

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"
Section 1.1-1.4

2 DTU Informatics, Technical University of Denmark

Lecture schedule

1. **Introduction**
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.2 +(A) + B.1)
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
Machine learning and data modelling in practice
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
(Tan 8.1-8.3 + 8.5.7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project presentation

3 DTU Informatics, Technical University of Denmark

Blooms taxonomy

4 DTU Informatics, Technical University of Denmark

Learning objectives

1. Describe the major steps involved in data modeling from preparing the data, modeling the data to evaluating and disseminating the results.
(Knowledge)
2. Discuss key machine learning concepts such as feature extraction, cross-validation, generalization and over-fitting, prediction and curse of dimensionality.
(Comprehension)
3. Sketch how the data modeling methods work and describe their assumptions and limitations.
(Knowledge and Comprehension)
4. Match practical problems to standard data modeling problems such as regression, classification, density estimation, clustering and association mining.
(Comprehension and Application)
5. Apply the data modeling framework to a broad range of application domains in medical engineering, bio-informatics, chemistry, electrical engineering and computer science.
(Application)
6. Compute the results of the data modeling framework by use of Matlab, R or Python.
(Application)
7. Use visualization techniques and statistics to evaluate model performance, identify patterns and data issues.
(Analysis)
8. Combine and modify data modeling tools in order to analyze a data set of their own and disseminate the results of the analysis.
(Application, Analysis, Synthesis and Evaluation)

5 DTU Informatics, Technical University of Denmark

Examination

- 4 hours written examination (multiple choice)
- Concentrated on the same topics as the lectures and exercises and linked to the learning objectives.
- Group reports
- Report Deadlines:
 - Report 1: 30th September, **Feature extraction and visualization**
 - Report 2: 4th November, **Supervised learning: Classification and regression**
 - Report 3: 2nd December, **Unsupervised learning: Clustering and density est.**
- Final grade based on an overall assessment of the reports and written exam

6 DTU Informatics, Technical University of Denmark

Course elements

- Lectures**
 - Top-level presentation of concepts and their relation
- Problems**
 - Train the understanding of theory
 - Prepare for the practical exercises and reports
 - Prepare for exam
- Exercises**
 - Illustrate real applications
 - Use Matlab, R, or Python to understand theoretical properties of data
 - Prepare for report work

7 DTU Informatics, Technical University of Denmark

Relevant courses at Section for Cognitive Systems

```

graph TD
    A[02450 Introduction to machine learning and data mining] --> B[028xx Digital media engineering]
    A --> C[02451 Digital signal processing]
    A --> D[02457 Nonlinear signal processing]
    A --> E[02454 Introduction to cognitive psychology]
    C -.-> F[02453 Applied digital signal processing]
    C -.-> G[02460 Advanced Machine Learning]
    D -.-> F
    D -.-> G
    E -.-> G
    E -.-> H[02458 Cognitive modeling]
  
```

8 DTU Informatics, Technical University of Denmark

Pretest

The purpose of the pretest is to assess the students background and academic level in order to allow the teachers to adjust the presentation of the course material as well as measure student learning. This pretest will not be graded and will not influence exam results in any way. We do not expect you to be able to answer all questions.

9 DTU Informatics, Technical University of Denmark

Data modeling framework

DTU

Evaluation, interpretation, and visualization

- Data preparation
 - Feature extraction
 - Similarity measures
 - Summary statistics
 - Data visualization
- Data modeling**
 - Classification
 - Regression
 - Clustering
 - Density estimation
- Evaluation**
 - Anomaly detection
 - Decision making
 - Result visualization
 - Dissemination
- Result**

Domain knowledge

11 DTU Informatik, Technical University of Denmark

We are entering the era of big data

Every day, we create 2.5 quintillion (10^{18}) bytes of data — so much that 90% of the data in the world today has been created in the last two years alone.

Source: <http://www-01.ibm.com/software/data/bigdata/>

“We are drowning in information and starving for knowledge”
John Naisbitt, author of “Megatrends”

13 DTU Informatik, Danmarks Tekniske Universitet

Introduction to Machine Learning 22-01-2013

14 DTU Informatik, Danmarks Tekniske Universitet

Introduction to Machine Learning 22-01-2013

What is machine learning?

- Arthur Samuel (1959)
 - **Machine learning:** "Field of study that gives computers the ability to learn without being explicitly programmed"
 - Samuels wrote a checkers playing program
 - Had the program play 10000 games against itself
 - Work out which board positions were good and bad depending on wins/losses
- Tom Michell (1999)
 - **Well posed learning problem:** "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E."
 - The checkers example,
 - E = 10000 games
 - T = playing checkers
 - P = if you win or not

Source: http://www.holehouse.org/mlclass/08_Neural_Networks_Representation.html

16 DTU Informatik, Danmarks Tekniske Universitet

Introduction to Machine Learning 22-01-2013

Machine learning and data mining

18 DTU Informatik, Technical University of Denmark

19 DTU Informatik, Technical University of Denmark

Human abilities examined

It is true machines are not very smart. However for many tasks humans are even worse.

Fig. 1. Proportion of rulings in favor of the prisoners by ordinal position. Circled points indicate the first decision in each of the three decision sessions; tick marks on x axis denote every third case; dotted line denotes food break. Because unequal session lengths resulted in a low number of cases for some of the later ordinal positions, the graph is based on the first 95% of the data from each session.

(Source: Danziger, S., J. Levav, and L. Avnaim-Pesso, 2011:Extraneous factors in judicial decisions. Proc. Nat. Acad. Sci.)

20 DTU Informatik, Technical University of Denmark

Human abilities examined

The number of studies reporting comparisons of clinical and statistical predictions has increased to roughly two hundred, but the score in the contest between algorithms and humans has not changed. About 60% of the studies have shown significantly better accuracy for the algorithms. The other comparisons scored a draw in accuracy, but a tie is tantamount to a win for the statistical rules, which are normally much less expensive to use than expert judgement. No exception has been convincingly documented.

The range of predicted outcomes has expanded to cover medical variables such as longevity of cancer patients, the length of hospital stays, the diagnosis of cardiac disease, and the susceptibility of babies to sudden infant death syndrome; economic measures such as the prospect of success for new businesses, the evaluation of credit risks by banks, and the future career satisfaction of workers; questions of interest to government agencies, including assessment of the suitability of foster parents, the odds of recidivism among juvenile offenders, and the likelihood of other forms of violent behaviour; and the miscellaneous outcomes such as the evaluation of scientific presentations, the winners of football games, and the future prices of Bordeaux wine. Each of these domains entails a significant degree of uncertainty and unpredictability. We describe them as “low-validity environments”. In every case, the accuracy of experts was matched or exceeded by a simple algorithm. (Kahneman 2011)

21 DTU Informatik, Technical University of Denmark

Applications

- Chemistry**
 - Spectrometry
 - Chemical sensors
- Audio processing**
 - Spoken digit classification
 - Music genre classification
- Image processing**
 - Hand-written digit recognition
 - Image tagging and classification
 - Number plate recognition
- Informatics**
 - Collaborative filtering
 - Text corpus analysis
 - Spam filters
 - Computer games
- Biomedical**
 - Micro-array gene analysis
 - Medical Imaging
- Financial data mining**
 - Market predictions
- Climate data**
 - Weather forecast

23 DTU Informatik, Technical University of Denmark

Current trend in machine learning

Simple Machine-Learning methods

- + Lots of data
- + Fast computers

= Solution to many “hard” problems

24 DTU Informatics, Technical University of Denmark

Advanced applications of machine learning

Armies of Expensive Lawyers, Replaced by Cheaper Software (NYTimes, 2011)
“a lot of people who used to be allocated to conduct document review are no longer able to be billed out,” said Bill Herr, who as a lawyer at a major chemical company used to muster auditoriums of lawyers to read documents for weeks on end.



In August 2012, the [Google autonomous car] team announced that they have completed over 300,000 autonomous-driving miles (500,000 km) accident-free ... Four U.S. states have passed laws permitting autonomous cars as of December 2013 (Wikipedia)



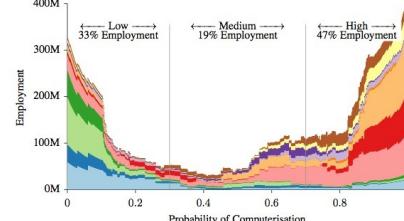
In 2011, Watson competed on Jeopardy! against former winners Brad Rutter and Ken Jennings. Watson received the first place prize of \$1 million (Wikipedia)



25 DTU Informatik, Danmarks Tekniske Universitet
Introduction to Machine Learning 22.01.2013

Machine learning as a disruptive technology

- A recent Oxford study suggest about 47% of all US jobs could be automated within two decades (Frey & Osborne, 2013)

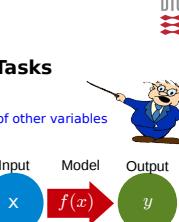


26 DTU Informatik, Danmarks Tekniske Universitet

Data Mining and Machine Learning Tasks

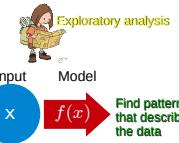
Predictive tasks (Supervised learning)

- Use some variables to predict unknown or future values of other variables
- Classification
 - Discrete output
(Determine which class a new data object belongs to)
- Regression
 - Continuous output
(Determine the output value from the input variables)



Descriptive tasks (Unsupervised learning)

- Find human-interpretable patterns that describe the data
- Clustering
 - Discover group structure in data
- Association rule discovery
 - Discover how data objects relate to each other
- Anomaly detection
 - Find data objects that are abnormal



28 DTU Informatics, Technical University of Denmark

Classification: Definition

- Given a collection of data objects (**training set**)
 - Each object has associated a number of features
 - Each object belongs to a certain class
- Define a **model** for the class given the other features
- Goal: Assign a class label to a **previously unseen object**

29 DTU Informatics, Technical University of Denmark

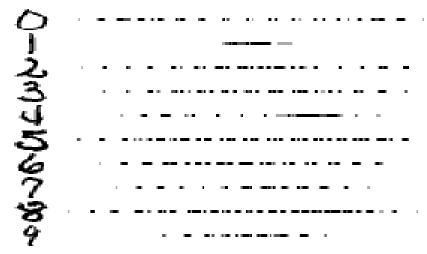
Classification: Example

Training set										Classify		
0	1	2	3	4	5	6	7	8	9	?	?	?
0	1	2	3	4	5	6	7	8	9	5	2	4
0	1	2	3	4	5	6	7	8	9	2	3	4
0	1	2	3	4	5	6	7	8	9	3	4	5
0	1	2	3	4	5	6	7	8	9	4	5	6
0	1	2	3	4	5	6	7	8	9	1	2	3

30 DTU Informatics, Technical University of Denmark

Classification: Example

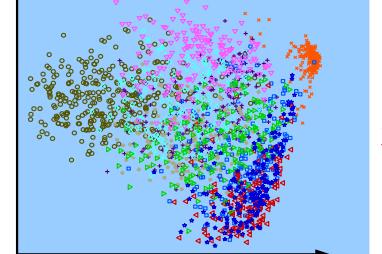
Data representation



31 DTU Informatics, Technical University of Denmark

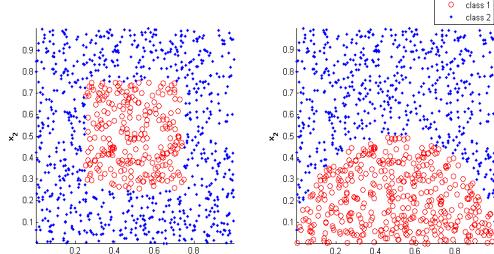
Classification: Example

Visualization



32 DTU Informatics, Technical University of Denmark

How would you characterize the two classes and how would you determine what class a new observation belongs to?



33 DTU Informatics, Technical University of Denmark

Regression: Definition

- Given a collection of data objects
 - Each object has associated a number of features
 - Each object has associated a **continuous valued variable**
- Define a **model** for the variable given the features
- Goal: Predict the value of the variable for a **previously unseen object**

34 DTU Informatics, Technical University of Denmark

Regression: Example

- Predict **sales amounts** of new product based on
 - advertising expenditure
- Predict **wind velocity** as a function of
 - temperature, humidity, and air pressure
- Predict the value of a **stock market index** based on
 - previous index time series and market indicators

35 DTU Informatics, Technical University of Denmark

Clustering: Definition

- Given a collection of data objects
 - Each object has associated a number of features
 - A measure of **similarity** between objects is defined
- Goal: **Group the objects** into clusters such that
 - Objects within each cluster are similar
 - Objects in separate clusters are less similar

36 DTU Informatics, Technical University of Denmark

Clustering: Definition

- Given a collection of data objects
 - Each object has associated a number of features
 - A measure of **similarity** between objects is defined
- Goal: **Group the objects** into clusters such that
 - Objects within each cluster are similar
 - Objects in separate clusters are less similar

37 DTU Informatics, Technical University of Denmark

Clustering: Example

Document clustering

- Goal**
 - Find groups of similar documents based on the words appearing in them
- Approach**
 - Identify frequently occurring words in each document
 - Define a similarity measure based on the word frequencies
 - Perform clustering to find groups of documents
- Motivation**
 - Use the clusters to relate a new document to existing documents
 - Better search algorithms: Return documents that are similar but do not have the exact search keywords

38 DTU Informatics, Technical University of Denmark

Association rule discovery: Definition

- Given a set of **records**
 - Each containing a number of **items from a set**
- Goal: Produce dependency rules
 - Predict the occurrence of an item based on occurrences of other items

39 DTU Informatics, Technical University of Denmark

Association rule discovery: Example

Market basket analysis

Training set

- {Bread, Coke, Milk}
- {Beer, Bread}
- {Beer, Coke, Diaper, Milk}
- {Beer, Bread, Diaper, Milk}
- {Coke, Milk}

Rules discovered

{Milk} \rightarrow {Coke}

{Diaper, Milk} \rightarrow {Beer}

40 DTU Informatics, Technical University of Denmark

Anomaly detection: Definition

- Given a collection of data objects
 - Each object has associated a number of features
- Detect which objects **deviate from normal** behaviour

41 DTU Informatics, Technical University of Denmark

Anomaly detection: Example

- Credit card **fraud detection**
 - Recognize dubious credit card transactions based on the transaction history of the card holder
- Detection of **outliers** in data measurements
 - Remove erroneous measurements due to misreading from an instrument
- Fault detection** in system health monitoring
 - Detect when a wind turbine performs poorly due to ice coating on blades

42 DTU Informatics, Technical University of Denmark

DTU


Discussion

Which of these activities are machine learning / data modeling tasks?

(*Tan Ex. 1.1*)

- a) Dividing the customers of a company according to their gender
- b) Dividing the customers of a company according to their profitability
- c) Computing the total sales of a company
- d) Sorting a student database based on student identification number
- e) Predicting the outcomes of tossing a (fair) pair of dice
- f) Predicting the future stock price of a company using historical data
- g) Monitoring the heart rate of a patient for abnormalities
- h) Monitoring seismic waves for earthquake activities
- i) Extracting the frequencies of a sound wave

43 DTU Informatics, Technical University of Denmark

A slide titled "Group discussion" featuring a cartoon illustration of three people and the DTU logo.



Exercises

All the exercises will be in Matlab, Python and R.

- Two options for running Matlab on your computer:
 - 1)Install Matlab on your computer and run it using internet connection to a license server: <http://www.gbar.dtu.dk/downloads/#>
 - 2) Run Matlab on the GBAR from your computer using thinlinc:
<http://gbar.dtu.dk/ThinLinc>
- Python and R is freely available
(exercise 1 today will guide you through how to install the programs)

You should form groups of 2-3 people for the exercises and for the 3 (group) reports.

Each group will give feedback to the teachers on the lectures and exercises of one of the course weeks.

45 DTU Informatics, Technical University of Denmark

Lecture schedule	
Report 1	<ul style="list-style-type: none">1. Introduction (Tan 1.1-1.4) Data: Feature extraction and visualization2. Data and feature extraction (Tan 2.1-2.2 +A(A)+ B.1)3. Measures of similarity and summary statistics (Tan 2.4 + 3.1-3.2 + C1-C2)4. Data visualization (Tan 3.3) Supervised learning: Classification and regression5. Decision trees and linear regression (Tan 4.1-4.3 + D)6. Overfitting and performance evaluation (Tan 4.4-4.6)7. Nearest neighbor, naive Bayes, and artificial neural networks (Tan 5.2-5.4)
Report 2	<ul style="list-style-type: none">8. Ensemble methods and multi class classifiers (Tan 5.6-5.8) Unsupervised learning: Clustering and density estimation9. K-means and hierarchical clustering (Tan 8.1-8.3 + 8.5.7)10. Mixture models and association mining (Tan 9.2.2 + 6.1-6.3)11. Density estimation and anomaly detection (Tan 10.1-10.4) Machine learning and data modelling in practice12. Putting it all together: Summary and overview13. Mini project presentation

Today's exercise:

- Form groups for project work**
For the projects 1, 2 and 3 you must work in groups of 2-3 students. The sooner you form the groups the better. Once you have formed groups please send me an email at tuh@dtu.dk with your information in the format:
 - Alice Hansen (s113589) aliceh@gmail.com
 - Bob Carlsen (s095458) bobc@yahoo.com
 - Charlie Jensen (s125199) charliej@mail.com
 - Please send the mail before 9th September
- Find a dataset to analyze throughout the course:**
Each group will find a dataset of their own. Either your own dataset, a dataset you find yourself or for instance taken from one of the resources given in the guideline "*FindingADatasetForTheReport.pdf*" on Campusnet. The 3 group reports will be based on the dataset that in turn will be analyzed by the various approaches taught during the course. Once you have found a dataset you need to have the dataset approved by me for the course.
 - Deadline for finding and having the dataset approved: 16th September
 - First report must be uploaded to campusnet 30th September. Notice file upload is closed after 13:00.
- Familiarize yourself with Matlab, Python or R:**
Todays exercise is a brush up course on either Matlab, Python or R targeted for those not familiar with these programming language. We recommend you use Matlab unless you are much more familiar with Python or R.
- Exercises are in building 358 rooms 006, 029, 031 and 032.**
Instructor with expertise in Python in rooms 006 & 029 and instructor with expertise in R in 029. All instructors have expertise in Matlab.

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"
Section 2.1-2.3 + B1 (+ A)

Feedback Groups of the day:
Rasmus Olsen
Oliver Naaby
Jesper Baltzersen
Steen Barkholt
Mads Slotsbo
Signe Bendsen Brusgaard
Amalie Sofie Ekstrand

If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + (A) + B.1)
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
(Tan 8.1-8.3)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project

Data modeling framework

Todays learning objectives:
Understand types of data, their attributes and data issues.
Be able to apply principal component analysis for data visualization and feature extraction.

What is data?

- Collection of data objects and their attributes
- An attribute is a property or characteristic of an object
 - Also known as variable, field, characteristic, or feature
- Collection of attributes describe an object
 - Also known as record, point, case, sample, entity, or instance

Attributes

ID	Age	Gender	Name
1	31	F	Alex
2	24	M	Ben
3	52	F	Cindy
4	35	M	Dan
5	58	M	Eric
6	46	F	Fay
7	42	M	George

Data objects

Discrete / continuous attributes

- **Discrete**
 - Finite (or countably infinite) set of values
 - Examples:
 - Zip codes
 - Counts
 - Set of words in a collection of documents
 - Often represented as integer variables
- **Continuous**
 - Has real numbers as attribute values
 - Examples:
 - Temperature
 - Height
 - Weight
 - Often represented as floating point variables

Types of attributes

- **Nominal:** Objects belong to a category (Equal / Not equal)
 - ID numbers
 - Eye color
 - Zip codes
- **Ordinal:** Objects can be ranked (Greater than / Less than)
 - Taste of potato chips on a scale from 1-10
 - Grades
 - Height in {short, medium, tall}
- **Interval:** Distance between objects can be measured (Addition / Subtraction)
 - Calendar dates
 - Temperature in Fahrenheit and Celcius
- **Ratio:** Zero means absence of what is measured (Multiplication / Division)
 - Length
 - Time
 - Counts
 - Temperature in Kelvin

Qualitative

Discussion

• Classify the following attributes

- a) Military rank
- b) Angles measured in degrees
- c) A persons year of birth
- d) A persons age in years
- e) Coat check number
- f) Distance from center of campus
- g) Number of patients in a hospital
- h) Sea level

Discrete

- Finite (or countably infinite) set of values

Continuous

- Real number

Nominal (Equal / Not equal)

Ordinal (Greater than / Less than)

Interval (Addition / Subtraction)

Ratio (Multiplication / Division)

• Zero means absence of what is measured

Types of data sets

- **Record data**
 - Collection of data objects and their attributes
 - Representation: Table
- **Relational data**
 - Collection of data objects and their relation
 - Representation: Graph
- **Ordered data**
 - Ordered collection of data objects
 - Representation: Sequence

Record data example: Market basket data

- Transaction data table

ID	Items
1	Bread, Soda, Milk
2	Beer, Bread
3	Beer, Soda, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Soda, Diaper, Milk

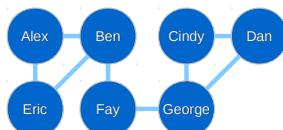
- Matrix

ID	Bread	Soda	Milk	Beer	Diaper
1	1	1	1	0	0
2	1	0	0	1	0
3	0	1	1	1	1
4	1	0	1	1	1
5	0	1	1	0	1

11 DTU Informatics, Technical University of Denmark

Relational data example: Who knows who?

- Graph



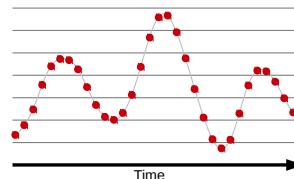
- Matrix

A	B	C	D	E	F	G
0	1	0	1	0	0	0
1	0	0	0	1	1	0
0	0	0	1	0	0	1
0	0	1	0	0	0	1
1	1	0	0	0	0	0
0	1	0	0	0	0	1
0	0	1	0	1	0	0

12 DTU Informatics, Technical University of Denmark

Ordered data example: Time series

- Sequence



- Matrix

Time	Value
0	1.3
1	1.8
2	2.5
3	3.6
4	4.4
5	4.7
6	4.6
7	4.3
8	2.4
9	2.1
10	2.0
11	2.3
12	3.1

13 DTU Informatics, Technical University of Denmark

Data quality

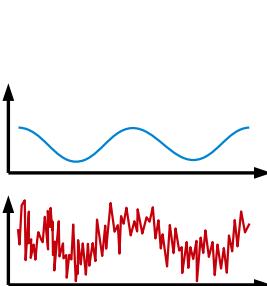
- Data is of high quality if they
 - Are fit for their intended use
 - Correctly represent the phenomena they correspond to
- Examples of quality problems
 - Noise
 - Outliers
 - Missing values



14 DTU Informatics, Technical University of Denmark

Noise

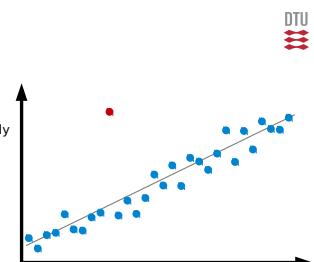
- Definition
 - Unwanted perturbation to a signal
 - Unwanted data
- Reasons for noise
 - Fundamental limits in measurement accuracy
 - Interference from other signals
 - Measurement of attributes not related to the data modeling task
- Handling noise
 - Exclude noisy attributes
 - Remove noise by filtering
 - Include a model of the noise



15 DTU Informatics, Technical University of Denmark

Outliers

- Definition
 - Data objects which are significantly different from most others
- Reasons for outliers
 - Measurement error
 - Natural property of data
- Handling outliers
 - Identify & exclude outliers
 - Model the outliers



16 DTU Informatics, Technical University of Denmark

Missing values

- Definition
 - No value is stored for an attribute in a data object
- Reasons for missing values
 - Information is not collected
 - People decline to give their age
 - Attribute is not applicable
 - Annual income is not applicable to children
- Handling missing values
 - Eliminate data objects
 - Estimate missing values (e.g. an average)
 - Ignore the missing value in analysis

ID	Age	Gender	Name
1	31	F	Alex
2	(?)	M	Ben
3	52	F	Cindy
4	35	(?)	Dan
5	(?)	M	Eric
6	(?)	F	Fay
7	42	M	(?)

17 DTU Informatics, Technical University of Denmark



Discussion

- A group of people were asked to write how many children they have
 - Their response was this

3 1 NONE 2 7 3 15 2 1 3 2 zero 4 0 1

- A research assistant typed the results into a table
 - His table looked like this

Children	3	1	0	2	7	5	15	0	1	3	-2	0	0	0	1

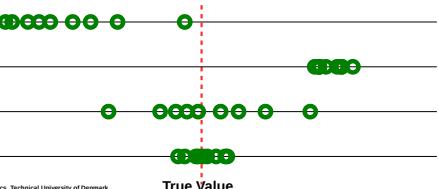
- Are there any data quality issues?
 - Noise?
 - Outliers?
 - Missing values?
- Why have these issues occurred, and how should they be handled?

18 DTU Informatics, Technical University of Denmark

Precision, Bias and Accuracy

Assume we make repeated measurements of the same underlying quantity and use this set of values to calculate a mean value (average) that serves as our estimate of the true value.

- **Definition 2.3 (Precision):** The closeness of repeated measurements (of the same quantity) to one another (often measured by standard deviation)
- **Definition 2.4 (Bias):** A systematic variation of measurements from the quantity being measured.
- **Definition 2.5 (Accuracy):** The closeness of measurements to the true value of the quantity being measured.



19 DTU Informatics, Technical University of Denmark

True Value

Data preprocessing and dimensionality reduction

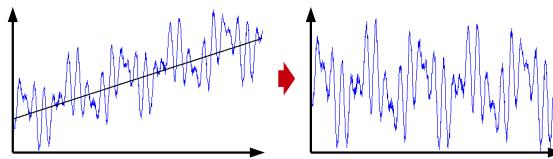
- Aggregation**
 - Combining several attributes into a single attribute
- Sampling**
 - Selecting a representative subset of data points
- Dimensionality reduction**
 - Project data to a low-dimensional subspace (**explained today**)
- Feature subset selection**
 - Choose a subset of attributes
- Feature extraction**
 - Create new features from existing attributes
- Discretization and binarization**
 - Reduce continuous attributes to discrete
- Attribute transformation**
 - Apply a fixed transformation to an attribute

20 DTU Informatics, Technical University of Denmark

Filtering

- Eliminating, suppressing, or attenuating certain aspects of the data
 - Noise removal in audio signals
 - Elimination of common words in text documents
 - Removal of background in images
 - Removal of examples which are corrupted
 - De-trending data (if it is not stationary)

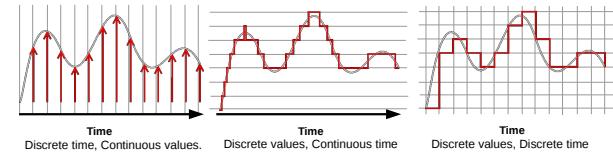
Example of de-trending data



21 DTU Informatics, Technical University of Denmark

From Analog to Digital

- Data are often **analog** and needs to be converted to a **digital** format



Discrete time, Continuous values. Discrete values, Continuous time. Discrete values, Discrete time

- Some data are "born" digital
 - On/off sensors
 - Questionnaires
 - User ratings



22 DTU Informatics, Technical University of Denmark

Chemical sensor data

- Example: The nano-nose



A	B	C	D	E	F	G	H	I	J	K	L
1	Name	Concentration									
2	Sample Type		0.5	21	6	11.02631	64.13482	21.40505	5.57374	1.17910	
3	Water		2000								
4	Water		3000	94.5	17	5	3.604615	63.25577	19.05086	4.988	1.822444
5	Water		5000	92	19	5	11.02709	62.37698	19.83532	5.12349	1.566855
6	Water		6000	53	2.5	2.5	4.663865	34.01769	6.75241	1.611996	0.055663
7	Water		4600	51	7.5	2.5	4.663865	34.01769	6.75241	1.611996	0.055663
8	Water		2000	50	4	2.5	4.663865	34.01769	6.75241	1.611996	0.055663
9	Water		2000	27.5	4	2.5	2.2	18.37597	2.798553	0.9383	0.020985
10	Water		2000	23.5	4.5	2.5	2.2	18.37597	2.798553	0.9383	0.020985
11	Water		2000	27	4	3.5	3.1	17.47115	2.6465	0.648483	0
12	Water		1350	13.5	8	0	1.700866	10.15383	1.070215	0.431933	0
13	Water		3150	13	2	0.5	2.650866	8.841545	0.620060	0.115	0
14	Water		1150	13.5	2.5	0.5	2.650866	8.841545	0.620060	0.115	0
15	Water		1350	4.5	0.5	0	0.000013	0.000013	0.000013	0	0
16	Water		375	4	0.5	0	0.000013	0.000013	0.000013	0	0
17	Water		375	3.5	0.5	0	0.000013	0.000013	0.000013	0	0
18	Never		160	1.5	n	0	0	0.147663	0.115934	0.144	n

23 DTU Informatics, Technical University of Denmark

Bag of words

- First three sentences on wikipedia.org

- The bag-of-words model is a simplifying assumption used in natural language processing and information retrieval
- In this model, a text (such as a sentence or a document) is represented as an unordered collection of words, disregarding grammar and even word order
- The bag-of-words model is used in some methods of document classification



24 DTU Informatics, Technical University of Denmark

Bag of words

- First three sentences on wikipedia.org

- The bag-of-words model is a simplifying assumption used in natural language processing and information retrieval
- In this model, a text (such as a sentence or a document) is represented as an unordered collection of words, disregarding grammar and even word order
- The bag-of-words model is used in some methods of document classification

• We will treat **this text** as a data set and create a bag-of-words model of it



25 DTU Informatics, Technical University of Denmark

Bag of words

- Elimination of common words (so-called stop words)
 - The bag-of-words model is a simplifying assumption used in natural language processing and information retrieval
 - In this model, a text (such as a sentence or a document) is represented as an unordered collection of words, disregarding grammar and even word order
 - The bag-of-words model is used in some methods of document classification

bag-of-words model
document classification

26 DTU Informatics, Technical University of Denmark

Bag of words

- Stemming

Word	Sentence	1	2	3
bag-of-word*	1	1	1	
model*	1	1	1	
simil*	1			
assum*	1			
natur*	1			
languag*	1			
process*	1			
information*	1			
retriev	1			
text*	1			
sentence*	1			
document*	1			
represent*	1			
unorder*	1			
collect*	1			
word*	2			
disregard*	1			
grammar*	1			
order*	1			
method*	1			
classif*	1			

28 DTU Informatics, Technical University of Denmark

Image representation

- Example: Handwritten digits

- Preprocessing
 - Digitalization
 - Centering
 - Rotation
 - Scaling



$$M_0 = \begin{bmatrix} 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & \dots & 0.3 & 1 & 0.2 & 0 & \dots \\ \vdots & & & & & & \vdots \\ 0 & \dots & 0 & 0 & 0 & 0 & \dots \end{bmatrix}$$

- Vectorization

$$x_0 = [0 \dots 0.3 1 0.2 0 \dots 0]^T$$

- Matrix representation of data set

$$X = \begin{bmatrix} \cdots & x_1 & \cdots \\ \cdots & x_2 & \cdots \\ \vdots & & \vdots \\ \cdots & x_N & \cdots \end{bmatrix}$$

If each image is 16 × 16 pixels then X is a $N \times 256$ matrix.

09/09/14

Vector space representation

- All these data objects have a vector space representation

$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$

30 DTU Informatics, Technical University of Denmark 09/09/14

Vectors and matrices

- Data represented as vectors and matrices
 - Linear algebra useful for manipulating and analyzing data
- We will derive the Principal Component Analysis (PCA) and discuss the Singular Value Decomposition (SVD)
 - First a (brief) highlight of linear algebra
 - PCA is very important for data visualization

31 DTU Informatics, Technical University of Denmark 09/09/14

Vectors and matrices

- Common matrix notation

$$A, A, \bar{A} \quad A = \begin{bmatrix} a_{1,1} & \cdots & a_{1,M} \\ \vdots & & \vdots \\ a_{N,1} & \cdots & a_{N,M} \end{bmatrix} \in \mathbb{R}^{N \times M}$$

- Common vector notation

$$x, \underline{x}, \vec{x} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N$$

32 DTU Informatics, Technical University of Denmark 09/09/14

Matrix Multiplication

- Two matrices can be multiplied $AB = C$
 - if the number of columns in the first equals the number of rows in the second

$A \times B = C$
 $L \times M \quad M \times N \quad L \times N$

$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot & a & \cdot \\ \cdot & \cdot & \cdot & b & \cdot \\ \cdot & \cdot & \cdot & c & \cdot \\ \cdot & \cdot & \cdot & d & \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & x_{3,4} & \cdot \end{bmatrix}$

$x_{3,4} = 1 \cdot a + 2 \cdot b + 3 \cdot c + 4 \cdot d$
 $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = ?$

33 DTU Informatics, Technical University of Denmark 09/09/14

Matrix transpose

- The transpose of a matrix

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad A^\top = \begin{bmatrix} 1 & 3 & 7 \\ 2 & 4 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

- Transpose of a sum

$$(A + B)^\top = A^\top + B^\top$$

- Transpose of a product

$$(AB)^\top = B^\top A^\top$$

$$(Ax)^\top y = x^\top A^\top y = x^\top (A^\top y)$$

34 DTU Informatics, Technical University of Denmark 09/09/14

The identity matrix

- Ones on the diagonal and zeros everywhere else

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad I^\top = I$$

- Multiplying by the identity does not change anything

$$IA = A$$

$$I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$I_2 A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

35 DTU Informatics, Technical University of Denmark 09/09/14

Matrix inverses

- For a square matrix, the inverse satisfies

$$AA^{-1} = A^{-1}A = I$$

- Inverse of a product of square matrices

$$(AB)^{-1} = B^{-1}A^{-1}$$

- Transpose of inverse

$$(A^{-1})^\top = (A^\top)^{-1}$$

36 DTU Informatics, Technical University of Denmark 09/09/14

Norms

- The norm of a vector is usually written as

$$\|x\|_p = \left(\sum_{n=1}^N |x_n|^p \right)^{1/p}$$

- Of particular interest is the 1-norm and the 2-norm

$$\|x\|_1 = \sum_{n=1}^N |x_n| \quad \|x\|_2 = \sqrt{\sum_{n=1}^N |x_n|^2} = \sqrt{x^\top x}$$

- The Frobenius norm of a matrix

$$\|X\|_F^2 = \sum_{i,j} x_{i,j}^2 = \text{trace}(XX^T) = \text{trace}(X^T X)$$

Where trace takes the sum of the diagonal elements, i.e.

$$\text{trace}(A) = \sum_{i=1}^N a_{i,i}$$

37 DTU Informatics, Technical University of Denmark 09/09/14

Discussion

- I want to go from **61st and West End Ave.** to **110th and Central Park West**
- I have created a coordinate system and a vector to my destination

- Which vector norm should I use to measure the length of my trip?

$$\|x\|_p = \left(\sum_{n=1}^N |x_n|^p \right)^{1/p}$$

$$\|x\|_1 = \sum_{n=1}^N |x_n|$$

$$\|x\|_2 = \sqrt{\sum_{n=1}^N |x_n|^2} = \sqrt{x^\top x}$$

38 DTU Informatics, Technical University of Denmark 09/09/14

Vector spaces

- Vector spaces can be of arbitrary size
- Typically defined using a matrix of basis vectors
- The basis vectors must be **linearly independent**

$$Vx = 0 \Rightarrow x = 0$$

$$V = \begin{bmatrix} | & | & | \\ v_1 & v_2 & \cdots & v_N \\ | & | & | \end{bmatrix}$$

39 DTU Informatics, Technical University of Denmark 09/09/14

Vector spaces

- Often the vectors are taken to be **mutually orthogonal** and of **unit length**

$$v_i^\top v_j = 0 \quad \|v_i\|_2 = 1$$

- This defines an orthonormal basis for the vector space

- An orthonormal matrix satisfies

$$V^\top V = I, \quad V^\top = V^{-1}$$

40 DTU Informatics, Technical University of Denmark 09/09/14

Subspaces

41 DTU Informatics, Technical University of Denmark 09/09/14

Projection

- Projection onto a vector

42 DTU Informatics, Technical University of Denmark 09/09/14

- Angle between vectors

$$\cos(\phi) = \frac{\mathbf{v}^\top \mathbf{x}}{\|\mathbf{x}\|_2 \|\mathbf{v}\|_2}$$

- Length of projection

$$p = \|\mathbf{x}\|_2 \cos(\phi) = \frac{\mathbf{v}^\top \mathbf{x}}{\|\mathbf{v}\|_2}$$

- Projection onto unit vector

$$p = \mathbf{v}^\top \mathbf{x}$$

09/09/14

Projection onto subspace

- **Projection onto a subspace**
 - Subspace defined by an orthonormal basis matrix
 - Projection given by

$$\mathbf{x}^\top V$$

- **Example:** Projection of 3-D vector onto the (x,z) plane

$$V = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

$$\mathbf{x}^\top V = \begin{bmatrix} x & y & z \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} x & z \end{bmatrix}$$

43 DTU Informatics, Technical University of Denmark 09/09/14

Visualizing high dimensional data

- Humans are good at perceiving data in 2-D
 - However we cannot visualize high dimensional data
- We can project high dimensional data to a lower dimensional subspace
- But what is a good projection?
 - Account for as much of the variation as possible

44 DTU Informatics, Technical University of Denmark 09/09/14

Singular Value Decomposition (SVD)

(Eugenio Beltrami & Camille Jordan, independently, 1873-1874)

Any N x M matrix can be decomposed as follows:

$$X = U \Sigma V^\top$$

X is $N \times M$. U is $N \times N$ (Orthonormal). Σ is $N \times M$ (Diagonal). V is $M \times M$ (Orthonormal).

$$U = \begin{bmatrix} | & \cdots & | \\ u_1 & \cdots & u_N \\ | & \cdots & | \end{bmatrix} \quad V = \begin{bmatrix} | & \cdots & | \\ v_1 & \cdots & v_M \\ | & \cdots & | \end{bmatrix}$$

$\sigma_1, \dots, \sigma_M$ is known as the singular values

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_M \\ 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}$$

$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_M \geq 0$

If $i \neq j$: $\Sigma_{i,j} = 0$, $U^\top U = I_{N \times N}$, $V^\top V = I_{M \times M}$

45 DTU Informatics, Technical University of Denmark 09/09/14

Principal component analysis (PCA)

(Karl Pearson, 1901)

- 1) Subtract the mean from each attribute
- 2) Apply singular value decomposition (SVD)

$$X = U \Sigma V^\top$$

X is $N \times M$. U is $N \times N$. Σ is $N \times M$. V is $M \times M$.

- 3) Select first K columns of V (the PCA projection operation) and first K columns of Σ .

$$\hat{X} = U \Sigma_{(K)}$$

\hat{X} is $N \times K$ (PCA components or PCA projection of the data).

$$(X^\top X)v_k = \sigma_k^2 v_k$$

v_k is $M \times K$ (PCA loadings).

46 DTU Informatics, Technical University of Denmark 09/09/14

Principal component analysis

$$X = U \Sigma V^\top$$

- Entries in the diagonal matrix Σ are called **singular values**
 - They are sorted (largest first)
 - Indicate how much variability is explained by the corresponding component
 - 1st component explains most of the variability
 - 2nd component explains most of the remaining variability
 - Etc.

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N \end{bmatrix} \quad \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_N \geq 0$$

- Singular value spectrum

48 DTU Informatics, Technical University of Denmark 09/09/14



- Show that

$$\|X\|_F^2 = \sum_i \sigma_i^2$$

where
 $\sigma_i = \sqrt{\lambda_i}$

Fraction of the variation in the data explained by the i^{th} principal component is given by:

$$\frac{\sigma_i^2}{\sum_j \sigma_j^2}$$

And by the first K principal components

$$\frac{\sum_{i=1}^K \sigma_i^2}{\sum_j \sigma_j^2}$$

Hints:

$$X = U\Sigma V^T$$

$$\|X\|_F^2 = \text{trace}(XX^T)$$

$$\text{trace}(AB) = \text{trace}(BA)$$

Fishers Iris Data

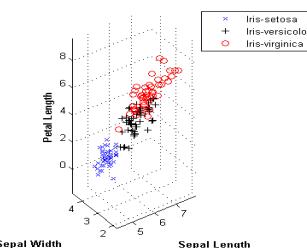


Three types of flowers:
Iris Setosa, Iris Versicolor, Iris Virginica

Flower ID	Attribute			Petal Width
	Sepal Length	Sepal Width	Petal Length	
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
.
.
150	5.9	3.0	5.1	1.8

We will presently consider the first 3 attributes, i.e. Sepal length, Sepal Width and Petal Length.

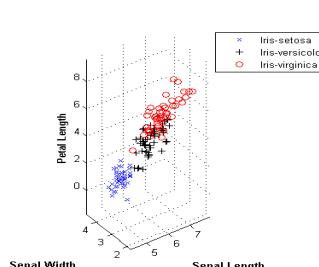
3D scatter plot of the data



What fraction of the total variation in the data will the first principal component account for?



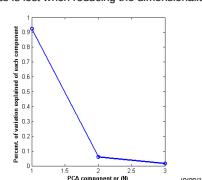
3D scatter plot of the data



- Subtract the mean
- Apply singular value decomposition (SVD)

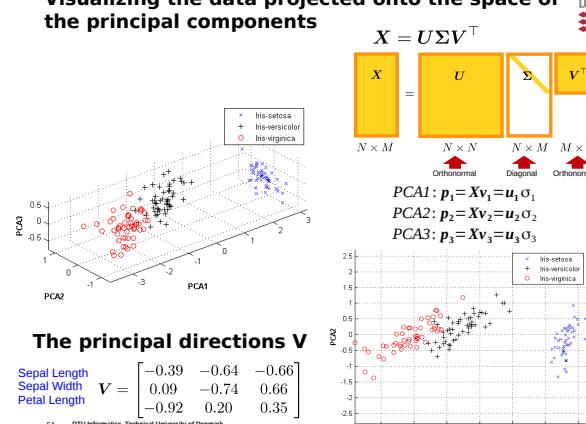
$$X = U\Sigma V^T$$

Evaluate the singular values to determine how much of the dynamics is lost when reducing the dimensionality



What fraction of the total variation in the data will the first principal component account for?

Visualizing the data projected onto the space of the principal components



The principal directions V

$$V = \begin{bmatrix} -0.39 & -0.64 & -0.66 \\ 0.09 & -0.74 & 0.66 \\ -0.92 & 0.20 & 0.35 \end{bmatrix}$$

$$X = U\Sigma V^T$$

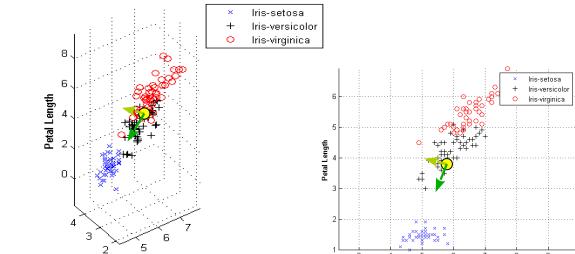
$$\text{PCA1: } p_1 = Xv_1 = u_1\sigma_1$$

$$\text{PCA2: } p_2 = Xv_2 = u_2\sigma_2$$

$$\text{PCA3: } p_3 = Xv_3 = u_3\sigma_3$$

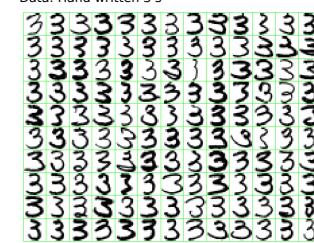
$$\mu = \begin{bmatrix} 5.8 \\ 3.1 \\ 3.8 \end{bmatrix}, \quad v_1 = \begin{bmatrix} -0.39 \\ 0.09 \\ -0.92 \end{bmatrix}, \quad v_2 = \begin{bmatrix} -0.64 \\ -0.74 \\ 0.20 \end{bmatrix}$$

Sepal Length
Sepal Width
Petal Length



Visualization of hand written digits

- Data: Hand written 3's



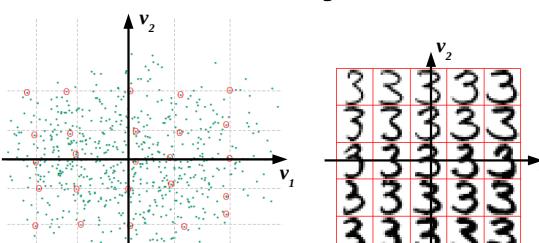
- Data matrix

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

If each image is 16×16 pixels then X is a $N \times 256$ matrix.

$$X = U\Sigma V^T$$

Visualization of hand written digits



(Figures taken from "Elements of Statistical learning" by Trevor Hastie, Robert Tibshirani, Jerome Friedman)

What is the dynamics captured by the first two principal components?

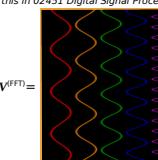
Data and Domain driven feature extraction

PCA is an example of a data driven approach for feature extraction
i.e., we define from data the features extracted in terms of the projections $V^{(PCA)}$ that preserve most of the variation in the data

$$X = U\Sigma V^T$$

The fourier transform is an example of a domain driven approach for feature extraction

i.e., in the analysis of sound good features are often to use spectral representations. These can be found by projecting the data using the so-called fourier transform matrix $V^{(FFT)}$ where the components are defined as specific frequencies such that the projection of the data onto these frequencies defines the extend to which these frequencies are present in the data. (you can learn much more about this in 02451 Digital Signal Processing)



02450 Introduction to machine learning and data mining (Preliminary Version)

DTU Informatics
Department of Informatics and Mathematical Modeling

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"
Section 2.4 + 3.1-3.2 + C1-C2

Groups of the day:
TBA

Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. **Measures of similarity and summary statistics**
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
(Tan 8.1-8.3+8.5.7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project

Data modeling framework

Evaluation, interpretation, and visualization

Data
Data preparation
• Feature extraction
• Similarity measures
• Summary statistics
• Data visualization

Data modeling
• Classification
• Regression
• Clustering
• Density estimation

Evaluation
• Anomaly detection
• Decision making
• Result visualization
• Dissemination

Result

Domain knowledge

Todays learning objectives:
Be able to calculate various measures of similarity and dissimilarity.
Understand how various summary statistics are calculated and can be interpreted.
Explain and apply Bayes theorem
Understand the normal and multi-variate normal distribution and the role of the covariance matrix

Example: Principal component analysis of images

- 1000 images, 86 x 86 pixels, 3 RGB intensities
Tamara Berg "Faces in the wild"

Preprocessing

- Each image
 - 86 x 86 pixels
 - 3 RGB intensities
- Split image into red, green, and blue color channels

Preprocessing

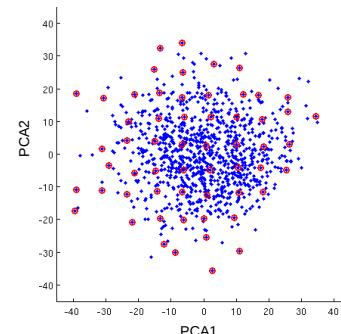
- Concatenate all pixel color values in one long vector
 - 86 x 86 x 3 = 22,188
 - Image is now represented as a 22,188 dimensional vector
- Put all 1000 images into a big matrix
 - 1000 x 22,188

Principal component analysis (PCA)

1. Subtract the mean
2. Compute the singular value decomposition (SVD)
 - Orthogonal linear transformation
 - Transforms data to a new coordinate system
 - Greatest variance along the first axis
 - Second greatest variance along the second axis
 - Etc.
- Plot data in the transformed coordinate system
 - Corresponds to looking at data from an angle where it is most spread out

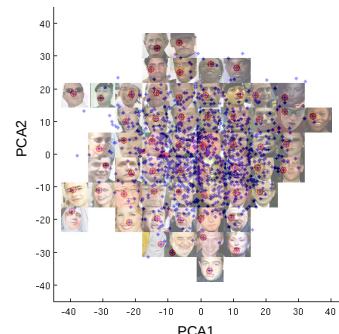
PCA of face images

PCA of face images



10 DTU Informatics, Technical University of Denmark

PCA of face images



11 DTU Informatics, Technical University of Denmark

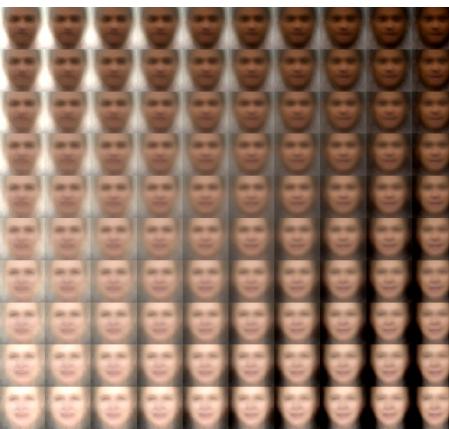
Discussion



12 DTU Informatics, Technical University of Denmark



- What information do the two principal axes capture?



13 DTU Informatics, Technical University of Denmark



- What information do the two principal axes capture?



14 DTU Informatics, Technical University of Denmark

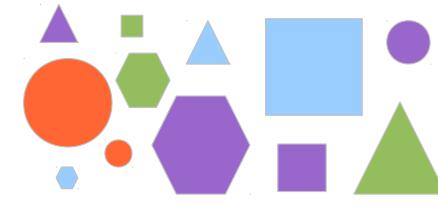
Similarity / Dissimilarity

• Definition

- A numerical measure of *how alike/different* two data objects are
- Often defined on the interval [0,1]

• Similarity / dissimilarity between two data objects

$$s(x, y), \quad d(x, y)$$



15 DTU Informatics, Technical University of Denmark

Dissimilarity measures

• Euclidean distance (2-norm)

$$d_2(x, y) = \sqrt{\sum_{n=1}^N (x_n - y_n)^2}$$

• Minkowski distance (p-norm)

$$d_p(x, y) = \left(\sum_{n=1}^N |x_n - y_n|^p \right)^{1/p}$$

Properties of dissimilarity measures

• Positivity

$$\begin{aligned} d(x, y) &\geq 0 \\ d(x, y) &= 0 \Leftrightarrow x = y \end{aligned}$$

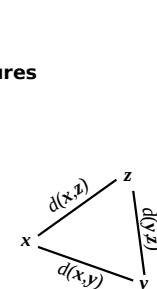
• Symmetry

$$d(x, y) = d(y, x)$$

• Triangle inequality

$$d(x, z) \leq d(x, y) + d(y, z)$$

• Measures that satisfy all these properties: **Metrics**

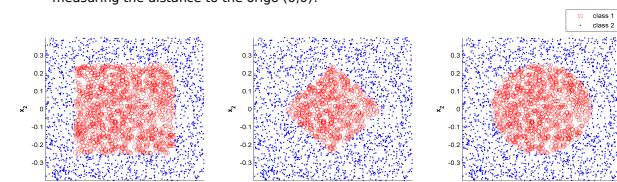


16 DTU Informatics, Technical University of Denmark

17 DTU Informatics, Technical University of Denmark

Minkowski distance

Which Minkowski distance (p -norm) can be used to separate the two classes, by measuring the distance to the origo (0,0)?



$$d_p(x, y) = \left(\sum_{n=1}^N |x_n - y_n|^p \right)^{1/p}$$

18 DTU Informatics, Technical University of Denmark

Binary similarity measures

- Simple matching coefficient (SMC)

- Symmetric: Counts present and absent attributes equally

$$\text{SMC}(x, y) = \frac{f_{00} + f_{11}}{K}$$

- Jaccard coefficient

- Asymmetric: Counts only present attributes

$$J(x, y) = \frac{f_{11}}{K - f_{00}}$$

K : Total number of attributes

f_{00} : Number of attributes where $x_k = y_k = 0$

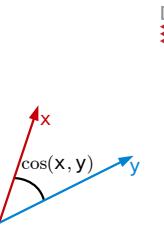
f_{11} : Number of attributes where $x_k = y_k = 1$

19 DTU Informatics, Technical University of Denmark

Continuous similarity measures

- Cosine similarity

$$\cos(x, y) = \frac{x^\top y}{\|x\|_2 \|y\|_2}$$



- Extended Jaccard coefficient

$$\text{EJ}(x, y) = \frac{x^\top y}{\|x\|_2^2 + \|y\|_2^2 - x^\top y}$$

20 DTU Informatics, Technical University of Denmark

Empirical statistics

- Empirical mean

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

- Empirical covariance

$$\text{cov}(x, y) = \frac{1}{N-1} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

- Empirical variance

$$\text{var}(x) = \text{cov}(x, x) = \frac{1}{N-1} \sum_{n=1}^N (x_n - \bar{x})^2$$

- Empirical standard deviation

$$\text{std}(x) = \sqrt{\text{var}(x)}$$

22 DTU Informatics, Technical University of Denmark

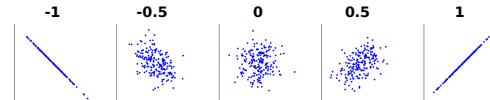
Correlation

- Measure of linear relation

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\text{std}(x)\text{std}(y)}$$

- A correlation of **1** or **-1** means there is a perfect linear relation

$$x_k = ay_k + b$$



23 DTU Informatics, Technical University of Denmark

Calculate the SMC, Jaccard, Cosine and Extended Jaccard similarity between customer 1 and customer 2 in the market basket data below.

ID	Bread	Soda	Milk	Beer	Diaper
1	1	1	1	0	0
2	0	1	1	0	1

$$\text{SMC}(x, y) = \frac{f_{00} + f_{11}}{K}$$

$$J(x, y) = \frac{f_{11}}{K - f_{00}}$$

$$\cos(x, y) = \frac{x^\top y}{\|x\|_2 \|y\|_2}$$

$$\text{EJ}(x, y) = \frac{x^\top y}{\|x\|_2^2 + \|y\|_2^2 - x^\top y}$$

K : Total number of attributes

f_{00} : Number of attributes where $x_k = y_k = 0$

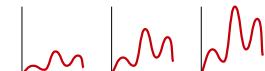
f_{11} : Number of attributes where $x_k = y_k = 1$

21 DTU Informatics, Technical University of Denmark

Invariance

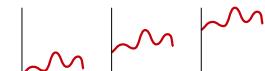
- Scale

$$d(x, y) = d(\alpha x, y)$$



- Translation

$$d(x, y) = d(\beta + x, y)$$

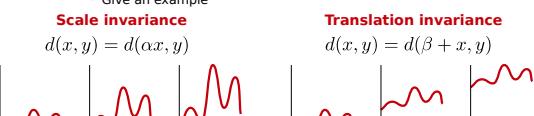


24 DTU Informatics, Technical University of Denmark

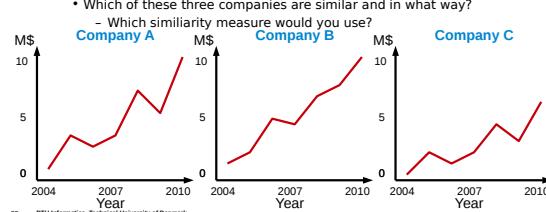
Discussion

- When would a **scale invariant** similarity measure be useful
 - Give an example

- When would a **translation invariant** similarity measure be useful
 - Give an example



- Which of these three companies are similar and in what way?
- Which similarity measure would you use?



25 DTU Informatics, Technical University of Denmark

Issues in proximity calculation

- What to do when attributes have different

- Scale?

- Standardize attributes
- Use scale invariant similarity measure

- Type?

- Compute similarities for each attribute and combine

- Importance?

- Compute a weighted similarity measure

26 DTU Informatics, Technical University of Denmark

Standardization

- Attributes have different scales.

- Example:

- Number of children ~ 0-5
- Age ~ 0-100 years
- Annual income ~ 0-50.000 €

- Unless we do something, Annual income will dominate

- Standardization: Subtract mean and divide by standard deviation

$$x_k^* = \frac{x_k - \bar{x}_k}{\text{std}(x_k)}$$

27 DTU Informatics, Technical University of Denmark

Combining heterogeneous attributes

- Attributes have different type

- Example:
 - Age:** Continuous
 - Education:** Binary
 - Primary (yes/no)
 - Secondary (yes/no)
 - Tertiary (yes/no)

- Similarity measure must handle **continuous** and **binary** features

- Compute similarities for each attribute and combine

$$s_{\text{Age}} = d_1(x_{\text{Age}}, y_{\text{Age}})^{-1} \quad s_{\text{Edu.}} = \text{SMC}(x_{\text{Edu.}}, y_{\text{Edu.}})$$

$$s(x, y) = \frac{1}{2}(s_{\text{Age}} + s_{\text{Edu.}})$$

Weighting attributes by importance

- Attributes have different importance

- Example:
 - Age:** Very important
 - Education:** Less important

- Similarity measure must take **importance** into account

- Introduce **importance weights** for each attribute

$$s_{\text{Age}} = d_1(x_{\text{Age}}, y_{\text{Age}})^{-1} \quad s_{\text{Edu.}} = \text{SMC}(x_{\text{Edu.}}, y_{\text{Edu.}})$$

$$s(x, y) = 0.99 \cdot s_{\text{Age}} + 0.01 \cdot s_{\text{Edu.}}$$



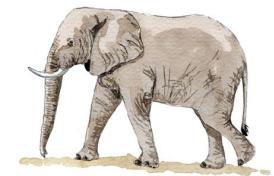
Discussion

- The following **attributes** are measured for a herd of elephants

- Weight
- Height
- Tusk length
- Trunk length
- Ear area
- Gender

- Based on these measurements

- How would you evaluate how similar elephants are?
- Justify your answer



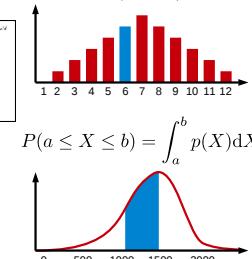
Probabilities

- Discrete: Probability mass

- Example: The sum of two dice

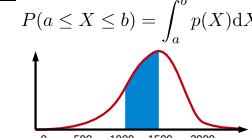


$$P(X = v)$$



- Continuous: Probability density

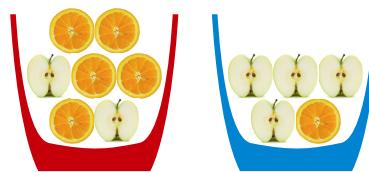
- Example: Lifetime of a light bulb



Probabilities

- What is the probability of an **orange** if the bowl is **red**?

- What is the probability of the **red** bowl if the fruit is **orange**?



Probabilities

- Basic rules of probability

- Sum rule

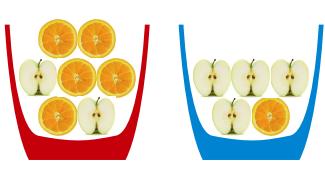
$$p(x, y) = \sum_y p(x|y)p(y)$$

- Product rule

$$p(x, y) = p(x|y)p(y)$$

- Bayes' rule

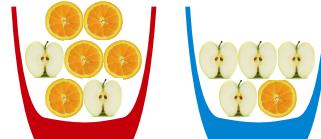
$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$



Probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

2	5
4	1



Probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

2	5
4	1
6	12

Probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

2	5
4	1
6	12

$$p(o|r) = \frac{p(r, o)}{p(r)} = \frac{5/12}{7/12} = \frac{5}{7}$$

Probabilities

- What is the probability of an orange if the bowl is red?
- What is the probability of the red bowl if the fruit is orange?

2	5
4	1
6	12

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$p(r|o) = \frac{p(r,o)}{p(o)} = \frac{5/12}{6/12} = 5/6$$

37 DTU Informatics, Technical University of Denmark

Probabilities

- What is the probability of an orange if the bowl is red?
- What is the probability of the red bowl if the fruit is orange?

2	5
4	1
6	12

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$\begin{aligned} p(r|o) &= \frac{p(r,o)}{p(o)} = \frac{5/12}{6/12} = 5/6 \\ &= \frac{p(o|r)p(r)}{p(o)} \end{aligned}$$

38 DTU Informatics, Technical University of Denmark

Probabilities

- What is the probability of an orange if the bowl is red?
- What is the probability of the red bowl if the fruit is orange?

2	5
4	1
6	12

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$\begin{aligned} p(r|o) &= \frac{p(r,o)}{p(o)} = \frac{5/12}{6/12} = 5/6 \\ &= \frac{p(o|r)p(r)}{p(o)} \\ &= \frac{5/7 \cdot 7/12}{6/12} = 5/6 \end{aligned}$$

39 DTU Informatics, Technical University of Denmark



Medical test

A medical test for a given disease

- Correctly identifies the disease 99% of the time (true positives), and
- Incorrectly turns out positive 2% of the time (false positives).

You know that

- 1% of the population suffers from the disease.

You go to the doctor to get tested, and the test turns out to be positive.

What is the probability you have the disease?

Hints:

- Identify from the text:
 $p(\text{Positive}|\text{Disease})$
 $p(\text{Positive}|\text{No Disease})$
 $p(\text{Disease})$
 $p(\text{No Disease})$
- Use the basic rules of probability given to the right to find:
 $p(\text{Disease}|\text{Positive})$

$$p(x) = \sum_y p(x,y)$$

$$p(x,y) = p(x|y)p(y)$$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

40 DTU Informatics, Technical University of Denmark



Frequency and mode

- Frequency:** Percentage of time a value occurs
 - Example: Given the attribute **Gender** and a representative population of people, the value **Female** occurs about 50% of the time
- Mode:** The most frequent attribute value
 - Example: Given the attribute **Operating System** and a representative population of computers, the value **Microsoft Windows** is the mode
- The notions of frequency and mode are typically used with categorical data

41 DTU Informatics, Technical University of Denmark



Percentiles

- Percentiles:** Given an ordinal or continuous attribute **x** and a number **p** between 0 and 100, the **p**th percentile is a value **x_p** of **x** such that **p** percent of the observed values of **x** are less than **x_p**.
 - Example: The 10th percentile of **x** is the value **x_{10%}**, such that 10% of all values are less than **x_{10%}**.
- Median:** The 50th percentile
 - Sort the numbers and take the middle value (if there are an even number of values, average the two middle values)

42 DTU Informatics, Technical University of Denmark

Measures of location

Mean: Average

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

Median: Sort the numbers

$$\text{median}(x) = \begin{cases} x_{R+1} & \text{if } N \text{ is odd } (N = 2R + 1) \\ \frac{1}{2}(x_R + x_{R+1}) & \text{if } N \text{ is even } (N = 2R) \end{cases}$$

43 DTU Informatics, Technical University of Denmark



Discussion

- What are the frequencies and the mode of these numbers and what is the mean and median?

0, 1, 1, 3, 5, 590

- Explain also what to be careful about when using the mean and median

- Frequency:** Percentage of time a value occurs
- Mode:** The most frequent attribute value
- Mean:** Average

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

Median: Sort the numbers

$$\text{median}(x) = \begin{cases} x_{R+1} & \text{if } N \text{ is odd } (N = 2R + 1) \\ \frac{1}{2}(x_R + x_{R+1}) & \text{if } N \text{ is even } (N = 2R) \end{cases}$$

44 DTU Informatics, Technical University of Denmark



Measures of spread

Range

$$\text{range}(x) = \max(x) - \min(x)$$

Variance

$$\text{variance}(x) = s_x^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \bar{x})^2$$

Absolute average deviation (AAD)

$$\text{AAD}(x) = \frac{1}{N} \sum_{n=1}^N |x_n - \bar{x}|$$

Median absolute deviation (MAD)

$$\text{MAD}(x) = \text{median}\{|x_1 - \bar{x}|, \dots, |x_N - \bar{x}|\}$$

Interquartile range (IQR)

$$\text{IQR}(x) = x_{75\%} - x_{25\%}$$

45 DTU Informatics, Technical University of Denmark

Expected values

- Discrete random variable

$$\mathbb{E}[g(X)] = \sum_i g(x_i)P(X = x_i)$$

- Continuous random variable

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(X)p(X)dX$$

46 DTU Informatics, Technical University of Denmark

Statistics

- Mean

$$\bar{x} = \mathbb{E}[x]$$

- Covariance

$$\text{cov}(x, y) = \mathbb{E}[(x - \bar{x})(y - \bar{y})]$$

- Variance

$$\text{var}(x) = \text{cov}(x, x) = \mathbb{E}[(x - \bar{x})^2]$$

- Standard deviation

$$\text{std}(x) = \sqrt{\text{var}(x)}$$

47 DTU Informatics, Technical University of Denmark

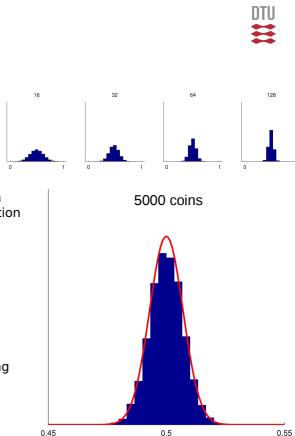
Normal distribution

- Central limit theorem

- The mean of a large number of random variables will tend to a Normal distribution irrespective of the distribution of the random variables
(Under certain conditions)

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

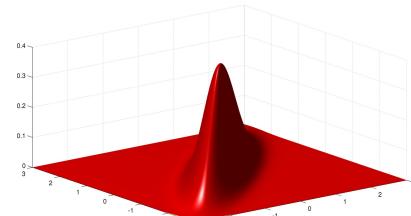
- Example: Proportion of heads when flipping - 1 coin, 2 coins, 4 coins etc.



48 DTU Informatics, Technical University of Denmark

Multivariate Normal distribution

$$p(x) = \frac{1}{(2\pi)^{k/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$



49 DTU Informatics, Technical University of Denmark

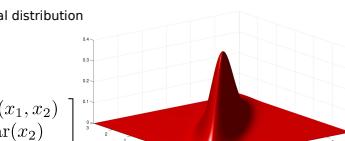
Multivariate Normal distribution

$$p(x) = \frac{1}{(2\pi)^{k/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

- Example: 2-dimensional Normal distribution

$$\mu = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

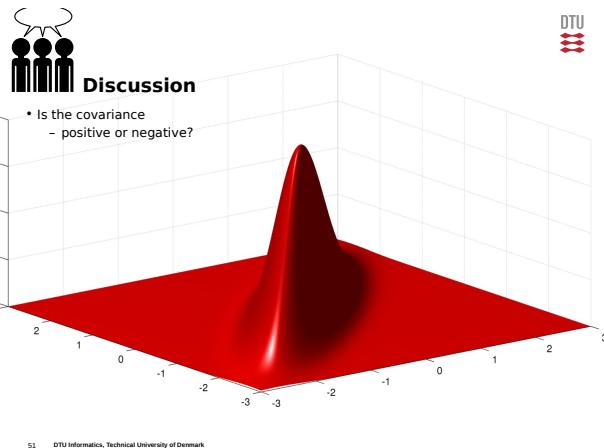
$$\Sigma = \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{var}(x_2) \end{bmatrix}$$



50 DTU Informatics, Technical University of Denmark

Discussion

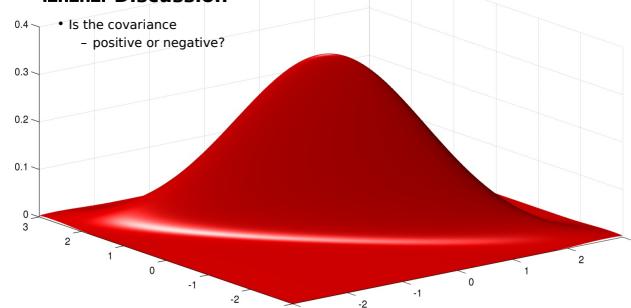
- Is the covariance positive or negative?



51 DTU Informatics, Technical University of Denmark

Discussion

- Is the covariance positive or negative?



52 DTU Informatics, Technical University of Denmark

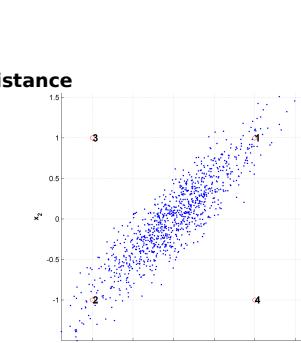
The Mahalanobis distance

- How far are x_1 and x_2 apart?

- $\text{mahalanobis}(x_1, x_2) = 4.2$
- $d_{\text{Euclidean}}(x_1, x_2)^2 = 8.0$

- How far are x_3 and x_4 apart?

- $\text{mahalanobis}(x_3, x_4) = 80$
- $d_{\text{Euclidean}}(x_3, x_4)^2 = 8.0$



$$\text{mahalanobis}(x, y) = (x - y)^\top \Sigma^{-1} (x - y)$$

$$d_{\text{Euclidean}}(x, y)^2 = (x - y)^\top \Sigma^{-1} (x - y)$$

53 DTU Informatics, Technical University of Denmark

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

$f(x+\Delta x) = \sum_{i=1}^n f_i(x)$

$\int_a^b \Theta + \Omega \int \delta e^{ix} = [2.7182818284] \sum!$

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"

Section 3.3

Feedback Groups of the day:
Søren Trads Steen
Janus Nortoft Jensen

Veli Kerim Celik
Fadi Abdul Halim
Ziad Bahawan

Andreas Jacobsen
Asger Anker Sørensen
Phong Le Trung

Anders Ulrik Ellisen
Rasmus Sten Andersen
Rasmus Lund Bendtsen

If possible, please (1) stay to give me feedback after the second exercise today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

INTRODUCTION TO DATA MINING
PANG-NING TAN
MICHAEL STEINBACH
VIPIN KUMAR

Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. **Data visualization**
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
(Tan 8.1-8.3, 8.5.7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project

Probabilities (revisited from last week)

- Basic rules of probability
 - Sum rule
- Product rule
- $p(x, y) = p(x|y)p(y)$
- Bayes' rule

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

4 DTU Informatics, Technical University of Denmark

Probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

		2	5	7	
		4	1	5	
2	6	6	6	12	

5 DTU Informatics, Technical University of Denmark

Probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

		2	5	7	
		4	1	5	
2	6	6	6	12	

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$p(r|o) = \frac{p(r,o)}{p(o)} = \frac{5/12}{5/12} = 5/5 = 1$$

$$= \frac{p(o|r)p(r)}{p(o)}$$

$$= \frac{5/7 \cdot 7/12}{5/12} = 1/2$$

6 DTU Informatics, Technical University of Denmark

The news paper "Slam the Glam"

News media agency Reuters Bureau sends stories to the Economist:

- 80% of the news stories from Reuters are positive and 20% of the news stories are negative.
- 90% of the negative news stories are published in the Economist while only 5% of the positive stories are published.

What is the probability that a story from Reuters is published in the Economist given it is positive?

Hints:

- Sum Rule
- Product rule
- Bayes theorem

$$p(x) = \sum_y p(x,y)$$

$$p(x,y) = p(x|y)p(y)$$

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

7 DTU Informatics, Technical University of Denmark

Data modeling framework

Evaluation, interpretation, and visualization

Data

Evaluation
• Anomaly detection
• Decision making
• Feature extraction
• Feature selection
• Density estimation
• Dissemination
• Feature extraction
• Feature selection
• Density estimation

Domain knowledge

Result

Todays learning objectives:
Be able to understand and apply a wide range of data visualization approaches
Understand the ACCENT principle and Tufte's guidelines

8 DTU Informatics, Technical University of Denmark

Visualization

- A main function of the brain is to process visual information
- We can exploit this capacity using visualization of the data in order to:
 - Detect new patterns, i.e. exploratory data analysis
 - **Disseminate results, i.e. visualizations/plots in written work (todays lecture)**
- Since we are making use of the brains visual system we should take into account how the visual system work

Gun deaths in Florida

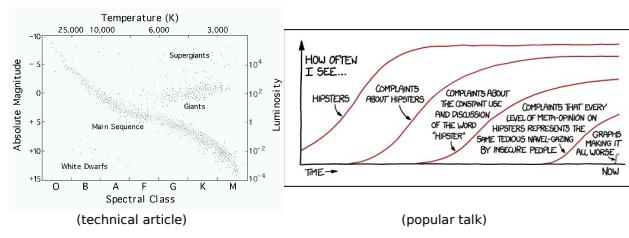
Number of murders committed using firearms

Source: Florida Department of Law Enforcement

9 DTU Informatics, Technical University of Denmark

Visualization

- A visualization is an attempt to communicate/convince a reader about a particular point supported by the data.
- Keep the audience in mind (compare a popular talk to a technical article)



12 DTU Informatics, Technical University of Denmark

(Sources: <http://cass.ucsd.edu/>, <http://xkcd.com>)

Three questions to ask of any visualization

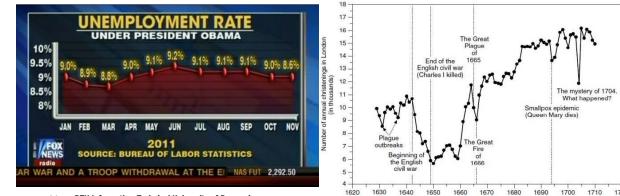
(From <http://junkcharts.typepad.com/> (JC))

- What is the **Question**?
 - Any visualization should attempt to answer a particular question or argue a particular point.
 - Without a question there is no need for the visualization
 - What does the **Data** say?
 - The data should be relevant for the question.
 - Applying appropriate transformations (not to few, not to many)
 - discarding irrelevant attributes
 - remove outliers
 - What does the **Visual** say?
 - The visual element being chosen should address the question being asked using the data
 - Type of plot, use of colors, use of other graphical elements to help/confuse the reader
- This talk will mainly consider the **visual** question

13 DTU Informatics, Technical University of Denmark

Visualization as technical writing

- Just as there is no recipe for writing, there are no definite recipes for plotting
- However there are guidelines, often more apparent for writing than for illustrations
 - The purpose of the text is to communicate an idea (*vs. plots has a purpose*)
 - Be grammatically correct (*vs. elementary "rules" of good plotting*)
 - Ensure the text is readable (*vs. labels, legends or lines nobody can read*)
 - Avoid long/complicated paragraph (*vs. plots that are overly complicated*)
 - Dont lie or gloss over inconvenient facts. (*vs. distort data in a plot*)



14 DTU Informatics, Technical University of Denmark

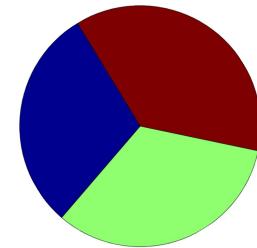
Common difficulties for visualizations

- Representation:** How will you map objects, attributes, and relations to visual elements?
 - Positions, lengths, colors, areas, orientation
- Arrangement:** How will you display the visual elements?
 - Viewpoint, transparency, separation, grouping
- Selection:** How will you handle a large number of attributes and data objects?
 - Display a subset, focus on a region of interest, show summaries

16 DTU Informatics, Technical University of Denmark

Representation

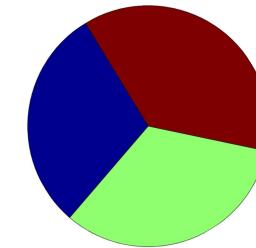
- Area represents proportion**
 - Which is smallest, middle, and largest?
 - What are the proportions approximately?



17 DTU Informatics, Technical University of Denmark

Representation

- Area represents proportion**
 - Which is smallest, middle, and largest?
 - What are the proportions approximately?



18 DTU Informatics, Technical University of Denmark

Arrangement

- Placement of visual elements**
 - Can make a great difference in how easy it is to interpret data

Example

1	0	1	0	1	1	0
2	1	0	1	0	0	1
3	0	1	0	1	1	0
4	1	0	1	0	0	1
5	0	1	0	1	1	0
6	1	0	1	0	0	1
7	0	1	0	1	1	0
8	1	0	1	0	0	1
9	0	1	0	1	1	0

→

6	1	3	2	5	4
4	1	1	0	0	0
2	1	1	0	0	0
6	1	1	0	0	0
8	1	1	0	0	0
5	0	0	1	1	1
3	0	0	1	1	1
9	0	0	1	1	1
1	0	0	1	1	1
7	0	0	1	1	1

19 DTU Informatics, Technical University of Denmark

Selection

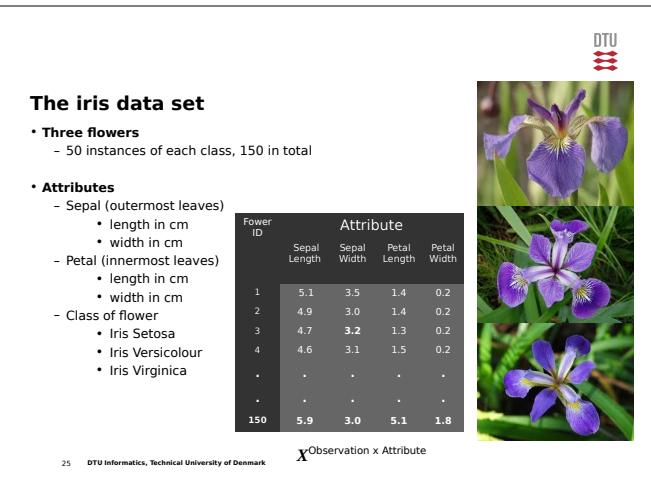
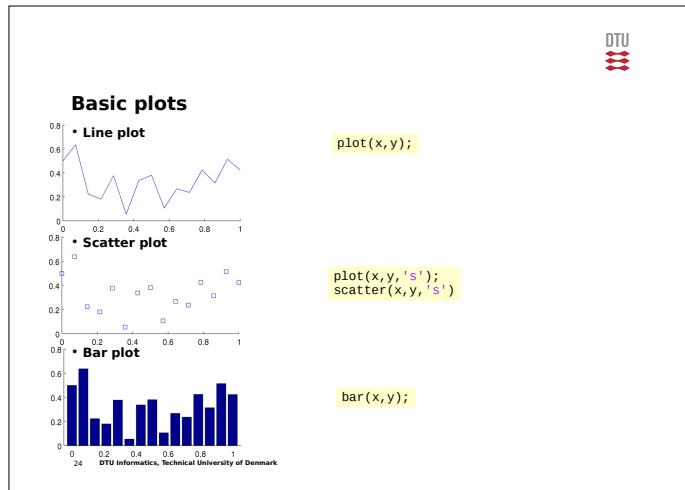
- Elimination or de-emphasis of certain objects or attributes
- A subset of **attributes**
 - Why?** A graph can only show so many attributes - focus on the relevant
 - How?**
 - Dimensionality reduction
 - Plot pairs of attributes
- A subset of **objects**
 - Why?** A graph can only show so many objects - focus on the relevant
 - How?**
 - Random sampling
 - Display of region of interest
 - Use density estimation

20 DTU Informatics, Technical University of Denmark

Types of plots

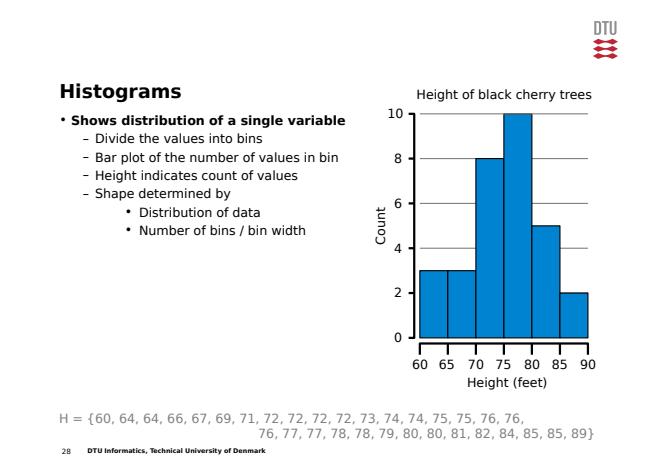
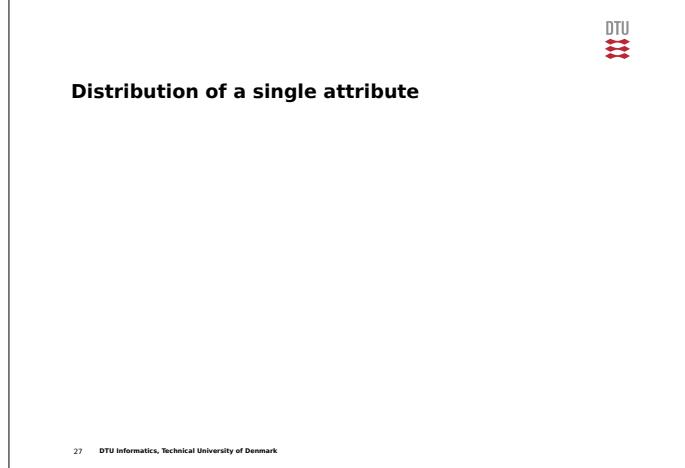
- Distribution of a single attribute**
 - Histogram
 - Empirical cumulative distribution
 - Percentile plots
 - Box plot
- Relation between attributes**
 - 2D histogram
 - Heat maps and contour plots
 - Scatter plots
- Visualization of high-dimensional objects**
 - Matrix plots
 - Parallel coordinates
 - Star plots

21 DTU Informatics, Technical University of Denmark

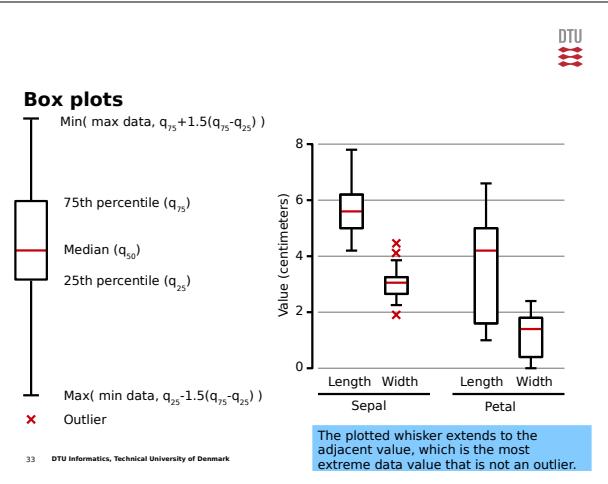
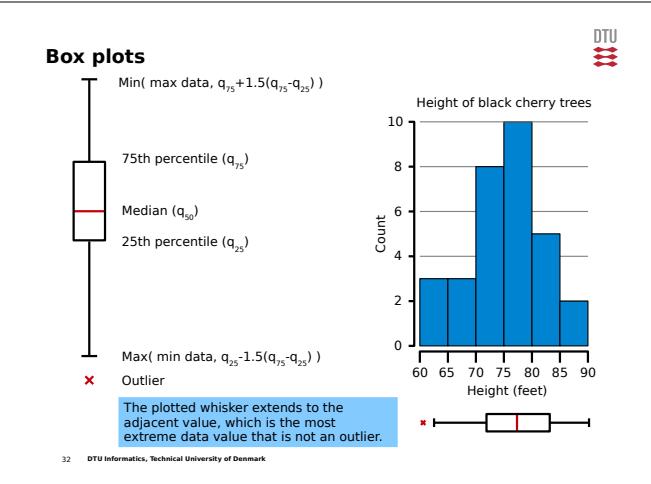
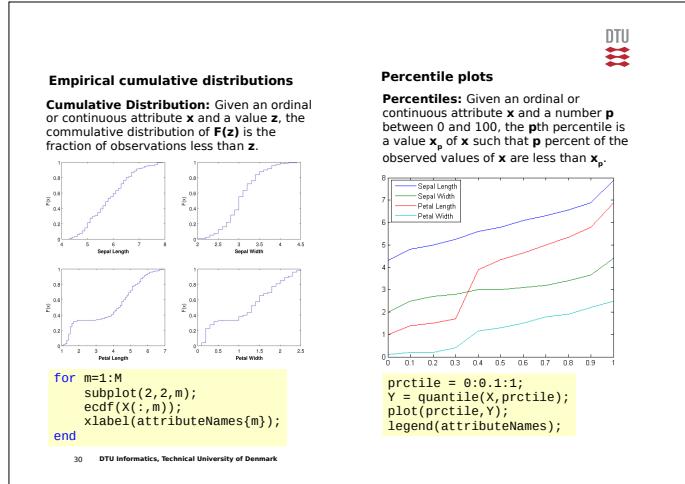
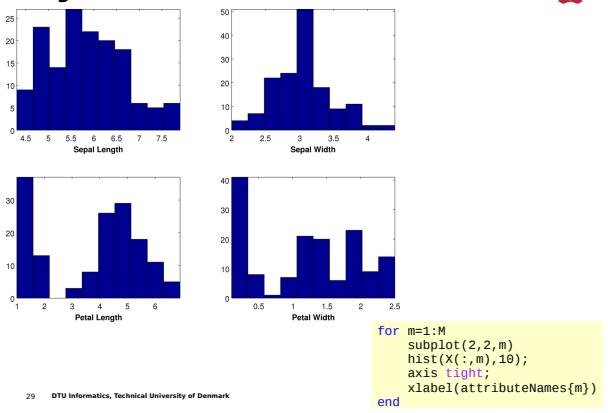


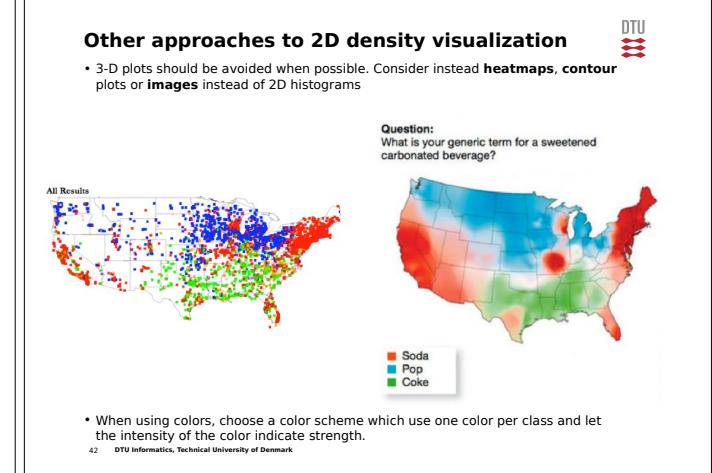
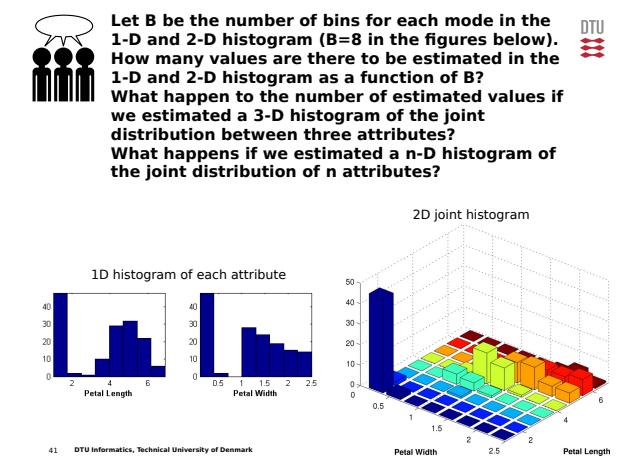
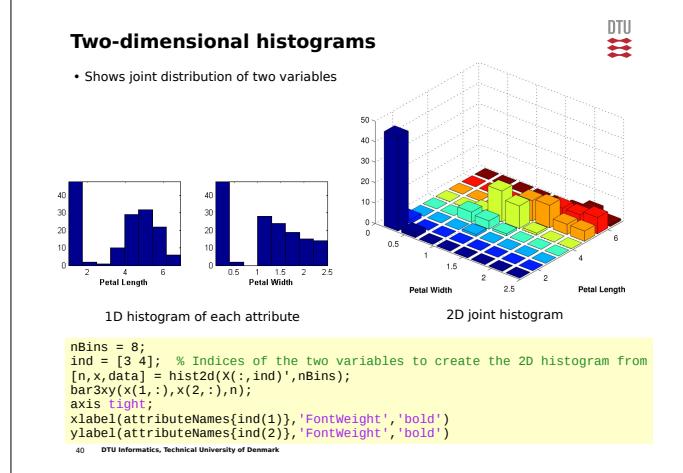
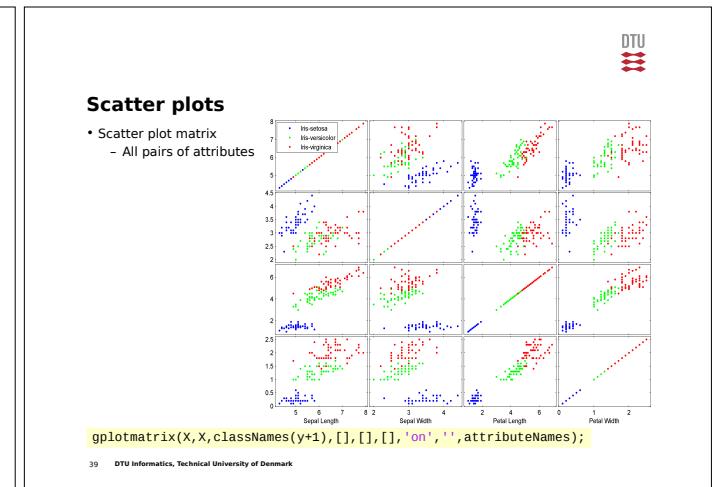
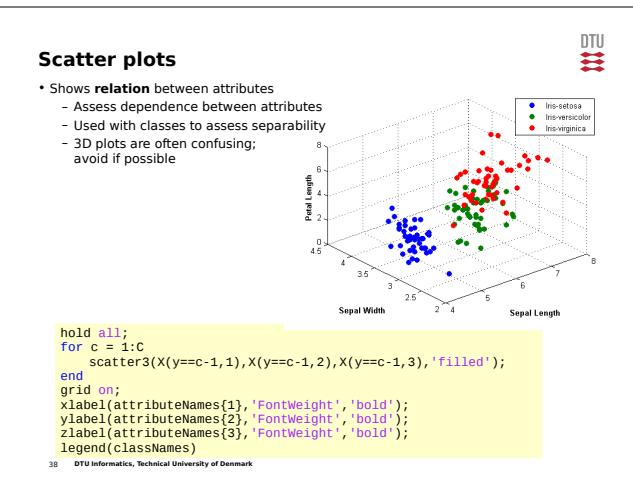
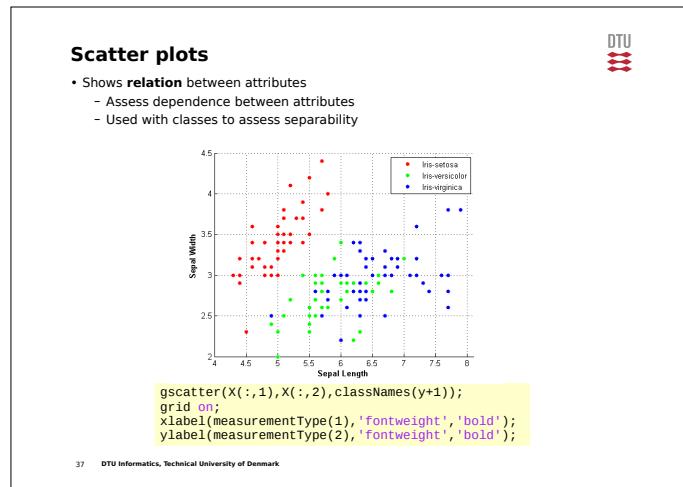
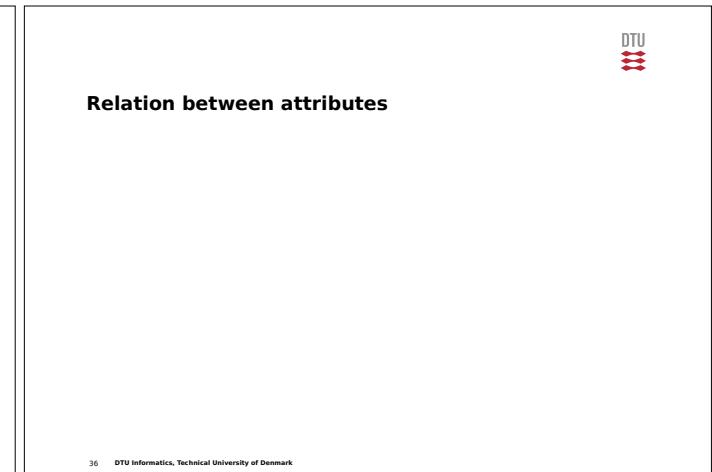
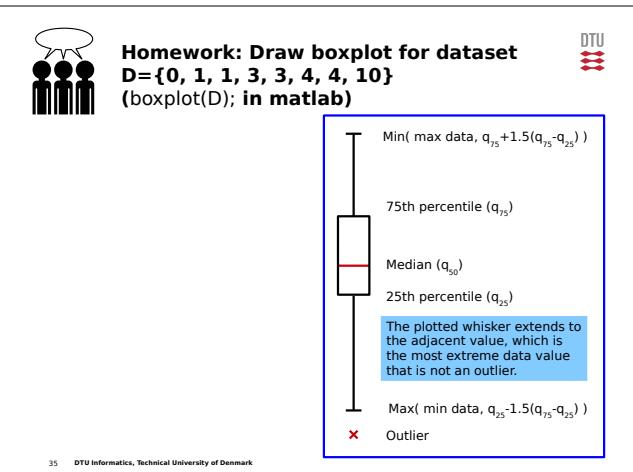
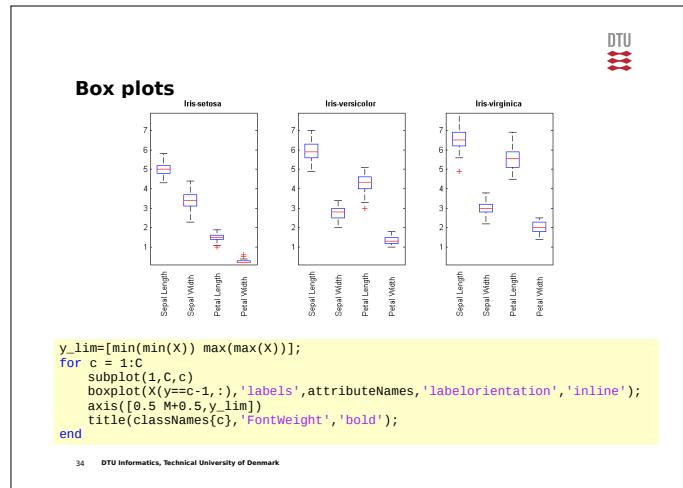
Matlab var.	Type	Size	Description
X	Numeric	$N \times M$	Data matrix: The rows correspond to N data objects, each of which contains M attributes.
y	Numeric	$N \times 1$	Class index: For each data object, y contains a class index, $y \in \{0, 1, \dots, C-1\}$, where C is the total number of classes.
classNames	Cell array	$C \times 1$	Class names: Name (string) for each of the C classes.
attributeNames	Cell array	$M \times 1$	Attribute names: Name (string) for each of the M attributes.
N	Numeric	Scalar	Number of data objects.
M	Numeric	Scalar	Number of attributes.
C	Numeric	Scalar	Number of classes.

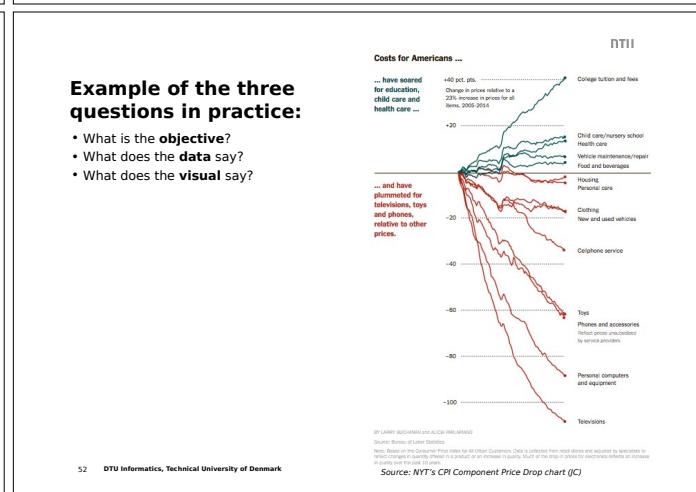
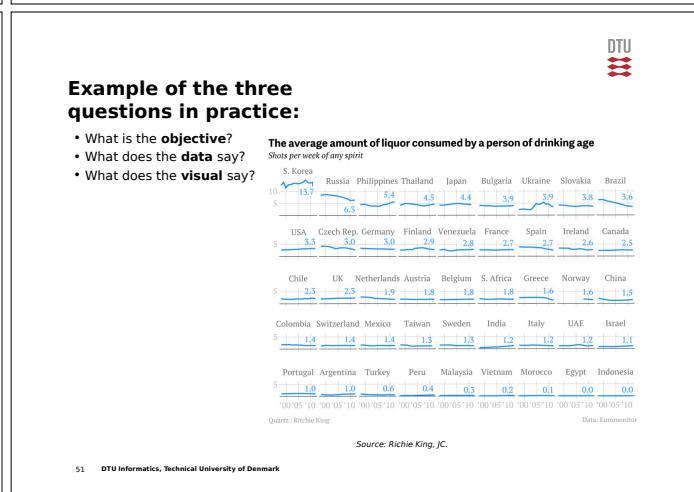
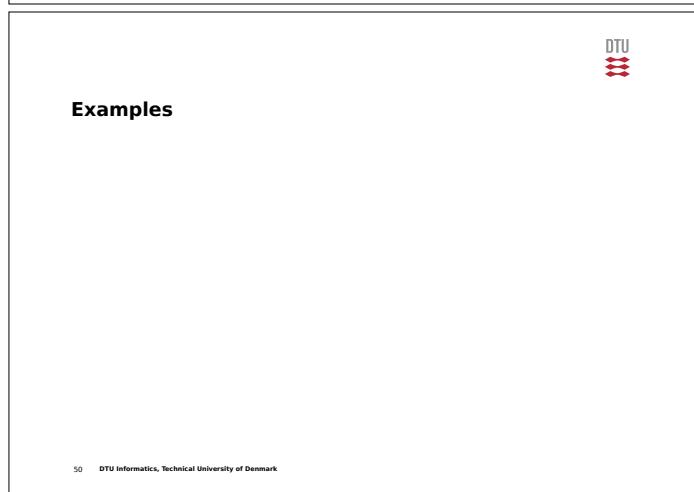
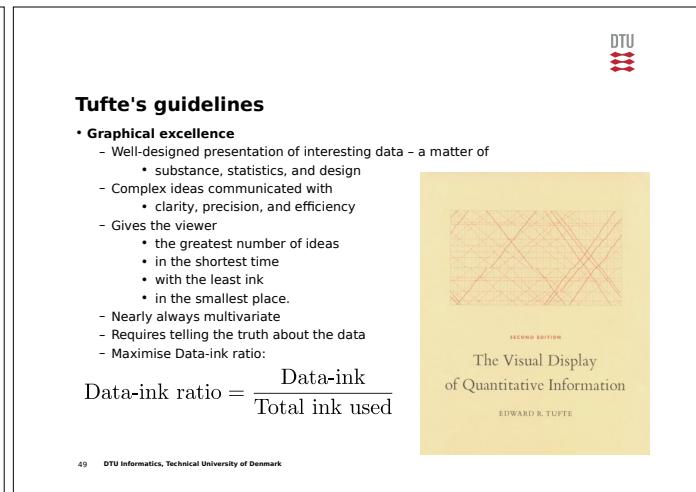
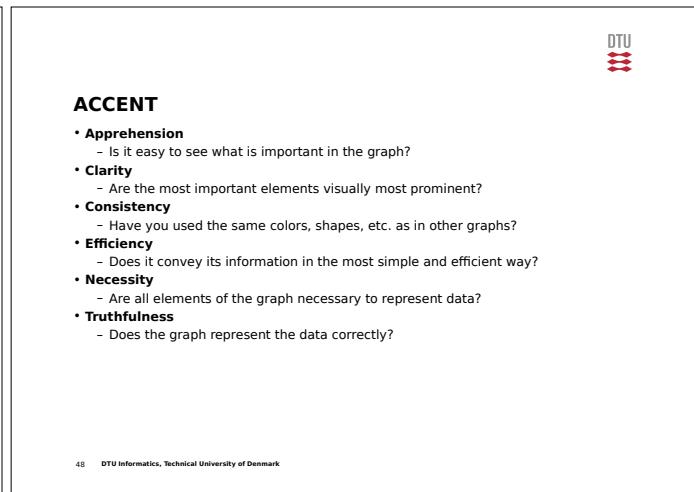
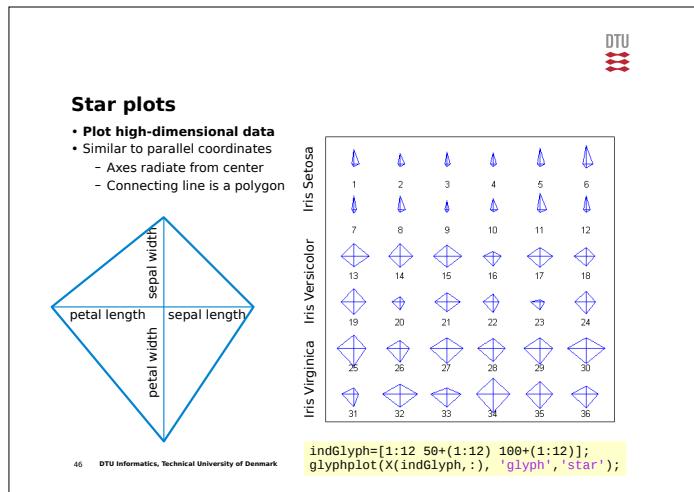
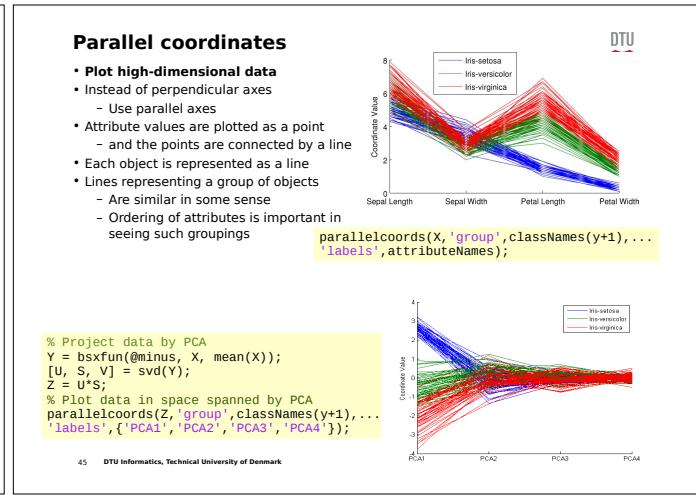
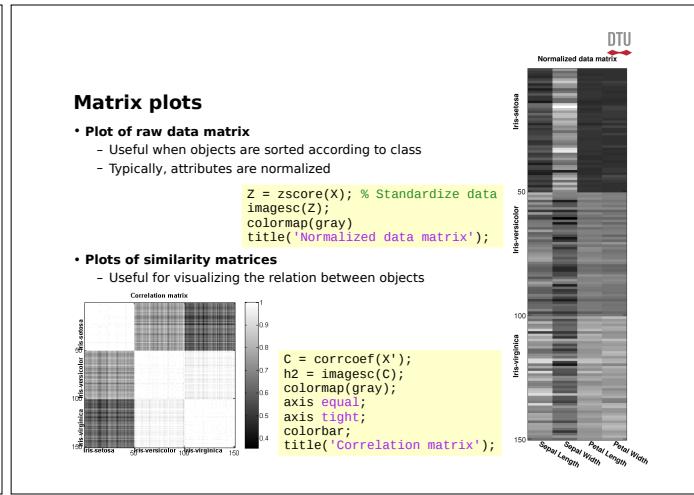
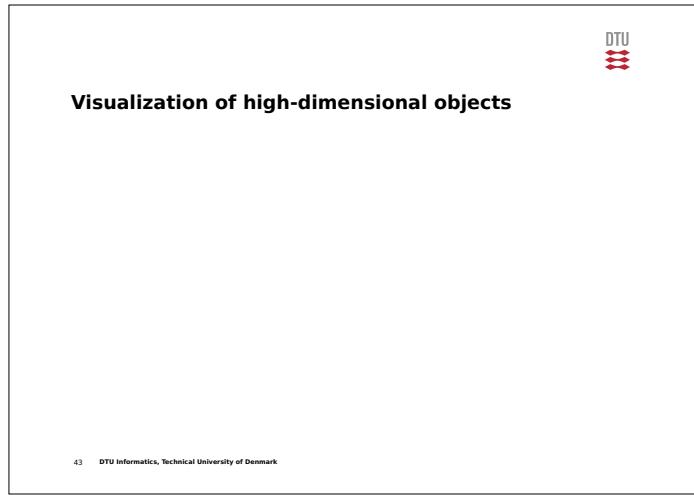
26 DTU Informatics, Technical University of Denmark



Histograms of the Iris data attributes

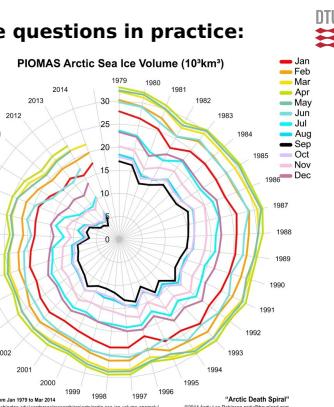






Example of the three questions in practice:

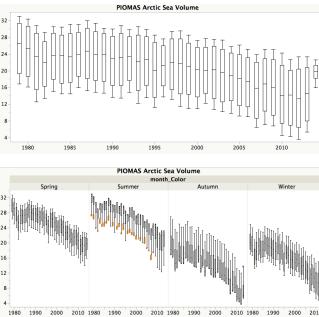
- What is the **objective**?
- What does the **data** say?
- What does the **visual** say?



53 DTU Informatics, Technical University of Denmark

Example of the three questions in practice:

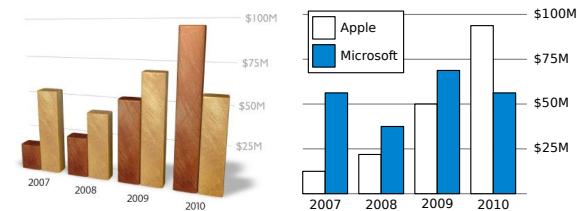
- What is the **objective**?
- What does the **data** say?
- What does the **visual** say?



54 DTU Informatics, Technical University of Denmark

Apprehension: Is it easy to see what is important in the graph?
Clarity: Are the most important elements visually most prominent?
Consistency: Have you used the same colors, shapes, etc. as in other graphs?
Efficiency: Does it convey its information in the most simple and efficient way?
Necessity: Are all elements of the graph necessary to represent data?
Truthfulness: Does the graph represent the data correctly?

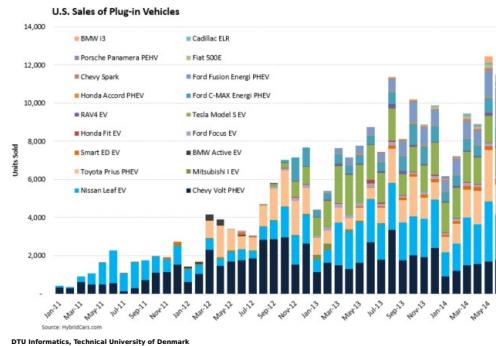
- Compare the graphs using ACCENT



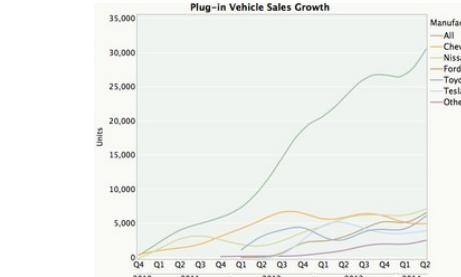
55 DTU Informatics, Technical University of Denmark

Problems with bar charts

(Source: JC)



The problem with bar charts (cont.)

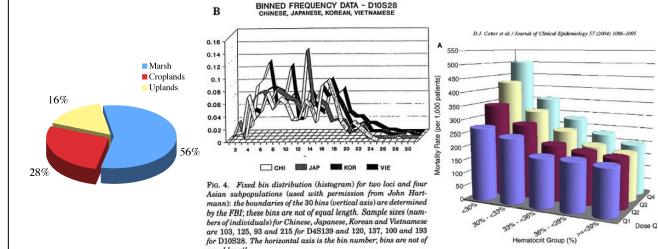


- Bar charts use a lot of space to convey little information. If one wish to convey a lot of information bar charts should be avoided

58 DTU Informatics, Technical University of Denmark

The horroccabinet

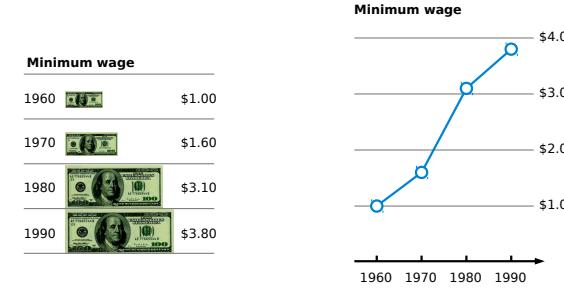
- No matter who you are writing to (even the Pointy-haired boss or 3rd grade students) you should never, ever use the **trivial pursuit**, **spaghetti** or **stairs** plot



59 DTU Informatics, Technical University of Denmark

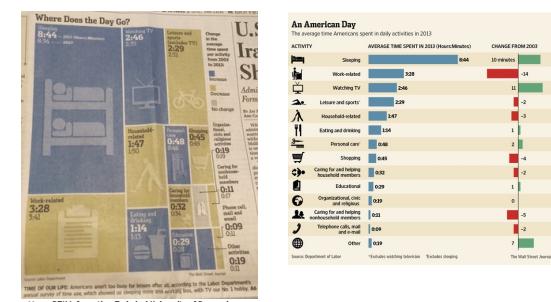
Apprehension: Is it easy to see what is important in the graph?
Clarity: Are the most important elements visually most prominent?
Consistency: Have you used the same colors, shapes, etc. as in other graphs?
Efficiency: Does it convey its information in the most simple and efficient way?
Necessity: Are all elements of the graph necessary to represent data?
Truthfulness: Does the graph represent the data correctly?

- Compare the graphs using ACCENT



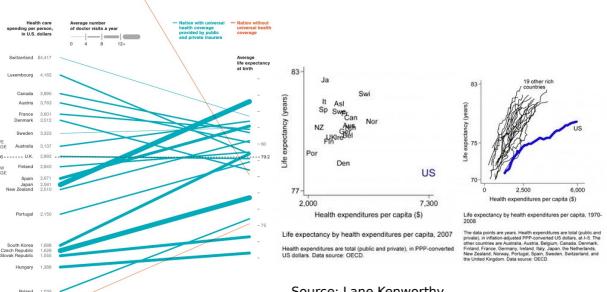
Apprehension: Is it easy to see what is important in the graph?
Clarity: Are the most important elements visually most prominent?
Consistency: Have you used the same colors, shapes, etc. as in other graphs?
Efficiency: Does it convey its information in the most simple and efficient way?
Necessity: Are all elements of the graph necessary to represent data?
Truthfulness: Does the graph represent the data correctly?

- Compare the graphs using ACCENT (Source: WSJ & JC)



Apprehension: Is it easy to see what is important in the graph?
Clarity: Are the most important elements visually most prominent?
Consistency: Have you used the same colors, shapes, etc. as in other graphs?
Efficiency: Does it convey its information in the most simple and efficient way?
Necessity: Are all elements of the graph necessary to represent data?
Truthfulness: Does the graph represent the data correctly?

- Compare these graphs using ACCENT

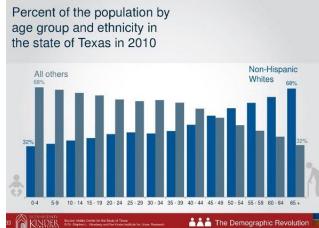


Source: Lane Kenworthy.



Apprehension: Is it easy to see what is important in the graph?
Clarity: Are the most important elements visually most prominent?
Consistency: Have you used the same colors, shapes, etc. as in other graphs?
Efficiency: Does it convey its information in the most simple and efficient way?
Necessity: Are all elements of the graph necessary to represent data?
Truthfulness: Does the graph represent the data correctly?

- Compare these graphs using ACCENT (Source: JC)



65 DTU Informatics, Technical University of Denmark



Imagine you have a dataset where some of the attributes are numeric given in the matrix X but you also have a categorical attribute given by TXT (see below). You would like to carry out a PCA on the data taking both the numeric and categorical attributes into account. How would you proceed?

Age	Height	Weight	Nationality	
-0.2248	0.4762	-0.2097	'Sweden'	0 0 1
-0.5890	0.8620	0.6252	'Sweden'	0 0 1
-0.2938	-1.3610	0.1832	'Sweden'	0 0 1
-0.8940	0.2040	0.0988	'Sweden'	0 0 1
-1.1201	0.8487	0.9492	'Norway'	0 0 1
2.5260	0.3349	0.3071	'Norway'	0 0 1
1.6555	0.5526	0.1352	'Norway'	0 0 1
0.3075	1.0391	0.5152	'Norway'	0 0 1
1.2216	1.1146	0.1414	'Norway'	0 0 1
0.8655	1.2607	-0.9415	'Sweden'	0 0 1
-0.1765	0.6601	-0.1623	'Norway'	0 0 1
0.7914	-0.0671	-0.1461	'Denmark'	1 0 0
-1.3320	-0.1957	-0.5320	'Denmark'	1 0 0
-2.2326	-0.3016	-0.2111	'Sweden'	0 0 1
-1.4491	-0.3031	-0.8757	'Sweden'	0 0 1
0.3335	0.0230	-0.4838	'Sweden'	0 0 1
0.3914	0.0511	-0.7120	'Denmark'	1 0 0
0.4517	0.8261	-1.1742	'Sweden'	0 0 1
-0.1303	1.5270	-0.1922	'Norway'	0 1 0
0.1837	0.4669	-0.2741	'Denmark'	1 0 0

X=[X tmp, attributeNames_tmp]=categorical2numeric(TXT); X=[X X_tmp]; attributeNames=[attributeNames; attributeNames_tmp];

One-out-of-K coding

66 DTU Informatics, Technical University of Denmark



Some guidelines I follow for technical writing

- Everything should be vector graphics at all times (PDF or EPS; the zoom-in test)
- Print out the graphics and ensure it is readable
- Never use a piechart. Never use something worse than a pie chart
- Avoid 3D plots if at all possible (heat-maps or contour plots are preferable)
- Under no circumstances must you add a 3D effect to a plot
- Excel is not for serious work
- Always use a white background. Don't add colors unless it improves the plot
- The best plots are often "`plot(x,y,'.')`". Can the point be made like this?
- Turn off the box around the plot ("`box off`" in matlab)
- Consider the scale/location of axis: "`axis equal`"; if axis are of the same type
- Automate as much as possible (labels, legends, axis). Avoid the Matlab plot editor
- Try to avoid graphics programs. I prefer **Matlab+matlabfrag+Latex+Tikz**
- Use image captions. Add a "*mini-conclusion*" to the caption to tell a reader what he should/can take away from a graphics if it is not apparent
- Ask others about your plot without explaining them at first
- Be aware a typical reader won't read your main text
- **Use common sense**

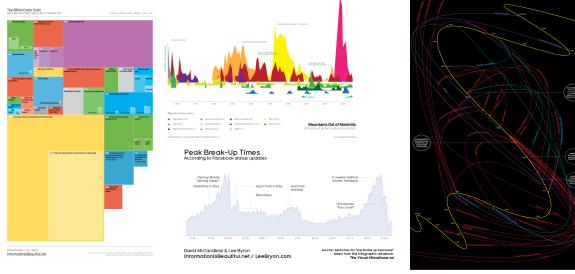
67 DTU Informatics, Technical University of Denmark



Making good data visualizations is an art

For some interesting data visualizations see also

http://www.ted.com/talks/david_mccandless_the_beauty_of_data_visualization.html
<http://www.informationisbeautiful.net/>
<http://www.junkcharts.typepad.com/>



68 DTU Informatics, Technical University of Denmark

02450 Introduction to machine learning and data mining

$$f(x+\Delta x) = \sum_{i=1}^n f_i(x) + \delta e^{i\pi} = \Theta + \Omega \int_a^b \delta e^{ix} dx = [2.7182818284] \sum_{i=1}^n \delta e^{i\pi},$$

DTU Informatics
Department of Informatics and Mathematical Modeling

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"

Section 4.1-4.3 + Appendix D

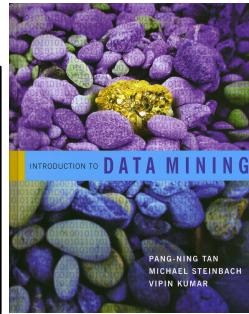
Feedback Groups of the day:
 Anne Brandt Nielsen
 Julie Sofie Bjørn Olsen
 Emilie Lundbye Dalsgaard

 Harri Laine
 Istvan Szonyi

 Astrid Ottosen
 Nanna Elkjaer Møller
 Zenita Omanovic

 Andreas Olsen
 Shah Nazar
 Qian Junging

 If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email to contact me at the exercises next week with feedback/suggestions on the exercises for today.



2 DTU Informatics, Technical University of Denmark

Lecture schedule

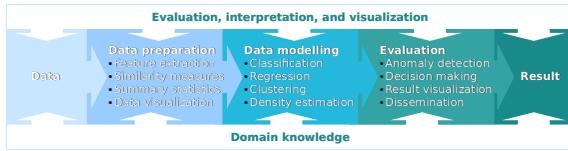
1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.2 + (A) + B.1)
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
(Tan 8.1-8.3 + 8.5-7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project presentation

Report 1

Report 2

3 DTU Informatics, Technical University of Denmark

Data modeling framework

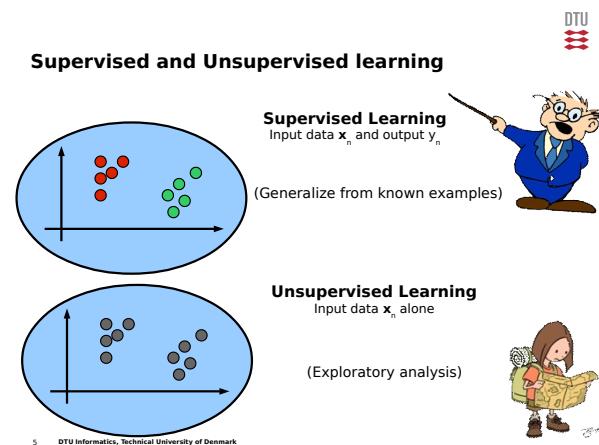


After today you should be able to:

Explain what supervised learning is
 Explain the difference between classification and regression
 Be able to evaluate classifiers in terms of the confusion matrix, error rate and accuracy
 Understand the principles behind decision trees and Hunt's algorithm
 Apply and interpret decision trees, linear regression and logistic regression

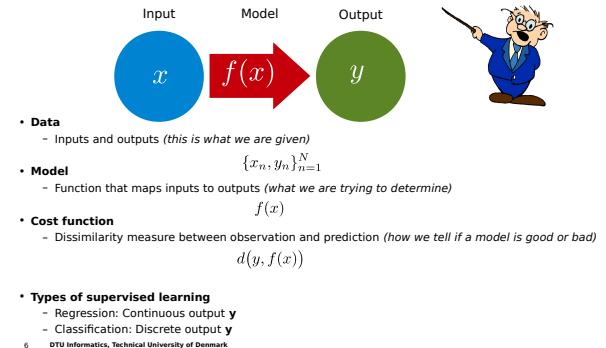
4 DTU Informatics, Technical University of Denmark

Supervised and Unsupervised learning



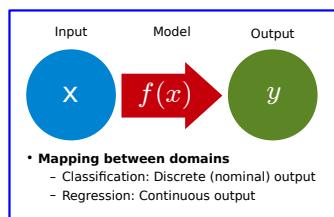
5 DTU Informatics, Technical University of Denmark

Supervised learning



6 DTU Informatics, Technical University of Denmark

 Give an example of a classification and a regression problem and explain what the model $f(x)$ can be used for.



8 DTU Informatics, Technical University of Denmark

Classification

- **Definition:** Learning a function that maps a data object to a discrete class
- **Why classify?**
 - Descriptive modeling
 - Explain / understand the relation between attributes and class
 - Predictive modeling
 - Predict the class of a new data object

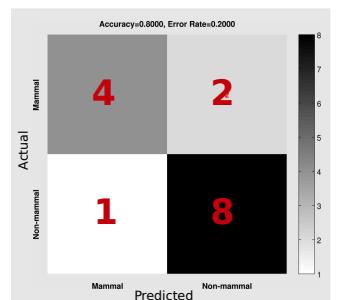
9 DTU Informatics, Technical University of Denmark

Confusion matrix

- Visualization of actual versus predicted class labels
- **Accuracy**
 (Number of correctly predicted observations divided by the total number of observations)

$$\frac{4 + 8}{4 + 2 + 1 + 8} = 80\%$$
- **Error rate**
 (Number of in-correctly predicted observations divided by the total number of observations)

$$\frac{2 + 1}{4 + 2 + 1 + 8} = 20\%$$



10 DTU Informatics, Technical University of Denmark

Decision trees

- Remember the game “20 questions to the professor”? (see also www.20q.net)

Q1. Is it an Animal? Yes.
 Q2. Can you hold it? No.
 Q3. Does it live in groups (gregarious)? Yes.
 Q4. Are there many different sorts of it? No.
 Q5. Can it jump? Yes.
 Q6. Does it eat seeds? No.
 Q7. Is it black and white? Sometimes.
 Q8. Is it black and white? No.
 Q9. Does it have paws? Yes.
 Q10. Can you see it in a zoo? Yes.
 Q11. Does it roar? Yes.
 Q12. Is it worth a lot of money? Yes.
 Q13. Does it have spots? Yes.
 Q14. Is it multicoloured? Yes.
 Q15. Can you make money by selling it? Yes.
 Q16. Does it live in the jungle? Yes.
 Q17. I guessed that it was a leopard? Wrong.
 Q18. Does it like to play? Yes.
 Q19. I guessed that it was a cheetah? Wrong.
 Q20. I am guessing that it is a siberian tiger? Correct.

11 DTU Informatics, Technical University of Denmark

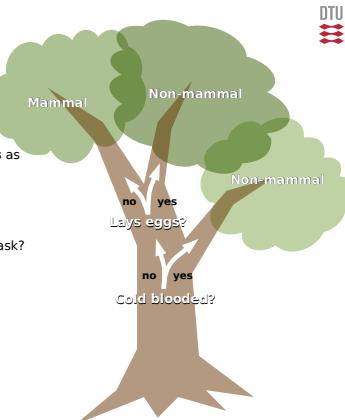
Decision trees

- Ask a series of questions until a conclusion is reached

- Example:** Classify vertebrates as
 - Mammal or
 - Non-mammal

- Learning task**

- Which questions should we ask?



12 DTU Informatics, Technical University of Denmark

Hunts algorithm

- Assign all data objects to the root



13 DTU Informatics, Technical University of Denmark

Hunts algorithm

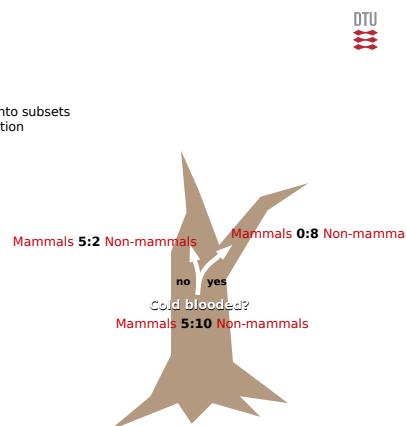
- Select an attribute test condition
- Find a good question to ask



14 DTU Informatics, Technical University of Denmark

Hunt's Algorithm

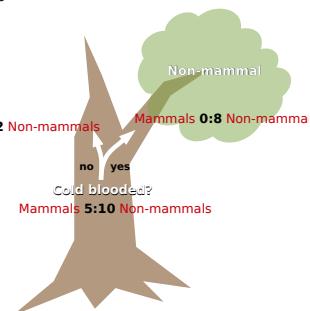
- Partition the data objects into subsets according to the test condition



15 DTU Informatics, Technical University of Denmark

Hunts algorithm

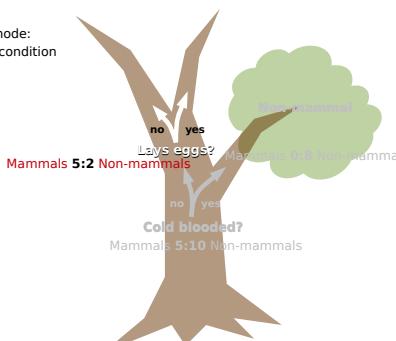
- If all data objects belong to the same class
- Create a leaf node



16 DTU Informatics, Technical University of Denmark

Hunts algorithm

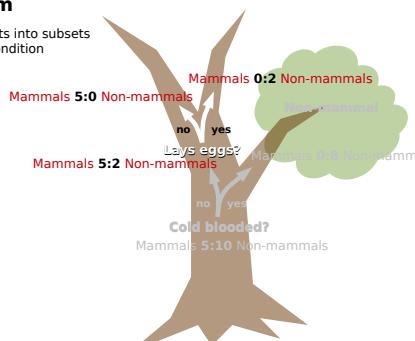
- Repeat for each non-leaf node:
- Select an attribute test condition



17 DTU Informatics, Technical University of Denmark

Hunts algorithm

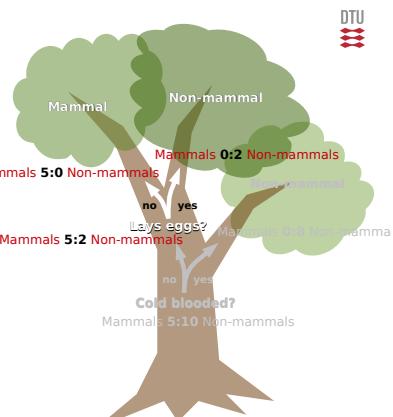
- Partition the data objects into subsets according to the test condition



18 DTU Informatics, Technical University of Denmark

Hunts algorithm

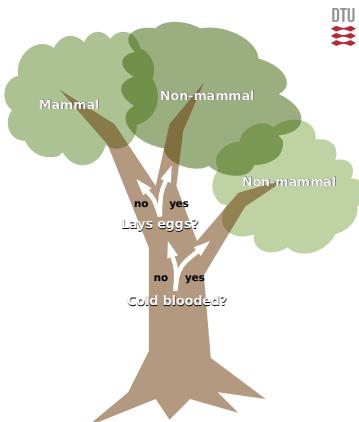
- If all data objects belong to the same class
- Create a leaf node



19 DTU Informatics, Technical University of Denmark

Hunt's algorithm

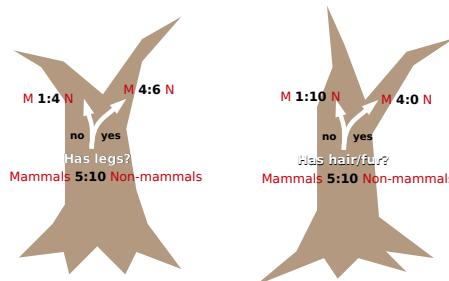
- But how do we find the **best question** at each step?



20 DTU Informatics, Technical University of Denmark

Selecting the best split

- Which of these two questions is best and why?



21 DTU Informatics, Technical University of Denmark

Selecting the best split

- Consider a large number of possible splits
- Compute a measure of impurity after the proposed split
 - For each new branch of the tree
 - Compute weighted average impurity
- Choose split that reduces impurity most

22 DTU Informatics, Technical University of Denmark



Selecting the best split: Impurity measures

- Compute the purity gain, Δ

$$\text{Entropy}(t) = - \sum_{i=0}^{c-1} p(i|t) \log_2 p(i|t)$$

$$\Delta = ?$$

$$\Delta = ?$$

$$\text{Gini}(t) = 1 - \sum_{i=0}^{c-1} (p(i|t))^2$$

$$\Delta = ?$$

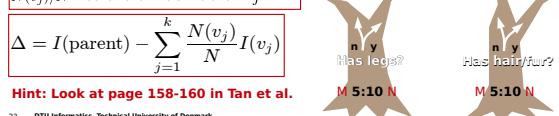
$$\Delta = ?$$

$$\text{Class. error}(t) = 1 - \max_i p(i|t)$$

$$\Delta = ?$$

$$\Delta = ?$$

$p(i|t)$ Fraction of objects that belong to class i
 $N(v_j)/N$ Fraction of animals in branch v_j



Hint: Look at page 158-160 in Tan et al.

23 DTU Informatics, Technical University of Denmark



Selecting the best split: Impurity measures

- Compute the purity gain, Δ

$$\text{Entropy}(t) = - \sum_{i=0}^{c-1} p(i|t) \log_2 p(i|t)$$

$$\Delta(\text{Parent}) = -5/15 \log(5/15) - 10/15 \log(10/15)$$

$$= -1/3 \log(1/3) - 2/3 \log(2/3)$$

$$\text{Left} = -1/3 \log(1/11) - 1/11 \log(1/11)$$

$$\text{Right} = -2/3 \log(4/10) - 4/10 \log(4/10)$$

$$\Delta = 0.9183 \cdot 5/15 + 0.7219 \cdot 10/15 = 0.9710$$

$$-0.0363$$

$$\text{Gini}(t) = 1 - \sum_{i=0}^{c-1} (p(i|t))^2$$

$$\Delta(\text{Parent}) = 1 - (5/15)^2 - (10/15)^2$$

$$= 1 - 1/3 - 2/3$$

$$\text{Left} = 1 - (1/11)^2 - (10/11)^2$$

$$\text{Right} = 1 - (4/10)^2 - (6/10)^2$$

$$\Delta = 0.4444 \cdot 5/15 + 0.3200 \cdot 10/15 = 0.4800$$

$$-0.0177$$

$$\text{Class. error}(t) = 1 - \max_i p(i|t)$$

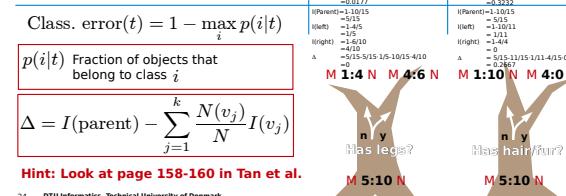
$$\Delta(\text{Parent}) = 1 - (5/15)^2 - (10/15)^2$$

$$= 1 - 1/3 - 2/3$$

$$\text{Left} = 1 - (1/11)^2 - (10/11)^2$$

$$\text{Right} = 1 - (4/10)^2 - (6/10)^2$$

$$\Delta = 0.4444 \cdot 11/15 + 0.1653 \cdot 4/15 = 0.3232$$



Hint: Look at page 158-160 in Tan et al.

24 DTU Informatics, Technical University of Denmark

Which splits to consider

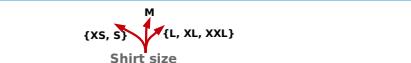
- Binary



- Nominal



- Ordinal

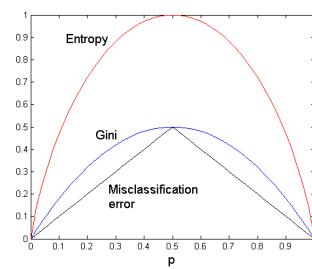


- Continuous



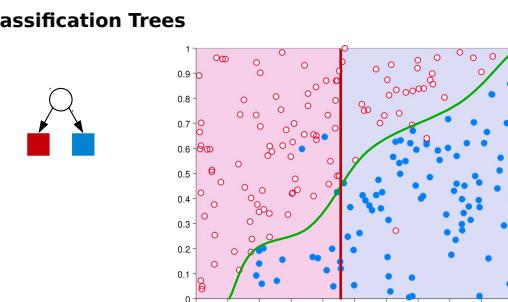
25 DTU Informatics, Technical University of Denmark

For a two class problem



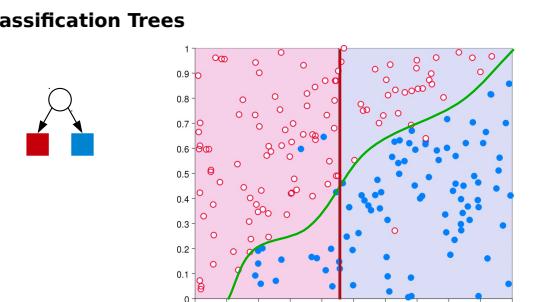
26 DTU Informatics, Technical University of Denmark

Classification Trees

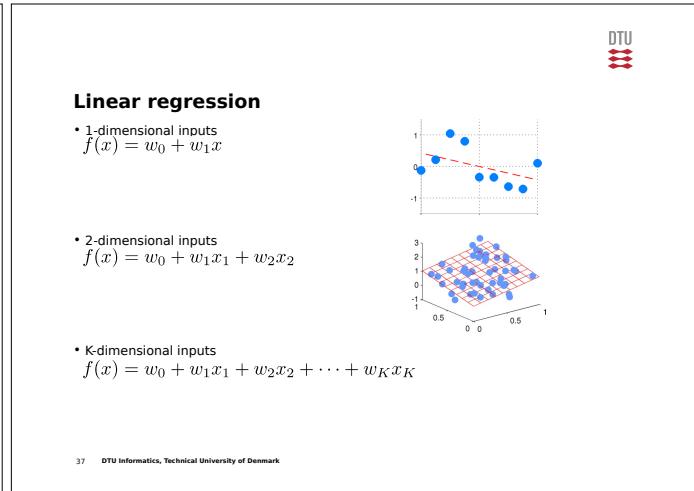
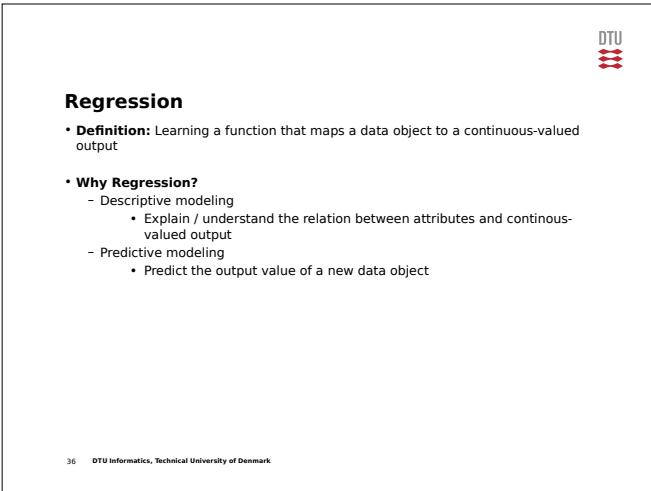
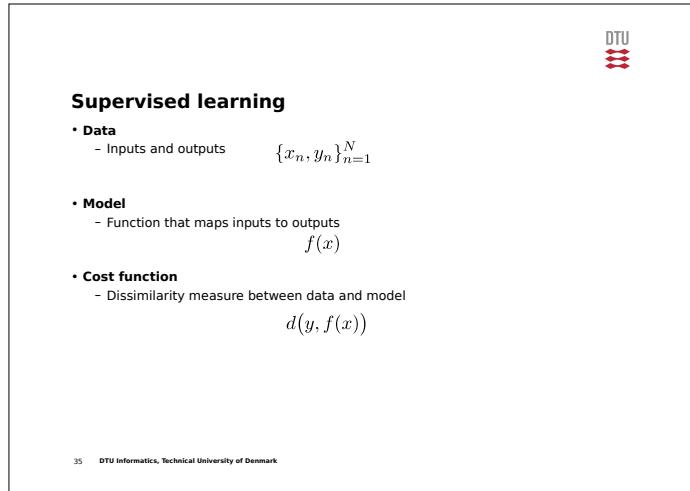
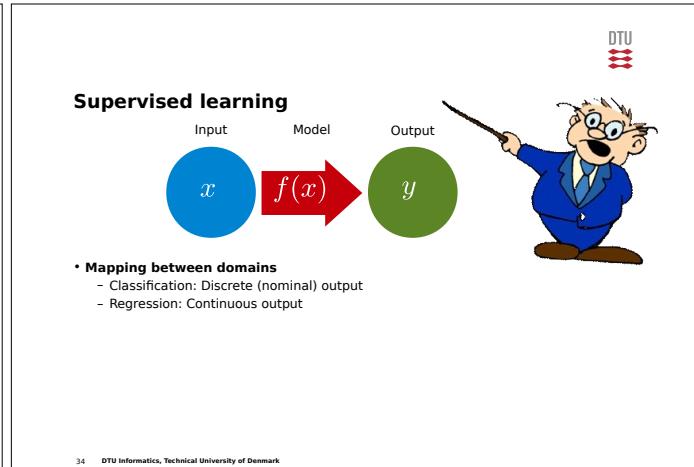
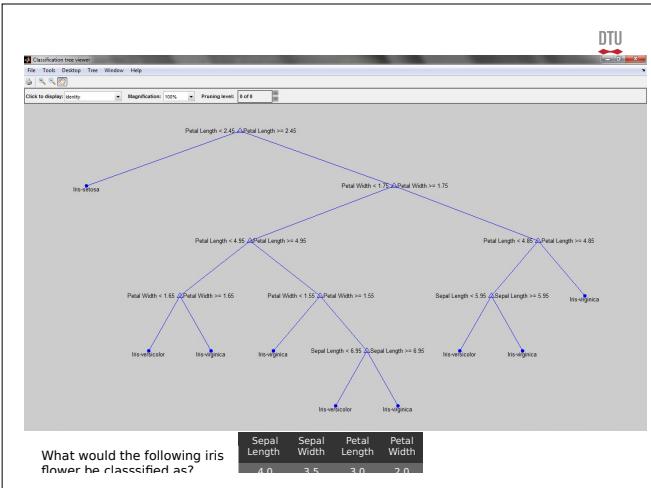
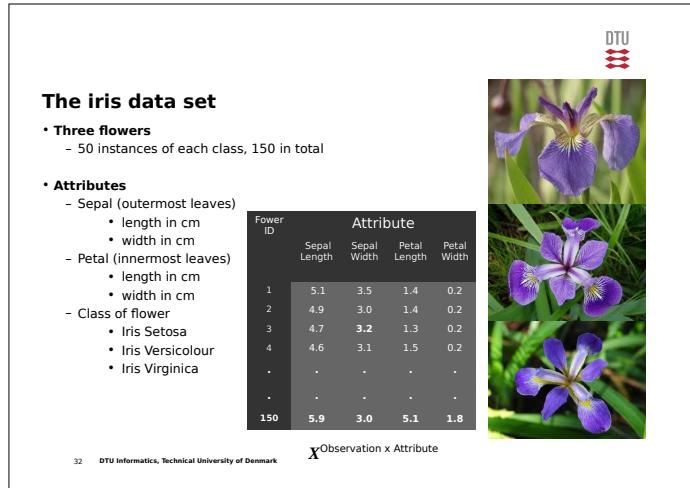
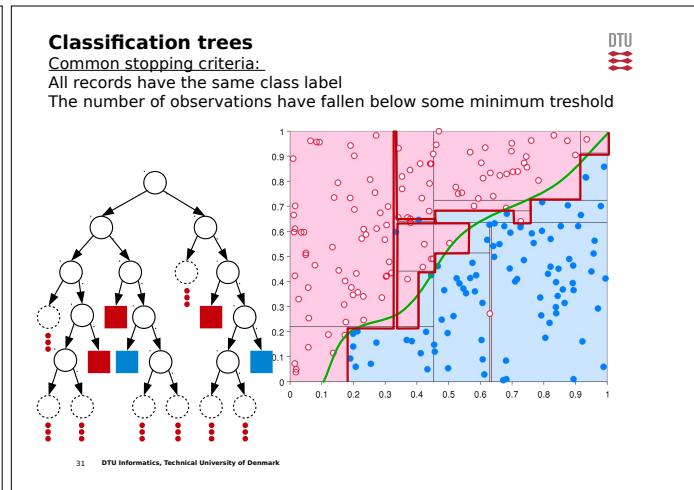
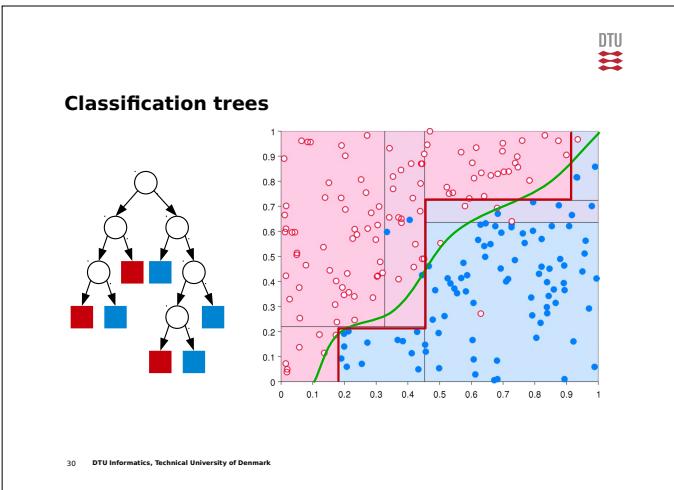
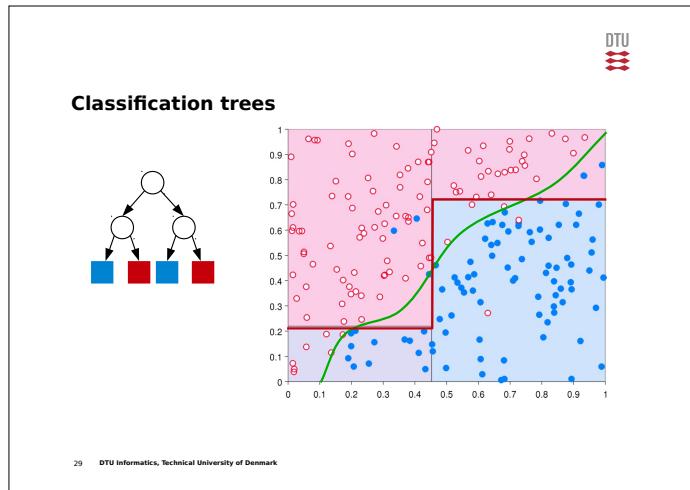


27 DTU Informatics, Technical University of Denmark

Classification Trees



28 DTU Informatics, Technical University of Denmark



Linear regression

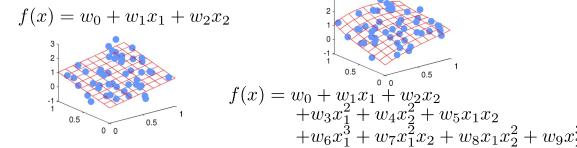
- K-dimensional inputs
 $f(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Kx_K$
- Non-linearly transformed inputs
 $f(x) = w_0 + w_1x + w_2x^2 + \dots + w_Kx^K$
 $f(x) = w_0 + w_1\sin(x) + w_2\cos(x)$

38 DTU Informatics, Technical University of Denmark

Linear regression

- K-dimensional inputs
 $f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Kx_K$
- Non-linearly transformed inputs
 $f(x) = w_0 + w_1x + w_2x^2 + \dots + w_Kx^K$
 $f(x) = w_0 + w_1\sin(x) + w_2\cos(x)$

Example



39 DTU Informatics, Technical University of Denmark

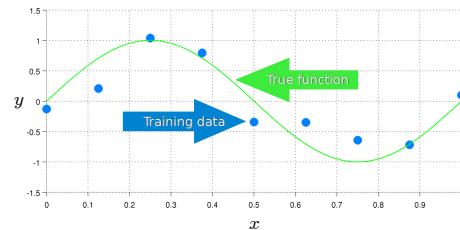
Vector notation

- The linear model can be written compactly using vector notation
 $f(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Kx_K$
 $= \sum_{k=0}^K w_k x_k = \mathbf{x}^\top \mathbf{w}$

- where $x_0 = 1$

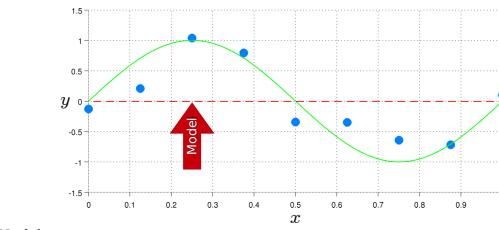
40 DTU Informatics, Technical University of Denmark

Linear regression



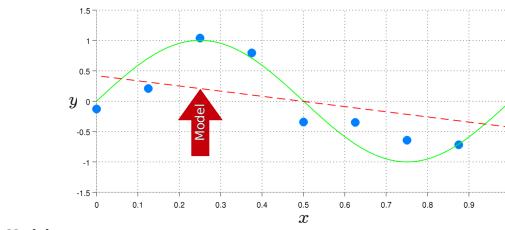
41 DTU Informatics, Technical University of Denmark

Linear regression



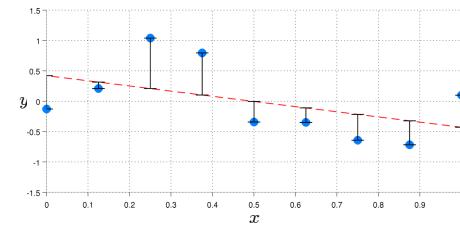
42 DTU Informatics, Technical University of Denmark

Linear regression



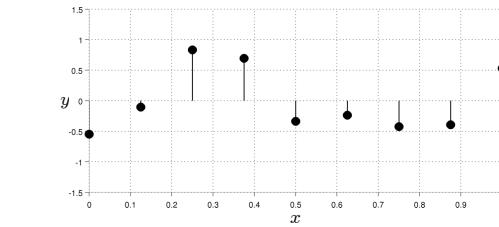
43 DTU Informatics, Technical University of Denmark

Residual error



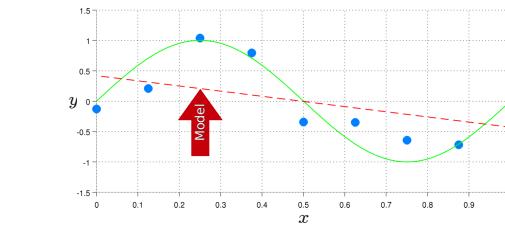
44 DTU Informatics, Technical University of Denmark

Residual error

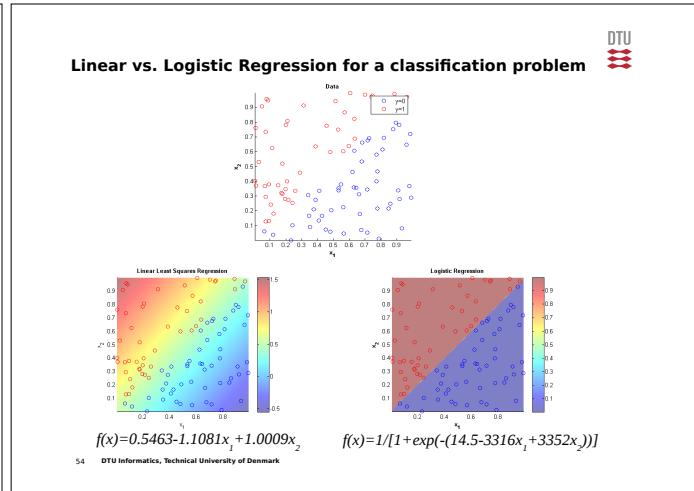
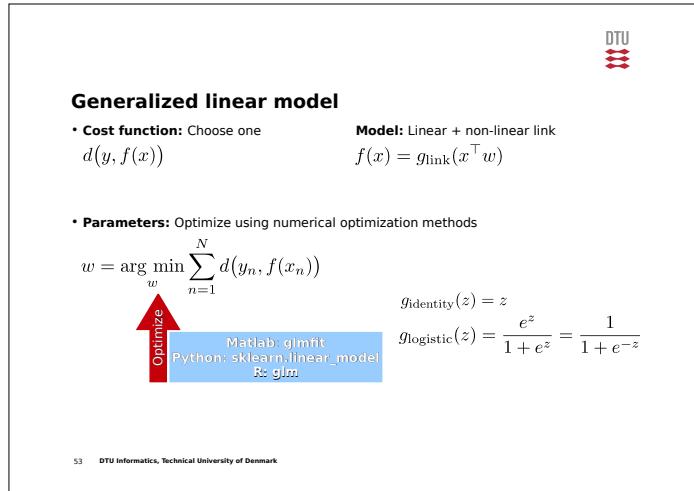
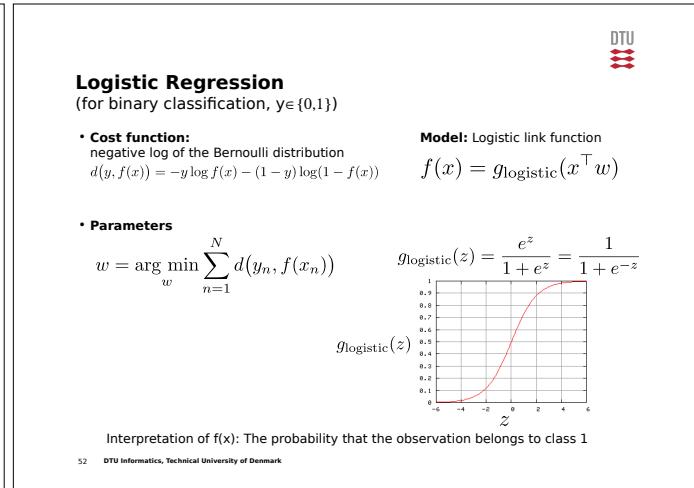
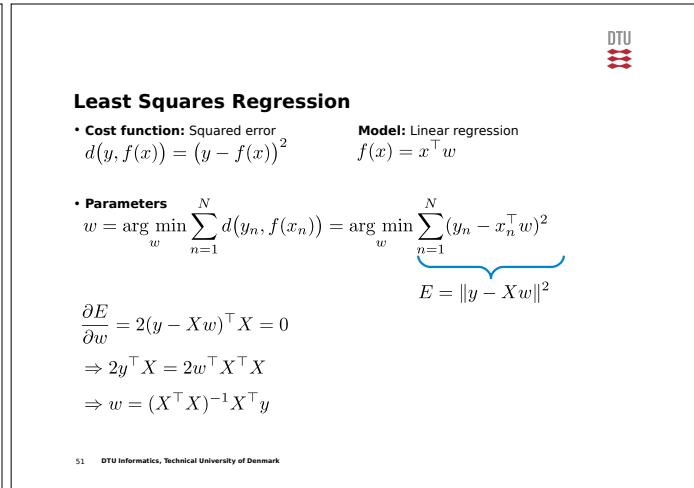
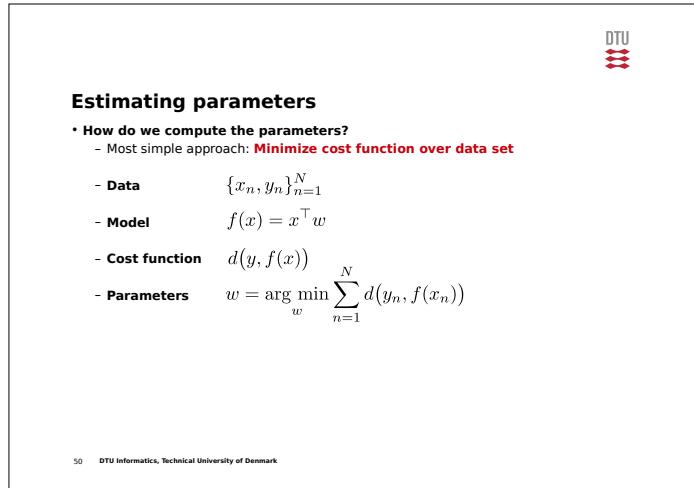
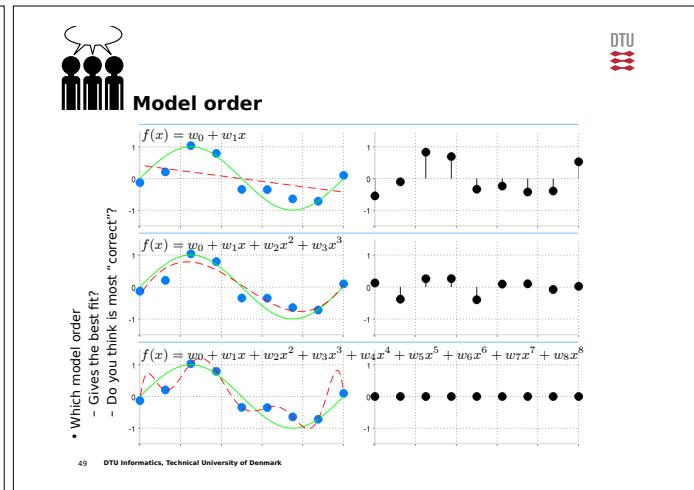
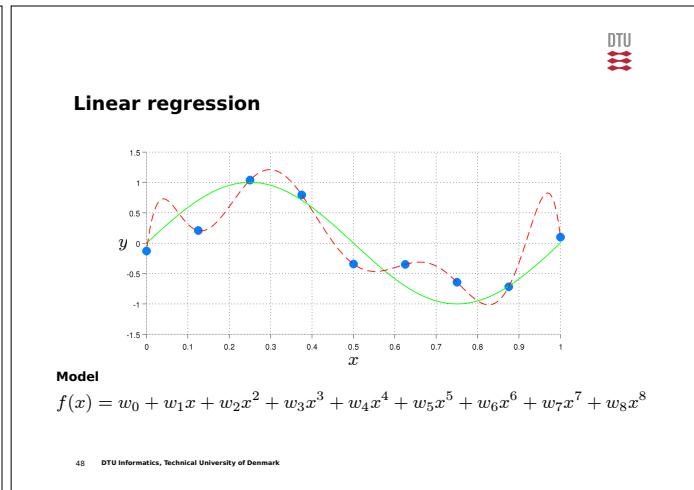
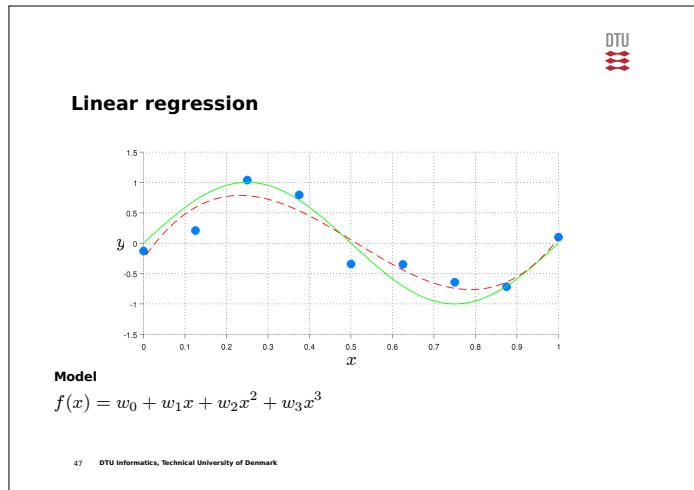


45 DTU Informatics, Technical University of Denmark

Linear regression



46 DTU Informatics, Technical University of Denmark



02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

$f(x+\Delta x) = \sum_i \frac{(\Delta x)^i}{i!} f^{(i)}(x)$

$\Delta \int_a^b \Theta + \Omega \int \delta e^{ix} = [2.7182818284] \sum!$

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"

Section 4.4-4.6

Feedback Groups of the day:
Pietro Davide Doldo, Josephine Robert, Late Vitolina
Jakob Frederiksen, Rasmus Hansen
Lise Oppermann Aagesen, Sophie Nielsen, Isabel Amalia Jespersen
Andrew Marc, Bjarki Sigurjonsson
Rasmus Lau Petersen, Daniel Tolboe Handler
Wen Hsin Li, Robert Lyck, Vivi Halla-aho
Christian Fiksdal Brandes, Anna Emilie Henius

If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

Lecture schedule

1. Introduction (Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction (Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics (Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization (Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression (Tan 4.1-4.6 + D)
6. **Oversampling and performance evaluation** (Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks (Tan 5.2-5.4)
Machine learning and data modelling in practice
8. Ensemble methods and multi class classifiers (Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering (Tan 8.1-8.3+8.5.7)
10. Mixture models and association mining (Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection (Tan 10.1-10.4)
12. Putting it all together: Summary and overview
13. Mini project

3 DTU Informatics, Technical University of Denmark

Data modeling framework

After today you should be able to:

- Explain the difference between training and test (generalization) error
- Explain how cross-validation can be used for model selection
- Apply forward and backward selection
- Prune decision trees
- Understand the Bias-Variance tradeoff as illustrated for regularized least squares estimation
- Test the significance of classifiers

4 DTU Informatics, Technical University of Denmark

Supervised learning

Input x → Model $f(x)$ → Output y

Mapping between domains

- Classification: Discrete output
- Regression: Continuous output

5 DTU Informatics, Technical University of Denmark

Linear regression

- Bad fit
- Too simple model**

$f(x) = w_0 + w_1 x$

6 DTU Informatics, Technical University of Denmark

Linear regression

- Reasonable fit
- Reasonable model**

$f(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$

7 DTU Informatics, Technical University of Denmark

Linear regression

- Perfect fit
- Too complex model**

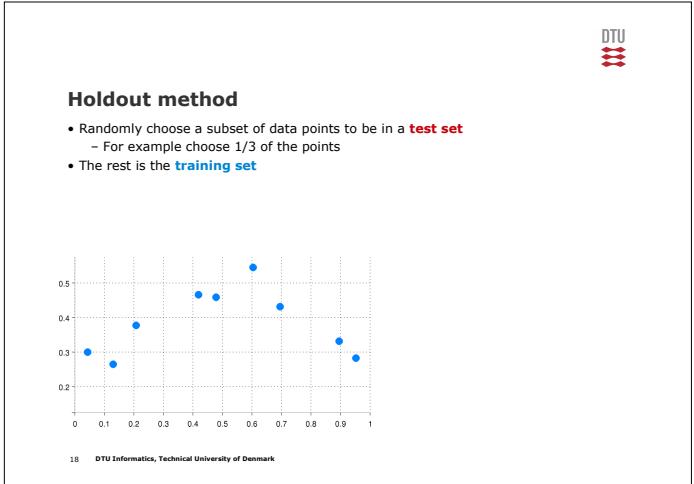
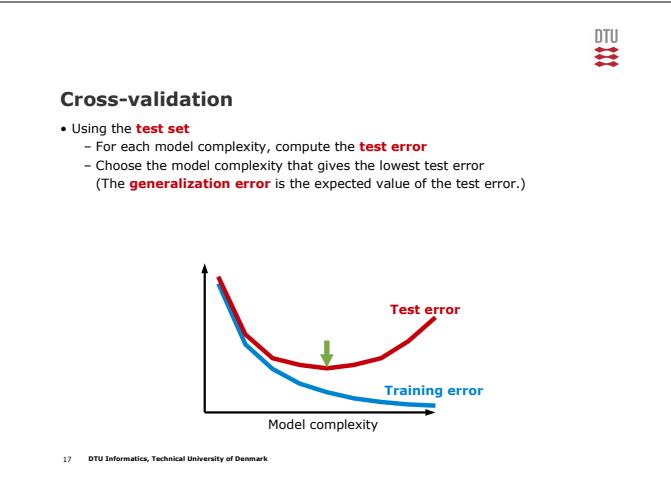
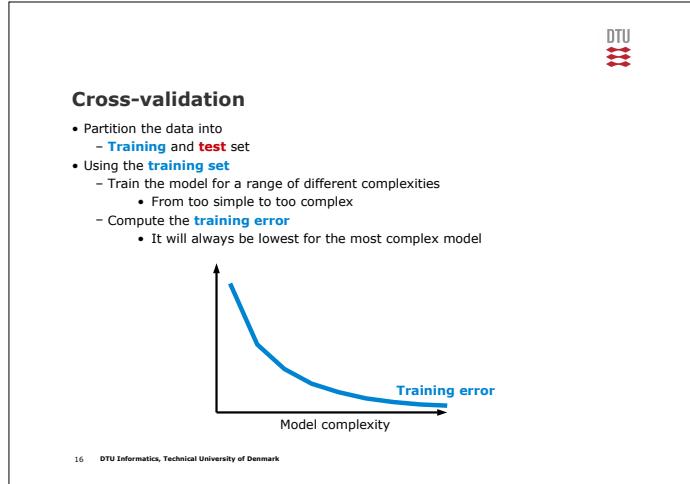
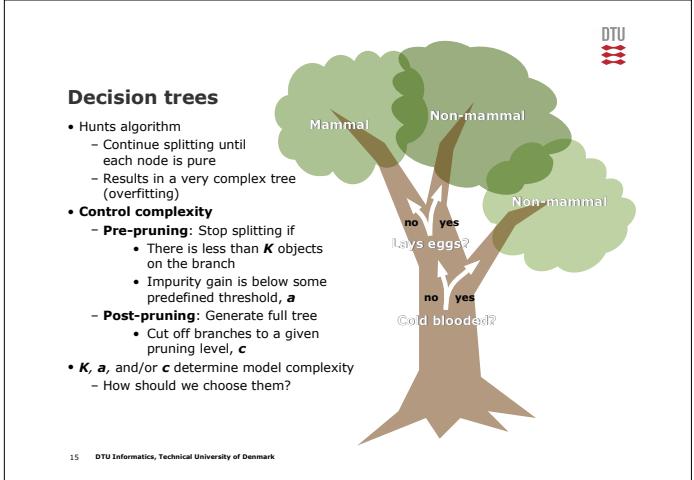
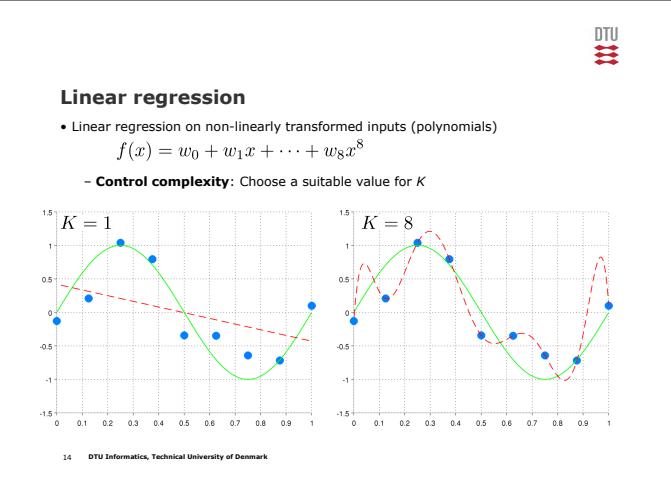
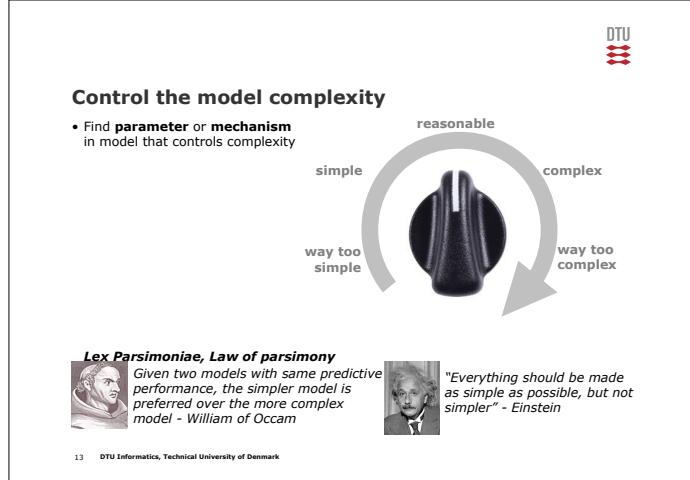
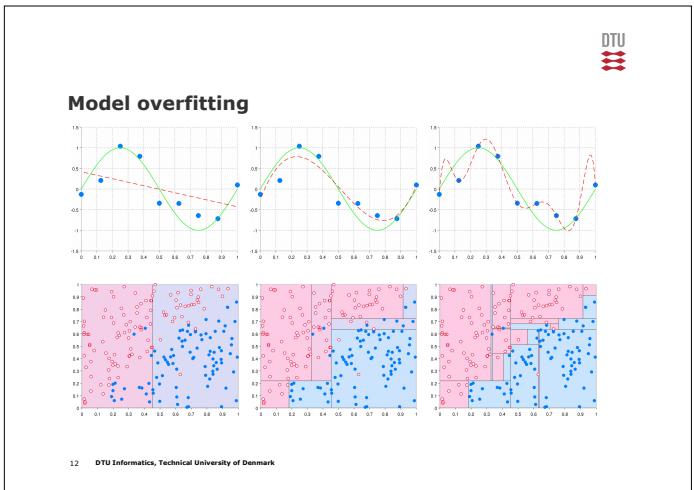
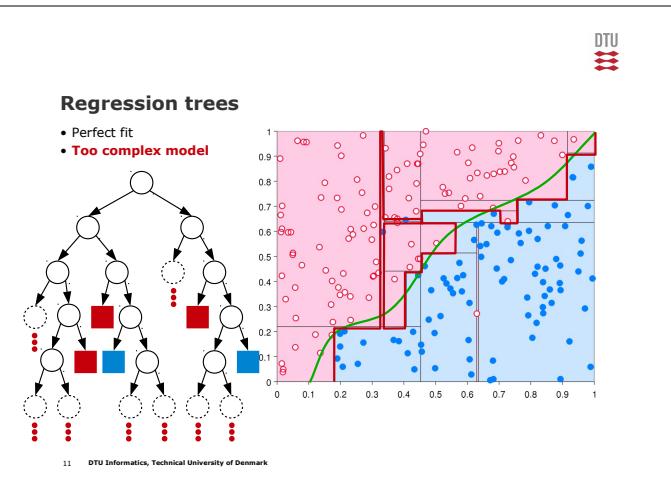
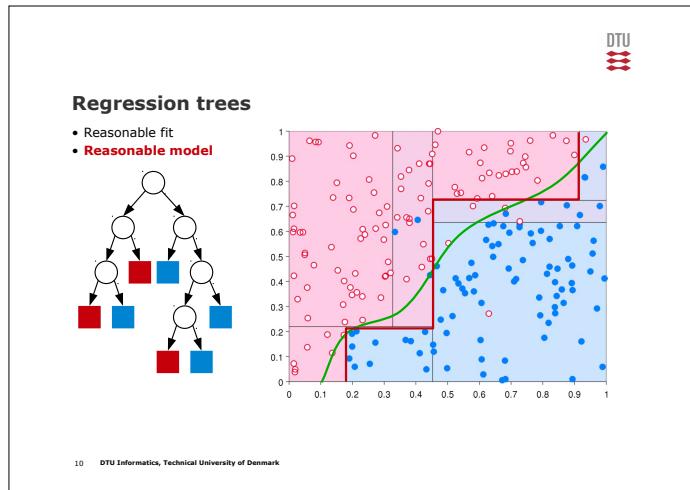
$f(x) = w_0 + w_1 x + \dots + w_8 x^8$

8 DTU Informatics, Technical University of Denmark

Regression trees

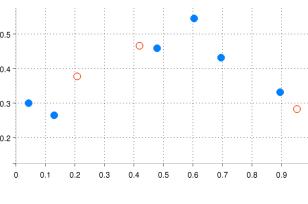
- Bad fit
- Too simple model**

9 DTU Informatics, Technical University of Denmark



Holdout method

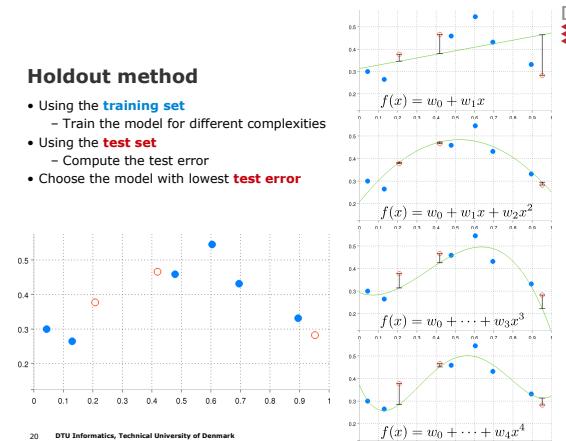
- Randomly choose a subset of data point to be in a **test set**
 - For example choose 1/3 of the points
- The rest is the **training set**



19 DTU Informatics, Technical University of Denmark

Holdout method

- Using the **training set**
 - Train the model for different complexities
- Using the **test set**
 - Compute the test error
- Choose the model with lowest **test error**



20 DTU Informatics, Technical University of Denmark

Holdout method

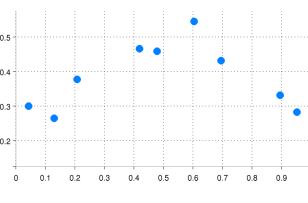
- Using the **training set**
 - Train the model for different complexities
- Using the **test set**
 - Compute the test error
- Choose the model with lowest **test error**



21 DTU Informatics, Technical University of Denmark

Leave-one-out

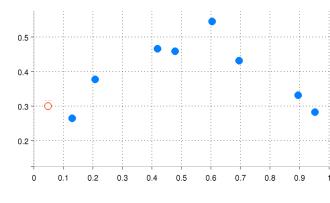
- Choose the first data point as a **test set**
- The rest is the **training set**



22 DTU Informatics, Technical University of Denmark

Leave-one-out

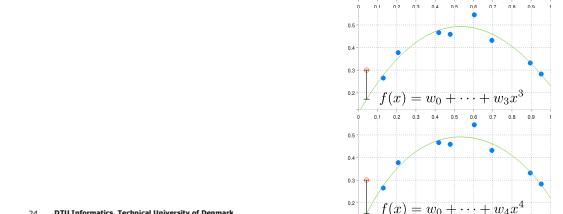
- Choose the first data point as a **test set**
- The rest is the **training set**



23 DTU Informatics, Technical University of Denmark

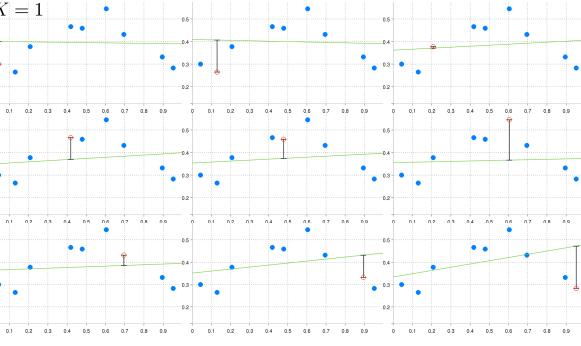
Leave-one-out

- Using the **training set**
 - Train the model for different complexities
- Using the **test set**
 - Compute the test error
- Repeat for all data points**
 - All data points get to be test set
 - Compute **average test error**



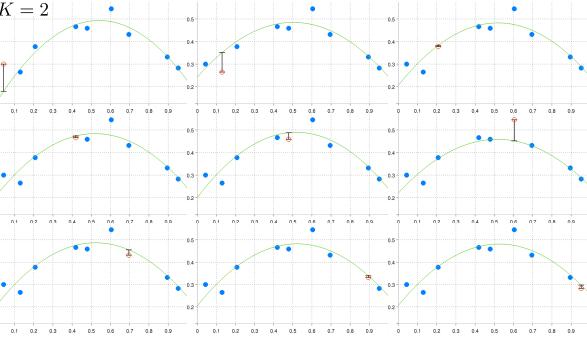
24 DTU Informatics, Technical University of Denmark

Leave-one-out



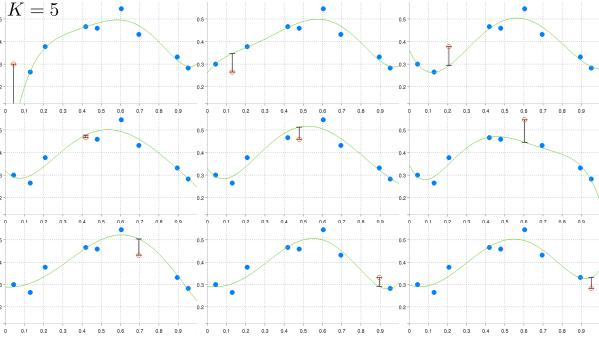
25 DTU Informatics, Technical University of Denmark

Leave-one-out



26 DTU Informatics, Technical University of Denmark

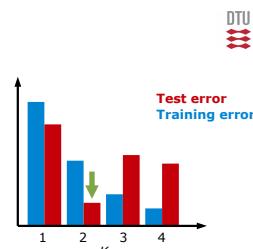
Leave-one-out cross-validation



27 DTU Informatics, Technical University of Denmark

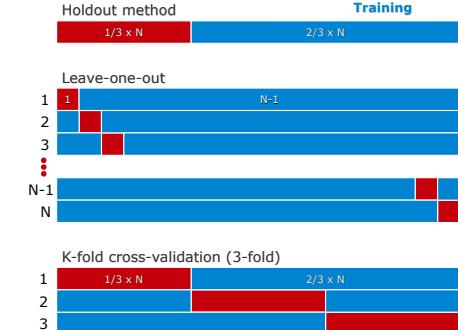
Leave-one-out

- Using the **training set**
 - Train the model for different complexities
- Using the **test set**
 - Compute the test error
- Repeat for all data points**
 - All data points get to be test set
 - Compute **average test error**



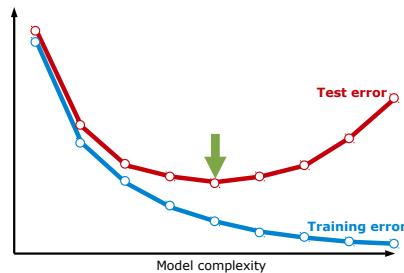
28 DTU Informatics, Technical University of Denmark

Cross-validation methods



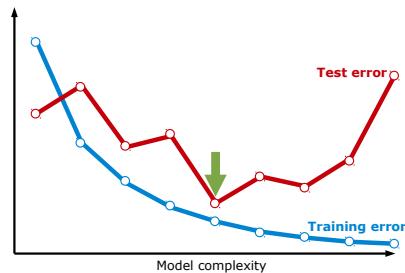
29 DTU Informatics, Technical University of Denmark

Training and test error



31 DTU Informatics, Technical University of Denmark

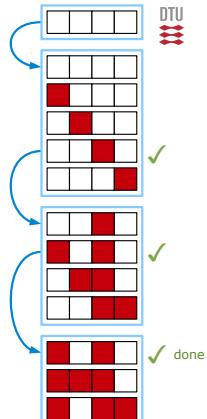
Training and test error



32 DTU Informatics, Technical University of Denmark

Sequential feature selection

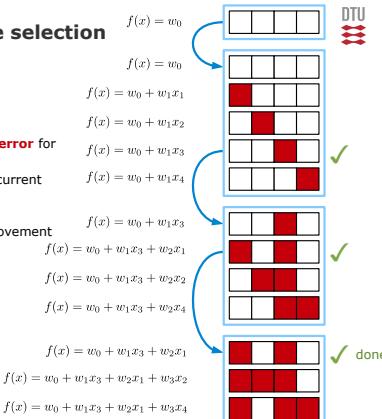
- Forward selection**
 - Start with no features
 - Compute **cross-validation error** for
 - Current feature subset
 - All subsets equal to the current + one added feature
 - Choose best subset
 - Repeat until no further improvement



34 DTU Informatics, Technical University of Denmark

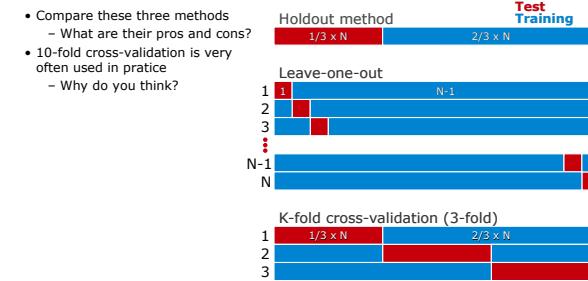
Sequential feature selection

- Forward selection**
 - Start with no features
 - Compute **cross-validation error** for
 - Current feature subset
 - All subsets equal to the current + one added feature
 - Choose best subset
 - Repeat until no further improvement



35 DTU Informatics, Technical University of Denmark

Cross-validation methods



30 DTU Informatics, Technical University of Denmark

Feature subset selection

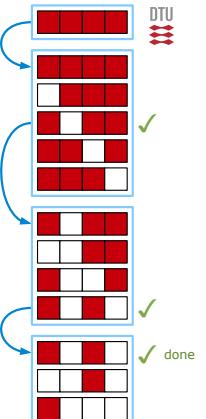
- Let's say we want to do linear regression
 - We have a large number of attributes
- x_1, x_2, \dots, x_M
- Using all attributes results in a too complex model
 - Control complexity:** Choose a subset of attributes
 - Small subset = Simple model
 - Large subset = Complex model
- How many different ways can we choose a subset?**
 - How many models must be compared for
 - M=4
 - M=10
 - M=100

$$\begin{aligned}f(x) &= w_0 \\f(x) &= w_0 + w_1x_1 + w_2x_{27} + w_3x_{88} \\f(x) &= w_0 + w_1x_{19} + w_2x_{76} \\f(x) &= w_0 + w_1x_1 + w_2x_{27} + w_3x_{19} \\f(x) &= w_0 + w_1x_{27} + w_2x_{88}\end{aligned}$$

33 DTU Informatics, Technical University of Denmark

Sequential feature selection

- Backward selection**
 - Start with all features
 - Compute **cross-validation error** for
 - Current feature subset
 - All subsets equal to the current - one removed feature
 - Choose best subset
 - Repeat until no further improvement



36 DTU Informatics, Technical University of Denmark



Feature subset selection

- How many models do we maximally have to evaluate by forward or backward selection? x_1, x_2, \dots, x_M

- M=4
- M=10
- M=100

37 DTU Informatics, Technical University of Denmark

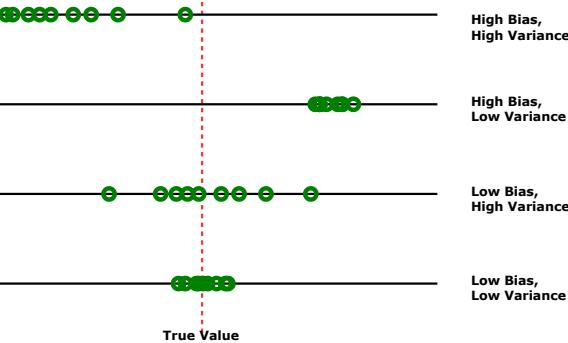


Bias-Variance tradeoff

Consider learning the parameter w for each of N cross-validations, denoted $w^{(i)}$ and consider the squared error between the estimated and true parameter

$$\begin{aligned}\bar{w} &= \frac{1}{N} \sum_{i=1}^N w^{(i)} \\ \mathbb{E}[(w^{\text{true}} - w^{(i)})^2] &= \frac{1}{N} \sum_{i=1}^N (w^{\text{true}} - w^{(i)})^2 \\ &= \frac{1}{N} \sum_{i=1}^N (w^{\text{true}} - \bar{w} + \bar{w} - w^{(i)})^2 \\ &= \frac{1}{N} \sum_{i=1}^N (w^{\text{true}} - \bar{w}) + \frac{1}{N} \sum_{i=1}^N (\bar{w} - w^{(i)}) + \frac{1}{N} \sum_{i=1}^N 2(w^{\text{true}} - \bar{w})(\bar{w} - w^{(i)}) \\ &= \underbrace{\frac{1}{N} \sum_{i=1}^N (w^{\text{true}} - \bar{w})}_{\text{Bias}} + \underbrace{\frac{1}{N} \sum_{i=1}^N (\bar{w} - w^{(i)})}_{\text{Variance}}\end{aligned}$$

38 DTU Informatics, Technical University of Denmark



39 DTU Informatics, Technical University of Denmark

Regularized Least Squares Regression

- Cost function:** Squared error+ridge penalty **Model:** Linear regression

$$E = \sum_{n=1}^N (y_n - f(x_n))^2 + \lambda w^T w \quad f(x) = x^T w$$

- Parameters** $w = \arg \min_w \left[\sum_{n=1}^N (y_n - x_n^T w)^2 + \lambda w^T w \right]$

$$\begin{aligned}\frac{\partial E}{\partial w} &= 2(y - Xw)^T X + 2\lambda w^T = 0 \\ \Rightarrow 2y^T X &= 2w^T(X^T X + \lambda I) \\ \Rightarrow w &= (X^T X + \lambda I)^{-1} X^T y\end{aligned}$$

40 DTU Informatics, Technical University of Denmark



What do you think happens when $\lambda \rightarrow \infty$?
And how can we select the optimal value of λ ?

$$\begin{aligned}w &= \arg \min_w \left[\sum_{n=1}^N (y_n - x_n^T w)^2 + \lambda w^T w \right] \\ \frac{\partial E}{\partial w} &= 2(y - Xw)^T X + 2\lambda w^T = 0 \\ \Rightarrow 2y^T X &= 2w^T(X^T X + \lambda I) \\ \Rightarrow w &= (X^T X + \lambda I)^{-1} X^T y\end{aligned}$$

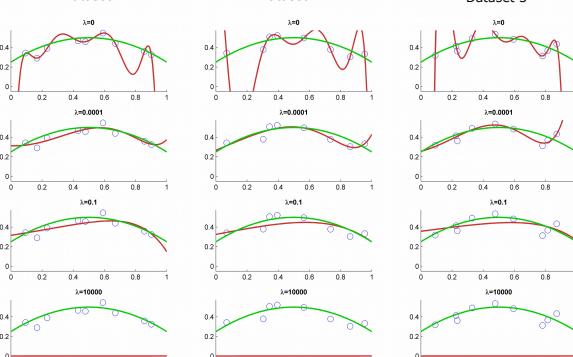
41 DTU Informatics, Technical University of Denmark



Dataset 1

Dataset 2

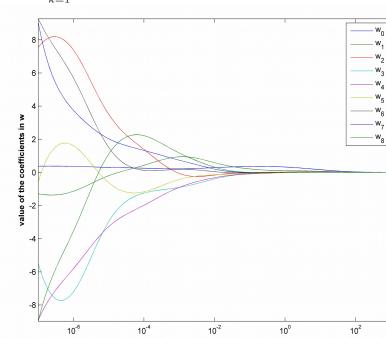
Dataset 3



By regularization we can tradeoff bias and variance, in particular we can hope to substantially reduce variance without introducing too much bias!

42 DTU Informatics, Technical University of Denmark

$$\begin{aligned}f(x) &= w_0 + w_1 x^1 + w_2 x^2 + w_3 x^3 + w_4 x^4 + w_5 x^5 + w_6 x^6 + w_7 x^7 + w_8 x^8 \\ &= \sum_{k=1}^8 w_k x^k\end{aligned}$$

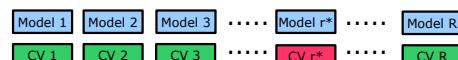


43



Imagine we estimate R different models and select the model r^* with the lowest cross-validation error as the best model. Will the estimated cross-validation error for this model be a correct estimate of how well the model generalizes to new data?

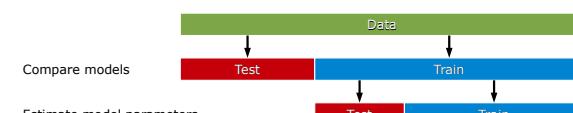
(i.e., is the obtained cross-validation error an unbiased estimator of the generalization error?)



44 DTU Informatics, Technical University of Denmark



Multi-level cross-validation



45 DTU Informatics, Technical University of Denmark

Evaluating the performance of classifiers

- Compare classifier performance to random guessing
- i.e., evaluate how significantly the classifier performs relative to random guessing
- Derive confidence interval for accuracy of classifiers
- i.e., confidence intervals are used to indicate the reliability of an estimate.
- Compare the performance of two classifiers
- i.e., is one classifier significantly better than another classifier.

Evaluating the significance of a classifier

Imagine we train a classifier on a balanced two class problem with N observations and correctly classify L of the observations. We want to know if our classifier is better than random guessing. I.e. We would like to reject the null hypothesis:

H_0 : Our classifier is guessing at random

If we were random guessing the distribution of classifying L observations correctly would follow a binomial with $p=\frac{1}{2}$.

$$p(L|N, p) = \frac{N!}{L!(N-L)!} p^L (1-p)^{N-L}$$

The probability by random of classifying L or more observations correctly by random is then given by

$$\begin{aligned} p(l \geq L|N, p) &= \sum_{l=L}^N \frac{N!}{l!(N-l)!} p^l (1-p)^{N-l} \\ &= 1 - \sum_{l=0}^{L-1} \frac{N!}{l!(N-l)!} p^l (1-p)^{N-l} \end{aligned}$$

Consider two classification problems.

Problem 1: $N=10, L=8$
Problem 2: $N=100, L=60$

$$\begin{aligned} p(1 \geq 8 | 10, 0.5) &= 0.0547 \\ p(1 \geq 60 | 100, 0.5) &= 0.0284 \end{aligned}$$

Confidence interval for the accuracy of a classifier

The empirical accuracy has mean $p=L/N$ and variance $p(1-p)/N$. For large N we can approximate the binomial distributed by the normal distribution to obtain confidence intervals for the accuracy, i.e.

$$\begin{aligned} p\left(Z_{\alpha/2} \leq \frac{\text{acc} - p}{\sqrt{p(1-p)/N}} \leq Z_{1-\alpha/2}\right) &= 1 - \alpha \\ Z_{1-\alpha/2} &= -Z_{\alpha/2} \\ p\left(\frac{Z_{\alpha/2} \leq \frac{\text{acc} - p}{\sqrt{p(1-p)/N}} \leq Z_{1-\alpha/2}}{\sqrt{p(1-p)/N}}\right) &= p\left(\left[\frac{\text{acc} - p}{\sqrt{p(1-p)/N}}\right]^2 \leq Z_{\alpha/2}^2\right) \\ \left[\frac{\text{acc} - p}{\sqrt{p(1-p)/N}}\right]^2 &= Z_{\alpha/2}^2 \\ 0 &= (N + Z_{\alpha/2}^2)p^2 - (2N \text{ acc} + Z_{\alpha/2}^2)p + N \text{ acc}^2 \\ p &\in [\mu - d, \mu + d] \\ \mu &= \frac{2N \text{ acc} + Z_{\alpha/2}^2}{2(N + Z_{\alpha/2}^2)}, \quad d = \frac{Z_{\alpha/2}\sqrt{Z_{\alpha/2}^2 + 4N \text{ acc} - 4N \text{ acc}^2}}{2(N + Z_{\alpha/2}^2)} \end{aligned}$$

Consider a model with accuracy 80%, the confidence interval as a function of number of observations N

N	20	50	100	500	1000	5000
p	[0.584; 0.919]	[0.670; 0.888]	[0.711; 0.867]	[0.763; 0.833]	[0.774; 0.824]	[0.789; 0.811]

Comparing the performance of two classifiers based on k-fold cross-validation using the same splits for the two models

Let e_j and \bar{e}_j be the error rates at cross validation split j for model 1 and model 2 respectively and $d_j = e_j - \bar{e}_j$

$$\begin{aligned} \bar{d} &= \frac{1}{k} \sum_{j=1}^k d_j && \text{Variance of the mean value} \\ \bar{\sigma}^2 &= \frac{\sum_{j=1}^k (d_j - \bar{d})^2}{k(k-1)} \end{aligned}$$

$$d^{\text{cv}} \in [\bar{d} - t_{1-\alpha/2, k-1} \bar{\sigma}, \bar{d} + t_{1-\alpha/2, k-1} \bar{\sigma}]$$

Consider two classifiers with the following error rates across 5-cross validation splits.

Classifier 1: $e_j = [0.2 \ 0.1 \ 0.3 \ 0.2 \ 0.4]$

$$d = e_j - \bar{e}_j = [0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.3]$$

Classifier 2: $e_j = [0.1 \ 0.1 \ 0.2 \ 0.2 \ 0.1]$

$$\bar{d} = 0.1, \quad \sigma^2 = 0.003,$$

$$\begin{aligned} d^{\text{cv}} &\in [0.1 - 2.7764 \cdot \sqrt{0.003}; 0.1 + 2.7764 \cdot \sqrt{0.003}] \\ &= [-0.0521; 0.2521] \end{aligned}$$

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"

Section 5.2-5.4

Feedback Groups of the day:

- Mathias Gæde, Matthias S. Alan Larsen, Frederik Kirk
- Søren Howe Gersager, Ahmet Yıldırım,
- Christopher Frederick Wilm Schenk
- Luke Kristensen, Michael Schmidt Nielsen
- Jamie Neubert Pedersen, Ulf Aslak Jensen, Daniel Heestermann Svendsen
- Thomas Kristensen, Frederik Wedel-Henien
- Karol Dzitkowski, Marco Beccatini
- Anders Rahbek

If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

Lecture schedule

- Introduction (Tan 1.1-1.4)
Data: Feature extraction and visualization
- Data and feature extraction (Tan 2.1-2.3 + BI (+ A))
- Measures of similarity and summary statistics (Tan 2.4 + 3.1-3.2 + C1-C2)
- Data visualization (Tan 3.3)
Supervised learning: Classification and regression
- Decision trees and linear regression (Tan 4.1-4.3 + D)
- Oversampling and performance evaluation (Tan 4.4-4.6)
- Nearest neighbor, naive Bayes, and artificial neural networks (Tan 5.2-5.4)
Machine learning and data modelling in practice
- Putting it all together: Summary and overview
- Mini project

8. Ensemble methods and multi class classifiers (Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.

9. K-means and hierarchical clustering (Tan 8.1-8.3)

10. Mixture models and association mining (Tan 9.2.2 + 6.1-6.3)

11. Density estimation and anomaly detection (Tan 10.1-10.4)
Machine learning and data modelling in practice

12. Putting it all together: Summary and overview

3 DTU Informatics, Technical University of Denmark

Data modeling framework

Evaluation, interpretation, and visualization

Data

- Data preparation
 - Feature extraction
 - Similarity measures
 - Summary statistics
 - Data visualization
- Domain knowledge**

Data modelling

- Classification
 - Regression
 - Clustering
 - Density estimation

Evaluation

- Anomaly detection
- Decision making
- Result visualization
- Dissemination

Result

After today you should be able to:

Explain how K-Nearest Neighbors can be used to classify data
Account for the assumptions made in Naïve Bayes
Apply Bayes theorem to obtain the class posterior likelihood
Understand the principle behind artificial neural networks (ANN) and how ANN can be used for classification and regression.

4 DTU Informatics, Technical University of Denmark

Data modeling framework

Evaluation, interpretation, and visualization

Data

- Data preparation
 - Feature extraction
 - Similarity measures
 - Summary statistics
 - Data visualization
- Domain knowledge**

Data modelling

- Classification
 - Regression
 - Clustering
 - Density estimation

Evaluation

- Anomaly detection
- Decision making
- Result visualization
- Dissemination

Result

Feedback on projects:

To find out which teacher has given feedback to your project, download "Instructions for group project 1.2 and 3/02450_project_1_2014fall_students.xlsx" from campusnet and find the teacher who has graded your project after the lecture.
Please ensure your student ID is on the list to ensure you will get credit.
Otherwise contact the relevant exercise teacher or me.

5 DTU Informatics, Technical University of Denmark

Classify gender based on height and weight

	Height	Weight	Gender
1	183	90	Male
2	180	75	Male
3	170	85	Male
4	185	83	Male
5	159	59	Female
6	167	75	Female
7	165	68	Female
8	175	72	Female
9	171	82	?

6 DTU Informatics, Technical University of Denmark

Nearest neighbor classifier

- 1 nearest neighbor

7 DTU Informatics, Technical University of Denmark

Nearest neighbor classifier

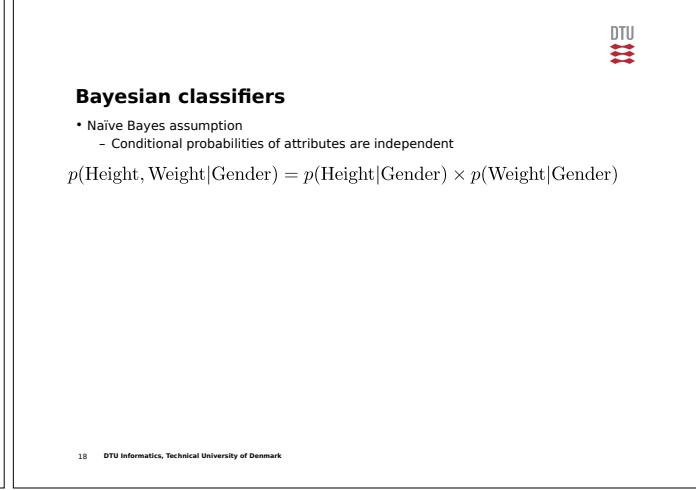
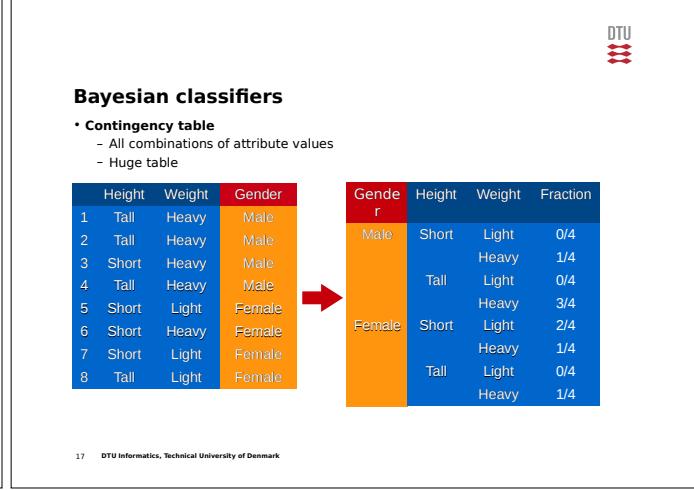
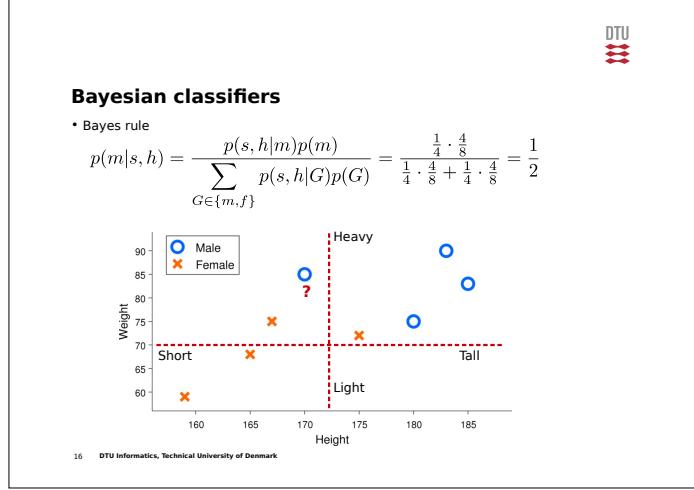
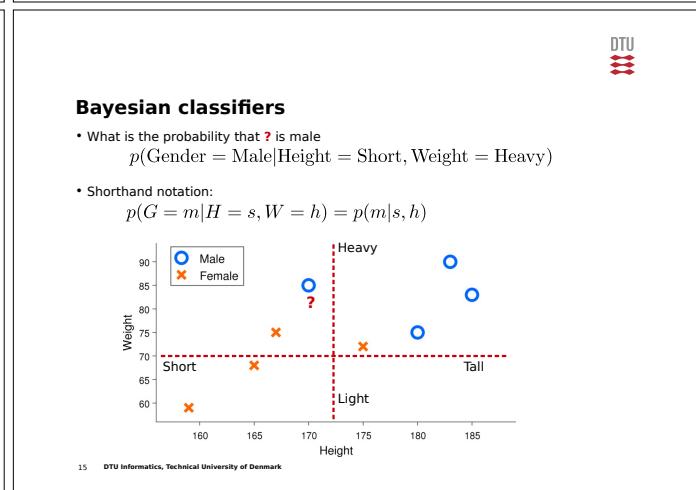
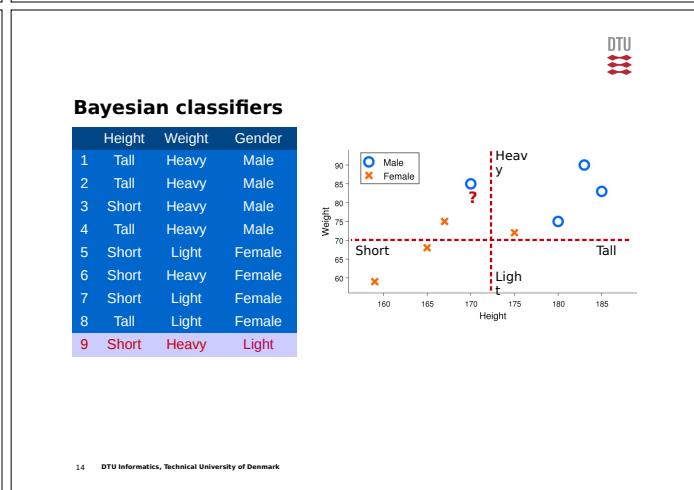
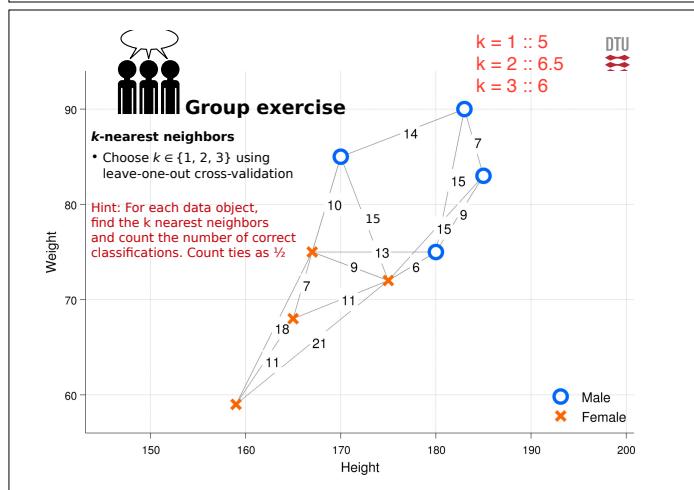
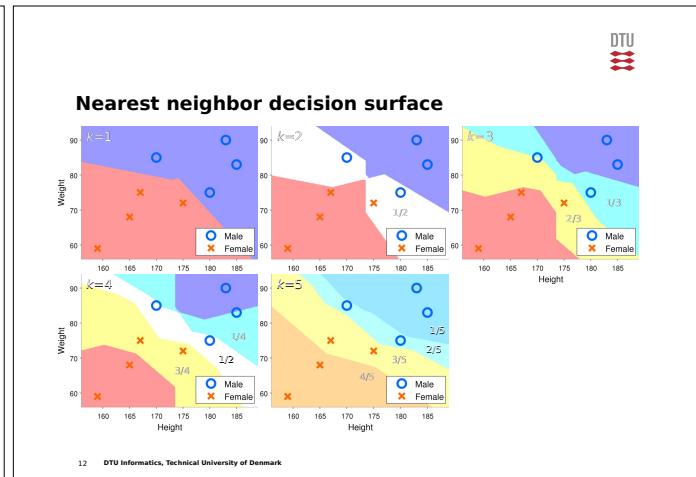
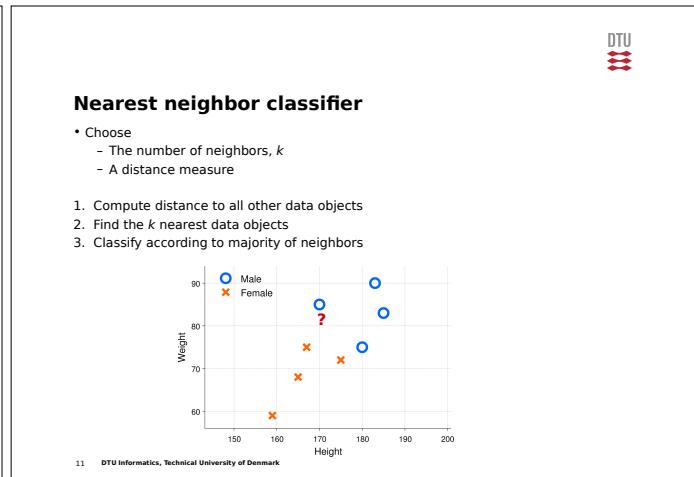
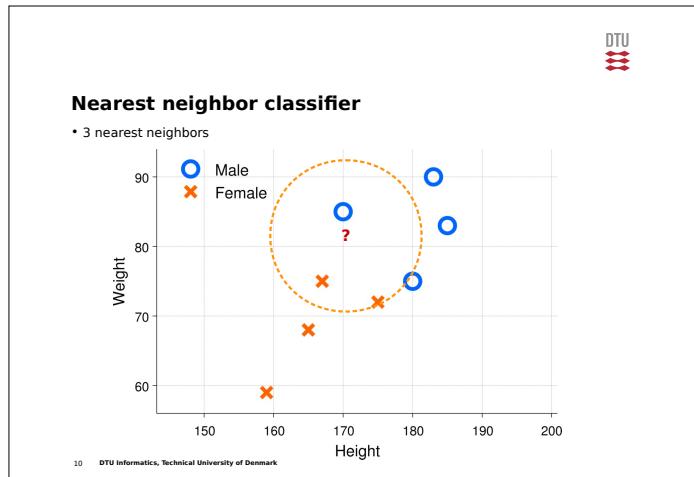
- 1 nearest neighbor

8 DTU Informatics, Technical University of Denmark

Nearest neighbor classifier

- 2 nearest neighbors

9 DTU Informatics, Technical University of Denmark

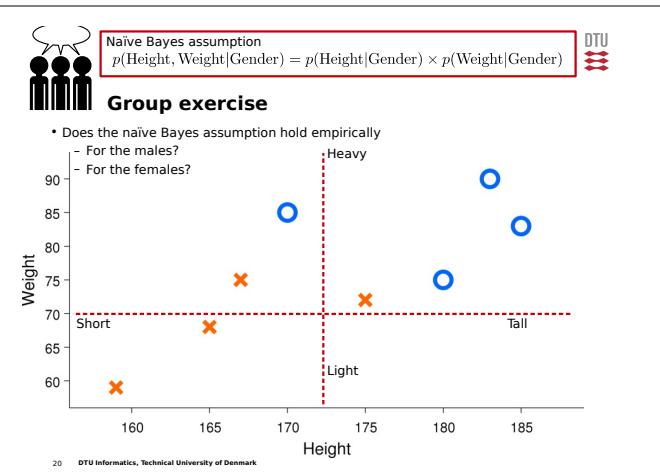


Bayesian classifiers

- Naive Bayes classifier**

$$p(m|s, h) = \frac{p(s|m)p(h|m)p(m)}{\sum_{G \in \{m,f\}} p(s|G)p(h|G)p(G)} = \frac{\frac{1}{4} \cdot \frac{4}{4} \cdot \frac{4}{8}}{\frac{1}{4} \cdot \frac{4}{4} \cdot \frac{4}{8} + \frac{3}{4} \cdot \frac{2}{4} \cdot \frac{4}{8}} = \frac{2}{5}$$

19 DTU Informatics, Technical University of Denmark



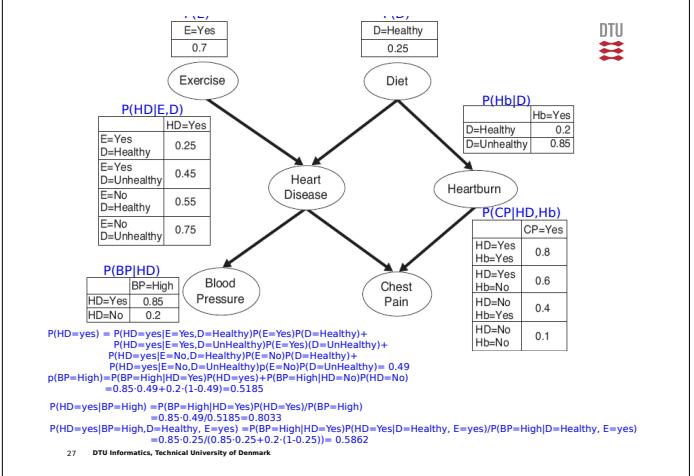
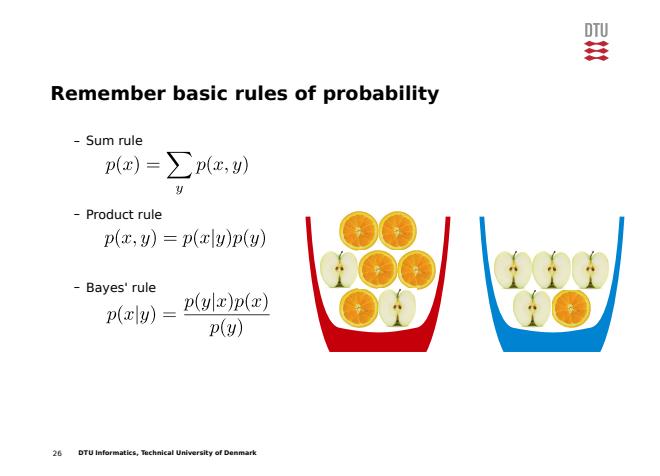
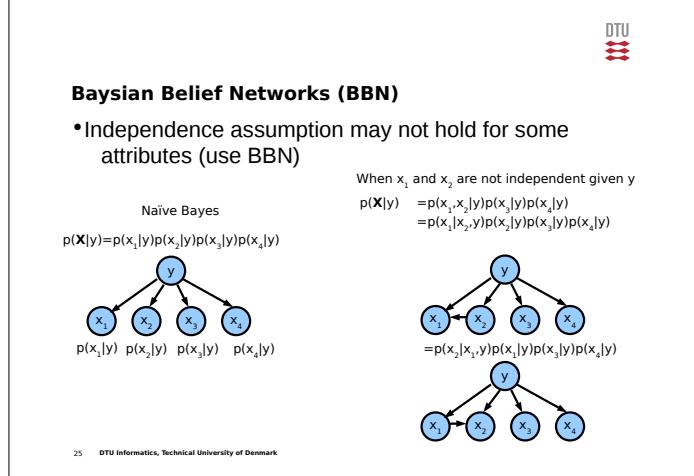
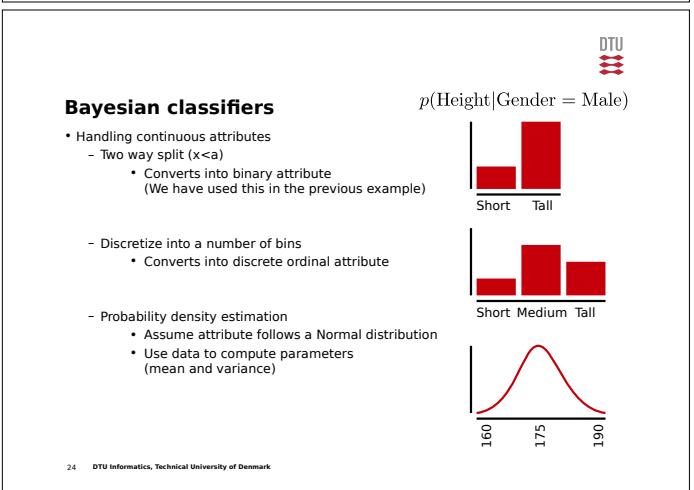
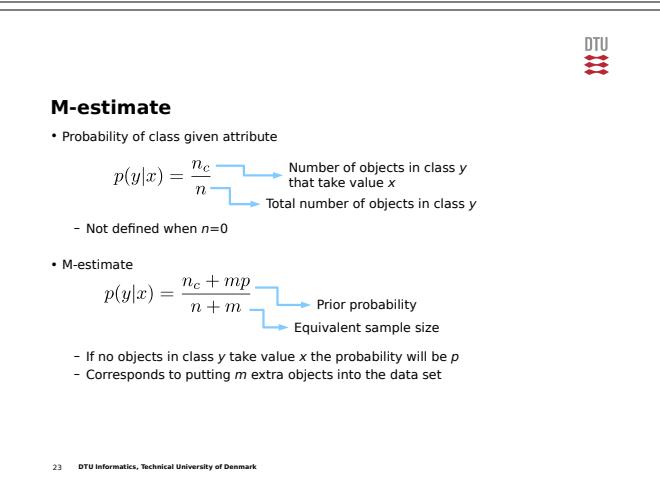
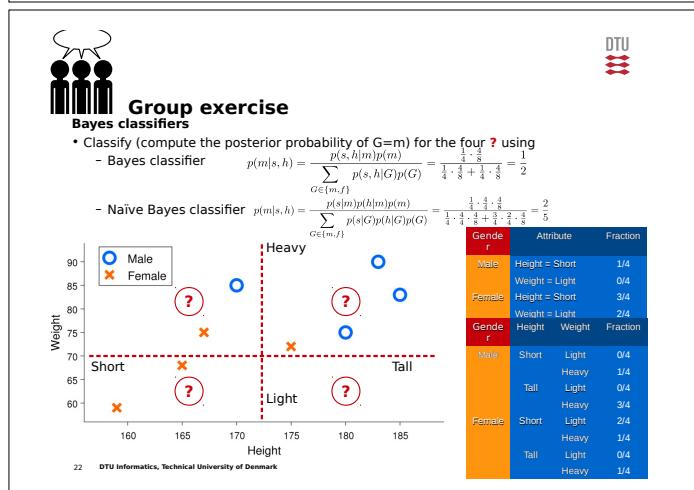
Bayesian classifiers

- Naive Bayes contingency table**
 - Only counts for each attribute
 - Small table

Height	Weight	Gender	
1	Tall	Heavy	Male
2	Tall	Heavy	Male
3	Short	Heavy	Male
4	Tall	Heavy	Male
5	Short	Light	Female
6	Short	Heavy	Female
7	Short	Light	Female
8	Tall	Light	Female

Gender	Attribute	Fraction
Male	Height=Short	1/4
	Weight=Light	0/4
Female	Height=Short	3/4
	Weight=Light	2/4

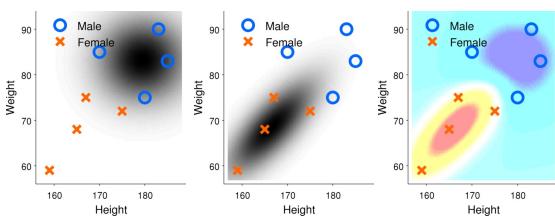
21 DTU Informatics, Technical University of Denmark



Bayesian classification by the multivariate normal distribution

Continuous density estimation $P(x|y=c) = \frac{1}{(2\pi)^{M/2} \det(\Sigma_c)^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_c)^\top \Sigma_c (x - \mu_c)\right)$

- Fit a Normal distribution to each class
 - Compute class mean and covariance
 - Classify using Bayes rule as before

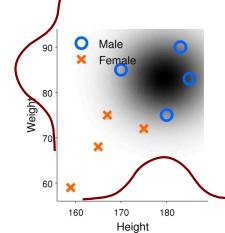


28 DTU Informatics, Technical University of Denmark



Group exercise

- What does the Naive Bayes assumption of independence of the attributes correspond to in terms of the parameters of the multivariate normal distribution?

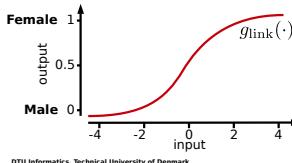
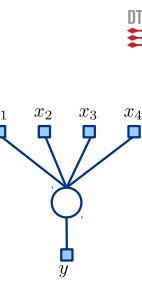


29 DTU Informatics, Technical University of Denmark

Artificial neural networks (ANN)

- Remember the generalized linear model?

- Data $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Model $f(\mathbf{x}) = g_{\text{link}}(\mathbf{x}^\top \mathbf{w})$
- Cost function $d(y, f(\mathbf{x}))$
- Parameters $\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N d(y_n, f(\mathbf{x}_n))$

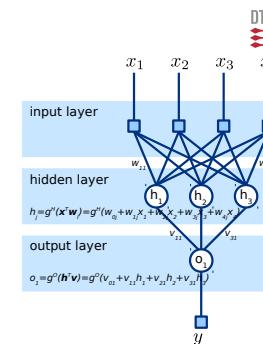


31 DTU Informatics, Technical University of Denmark

Artificial neural networks

Feed forward network

- Each "neuron"
 - Computes a non-linear function of the sum of its inputs
 - Is just like a generalized linear model
 - Has its own set of parameters
- Modeling choices
 - Cost function
 - Non-linearities
 - Number of neurons and hidden layers
 - Selection of inputs
- Parameter estimation using numerical optimization methods
- Very flexible model: Can easily overfit



32 DTU Informatics, Technical University of Denmark

Midterm practice test

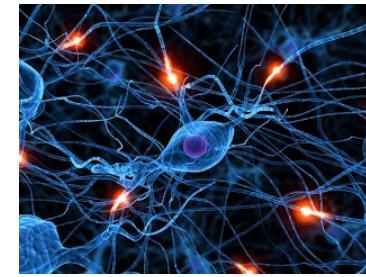
The midterm practice test is used solely for you to test your knowledge and for me to see how well you have understood the covered material so far.

The test **does not** count towards your grade for this course.

34 DTU Informatics, Technical University of Denmark

Artificial neural networks (ANN)

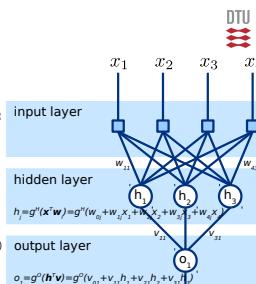
- The human brain contain in the order of 10^{11} neurons with 10^{15} connections
- Artificial Neural Networks (McCulloch & Pitts, 1943) are inspired by the architecture of the human brain



Artificial Neural Networks

- The ANN we will consider in the exercises:

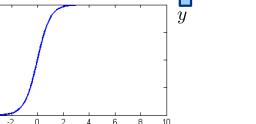
- Data $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Model $f(\mathbf{x}) = g^O(v_{10} + \sum_i g^H(x^\top w_i))$
- Cost function $d(y, f(\mathbf{x}))$
- Parameters $\mathbf{W}, \mathbf{v} = \arg \min_{\mathbf{W}, \mathbf{v}} \sum_{n=1}^N d(y_n, f(\mathbf{x}_n))$



- The implementation (in Matlab) uses:

$$\begin{aligned} d(y, f(\mathbf{x})) &= (y - f(\mathbf{x}))^2 \\ g^O(t) &= t \\ g^H(t) &= \tanh(t) \\ o_1 &= g^O(h^T v) = g^O(v_{o1} + v_{o2}h_1 + v_{o3}h_2 + v_{o4}h_3) \end{aligned}$$

33 DTU Informatics, Technical University of Denmark



02450 Introduction to machine learning and data mining

$$f(x+\Delta x) = \sum_{i=1}^n (\Delta x_i)^{\alpha_i} f^{\alpha_i}(x)$$

$$\Delta \int_a^b \Theta^{17} + \Omega^{\delta} e^{i\pi} = [2.7182818284 \sum_{i=1}^n \delta_i]$$

DTU Informatics
Department of Informatics and Mathematical Modeling

Midterm test: Question 9

When carrying out a principal component analysis of a dataset with four attributes we obtain the following singular values $\sigma_1=4$, $\sigma_2=2$, $\sigma_3=1$, and $\sigma_4=0$. Which one of the following statements is wrong?

- 1: The first principal component accounts for more than 60% of the variation in the data.
- 2: The third principal component accounts for less than 5 % of the variation in the data.
- 3: The second principal component accounts for more than 20 % of the variation in the data.
- 4: The data can be perfectly represented in a three dimensional sub-space.

7 DTU Informatics, Technical University of Denmark

Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project

10 DTU Informatics, Technical University of Denmark

Question 5

Consider a data set of four features: A, B, C, and D that are applied in a classification algorithm. The table below shows the cross-validated Error rate when using different combinations of the features.
We will apply a forward feature selection algorithm. Which feature set will the selection algorithm choose?

- C
- A and B
- A and B and C
- B and C and D
- Don't Know

Feature(s)	Error rate
A	0.40
B	0.45
C	0.33
D	0.42
A and B	0.20
A and C	0.25
A and D	0.34
B and C	0.29
B and D	0.42
C and D	0.40
A and B and C	0.13
A and B and D	0.17
B and C and D	0.10
A and C and D	0.15
A and B and C and D	0.28

4 DTU Informatics, Technical University of Denmark

Midterm test: Question 6

When training a decision tree we will use the classification error as impurity measure $I(t)$ given by $I(t) = 1 - \max_i p_i(t)$ where $p_i(t)$ denotes the fraction of data objects belonging to class i at a given node t . We will use Hunt's algorithm to grow the tree and recall that the purity gain is given by:

$$\Delta = I(\text{Parent}) - \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j)$$

where N is the total number of data objects at the parent node, k is the number of child nodes and $N(v_j)$ is the number of data objects associated with the child node, v_j . We will consider classification of Iris flowers into Iris-Setosa, Iris-Virginica and Iris-Versicolor. At a potential split we have:

- Before the split: 5 Iris-Setosa, 10 Iris-Virginica and 10 Iris Versicolor.

• After the split:

- 0 Iris-Setosa, 8 Iris-Virginica and 2 Iris-Versicolor in the left node.
- 5 Iris-Setosa, 2 Iris-Virginica and 8 Iris-Versicolor in the right node.

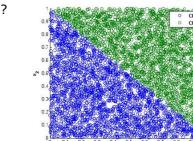
Which statement is correct?

- 1: The purity gain is 3/5.
- 2: The purity gain is 3/15.
- 3: The purity gain is 6/25.
- 4: The purity gain is 7/15.

6 DTU Informatics, Technical University of Denmark

Midterm test: Question 12

Consider the classification problem given in Figure 5 where x_1 and x_2 are used as features for a logistic regression classifier and a decision tree. The considered logistic regression models all include the constant term w_0 . Which one of the following statements is wrong?



- 1: The two classes can be perfectly separated by a logistic regression model using x_1 and x_2 as features.
- 2: A decision tree with less than five nodes can perfectly separate the classes using only x_1 and x_2 as features.
- 3: A logistic regression model can perfectly separate the two classes using only the feature t given by $t = x_1 + x_2$.
- 4: In logistic regression the probability that each observation belong to the two classes can be derived from the logit link function.

8 DTU Informatics, Technical University of Denmark

Reading material

Tan, Steinbach and Kumar "Introduction to Data Mining"

Section 5.6-5.8

Feedback Groups of the day:

Jakob Kaufmann Jespersen

Morten Bjelbo Thomsen, Rolf s123490

Arthur Katovsky, Anders Emil Nielsen

Mette Nielsen, Henrik Thistrup

Utku Norman, Darren Williams, Jose Esteves

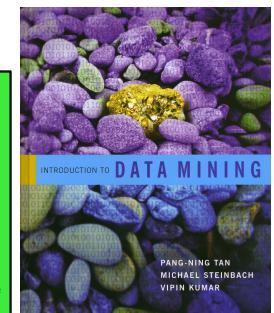
Henry Henry, Dominik Otto, Felix Ngajieh Ekwoge

Simon B. Hemmingsen,

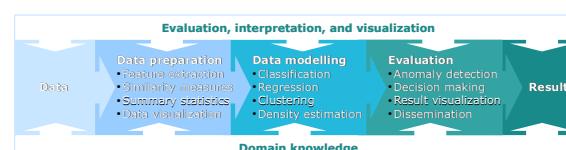
Thomas Kjæller Hammerbak

Minh Haw Truong

If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at my exercises next week with feedback/suggestions on the exercises for today.



Data modeling framework



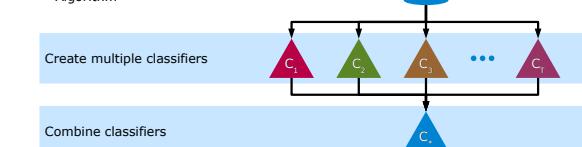
After today you should be able to:

Explain the principle behind boosting and bagging and apply it to improve classifiers
Be able to address issues of class-imbalances by resampling
Understand the definition of Precision, Recall, ROC and AUC
Be able to extend binary classifiers to multi-class classification

11 DTU Informatics, Technical University of Denmark

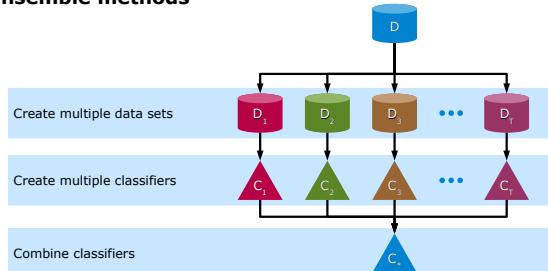
Ensemble methods

- Combine multiple (weak) classifiers into one (strong) classifier
- Each classifier trained using different variations of
 - Data set
 - Input attributes
 - Class labels
 - Algorithm



12 DTU Informatics, Technical University of Denmark

Ensemble methods



13 DTU Informatics, Technical University of Denmark

Why ensemble methods?

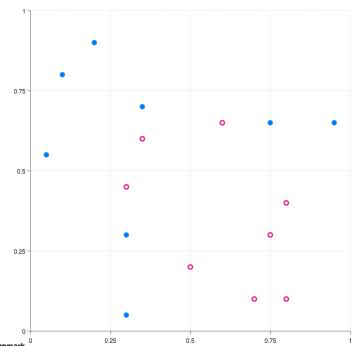
- Can improve classification algorithms in terms of
 - Better classification accuracy
 - Increased stability
 - Reduced variance
 - Less overfitting
- Consider $N=25$ independent classifiers for binary classification, each with error rate $\epsilon=0.35$

$$\begin{aligned} \epsilon_* &= \sum_{n=\lceil \frac{N}{2} \rceil}^N \binom{N}{n} \epsilon^n (1-\epsilon)^{N-n} \\ &= \sum_{n=13}^{25} \binom{25}{n} 0.35^n (1-0.35)^{25-n} = 0.06 \end{aligned}$$

14 DTU Informatics, Technical University of Denmark

Data example

- Classification using logistic regression



15 DTU Informatics, Technical University of Denmark

Bagging

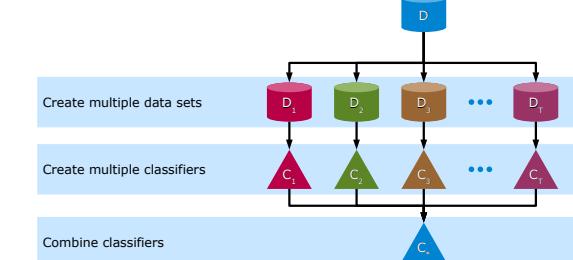
- New training data sets drawn randomly from pool with replacement

Pool of training data	1	2	3	4	5	6	7	8	9	10
New training data sets	3	5	4	3	9	7	9	5	1	1
	5	8	2	6	2	3	8	3	5	1
	1	7	4	1	10	6	10	8	8	7
	4	3	8	5	2	4	7	10	10	8

16 DTU Informatics, Technical University of Denmark

Bagging

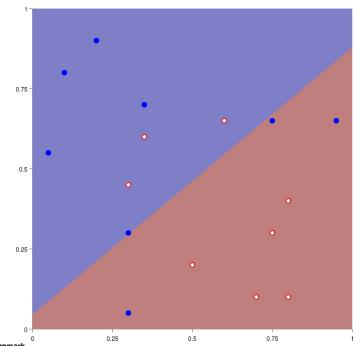
- Single classifier



17 DTU Informatics, Technical University of Denmark

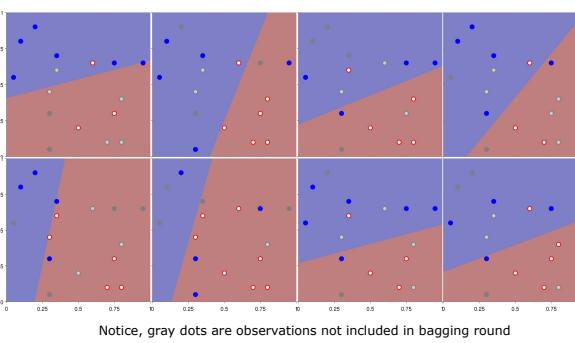
Bagging

- Single classifier
 - Logistic regression
 - Two features, (x, y)



18 DTU Informatics, Technical University of Denmark

Bagging

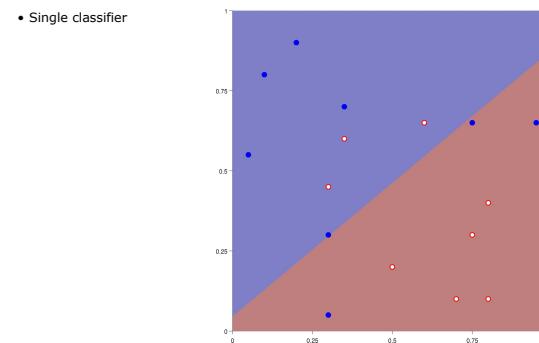


Notice, gray dots are observations not included in bagging round

19 DTU Informatics, Technical University of Denmark

Bagging

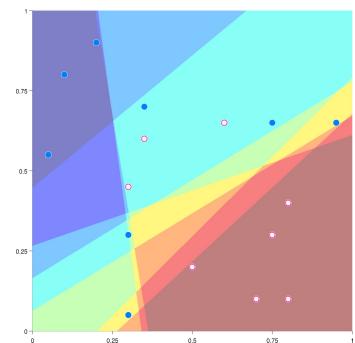
- Single classifier



20 DTU Informatics, Technical University of Denmark

Bagging

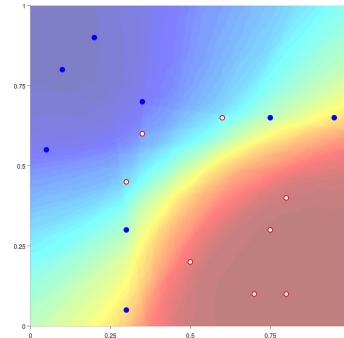
- Majority voting of 10 bagged classifiers



21 DTU Informatics, Technical University of Denmark

Bagging

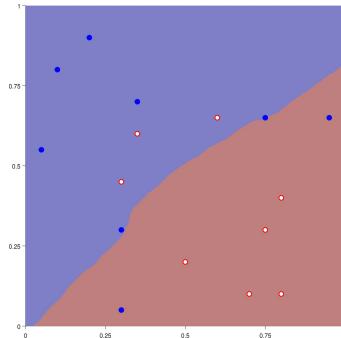
- Majority voting of 100 bagged classifiers



22 DTU Informatics, Technical University of Denmark

Bagging

- Decision boundary of 100 bagged classifier ensemble



23 DTU Informatics, Technical University of Denmark

Boosting

Pool of training data

1	2	3	4	5	6	7	8	9	10
.1	.1	.1	.1	.1	.1	.1	.1	.1	.1

Weights

3	5	4	3	9	7	9	5	1	1
---	---	---	---	---	---	---	---	---	---

New training data set

C ₁

Train classifier

Boosting

Pool of training data	1	2	3	4	5	6	7	8	9	10
Weights	.1	.1	.1	.1	.1	.1	.1	.1	.1	.1
New training data set	3	5	4	3	9	7	9	5	1	1
Train classifier	C ₁									
Classify all data objects	1✓	2✗	3✓	4✗	5✓	6✗	7✓	8✓	9✓	10✓

25 DTU Informatics, Technical University of Denmark

Boosting

Pool of training data	1	2	3	4	5	6	7	8	9	10
Weights	.1	.1	.1	.1	.1	.1	.1	.1	.1	.1
New training data set	3	5	4	3	9	7	9	5	1	1
Train classifier	C ₁									
Classify all data objects	1✓	2✗	3✓	4✗	5✓	6✗	7✓	8✓	9✓	10✓
Update weights	.07	.17	.07	.17	.07	.17	.07	.07	.07	.07

26 DTU Informatics, Technical University of Denmark

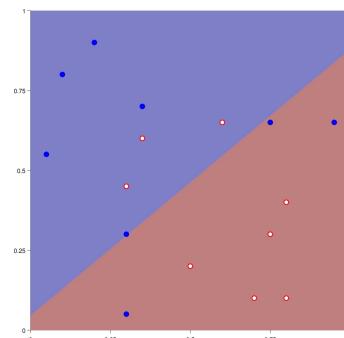
Boosting

Pool of training data	1	2	3	4	5	6	7	8	9	10
Weights	.1	.1	.1	.1	.1	.1	.1	.1	.1	.1
New training data set	3	5	4	3	9	7	9	5	1	1
Train classifier	C ₁									
Classify all data objects	1✓	2✗	3✓	4✗	5✓	6✗	7✓	8✓	9✓	10✓
Update weights	.07	.17	.07	.17	.07	.17	.07	.07	.07	.07
New training data set	6	4	7	3	2	4	10	2	5	6
Train classifier	C ₂									

27 DTU Informatics, Technical University of Denmark

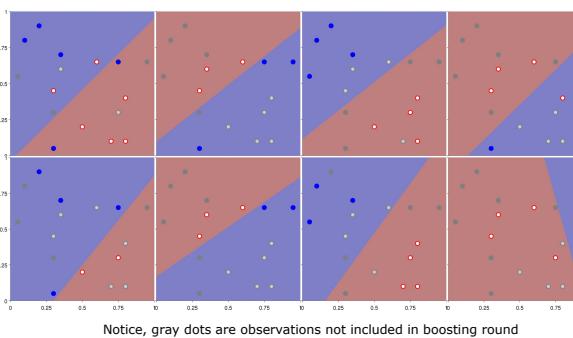
Boosting

- Single classifier



28 DTU Informatics, Technical University of Denmark

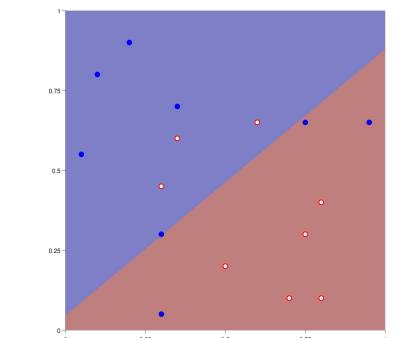
Boosting



29 DTU Informatics, Technical University of Denmark

Boosting

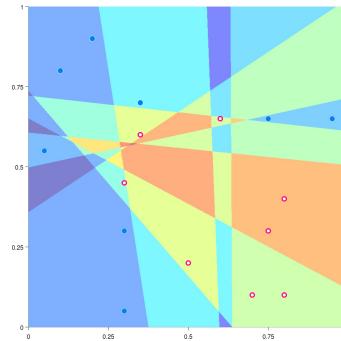
- Single classifier



30 DTU Informatics, Technical University of Denmark

Boosting

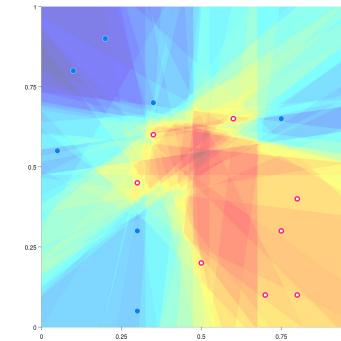
- Majority voting of 10 boosted classifiers



31 DTU Informatics, Technical University of Denmark

Boosting

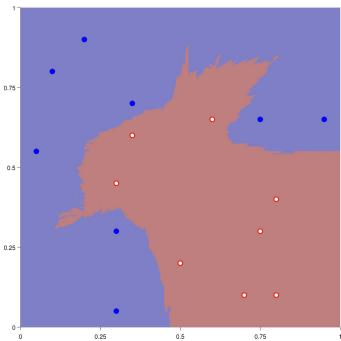
- Majority voting of 100 boosted classifiers



32 DTU Informatics, Technical University of Denmark

Boosting

- Decision boundary of 100 boosted classifier ensemble



33 DTU Informatics, Technical University of Denmark

Class imbalance problem

- Many data sets have **imbalanced class distributions**
 - Example: Detection of defects that only occur rarely (e.g. 1/1,000,000)
 - Danger: Algorithm that says nothing is defect will be 99.999% correct
- Solution approaches**
 - Resample to balance data sets
 - Modify existing classification algorithms
 - Measure performance in a way that takes balance into account

34 DTU Informatics, Technical University of Denmark

Resampling balanced data

- New sample has equal number of data objects from each class
- Approaches**
 - Oversampling** majority class: Throws out potentially useful data
 - Oversampling** minority class: Increase data size and computational burden
 - Somewhere in between...**

Imbalanced training data	1 2 3 4 5 6 7 8 9 10
Oversampling	1 2 3 4 5 1 2 3 6 6
Undersampling	3 5 6 8
Somewhere in between	3 5 4 3 9 6 6 8 8 8

35 DTU Informatics, Technical University of Denmark

Confusion matrix

		Predicted	
		positive	negative
Actual	positive	TP True Positive	FN False Negative
	negative	FP False Positive	TN True Negative

36 DTU Informatics, Technical University of Denmark

Precision and recall

- Precision**
 - Fraction of true positive among objects predicted to be positive

$$p = \frac{TP}{TP + FP}$$

- Recall**
 - Fraction of objects predicted to be positive among all positive objects

$$r = \frac{TP}{TP + FN}$$

		Predicted	
		positive	negative
Actual	positive	TP True Positive	FN False Negative
	negative	FP False Positive	TN True Negative



37 DTU Informatics, Technical University of Denmark



Group exercise

- You consider two different classifiers, on a test set with 20 positive objects

- **Classifier 1** detects 54 positives of which 18 are actually positive

- **Classifier 2** detects 16 positives of which 14 are actually positive

- Compute the **precision** and **recall** for the two classifiers

- Which classifier (if any) is the best?
- Which would you use if the objective is to detect credit card fraud (consider what is most costly – **missing** or **falsely detecting** a positive)

Precision

- Fraction of true positive among objects predicted to be positive
 $p = \frac{TP}{TP + FP}$

Recall

- Fraction of objects predicted to be positive among all positive objects
 $r = \frac{TP}{TP + FN}$

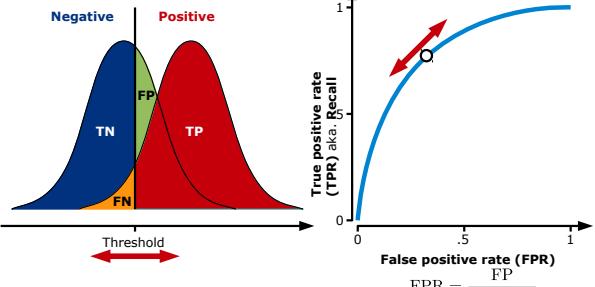
Predicted

		Predicted	
		positive	negative
Actual	positive	TP True Positive	FN False Negative
	negative	FP False Positive	TN True Negative

38 DTU Informatics, Technical University of Denmark

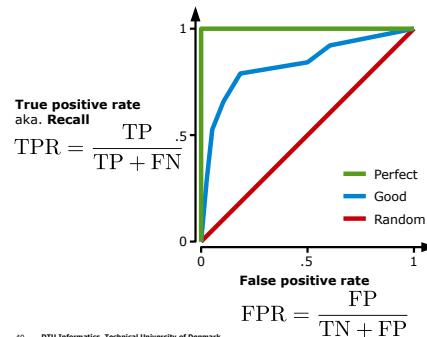
Receiver operating characteristic

$$TPR = \frac{TP}{TP + FN}$$



39 DTU Informatics, Technical University of Denmark

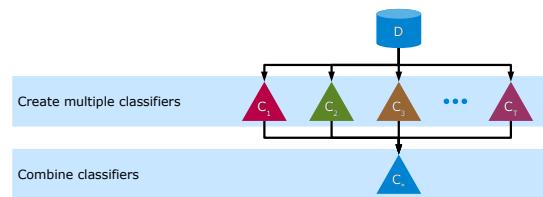
Receiver operating characteristic



40 DTU Informatics, Technical University of Denmark

Multiclass problems

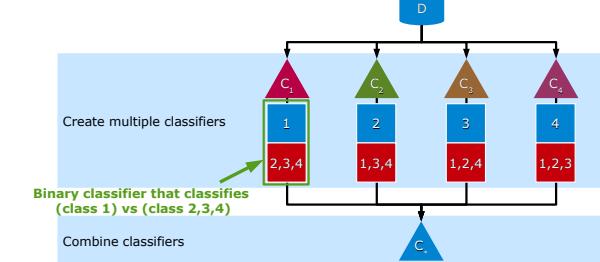
- Combine multiple **binary** classifiers into one **multiclass** classifier
 - Each classifier trained using different class labels



41 DTU Informatics, Technical University of Denmark

Multiclass problems

- Combine multiple **binary** classifiers into one **multiclass** classifier
 - Each classifier trained using different class labels



42 DTU Informatics, Technical University of Denmark

Multiclass problems

- Classification algorithms designed for **binary classification** can also be applied to **multiclass problems**
 - Train a number of classifiers and make final classification by **majority voting**

1-against-rest

- K classifiers

1	2	3	4
2,3,4	1,3,4	1,2,4	1,2,3

1-against-1

- $K(K-1)/2$ classifiers

1	1	1	2	2	3
2	3	4	3	4	4

Error correcting output coding

- Robustness against errors

1	1,4	1,3	1,3,4	1,2	1,2,4	1,2,3
2,3,4	2,3	2,4	2	3,4	3	4

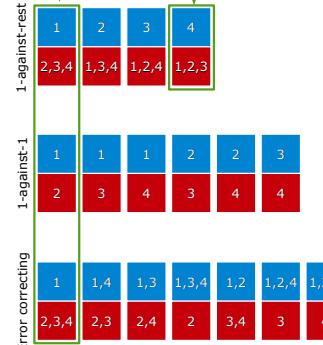
43 DTU Informatics, Technical University of Denmark



Group exercise

- Assuming that **all classifiers are perfect**
 - Use the three multiclass methods to classify a data object belonging to class 1

Hint: Go through the classifiers and add up the votes, giving one vote to each of the winning classes. (In 1-against-1 give 0.5 to each class when class 1 is not included.)
- Assuming that **the first classifier is broken** and makes a wrong classification
 - Repeat the exercise



44 DTU Informatics, Technical University of Denmark

Remember one-out-of-K coding

Nationality

	Denmark	Norway	Sweden
'Sweden'	0	0	1
'Sweden'	0	0	1
'Sweden'	0	0	1
'Norway'	0	1	0
'Norway'	0	1	0
'Norway'	0	1	0
'Norway'	0	1	0
'Norway'	0	1	0
'Norway'	0	1	0
'Norway'	0	1	0
'Denmark'	1	0	0
'Denmark'	1	0	0
'Sweden'	0	0	1
'Sweden'	0	0	1
'Denmark'	1	0	0
'Denmark'	1	0	0
'Sweden'	0	0	1
'Sweden'	0	0	1
'Norway'	0	1	0
'Denmark'	1	0	0

One-out-of-K coding

45 DTU Informatics, Technical University of Denmark

Extending classifiers to handle multi-class problems

Logistic regression

$$\{x_n, y_n\}_{n=1}^N$$

$$f(x) = \text{logit}(x^\top w) = \frac{1}{1 + \exp(-x^\top w)}$$

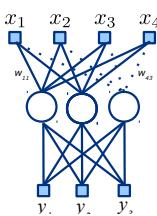


Multinomial Regression

$$\{x_n, y_n\}_{n=1}^N$$

$$f_c(x) = \text{softmax}(x^\top W) = \frac{\exp(x^\top w_c)}{\sum_{c'} \exp(x^\top w_{c'})}$$

$$w_C = 0$$



46 DTU Informatics, Technical University of Denmark



- Show that for a binary (i.e., two class problem) multinomial regression is equivalent to logistic regression if $w_0 = 0$.

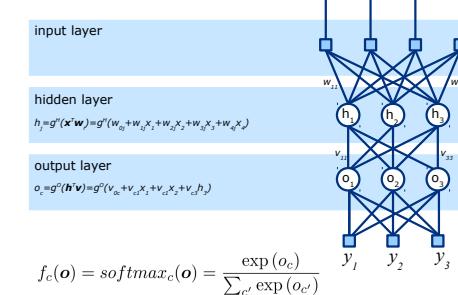
Hint: consider the softmax function for two classes:

$$f_c(x) = \text{softmax}_c(x^\top W) = \frac{\exp(x^\top w_c)}{\sum_{c'} \exp(x^\top w_{c'})}$$

and show that this is equivalent to the logit function:

$$f(x) = \text{logit}(x^\top w) = \frac{1}{1 + \exp(-x^\top w)}$$

Artificial Neural Networks for multiclass classification



48 DTU Informatics, Technical University of Denmark

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

$$f(x+\Delta x) = \sum_{i=1}^n f_i(x)$$

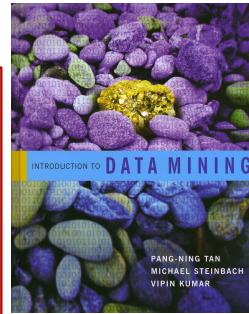
$$\Delta \int_a^b \Theta + \Omega \int \delta e^{i\pi} = \sqrt{17} \cdot [2.7182818284] \sum_{i=1}^n,$$

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"

Section 8.1-8.3+8.5.7

Feedback Groups of the day:
 Carlos Roncal, Antonio Martin Rodriguez
 Julian Bellino, Simon Thordal
 Siham Bargadouch, Ali Ihsan, Abiramy Mohanarajah
 Kristian K. Rüsgaard, Tryfon Tzanetis Ioanna Psylla
 Michael Raagaard, Mathias Kaas-Olsen
 Joachim Blom Hansen, Roman Piper
 Karol Dzitkowski, Marco Becattini
Thorbjørn Wolf, Mick Neupart
 If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

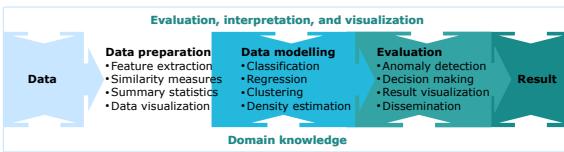


Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
- Supervised learning: Classification and regression**
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
- Unsupervised learning: Clustering and density est.**
9. **K-means and hierarchical clustering**
(Tan 8.1-8.3+8.5.7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
- Machine learning and data modelling in practice**
12. Putting it all together: Summary and overview
13. Mini project

3 DTU Informatics, Technical University of Denmark

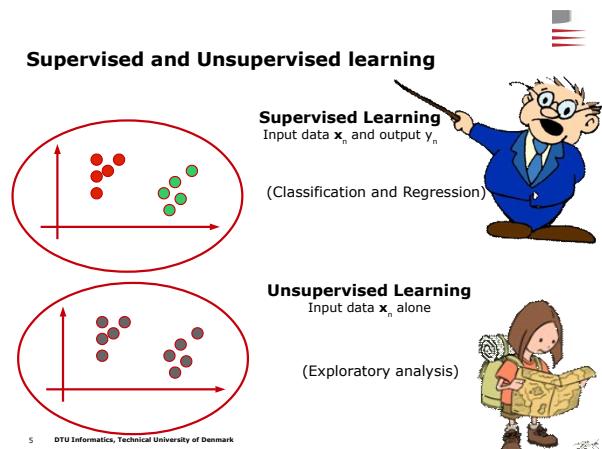
Data modeling framework



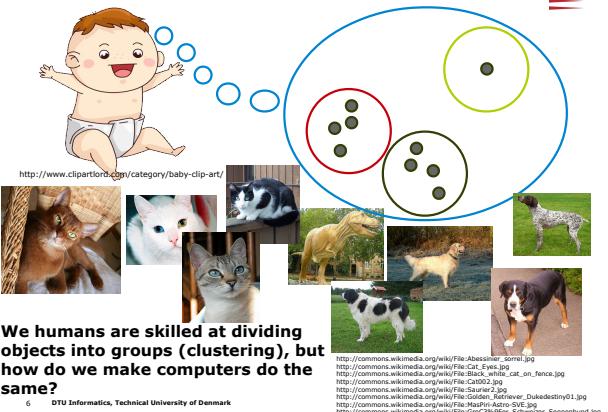
After today you should be able to:
 Discuss the aims of unsupervised learning
 Explain the principles behind hierarchical clustering and k-means clustering
 Apply hierarchical clustering and k-means clustering
 Evaluate the quality of the clustering using supervised measures of cluster validity

4 DTU Informatics, Technical University of Denmark

Supervised and Unsupervised learning



Imagine you observe the world for the first time!



5 DTU Informatics, Technical University of Denmark

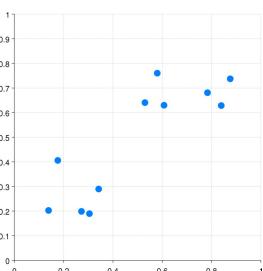
Unsupervised learning

- **Supervised learning**
 - Use the data to learn the output values
- **Unsupervised learning**
 - No output variables available
 - Sometimes called exploratory analysis
 - What learn from the data?
 - Structure
 - Regularities
 - Hidden information
 - Etc.

Clustering

- Divide data into groups (subsets/clusters) that are
 - **Meaningful:** Capture the natural structure of the data
 - **Useful:** Depends on purpose
- Observations in the same cluster are **similar in some sense**
- Unsupervised classification

Clustering

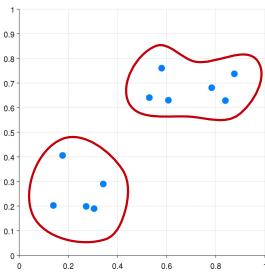


7 DTU Informatics, Technical University of Denmark

8 DTU Informatics, Technical University of Denmark

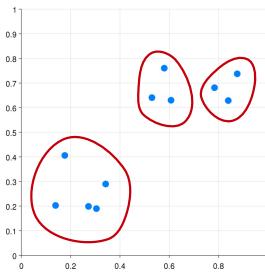
9 DTU Informatics, Technical University of Denmark

Clustering



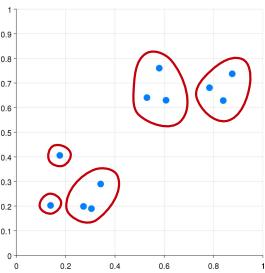
10 DTU Informatics, Technical University of Denmark

Clustering



11 DTU Informatics, Technical University of Denmark

Clustering

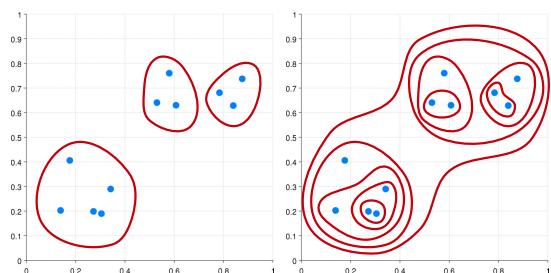


12 DTU Informatics, Technical University of Denmark

Partitional / hierarchical clustering

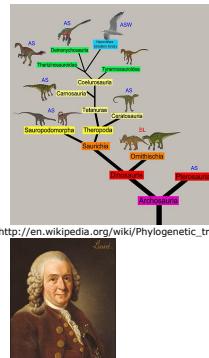
Partitional

Hierarchical



13 DTU Informatics, Technical University of Denmark

Phylogenetic trees may be considered a type of hierarchical clustering

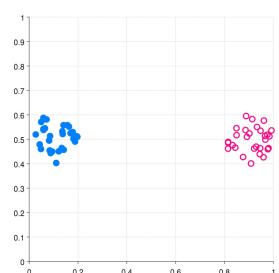


http://en.wikipedia.org/wiki/Phylogenetic_tree

Types of clustering

Well-separated

- Each point is closer to all points in its cluster than any point in another cluster

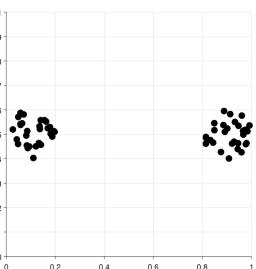


16 DTU Informatics, Technical University of Denmark

Types of clustering

Well-separated

- Each point is closer to all points in its cluster than any point in another cluster

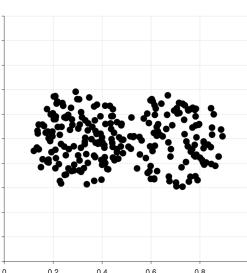


15 DTU Informatics, Technical University of Denmark

Types of clustering

Center-based

- Each point is closer to the center of its cluster than to the center of any other cluster

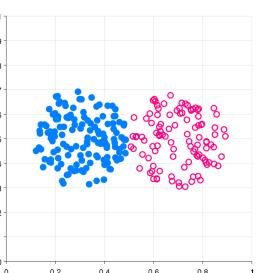


17 DTU Informatics, Technical University of Denmark

Types of clustering

Center-based

- Each point is closer to the center of its cluster than to the center of any other cluster

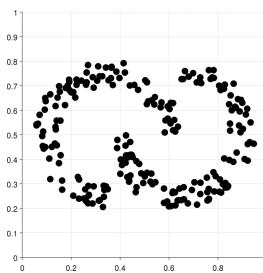


18 DTU Informatics, Technical University of Denmark

Types of clustering

Contiguity-based

- Each point is closer to at least one point in its cluster than to any point in another cluster



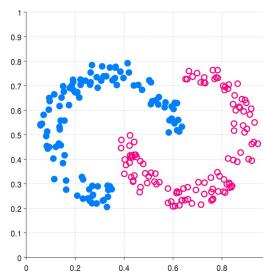
19 DTU Informatics, Technical University of Denmark

Types of clustering

Contiguity-based

- Each point is closer to at least one point in its cluster than to any point in another cluster

Clustering based on neighbor distance

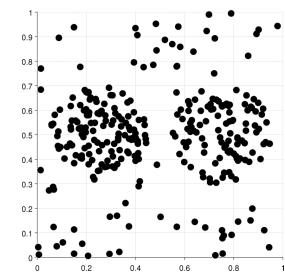


20 DTU Informatics, Technical University of Denmark

Types of clustering

Density-based

- Clusters are regions of high density separated by regions of low density

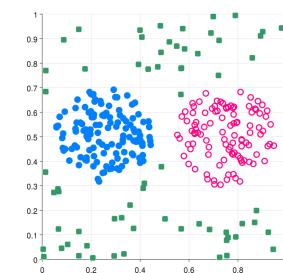


21 DTU Informatics, Technical University of Denmark

Types of clustering

Density-based

- Clusters are regions of high density separated by regions of low density



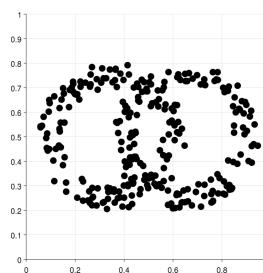
22 DTU Informatics, Technical University of Denmark

Types of clustering

Conceptual clusters

- Points in a cluster share some general property that derives from the entire set of points

Requires some intelligence to parse

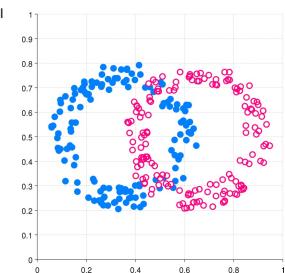


23 DTU Informatics, Technical University of Denmark

Types of clustering

Conceptual clusters

- Points in a cluster share some general property that derives from the entire set of points



24 DTU Informatics, Technical University of Denmark

Group exercise

Using the five criteria

- How will these points be clustered?
- How many clusters?

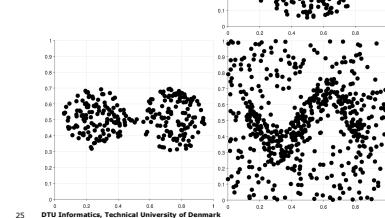
- Well-separated**
- Each point is closer to all points in its cluster than any point in another cluster

- Center-based**
- Each point is closer to the center of its cluster than to the center of any other cluster

- Contiguity-based**
- Each point is closer to at least one point in its cluster than to any point in another cluster

- Density-based**
- Clusters are regions of high density separated by regions of low density

- Conceptual clusters**
- Points in a cluster share some general property that derives from the entire set of points



25 DTU Informatics, Technical University of Denmark

K-means clustering

Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change

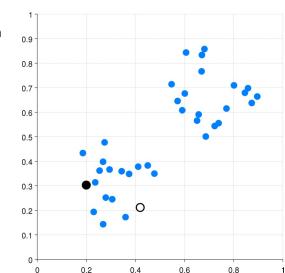
K-means clustering

Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change



26 DTU Informatics, Technical University of Denmark

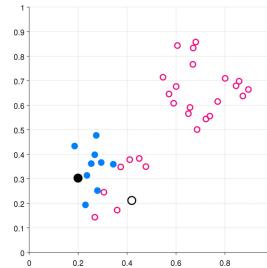
K-means clustering

Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change



28 DTU Informatics, Technical University of Denmark

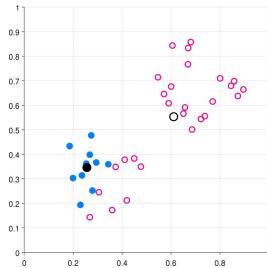
K-means clustering

Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change



29 DTU Informatics, Technical University of Denmark

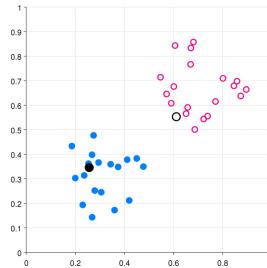
K-means clustering

Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change



30 DTU Informatics, Technical University of Denmark

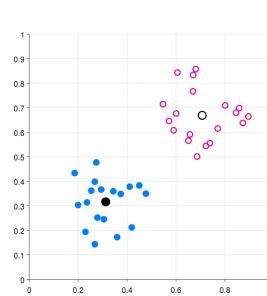
K-means clustering

Select K points as initial centroids

Repeat

- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change



31 DTU Informatics, Technical University of Denmark

K-means clustering

How do I

- Find the closest centroid?
 - Use a suitable **dissimilarity/similarity measure**
- Compute the cluster centroids
 - Depends on dissimilarity/similarity measure
 - For example, for Euclidean distance the mean is optimal (See Section 8.2.6 in Tan et al.)

To find the center of a cluster we minimize the L1 error.
L1 is minimal when $c_1 = \text{mean of } X_1 - c_1$ distances

32 DTU Informatics, Technical University of Denmark

Group exercise

Using pen-and-paper k-means, cluster the following data objects

- Number of clusters
 - $K=2$
- Distance measure
 - Euclidean
- Computation of centroid
 - Mean of cluster members
- Initial centroids
 - For example the first two data objects
 - In case of any ties, flip a coin to decide
- Data objects

$$x = \{42, 60, 17, 48, 12\}$$

Select K points as initial centroids

Repeat

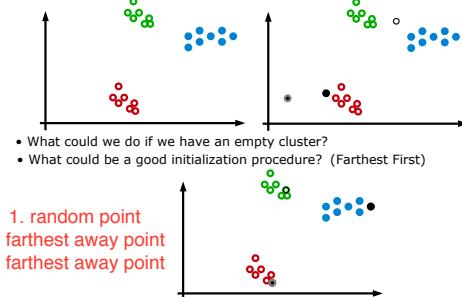
- Form K clusters by assigning each point to its closest centroid
- Recompute the centroids of each cluster

Until centroids do not change

33 DTU Informatics, Technical University of Denmark



How will the data (top-left diagram) be clustered given the initialization of the three centroids shown at the right and at the bottom?



34 DTU Informatics, Technical University of Denmark

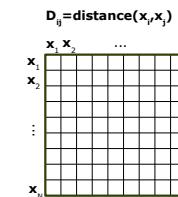
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



35 DTU Informatics, Technical University of Denmark

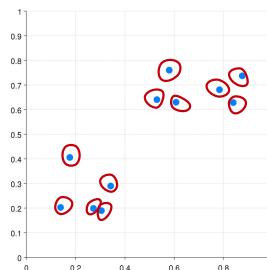
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



36 DTU Informatics, Technical University of Denmark

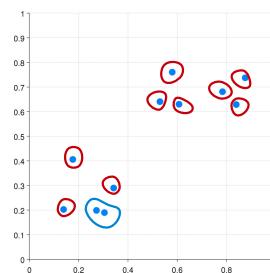
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



37 DTU Informatics, Technical University of Denmark

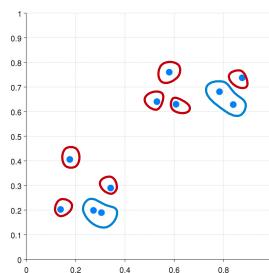
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



38 DTU Informatics, Technical University of Denmark

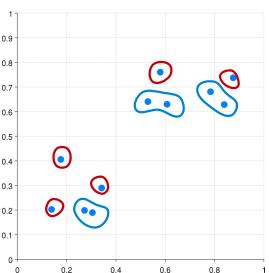
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



39 DTU Informatics, Technical University of Denmark

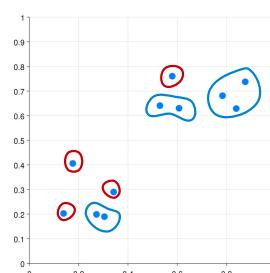
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



40 DTU Informatics, Technical University of Denmark

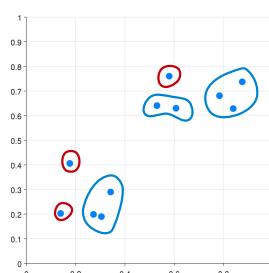
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



41 DTU Informatics, Technical University of Denmark

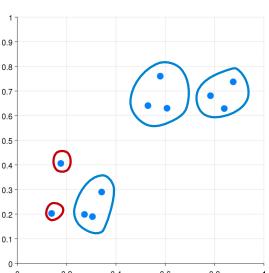
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



42 DTU Informatics, Technical University of Denmark

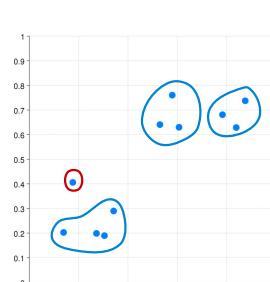
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



43 DTU Informatics, Technical University of Denmark

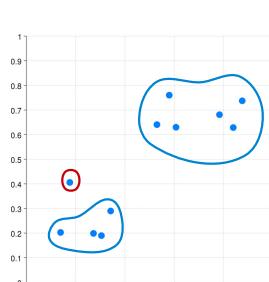
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



44 DTU Informatics, Technical University of Denmark

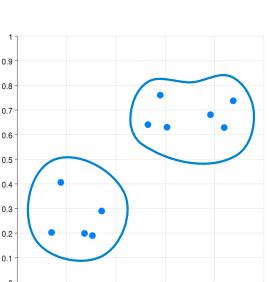
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains



45 DTU Informatics, Technical University of Denmark

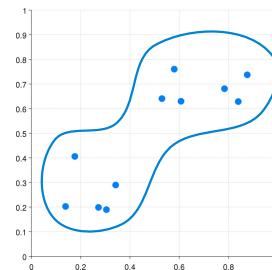
Agglomerative hierarchical clustering

Compute the proximity matrix

Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

Until only one cluster remains

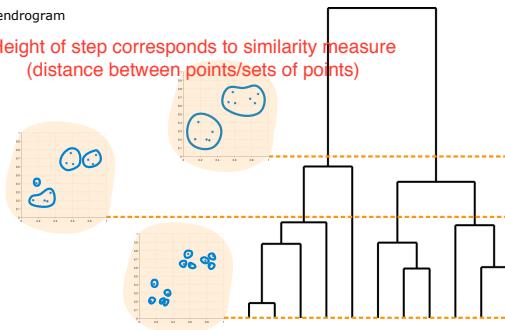


46 DTU Informatics, Technical University of Denmark

Agglomerative hierarchical clustering

• Dendrogram

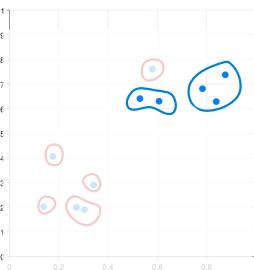
Height of step corresponds to similarity measure
(distance between points/sets of points)



47 DTU Informatics, Technical University of Denmark

Similarity between clusters

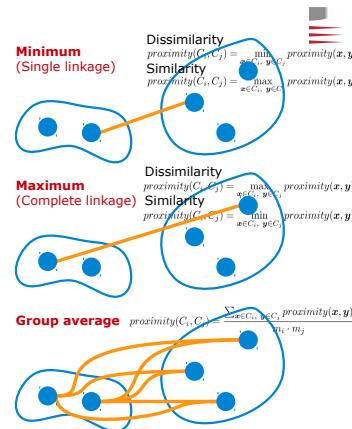
- The key operation in agglomerative hierarchical clustering is measuring similarity between clusters



48 DTU Informatics, Technical University of Denmark

Proximity between clusters

- Can be computed using proximity between objects
- Notice we need different definition if we are given a similarity or dissimilarity measure
- In our example before we used Euclidean distance as proximity measure; i.e. it is the first definition which is relevant (dissimilarity)

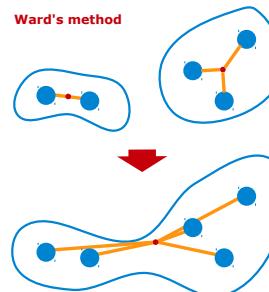


C_i : Observations in cluster i
 C_j : Observations in cluster j
 m_i : Number of observations in cluster i
 m_j : Number of observations in cluster j

49 DTU Informatics, Technical University of Denmark

Similarity between clusters

- Increase in sum of squared error after merging the two clusters should be as small as possible

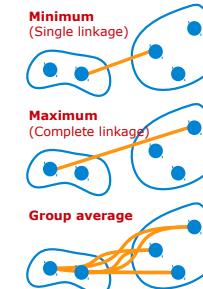


50 DTU Informatics, Technical University of Denmark



Group exercise

Can the choice of linkage be related to the notion of what constitutes clusters?



51 DTU Informatics, Technical University of Denmark

- Well-separated**
 - Each point is closer to all points in its cluster than any point in another cluster
- Center-based**
 - Each point is closer to the center of its cluster than to the center of any other cluster
- Contiguity-based**
 - Each point is closer to at least one point in its cluster than to any point in another cluster
- Density-based**
 - Clusters are regions of high density separated by regions of low density
- Common class**
 - Points in a cluster share some general property that derives from the entire set of points



Group exercise

Using pen-and-paper agglomerative hierarchical clustering, cluster the following data objects and draw the dendrogram

- Distance measure
 - Euclidean
- Similarity between clusters
 - Minimum (Single linkage)

Data objects

$$x = \{42, 60, 17, 48, 12\}$$

Compute the proximity matrix
Repeat

- Merge the two closest clusters
- Update the proximity matrix to reflect the proximity between the new cluster and the original clusters

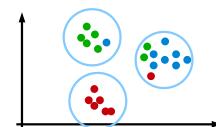
Until only one cluster remains

Minimum (Single linkage)

52 DTU Informatics, Technical University of Denmark

Cluster Purity measures when class labels are available

Motivation: Evaluate the extent to which manual classification process can be automatically produced by cluster analysis by comparing clustering to "ground truth"



54 DTU Informatics, Technical University of Denmark

Supervised measures of cluster validity

If we have (supervised) class labels, they can be used to measure the quality of the clustering

• Cluster purity measures

- (recall these from decision trees)
 - Entropy
 - Gini
 - Class error

• Class label accuracy measures

- Precision
- Recall
- F-measure

• Binary similarity measures

- Simple matching coefficient
- Jaccard coefficient

55 DTU Informatics, Technical University of Denmark

Supervised measures of cluster validity

Cluster purity measures

- Entropy

$$\text{Entropy}(i) = - \sum_j p_{ij} \log_2 p_{ij}$$

- Gini

$$\text{Gini}(i) = 1 - \sum_j p_{ij}^2$$

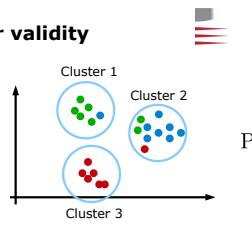
- Class error

$$\text{ClassError}(i) = 1 - \max_j p_{ij}$$

- Purity

$$\text{Purity}(i) = \max_j p_{ij} \quad p_{ij} : \text{Probability that member of cluster } i \text{ belongs to class } j \quad p_{ij} = \frac{m_{ij}}{m_i}$$

What is the Entropy, Gini, Class. Error and Purity for each of the three clusters given above?



Notice that $\theta \log(0)=0$

56 DTU Informatics, Technical University of Denmark

Supervised measures of cluster validity

Class label accuracy measures

- Precision

$$\text{Precision}(i, j) = \frac{m_{ij}}{m_i} = p_{ij}$$

- Recall

$$\text{Recall}(i, j) = \frac{m_{ij}}{m_j}$$

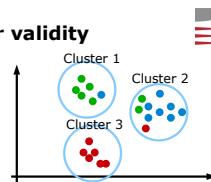
- F-measure

$$F(i, j) = \frac{2 \cdot \text{Precision}(i, j) \cdot \text{Recall}(i, j)}{\text{Precision}(i, j) + \text{Recall}(i, j)} = \frac{2m_{ij}}{m_i + m_j}$$

What is the Precision, recall and F-measure of the blue class in cluster 2?

$$\begin{aligned} \text{Precision(blue,cluster 2)} &= 7/10 \\ \text{Recall(blue,cluster 2)} &= 7/8 \\ F(\text{blue,cluster 2}) &= 2 \cdot 7/(8+10) = 7/9 \end{aligned}$$

57 DTU Informatics, Technical University of Denmark



Supervised measures of cluster validity

Binary similarity measures

- Simple matching coefficient (SMC)/Rand statistic

$$\text{SMC}(x, y) = \frac{f_{00} + f_{11}}{K}$$

- Jaccard coefficient

$$J(x, y) = \frac{f_{11}}{K - f_{00}}$$

K : Total number of pairs of objects, $N \cdot (N-1)/2$

f_{00} : Number of object pairs in different class assigned to different clusters

f_{11} : Number of object pairs in same class assigned to same cluster

In our example we find:

$$K = 22 \cdot (22-1)/2 = 231$$

$$f_{11} = (5 \cdot (5-1)/2 + 1 \cdot (1-1)/2)_{c1} + (7 \cdot (7-1)/2 + 2 \cdot (2-1)/2 + 1 \cdot (1-1)/2)_{c2} + (6 \cdot (6-1)/2)_{c3} = 10 + 22 + 15 = 47$$

$$f_{00} = (5 \cdot (7+1) + 1 \cdot (2+1) + 0 \cdot (2+7))_{c1 \rightarrow c2} + (5 \cdot 6 + 1 \cdot 6 + 0 \cdot 0)_{c1 \rightarrow c3} + (2 \cdot 6 + 7 \cdot 6 + 1 \cdot 0)_{c2 \rightarrow c3} = 43 + 36 + 54 = 133$$

$$\text{SMC} = (47+133)/231 = 180/231$$

$$\text{Jaccard} = 47/(231-133) = 47/98$$

58 DTU Informatics, Technical University of Denmark

Exam question examples

QUESTION I:

We have a one dimensional data set of size $N = 5$ with data examples
 $x_1 = 1, x_2 = 3, x_3 = 6, x_4 = 7$ and $x_5 = 12$.

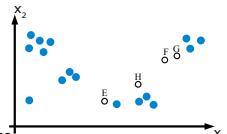
We run hierarchical clustering with a Euclidean dissimilarity between data points using group average linkage. We will use the following notation to summarize the dendrogram:
 $(x|y)$ means that x and y are joined in the binary tree. x and y can themselves be binary trees. What is the order we build the tree?

- A. $12(34)5 \rightarrow (12)(34)5 \rightarrow (((12)(34))5 \rightarrow (((12)(34))5))$.
- B. $12(34)5 \rightarrow 12((34)5) \rightarrow (12)((34)5) \rightarrow (12((34)5))$.
- C. $12(34)5 \rightarrow (12)(34)5 \rightarrow (12)((34)5) \rightarrow ((12)((34)5))$.
- D. $12(34)5 \rightarrow 12((34)5) \rightarrow 1(2((34)5)) \rightarrow (1(2((34)5)))$.

QUESTION II:

Consider the clustering problem given to the right where blue dots are observations and black circles are the initial position of four centroids denoted E, F, G and H used to cluster the data by k-means using Euclidean distances as dissimilarity. Upon convergence of the k-means algorithm which one of the following statements is wrong?

- A. Cluster formed by centroid F will be empty.
- B. Cluster formed by centroid E will contain 10 observations.
- C. Clusters formed by centroid H will contain 4 observations.
- D. Cluster formed by centroid G will contain 3 observations.



(Solution: QI: A, QII: B)

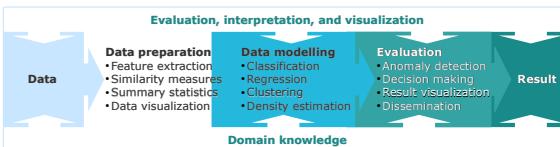
59 DTU Informatics, Technical University of Denmark

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

$$f(x+\Delta x) = \sum_{i=1}^n f_i(x) + \delta e^{i\pi} = \Theta + \Omega \int_a^x \delta e^{it} dt = [2.7182818284] \sum_{i=1}^n \delta e^{i\pi},$$

Data modeling framework

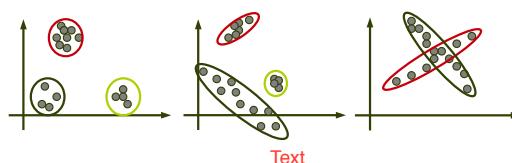


After today you should be able to:
 Explain the role of the parameters in the Gaussian Mixture Model (GMM) and how the parameters are updated using the EM-algorithm
 Explain why cross-validation can be used for GMM
 Describe the Apriori principle in association mining and explain how this can be used for efficient estimation of association rules.
 Calculate support and confidence of association rules

4 DTU Informatics, Technical University of Denmark

Group exercise

- What are the clusters below and what characterize each cluster?
- Is **k-means** well suited for modeling the clusters below?
 - Will it always find the optimum solution?
 - Can it model the sizes of the clusters?
 - Can it model the shape of the clusters?
 - How can we determine the number of clusters?



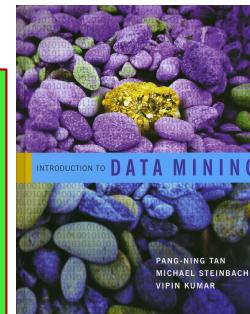
7 DTU Informatics, Technical University of Denmark

Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"
Section 9.2.2 + 6.1-6.3

Feedback Groups of the day:
 Sofie Knosgaard, Sarah Jørgensen,
 Sara Løansen
 Eleftherios-Marios Kotsonis-Tzannes, Mindaugas Venckus, Saad Eddin Al Orjany
 Johannes Sanders, Øystein Monsen, Tomasz Kamiński
 Louise Jørgensen, Ida Fillingsnes
 Anders Ulrik Eliesen, Rasmus Sten Andersen, Rasmus Bendtsen
 Løke Kristensen, Michael Nielsen
 Søren Trads Steen, Janus Nortoft Jensen

If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

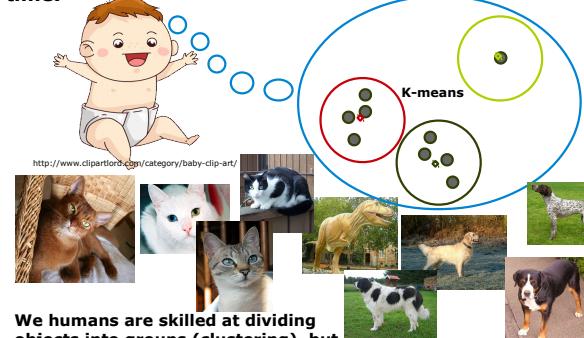


Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
Machine learning and data modelling in practice
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project

3 DTU Informatics, Technical University of Denmark

Imagine (again) you observe the world for the first time!

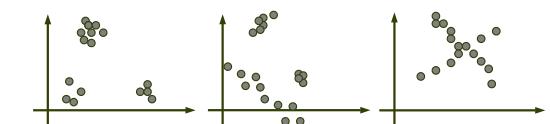


We humans are skilled at dividing objects into groups (clustering), but how do we make computers do the same?

5 DTU Informatics, Technical University of Denmark

Group exercise

- What are the clusters below and what characterize each cluster?
- Is **k-means** well suited for modeling the clusters below?
 - Will it always find the optimum solution?
 - Can it model the sizes of the clusters?
 - Can it model the shape of the clusters?
 - How can we determine the number of clusters?

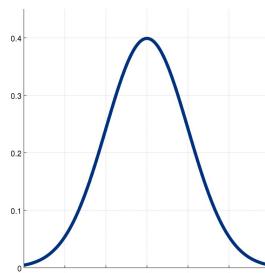


6 DTU Informatics, Technical University of Denmark

Normal distribution

- Probability density function describes the relative chance of a given value to occur
- Normal distribution characterized by
 - Mean
 - Variance

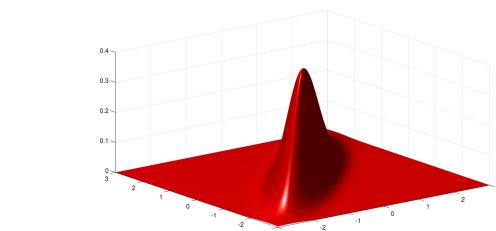
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



8 DTU Informatics, Technical University of Denmark

Multivariate Normal distribution

$$p(x) = \frac{1}{(2\pi)^{k/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right)$$



9 DTU Informatics, Technical University of Denmark

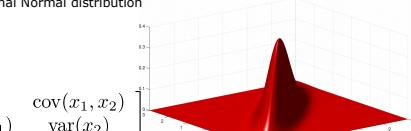
Multivariate Normal distribution

$$p(x) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

- Example: 2-dimensional Normal distribution

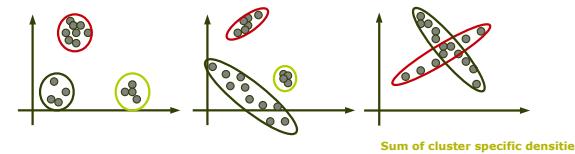
$$\mu = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{var}(x_2) \end{bmatrix}$$



10 DTU Informatics, Technical University of Denmark

The Gaussian Mixture Model (GMM)



- Different locations $\mu_{(k)}$
- Different shape $\Sigma_{(k)}$
- Different sizes w_k

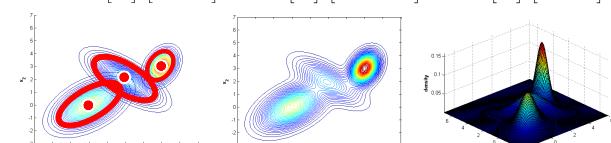
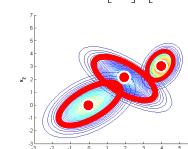
$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \mu_{(k)}, \Sigma_{(k)})$$

$$\text{s.t. } \sum_{k=1}^K w_k = 1, w_k \geq 0$$

11 DTU Informatics, Technical University of Denmark

GMM example

$$p(\mathbf{x}) = 0.5\mathcal{N}(\mathbf{x} | \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}) + 0.2\mathcal{N}(\mathbf{x} | \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & -0.7 \\ -0.7 & 1 \end{bmatrix}) + 0.3\mathcal{N}(\mathbf{x} | \begin{bmatrix} 4 \\ 3 \end{bmatrix}, \begin{bmatrix} 0.2 & 0.1 \\ 0.1 & 0.5 \end{bmatrix})$$



$\mu_{(k)}$: Cluster center (prototypical example in cluster)
 $\Sigma_{(k)}$: Shape of the cluster
 w_k : Relative size/density of the cluster

12 DTU Informatics, Technical University of Denmark



Group exercise

- Consider the Gaussian mixture model (GMM)

$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \mu_{(k)}, \Sigma_{(k)}) \quad \text{s.t. } \sum_{k=1}^K w_k = 1, w_k \geq 0$$

- What is the value of the integral?

$$\int p(\mathbf{x}) d\mathbf{x} = 1$$

Apply sum rule!

13 DTU Informatics, Technical University of Denmark

Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

- Expectation

- For each object, calculate the probability of belonging to each distribution

- Maximization

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change

E-step

$$p(z_n = k | \mathbf{x}_n) = \frac{w_k \mathcal{N}(\mathbf{x}_n | \mu_{(k)}, \Sigma_{(k)})}{\sum_{k=1}^K w_k \mathcal{N}(\mathbf{x}_n | \mu_{(k)}, \Sigma_{(k)})}$$

M-step

$$\begin{aligned} N_k &= \sum_{n=1}^N p(z_n = k | \mathbf{x}_n) \\ \mu_{(k)} &= \frac{1}{N_k} \sum_{n=1}^N \mathbf{x}_n p(z_n = k | \mathbf{x}_n) \\ w_k &= \frac{N_k}{N} \\ \Sigma_{(k)} &= \frac{1}{N_k} \sum_{n=1}^N (\mathbf{x}_n - \mu_{(k)}) (\mathbf{x}_n - \mu_{(k)})^\top p(z_n = k | \mathbf{x}_n) \end{aligned}$$

14 DTU Informatics, Technical University of Denmark

Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

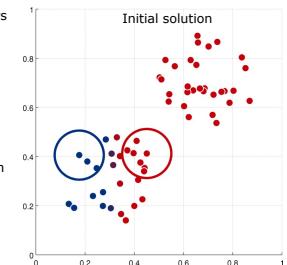
- Expectation

- For each object, calculate the probability of belonging to each distribution

- Maximization

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



15 DTU Informatics, Technical University of Denmark

Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

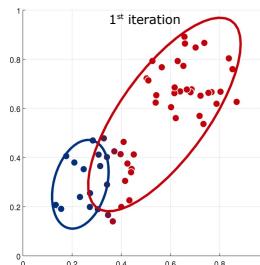
- Expectation

- For each object, calculate the probability of belonging to each distribution

- Maximization

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



16 DTU Informatics, Technical University of Denmark

Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

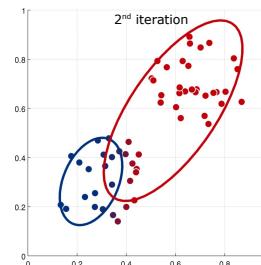
- Expectation

- For each object, calculate the probability of belonging to each distribution

- Maximization

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



17 DTU Informatics, Technical University of Denmark

Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

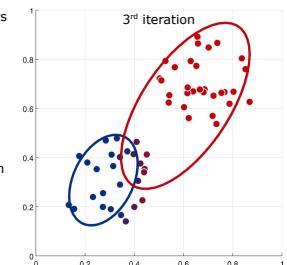
- Expectation

- For each object, calculate the probability of belonging to each distribution

- Maximization

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



18 DTU Informatics, Technical University of Denmark

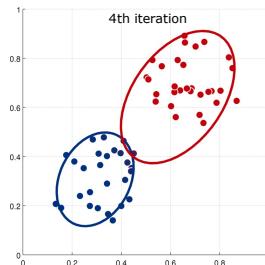
Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

- **Expectation**
 - For each object, calculate the probability of belonging to each distribution
- **Maximization**
 - For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



19 DTU Informatics, Technical University of Denmark

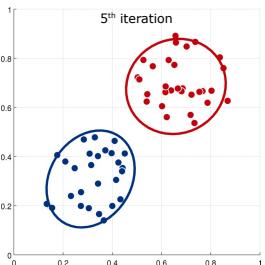
Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

- **Expectation**
 - For each object, calculate the probability of belonging to each distribution
- **Maximization**
 - For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



20 DTU Informatics, Technical University of Denmark

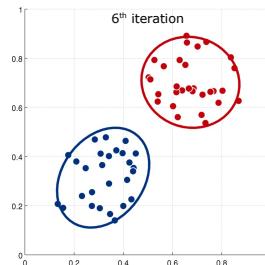
Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

- **Expectation**
 - For each object, calculate the probability of belonging to each distribution
- **Maximization**
 - For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change

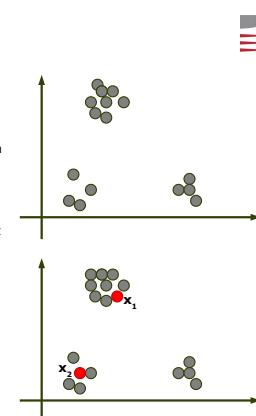


21 DTU Informatics, Technical University of Denmark



Group exercise

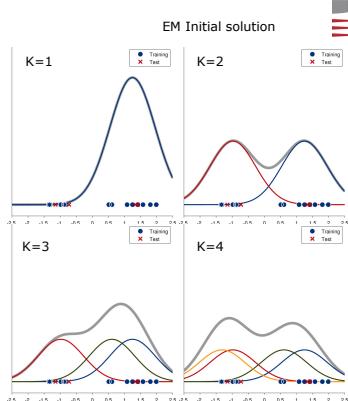
- Consider the data to the right with 16 observations.
 - What would ideally happen if we used a GMM with K=16 clusters to model the data?
- Imagine we have two **test observations**, denoted x_1 and x_2 (red points) that are not used for training.
 - What happens to $p(x_1)$ and $p(x_2)$ if we use K=3 and K=14 clusters?



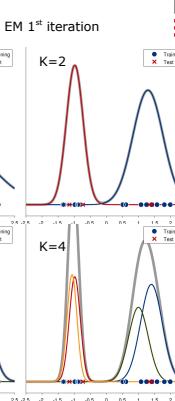
22 DTU Informatics, Technical University of Denmark

Mixture models

- Selecting complexity using crossvalidation



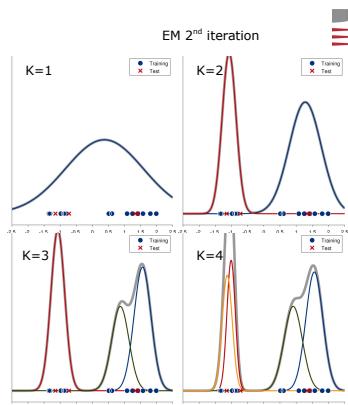
23 DTU Informatics, Technical University of Denmark



24 DTU Informatics, Technical University of Denmark

Mixture models

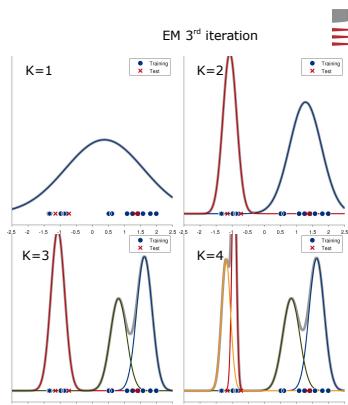
- Selecting complexity using crossvalidation



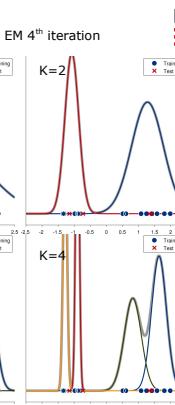
25 DTU Informatics, Technical University of Denmark

Mixture models

- Selecting complexity using crossvalidation



26 DTU Informatics, Technical University of Denmark



27 DTU Informatics, Technical University of Denmark

Mixture models

- Selecting complexity using crossvalidation

EM 5th iteration

K=1

K=2

K=3

K=4

The figure consists of three vertically stacked plots, each titled with a different value of K : $K=1$, $K=2$, and $K=4$. Each plot shows a set of data points (blue dots) and a fitted Gaussian mixture model (grey lines). The x-axis ranges from -2.5 to 2.5. The y-axis represents probability density.

- Plot K=1:** Shows a single broad Gaussian curve fitted to the data. The legend indicates blue dots for 'Training' and red 'x' marks for 'Test'. The curve is centered around 0.5.
- Plot K=2:** Shows two narrower Gaussian curves. The legend indicates blue dots for 'Training' and red 'x' marks for 'Test'. One peak is centered around -1.0 and the other around 1.0.
- Plot K=4:** Shows four very narrow Gaussian curves. The legend indicates blue dots for 'Training' and red 'x' marks for 'Test'. The peaks are centered around -1.5, -0.5, 0.5, and 1.5.

Each plot has a legend in the top right corner with a blue dot labeled 'Training' and a red 'x' labeled 'Test'.

DTU

Mixture models

- Selecting complexity using crossvalidation

Number of mixture components	Test log likelihood
1	-6.5
2	-4.5
3	-4.8
4	-14.5

30 DTU Informatics, Technical University of Denmark

K-means versus GMM

K-means

- No guarantee of optimal solution
- Does not model shape of clusters
- Does not model the size of clusters
- Difficult to assess the number of clusters to use particularly when there is no ground truth

Gaussian mixture model (GMM)

- No guarantee of optimal solution (even more local minima issues due to the additional model parameters)
- Models shape of cluster as ellipsoid
- Models the size of clusters
- Possible to estimate the number of components by cross-validation

<http://siguemeandosalvajecito.es/comics/merchandising/excultura-de-superman-vs-muhammad-ali/>

31 DTU Informatics, Technical University of Denmark

K-means versus GMM

K-means

- No guarantee of optimal solution
- Does not model shape of clusters
- Does not model the size of clusters
- Difficult to assess the number of clusters to use particularly when there is no ground truth

Gaussian mixture model (GMM)

- No guarantee of optimal solution (even more local minima issues due to the additional model parameters)
- Models shape of cluster as ellipsoid
- Models the size of clusters
- Possible to estimate the number of components by cross-validation

<http://sagecomunicacion.com.es/comics/merchandising/escultura-de-superman-vs-muhammad-ali/>

32

DTU Informatics, Technical University of Denmark

Association Mining

Association mining Agarwal

Ca. 24 900 resultater (0,67 sek.)

ip Seg efter resultater på Dansk alone. Du kan sende direkte sprogetstilling i Indstillingen for Scholia

finding association rules between sets of items in large databases

[Agarwal - 1994] ACM SIGMOD Record, 1994 - dl.acm.org

Abstract We introduce a new database of customer transactions in which a transaction consists of items purchased by a customer in a visit. We present an efficient algorithm that generates itemsets of size k from such a database. The algorithm can incorporate ...
[Agarwal - 1994] ACM SIGMOD Record, 1994 - dl.acm.org

» ip Fast algorithms for mining association rules

[Agarwal - 1994] Proc. 20th Int. Conf. Very Large Data Bases, VLDB 1994 - dl.acm.org

Abstract We introduce a new database of customer transactions in which a transaction consists of items purchased by a customer in a visit. We present two new algorithms for solving this problem that ...
[Agarwal - 1994] ACM SIGMOD Record, 1994 - dl.acm.org

» ip Mining quantitative association rules in large relational tables

[Biskup, R. Agarwal - 1995] ACM SIGMOD Record, 1995 - dl.acm.org

Abstract We introduce a new database of customer transactions in which a transaction consists of itemsets of size k obtained from large relational tables containing both quantitative and categorical attributes. An example of such an association rule is "10% of married people between age 50 and 60 have at least 2 cars". We deal ...
[Biskup, R. Agarwal - 1995] ACM SIGMOD Record, 1995 - dl.acm.org

» ip Mining generalized association rules

[Lifshitz, R. Agarwal - 1995 - www.csail.mit.edu/com

Abstract We introduce the problem of **mining generalized association rules**. Given a

large database of transactions where each transaction consists of a set of items and a

33 DTU Informatics, Technical University of Denmark

Association rule discovery: Definition

- Given a set of **records**
 - Each containing a number of **items from a set**
- **Goal:** Produce dependency rules
 - Predict the occurrence of an item based on occurrences of other items

34 DTU Informatics, Technical University of Denmark

Association rule discovery: Example

Market basket analysis

<u>Training set</u>	<u>Rules discovered</u>
1. {Bread, Soda, Milk} 2. {Beer, Bread} 3. {Beer, Soda, Diaper, Milk} 4. {Beer, Bread, Diaper, Milk} 5. {Soda, Diaper, Milk}	{Milk} \rightarrow {Soda} {Diaper, Milk} \rightarrow {Beer}

35

DTU Informatics, Technical University of Denmark

Market basket data

- Representation as

Transaction table

ID	Items
1	Bread, Soda, Milk
2	Beer, Bread
3	Beer, Soda, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Soda, Diaper, Milk

Data matrix

ID	Bread	Soda	Milk	Beer	Diaper
1	1	1	1	0	0
2	1	0	0	1	0
3	0	1	1	1	1
4	1	0	1	1	1
5	0	1	1	0	1

36

DTU Informatics, Technical University of Denmark

Association analysis

- Itemset**

- For example {Bread, Soda, Milk}, {Milk, Diaper}, {}

- Support** for an itemset X

- Percentage of transactions that contain X

- Association rule**

- Expression of the form: $X \rightarrow Y$
where X and Y are disjoint item sets

- Support** for an association rule $X \rightarrow Y$

- Percentage of transactions that contain $X \cup Y$
 $s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N} = P(X, Y)$

- Confidence** for an association rule $X \rightarrow Y$

- Percentage of transactions containing X that also contain Y

$$c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)} = \frac{P(Y|X)}{P(X)} = P(Y|X)$$

37 DTU Informatics, Technical University of Denmark

sigma = number of times



Group exercise

- What is the **support** for
- {Bread}
- {Milk, Diaper}
- What is the **support and confidence** for
- {Bread} \rightarrow {Milk}
- {} \rightarrow {Milk}
- Find an **itemset** with
- 0% support
- 100% support
- Find an **association rule** with
- 0% confidence
- 100% confidence

ID	Bread	Soda	Milk	Beer	Diaper
1	1	1	1	0	0
2	1	0	0	1	0
3	0	1	1	1	1
4	1	0	1	1	1
5	0	1	1	0	1

Itemset

- For example {Bread, Soda, Milk}, {Milk, Diaper}, {}
- Support** for an itemset X

- Percentage of transactions that contain X

Association rule

- Expression of the form: $X \rightarrow Y$

where X and Y are disjoint item sets

Support for an association rule $X \rightarrow Y$

- Percentage of transactions that contain $X \cup Y$

Confidence for an association rule $X \rightarrow Y$

- Percentage of transactions containing X that also contain Y

Association rule mining

- Find all association rules that have

- **Support \geq minsup**

- **Confidence \geq minconf**

- Approach**

- Frequent itemset generation**

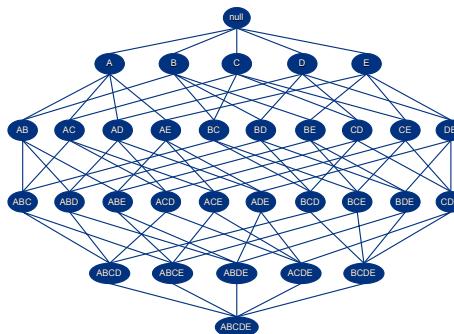
- Generate a list of all **itemsets** with **Support \geq minsup**

- Association rule generation**

- Generate all **association rules** with **Confidence \geq minconf**

39 DTU Informatics, Technical University of Denmark

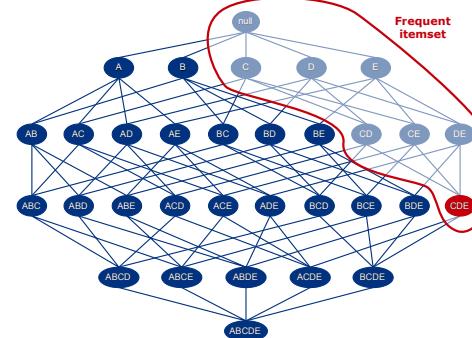
Frequent itemset generation



How many different itemsets can be created for a problem with a total of D items?

40 DTU Informatics, Technical University of Denmark

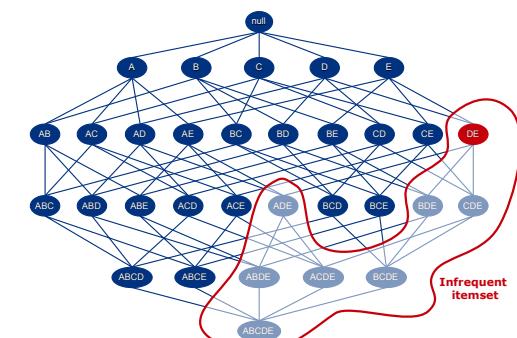
Frequent itemset generation



If an itemset is frequent, then all of its subsets must also be frequent

41 DTU Informatics, Technical University of Denmark

Frequent itemset generation



If an itemset is infrequent, then all of its supersets must also be infrequent

42 DTU Informatics, Technical University of Denmark

Group exercise

How many possible itemsets are there in the market basket below?
What are all the itemsets with support $\geq 35\%$?

	Juice	Milk	Beer	Cheese	Chocolate	Yoghurt	Sugar	Flour	Egg	Wine
Customer 1	0	0	1	0	0	0	1	0	1	0
Customer 2	1	1	0	0	0	1	0	1	1	1
Customer 3	0	1	0	1	1	0	0	0	0	1
Customer 4	1	1	0	0	0	1	0	0	0	1
Customer 5	1	0	1	0	0	0	0	0	1	0
Customer 6	1	0	0	0	0	0	0	0	1	0
Customer 7	1	1	0	0	1	1	0	0	0	1
Customer 8	0	1	0	1	0	1	1	1	0	1
Customer 9	1	1	0	0	1	1	0	0	0	0
Customer 10	0	0	1	0	0	1	0	0	1	1

Group exercise

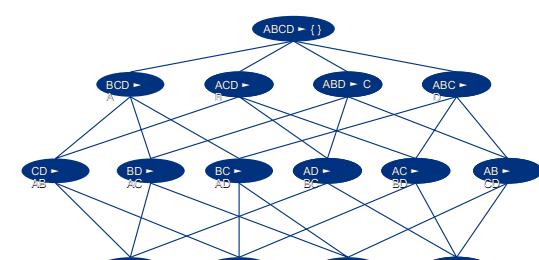
How many possible itemsets are there in the market basket below?
What are all the itemsets with support $\geq 35\%$?

	Juice	Milk	Beer	Cheese	Chocolate	Yoghurt	Sugar	Flour	Egg	Wine
Customer 1	0	0	1	0	0	0	1	0	1	0
Customer 2	1	1	0	0	0	1	0	1	1	1
Customer 3	0	1	0	1	1	0	0	0	0	1
Customer 4	1	1	0	0	0	1	0	0	0	1
Customer 5	1	0	1	0	0	0	0	0	1	0
Customer 6	1	0	0	0	0	0	0	0	1	0
Customer 7	1	1	0	0	1	1	0	0	0	1
Customer 8	0	1	0	1	0	0	1	1	0	1
Customer 9	1	1	0	0	1	1	0	0	0	0
Customer 10	0	0	1	0	0	1	0	0	1	1

$2^{10} = 1024$ itemsets, itemsets with support $\geq 35\%$ are:
 - {Juice}, {Milk}, {Yoghurt}, {Egg}, {Wine}
 - {Juice, Milk}, {Juice, Yoghurt}, {Milk, Yoghurt}, {Wine, Milk}, {Wine, Yoghurt}
 - {Juice, Milk, Yoghurt}

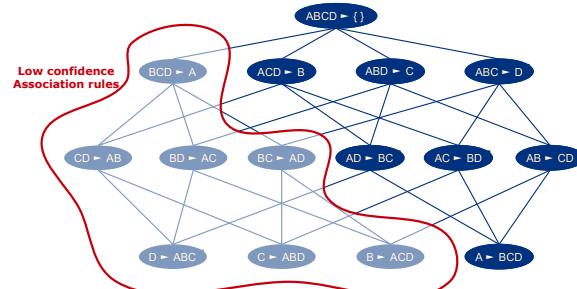
43 DTU Informatics, Technical University of Denmark

Association rule generation



44 DTU Informatics, Technical University of Denmark

Association rule generation



46 DTU Informatics, Technical University of Denmark

Results for market basket example

Itemset	Support	Association rule	Support	Confidence
Milk	80%	{ } → Milk	80%	80%
Bread	60%	Soda → Milk	60%	100%
Soda	60%	Diaper → Milk	60%	100%
Beer	60%	Soda, Diaper → Milk	40%	100%
Diaper	60%	Beer, Diaper → Milk	40%	100%
Diaper Milk	60%	Beer, Milk → Diaper	40%	100%
Soda Milk	60%			
Bread Beer	40%			
Bread Milk	40%			
Soda Diaper	40%			
Beer Diaper	40%			
Beer Milk	40%			
Soda Diaper Milk	40%			
Beer Diaper Milk	40%			

ID	Bread	Soda	Milk	Beer	Diaper
1	1	1	1	0	0
2	1	0	0	1	0
3	0	1	1	1	1
4	1	0	1	1	1
5	0	1	1	0	1

Binarize data according to percentiles

AttributeNames=	Attribute 1	Attribute 2	Attribute 3	AttributeBinned=					
				0-50 %	50-100 %	5-10 %	33.33-50 %	66.67-100 %	0-50 %
0.3689	0.9827	0.6999		1	0	0	1	0	1
0.4607	0.7302	0.6385		0	1	0	1	0	1
0.9816	0.3439	0.0336		0	1	0	1	0	1
0.1564	0.5841	0.0688		1	0	1	0	1	0
0.8555	0.1078	0.3196		0	1	0	0	1	0
0.6448	0.9063	0.5309		0	1	0	0	1	0
0.3763	0.8797	0.6544		1	0	0	1	0	1
0.1909	0.8178	0.4076		1	0	0	0	1	0
0.4283	0.2607	0.8200		0	1	1	0	0	1
X=	0.4820	0.7184	Xbinary=	0	1	0	1	0	1
	0.1206	0.0225		1	0	1	0	0	1
	0.5895	0.4253		0	1	0	1	1	0
	0.2262	0.3127		1	1	0	0	0	0
	0.3846	0.1615		0	0	1	0	0	0
	0.5830	0.1788		1	0	1	0	0	1
	0.2518	0.4229		1	0	1	0	0	1
	0.2904	0.0942		0	1	0	1	0	0
	0.6171	0.5985		1	1	0	0	0	1
	0.2653	0.2665		1	0	0	1	0	1
	0.8244	0.6959		0	1	0	1	1	0

49 DTU Informatics, Technical University of Denmark

Further reading (not required)

Association Mining

- Rakesh Agrawal and Ramakrishnan Srikant "Fast Algorithms for Mining Association Rules", Proc. 20th Int. Conf. Very Large Data Bases, 1994
- Rakesh Agrawal, Tomasz Imielinski and Arun Swami "Mining association rules between sets of items in large databases" Proceedings of the 1993 ACM SIGMOD international conference on Management of data

Gaussian Mixture Model

- Christopher M. Bishop "Pattern Recognition and Machine Learning", Chapter 9, Springer 2006



http://commons.wikimedia.org/wiki/File:Old_book_bindings.jpg

50 DTU Informatics, Technical University of Denmark

Group exercise

- How can we do association mining for continuous data?

Attribute 1	Attribute 2	Attribute 3
0.3689	0.9827	0.6999
0.4607	0.7302	0.6385
0.9816	0.3439	0.0336
0.1564	0.5841	0.0688
0.8555	0.1078	0.3196
0.6448	0.9063	0.5309
0.3763	0.8797	0.6544
0.1909	0.8178	0.4076
0.4283	0.2607	0.8200
X=	0.4820	0.5944
	0.1206	0.0225
	0.5895	0.4253
	0.2262	0.3127
	0.3846	0.1615
	0.5830	0.1788
	0.2518	0.4229
	0.2904	0.0942
	0.6171	0.5985
	0.2653	0.2665
	0.8244	0.6959

48 DTU Informatics, Technical University of Denmark

Exam question examples

- QI: Consider the market basket given to the right where 10 customers have purchased various types of fruits in a grocery store. Disregarding the empty set, what are all frequent itemsets with support $\geq 40\%$
- A: {Orange}, {Apple}, {Banana}, {Pineapple}
- B: {Orange}, {Apple}, {Banana}, {Pineapple}, {Orange, Apple}, {Orange, Pineapple}
- C: {Orange}, {Apple}, {Banana}, {Pineapple}, {Orange, Apple}, {Orange, Pineapple}, {Apple, Pineapple}, {Orange, Apple, Pineapple}
- D: {Orange}, {Apple}, {Banana}, {Pineapple}, {Orange, Apple}, {Orange, Pineapple}, {Apple, Pineapple}, {Orange, Apple, Pineapple}

QII: What is the confidence of the decision rules {Orange, Apple} \rightarrow {Pear}

- A: 25%
B: 50%
C: 75%
D: 100%

Correct answer QI: B, QII: B

51 DTU Informatics, Technical University of Denmark

02450 Introduction to machine learning and data mining

DTU Informatics
Department of Informatics and Mathematical Modeling

$$f(x+\Delta x) = \sum_{i=1}^n f_i(x)$$

$$\Delta \int_a^b \Theta + \Omega \int \delta e^{ix} = [2.7182818284] \sum_{i=1}^n \epsilon^i$$

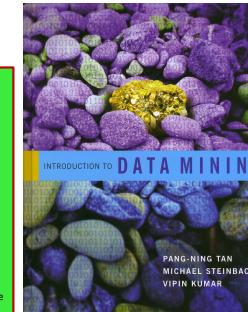
Reading material

Tan, Steinbach and Kumar
"Introduction to Data Mining"

Section 10.1-10.4

Feedback Groups of the day:
 Wen Hsin Li, Robert Lyck, Viivi Halla-aho
 Oliver Naaby, Rasmus Olsen
 Simon Hemmingsen, Thomas Hammerbak
 Henrik Thirstrup, Mette Marie Nielsen
 Thomas Kristensen
 Nicolai Pedersen, Jeppe Thagaard, Mathias Brade, Christian Kragh
 Oscar Petersen, Søren Vorne Nielsen
 Andreas Jacobsen, Asger Anker Sørensen, Phong Le Trung

If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

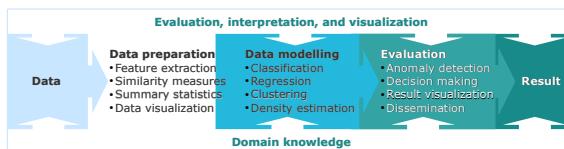


Lecture schedule

1. Introduction
(Tan 1.1-1.4)
Data: Feature extraction and visualization
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
Supervised learning: Classification and regression
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
Machine learning and data modelling in practice
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
Unsupervised learning: Clustering and density est.
9. K-means and hierarchical clustering
(Tan 8.1-8.3+8.5.7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
- 11. Density estimation and anomaly detection**
(Tan 10.1-10.4)
Machine learning and data modelling in practice
12. Putting it all together: Summary and overview
13. Mini project

3 DTU Informatics, Technical University of Denmark

Data modeling framework

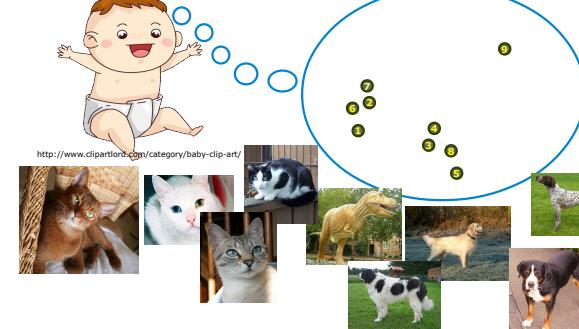


After today you should be able to:

Understand and apply a variety of outlier/anomaly detection approaches, including density estimation approaches and proximity-based techniques.

4 DTU Informatics, Technical University of Denmark

Imagine (yet again) you observe the world for the first time!



How do we detect anomalous objects
(i.e., the dinosaur in the world of cats and dogs)

5 DTU Informatics, Technical University of Denmark

Anomaly detection: Definition

- Given a collection of data objects
 - Each object has associated a number of features
- Detect which objects **deviate from normal** behaviour

6 DTU Informatics, Technical University of Denmark

Anomaly detection: Example

- Credit card **fraud detection**
 - Recognize dubious credit card transactions based on the transaction history of the card holder
- Network **intrusion detection**
 - Detect hacker attacks, web crawlers etc.
- Ecological disturbances**
 - Detect hurricanes, floods droughts, heat waves and fires
- Health and medicine monitoring**
 - Detect abnormal behaviour in populations and patients
- Fault detection in industry systems**
 - Detect when a wind turbine performs poorly due to ice coating on blades
- Detection of **outliers** in data measurements
 - Remove erroneous measurements due to misreading from an instrument

7 DTU Informatics, Technical University of Denmark



Group exercise

- Come up with **your own definition** of an outlier / anomaly
- How can we detect outliers using some of the methods you have already learned in the course?

8 DTU Informatics, Technical University of Denmark



Group exercise

- Come up with **your own definition** of an outlier / anomaly
- How can we detect outliers using some of the methods you have already learned in the course?

Hawkins' definition of an outlier: An outlier is an observation that differs so much from other observations as to arouse suspicion that it was generated by a different mechanism



Probabilistic definition of an outlier: An outlier is an object that has a low probability with respect to a probability distribution model of the data.

9 DTU Informatics, Technical University of Denmark

Data example I: Cats, Dogs and Dinosaurs

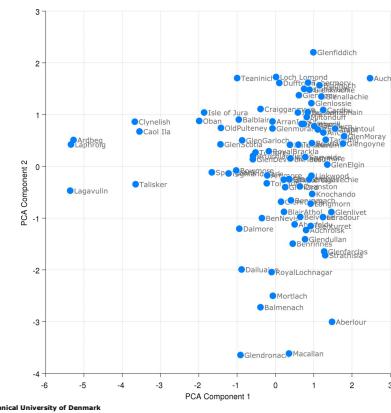


Data example II: Whisky

- 86 types of Scotch whisky
- Human ratings 1-5
- 12 taste categories
 - body, sweetness, smoky, medicinal, tobacco, honey, spicy, winey, nutty, malty, fruity, floral

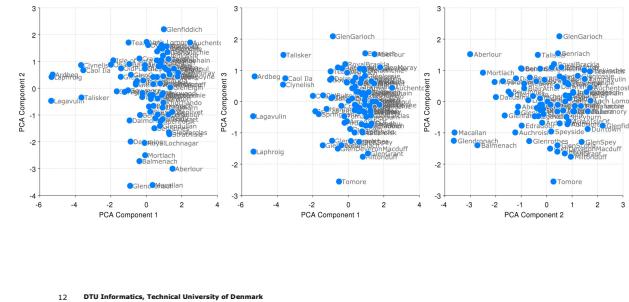
10 DTU Informatics, Technical University of Denmark

PCA plot



11 DTU Informatics, Technical University of Denmark

PCA plot

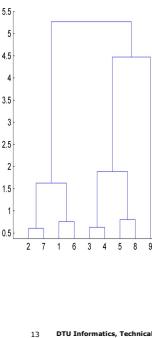


12 DTU Informatics, Technical University of Denmark

Dendrogram

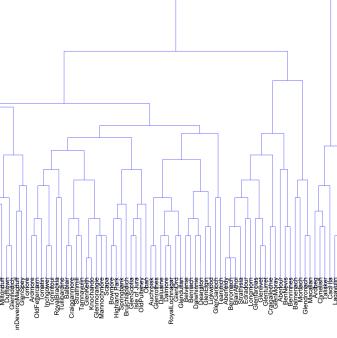
- Dendograms can be used to visualize relative distances between the observations

Data I: Cats, Dogs and Dinosaurs



13 DTU Informatics, Technical University of Denmark

Data II: Whisky



Approaches to anomaly detection

• Density-based techniques

- Estimate the density of data objects
- Outliers are:
 - Data objects in low density area

Approaches we will consider:

- Univariate normal distribution
- Kernel density estimation

• Proximity-based techniques

- Measure the distance between data objects
- Outliers are:
 - Data objects far from the other data objects

Approaches we will consider:

- Mahalanobis distance to center of data
- Distance to Kth nearest neighbour
- Inverse average distance to K nearest neighbours (KNN density)
- Average relative KNN density

Normal distribution

- Map attribute to standard Normal variable

$$z = \frac{x - \mu}{\sigma}$$

- Choose a threshold

c	p(z >c)
1.0	0.3173
1.5	0.1336
2.0	0.0455
2.5	0.0124
3.0	0.0027
3.5	0.0005
4.0	0.0001

Medicinal

$$\mu = 1.55$$

$$\sigma = 0.99$$

Normal distribution

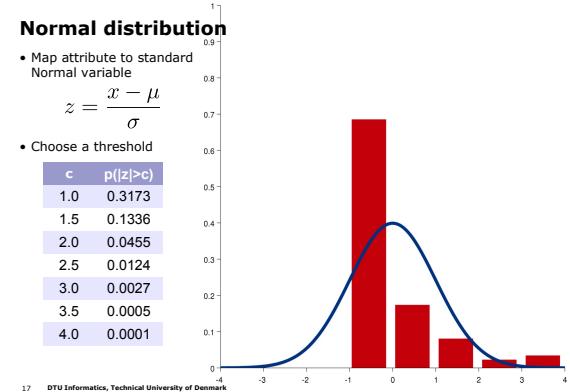
- Map attribute to standard Normal variable

$$z = \frac{x - \mu}{\sigma}$$

- Choose a threshold

c	p(z >c)
1.0	0.3173
1.5	0.1336
2.0	0.0455
2.5	0.0124
3.0	0.0027
3.5	0.0005
4.0	0.0001

Medicinal: z-score



17 DTU Informatics, Technical University of Denmark

Normal distribution

- Map attribute to standard Normal variable

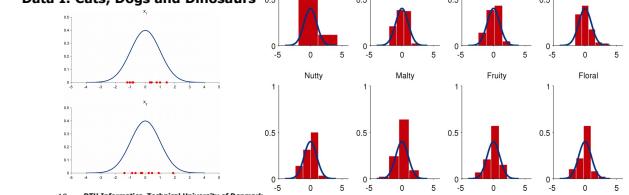
- Choose a threshold

$$z = \frac{x - \mu}{\sigma}$$

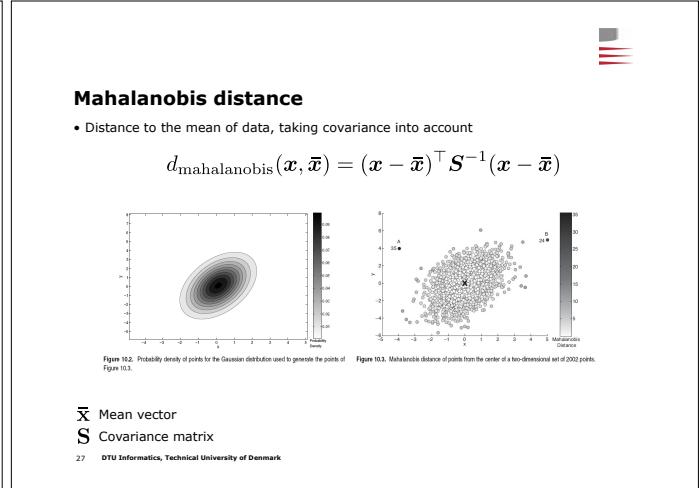
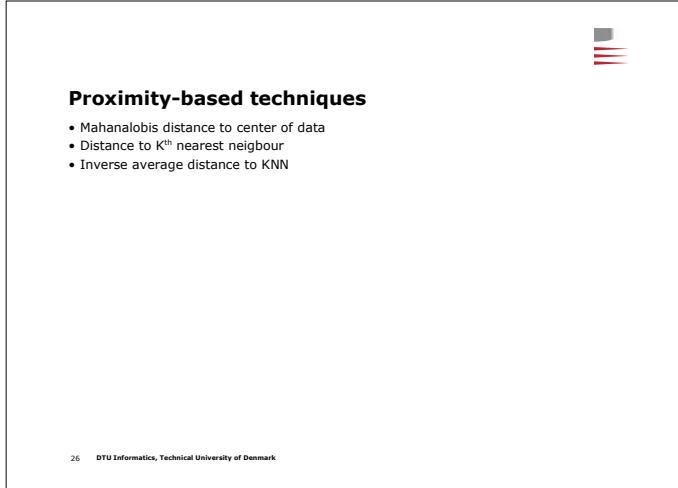
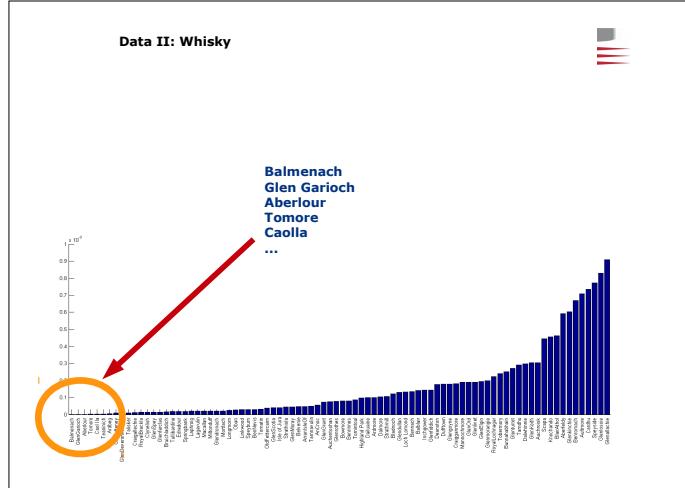
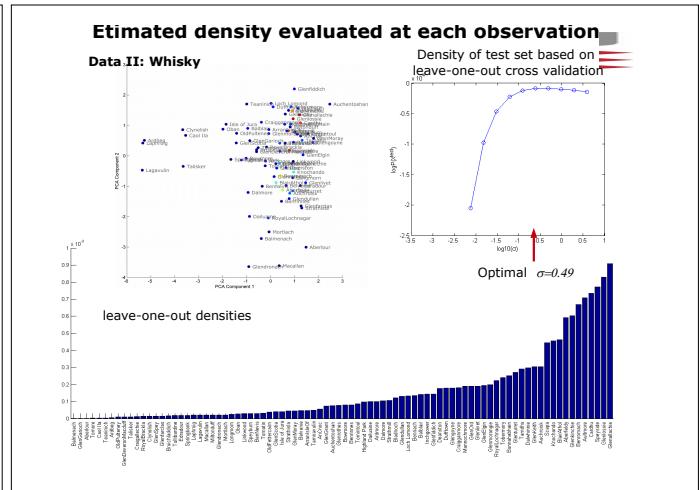
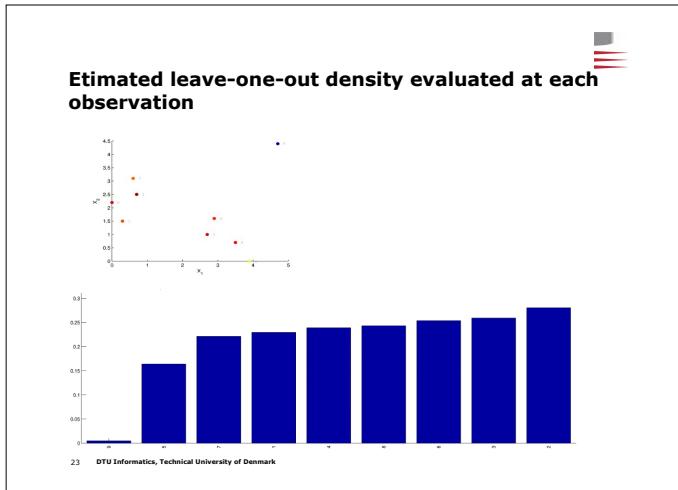
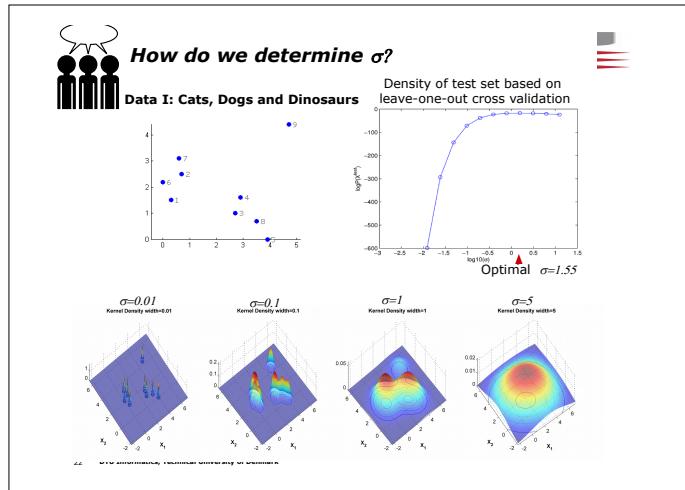
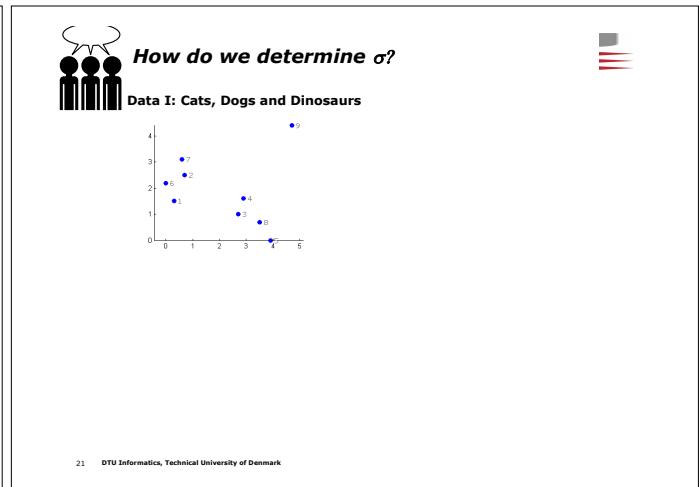
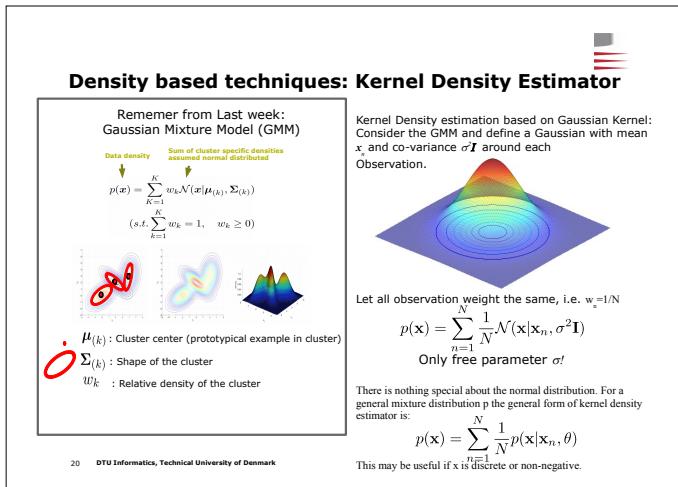
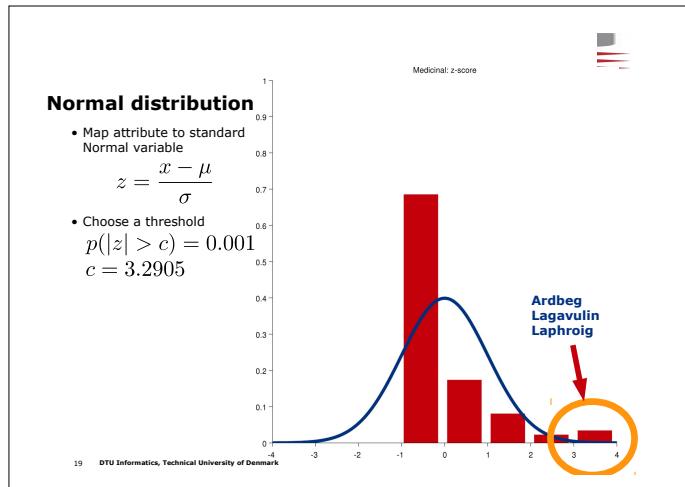
$$p(|z| > c) = 0.001$$

$$c = 3.2905$$

Data I: Cats, Dogs and Dinosaurs



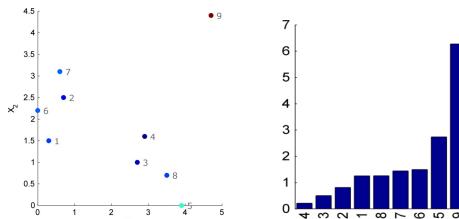
18 DTU Informatics, Technical University of Denmark



Mahalanobis distance

- Distance to the mean of data, taking covariance into account

Data I: Cats , dogs and dinosaurs

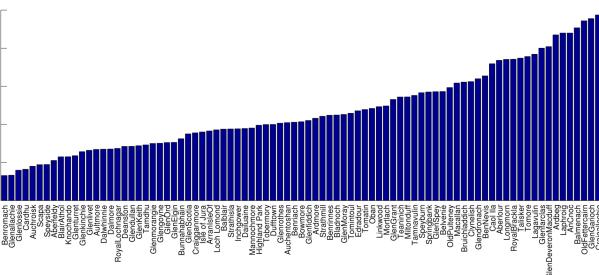


28 DTU Informatics, Technical University of Denmark

Mahalanobis distance

- Distance to the mean of data, taking covariance into account

Data II: Whisky

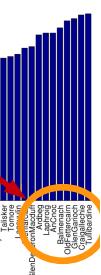


29 DTU Informatics, Technical University of Denmark

Mahalanobis distance

- Distance to the mean of data, taking covariance into account

Tullibardine
Craigellachie
Glen Garioch
Old Fettercairn
Balmenach
...



30 DTU Informatics, Technical University of Denmark

Distance to k-nearest neighbor

- Measure the distance to the k'th nearest neighbor

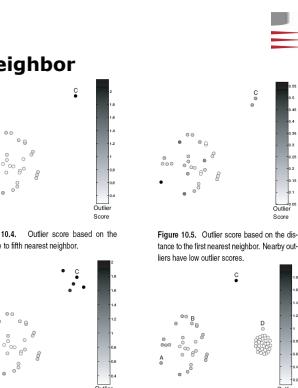


Figure 10.4. Outlier score based on the distance to 10th nearest neighbor.

Figure 10.5. Outlier score based on the distance to the 1st nearest neighbor. Nearby outliers have low outlier scores.

Figure 10.6. Outlier score based on distance to the 10th nearest neighbor. A small cluster becomes an outlier.

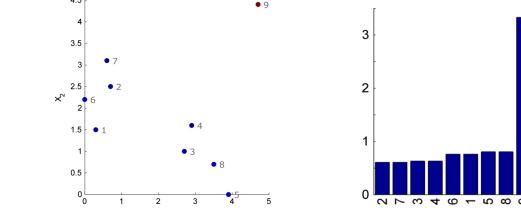
Figure 10.7. Outlier score based on the distance to the 10th nearest neighbor. Clusters of differing density.

31 DTU Informatics, Technical University of Denmark

Distance to k-nearest neighbor

- Measure the distance to the 1st nearest neighbor

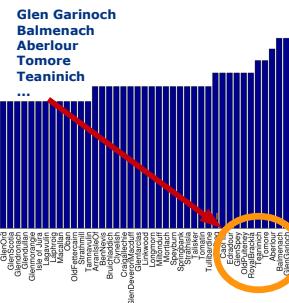
Data I: Cats , dogs and dinosaurs



32 DTU Informatics, Technical University of Denmark

Distance to k-nearest neighbor

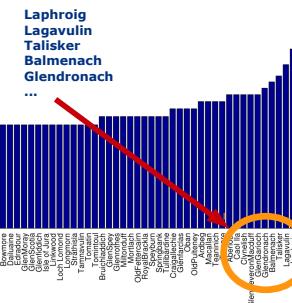
- Measure the distance to the 1st nearest neighbor



34 DTU Informatics, Technical University of Denmark

Distance to k-nearest neighbor

- Measure the distance to the 5th nearest neighbor

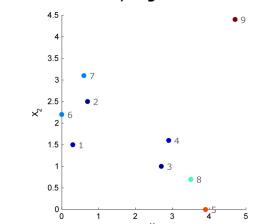


35 DTU Informatics, Technical University of Denmark

Distance to kth Nearest neighbour

- Measure the distance to the 5th nearest neighbor

Data I: Cats , dogs and dinosaurs



33 DTU Informatics, Technical University of Denmark

Inverse distance density estimation

- Distance based measure of density

- Density is inverse proportional to average distance to k nearest neighbors
- Density is low if nearest neighbors are far away

$$\text{density}(\mathbf{x}, k) = \left(\frac{1}{k} \sum_{\mathbf{y} \in N(\mathbf{x}, k)} \text{distance}(\mathbf{x}, \mathbf{y}) \right)^{-1}$$

- Relative density

- Density compared to density at nearest neighbors

$$\text{average relative density}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{\mathbf{y} \in N(\mathbf{x}, k)} \text{density}(\mathbf{y}, k)}$$

$N(\mathbf{x}, k)$ The set of k nearest neighbors

36 DTU Informatics, Technical University of Denmark

Consider the pairwise distance matrix given to the left. What is the density and average relative density of the first observation for k=3?

$d(\mathbf{x}_i, \mathbf{x}_j)$	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7
\mathbf{x}_1	0	2.5	2.4	4.0	0.8	0.6	3.3
\mathbf{x}_2	2.5	0	0.6	1.6	2.9	3.0	1.1
\mathbf{x}_3	2.4	0.6	0	1.9	3.0	2.7	1.0
\mathbf{x}_4	4.0	1.6	1.9	0	4.5	4.6	3.8
\mathbf{x}_5	0.8	2.9	3.0	4.5	0	1.1	3.9
\mathbf{x}_6	0.6	3.0	2.7	4.6	1.1	0	3.8
\mathbf{x}_7	3.3	1.1	1.0	3.8	3.9	3.8	0

$$\text{density}(\mathbf{x}, k) = \left(\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{distance}(\mathbf{x}, y) \right)^{-1}$$

$$\text{average relative density}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{density}(y, k)}$$

Consider the pairwise distance matrix given to the left. What is the density and average relative density of the first observation for k=3?

$d(\mathbf{x}_i, \mathbf{x}_j)$	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7
\mathbf{x}_1	0	2.5	2.4	4.0	0.8	0.6	3.3
\mathbf{x}_2	2.5	0	0.6	1.6	2.9	3.0	1.1
\mathbf{x}_3	2.4	0.6	0	1.9	3.0	2.7	1.0
\mathbf{x}_4	4.0	1.6	1.9	0	4.5	4.6	3.8
\mathbf{x}_5	0.8	2.9	3.0	4.5	0	1.1	3.9
\mathbf{x}_6	0.6	3.0	2.7	4.6	1.1	0	3.8
\mathbf{x}_7	3.3	1.1	1.0	3.8	3.9	3.8	0

$$\text{density}(\mathbf{x}, k) = \left(\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{distance}(\mathbf{x}, y) \right)^{-1}$$

$$\text{average relative density}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{density}(y, k)}$$

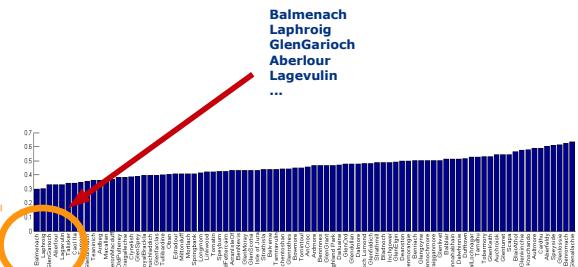
$$\begin{aligned}\text{density}(\mathbf{x}_1, 3) &= [1/3(0.6+0.8+2.4)]^{-1}=3/3.8 \\ \text{density}(\mathbf{x}_2, 3) &= [1/3(0.6+1.1+2.7)]^{-1}=3/4.4 \\ \text{density}(\mathbf{x}_3, 3) &= [1/3(0.8+1.1+2.9)]^{-1}=3/4.8 \\ \text{density}(\mathbf{x}_4, 3) &= [1/3(0.6+1.0+1.9)]^{-1}=3/3.5\end{aligned}$$

$$\text{Av. Rel. Density}(\mathbf{x}_1, 3) = [3/3.8]/[1/3(3/4.4+3/4.8+3/3.5)]=1.0945$$

37 DTU Informatics, Technical University of Denmark

Inverse distance density estimation

- KNN density (5 nearest neighbors)



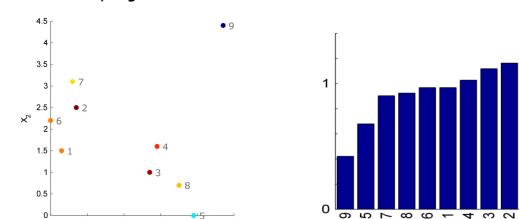
Results using different methods

- Univariate Normal distribution
- Distance to 5th nearest neighbor
- Kernel Density Estimation
- Mahalanobis distance
- Distance to nearest neighbor
- Common: Balmenach, Glen Garioch, Laphroig, Aberlour, Tomore, Lagavulin, Craigallechie

Average Relative density

- Average Relative KNN density (5 nearest neighbors)

Data I: Cats , dogs and dinosaurs



Example of exam questions

Q1: What is the average relative density for observation 2 (i.e. \mathbf{x}_2) for k=2 nearest neighbours?

- A:1/5
B:3/10
C:/10
D:1

$d(\mathbf{x}_i, \mathbf{x}_j)$	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5
\mathbf{x}_1	0	2.0	0.2	0.9	0.2
\mathbf{x}_2	2.0	0	1.5	0.5	2.0
\mathbf{x}_3	0.2	1.5	0	1.2	1.4
\mathbf{x}_4	0.9	0.5	1.2	0	1.0
\mathbf{x}_5	0.2	2.0	1.4	1.0	0

$$\text{density}(\mathbf{x}, k) = \left(\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{distance}(\mathbf{x}, y) \right)^{-1}$$

$$\text{avg. rel. den.}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{density}(y, k)}$$

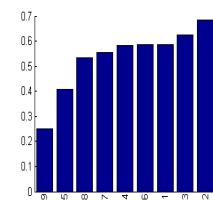
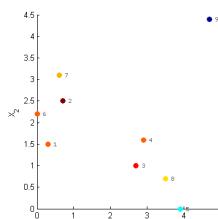
Common: Balmenach, Glen Garioch, Laphroig, Aberlour, Tomore, Lagavulin, Craigallechie

44 DTU Informatics, Technical University of Denmark

Inverse distance density estimation

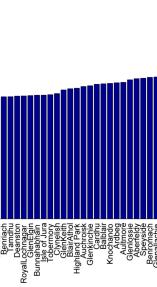
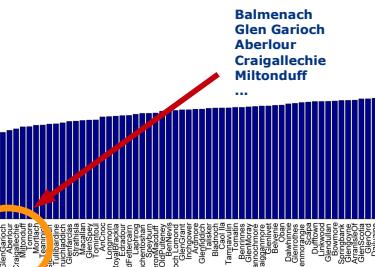
- KNN density (5 nearest neighbors)

Data I: Cats , dogs and dinosaurs



Average relative density

- Average relative KNN density (5 nearest neighbors)



43 DTU Informatics, Technical University of Denmark

©