and the data.
- Formally maximum a-posteriori estimation is defined as follows:
find the parameter value θ that maximizes

$$f(\theta \mid D) = \frac{f(D \mid \theta)f(\theta)}{f(D)} = \frac{f(D \mid \theta)f(\theta)}{\int_{-\infty}^{\infty} f(D \mid \theta)f(\theta) d\theta}$$

- As a comparison: maximum likelihood estimation maximizes

$$f(D \mid \theta)$$

- Note that $f(D)$ need not be computed: It is the same for all parameter values and since we are only interested in the value θ that maximizes $f(\theta \mid D)$ and not the *value of* $f(\theta \mid D)$, we can treat it as a normalization constant.

Maximum A-posteriori Estimation: Example

Given: a family of **binomial distributions**

$$f_X(x; \theta, n) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}.$$

Desired: an estimator for the unknown parameter θ .

a) **Uniform prior:** $f(\theta) = 1, \quad 0 \leq \theta \leq 1.$

$$f(\theta \mid D) = \gamma \binom{n}{x} \theta^x (1 - \theta)^{n-x} \cdot 1 \quad \Rightarrow \quad \hat{\theta} = \frac{x}{n}$$

b) **Tendency towards $\frac{1}{2}$:** $f(\theta) = 6\theta(1 - \theta), \quad 0 \leq \theta \leq 1.$

$$\begin{aligned} f(\theta \mid D) &= \gamma \binom{n}{x} \theta^x (1 - \theta)^{n-x} \cdot \theta(1 - \theta) = \gamma \binom{n}{x} \theta^{x+1} (1 - \theta)^{n-x+1} \\ &\Rightarrow \quad \hat{\theta} = \frac{x+1}{n+2} \end{aligned}$$

Excursion: Dirichlet's Integral

- For computing the normalization factors of the probability density functions that occur with polynomial distributions, **Dirichlet's Integral** is helpful:

$$\int_{\theta_1} \cdots \int_{\theta_k} \prod_{i=1}^k \theta_i^{x_i} d\theta_1 \cdots d\theta_k = \frac{\prod_{i=1}^k \Gamma(x_i + 1)}{\Gamma(n + k)}, \quad \text{where } n = \sum_{i=1}^k x_i$$

and the Γ -function is the so-called **generalized factorial**:

$$\Gamma(x) = \int_0^\infty e^{-t} t^{x-1} dt, \quad x > 0,$$

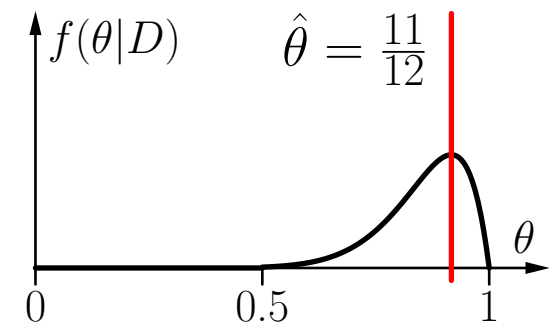
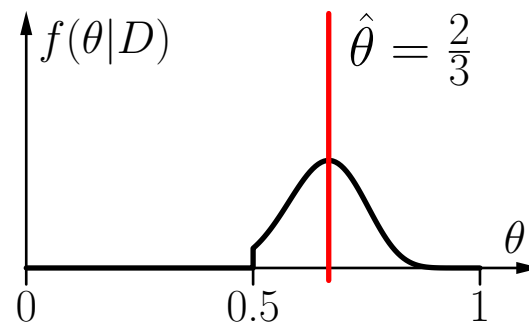
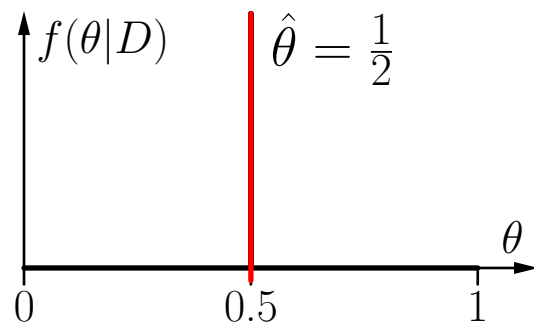
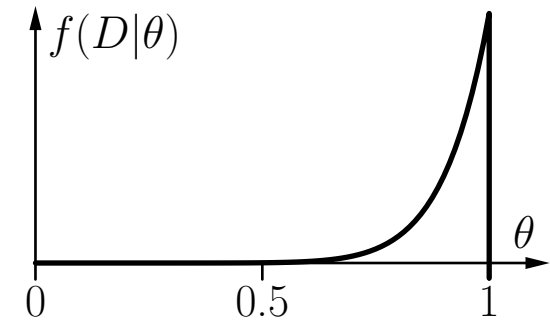
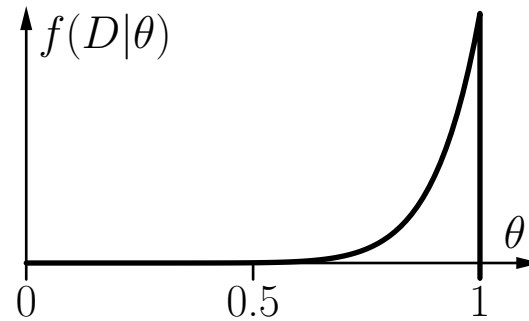
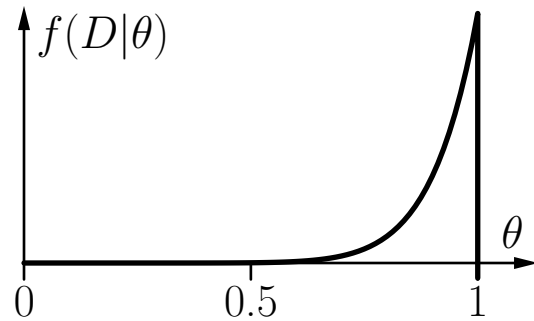
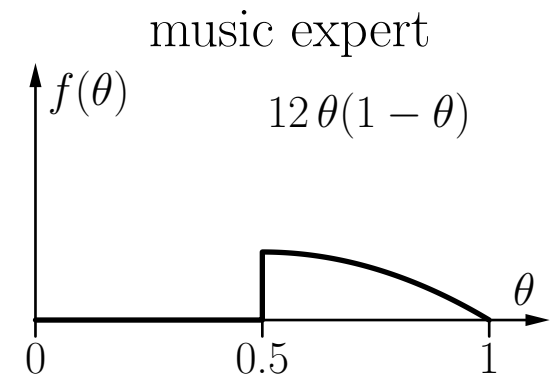
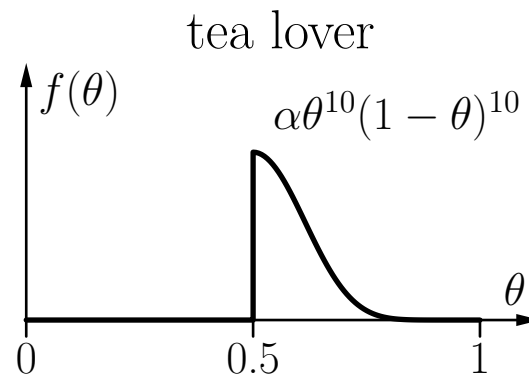
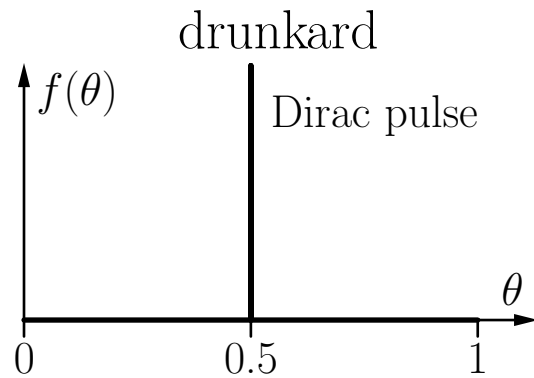
which satisfies

$$\Gamma(x + 1) = x \cdot \Gamma(x), \quad \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(1) = 1.$$

- Example:** the normalization factor α for the binomial distribution prior $f(\theta) = \alpha \theta^2(1 - \theta)^3$ is

$$\alpha = \frac{1}{\int_{\theta} \theta^2(1 - \theta)^3 d\theta} = \frac{\Gamma(5 + 2)}{\Gamma(2 + 1) \Gamma(3 + 1)} = \frac{6!}{2! 3!} = \frac{720}{12} = 60.$$

Maximum A-posteriori Estimation: Example



Inductive Statistics:

Interval Estimation

Interval Estimation

- In general the estimated value of a parameter will differ from the true value.
- It is desirable to be able to make an assertion about the possible deviations.
- The simplest possibility is to state not only a point estimate, but also the standard deviation of the estimator:

$$t \pm D(T) = t \pm \sqrt{D^2(T)}.$$

- A better possibility is to find intervals that contain the true value with high probability. Formally they can be defined as follows:

Let $A = g_A(X_1, \dots, X_n)$ and $B = g_B(X_1, \dots, X_n)$ be two statistics with

$$P(A < \theta < B) = 1 - \alpha, \quad P(\theta \leq A) = \frac{\alpha}{2}, \quad P(\theta \geq B) = \frac{\alpha}{2}.$$

Then the random interval $[A, B]$ (or an instantiation $[a, b]$ of this interval) is called $(1 - \alpha) \cdot 100\%$ **confidence interval** for θ . The value $1 - \alpha$ is called **confidence level**.

Interval Estimation

- This definition of a confidence interval is not specific enough:
 A and B are not uniquely determined.
- Common solution: Start from a point estimator T for the unknown parameter θ and define A and B as functions of T :

$$A = h_A(T) \quad \text{and} \quad B = h_B(T).$$

- Instead of $A \leq \theta \leq B$ consider the corresponding event w.r.t. the estimator T , that is, $A^* \leq T \leq B^*$.
- Determine $A = h_A(T)$ and $B = h_B(T)$ from the inverse functions $A^* = h_A^{-1}(\theta)$ and $B^* = h_B^{-1}(\theta)$.

$$\begin{aligned} \text{Procedure: } P(A^* < T < B^*) &= 1 - \alpha \\ \Rightarrow P(h_A^{-1}(\theta) < T < h_B^{-1}(\theta)) &= 1 - \alpha \\ \Rightarrow P(h_A(T) < \theta < h_B(T)) &= 1 - \alpha \\ \Rightarrow P(A < \theta < B) &= 1 - \alpha. \end{aligned}$$

Interval Estimation: Example

Given: a family of **uniform distributions** on the interval $[0, \theta]$, i.e.

$$f_X(x; \theta) = \begin{cases} \frac{1}{\theta}, & \text{if } 0 \leq x \leq \theta, \\ 0, & \text{otherwise.} \end{cases}$$

Desired: a confidence interval for the unknown parameter θ .

- Start from the unbiased point estimator $U = \frac{n+1}{n} \max\{X_1, \dots, X_n\}$:

$$P(U \leq B^*) = \int_0^{B^*} f_U(u; \theta) \, du = \frac{\alpha}{2}$$

$$P(U \geq A^*) = \int_{A^*}^{\frac{n+1}{n}\theta} f_U(u; \theta) \, du = \frac{\alpha}{2}$$

- From the study of point estimators we know

$$f_U(u; \theta) = \frac{n^{n+1}}{(n+1)^n} \frac{u^{n-1}}{\theta^n}.$$

Interval Estimation: Example

- Solving the integrals gives us

$$B^* = \sqrt[n]{\frac{\alpha}{2}} \frac{n+1}{n} \theta \quad \text{and} \quad A^* = \sqrt[n]{1 - \frac{\alpha}{2}} \frac{n+1}{n} \theta,$$

that is,

$$P \left(\sqrt[n]{\frac{\alpha}{2}} \frac{n+1}{n} \theta < U < \sqrt[n]{1 - \frac{\alpha}{2}} \frac{n+1}{n} \theta \right) = 1 - \alpha.$$

- Computing the inverse functions leads to

$$P \left(\frac{U}{\sqrt[n]{1 - \frac{\alpha}{2}} \frac{n+1}{n}} < \theta < \frac{U}{\sqrt[n]{\frac{\alpha}{2}} \frac{n+1}{n}} \right) = 1 - \alpha,$$

that is,

$$A = \frac{U}{\sqrt[n]{1 - \frac{\alpha}{2}} \frac{n+1}{n}} \quad \text{and} \quad B = \frac{U}{\sqrt[n]{\frac{\alpha}{2}} \frac{n+1}{n}}.$$

Inductive Statistics: Hypothesis Testing

Hypothesis Testing

- A **hypothesis test** is a statistical procedure with which a decision is made between two contrary hypothesis about the process that generated the data.
- The two hypotheses can refer to
 - the value of a parameter (**Parameter Test**),
 - a distribution assumption (**Goodness-of-Fit Test**),
 - the dependence of two attributes (**Dependence Test**).
- One of the two hypothesis is preferred, that is, in case of doubt the decision is made in its favor. (One says that it gets the “benefit of the doubt”.)
- The preferred hypothesis is called the **Null Hypothesis** H_0 , the other hypothesis is called the **Alternative Hypothesis** H_a .
- Intuitively: the null hypothesis H_0 is put on trial.
Only if the evidence is strong enough, it is convicted (i.e. rejected).
If there is doubt, however, it is acquitted (i.e. accepted).

Hypothesis Testing

- The test decision is based on a **test statistic**, that is, a function of the sample values.
- The null hypothesis is rejected if the value of the test statistic lies inside the so-called **critical region** C .
- Developing a hypothesis test consists in finding the critical region for a given test statistic and significance level (see below).
- The test decision may be wrong. There are two possible types of errors:
 - Type 1:** The null hypothesis H_0 is rejected, even though it is correct.
 - Type 2:** The null hypothesis H_0 is accepted, even though it is false.
- Type 1 errors are considered to be more severe, since the null hypothesis gets “the benefit of the doubt”.
- Therefore it is tried to restrict the probability of a type 1 error to a certain maximum α . This maximum value α is called **significance level**.

Parameter Test

- In a parameter test the contrary hypotheses refer to the value of a parameter, for example (one-sided test):

$$H_0 : \theta \geq \theta_0, \quad H_a : \theta < \theta_0.$$

- For such a test usually a point estimator T is chosen as the test statistic.
- The null hypothesis H_0 is rejected if the value t of the point estimator does not exceed a certain value c , the so-called **critical value** (i.e. $C = (-\infty, c]$).
- Formally the critical value c is determined as follows: We consider

$$\beta(\theta) = P_\theta(H_0 \text{ is rejected}) = P_\theta(T \in C),$$

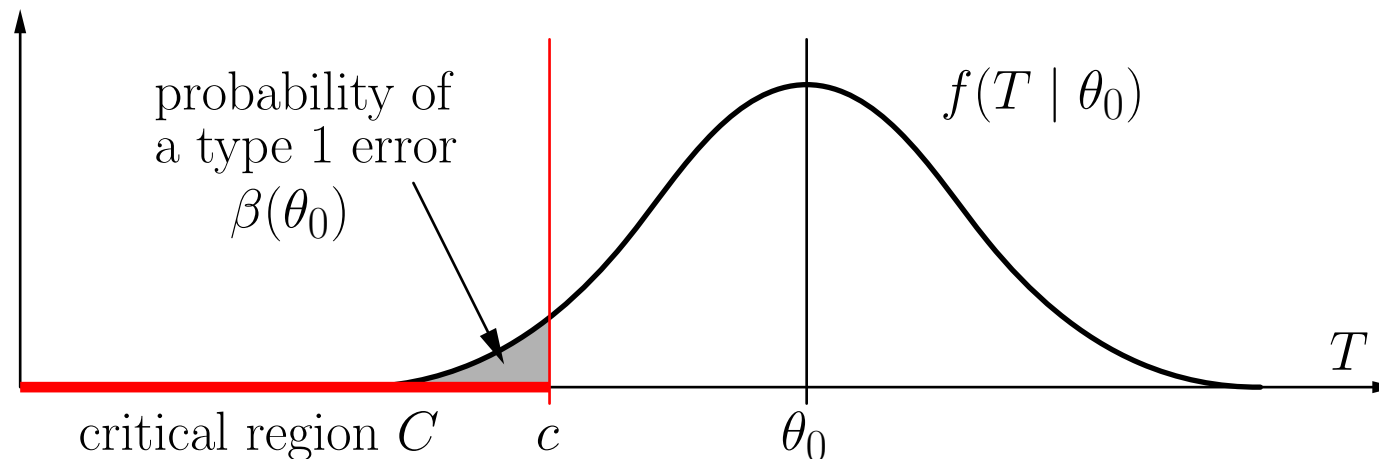
the so-called **power** β of the test.

- The power must be not exceed the significance level α for values θ satisfying H_0 :

$$\max_{\theta: \theta \text{ satisfies } H_0} \beta(\theta) \leq \alpha. \quad (\text{here: } \beta(\theta_0) \leq \alpha)$$

Parameter Test: Intuition

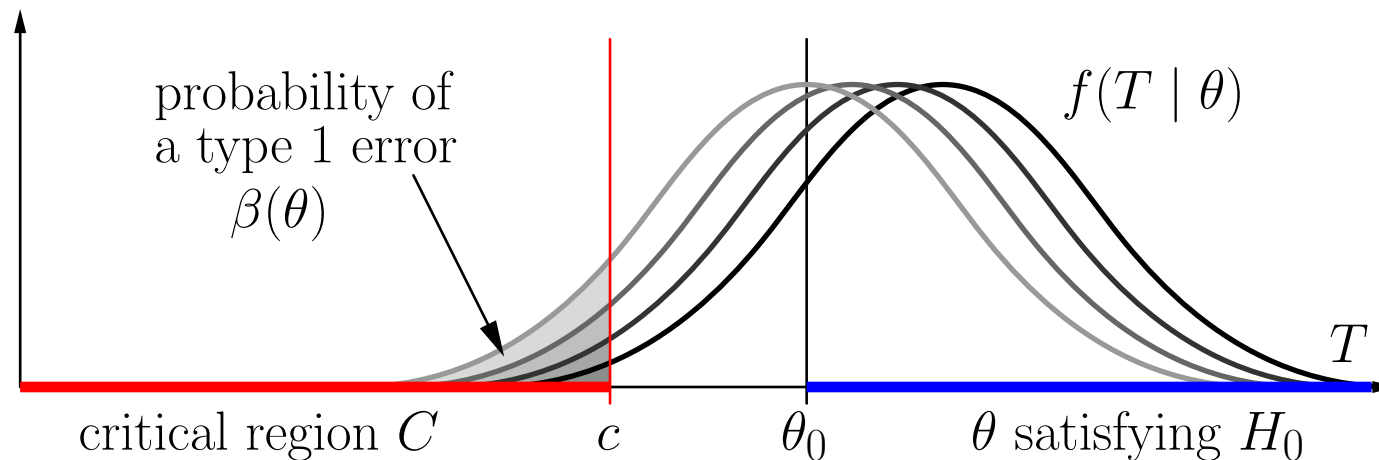
- The probability of a type 1 error is the area under the estimator's probability density function $f(T | \theta_0)$ to the left of the critical value c . (Note: This example illustrates $H_0 : \theta \geq \theta_0$ and $H_a : \theta < \theta_0$.)



- Obviously the probability of a type 1 error depends on the location of the critical value c : higher values mean a higher error probability.
- Idea: Choose the location of the critical value so that the maximal probability of a type 1 error equals α , the chosen significance level.

Parameter Test: Intuition

- What is so special about θ_0 that we use $f(T | \theta_0)$?



- In principle, all θ satisfying H_0 have to be considered, that is, all density functions $f(T | \theta)$ with $\theta \geq \theta_0$.
- Among these values θ , the one with the highest probability of a type 1 error (i.e., the one with the highest power $\beta(\theta)$) determines the critical value.
Intuitively: we consider the **worst possible case**.

Parameter Test: Example

- We consider a one-sided test of the expected value μ of a normal distribution $N(\mu, \sigma^2)$ with known variance σ^2 , i.e., we consider the hypotheses

$$H_0 : \mu \geq \mu_0, \quad H_a : \mu < \mu_0.$$

- As a test statistic we use the standard point estimator for the expected value

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

This point estimator has the probability density

$$f_{\bar{X}}(x) = N\left(x; \mu, \frac{\sigma^2}{n}\right).$$

- Therefore it is (with the $N(0, 1)$ -distributed random variable Z)

$$\alpha = \beta(\mu_0) = P_{\mu_0}(\bar{X} \leq c) = P\left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \leq \frac{c - \mu_0}{\sigma/\sqrt{n}}\right) = P\left(Z \leq \frac{c - \mu_0}{\sigma/\sqrt{n}}\right).$$

Parameter Test: Example

- We have as a result that

$$\alpha = \Phi \left(\frac{c - \mu_0}{\sigma / \sqrt{n}} \right),$$

where Φ is the distribution function of the standard normal distribution.

- The distribution function Φ is tabulated, because it cannot be represented in closed form. From such a table we retrieve the value z_α satisfying $\alpha = \Phi(z_\alpha)$.
- Then the critical value is

$$c = \mu_0 + z_\alpha \frac{\sigma}{\sqrt{n}}.$$

(Note that the value of z_α is negative due to the usually small value of α . Typical values are $\alpha = 0.1$, $\alpha = 0.05$ or $\alpha = 0.01$.)

- H_0 is rejected if the value \bar{x} of the point estimator \bar{X} does not exceed c , otherwise it is accepted.

Parameter Test: Example

- Let $\sigma = 5.4$, $n = 25$ and $\bar{x} = 128$. We choose $\mu_0 = 130$ and $\alpha = 0.05$.
- From a standard normal distribution table we retrieve $z_{0.05} \approx -1.645$ and get

$$c_{0.05} \approx 130 - 1.645 \frac{5.4}{\sqrt{25}} \approx 128.22.$$

Since $\bar{x} = 128 < 128.22 = c$, we reject the null hypothesis H_0 .

- If, however, we had chosen $\alpha = 0.01$, it would have been (with $z_{0.01} \approx -2.326$):

$$c_{0.01} \approx 130 - 2.326 \frac{5.4}{\sqrt{25}} \approx 127.49$$

Since $\bar{x} = 128 > 127.49 = c$, we would have accepted the null hypothesis H_0 .

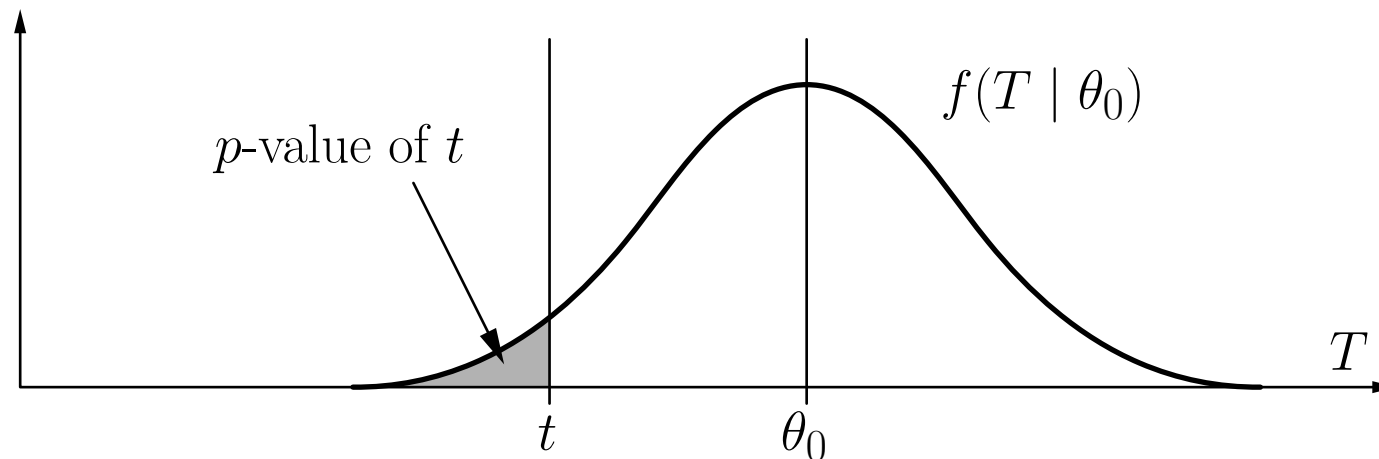
- Instead of fixing a significance level α one may state the so-called **p-value**

$$p = \Phi \left(\frac{128 - 130}{5.4/\sqrt{25}} \right) \approx 0.032.$$

For $\alpha \geq p = 0.032$ the null hypothesis is rejected, for $\alpha < p = 0.032$ accepted.

Parameter Test: p-value

- Let t be the value of the test statistic T that has been computed from a given data set.
(Note: This example illustrates $H_0 : \theta \geq \theta_0$ and $H_a : \theta < \theta_0$.)



- The **p-value** is the probability that a value of t or less can be observed for the chosen test statistic T .
- The p -value is a **lower limit for the significance level α** that may have been chosen if we wanted to reject the null hypothesis H_0 .

Parameter Test: p -value

Attention: p -values are often misused or misinterpreted!

- A low p -value does **not** mean that the result is very reliable!

All that matters for the test is whether the computed p -value is **below the chosen significance level or not**.

(A low p -value could just be a chance event, an accident!)

- The significance level may **not** be chosen **after** computing the p -value, since we tend to choose lower significance levels if we know that they are met. Doing so would undermine the reliability of the procedure!

- Stating p -values is only a convenient way of avoiding a fixed significance level. (Since significance levels are a matter of choice and thus user-dependent.)

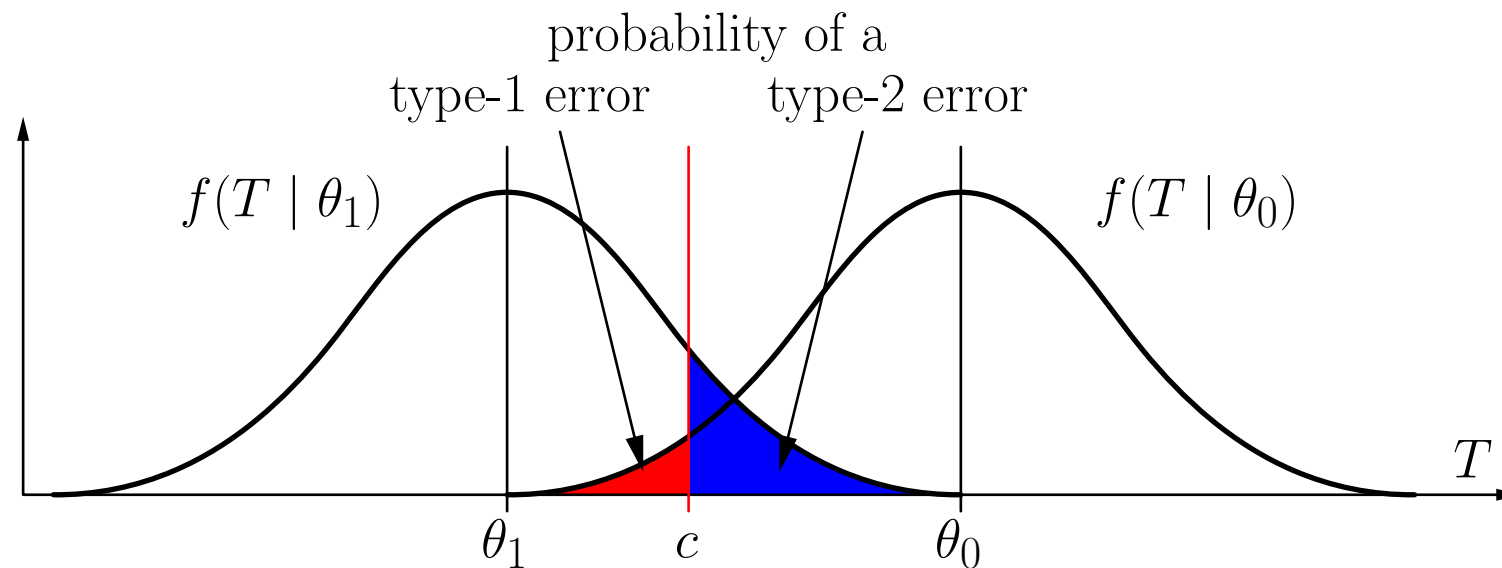
However: A significance level must still be chosen **before** a reported p -value is looked at.

Relevance of the Type-2 Error

- Reminder: There are two possible types of errors:
 - Type 1:** The null hypothesis H_0 is rejected, even though it is correct.
 - Type 2:** The null hypothesis H_0 is accepted, even though it is false.
- Type-1 errors are considered to be more severe, since the null hypothesis gets “the benefit of the doubt”.
- However, **type-2 errors should not be neglected** completely:
 - It is always possible to achieve a vanishing probability of a type-1 error: Simply accept the null hypothesis in all instances, regardless of the data.
 - Unfortunately such an approach maximizes the type-2 error.
- Generally, **type-1 and type-2 errors are complementary quantities**:
The lower we require the type-1 error to be (the lower the significance level), the higher will be the probability of a type-2 error.

Relationship between Type-1 and Type-2 Error

- Suppose there are only two possible parameter values θ_0 and θ_1 with $\theta_1 < \theta_0$. (That is, we have $H_0 : \theta = \theta_0$ and $H_a : \theta = \theta_1$.)



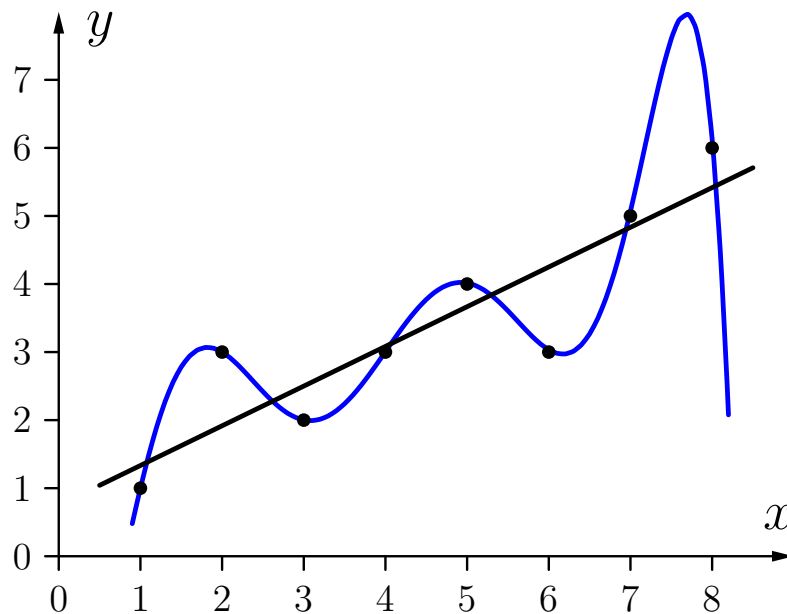
- Lowering the significance level α moves the critical value c to the left: lower type-1 error (red), but higher type-2 error (blue).
- Increasing the significance level α moves the critical value c to the right: higher type-1 error (red), but lower type-2 error (blue).

Inductive Statistics: Model Selection

Model Selection

- Objective: select the model that best fits the data, **taking the model complexity into account.**

The more complex the model, the better it usually fits the data.



black line:
regression line
(2 free parameters)

blue curve:
7th order regression polynomial
(8 free parameters)

- The blue curve fits the data points perfectly, *but it is not a good model.*

Information Criteria

- There is a **tradeoff between model complexity and fit to the data**.

Question: How much better must a more complex model fit the data in order to justify the higher complexity?

- One approach to quantify the tradeoff: **Information Criteria**

Let M be a model and Θ the set of free parameters of M . Then:

$$\text{IC}_{\kappa}(M, \Theta \mid D) = -2 \ln P(D \mid M, \Theta) + \kappa |\Theta|,$$

where D are the sample data and κ is a parameter.

Special cases:

- **Akaike Information Criterion** (AIC): $\kappa = 2$
 - **Bayesian Information Criterion** (BIC): $\kappa = \ln n$,
where n is the sample size.
- The lower the value of these measures, the better the model.

Minimum Description Length

- Idea: Consider the transmission of the data from a sender to a receiver.

Since the transmission of information is costly,
the length of the message to be transmitted should be minimized.

- A good model of the data can be used to transmit the data with fewer bits.

However, the receiver does not know the model the sender used
and thus cannot decode the message.

Therefore: if the sender uses a model, he/she has to transmit the model, too.

- **description length = length of model description
+ length of data description**

(A more complex model increases the length of the model description,
but reduces the length of the data description.)

- The model that leads to the smallest total description length is the best.

Minimum Description Length: Example

- Given: a one-dimensional sample from a polynomial distribution.
- Question: are the probabilities of the attribute values sufficiently different to justify a non-uniform distribution model?
- Coding using **no model** (equal probabilities for all values):

$$l_1 = n \log_2 k,$$

where n is the sample size and k the number of attribute values.

- Coding using a **polynomial distribution model**:

$$l_2 = \underbrace{\log_2 \frac{(n+k-1)!}{n!(k-1)!}}_{\text{model description}} + \underbrace{\log_2 \frac{n!}{x_1! \dots x_k!}}_{\text{data description}}$$

(Idea: Use a codebook with one page per configuration, i.e. frequency distribution (model) and specific sequence (data), and transmit the page number.)

Minimum Description Length: Example

Some details about the codebook idea:

- **Model Description:**

There are n objects (the sample cases) that have to be partitioned into k groups (one for each attribute value). (Model: distribute n balls on k boxes.)

$$\text{Number of possible distributions: } \frac{(n + k - 1)!}{n!(k - 1)!}$$

Idea: number of possible sequences of $n + k - 1$ objects (n balls and $k - 1$ box walls) of which n (the objects) and $k - 1$ (the box walls) are indistinguishable.

- **Data Description:**

There are k groups of objects with x_i , $i = 1, \dots, k$, elements in them. (The values of the x_k are known from the model description.)

$$\text{Number of possible sequences: } \frac{n!}{x_1! \dots x_k!}$$

Summary Statistics

Statistics has two main areas:

- **Descriptive Statistics**

- Display the data in tables or charts.
- Summarize the data in characteristic measures.
- Reduce the dimensionality of the data with principal component analysis.

- **Inductive Statistics**

- Use probability theory to draw inferences about the process that generated the data.
- Parameter Estimation
- Hypothesis Testing
- Model Selection

A short script in German can be found at

<http://fuzzy.cs.uni-magdeburg.de/studium/ida/>

Regression

Regression

- **General Idea of Regression**
 - Method of least squares
- **Linear Regression**
 - An illustrative example
- **Polynomial Regression**
 - Generalization to polynomial functional relationships
- **Multivariate Regression**
 - Generalization to more than one function argument
- **Logistic Regression**
 - Generalization to non-polynomial functional relationships
 - An illustrative example
- **Summary**

Regression

Also known as: **Method of Least Squares** (Carl Friedrich Gauß)

- Given:
- A data set of data tuples (one or more input values and one output value).
 - A hypothesis about the functional relationship between output and input values.
- Desired:
- A parameterization of the conjectured function that minimizes the sum of squared errors (“best fit”).

Depending on

- the hypothesis about the functional relationship and
- the number of arguments to the conjectured function

different types of regression are distinguished.

Reminder: Function Optimization

Task: Find values $\vec{x} = (x_1, \dots, x_m)$ such that $f(\vec{x}) = f(x_1, \dots, x_m)$ is optimal.

Often feasible approach:

- A necessary condition for a (local) optimum (maximum or minimum) is that the partial derivatives w.r.t. the parameters vanish (Pierre Fermat).
- Therefore: (Try to) solve the equation system that results from setting all partial derivatives w.r.t. the parameters equal to zero.

Example task: Minimize $f(x, y) = x^2 + y^2 + xy - 4x - 5y$.

Solution procedure:

1. Take the partial derivatives of the objective function and set them to zero:

$$\frac{\partial f}{\partial x} = 2x + y - 4 = 0, \quad \frac{\partial f}{\partial y} = 2y + x - 5 = 0.$$

2. Solve the resulting (here: linear) equation system: $x = 1, \quad y = 2$.

Linear Regression

- Given: data set $((x_1, y_1), \dots, (x_n, y_n))$ of n data tuples
- Conjecture: the functional relationship is linear, i.e., $y = g(x) = a + bx$.

Approach: Minimize the sum of squared errors, i.e.

$$F(a, b) = \sum_{i=1}^n (g(x_i) - y_i)^2 = \sum_{i=1}^n (a + bx_i - y_i)^2.$$

Necessary conditions for a minimum:

$$\frac{\partial F}{\partial a} = \sum_{i=1}^n 2(a + bx_i - y_i) = 0 \quad \text{and}$$

$$\frac{\partial F}{\partial b} = \sum_{i=1}^n 2(a + bx_i - y_i)x_i = 0$$

Linear Regression

Result of necessary conditions: System of so-called **normal equations**, i.e.

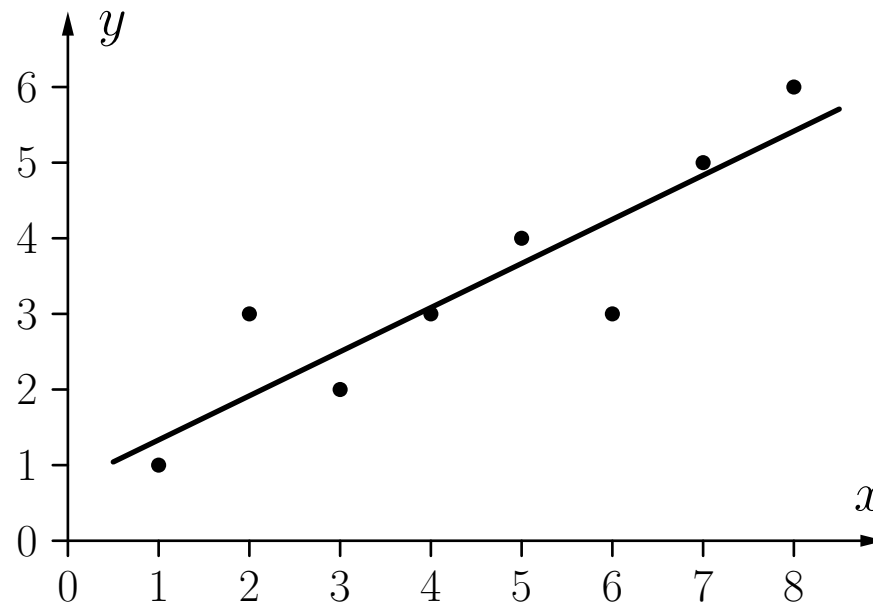
$$\begin{aligned} na + \left(\sum_{i=1}^n x_i \right) b &= \sum_{i=1}^n y_i, \\ \left(\sum_{i=1}^n x_i \right) a + \left(\sum_{i=1}^n x_i^2 \right) b &= \sum_{i=1}^n x_i y_i. \end{aligned}$$

- Two linear equations for two unknowns a and b .
- System can be solved with standard methods from linear algebra.
- Solution is unique unless all x -values are identical.
- The resulting line is called a **regression line**.

Linear Regression: Example

x	1	2	3	4	5	6	7	8
y	1	3	2	3	4	3	5	6

$$y = \frac{3}{4} + \frac{7}{12}x.$$



Least Squares and Maximum Likelihood

A regression line can be interpreted as a **maximum likelihood estimator**:

Assumption: The data generation process can be described well by the model

$$y = a + bx + \xi,$$

where ξ is normally distributed with mean 0 and (unknown) variance σ^2 (σ^2 independent of x , i.e. same dispersion of y for all x).

As a consequence we have

$$f(y \mid x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(y - (a + bx))^2}{2\sigma^2}\right).$$

With this expression we can set up the **likelihood function**

$$\begin{aligned} L((x_1, y_1), \dots, (x_n, y_n); a, b, \sigma^2) \\ = \prod_{i=1}^n f(x_i) f(y_i \mid x_i) = \prod_{i=1}^n f(x_i) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(y_i - (a + bx_i))^2}{2\sigma^2}\right). \end{aligned}$$

Least Squares and Maximum Likelihood

To simplify taking the derivatives, we compute the natural logarithm:

$$\begin{aligned} \ln L((x_1, y_1), \dots (x_n, y_n); a, b, \sigma^2) \\ &= \ln \prod_{i=1}^n f(x_i) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(y_i - (a + bx_i))^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^n \ln f(x_i) + \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (a + bx_i))^2 \end{aligned}$$

From this expression it becomes clear that (provided $f(x)$ is independent of a , b , and σ^2) maximizing the likelihood function is equivalent to minimizing

$$F(a, b) = \sum_{i=1}^n (y_i - (a + bx_i))^2.$$

Interpreting the method of least squares as a maximum likelihood estimator works also for the generalizations to polynomials and multilinear functions discussed next.

Polynomial Regression

Generalization to polynomials

$$y = p(x) = a_0 + a_1x + \dots + a_mx^m$$

Approach: Minimize the sum of squared errors, i.e.

$$F(a_0, a_1, \dots, a_m) = \sum_{i=1}^n (p(x_i) - y_i)^2 = \sum_{i=1}^n (a_0 + a_1x_i + \dots + a_mx_i^m - y_i)^2$$

Necessary conditions for a minimum: All partial derivatives vanish, i.e.

$$\frac{\partial F}{\partial a_0} = 0, \quad \frac{\partial F}{\partial a_1} = 0, \quad \dots, \quad \frac{\partial F}{\partial a_m} = 0.$$

Polynomial Regression

System of normal equations for polynomials

$$\begin{array}{ccccccc} na_0 & + & \left(\sum_{i=1}^n x_i \right) a_1 & + \dots + & \left(\sum_{i=1}^n x_i^m \right) a_m & = & \sum_{i=1}^n y_i \\ \left(\sum_{i=1}^n x_i \right) a_0 & + & \left(\sum_{i=1}^n x_i^2 \right) a_1 & + \dots + & \left(\sum_{i=1}^n x_i^{m+1} \right) a_m & = & \sum_{i=1}^n x_i y_i \\ \vdots & & & & & & \vdots \\ \left(\sum_{i=1}^n x_i^m \right) a_0 & + & \left(\sum_{i=1}^n x_i^{m+1} \right) a_1 & + \dots + & \left(\sum_{i=1}^n x_i^{2m} \right) a_m & = & \sum_{i=1}^n x_i^m y_i, \end{array}$$

- $m + 1$ linear equations for $m + 1$ unknowns a_0, \dots, a_m .
- System can be solved with standard methods from linear algebra.
- Solution is unique unless the points lie exactly on a polynomial of lower degree.

Multilinear Regression

Generalization to more than one argument

$$z = f(x, y) = a + bx + cy$$

Approach: Minimize the sum of squared errors, i.e.

$$F(a, b, c) = \sum_{i=1}^n (f(x_i, y_i) - z_i)^2 = \sum_{i=1}^n (a + bx_i + cy_i - z_i)^2$$

Necessary conditions for a minimum: All partial derivatives vanish, i.e.

$$\begin{aligned}\frac{\partial F}{\partial a} &= \sum_{i=1}^n 2(a + bx_i + cy_i - z_i) = 0, \\ \frac{\partial F}{\partial b} &= \sum_{i=1}^n 2(a + bx_i + cy_i - z_i)x_i = 0, \\ \frac{\partial F}{\partial c} &= \sum_{i=1}^n 2(a + bx_i + cy_i - z_i)y_i = 0.\end{aligned}$$

Multilinear Regression

System of normal equations for several arguments

$$\begin{aligned} na + \left(\sum_{i=1}^n x_i \right) b + \left(\sum_{i=1}^n y_i \right) c &= \sum_{i=1}^n z_i \\ \left(\sum_{i=1}^n x_i \right) a + \left(\sum_{i=1}^n x_i^2 \right) b + \left(\sum_{i=1}^n x_i y_i \right) c &= \sum_{i=1}^n z_i x_i \\ \left(\sum_{i=1}^n y_i \right) a + \left(\sum_{i=1}^n x_i y_i \right) b + \left(\sum_{i=1}^n y_i^2 \right) c &= \sum_{i=1}^n z_i y_i \end{aligned}$$

- 3 linear equations for 3 unknowns a , b , and c .
- System can be solved with standard methods from linear algebra.
- Solution is unique unless all data points lie on a straight line.

Multilinear Regression

General multilinear case:

$$y = f(x_1, \dots, x_m) = a_0 + \sum_{k=1}^m a_k x_k$$

Approach: Minimize the sum of squared errors, i.e.

$$F(\vec{a}) = (\mathbf{X}\vec{a} - \vec{y})^\top (\mathbf{X}\vec{a} - \vec{y}),$$

where

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nm} \end{pmatrix}, \quad \vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \text{and} \quad \vec{a} = \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_m \end{pmatrix}$$

Necessary condition for a minimum:

$$\nabla_{\vec{a}} F(\vec{a}) = \nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y})^\top (\mathbf{X}\vec{a} - \vec{y}) = \vec{0}$$

Multilinear Regression

- $\nabla_{\vec{a}} F(\vec{a})$ may easily be computed by remembering that the differential operator

$$\nabla_{\vec{a}} = \left(\frac{\partial}{\partial a_0}, \dots, \frac{\partial}{\partial a_m} \right)$$

behaves formally like a vector that is “multiplied” to the sum of squared errors.

- Alternatively, one may write out the differentiation componentwise.

With the former method we obtain for the derivative:

$$\begin{aligned} & \nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y})^\top (\mathbf{X}\vec{a} - \vec{y}) \\ &= (\nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y}))^\top (\mathbf{X}\vec{a} - \vec{y}) + ((\mathbf{X}\vec{a} - \vec{y})^\top (\nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y})))^\top \\ &= (\nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y}))^\top (\mathbf{X}\vec{a} - \vec{y}) + (\nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y}))^\top (\mathbf{X}\vec{a} - \vec{y}) \\ &= 2\mathbf{X}^\top (\mathbf{X}\vec{a} - \vec{y}) \\ &= 2\mathbf{X}^\top \mathbf{X}\vec{a} - 2\mathbf{X}^\top \vec{y} = \vec{0} \end{aligned}$$

Multilinear Regression

Necessary condition for a minimum therefore:

$$\begin{aligned}\nabla_{\vec{a}} F(\vec{a}) &= \nabla_{\vec{a}} (\mathbf{X}\vec{a} - \vec{y})^\top (\mathbf{X}\vec{a} - \vec{y}) \\ &= 2\mathbf{X}^\top \mathbf{X}\vec{a} - 2\mathbf{X}^\top \vec{y} \stackrel{!}{=} \vec{0}\end{aligned}$$

As a consequence we get the system of **normal equations**:

$$\mathbf{X}^\top \mathbf{X}\vec{a} = \mathbf{X}^\top \vec{y}$$

This system has a unique solution if $\mathbf{X}^\top \mathbf{X}$ is not singular. Then we have

$$\vec{a} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{y}.$$

$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is called the (Moore–Penrose) **pseudoinverse** of the matrix \mathbf{X} .

With the matrix-vector representation of the regression problem an extension to **multipolynomial regression** is straightforward:

Simply add the desired products of powers to the matrix \mathbf{X} .

Logistic Regression

Generalization to non-polynomial functions

Idea: Find transformation to linear/polynomial case.

Simple example: The function $y = ax^b$
can be transformed into $\ln y = \ln a + b \cdot \ln x$.

Special case: **logistic function**

$$y = \frac{Y}{1 + e^{a+bx}} \quad \Leftrightarrow \quad \frac{1}{\frac{Y}{y}} = \frac{1 + e^{a+bx}}{Y} \quad \Leftrightarrow \quad \frac{Y - y}{y} = e^{a+bx}.$$

Result: Apply so-called **Logit Transformation**

$$\ln \left(\frac{Y - y}{y} \right) = a + bx.$$

Logistic Regression: Example

x	1	2	3	4	5
y	0.4	1.0	3.0	5.0	5.6

Transform the data with

$$z = \ln \left(\frac{Y - y}{y} \right), \quad Y = 6.$$

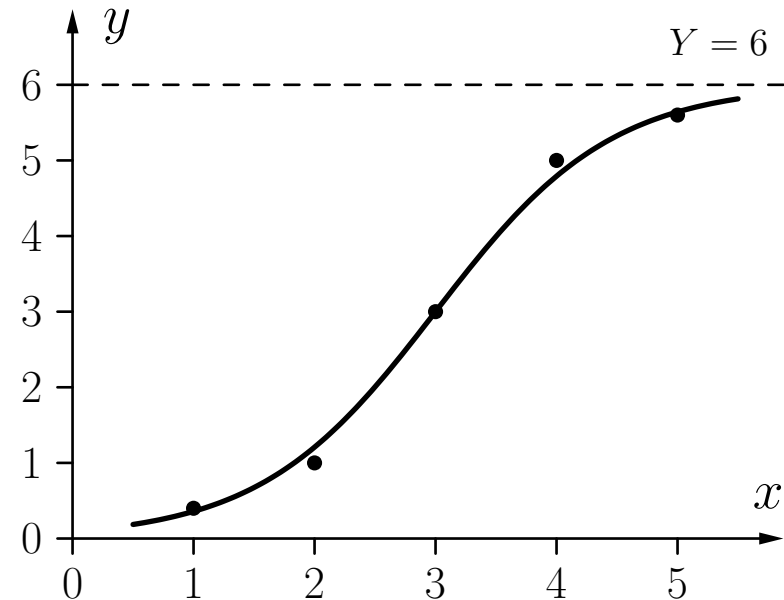
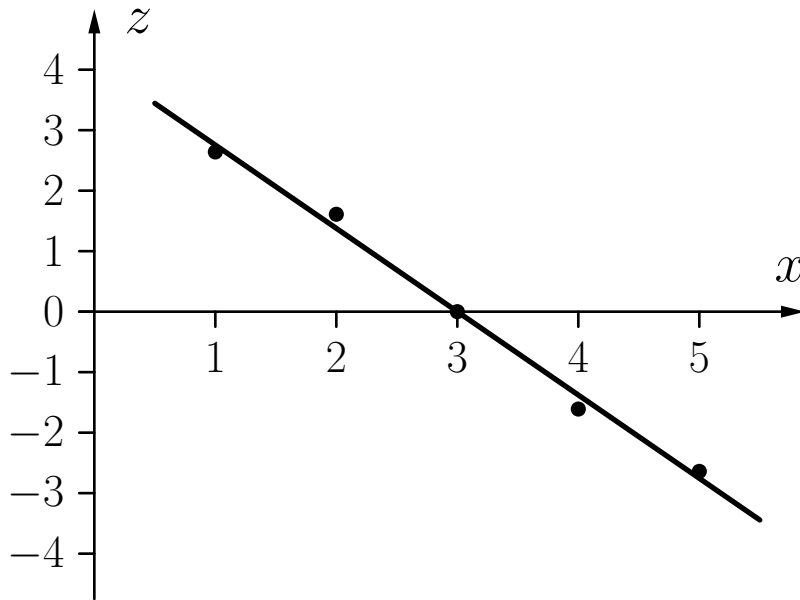
The transformed data points are

x	1	2	3	4	5
z	2.64	1.61	0.00	-1.61	-2.64

The resulting regression line is

$$z \approx -1.3775x + 4.133.$$

Logistic Regression: Example

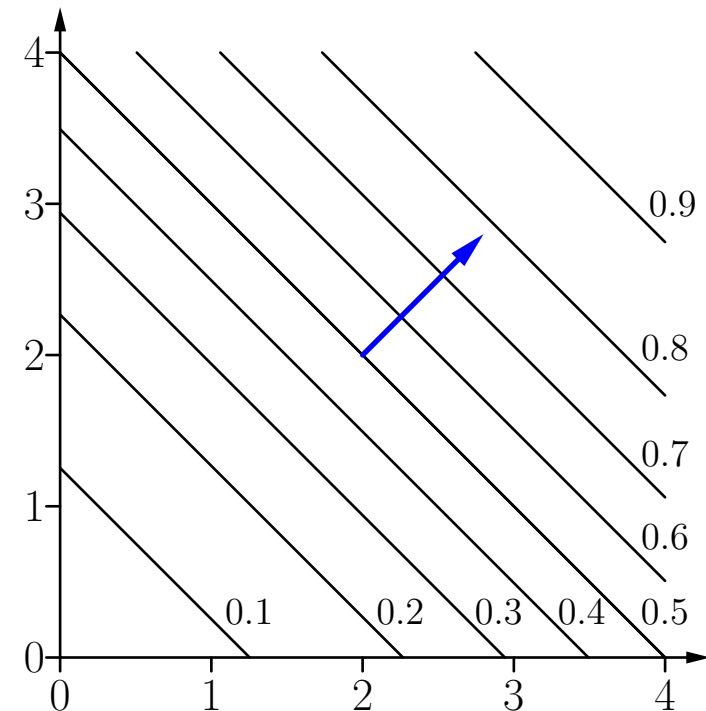
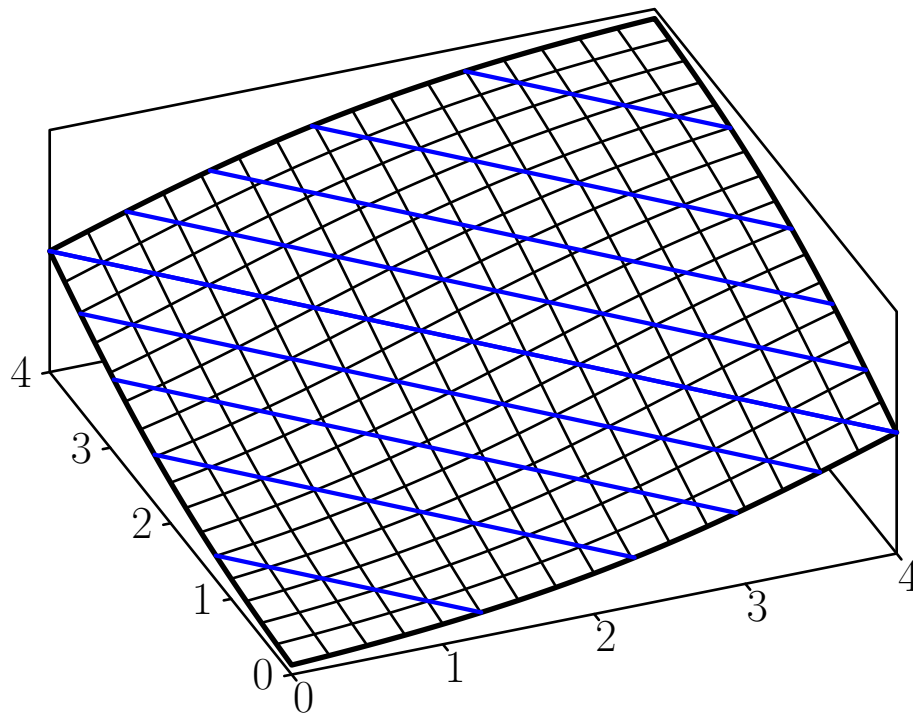


- **Attention:** The sum of squared errors is minimized only in the space the transformation maps to, not in the original space.
- Nevertheless this approach usually leads to very good results.
The result may be improved by a gradient descent in the original space.

Logistic Regression: Two-dimensional Example

Example logistic function for two arguments x_1 and x_2 :

$$y = \frac{1}{1 + \exp(4 - x_1 - x_2)} = \frac{1}{1 + \exp(4 - (1, 1)(x_1, x_2)^\top)}$$



Logistic Regression: Two Class Problems

- Let C be a class attrib., $\text{dom}(C) = \{c_1, c_2\}$, and \vec{X} an m -dim. random vector.
Let $P(C = c_1 \mid \vec{X} = \vec{x}) = p(\vec{x})$ and $P(C = c_2 \mid \vec{X} = \vec{x}) = 1 - p(\vec{x})$.
- **Given:** A set of data points $\mathbf{X} = \{\vec{x}_1, \dots, \vec{x}_n\}$ (realizations of \vec{X}), each of which belongs to one of the two classes c_1 and c_2 .
- **Desired:** A simple description of the function $p(\vec{x})$.
- **Approach:** Describe p by a logistic function:

$$p(\vec{x}) = \frac{1}{1 + e^{a_0 + \vec{a}\vec{x}}} = \frac{1}{1 + \exp\left(a_0 + \sum_{i=1}^m a_i x_i\right)}$$

Apply logit transformation to $p(x)$:

$$\ln\left(\frac{1 - p(\vec{x})}{p(\vec{x})}\right) = a_0 + \vec{a}\vec{x} = a_0 + \sum_{i=1}^m a_i x_i$$

The values $p(\vec{x}_i)$ may be obtained by kernel estimation.

Kernel Estimation

- **Idea:** Define an “influence function” (kernel), which describes how strongly a data point influences the probability estimate for neighboring points.

- Common choice for the kernel function: **Gaussian function**

$$K(\vec{x}, \vec{y}) = \frac{1}{(2\pi\sigma^2)^{\frac{m}{2}}} \exp\left(-\frac{(\vec{x} - \vec{y})^\top (\vec{x} - \vec{y})}{2\sigma^2}\right)$$

- Kernel estimate of probability density given a data set $\mathcal{X} = \{\vec{x}_1, \dots, \vec{x}_n\}$:

$$\hat{f}(\vec{x}) = \frac{1}{n} \sum_{i=1}^n K(\vec{x}, \vec{x}_i).$$

- Kernel estimation applied to a two class problem:

$$\hat{p}(\vec{x}) = \frac{\sum_{i=1}^n c(\vec{x}_i) K(\vec{x}, \vec{x}_i)}{\sum_{i=1}^n K(\vec{x}, \vec{x}_i)}.$$

(It is $c(\vec{x}_i) = 1$ if x_i belongs to class c_1 and $c(\vec{x}_i) = 0$ otherwise.)

Summary Regression

- **Minimize the Sum of Squared Errors**
 - Write the sum of squared errors as a function of the parameters to be determined.
- **Exploit Necessary Conditions for a Minimum**
 - Partial derivatives w.r.t. the parameters to determine must vanish.
- **Solve the System of Normal Equations**
 - The best fit parameters are the solution of the system of normal equations.
- **Non-polynomial Regression Functions**
 - Find a transformation to the multipolynomial case.
 - Logistic regression can be used to solve two class classification problems.

Bayes Classifiers

Bayes Classifiers

- **Probabilistic Classification and Bayes' Rule**
- **Naive Bayes Classifiers**
 - Derivation of the classification formula
 - Probability estimation and Laplace correction
 - Simple examples of naive Bayes classifiers
 - A naive Bayes classifier for the Iris data
- **Full Bayes Classifiers**
 - Derivation of the classification formula
 - Comparison to naive Bayes classifiers
 - A simple example of a full Bayes classifier
 - A full Bayes classifier for the Iris data
- **Summary**

Probabilistic Classification

- A classifier is an algorithm that assigns a class from a predefined set to a case or object, based on the values of descriptive attributes.
- An optimal classifier maximizes the probability of a correct class assignment.
 - Let C be a class attribute with $\text{dom}(C) = \{c_1, \dots, c_{n_C}\}$, which occur with probabilities p_i , $1 \leq i \leq n_C$.
 - Let q_i be the probability with which a classifier assigns class c_i . ($q_i \in \{0, 1\}$ for a deterministic classifier)
 - The probability of a correct assignment is

$$P(\text{correct assignment}) = \sum_{i=1}^{n_C} p_i q_i.$$

- Therefore the best choice for the q_i is

$$q_i = \begin{cases} 1, & \text{if } p_i = \max_{k=1}^{n_C} p_k, \\ 0, & \text{otherwise.} \end{cases}$$

Probabilistic Classification

- Consequence: An optimal classifier should assign the **most probable class**.
- This argument does not change if we take descriptive attributes into account.
 - Let $U = \{A_1, \dots, A_m\}$ be a set of descriptive attributes with domains $\text{dom}(A_k)$, $1 \leq k \leq m$.
 - Let $A_1 = a_1, \dots, A_m = a_m$ be an instantiation of the attributes.
 - An optimal classifier should assign the class c_i for which

$$P(C = c_i \mid A_1 = a_1, \dots, A_m = a_m) = \max_{j=1}^{n_C} P(C = c_j \mid A_1 = a_1, \dots, A_m = a_m)$$

- **Problem:** We cannot store a class (or the class probabilities) for every possible instantiation $A_1 = a_1, \dots, A_m = a_m$ of the descriptive attributes. (The table size grows exponentially with the number of attributes.)
- Therefore: **Simplifying assumptions are necessary.**

Bayes' Rule and Bayes' Classifiers

- Bayes' rule is a formula that can be used to “invert” conditional probabilities: Let X and Y be events, $P(X) > 0$. Then

$$P(Y | X) = \frac{P(X | Y) \cdot P(Y)}{P(X)}.$$

- Bayes' rule follows directly from the definition of conditional probability:

$$P(Y | X) = \frac{P(X \cap Y)}{P(X)} \quad \text{and} \quad P(X | Y) = \frac{P(X \cap Y)}{P(Y)}.$$

- Bayes' classifiers: Compute the class probabilities as

$$P(C = c_i | A_1 = a_1, \dots, A_m = a_m) = \frac{P(A_1 = a_1, \dots, A_m = a_m | C = c_i) \cdot P(C = c_i)}{P(A_1 = a_1, \dots, A_m = a_m)}.$$

- Looks unreasonable at first sight: Even more probabilities to store.

Naive Bayes Classifiers

Naive Assumption:

The descriptive attributes are conditionally independent given the class.

Bayes' Rule:

$$P(C = c_i \mid \omega) = \frac{P(A_1 = a_1, \dots, A_m = a_m \mid C = c_i) \cdot P(C = c_i)}{P(A_1 = a_1, \dots, A_m = a_m)} \quad \leftarrow p_0$$

Chain Rule of Probability:

$$P(C = c_i \mid \omega) = \frac{P(C = c_i)}{p_0} \cdot \prod_{k=1}^m P(A_k = a_k \mid A_1 = a_1, \dots, A_{k-1} = a_{k-1}, C = c_i)$$

Conditional Independence Assumption:

$$P(C = c_i \mid \omega) = \frac{P(C = c_i)}{p_0} \cdot \prod_{k=1}^m P(A_k = a_k \mid C = c_i)$$

Reminder: Chain Rule of Probability

- Based on the **product rule** of probability:

$$P(A \wedge B) = P(A \mid B) \cdot P(B)$$

(Multiply definition of conditional probability with $P(B)$.)

- Multiple application** of the product rule yields:

$$\begin{aligned} P(A_1, \dots, A_m) &= P(A_m \mid A_1, \dots, A_{m-1}) \cdot P(A_1, \dots, A_{m-1}) \\ &= P(A_m \mid A_1, \dots, A_{m-1}) \\ &\quad \cdot P(A_{m-1} \mid A_1, \dots, A_{m-2}) \cdot P(A_1, \dots, A_{m-2}) \\ &= \vdots \\ &= \prod_{k=1}^m P(A_k \mid A_1, \dots, A_{k-1}) \end{aligned}$$

- The scheme works also if there is already a condition in the original expression:

$$P(A_1, \dots, A_m \mid C) = \prod_{i=1}^m P(A_i \mid A_1, \dots, A_{i-1}, C)$$

Conditional Independence

- Reminder: **stochastic independence** (unconditional)

$$P(A \wedge B) = P(A) \cdot P(B)$$

(Joint probability is the product of the individual probabilities.)

- Comparison to the **product rule**

$$P(A \wedge B) = P(A \mid B) \cdot P(B)$$

shows that this is equivalent to

$$P(A \mid B) = P(A)$$

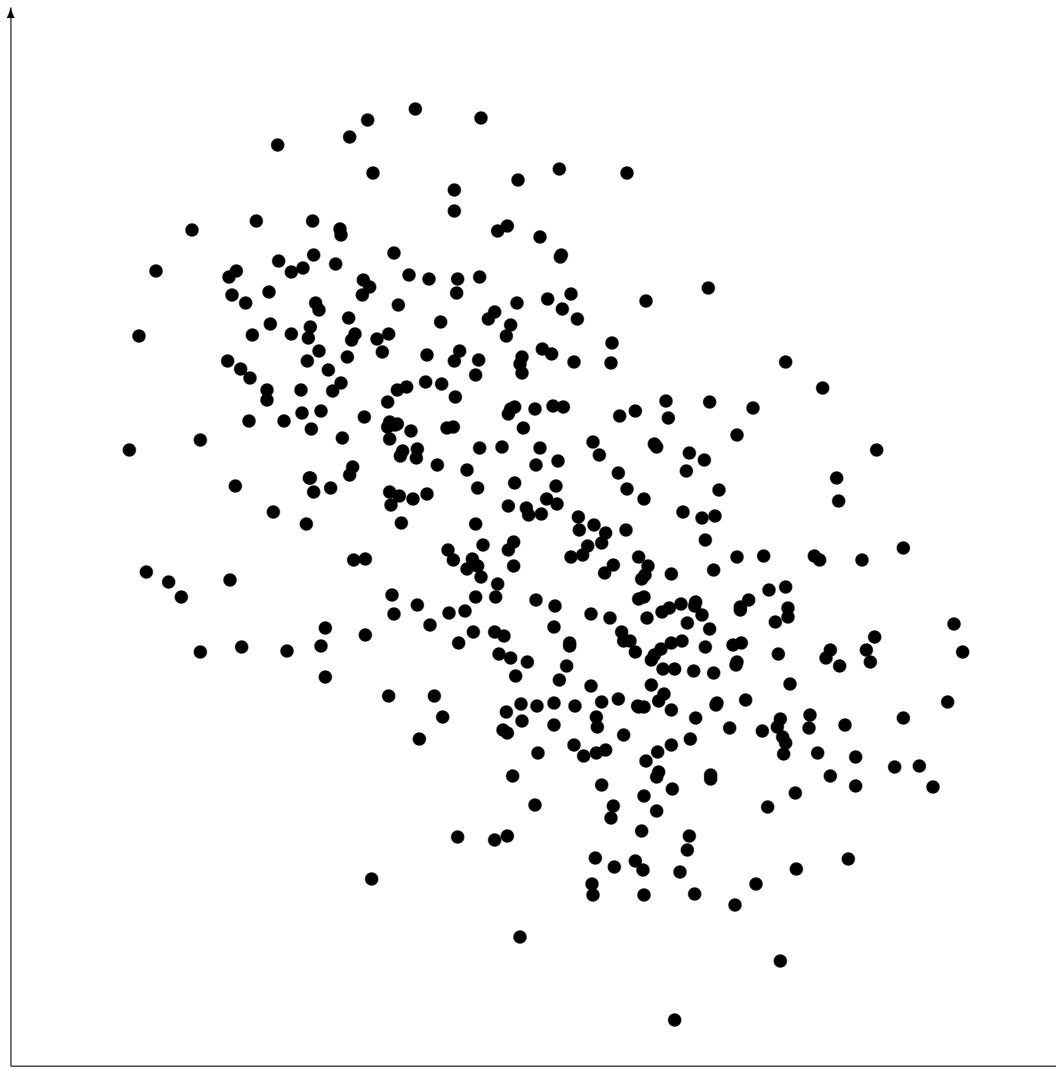
- The same formulae hold conditionally, i.e.

$$P(A \wedge B \mid C) = P(A \mid C) \cdot P(B \mid C) \quad \text{and}$$

$$P(A \mid B, C) = P(A \mid C).$$

- **Conditional independence allows us to cancel some conditions.**

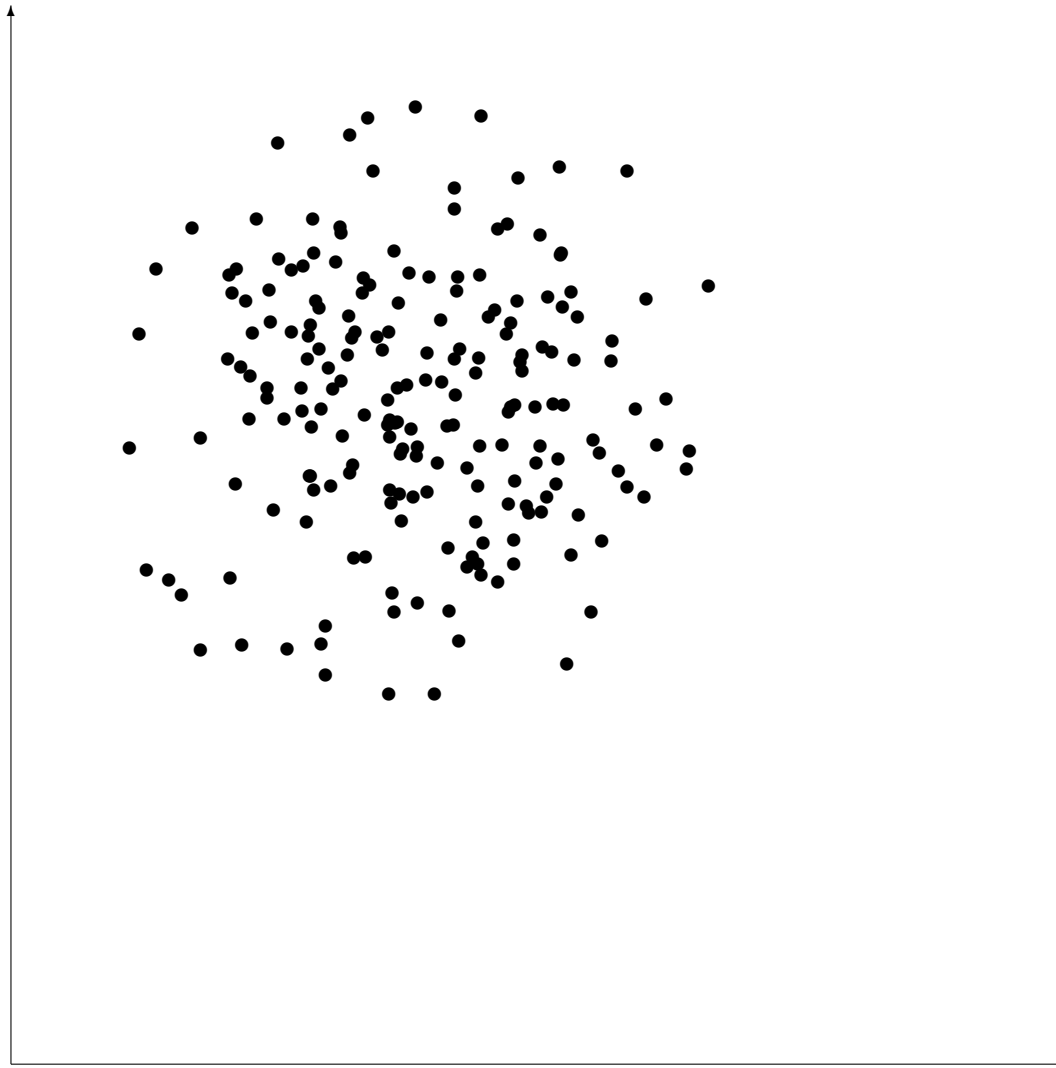
Conditional Independence: An Example



Group 1

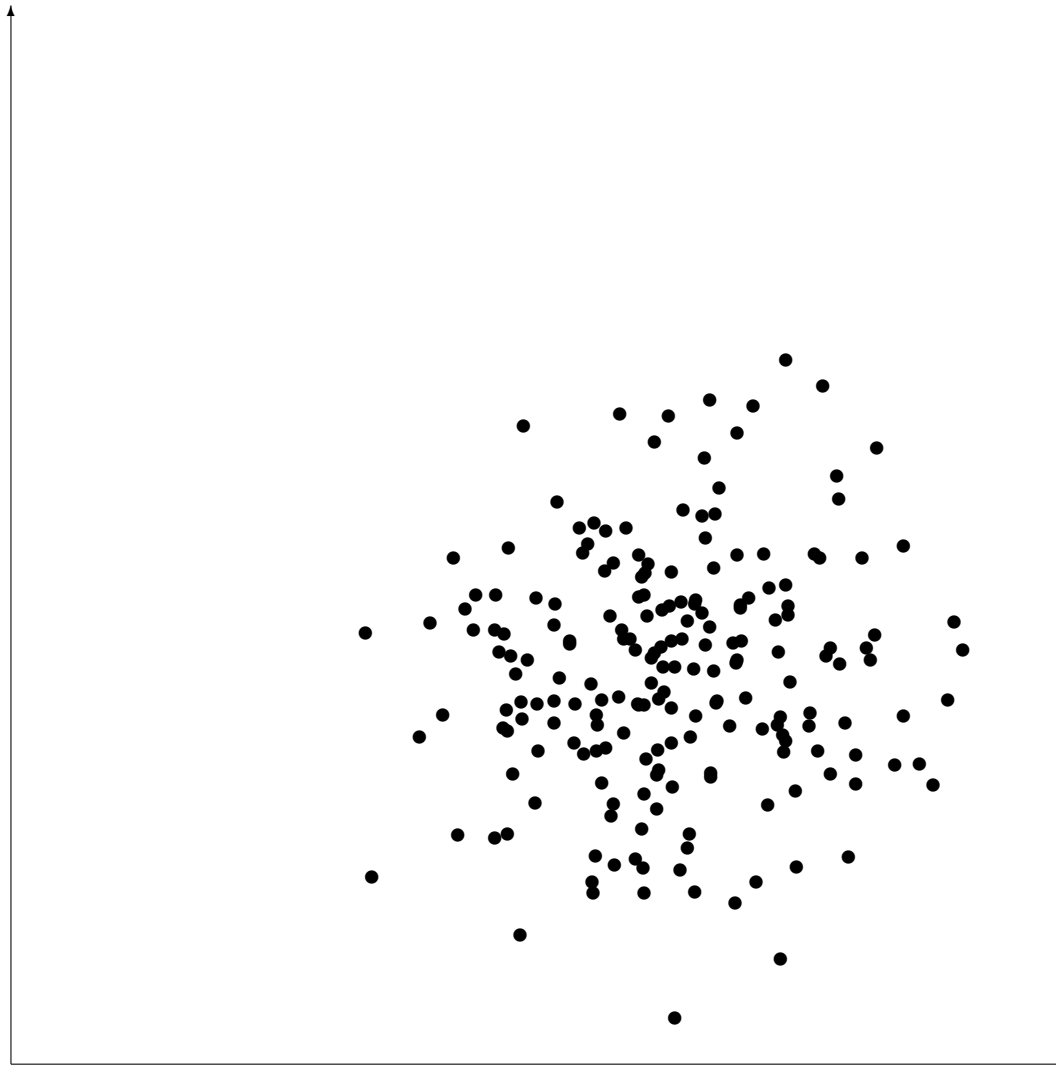
Group 2

Conditional Independence: An Example



Group 1

Conditional Independence: An Example



Group 2

Naive Bayes Classifiers

- Consequence: Manageable amount of data to store.
Store distributions $P(C = c_i)$ and $\forall 1 \leq k \leq m : P(A_k = a_k \mid C = c_i)$.
- It is not necessary to compute p_0 explicitly, because it can be computed implicitly by normalizing the computed values to sum 1.

Estimation of Probabilities:

- **Nominal/Symbolic Attributes**

$$\hat{P}(A_k = a_k \mid C = c_i) = \frac{\#(A_k = a_k, C = c_i) + \gamma}{\#(C = c_i) + n_{A_k} \gamma}$$

γ is called **Laplace correction**.

$\gamma = 0$: Maximum likelihood estimation.

Common choices: $\gamma = 1$ or $\gamma = \frac{1}{2}$.

Naive Bayes Classifiers

Estimation of Probabilities:

- **Metric/Numeric Attributes:** Assume a normal distribution.

$$P(A_k = a_k \mid C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_k(c_i)} \exp\left(-\frac{(a_k - \mu_k(c_i))^2}{2\sigma_k^2(c_i)}\right)$$

- Estimate of mean value

$$\hat{\mu}_k(c_i) = \frac{1}{\#(C = c_i)} \sum_{j=1}^{\#(C=c_i)} a_k(j)$$

- Estimate of variance

$$\hat{\sigma}_k^2(c_i) = \frac{1}{\xi} \sum_{j=1}^{\#(C=c_i)} (a_k(j) - \hat{\mu}_k(c_i))^2$$

$\xi = \#(C = c_i)$: Maximum likelihood estimation

$\xi = \#(C = c_i) - 1$: Unbiased estimation

Naive Bayes Classifiers: Simple Example 1

No	Sex	Age	Blood pr.	Drug
1	male	20	normal	A
2	female	73	normal	B
3	female	37	high	A
4	male	33	low	B
5	female	48	high	A
6	male	29	normal	A
7	female	52	normal	B
8	male	42	low	B
9	male	61	normal	B
10	female	30	normal	A
11	female	26	low	B
12	male	54	high	A

$P(\text{Drug})$	A	B
	0.5	0.5

$P(\text{Sex} \mid \text{Drug})$	A	B
male	0.5	0.5
female	0.5	0.5

$P(\text{Age} \mid \text{Drug})$	A	B
μ	36.3	47.8
σ^2	161.9	311.0

$P(\text{Blood Pr.} \mid \text{Drug})$	A	B
low	0	0.5
normal	0.5	0.5
high	0.5	0

A simple database and estimated (conditional) probability distributions.

Naive Bayes Classifiers: Simple Example 1

$$\begin{aligned} &P(\text{Drug A} \mid \text{male}, 61, \text{normal}) \\ &= c_1 \cdot P(\text{Drug A}) \cdot P(\text{male} \mid \text{Drug A}) \cdot P(61 \mid \text{Drug A}) \cdot P(\text{normal} \mid \text{Drug A}) \\ &\approx c_1 \cdot 0.5 \cdot 0.5 \cdot 0.004787 \cdot 0.5 = c_1 \cdot 5.984 \cdot 10^{-4} = 0.219 \end{aligned}$$

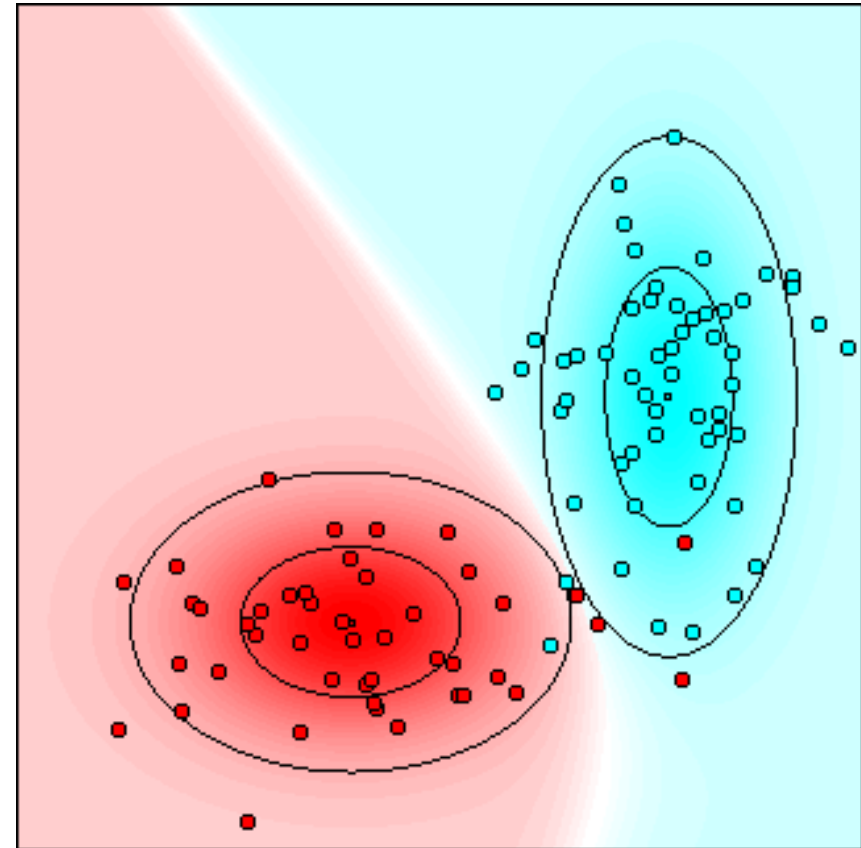
$$\begin{aligned} &P(\text{Drug B} \mid \text{male}, 61, \text{normal}) \\ &= c_1 \cdot P(\text{Drug B}) \cdot P(\text{male} \mid \text{Drug B}) \cdot P(61 \mid \text{Drug B}) \cdot P(\text{normal} \mid \text{Drug B}) \\ &\approx c_1 \cdot 0.5 \cdot 0.5 \cdot 0.017120 \cdot 0.5 = c_1 \cdot 2.140 \cdot 10^{-3} = 0.781 \end{aligned}$$

$$\begin{aligned} &P(\text{Drug A} \mid \text{female}, 30, \text{normal}) \\ &= c_2 \cdot P(\text{Drug A}) \cdot P(\text{female} \mid \text{Drug A}) \cdot P(30 \mid \text{Drug A}) \cdot P(\text{normal} \mid \text{Drug A}) \\ &\approx c_2 \cdot 0.5 \cdot 0.5 \cdot 0.027703 \cdot 0.5 = c_2 \cdot 3.471 \cdot 10^{-3} = 0.671 \end{aligned}$$

$$\begin{aligned} &P(\text{Drug B} \mid \text{female}, 30, \text{normal}) \\ &= c_2 \cdot P(\text{Drug B}) \cdot P(\text{female} \mid \text{Drug B}) \cdot P(30 \mid \text{Drug B}) \cdot P(\text{normal} \mid \text{Drug B}) \\ &\approx c_2 \cdot 0.5 \cdot 0.5 \cdot 0.013567 \cdot 0.5 = c_2 \cdot 1.696 \cdot 10^{-3} = 0.329 \end{aligned}$$

Naive Bayes Classifiers: Simple Example 2

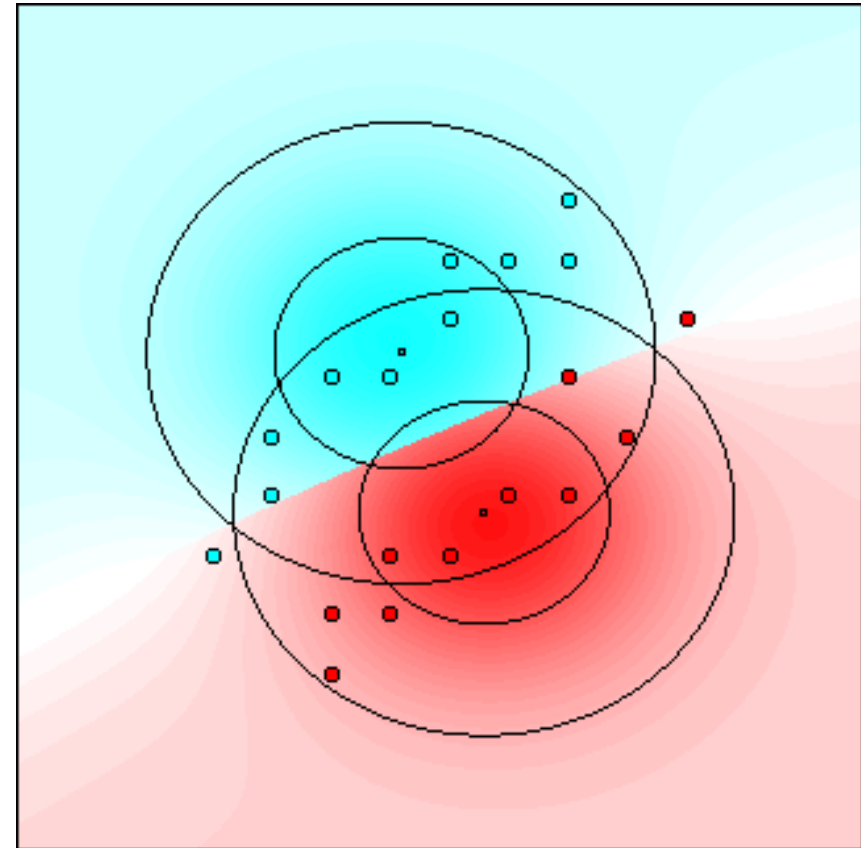
- 100 data points, 2 classes
- Small squares: mean values
- Inner ellipses: one standard deviation
- Outer ellipses: two standard deviations
- Classes overlap: classification is not perfect



Naive Bayes Classifier

Naive Bayes Classifiers: Simple Example 3

- 20 data points, 2 classes
- Small squares: mean values
- Inner ellipses: one standard deviation
- Outer ellipses: two standard deviations
- Attributes are not conditionally independent given the class



Naive Bayes Classifier

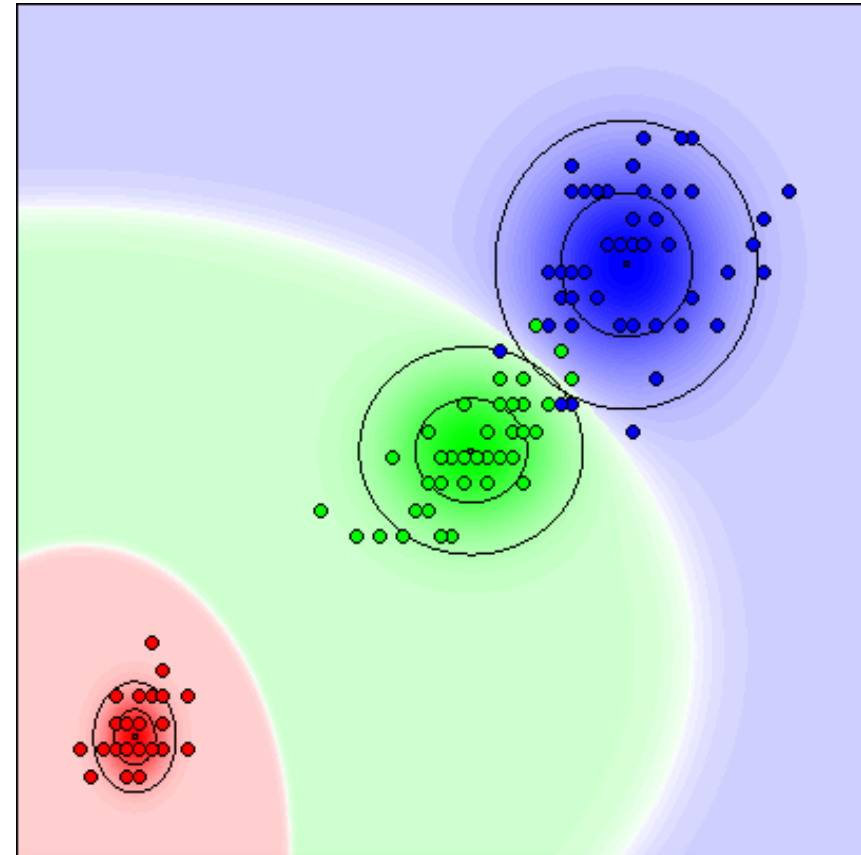
The Iris Data

pictures not available in online version

- First analyzed by Ronald Aylmer Fischer (famous statistician).
- 150 cases in total, 50 cases per Iris flower type.
- Measurements of sepal length and width and petal length and width (in cm).
- Most famous data set in pattern recognition and data analysis.

Naive Bayes Classifiers: Iris Data

- 150 data points, 3 classes
 - Iris setosa (red)
 - Iris versicolor (green)
 - Iris virginica (blue)
- Shown: 2 out of 4 attributes
 - sepal length
 - sepal width
 - petal length (horizontal)
 - petal width (vertical)
- 6 misclassifications on the training data (with all 4 attributes)



Naive Bayes Classifier

Full Bayes Classifiers

- Restricted to metric/numeric attributes (only the class is nominal/symbolic).
- **Simplifying Assumption:**
Each class can be described by a multivariate normal distribution.

$$f(A_1 = a_1, \dots, A_m = a_m \mid C = c_i) \\ = \frac{1}{\sqrt{(2\pi)^m |\mathbf{\Sigma}_i|}} \exp \left(-\frac{1}{2} (\vec{a} - \vec{\mu}_i)^\top \mathbf{\Sigma}_i^{-1} (\vec{a} - \vec{\mu}_i) \right)$$

$\vec{\mu}_i$: mean value vector for class c_i

$\mathbf{\Sigma}_i$: covariance matrix for class c_i

- Intuitively: Each class has a bell-shaped probability density.
- Naive Bayes classifiers: Covariance matrices are diagonal matrices.
(Details about this relation are given below.)

Full Bayes Classifiers

Estimation of Probabilities:

- Estimate of mean value vector

$$\hat{\vec{\mu}}_i = \frac{1}{\#(C = c_i)} \sum_{j=1}^{\#(C=c_i)} \vec{a}(j)$$

- Estimate of covariance matrix

$$\hat{\Sigma}_i = \frac{1}{\xi} \sum_{j=1}^{\#(C=c_i)} \left(\vec{a}(j) - \hat{\vec{\mu}}_i \right) \left(\vec{a}(j) - \hat{\vec{\mu}}_i \right)^\top$$

$\xi = \#(C = c_i)$: Maximum likelihood estimation

$\xi = \#(C = c_i) - 1$: Unbiased estimation

\vec{x}^\top denotes the transpose of the vector \vec{x} .

$\vec{x}\vec{x}^\top$ is the so-called **outer product** or **matrix product** of \vec{x} with itself.

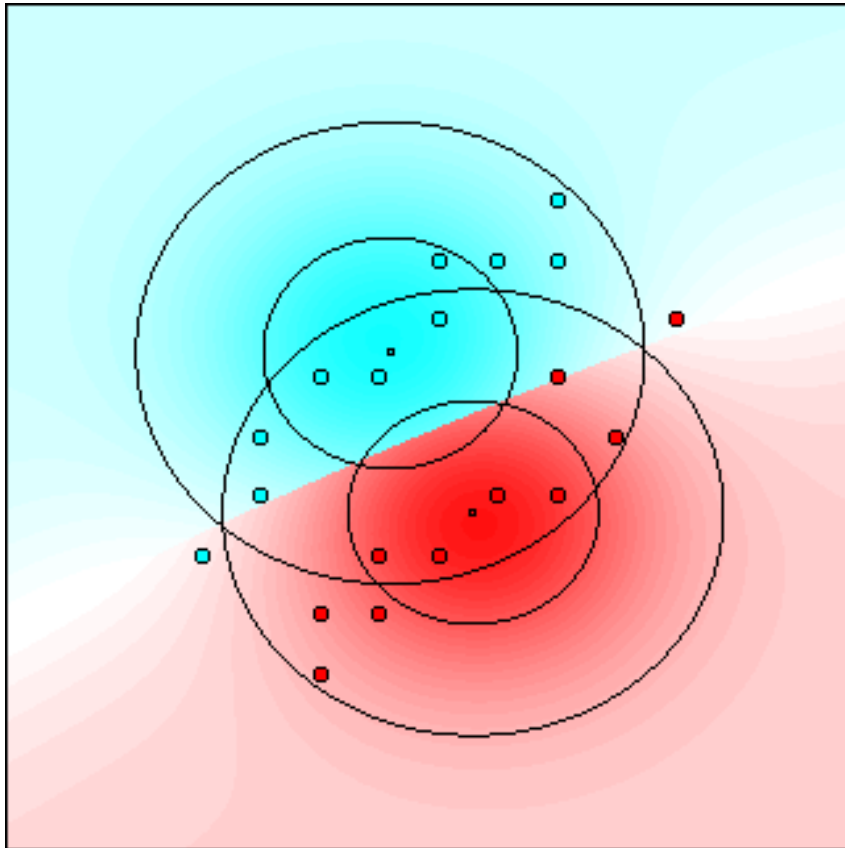
Comparison of Naive and Full Bayes Classifiers

Naive Bayes classifiers for metric/numeric data are equivalent to full Bayes classifiers with diagonal covariance matrices:

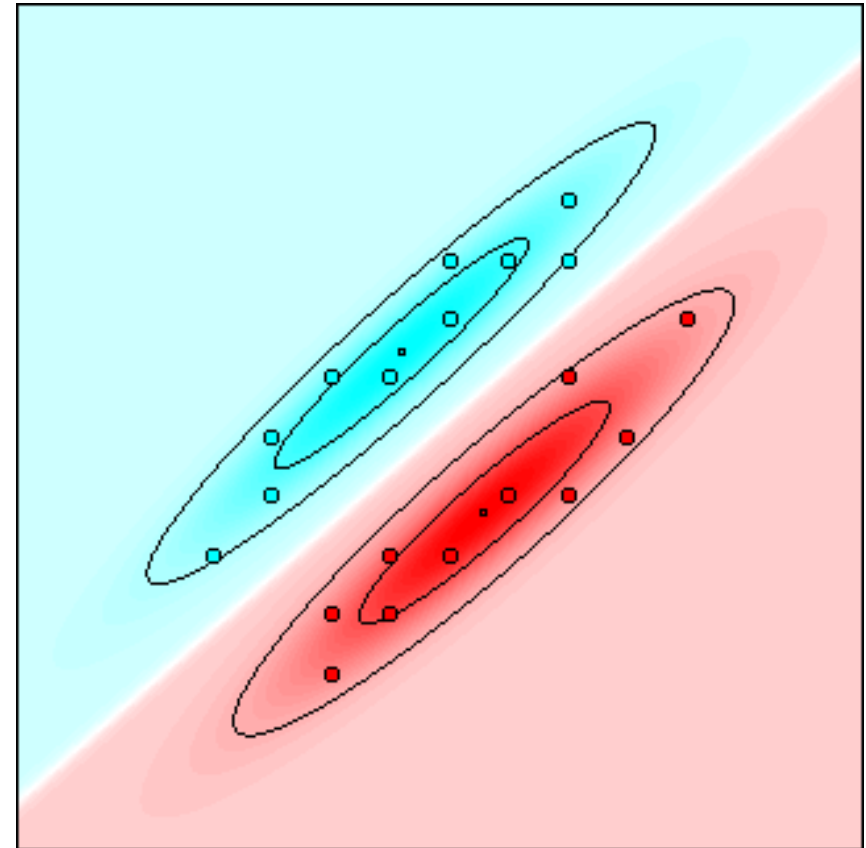
$$\begin{aligned} & f(A_1 = a_1, \dots, A_m = a_m \mid C = c_i) \\ &= \frac{1}{\sqrt{(2\pi)^m |\mathbf{\Sigma}_i|}} \cdot \exp\left(-\frac{1}{2}(\vec{a} - \vec{\mu}_i)^\top \mathbf{\Sigma}_i^{-1}(\vec{a} - \vec{\mu}_i)\right) \\ &= \frac{1}{\sqrt{(2\pi)^m \prod_{k=1}^m \sigma_{i,k}^2}} \cdot \exp\left(-\frac{1}{2}(\vec{a} - \vec{\mu}_i)^\top \text{diag}(\sigma_{i,1}^{-2}, \dots, \sigma_{i,m}^{-2}) (\vec{a} - \vec{\mu}_i)\right) \\ &= \frac{1}{\prod_{k=1}^m \sqrt{2\pi\sigma_{i,k}^2}} \cdot \exp\left(-\frac{1}{2} \sum_{k=1}^m \frac{(a_k - \mu_{i,k})^2}{\sigma_{i,k}^2}\right) \\ &= \prod_{k=1}^m \frac{1}{\sqrt{2\pi\sigma_{i,k}^2}} \cdot \exp\left(-\frac{(a_k - \mu_{i,k})^2}{2\sigma_{i,k}^2}\right) \hat{=} \prod_{k=1}^m f(A_k = a_k \mid C = c_i), \end{aligned}$$

where $f(A_k = a_k \mid C = c_i)$ are the density functions of a naive Bayes classifier.

Comparison of Naive and Full Bayes Classifiers



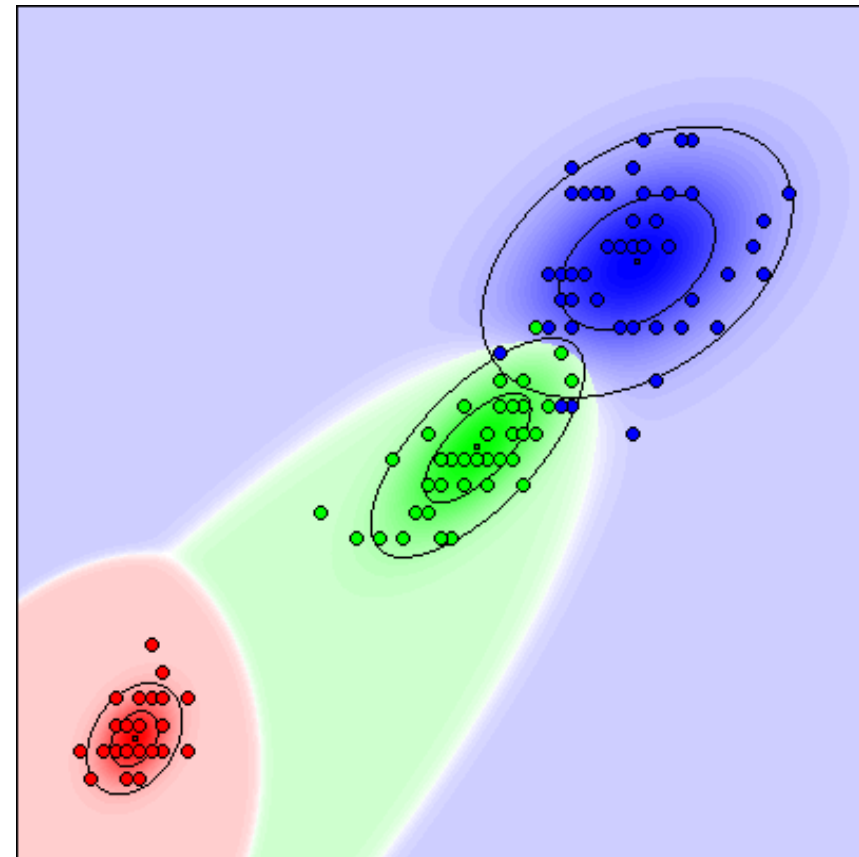
Naive Bayes Classifier



Full Bayes Classifier

Full Bayes Classifiers: Iris Data

- 150 data points, 3 classes
 - Iris setosa (red)
 - Iris versicolor (green)
 - Iris virginica (blue)
- Shown: 2 out of 4 attributes
 - sepal length
 - sepal width
 - petal length (horizontal)
 - petal width (vertical)
- 2 misclassifications on the training data (with all 4 attributes)



Full Bayes Classifier

Summary Bayes Classifiers

- **Probabilistic Classification:** Assign the most probable class.
- **Bayes' Rule:** “Invert” the conditional class probabilities.
- **Naive Bayes Classifiers**
 - Simplifying Assumption:
Attributes are conditionally independent given the class.
 - Can handle nominal/symbolic as well as metric/numeric attributes.
- **Full Bayes Classifiers**
 - Simplifying Assumption:
Each class can be described by a multivariate normal distribution.
 - Can handle only metric/numeric attributes.

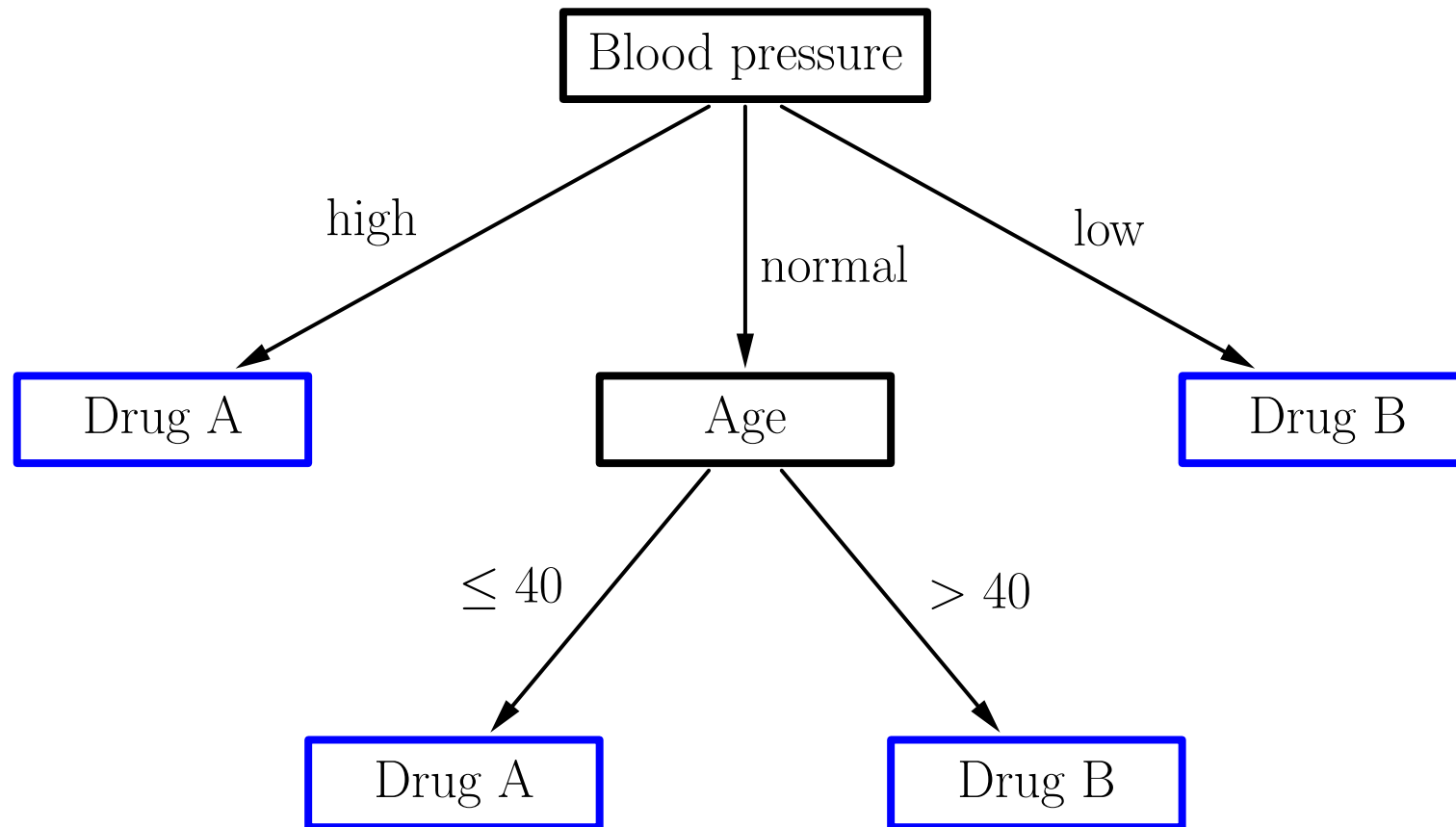
Decision and Regression Trees

Decision and Regression Trees

- **Classification with a Decision Tree**
- **Top-down Induction of Decision Trees**
 - A simple example
 - The general algorithm
 - Attribute selection measures
 - Treatment of numeric attributes and missing values
- **Pruning Decision Trees**
 - General approaches
 - A simple example
- **Regression Trees**
- **Summary**

A Very Simple Decision Tree

Assignment of a drug to a patient:



Classification with a Decision Tree

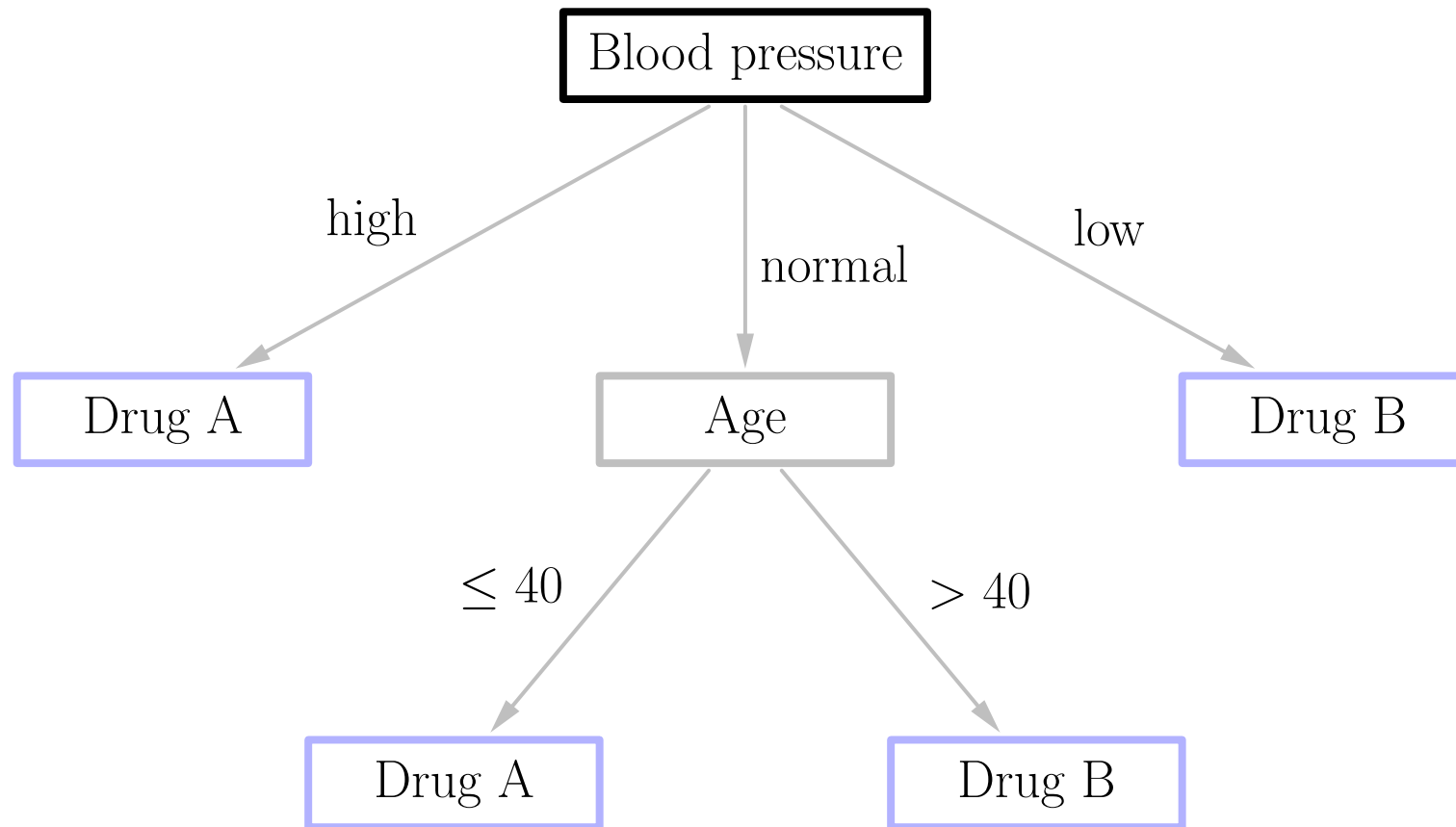
Recursive Descent:

- Start at the root node.
- If the current node is an **leaf node**:
 - Return the class assigned to the node.
- If the current node is an **inner node**:
 - Test the attribute associated with the node.
 - Follow the branch labeled with the outcome of the test.
 - Apply the algorithm recursively.

Intuitively: Follow the path corresponding to the case to be classified.

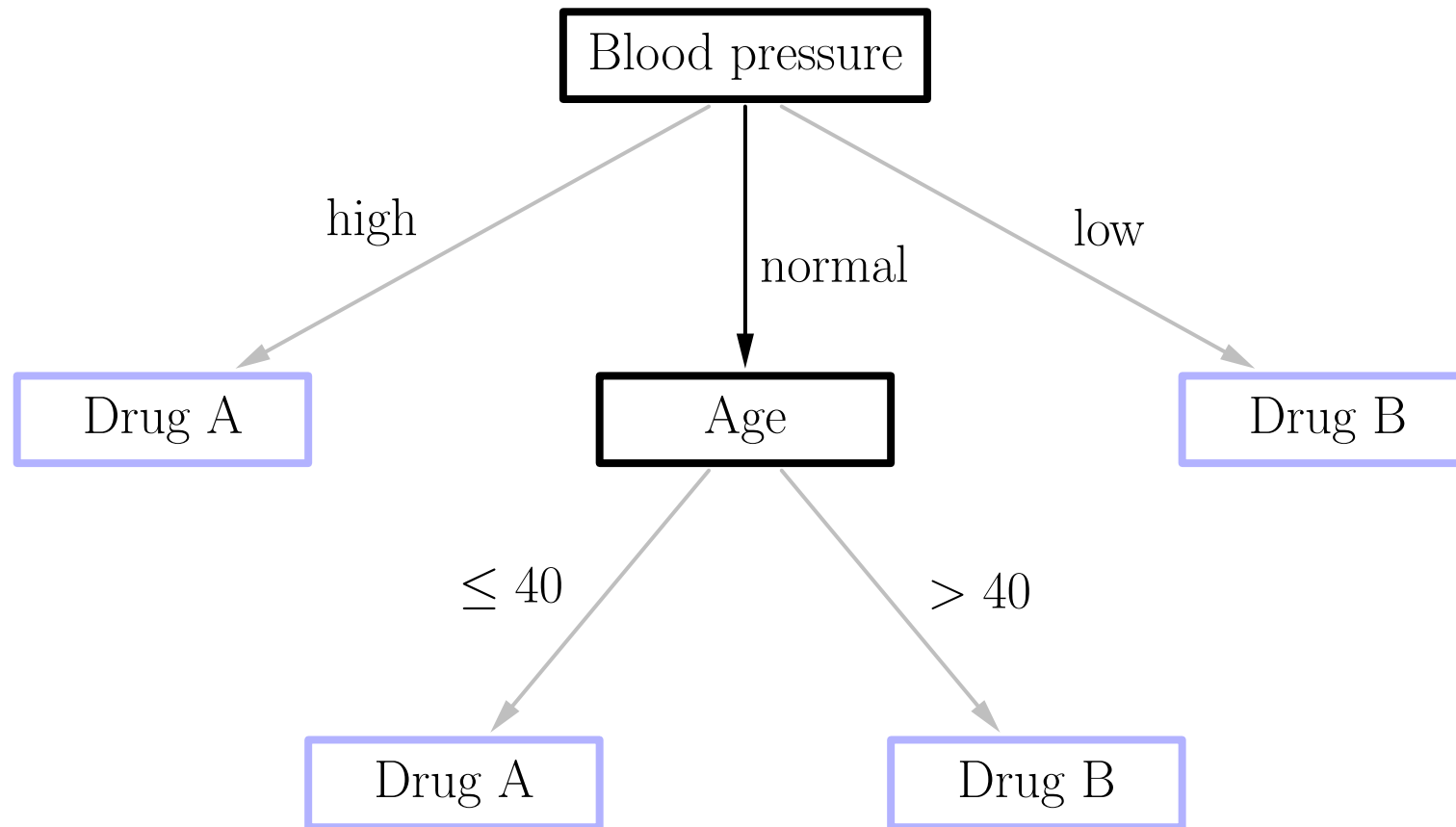
Classification in the Example

Assignment of a drug to a patient:



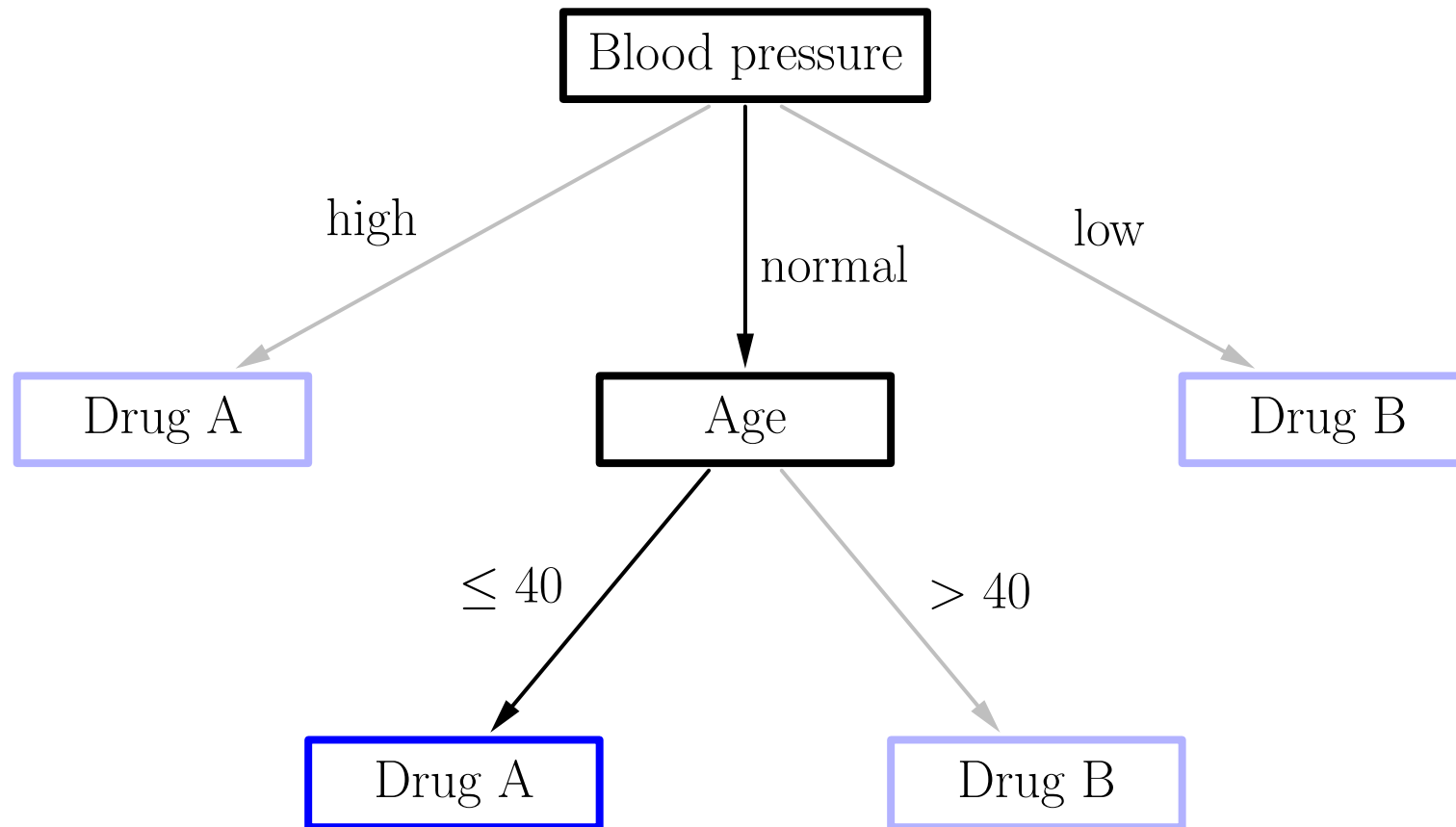
Classification in the Example

Assignment of a drug to a patient:



Classification in the Example

Assignment of a drug to a patient:



Induction of Decision Trees

- **Top-down approach**
 - Build the decision tree from top to bottom (from the root to the leaves).
- **Greedy Selection of a Test Attribute**
 - Compute an evaluation measure for all attributes.
 - Select the attribute with the best evaluation.
- **Divide and Conquer / Recursive Descent**
 - Divide the example cases according to the values of the test attribute.
 - Apply the procedure recursively to the subsets.
 - Terminate the recursion if
 - all cases belong to the same class
 - no more test attributes are available

Induction of a Decision Tree: Example

Patient database

- 12 example cases
- 3 descriptive attributes
- 1 class attribute

Assignment of drug

(without patient attributes)

always drug A or always drug B:

50% correct (in 6 of 12 cases)

No	Sex	Age	Blood pr.	Drug
1	male	20	normal	A
2	female	73	normal	B
3	female	37	high	A
4	male	33	low	B
5	female	48	high	A
6	male	29	normal	A
7	female	52	normal	B
8	male	42	low	B
9	male	61	normal	B
10	female	30	normal	A
11	female	26	low	B
12	male	54	high	A

Induction of a Decision Tree: Example

Sex of the patient

- Division w.r.t. male/female.

Assignment of drug

male: 50% correct (in 3 of 6 cases)

female: 50% correct (in 3 of 6 cases)

total: **50% correct** (in 6 of 12 cases)

No	Sex	Drug
1	male	A
6	male	A
12	male	A
4	male	B
8	male	B
9	male	B
3	female	A
5	female	A
10	female	A
2	female	B
7	female	B
11	female	B

Induction of a Decision Tree: Example

Age of the patient

- Sort according to age.
- Find best age split.
here: ca. 40 years

Assignment of drug

≤ 40 : A 67% correct (in 4 of 6 cases)

> 40 : B 67% correct (in 4 of 6 cases)

total: **67% correct** (in 8 of 12 cases)

No	Age	Drug
1	20	A
11	26	B
6	29	A
10	30	A
4	33	B
3	37	A
8	42	B
5	48	A
7	52	B
12	54	A
9	61	B
2	73	B

Induction of a Decision Tree: Example

Blood pressure of the patient

- Division w.r.t. high/normal/low.

Assignment of drug

high: A 100% correct (in 3 of 3 cases)

normal: 50% correct (in 3 of 6 cases)

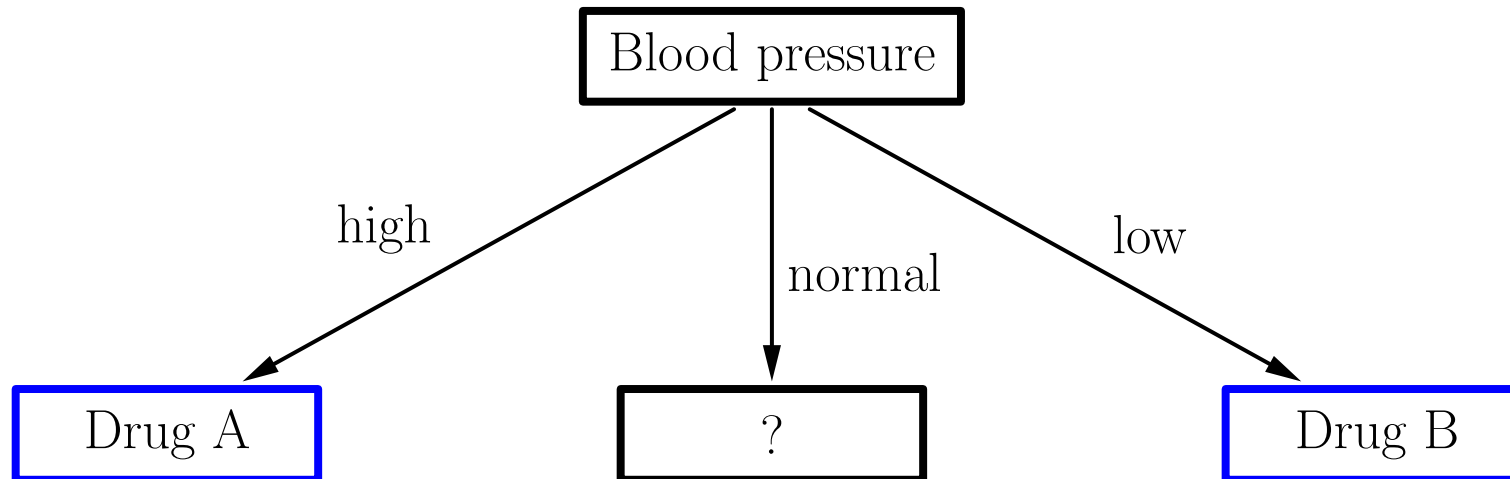
low: B 100% correct (in 3 of 3 cases)

total: **75% correct** (in 9 of 12 cases)

No	Blood pr.	Drug
3	high	A
5	high	A
12	high	A
1	normal	A
6	normal	A
10	normal	A
2	normal	B
7	normal	B
9	normal	B
4	low	B
8	low	B
11	low	B

Induction of a Decision Tree: Example

Current Decision Tree:



Induction of a Decision Tree: Example

Blood pressure and sex

- Only patients with normal blood pressure.
- Division w.r.t. male/female.

Assignment of drug

male: A 67% correct (2 of 3)

female: B 67% correct (2 of 3)

total: **67% correct** (4 of 6)

No	Blood pr.	Sex	Drug
3	high		A
5	high		A
12	high		A
1	normal	male	A
6	normal	male	A
9	normal	male	B
2	normal	female	B
7	normal	female	B
10	normal	female	A
4	low		B
8	low		B
11	low		B

Induction of a Decision Tree: Example

Blood pressure and age

- Only patients with normal blood pressure.
- Sort according to age.
- Find best age split.
here: ca. 40 years

Assignment of drug

≤ 40 : A 100% correct (3 of 3)

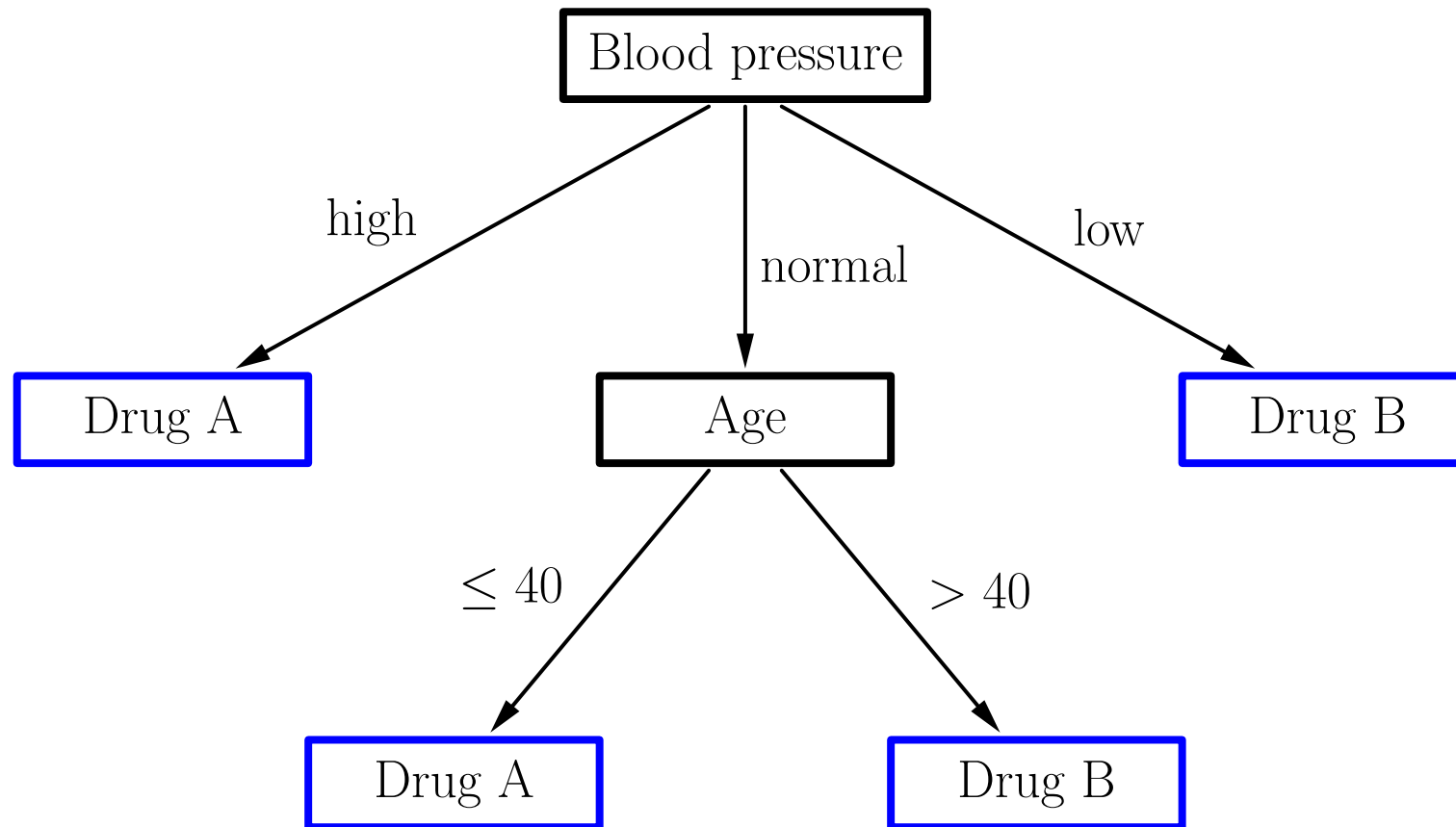
> 40 : B 100% correct (3 of 3)

total: **100% correct** (6 of 6)

No	Blood pr.	Age	Drug
3	high		A
5	high		A
12	high		A
1	normal	20	A
6	normal	29	A
10	normal	30	A
7	normal	52	B
9	normal	61	B
2	normal	73	B
11	low		B
4	low		B
8	low		B

Result of Decision Tree Induction

Assignment of a drug to a patient:



Decision Tree Induction: Notation

S	a set of case or object descriptions
C	the class attribute
$A^{(1)}, \dots, A^{(m)}$	other attributes (index dropped in the following)
$\text{dom}(C)$	$= \{c_1, \dots, c_{n_C}\}, \quad n_C$: number of classes
$\text{dom}(A)$	$= \{a_1, \dots, a_{n_A}\}, \quad n_A$: number of attribute values
$N_{..}$	total number of case or object descriptions i.e. $N_{..} = S $
$N_{i.}$	absolute frequency of the class c_i
$N_{.j}$	absolute frequency of the attribute value a_j
N_{ij}	absolute frequency of the combination of the class c_i and the attribute value a_j . It is $N_{i.} = \sum_{j=1}^{n_A} N_{ij}$ and $N_{.j} = \sum_{i=1}^{n_C} N_{ij}$.
$p_{i.}$	relative frequency of the class c_i , $p_{i.} = \frac{N_{i.}}{N_{..}}$
$p_{.j}$	relative frequency of the attribute value a_j , $p_{.j} = \frac{N_{.j}}{N_{..}}$
p_{ij}	relative frequency of the combination of class c_i and attribute value a_j , $p_{ij} = \frac{N_{ij}}{N_{..}}$
$p_{i j}$	relative frequency of the class c_i in cases having attribute value a_j , $p_{i j} = \frac{N_{ij}}{N_{.j}} = \frac{p_{ij}}{p_{.j}}$

Decision Tree Induction: General Algorithm

```
function grow_tree ( $S$  : set of cases) : node;  
begin  
     $best\_v :=$  WORTHLESS;  
    for all untested attributes  $A$  do  
        compute frequencies  $N_{ij}$ ,  $N_{i.}$ ,  $N_{.j}$  for  $1 \leq i \leq n_C$  and  $1 \leq j \leq n_A$ ;  
        compute value  $v$  of an evaluation measure using  $N_{ij}$ ,  $N_{i.}$ ,  $N_{.j}$ ;  
        if  $v > best\_v$  then  $best\_v := v$ ;  $best\_A := A$ ; end;  
    end  
    if  $best\_v =$  WORTHLESS  
    then create leaf node  $x$ ;  
        assign majority class of  $S$  to  $x$ ;  
    else create test node  $x$ ;  
        assign test on attribute  $best\_A$  to  $x$ ;  
        for all  $a \in \text{dom}(best\_A)$  do  $x.\text{child}[a] := \text{grow\_tree}(S|_{best\_A=a})$ ; end;  
    end;  
    return  $x$ ;  
end;
```

Evaluation Measures

- Evaluation measure used in the above example:
rate of correctly classified example cases.
 - Advantage: simple to compute, easy to understand.
 - Disadvantage: works well only for two classes.
- If there are more than two classes, the rate of misclassified example cases **neglects a lot of the available information.**
 - Only the majority class—that is, the class occurring most often in (a subset of) the example cases—is really considered.
 - The distribution of the other classes has no influence. However, a good choice here can be important for deeper levels of the decision tree.
- **Therefore:** Study also other evaluation measures. Here:
 - **Information gain** and its various normalizations.
 - χ^2 **measure** (well-known in statistics).

An Information-theoretic Evaluation Measure

Information Gain (Kullback and Leibler 1951, Quinlan 1986)

Based on Shannon Entropy $H = - \sum_{i=1}^n p_i \log_2 p_i$ (Shannon 1948)

$$\begin{aligned} I_{\text{gain}}(C, A) &= \overbrace{H(C)} - \overbrace{H(C|A)} \\ &= - \sum_{i=1}^{n_C} p_{i.} \log_2 p_{i.} - \sum_{j=1}^{n_A} p_{.j} \left(- \sum_{i=1}^{n_C} p_{i|j} \log_2 p_{i|j} \right) \end{aligned}$$

$H(C)$ Entropy of the class distribution (C : class attribute)

$H(C|A)$ *Expected entropy* of the class distribution
if the value of the attribute A becomes known

$H(C) - H(C|A)$ Expected entropy reduction or *information gain*

Interpretation of Shannon Entropy

- Let $S = \{s_1, \dots, s_n\}$ be a finite set of alternatives having positive probabilities $P(s_i)$, $i = 1, \dots, n$, satisfying $\sum_{i=1}^n P(s_i) = 1$.

- **Shannon Entropy:**

$$H(S) = - \sum_{i=1}^n P(s_i) \log_2 P(s_i)$$

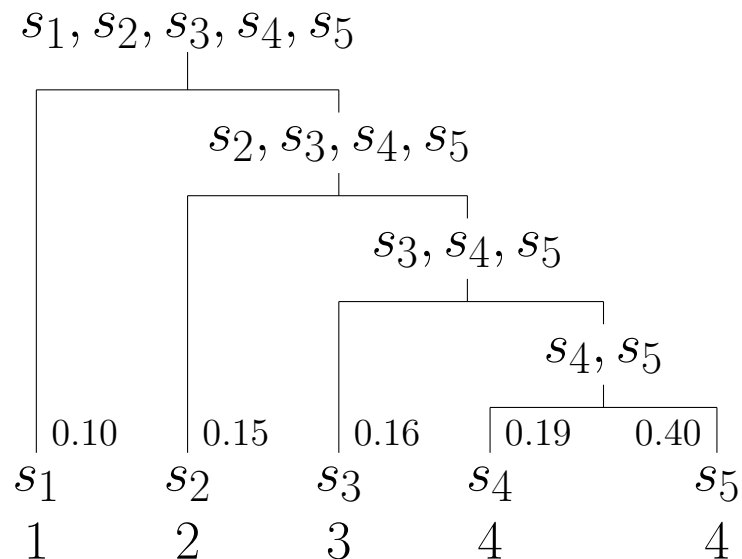
- Intuitively: **Expected number of yes/no questions that have to be asked in order to determine the obtaining alternative.**
 - Suppose there is an oracle, which knows the obtaining alternative, but responds only if the question can be answered with “yes” or “no”.
 - A better question scheme than asking for one alternative after the other can easily be found: Divide the set into two subsets of about equal size.
 - Ask for containment in an arbitrarily chosen subset.
 - Apply this scheme recursively \rightarrow number of questions bounded by $\lceil \log_2 n \rceil$.

Question/Coding Schemes

$$P(s_1) = 0.10, \quad P(s_2) = 0.15, \quad P(s_3) = 0.16, \quad P(s_4) = 0.19, \quad P(s_5) = 0.40$$

$$\text{Shannon entropy: } -\sum_i P(s_i) \log_2 P(s_i) = 2.15 \text{ bit/symbol}$$

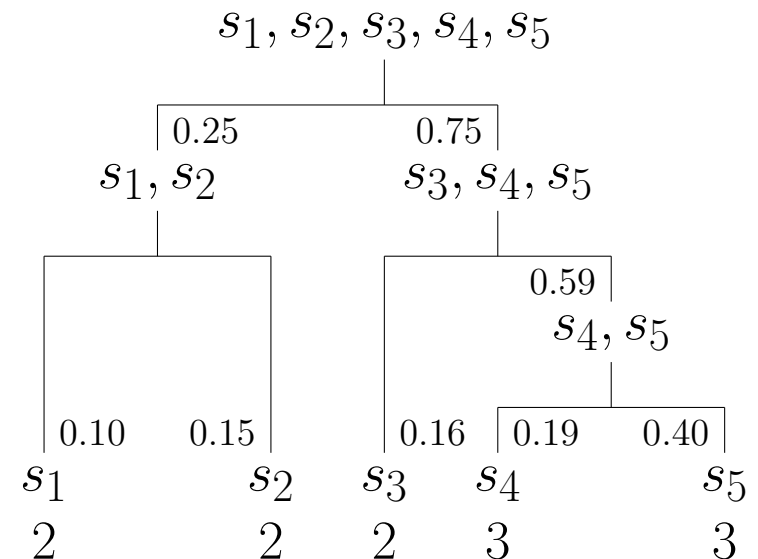
Linear Traversal



Code length: 3.24 bit/symbol

Code efficiency: 0.664

Equal Size Subsets



Code length: 2.59 bit/symbol

Code efficiency: 0.830

Question/Coding Schemes

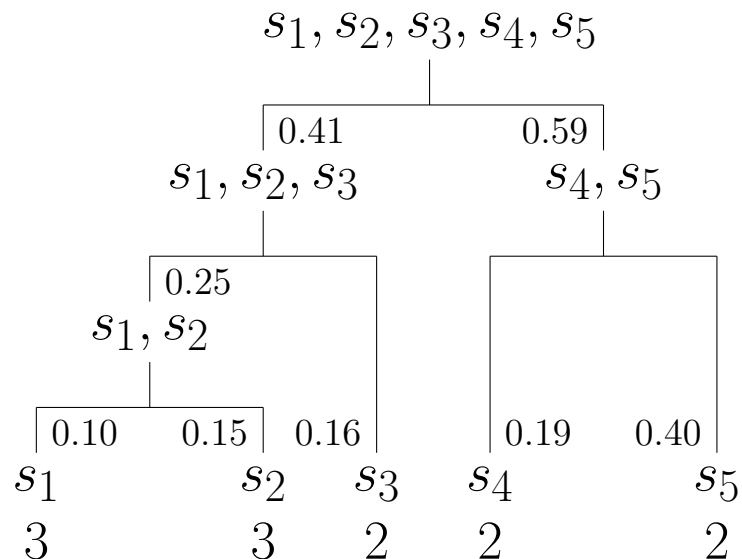
- Splitting into subsets of about equal size can lead to a bad arrangement of the alternatives into subsets → high expected number of questions.
- Good question schemes take the probability of the alternatives into account.
- **Shannon-Fano Coding** (1948)
 - Build the question/coding scheme top-down.
 - Sort the alternatives w.r.t. their probabilities.
 - Split the set so that the subsets have about equal *probability* (splits must respect the probability order of the alternatives).
- **Huffman Coding** (1952)
 - Build the question/coding scheme bottom-up.
 - Start with one element sets.
 - Always combine those two sets that have the smallest probabilities.

Question/Coding Schemes

$$P(s_1) = 0.10, \quad P(s_2) = 0.15, \quad P(s_3) = 0.16, \quad P(s_4) = 0.19, \quad P(s_5) = 0.40$$

$$\text{Shannon entropy: } -\sum_i P(s_i) \log_2 P(s_i) = 2.15 \text{ bit/symbol}$$

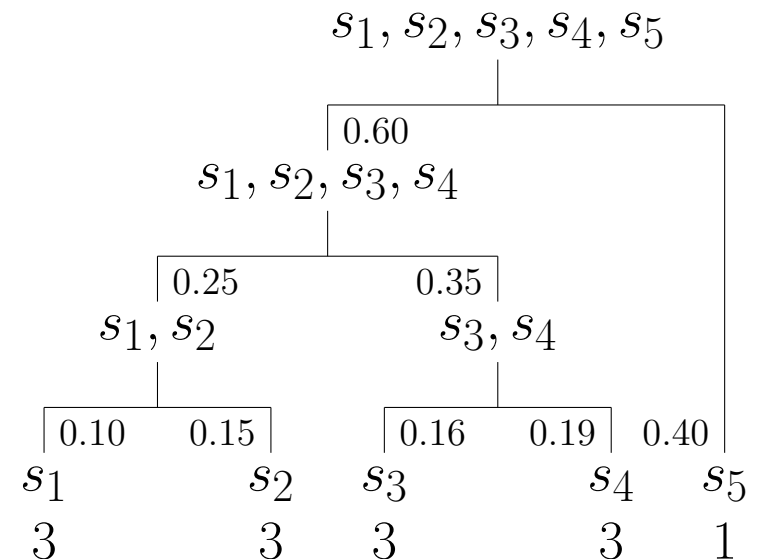
Shannon–Fano Coding (1948)



Code length: 2.25 bit/symbol

Code efficiency: 0.955

Huffman Coding (1952)



Code length: 2.20 bit/symbol

Code efficiency: 0.977

Question/Coding Schemes

- It can be shown that Huffman coding is optimal if we have to determine the obtaining alternative in a single instance.
(No question/coding scheme has a smaller expected number of questions.)
- Only if the obtaining alternative has to be determined in a sequence of (independent) situations, this scheme can be improved upon.
- Idea: Process the sequence not instance by instance, but combine two, three or more consecutive instances and ask directly for the obtaining combination of alternatives.
- Although this enlarges the question/coding scheme, the expected number of questions per identification is reduced (because each interrogation identifies the obtaining alternative for several situations).
- However, the expected number of questions per identification cannot be made arbitrarily small. Shannon showed that there is a lower bound, namely the Shannon entropy.

Interpretation of Shannon Entropy

$$P(s_1) = \frac{1}{2}, \quad P(s_2) = \frac{1}{4}, \quad P(s_3) = \frac{1}{8}, \quad P(s_4) = \frac{1}{16}, \quad P(s_5) = \frac{1}{16}$$

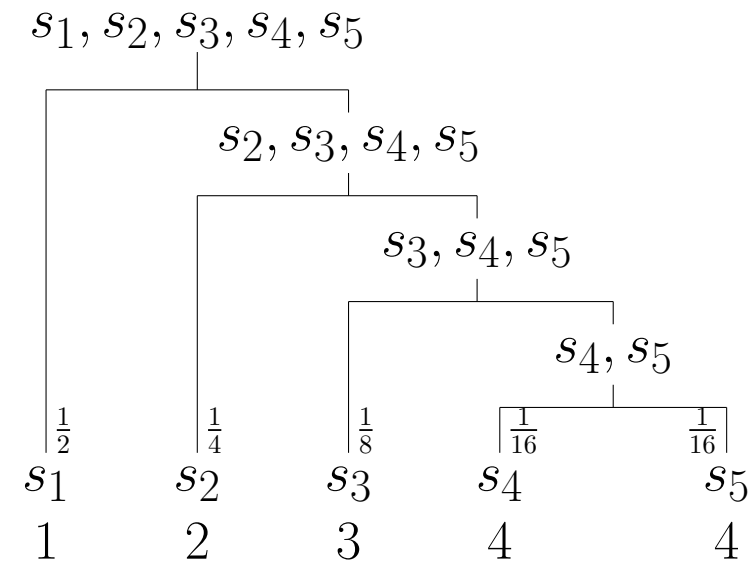
$$\text{Shannon entropy: } -\sum_i P(s_i) \log_2 P(s_i) = 1.875 \text{ bit/symbol}$$

If the probability distribution allows for a perfect Huffman code (code efficiency 1), the Shannon entropy can easily be interpreted as follows:

$$\begin{aligned} & -\sum_i P(s_i) \log_2 P(s_i) \\ &= \sum_i \underbrace{P(s_i)}_{\text{occurrence probability}} \cdot \underbrace{\log_2 \frac{1}{P(s_i)}}_{\text{path length in tree}}. \end{aligned}$$

In other words, it is the expected number of needed yes/no questions.

Perfect Question Scheme



Code length: 1.875 bit/symbol

Code efficiency: 1

Other Information-theoretic Evaluation Measures

Normalized Information Gain

- Information gain is biased towards many-valued attributes.
- Normalization removes / reduces this bias.

Information Gain Ratio (Quinlan 1986 / 1993)

$$I_{\text{gr}}(C, A) = \frac{I_{\text{gain}}(C, A)}{H_A} = \frac{I_{\text{gain}}(C, A)}{-\sum_{j=1}^{n_A} p_{.j} \log_2 p_{.j}}$$

Symmetric Information Gain Ratio (López de Mántaras 1991)

$$I_{\text{sgr}}^{(1)}(C, A) = \frac{I_{\text{gain}}(C, A)}{H_{AC}} \quad \text{or} \quad I_{\text{sgr}}^{(2)}(C, A) = \frac{I_{\text{gain}}(C, A)}{H_A + H_C}$$

Bias of Information Gain

- **Information gain is biased towards many-valued attributes**, i.e., of two attributes having about the same information content it tends to select the one having more values.
- The reasons are quantization effects caused by the finite number of example cases (due to which only a finite number of different probabilities can result in estimations) in connection with the following theorem:
- **Theorem:** Let A , B , and C be three attributes with finite domains and let their joint probability distribution be strictly positive, i.e., $\forall a \in \text{dom}(A) : \forall b \in \text{dom}(B) : \forall c \in \text{dom}(C) : P(A = a, B = b, C = c) > 0$. Then

$$I_{\text{gain}}(C, AB) \geq I_{\text{gain}}(C, B),$$

with equality obtaining only if the attributes C and A are conditionally independent given B , i.e., if $P(C = c \mid A = a, B = b) = P(C = c \mid B = b)$.

(A detailed proof of this theorem can be found, for example, in [Borgelt and Kruse 2002], p. 311ff.)

A Statistical Evaluation Measure

χ^2 Measure

- Compares the actual joint distribution with a **hypothetical independent distribution**.
- Uses absolute comparison.
- Can be interpreted as a difference measure.

$$\chi^2(C, A) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_A} N_{..} \frac{(p_{i.p.j} - p_{ij})^2}{p_{i.p.j}}$$

- Side remark: Information gain can also be interpreted as a difference measure.

$$I_{\text{gain}}(C, A) = \sum_{i=1}^{n_C} \sum_{j=1}^{n_A} p_{ij} \log_2 \frac{p_{ij}}{p_{i.p.j}}$$

Treatment of Numeric Attributes

General Approach: Discretization

- **Preprocessing I**
 - Form equally sized or equally populated intervals.
- **During the tree construction**
 - Sort the example cases according to the attribute's values.
 - Construct a binary symbolic attribute for every possible split (values: " \leq threshold" and " $>$ threshold").
 - Compute the evaluation measure for these binary attributes.
 - Possible improvements: Add a penalty depending on the number of splits.
- **Preprocessing II / Multisplits during tree construction**
 - Build a decision tree using only the numeric attribute.
 - Flatten the tree to obtain a multi-interval discretization.

Treatment of Missing Values

Induction

- Weight the evaluation measure with the fraction of cases with known values.
 - Idea: The attribute provides information only if it is known.
- Try to find a surrogate test attribute with similar properties (CART, Breiman *et al.* 1984)
- Assign the case to all branches, weighted in each branch with the relative frequency of the corresponding attribute value (C4.5, Quinlan 1993).

Classification

- Use the surrogate test attribute found during induction.
- Follow all branches of the test attribute, weighted with their relative number of cases, aggregate the class distributions of all leaves reached, and assign the majority class of the aggregated class distribution.

Pruning Decision Trees

Pruning serves the purpose

- to simplify the tree (improve interpretability),
- to avoid overfitting (improve generalization).

Basic ideas:

- Replace “bad” branches (subtrees) by leaves.
- Replace a subtree by its largest branch if it is better.

Common approaches:

- Reduced error pruning
- Pessimistic pruning
- Confidence level pruning
- Minimum description length pruning

Reduced Error Pruning

- Classify a set of new example cases with the decision tree.
(These cases must not have been used for the induction!)
- Determine the number of errors for all leaves.
- The number of errors of a subtree is the sum of the errors of all of its leaves.
- Determine the number of errors for leaves that replace subtrees.
- If such a leaf leads to the same or fewer errors than the subtree, replace the subtree by the leaf.
- If a subtree has been replaced, recompute the number of errors of the subtrees it is part of.

Advantage: Very good pruning, effective avoidance of overfitting.

Disadvantage: Additional example cases needed.

Pessimistic Pruning

- Classify a set of example cases with the decision tree.
(These cases may or may not have been used for the induction.)
- Determine the number of errors for all leaves and increase this number by a fixed, user-specified amount r .
- The number of errors of a subtree is the sum of the errors of all of its leaves.
- Determine the number of errors for leaves that replace subtrees (also increased by r).
- If such a leaf leads to the same or fewer errors than the subtree, replace the subtree by the leaf and recompute subtree errors.

Advantage: No additional example cases needed.

Disadvantage: Number of cases in a leaf has no influence.

Confidence Level Pruning

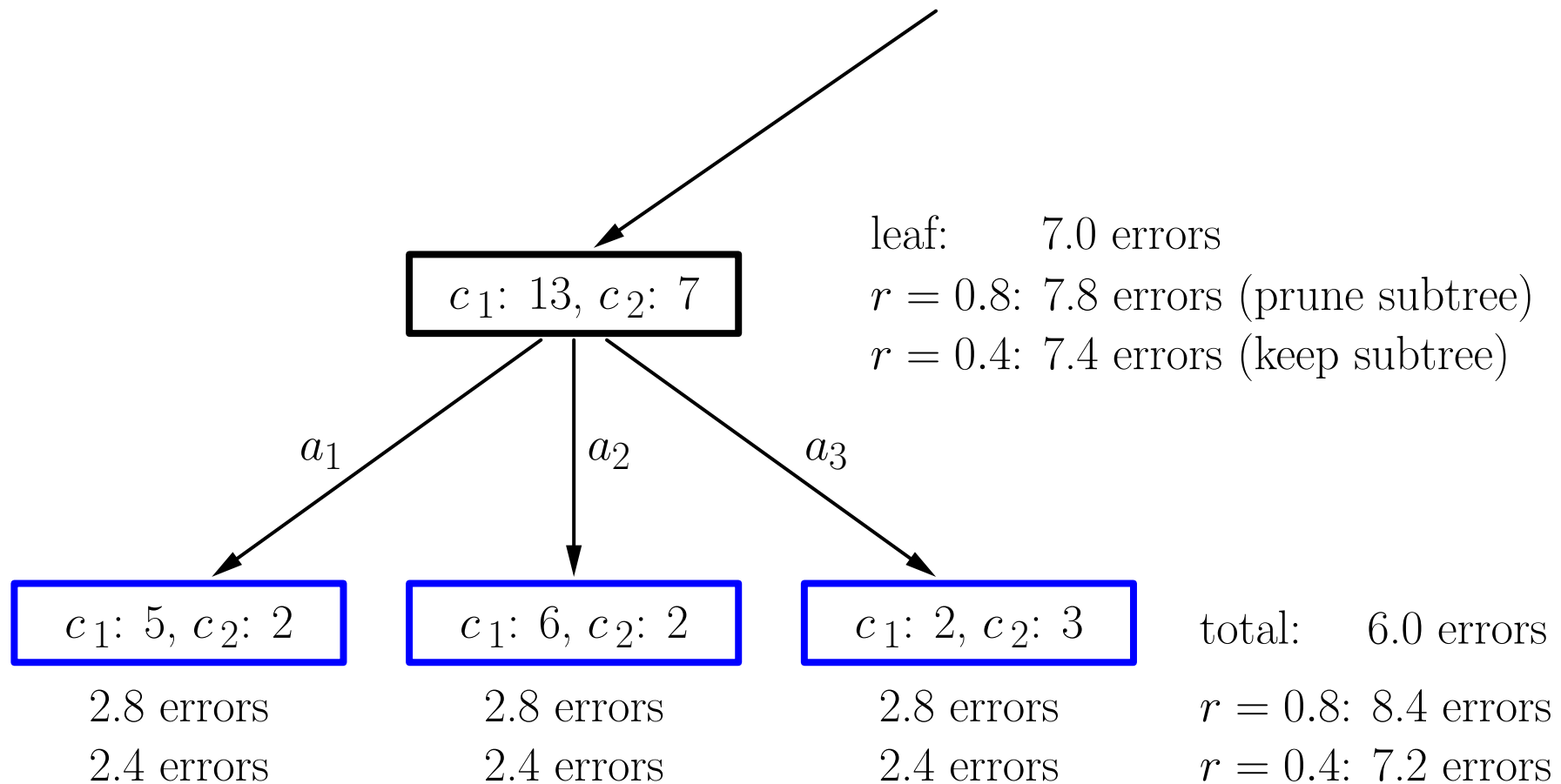
- Like pessimistic pruning, but the number of errors is computed as follows:
 - See classification in a leaf as a Bernoulli experiment (error / no error).
 - Estimate an interval for the error probability based on a user-specified confidence level α .
(use approximation of the binomial distribution by a normal distribution)
 - Increase error number to the upper level of the confidence interval times the number of cases assigned to the leaf.
 - Formal problem: Classification is not a random experiment.

Advantage: No additional example cases needed, good pruning.

Disadvantage: Statistically dubious foundation.

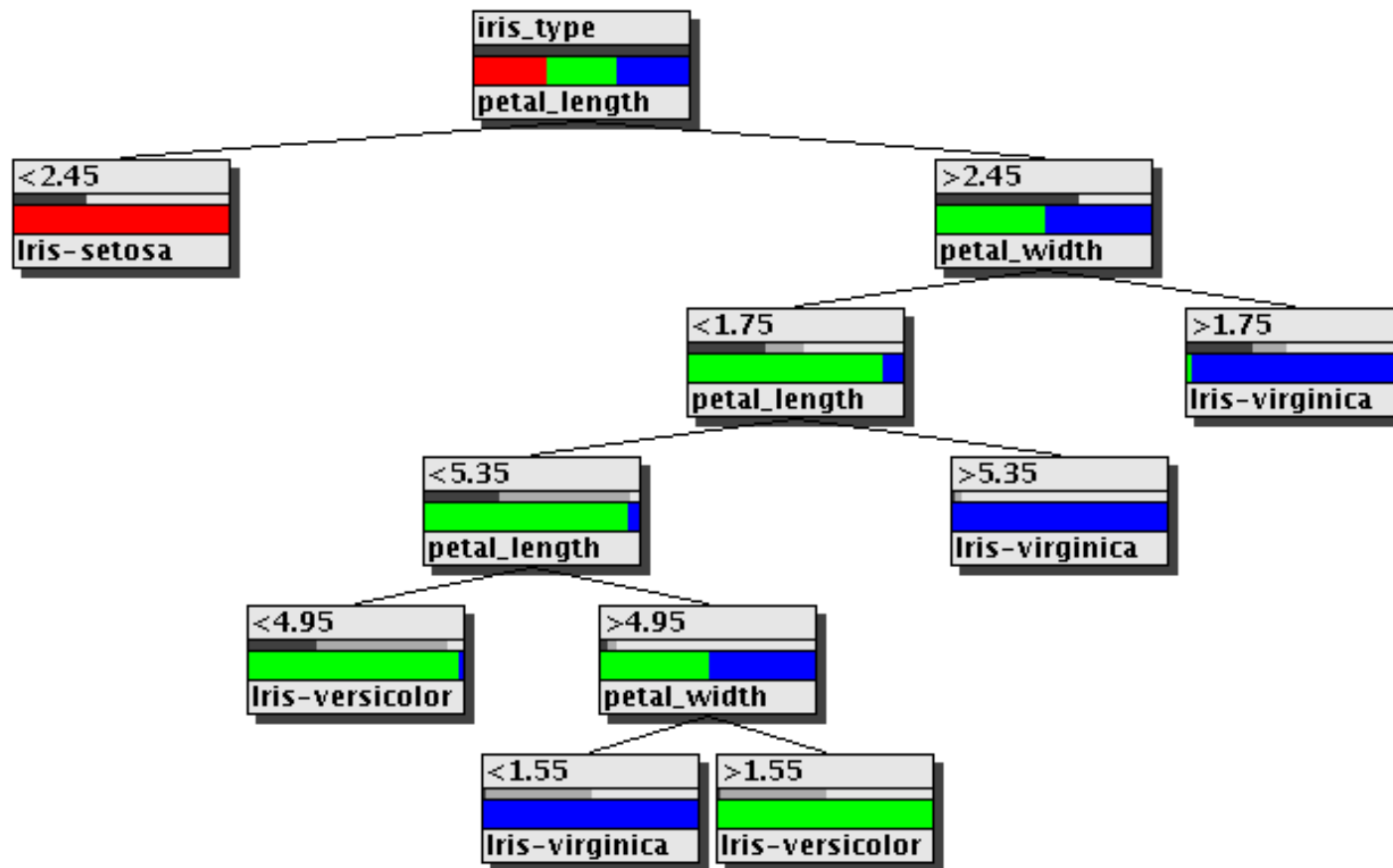
Pruning a Decision Tree: A Simple Example

Pessimistic Pruning with $r = 0.8$ and $r = 0.4$:



Decision Trees: An Example

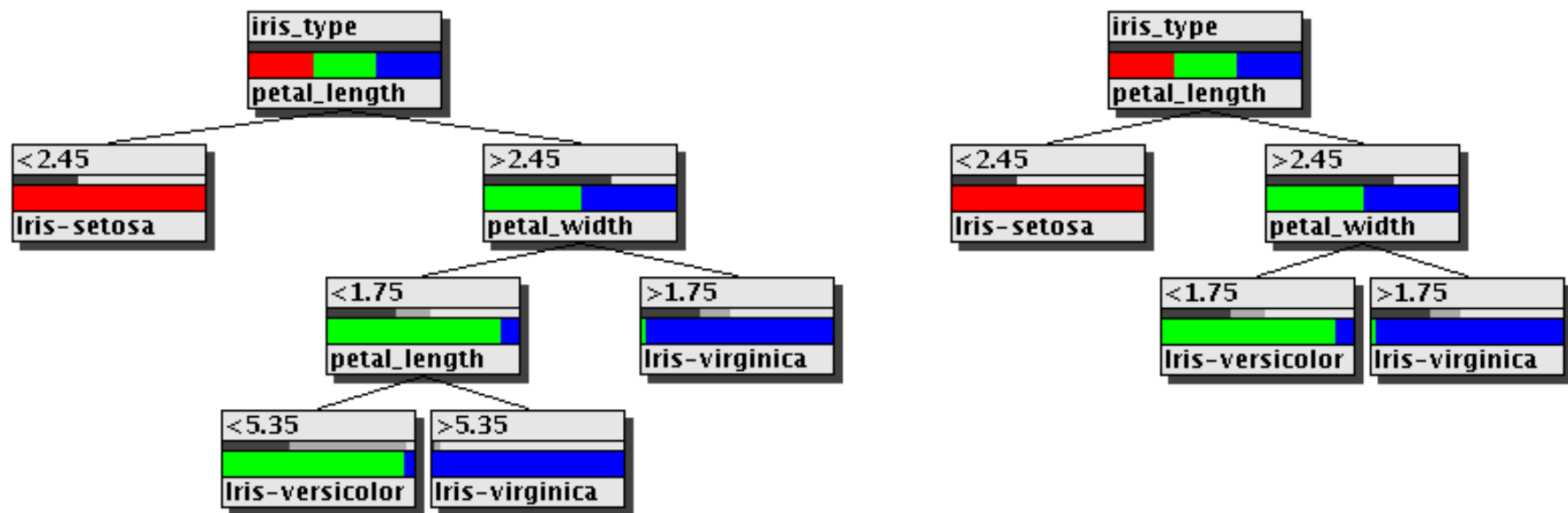
A decision tree for the Iris data
(induced with information gain ratio, unpruned)



Decision Trees: An Example

A decision tree for the Iris data

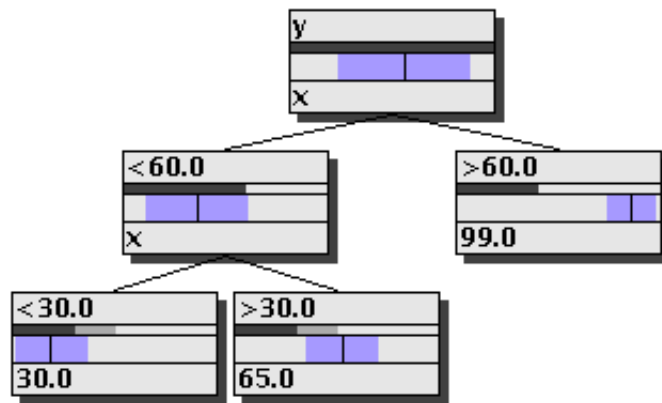
(pruned with confidence level pruning, $\alpha = 0.8$, and pessimistic pruning, $r = 2$)



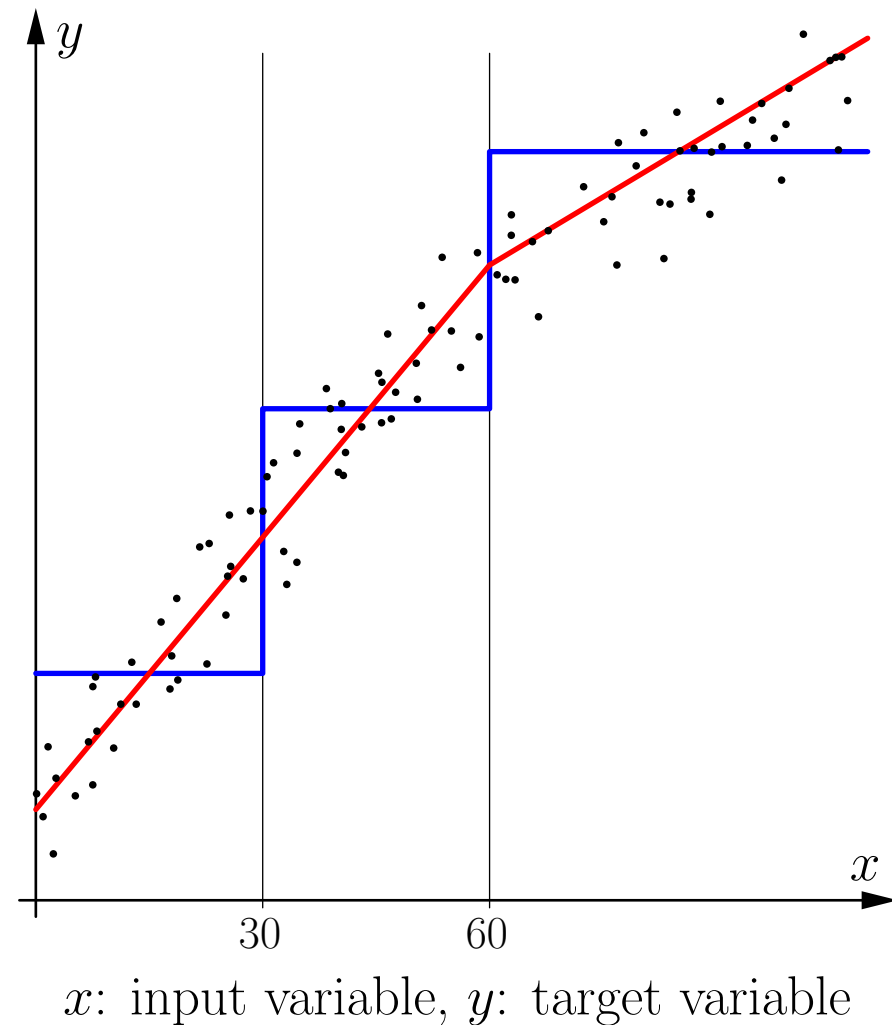
- Left: 7 instead of 11 nodes, 4 instead of 2 misclassifications.
- Right: 5 instead of 11 nodes, 6 instead of 2 misclassifications.
- The right tree is “minimal” for the three classes.

Regression Trees

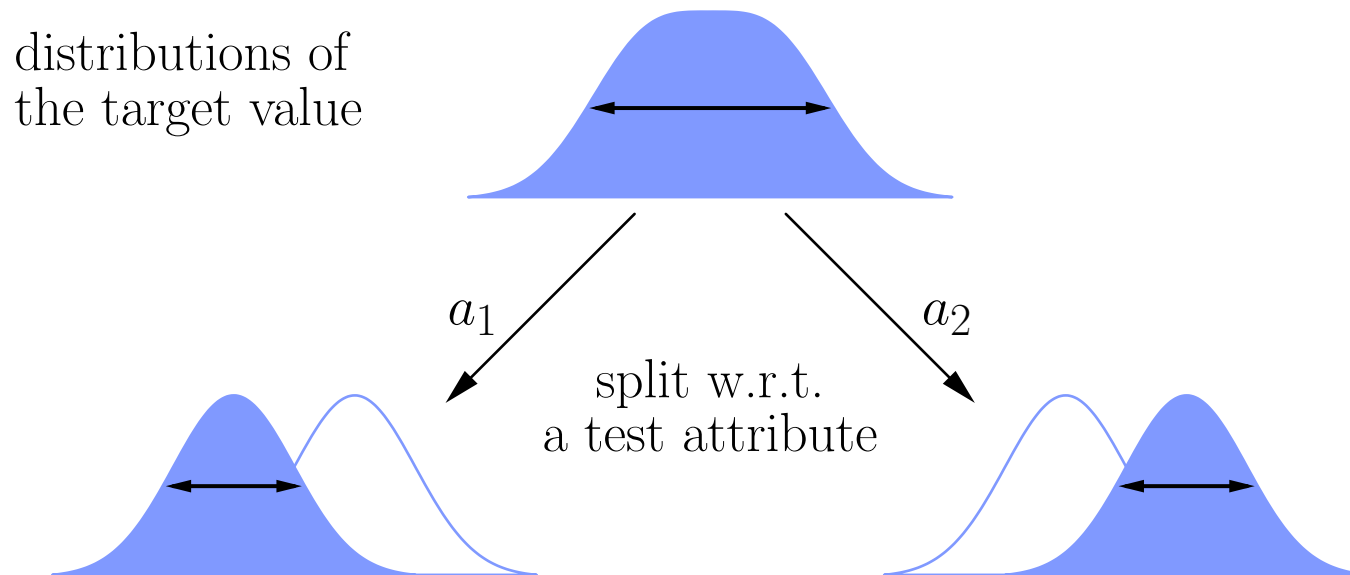
- Target variable is not a class, but a numeric quantity.
- Simple regression trees: predict constant values in leaves. (blue lines)



- More complex regression trees: predict linear functions in leaves. (red line)



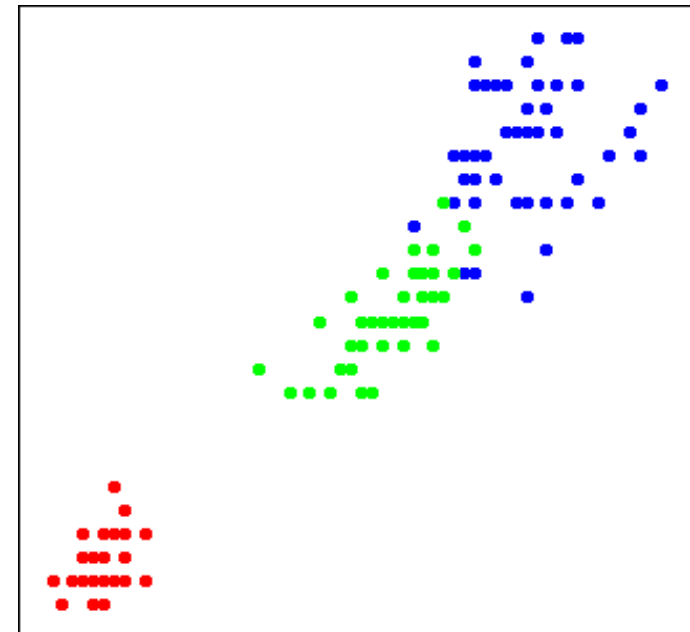
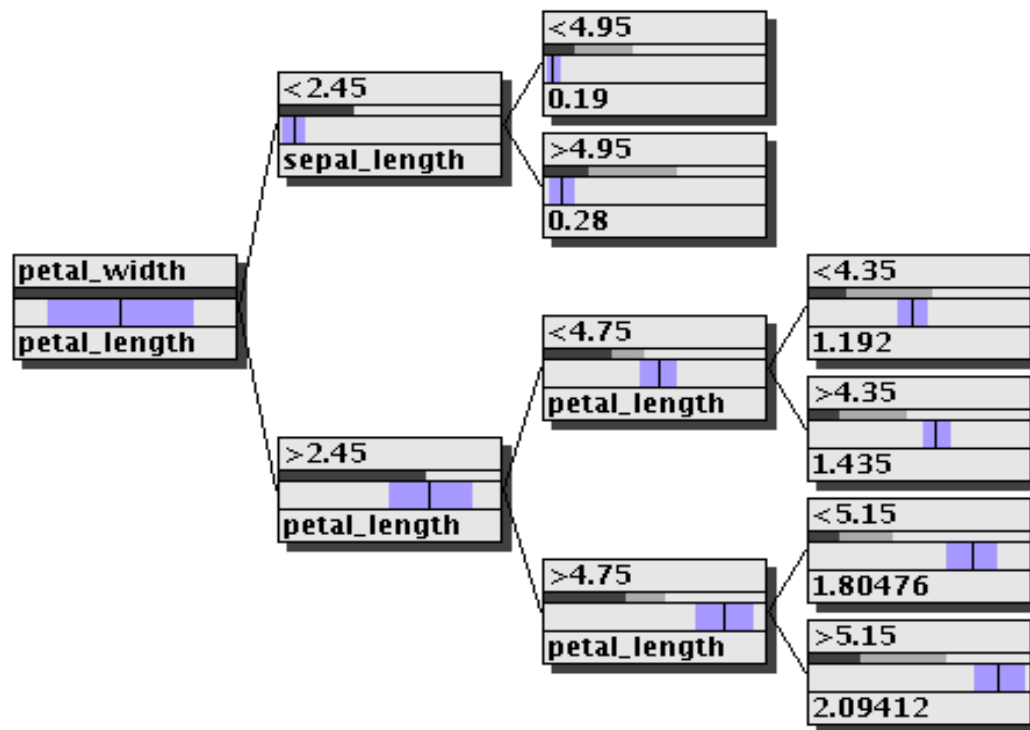
Regression Trees: Attribute Selection



- The variance / standard deviation is compared to the variance / standard deviation in the branches.
- The attribute that yields the highest reduction is selected.

Regression Trees: An Example

A regression tree for the Iris data (petal width)
(induced with reduction of sum of squared errors)



Summary Decision and Regression Trees

- **Decision Trees are Classifiers with Tree Structure**

- Inner node: Test of a descriptive attribute
- Leaf node: Assignment of a class

- **Induction of Decision Trees from Data**

(Top-Down Induction of Decision Trees, TDIDT)

- *Divide and conquer* approach / *recursive descent*
- *Greedy* selection of the test attributes
- Attributes are selected based on an *evaluation measure*, e.g. information gain, χ^2 measure
- Recommended: *Pruning* of the decision tree

- **Numeric Target: Regression Trees**

Classification Evaluation: Cross Validation

- General method to evaluate / predict the performance of classifiers.
- Serves the purpose to estimate the error rate on new example cases.
- Procedure of cross validation:
 - Split the given data set into n so-called *folds* of equal size (n -fold cross validation).
 - Combine $n - 1$ folds into a training data set, build a classifier, and test it on the n -th fold.
 - Do this for all n possible selections of $n - 1$ folds and average the error rates.
- Special case: Leave-1-out cross validation.
(use as many folds as there are example cases)
- Final classifier is learned from the full data set.

Clustering

Clustering

- **General Idea of Clustering**
 - Similarity and distance measures
- **Prototype-based Clustering**
 - Classical c -means clustering
 - Learning vector quantization
 - Fuzzy c -means clustering
 - Expectation maximization for Gaussian mixtures
- **Hierarchical Agglomerative Clustering**
 - Merging clusters: Dendrograms
 - Measuring the distance of clusters
 - Choosing the clusters
- **Summary**

General Idea of Clustering

- Goal: Arrange the given data tuples into **classes** or **clusters**.
- Data tuples assigned to the same cluster should be as similar as possible.
- Data tuples assigned to different clusters should be as dissimilar as possible.
- Similarity is most often measured with the help of a distance function.
(The smaller the distance, the more similar the data tuples.)
- Often: restriction to data points in \mathbb{R}^m (although this is not mandatory).

$d : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}_0^+$ is a **distance function** if it satisfies $\forall \vec{x}, \vec{y}, \vec{z} \in \mathbb{R}^m :$

- (i) $d(\vec{x}, \vec{y}) = 0 \iff \vec{x} = \vec{y},$
- (ii) $d(\vec{x}, \vec{y}) = d(\vec{y}, \vec{x})$ (symmetry),
- (iii) $d(\vec{x}, \vec{z}) \leq d(\vec{x}, \vec{y}) + d(\vec{y}, \vec{z})$ (triangle inequality).

Distance Functions

Illustration of distance functions: Minkowski Family

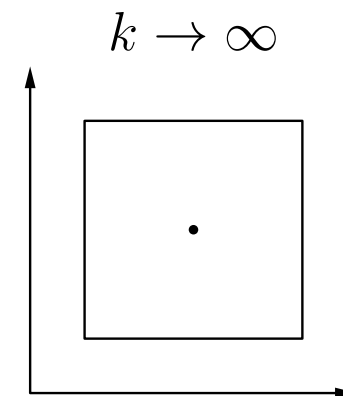
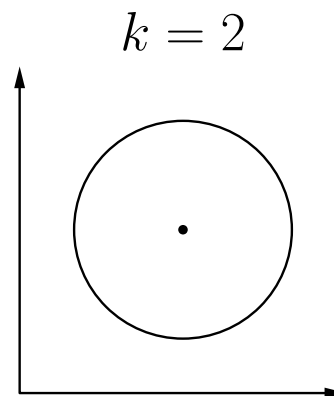
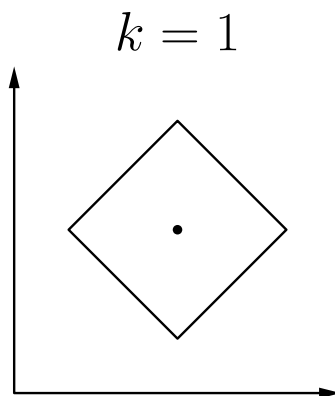
$$d_k(\vec{x}, \vec{y}) = \left(\sum_{i=1}^n (x_i - y_i)^k \right)^{\frac{1}{k}}$$

Well-known special cases from this family are:

$k = 1$: Manhattan or city block distance,

$k = 2$: Euclidean distance,

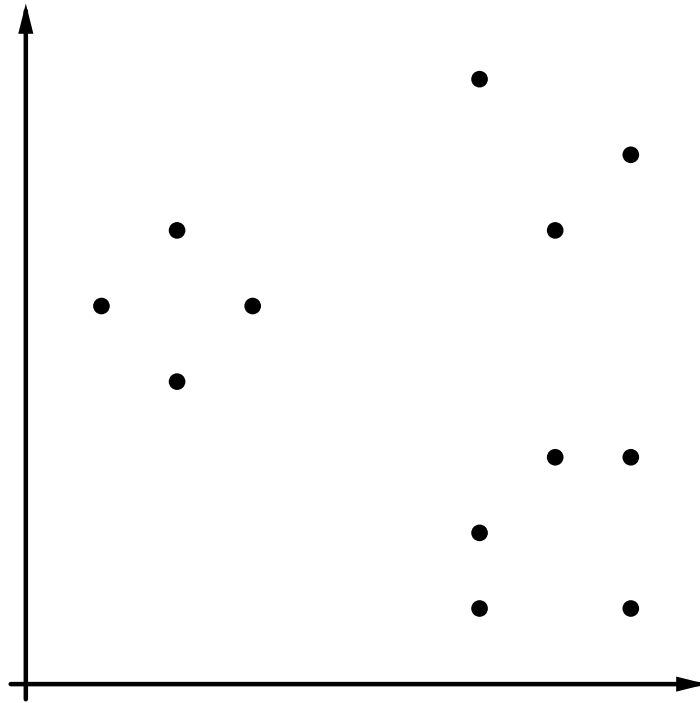
$k \rightarrow \infty$: maximum distance, i.e. $d_\infty(\vec{x}, \vec{y}) = \max_{i=1}^n |x_i - y_i|$.



c -Means Clustering

- Choose a number c of clusters to be found (user input).
- Initialize the cluster centers randomly
(for instance, by randomly selecting c data points).
- **Data point assignment:**
Assign each data point to the cluster center that is closest to it
(i.e. closer than any other cluster center).
- **Cluster center update:**
Compute new cluster centers as the mean vectors of the assigned data points.
(Intuitively: center of gravity if each data point has unit weight.)
- Repeat these two steps (data point assignment and cluster center update)
until the clusters centers do not change anymore.
- It can be shown that this scheme must converge,
i.e., the update of the cluster centers cannot go on forever.

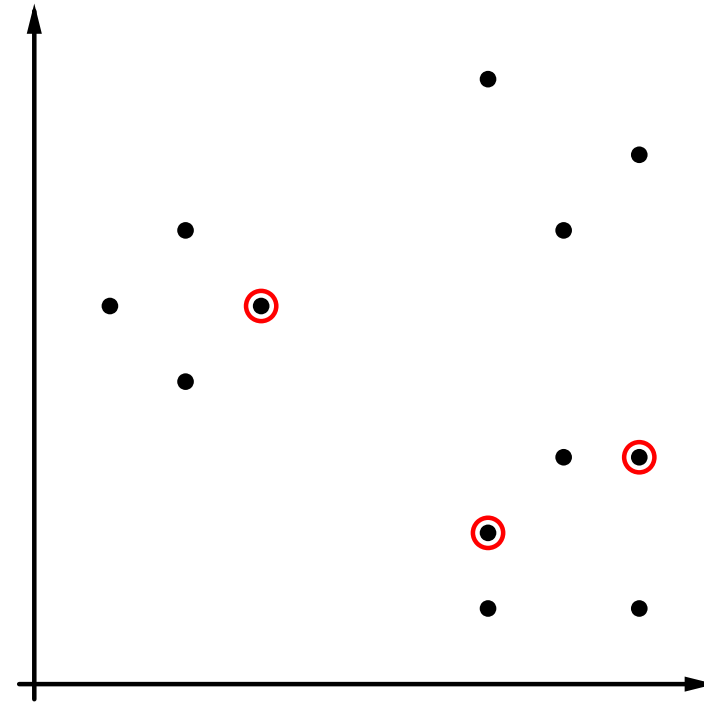
c -Means Clustering: Example



Data set to cluster.

Choose $c = 3$ clusters.

(From visual inspection, can be difficult to determine in general.)

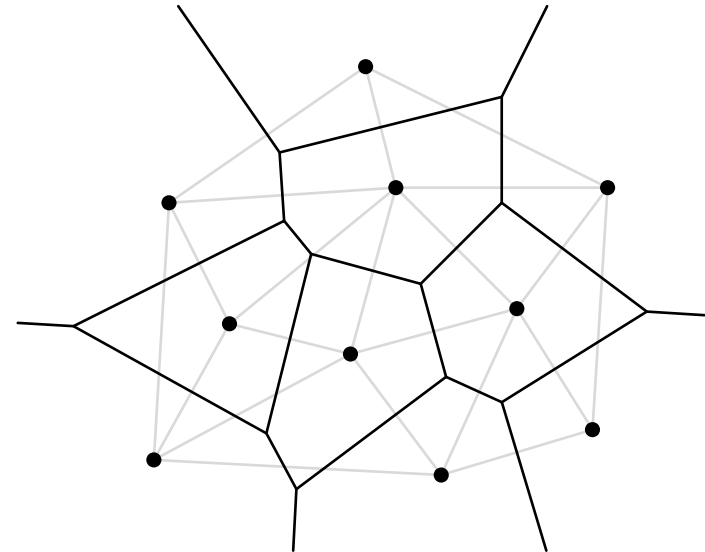
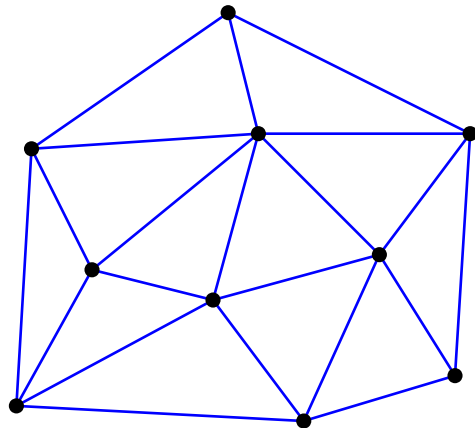


Initial position of cluster centers.

Randomly selected data points.

(Alternative methods include e.g. latin hypercube sampling)

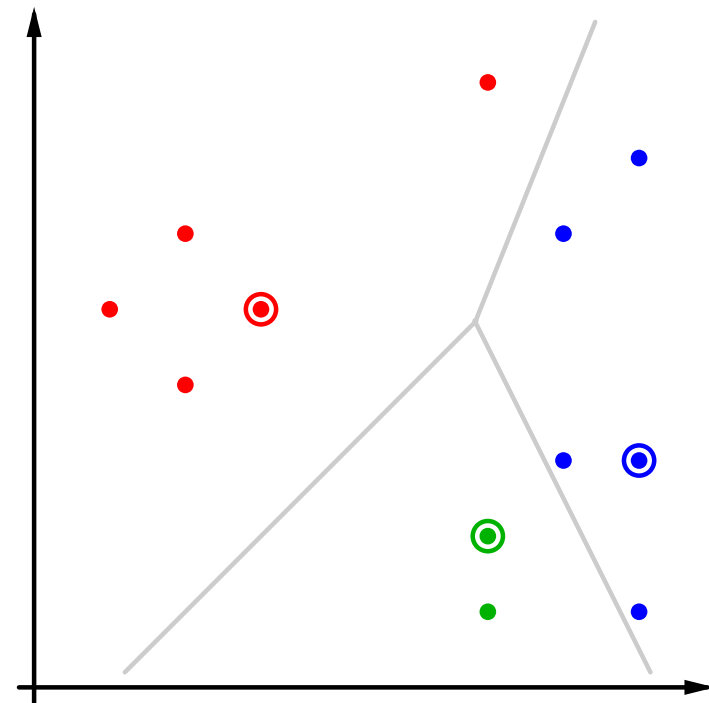
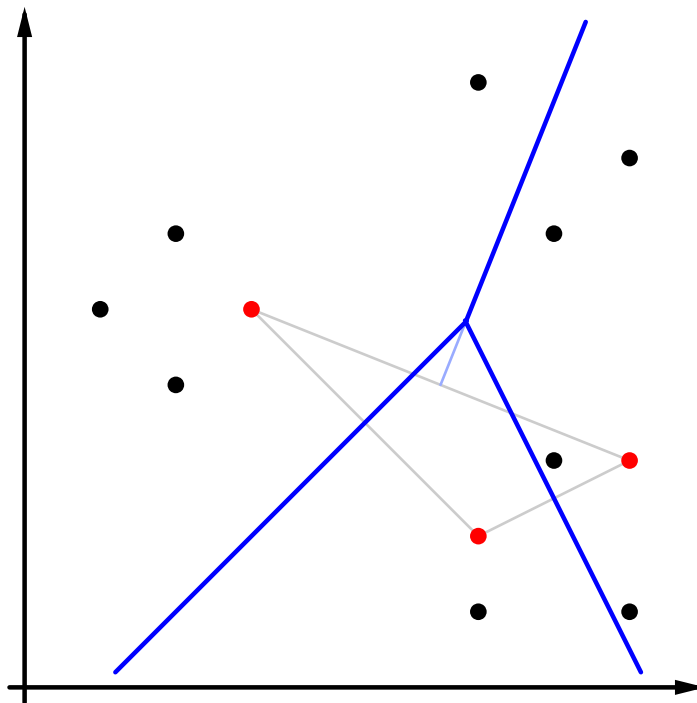
Delaunay Triangulations and Voronoi Diagrams



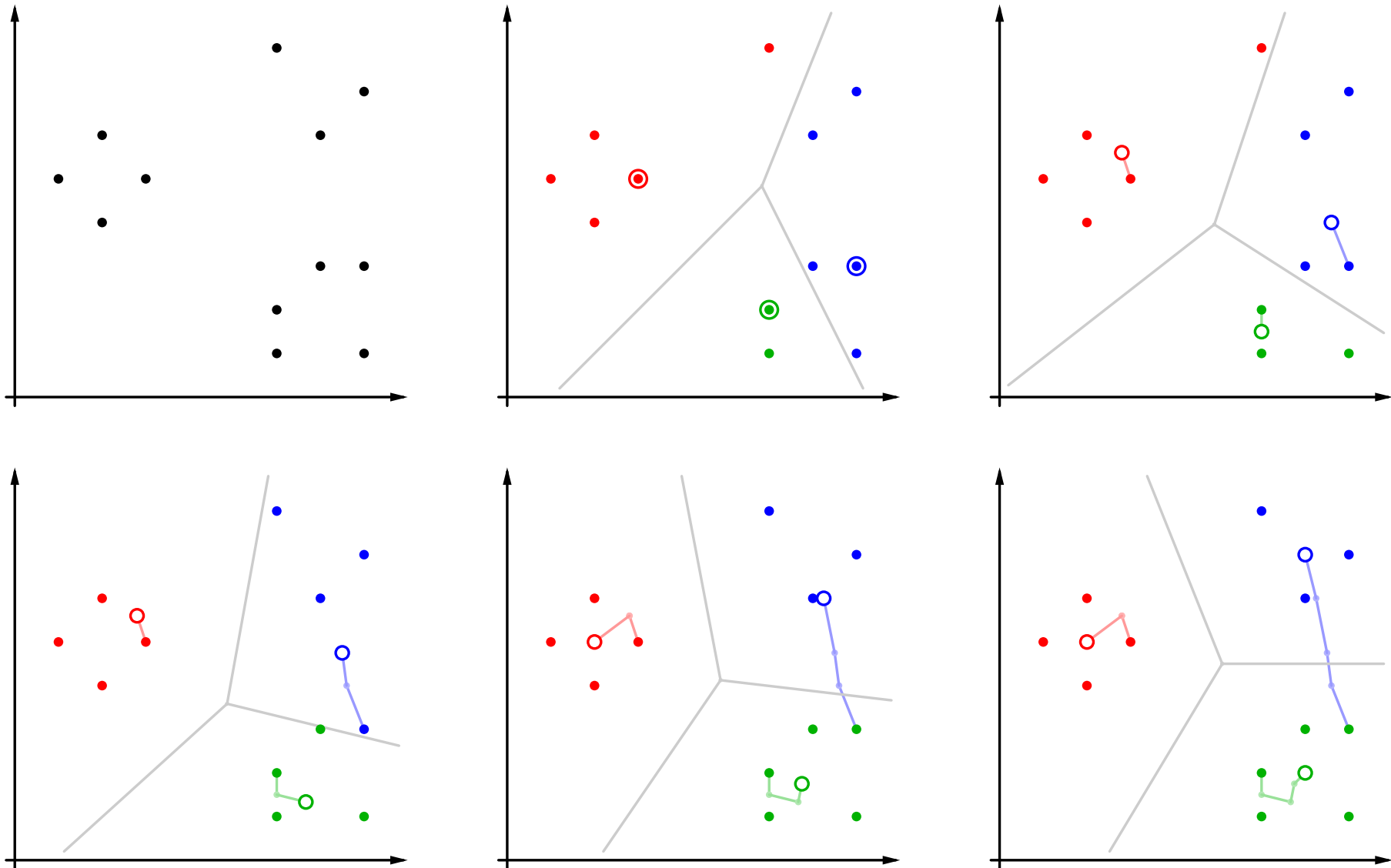
- Dots represent cluster centers (quantization vectors).
- Left: **Delaunay Triangulation**
(The circle through the corners of a triangle does not contain another point.)
- Right: **Voronoi Diagram**
(Midperpendiculars of the Delaunay triangulation: boundaries of the regions of points that are closest to the enclosed cluster center (Voronoi cells)).

Delaunay Triangulations and Voronoi Diagrams

- **Delaunay Triangulation:** simple triangle (shown in grey on the left)
- **Voronoi Diagram:** midperpendiculars of the triangle's edges (shown in blue on the left, in grey on the right)



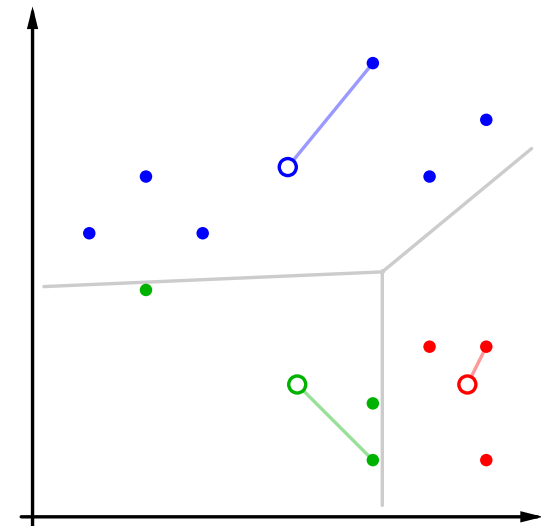
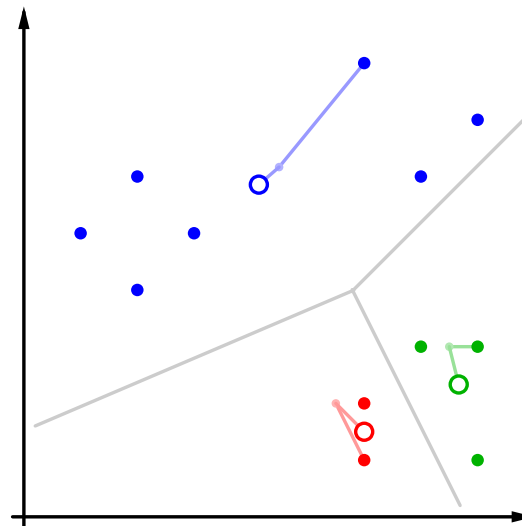
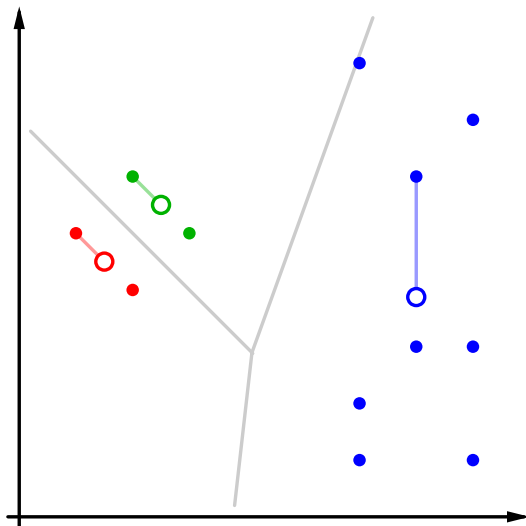
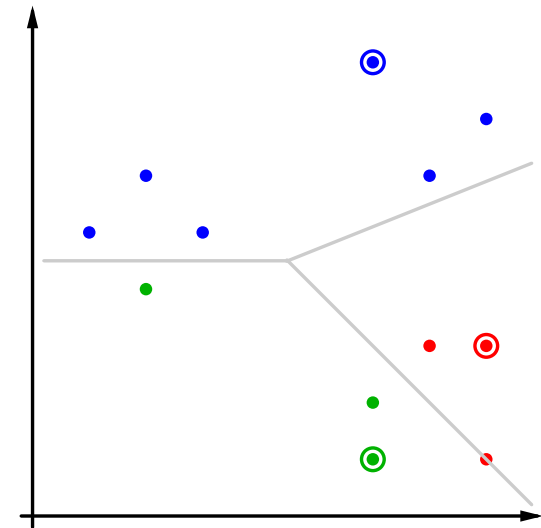
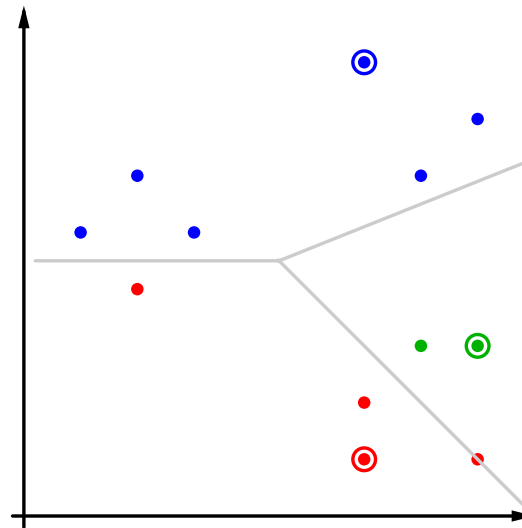
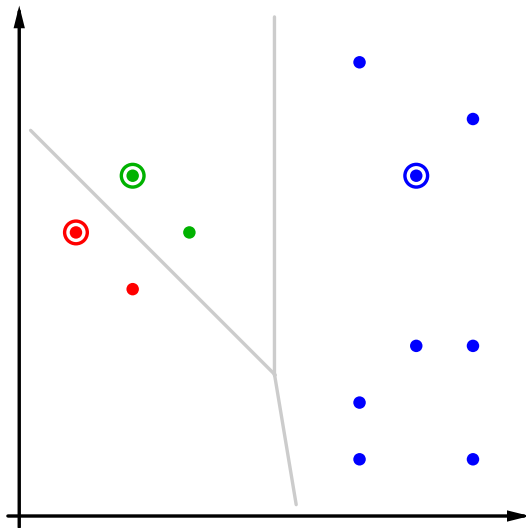
c -Means Clustering: Example



c -Means Clustering: Local Minima

- Clustering is successful in this example:
The clusters found are those that would have been formed intuitively.
- Convergence is achieved after only 5 steps.
(This is typical: convergence is usually very fast.)
- However: The clustering result is fairly **sensitive to the initial positions** of the cluster centers.
- With a bad initialization clustering may fail
(the alternating update process gets stuck in a local minimum).
- Fuzzy c -means clustering and the estimation of a mixture of Gaussians are much more robust (to be discussed later).
- Research issue: Can we determine the number of clusters automatically?
(Some approaches exist, but none of them is too successful.)

c -Means Clustering: Local Minima



Learning Vector Quantization

Adaptation of reference vectors / codebook vectors

- Like “online” *c*-means clustering (update after each data point).
- For each training pattern find the closest reference vector.
- Adapt only this reference vector (winner neuron).
- For classified data the class may be taken into account.
(reference vectors are assigned to classes)

Attraction rule (data point and reference vector have same class)

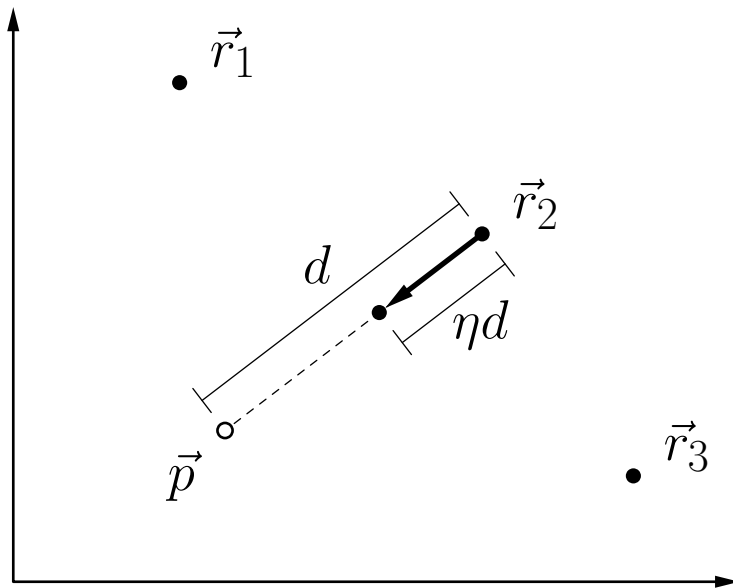
$$\vec{r}^{(\text{new})} = \vec{r}^{(\text{old})} + \eta(\vec{p} - \vec{r}^{(\text{old})}),$$

Repulsion rule (data point and reference vector have different class)

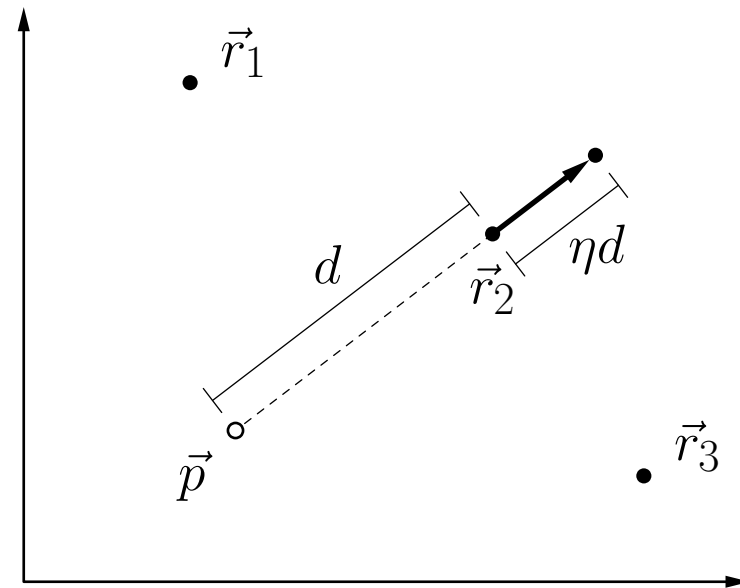
$$\vec{r}^{(\text{new})} = \vec{r}^{(\text{old})} - \eta(\vec{p} - \vec{r}^{(\text{old})}).$$

Learning Vector Quantization

Adaptation of reference vectors / codebook vectors



attraction rule

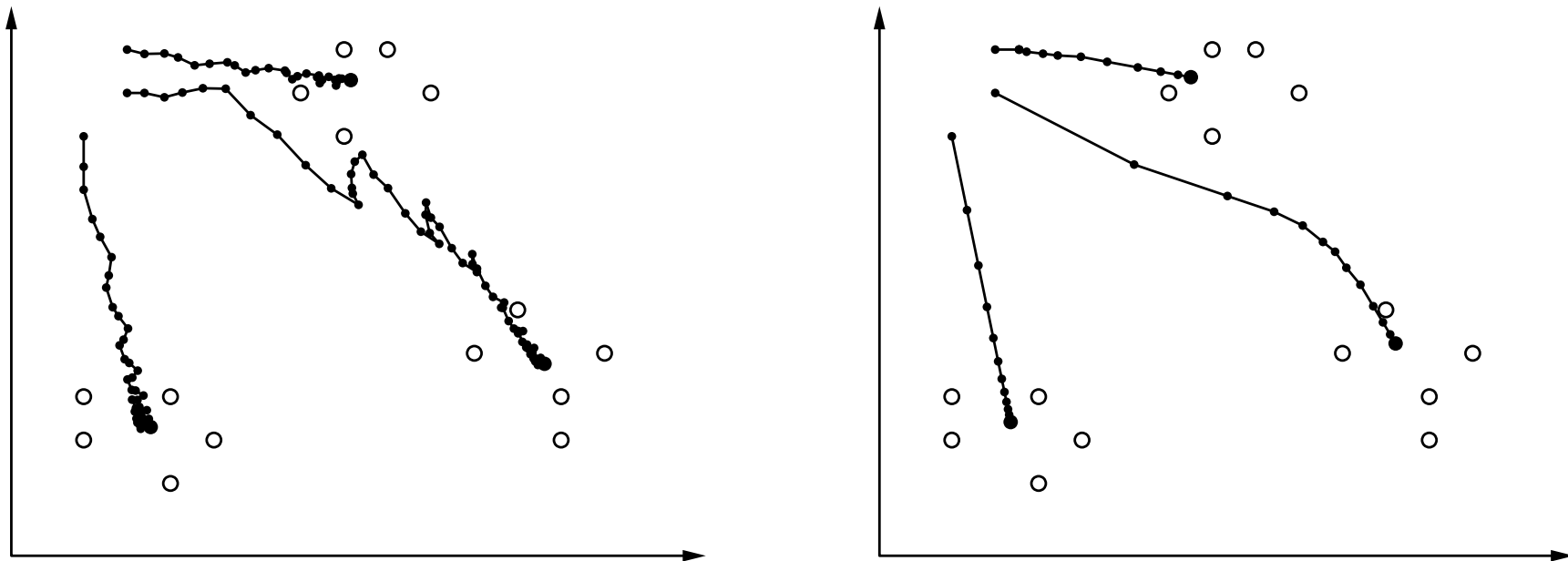


repulsion rule

- \vec{p} : data point, \vec{r}_i : reference vector
- $\eta = 0.4$ (learning rate)

Learning Vector Quantization: Example

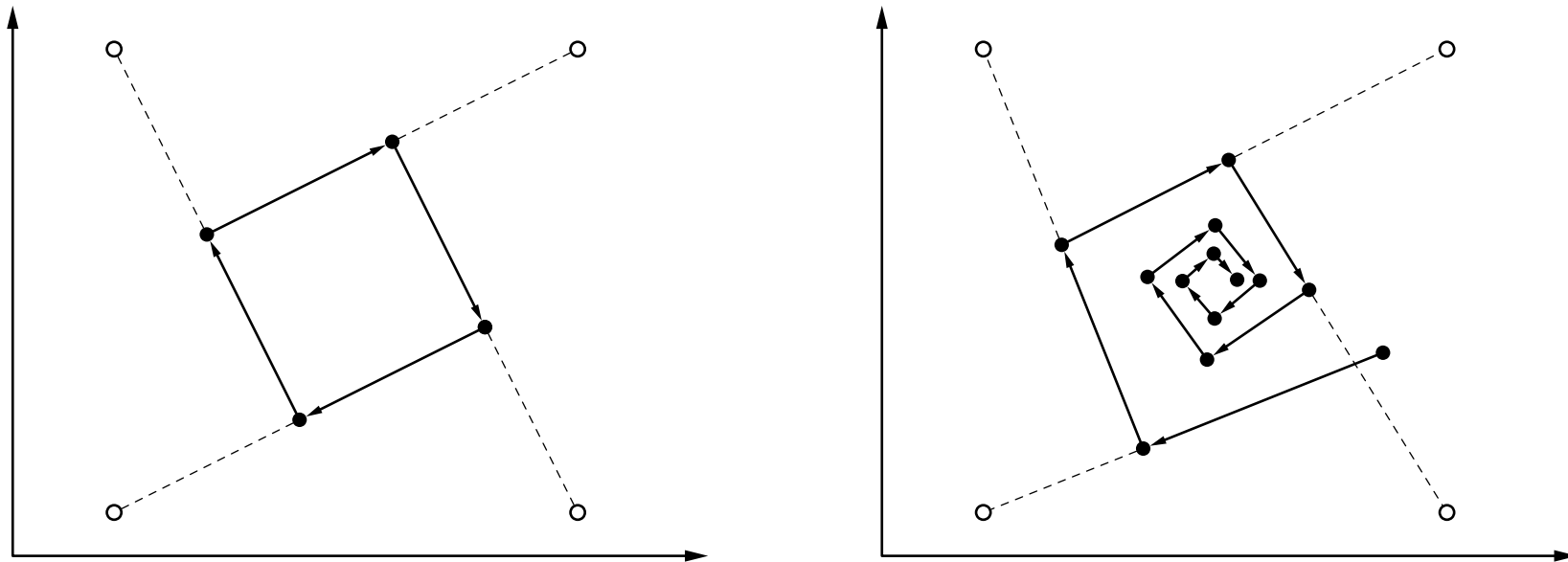
Adaptation of reference vectors / codebook vectors



- Left: Online training with learning rate $\eta = 0.1$,
- Right: Batch training with learning rate $\eta = 0.05$.

Learning Vector Quantization: Learning Rate Decay

Problem: fixed learning rate can lead to oscillations



Solution: **time dependent learning rate**

$$\eta(t) = \eta_0 \alpha^t, \quad 0 < \alpha < 1, \quad \text{or} \quad \eta(t) = \eta_0 t^\kappa, \quad \kappa < 0.$$

Fuzzy Clustering

- Allow degrees of membership of a datum to different clusters.
(Classical c -means clustering assigns data crisply.)

Objective Function: (to be minimized)

$$J(\mathbf{X}, \mathbf{B}, \mathbf{U}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^w d^2(\beta_i, \vec{x}_j)$$

- $\mathbf{U} = [u_{ij}]$ is the $c \times n$ fuzzy partition matrix,
 $u_{ij} \in [0, 1]$ is the membership degree of the data point \vec{x}_j to the i -th cluster.
- $\mathbf{B} = \{\beta_1, \dots, \beta_c\}$ is the set of cluster prototypes.
- w is the so-called “fuzzifier” (the higher w , the softer the cluster boundaries).

- Constraints:

$$\forall i \in \{1, \dots, c\} : \sum_{j=1}^n u_{ij} > 0 \quad \text{and} \quad \forall j \in \{1, \dots, n\} : \sum_{i=1}^c u_{ij} = 1.$$

Fuzzy and Hard Clustering

Relation to Classical c -Means Clustering:

- c -means clustering can be seen as optimizing the objective function

$$J(\mathbf{X}, \mathbf{B}, \mathbf{U}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij} d^2(\beta_i, \vec{x}_j),$$

where $\forall i, j : u_{ij} \in \{0, 1\}$ (i.e. hard assignment of the data points) and the cluster prototypes β_i consist only of cluster centers.

- To obtain a fuzzy assignment of the data points, it is not enough to extend the range of values for the u_{ij} to the unit interval $[0, 1]$: The objective function J is optimized for a hard assignment (each data point is assigned to the closest cluster center).
- **Necessary for degrees of membership:**
Apply a convex function $h : [0, 1] \rightarrow [0, 1]$ to the membership degrees u_{ij} .
Most common choice: $h(u) = u^w$, usually with $w = 2$.

Reminder: Function Optimization

Task: Find values $\vec{x} = (x_1, \dots, x_m)$ such that $f(\vec{x}) = f(x_1, \dots, x_m)$ is optimal.

Often feasible approach:

- A necessary condition for a (local) optimum (maximum or minimum) is that the partial derivatives w.r.t. the parameters vanish (Pierre Fermat).
- Therefore: (Try to) solve the equation system that results from setting all partial derivatives w.r.t. the parameters equal to zero.

Example task: Minimize $f(x, y) = x^2 + y^2 + xy - 4x - 5y$.

Solution procedure:

1. Take the partial derivatives of the objective function and set them to zero:

$$\frac{\partial f}{\partial x} = 2x + y - 4 = 0, \quad \frac{\partial f}{\partial y} = 2y + x - 5 = 0.$$

2. Solve the resulting (here: linear) equation system: $x = 1, \quad y = 2$.

Function Optimization with Constraints

Often a function has to be optimized subject to certain **constraints**.

Here: restriction to k **equality constraints** $C_i(\vec{x}) = 0, i = 1, \dots, k$.

Note: the equality constraints describe a subspace of the domain of the function.

Problem of optimization with constraints:

- The gradient of the objective function f may vanish outside the constrained subspace, leading to an unacceptable solution (violating the constraints).
- At an optimum *in the constrained subspace* the derivatives need not vanish.

One way to handle this problem are **generalized coordinates**:

- Exploit the dependence between the parameters specified in the constraints to express some parameters in terms of the others and thus reduce the set \vec{x} to a set \vec{x}' of independent parameters (*generalized coordinates*).
- Problem: Can be clumsy and cumbersome, if possible at all, because the form of the constraints may not allow for expressing some parameters as proper functions of the others.

Function Optimization with Constraints

A much more elegant approach is based on the following nice insights:

Let \vec{x}^* be a (local) optimum of $f(\vec{x})$ *in the constrained subspace*. Then:

- The gradient $\nabla_{\vec{x}} f(\vec{x}^*)$, if it does not vanish, must be **perpendicular** to the constrained subspace. (If $\nabla_{\vec{x}} f(\vec{x}^*)$ had a component in the constrained subspace, \vec{x}^* would not be a (local) optimum in this subspace.)
- The gradients $\nabla_{\vec{x}} C_j(\vec{x}^*)$, $1 \leq j \leq k$, must all be **perpendicular** to the constrained subspace, because they are constant, namely 0, in this subspace. Together they span the subspace perpendicular to the constrained subspace.
- Therefore it must be possible to find values λ_j , $1 \leq j \leq k$, such that

$$\nabla_{\vec{x}} f(\vec{x}^*) + \sum_{j=1}^s \lambda_j \nabla_{\vec{x}} C_j(\vec{x}^*) = 0.$$

If the constraints (and thus their gradients) are linearly independent, the values λ_j are uniquely determined. This equation can be used to **compensate the gradient** of $f(\vec{x}^*)$ so that it vanishes at \vec{x}^* .

Function Optimization: Lagrange Theory

As a consequence of these insights we obtain the

Method of Lagrange Multipliers:

Given:

- a function $f(\vec{x})$, which is to be optimized,
- k equality constraints $C_j(\vec{x}) = 0$, $1 \leq j \leq k$.

Procedure:

1. Construct the so-called **Lagrange function** by incorporating the equality constraints C_i , $i = 1, \dots, k$, with (unknown) **Lagrange multipliers** λ_i :

$$L(\vec{x}, \lambda_1, \dots, \lambda_k) = f(\vec{x}) + \sum_{i=1}^k \lambda_i C_i(\vec{x}).$$

2. Set the partial derivatives of the Lagrange function equal to zero:

$$\frac{\partial L}{\partial x_1} = 0, \quad \dots, \quad \frac{\partial L}{\partial x_m} = 0, \quad \frac{\partial L}{\partial \lambda_1} = 0, \quad \dots, \quad \frac{\partial L}{\partial \lambda_k} = 0.$$

3. (Try to) solve the resulting equation system.

Function Optimization: Lagrange Theory

Observations:

- Due to the representation of the gradient of $f(\vec{x})$ at a local optimum \vec{x}^* in the constrained subspace (see above) the gradient of L w.r.t. \vec{x} vanishes at \vec{x}^* .
→ The standard approach works again!
- If the constraints are satisfied, the additional terms have no influence.
→ The original task is not modified (same objective function).
- Taking the partial derivative w.r.t. a Lagrange multiplier reproduces the corresponding equality constraint:

$$\forall j; 1 \leq j \leq k : \quad \frac{\partial}{\partial \lambda_j} L(\vec{x}, \lambda_1, \dots, \lambda_k) = C_j(\vec{x}),$$

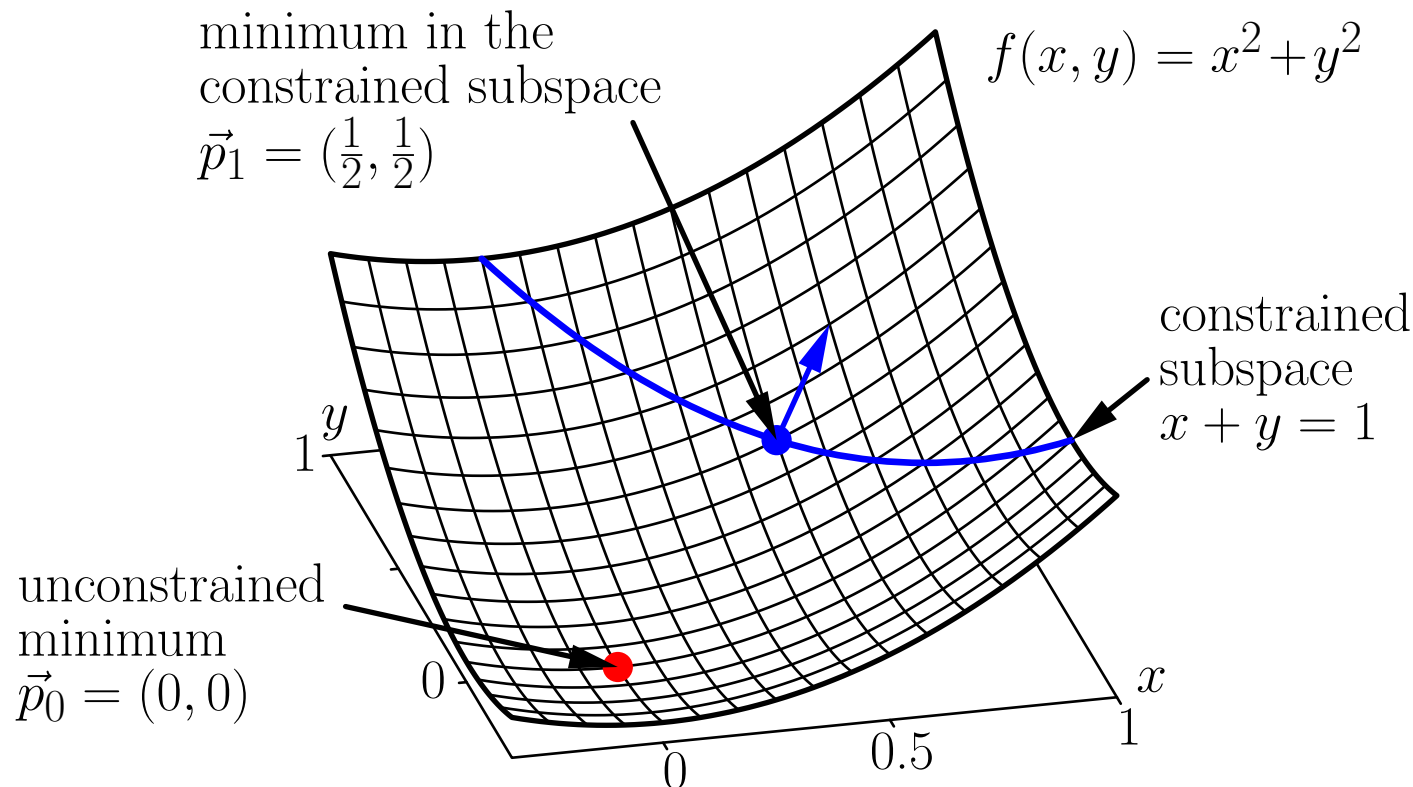
→ Constraints enter the equation system to solve in a natural way.

Remark:

- **Inequality** constraints can be handled with the **Kuhn–Tucker theory**.

Lagrange Theory: Example 1

Example task: Minimize $f(x, y) = x^2 + y^2$ subject to $x + y = 1$.



The unconstrained minimum is not in the constrained subspace, and at the minimum in the constrained subspace the gradient does not vanish.

Lagrange Theory: Example 1

Example task: Minimize $f(x, y) = x^2 + y^2$ subject to $x + y = 1$.

Solution procedure:

1. Rewrite the constraint, so that one side gets zero: $x + y - 1 = 0$.
2. Construct the Lagrange function by incorporating the constraint into the objective function with a Lagrange multiplier λ :

$$L(x, y, \lambda) = x^2 + y^2 + \lambda(x + y - 1).$$

3. Take the partial derivatives of the Lagrange function and set them to zero (necessary conditions for a minimum):

$$\frac{\partial L}{\partial x} = 2x + \lambda = 0, \quad \frac{\partial L}{\partial y} = 2y + \lambda = 0, \quad \frac{\partial L}{\partial \lambda} = x + y - 1 = 0.$$

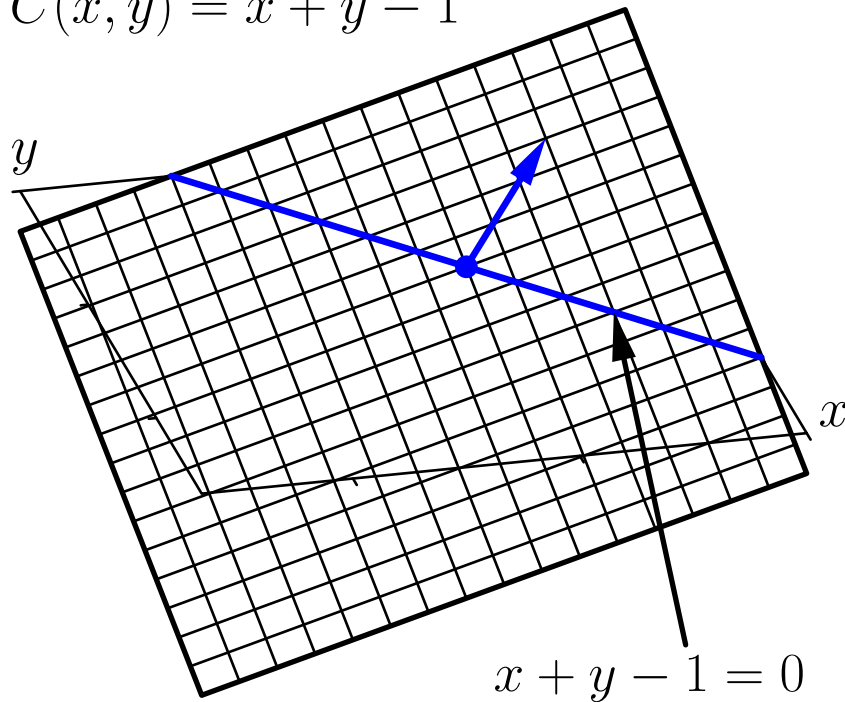
4. Solve the resulting (here: linear) equation system:

$$\lambda = -1, \quad x = y = \frac{1}{2}.$$

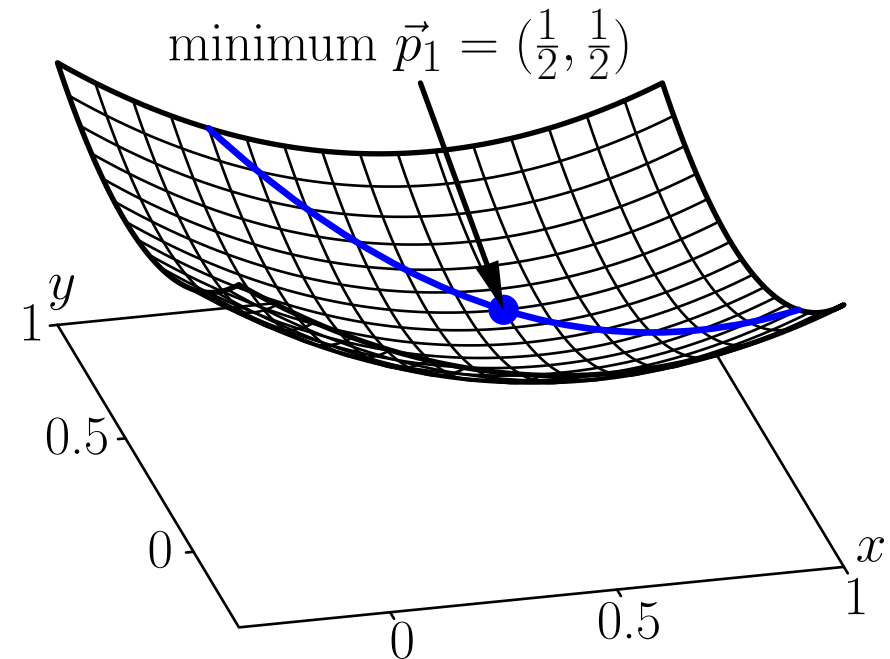
Lagrange Theory: Example 1

Example task: Minimize $f(x, y) = x^2 + y^2$ subject to $x + y = 1$.

$$C(x, y) = x + y - 1$$



$$L(x, y, -1) = x^2 + y^2 - (x + y - 1)$$



The gradient of the constraint is perpendicular to the constrained subspace.
The (unconstrained) minimum of the Lagrange function $L(x, y, \lambda)$
is the minimum of the objective function $f(x, y)$ in the constrained subspace.

Lagrange Theory: Example 2

Example task: Find the side lengths x, y, z of a box with maximum volume for a given area S of the surface.

Formally: Maximize $f(x, y, z) = xyz$
subject to $2xy + 2xz + 2yz = S$.

Solution procedure:

1. The constraint is $C(x, y, z) = 2xy + 2xz + 2yz - S = 0$.
2. The Lagrange function is

$$L(x, y, z, \lambda) = xyz + \lambda(2xy + 2xz + 2yz - S).$$

3. Taking the partial derivatives yields (in addition to the constraint):

$$\frac{\partial L}{\partial x} = yz + 2\lambda(y + z) = 0, \quad \frac{\partial L}{\partial y} = xz + 2\lambda(x + z) = 0, \quad \frac{\partial L}{\partial z} = xy + 2\lambda(x + y) = 0.$$

4. The solution is: $\lambda = -\frac{1}{4}\sqrt{\frac{S}{6}}, \quad x = y = z = \sqrt{\frac{S}{6}}$ (i.e., the box is a cube).

Fuzzy Clustering: Alternating Optimization

Objective function: (to be minimized)

$$J(\mathbf{X}, \mathbf{B}, \mathbf{U}) = \sum_{i=1}^c \sum_{j=1}^n u_{ij}^w d^2(\vec{x}_j, \beta_i)$$

Constraints:

$$\forall i \in \{1, \dots, c\} : \sum_{j=1}^n u_{ij} > 0 \quad \text{and} \quad \forall j \in \{1, \dots, n\} : \sum_{i=1}^c u_{ij} = 1.$$

- **Problem:** The objective function J cannot be minimized directly.
- Therefore: **Alternating Optimization**
 - Optimize membership degrees for fixed cluster parameters.
 - Optimize cluster parameters for fixed membership degrees.
(Update formulae are derived by differentiating the objective function J)
 - Iterate until convergence (checked, e.g., by change of cluster center).

Fuzzy Clustering: Alternating Optimization

First Step: Fix the cluster parameters.

Introduce Lagrange multipliers λ_j , $0 \leq j \leq n$, to incorporate the constraints $\forall j; 1 \leq j \leq n : \sum_{i=1}^c u_{ij} = 1$. This yields the Lagrange function (to be minimized)

$$L(\mathbf{X}, \mathbf{B}, \mathbf{U}, \Lambda) = \underbrace{\sum_{i=1}^c \sum_{j=1}^n u_{ij}^w d_{ij}^2}_{=J(\mathbf{X}, \mathbf{B}, \mathbf{U})} + \sum_{j=1}^n \lambda_j \left(1 - \sum_{i=1}^c u_{ij} \right),$$

A necessary condition for the minimum is that the partial derivatives of the Lagrange function w.r.t. the membership degrees vanish, i.e.,

$$\frac{\partial}{\partial u_{kl}} L(\mathbf{X}, \mathbf{B}, \mathbf{U}, \Lambda) = w u_{kl}^{w-1} d_{kl}^2 - \lambda_l \stackrel{!}{=} 0,$$

which leads to

$$\forall i; 1 \leq i \leq c : \forall j; 1 \leq j \leq n : \quad u_{ij} = \left(\frac{\lambda_j}{w d_{ij}^2} \right)^{\frac{1}{w-1}}.$$

Fuzzy Clustering: Alternating Optimization

Summing these equations over the clusters (in order to be able to exploit the corresponding constraints on the membership degrees), we get

$$1 = \sum_{i=1}^c u_{ij} = \sum_{i=1}^c \left(\frac{\lambda_j}{w d_{ij}^2} \right)^{\frac{1}{w-1}}.$$

Consequently the λ_j , $1 \leq j \leq n$, are

$$\lambda_j = \left(\sum_{i=1}^c \left(w d_{ij}^2 \right)^{\frac{1}{1-w}} \right)^{1-w}.$$

Inserting this into the equation for the membership degrees yields

$$\forall i; 1 \leq i \leq c : \forall j; 1 \leq j \leq n : \quad u_{ij} = \frac{d_{ij}^{\frac{2}{1-w}}}{\sum_{k=1}^c d_{kj}^{\frac{2}{1-w}}}.$$

This update formula results regardless of the distance measure.

Standard Fuzzy Clustering Algorithms

Fuzzy C-Means Algorithm: Euclidean distance

$$d_{\text{fcm}}^2(\vec{x}_j, \beta_i) = (\vec{x}_j - \vec{\mu}_i)^\top (\vec{x}_j - \vec{\mu}_i)$$

Necessary condition for a minimum: gradients w.r.t. cluster centers vanish.

$$\begin{aligned}\nabla_{\vec{\mu}_k} J_{\text{fcm}}(\mathbf{X}, \mathbf{B}, \mathbf{U}) &= \nabla_{\vec{\mu}_k} \sum_{i=1}^c \sum_{j=1}^n u_{ij}^w (\vec{x}_j - \vec{\mu}_i)^\top (\vec{x}_j - \vec{\mu}_i) \\ &= \sum_{j=1}^n u_{kj}^w \nabla_{\vec{\mu}_k} (\vec{x}_j - \vec{\mu}_k)^\top (\vec{x}_j - \vec{\mu}_k) \\ &= -2 \sum_{j=1}^n u_{kj}^w (\vec{x}_j - \vec{\mu}_k) \stackrel{!}{=} \vec{0}\end{aligned}$$

Resulting update rule for the cluster centers (**second step** of alt. optimization):

$$\forall i; 1 \leq i \leq c : \quad \vec{\mu}_i = \frac{\sum_{j=1}^n u_{ij}^w \vec{x}_j}{\sum_{j=1}^n u_{ij}^w}$$

Standard Fuzzy Clustering Algorithms

Gustafson–Kessel Algorithm: Mahalanobis distance

$$d_{\text{gk}}^2(\vec{x}_j, \beta_i) = (\vec{x}_j - \vec{\mu}_i)^\top \mathbf{C}_i^{-1} (\vec{x}_j - \vec{\mu}_i)$$

Additional constraints: $|\mathbf{C}_i| = 1$ (all cluster have unit size).

These constraints are incorporated again by Lagrange multipliers.

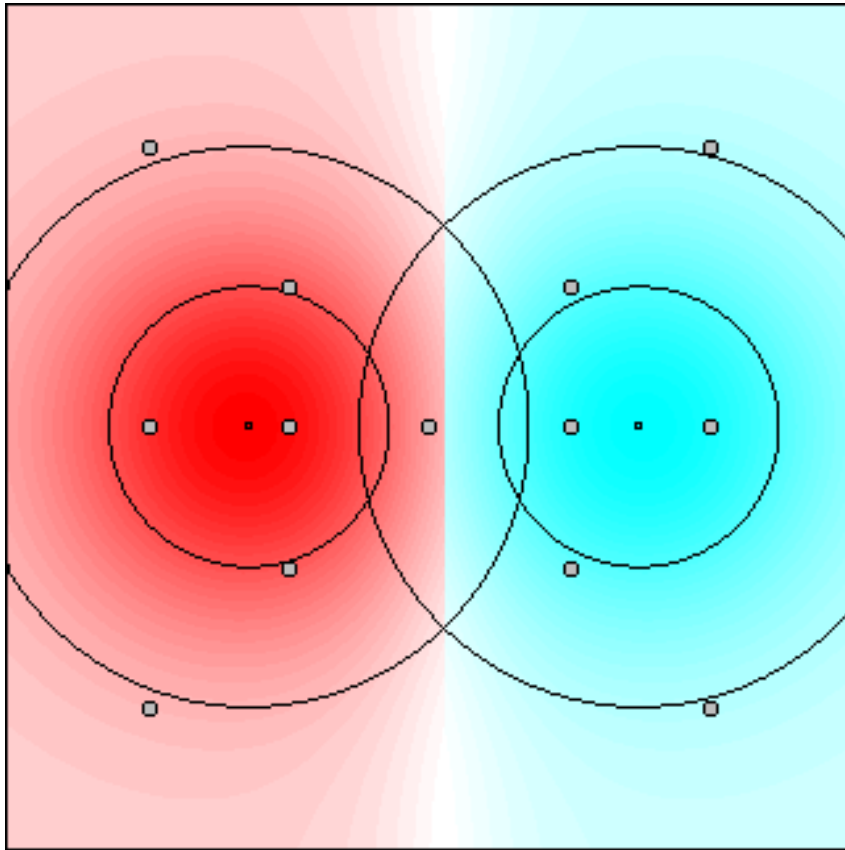
A similar derivation as for the fuzzy c -means algorithm yields the same update rule for the cluster centers:

$$\forall i; 1 \leq i \leq c : \quad \vec{\mu}_i = \frac{\sum_{j=1}^n u_{ij}^w \vec{x}_j}{\sum_{j=1}^n u_{ij}^w}$$

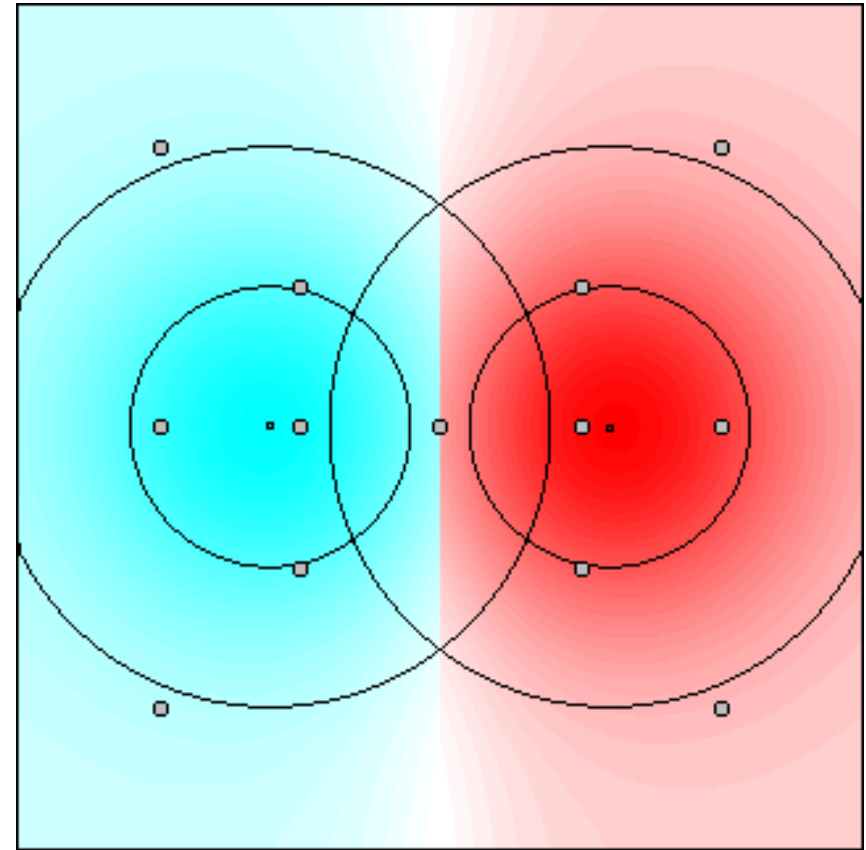
Update rule for the covariance matrices (m is the number of dimensions):

$$\mathbf{C}_i = \frac{1}{m \sqrt{|\boldsymbol{\Sigma}_i|}} \boldsymbol{\Sigma}_i \quad \text{where} \quad \boldsymbol{\Sigma}_i = \sum_{j=1}^n u_{ij}^w (\vec{x}_j - \vec{\mu}_i)(\vec{x}_j - \vec{\mu}_i)^\top.$$

Fuzzy Clustering: Overlapping Clusters

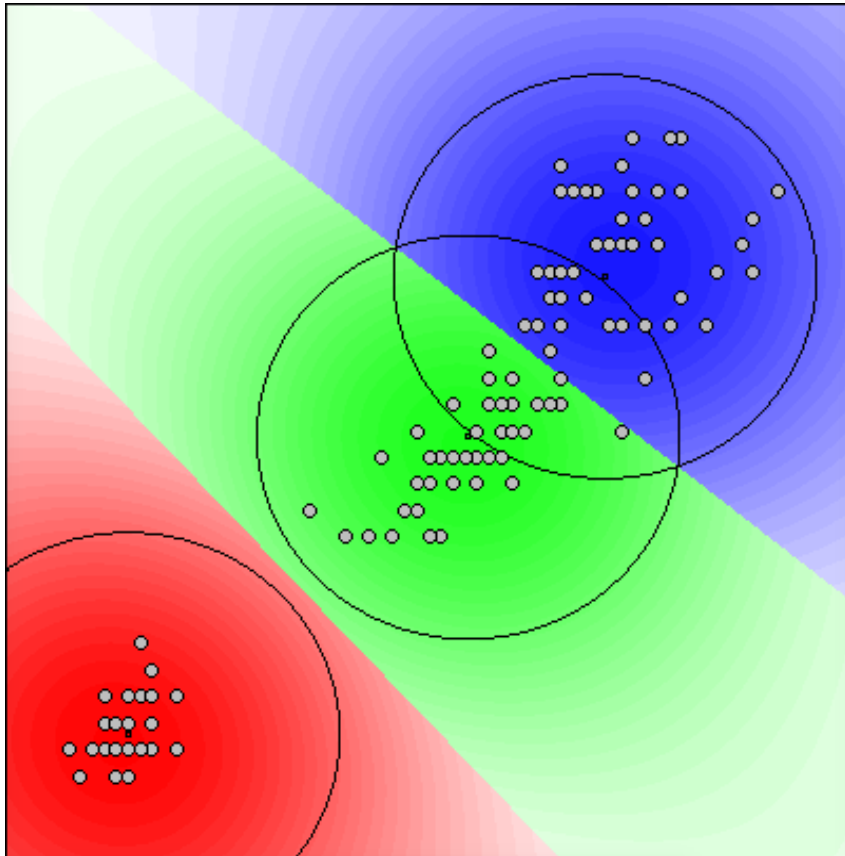


Classical c -Means

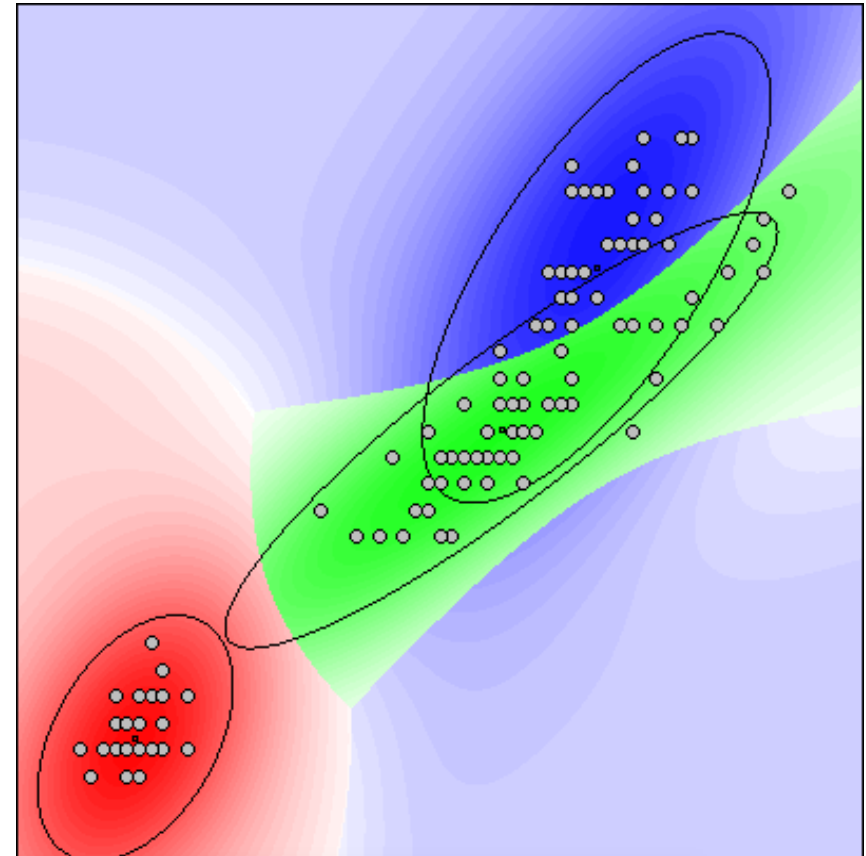


Fuzzy c -Means

Fuzzy Clustering of the Iris Data



Fuzzy c -Means



Gustafson-Kessel

Expectation Maximization: Mixture of Gaussians

- **Assumption:** Data was generated by sampling a set of normal distributions. (The probability density is a mixture of Gaussian distributions.)
- **Formally:** We assume that the probability density can be described as

$$f_{\vec{X}}(\vec{x}; \mathbf{C}) = \sum_{y=1}^c f_{\vec{X}, Y}(\vec{x}, y; \mathbf{C}) = \sum_{y=1}^c p_Y(y; \mathbf{C}) \cdot f_{\vec{X}|Y}(\vec{x}|y; \mathbf{C}).$$

\mathbf{C}	is the set of cluster parameters
\vec{X}	is a random vector that has the data space as its domain
Y	is a random variable that has the cluster indices as possible values (i.e., $\text{dom}(\vec{X}) = \mathbb{R}^m$ and $\text{dom}(Y) = \{1, \dots, c\}$)
$p_Y(y; \mathbf{C})$	is the probability that a data point belongs to (is generated by) the y -th component of the mixture
$f_{\vec{X} Y}(\vec{x} y; \mathbf{C})$	is the conditional probability density function of a data point given the cluster (specified by the cluster index y)

Expectation Maximization

- **Basic idea:** Do a maximum likelihood estimation of the cluster parameters.
- **Problem:** The likelihood function,

$$L(\mathbf{X}; \mathbf{C}) = \prod_{j=1}^n f_{\vec{X}_j}(\vec{x}_j; \mathbf{C}) = \prod_{j=1}^n \sum_{y=1}^c p_Y(y; \mathbf{C}) \cdot f_{\vec{X}|Y}(\vec{x}_j|y; \mathbf{C}),$$

is difficult to optimize, even if one takes the natural logarithm (cf. the maximum likelihood estimation of the parameters of a normal distribution), because

$$\ln L(\mathbf{X}; \mathbf{C}) = \sum_{j=1}^n \ln \sum_{y=1}^c p_Y(y; \mathbf{C}) \cdot f_{\vec{X}|Y}(\vec{x}_j|y; \mathbf{C})$$

contains the natural logarithms of complex sums.

- **Approach:** Assume that there are “hidden” variables Y_j stating the clusters that generated the data points \vec{x}_j , so that the sums reduce to one term.
- **Problem:** Since the Y_j are hidden, we do not know their values.

Expectation Maximization

- **Formally:** Maximize the likelihood of the “completed” data set (\mathbf{X}, \vec{y}) , where $\vec{y} = (y_1, \dots, y_n)$ combines the values of the variables Y_j . That is,

$$L(\mathbf{X}, \vec{y}; \mathbf{C}) = \prod_{j=1}^n f_{\vec{X}_j, Y_j}(\vec{x}_j, y_j; \mathbf{C}) = \prod_{j=1}^n p_{Y_j}(y_j; \mathbf{C}) \cdot f_{\vec{X}_j|Y_j}(\vec{x}_j|y_j; \mathbf{C}).$$

- **Problem:** Since the Y_j are hidden, the values y_j are unknown (and thus the factors $p_{Y_j}(y_j; \mathbf{C})$ cannot be computed).
- **Approach to find a solution nevertheless:**
 - See the Y_j as random variables (the values y_j are not fixed) and consider a probability distribution over the possible values.
 - As a consequence $L(\mathbf{X}, \vec{y}; \mathbf{C})$ becomes a random variable, even for a fixed data set \mathbf{X} and fixed cluster parameters \mathbf{C} .
 - Try to **maximize the expected value** of $L(\mathbf{X}, \vec{y}; \mathbf{C})$ or $\ln L(\mathbf{X}, \vec{y}; \mathbf{C})$ (hence the name **expectation maximization**).

Expectation Maximization

- **Formally:** Find the cluster parameters as

$$\hat{\mathbf{C}} = \underset{\mathbf{C}}{\operatorname{argmax}} E([\ln] L(\mathbf{X}, \vec{y}; \mathbf{C}) \mid \mathbf{X}; \mathbf{C}),$$

that is, maximize the expected likelihood

$$E(L(\mathbf{X}, \vec{y}; \mathbf{C}) \mid \mathbf{X}; \mathbf{C}) = \sum_{\vec{y} \in \{1, \dots, c\}^n} p_{\vec{Y} \mid \mathcal{X}}(\vec{y} \mid \mathbf{X}; \mathbf{C}) \cdot \prod_{j=1}^n f_{\vec{X}_j, Y_j}(\vec{x}_j, y_j; \mathbf{C})$$

or, alternatively, maximize the expected log-likelihood

$$E(\ln L(\mathbf{X}, \vec{y}; \mathbf{C}) \mid \mathbf{X}; \mathbf{C}) = \sum_{\vec{y} \in \{1, \dots, c\}^n} p_{\vec{Y} \mid \mathcal{X}}(\vec{y} \mid \mathbf{X}; \mathbf{C}) \cdot \sum_{j=1}^n \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, y_j; \mathbf{C}).$$

- Unfortunately, these functionals are still difficult to optimize directly.
- **Solution:** Use the equation as an iterative scheme, fixing \mathbf{C} in some terms (iteratively compute better approximations, similar to Heron's algorithm).

Excursion: Heron's Algorithm

- **Task:** Find the square root of a given number x , i.e., find $y = \sqrt{x}$.

- **Approach:** Rewrite the defining equation $y^2 = x$ as follows:

$$y^2 = x \quad \Leftrightarrow \quad 2y^2 = y^2 + x \quad \Leftrightarrow \quad y = \frac{1}{2y}(y^2 + x) \quad \Leftrightarrow \quad y = \frac{1}{2} \left(y + \frac{x}{y} \right).$$

- Use the resulting equation as an iteration formula, i.e., compute the sequence

$$y_{k+1} = \frac{1}{2} \left(y_k + \frac{x}{y_k} \right) \quad \text{with} \quad y_0 = 1.$$

- It can be shown that $0 \leq y_k - \sqrt{x} \leq y_{k-1} - y_n$ for $k \geq 2$.
Therefore this iteration formula provides increasingly better approximations of the square root of x and thus is a safe and simple way to compute it.
Ex.: $x = 2$: $y_0 = 1$, $y_1 = 1.5$, $y_2 \approx 1.41667$, $y_3 \approx 1.414216$, $y_4 \approx 1.414213$.
- Heron's algorithm converges very quickly and is often used in pocket calculators and microprocessors to implement the square root.

Expectation Maximization

- **Iterative scheme for expectation maximization:**

Choose some initial set \mathbf{C}_0 of cluster parameters and then compute

$$\begin{aligned}\mathbf{C}_{k+1} &= \operatorname{argmax}_{\mathbf{C}} E(\ln L(\mathbf{X}, \vec{y}; \mathbf{C}) \mid \mathbf{X}; \mathbf{C}_k) \\ &= \operatorname{argmax}_{\mathbf{C}} \sum_{\vec{y} \in \{1, \dots, c\}^n} p_{\vec{Y}|\mathcal{X}}(\vec{y}|\mathbf{X}; \mathbf{C}_k) \sum_{j=1}^n \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, y_j; \mathbf{C}) \\ &= \operatorname{argmax}_{\mathbf{C}} \sum_{\vec{y} \in \{1, \dots, c\}^n} \left(\prod_{l=1}^n p_{Y_l|\vec{X}_l}(y_l|\vec{x}_l; \mathbf{C}_k) \right) \sum_{j=1}^n \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, y_j; \mathbf{C}) \\ &= \operatorname{argmax}_{\mathbf{C}} \sum_{i=1}^c \sum_{j=1}^n p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k) \cdot \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, i; \mathbf{C}).\end{aligned}$$

- It can be shown that each EM iteration increases the likelihood of the data and that the algorithm converges to a local maximum of the likelihood function (i.e., EM is a safe way to maximize the likelihood function).

Expectation Maximization

Justification of the last step on the previous slide:

$$\begin{aligned}
 & \sum_{\vec{y} \in \{1, \dots, c\}^n} \left(\prod_{l=1}^n p_{Y_l | \vec{X}_l}(y_l | \vec{x}_l; \mathbf{C}_k) \right) \sum_{j=1}^n \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, y_j; \mathbf{C}) \\
 &= \sum_{y_1=1}^c \cdots \sum_{y_n=1}^c \left(\prod_{l=1}^n p_{Y_l | \vec{X}_l}(y_l | \vec{x}_l; \mathbf{C}_k) \right) \sum_{j=1}^n \sum_{i=1}^c \delta_{i, y_j} \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, i; \mathbf{C}) \\
 &= \sum_{i=1}^c \sum_{j=1}^n \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, i; \mathbf{C}) \sum_{y_1=1}^c \cdots \sum_{y_n=1}^c \delta_{i, y_j} \prod_{l=1}^n p_{Y_l | \vec{X}_l}(y_l | \vec{x}_l; \mathbf{C}_k) \\
 &= \sum_{i=1}^c \sum_{j=1}^n p_{Y_j | \vec{X}_j}(i | \vec{x}_j; \mathbf{C}_k) \cdot \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, i; \mathbf{C}) \\
 & \quad \underbrace{\sum_{y_1=1}^c \cdots \sum_{y_{j-1}=1}^c \sum_{y_{j+1}=1}^c \cdots \sum_{y_n=1}^c \prod_{l=1, l \neq j}^n p_{Y_l | \vec{X}_l}(y_l | \vec{x}_l; \mathbf{C}_k)}_{= \prod_{l=1, l \neq j}^n \sum_{y_l=1}^c p_{Y_l | \vec{X}_l}(y_l | \vec{x}_l; \mathbf{C}_k) = \prod_{l=1, l \neq j}^n 1 = 1} \\
 &= \prod_{l=1, l \neq j}^n \sum_{y_l=1}^c p_{Y_l | \vec{X}_l}(y_l | \vec{x}_l; \mathbf{C}_k) = \prod_{l=1, l \neq j}^n 1 = 1
 \end{aligned}$$

Expectation Maximization

- The probabilities $p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k)$ are computed as

$$p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k) = \frac{f_{\vec{X}_j, Y_j}(\vec{x}_j, i; \mathbf{C}_k)}{f_{\vec{X}_j}(\vec{x}_j; \mathbf{C}_k)} = \frac{f_{\vec{X}_j|Y_j}(\vec{x}_j|i; \mathbf{C}_k) \cdot p_{Y_j}(i; \mathbf{C}_k)}{\sum_{l=1}^c f_{\vec{X}_j|Y_j}(\vec{x}_j|l; \mathbf{C}_k) \cdot p_{Y_j}(l; \mathbf{C}_k)},$$

that is, as the relative probability densities of the different clusters (as specified by the cluster parameters) at the location of the data points \vec{x}_j .

- The $p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k)$ are the posterior probabilities of the clusters given the data point \vec{x}_j and a set of cluster parameters \mathbf{C}_k .
- They can be seen as **case weights** of a “completed” data set:
 - Split each data point \vec{x}_j into c data points (\vec{x}_j, i) , $i = 1, \dots, c$.
 - Distribute the unit weight of the data point \vec{x}_j according to the above probabilities, i.e., assign to (\vec{x}_j, i) the weight $p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k)$, $i = 1, \dots, c$.

Expectation Maximization: Cookbook Recipe

Core Iteration Formula

$$\mathbf{C}_{k+1} = \operatorname{argmax}_{\mathbf{C}} \sum_{i=1}^c \sum_{j=1}^n p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k) \cdot \ln f_{\vec{X}_j, Y_j}(\vec{x}_j, i; \mathbf{C})$$

Expectation Step

- For all data points \vec{x}_j :
Compute for each normal distribution the probability $p_{Y_j|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}_k)$ that the data point was generated from it
(ratio of probability densities at the location of the data point).
→ “weight” of the data point for the estimation.

Maximization Step

- For all normal distributions:
Estimate the parameters by standard maximum likelihood estimation using the probabilities (“weights”) assigned to the data points w.r.t. the distribution in the expectation step.

Expectation Maximization: Mixture of Gaussians

Expectation Step: Use Bayes' rule to compute

$$p_{C|\vec{X}}(i|\vec{x}; \mathbf{C}) = \frac{p_C(i; \mathbf{c}_i) \cdot f_{\vec{X}|C}(\vec{x}|i; \mathbf{c}_i)}{f_{\vec{X}}(\vec{x}; \mathbf{C})} = \frac{p_C(i; \mathbf{c}_i) \cdot f_{\vec{X}|C}(\vec{x}|i; \mathbf{c}_i)}{\sum_{k=1}^c p_C(k; \mathbf{c}_k) \cdot f_{\vec{X}|C}(\vec{x}|k; \mathbf{c}_k)}.$$

→ “weight” of the data point \vec{x} for the estimation.

Maximization Step: Use maximum likelihood estimation to compute

$$\varrho_i^{(t+1)} = \frac{1}{n} \sum_{j=1}^n p_{C|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}^{(t)}), \quad \vec{\mu}_i^{(t+1)} = \frac{\sum_{j=1}^n p_{C|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}^{(t)}) \cdot \vec{x}_j}{\sum_{j=1}^n p_{C|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}^{(t)})},$$

$$\text{and} \quad \Sigma_i^{(t+1)} = \frac{\sum_{j=1}^n p_{C|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}^{(t)}) \cdot (\vec{x}_j - \vec{\mu}_i^{(t+1)}) (\vec{x}_j - \vec{\mu}_i^{(t+1)})^\top}{\sum_{j=1}^n p_{C|\vec{X}_j}(i|\vec{x}_j; \mathbf{C}^{(t)})}$$

Iterate until convergence (checked, e.g., by change of mean vector).

Expectation Maximization: Technical Problems

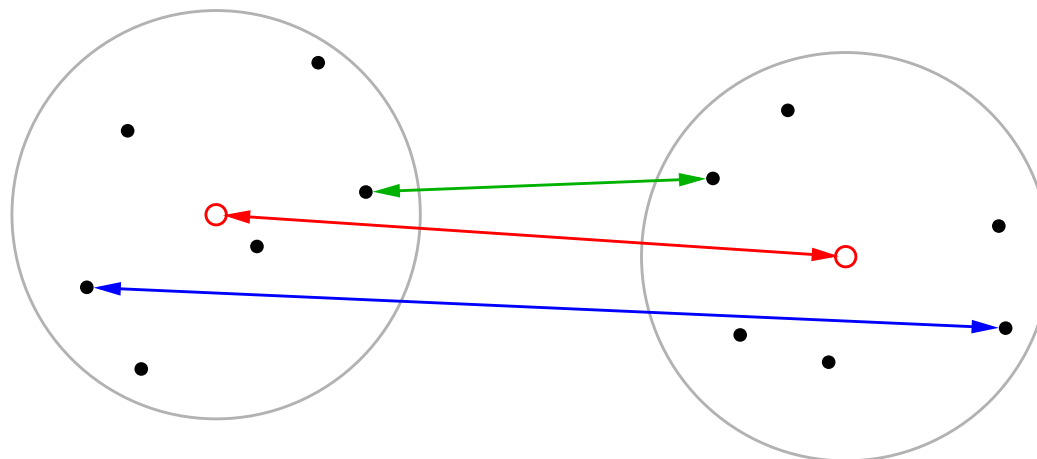
- If a fully general mixture of Gaussian distributions is used, the likelihood function is truly optimized if
 - all normal distributions except one are contracted to single data points and
 - the remaining normal distribution is the maximum likelihood estimate for the remaining data points.
- This undesired result is rare, because the algorithm gets stuck in a local optimum.
- Nevertheless it is recommended to take countermeasures, which consist mainly in reducing the degrees of freedom, like
 - Fix the determinants of the covariance matrices to equal values.
 - Use a diagonal instead of a general covariance matrix.
 - Use an isotropic variance instead of a covariance matrix.
 - Fix the prior probabilities of the clusters to equal values.

Hierarchical Agglomerative Clustering

- Start with every data point in its own cluster.
(i.e., start with so-called **singletons**: single element clusters)
- In each step merge those two clusters that are closest to each other.
- Keep on merging clusters until all data points are contained in one cluster.
- The result is a hierarchy of clusters that can be visualized in a tree structure
(a so-called **dendrogram** — from the Greek *δέντρον* (dendron): tree)
- **Measuring the Distances**
 - The distance between singletons is simply the distance between the (single) data points contained in them.
 - However: How do we compute the distance between clusters that contain more than one data point?

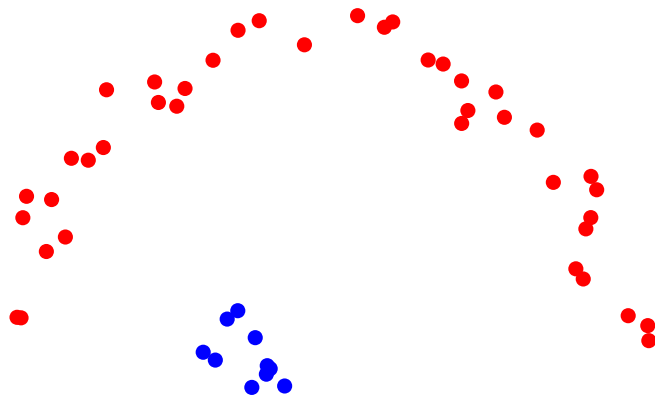
Measuring the Distance between Clusters

- **Centroid** (red)
Distance between the centroids (mean value vectors) of the two clusters.
- **Average Linkage**
Average distance between two points of the two clusters.
- **Single Linkage** (green)
Distance between the two closest points of the two clusters.
- **Complete Linkage** (blue)
Distance between the two farthest points of the two clusters.

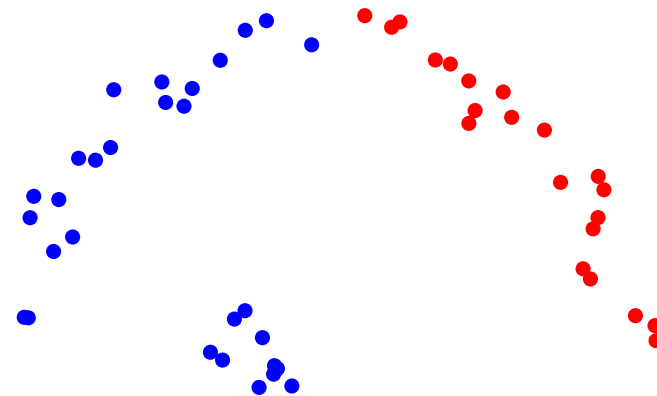


Measuring the Distance between Clusters

- Single linkage can “follow chains” in the data (may be desirable in certain applications).
- Complete linkage leads to very compact clusters.
- Average linkage also tends clearly towards compact clusters.



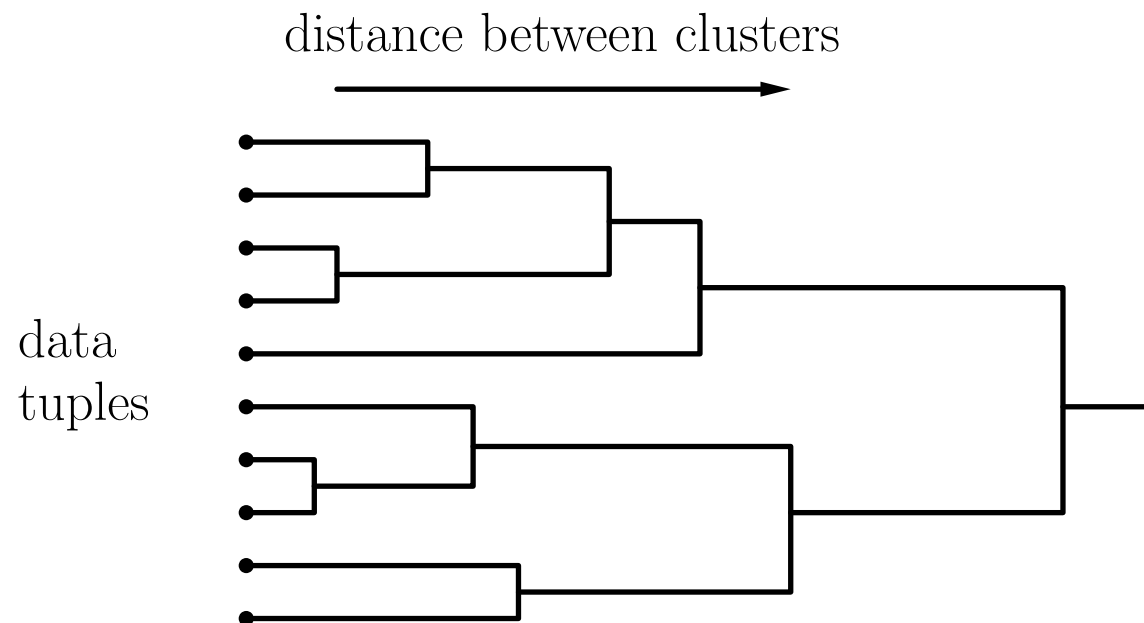
Single Linkage



Complete Linkage

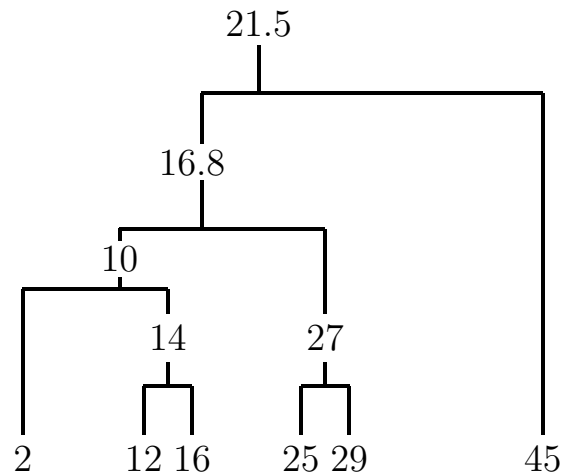
Dendrograms

- The cluster merging process arranges the data points in a binary tree.
- Draw the data tuples at the bottom or on the left (equally spaced if they are multi-dimensional).
- Draw a connection between clusters that are merged, with the distance to the data points representing the distance between the clusters.

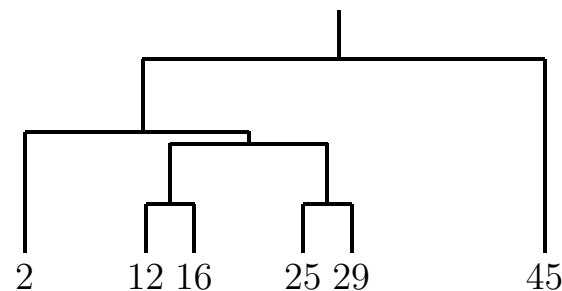


Hierarchical Agglomerative Clustering

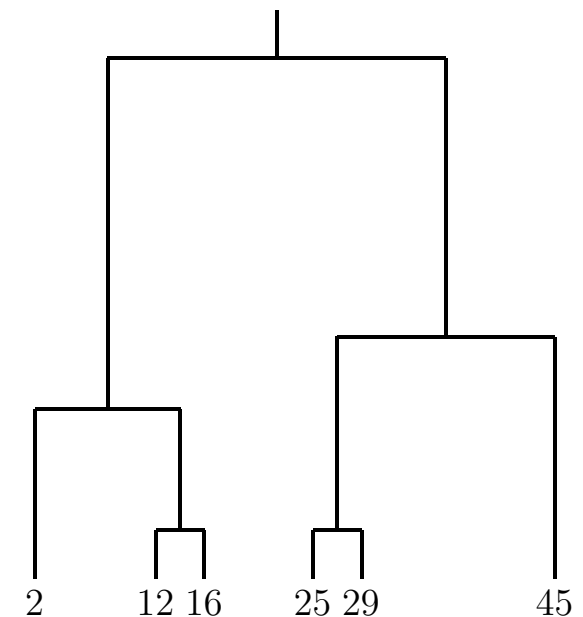
- Example: Clustering of the 1-dimensional data set $\{2, 12, 16, 25, 29, 45\}$.
- All three approaches to measure the distance between clusters lead to different dendrograms.



Centroid



Single Linkage



Complete Linkage

Implementation Aspects

- Hierarchical agglomerative clustering can be implemented by processing the matrix $\mathbf{D} = (d_{ij})_{1 \leq i, j \leq n}$ containing the pairwise distances of the data points. (The data points themselves are actually not needed.)
- In each step the rows and columns corresponding to the two clusters that are closest to each other are deleted.
- A new row and column corresponding to the cluster formed by merging these clusters is added to the matrix.
- The elements of this new row/column are computed according to

$$\forall k : \quad d_{k*} = d_{*k} = \alpha_i d_{ik} + \alpha_j d_{jk} + \beta d_{ij} + \gamma |d_{ik} - d_{jk}|$$

i, j	indices of the two clusters that are merged
k	indices of the old clusters that are <i>not</i> merged
$*$	index of the new cluster (result of merger)
$\alpha_i, \alpha_j, \beta, \gamma$	parameters specifying the method (single linkage etc.)

Implementation Aspects

- The parameters defining the different methods are
(n_i, n_j, n_k are the numbers of data points in the clusters):

method	α_i	α_j	β	γ
centroid method	$\frac{n_i}{n_i+n_j}$	$\frac{n_j}{n_i+n_j}$	$-\frac{n_i n_j}{n_i+n_j}$	0
median method	$\frac{1}{2}$	$\frac{1}{2}$	$-\frac{1}{4}$	0
single linkage	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$
complete linkage	$\frac{1}{2}$	$\frac{1}{2}$	0	$+\frac{1}{2}$
average linkage	$\frac{n_i}{n_i+n_j}$	$\frac{n_j}{n_i+n_j}$	0	0
Ward's method	$\frac{n_i+n_k}{n_i+n_j+n_k}$	$\frac{n_j+n_k}{n_i+n_j+n_k}$	0	0

Choosing the Clusters

- **Simplest Approach:**

- Specify a minimum desired distance between clusters.
- Stop merging clusters if the closest two clusters are farther apart than this distance.

- **Visual Approach:**

- Merge clusters until all data points are combined into one cluster.
- Draw the dendrogram and find a good cut level.
- Advantage: Cut need not be strictly horizontal.

- **More Sophisticated Approaches:**

- Analyze the sequence of distances in the merging process.
- Try to find a step in which the distance between the two clusters merged is considerably larger than the distance of the previous step.
- Several heuristic criteria exist for this step selection.

Summary Clustering

- **Prototype-based Clustering**

- Alternating adaptation of data point assignment and cluster parameters.
- Online or batch adaptation of the cluster center.
- Crisp or fuzzy/probabilistic assignment of a datum to a cluster.
- Local minima can pose a problem.
- Fuzzy/probabilistic approaches are usually more robust.

- **Hierarchical Agglomerative Clustering**

- Start with singletons (one element clusters).
- Always merge those clusters that are closest.
- Different ways to measure the distance of clusters.
- Cluster hierarchy can be depicted as a dendrogram.

Software

Software for

- Multipolynomial and Logistic Regression,
- Bayes Classifier Induction (naive and full),
- Decision and Regression Tree Induction,
- Artificial Neural Networks (MLPs, RBFNs),
- Learning Vector Quantization,
- Fuzzy and Probabilistic Clustering,
- Association Rule Induction and Frequent Item Set Mining
- Frequent Subgraph Mining / Molecular Fragment Mining

can be found at

<http://www.borgelt.net/software.html>