# Knowledge Discovery and Data Mining 1 (VO) (707.003)

Denis Helic

KTI, TU Graz

Oct 1, 2015

# Lecturer

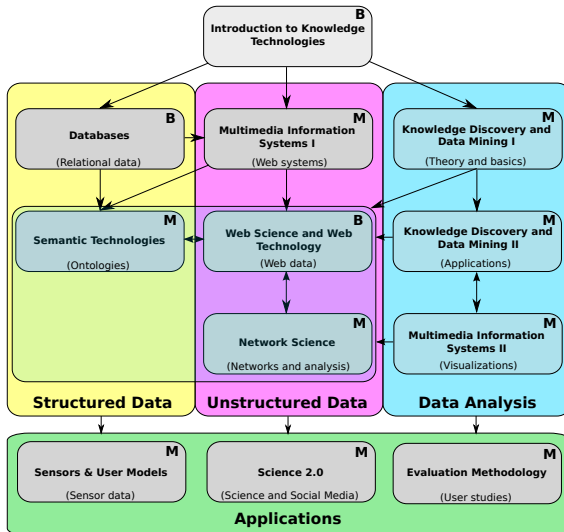| | |
|---:|:---:|
| Name: | Denis Helic |
| Office: | IWT, Inffeldgasse 13, 5th Floor, Room 070 |
| Office hours: | Tuesday from 12 til 13 |
| Phone: | +43-316/873-30610 |
| email: | dhelic@tugraz.at |

# Lecturer

|  |  |
|---:|---:|
| Name: | Roman Kern |
| Office: | IWT, Inffeldgasse 13, 6th Floor, Room 072 |
| Office hours: | By appointment |
| Phone: | +43-316/873-30860 |
| email: | rkern@know-center.at |

# Language

- Lectures in English
- Communication in German/English
- If in German: please informally (Du)!
- Examination: German/English

# Outline

1. Welcome and Introduction

2. Course Organization

3. Motivation

4. Course Overview

5. Course Highlights

# Teaching @ KTI



+ Projects, Bachelor Thesis, Master Projects, Master Thesis, PhD Thesis

# Course context

- Knowledge Discovery and Data Mining 1 (VO) (707.003)
- Obligatory course Master Software Development and Business (1st Semester)
- Elective course in subject catalogue "Knowledge Technologies"
- Computer Science, Telematics

# Course context

- Knowledge Discovery and Data Mining 1 (KU) (707.004)
- Free course
- An add-on for the theoretical part
- Highly suggested

# Goals of the course

- The overall goal of KDDM and related courses is to learn how to discover patterns and models in data. Discovered patterns need to be:
  - (i) Valid: hold for new data with high probability
  - (ii) Useful: we can base further actions on them
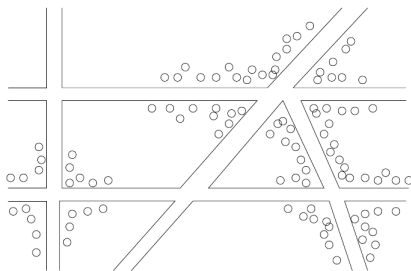  - (iii) Unexpected: non-obvious
  - (iv) Understandable: humans can intepret them

# Goals of the course: patterns example

**1854 Broad Street cholera outbreak**

Extracting clusters of cholera outbreak in the city of London in 1854. The cases clusterd around some intersections of roads in London. These had contaminated water wells.

- http://en.wikipedia.org/wiki/1854_Broad_Street_cholera_outbreak

# Goals of the course: patterns example

# Goals of the course

- Specific goals of this course are to learn about **two** (three) cruicial elements of that discovery:

    (i) **Tools: mathematical tools such as probability theory, linear algebra, information theory, and statistical inference**
    (ii) Infrastructure: models of computation for large data (KDDM2 has also a practical part)
    (iii) **Process: steps that are needed to discover patterns**

- I assume here that you already know how to program and develop software

# Goals of the course

- Student goals: to pass the examination
- Bonus goal for all: to have fun!

# Course Calendar

- 01.10.2015: Course organization, Introduction and Motivation (Denis)
- 08.10.2015: Preprocessing (Roman)
- 15.10.2015: Feature Extraction (Roman)
- 22.10.2015: **Partial Exam 1** / Project presentations (KU)
- 29.10.2015: Feature Engineering (Roman)

# Course Calendar

- 05.11.2015: Data Matrices (Denis)
- 12.11.2015: Principal Component Analysis (Denis)
- 19.11.2015: SVD and Latent Semantic Analysis (Denis)
- 26.11.2015: Recommender Systems: Matrix Factorization (Denis)
- 03.12.2015: **Partial Exam 2** / Project presentations (KU)

# Course Calendar

- 10.12.2015: Classification (Denis)
- 17.12.2015: Clustering (Roman)
- 14.01.2016: Pattern Mining (Roman)
- 21.01.2016: **Partial Exam 3** / Project presentations (KU)
- 28.01.2016: Examination

# Course Logistics

- Course website:
  http://kti.tugraz.at/staff/denis/courses/kddm1
- Slides will be made available on the course website
- Additional readings, references, links, etc. also on the website
- We expect that you have basic knowledge in **probability theory** and **linear algebra**
- To freshen the knowledge in the first midterm there will be questions from these topics!
- Please check the homepage for details
- As a side note: we also expect that you know how to **program** (relevant for the practical part)

# Grading

- Partial examinations within the class
- Written examination at the end of the class
- Two additional examinations in winter semester
- Three examinations in summer semester
- Examination material: contents of slides
- In class we will discuss certain types of exam questions

# Partial Examinations

- 3 written examinations
- In the beginning of a lecture: 30 minutes
- Each partial examination 2 questions
- Difficulty adjusted to solve both problems in approx. 25 minutes
- Max 15 points for each question
- Total points: 90

# Examination

- Written examination 90 minutes
- 4 questions
- Difficulty adjusted to solve all four in approx. 70-80 minutes
- Max 20 points for each question
- Total points: 80

# Grading

- 0-40 points: 5
- 41-50 points: 4
- 51-60 points: 3
- 61-70 points: 2
- 71-80 points: 1

# KU Organization

- KU organization on 22.10.2015

# Questions?

- Raise them now ($+1$ $+1$)
- Ask after the lecture ($+1$)
- Visit me in the office hours ($+1$)
- Send me an e-mail ($\pm 1$)
- As a side note: you should(!) interrupt me immediately ($+1$ $+1$ $+1$) and ask any question you might have during the lecture

# How much information is being produced?

- Study at Berkley: `http://www2.sims.berkeley.edu/research/projects/how-much-info-2003/`
- The World Wide Web contains about 170 terabytes of information on its surface
- This is seventeen times the size of the Library of Congress print collections
- Instant messaging generates five billion messages a day
- This is 274 terabytes a year

# How much information is being produced?

- Email generates about 400,000 terabytes of new information each year worldwide
- P2P file exchange on the Internet is growing rapidly
- That was 2003, what do we have today?
- In 1993: 100.000G transferred over the Internet per year

# How much information is being produced?

- Email generates about 400,000 terabytes of new information each year worldwide
- P2P file exchange on the Internet is growing rapidly
- That was 2003, what do we have today?
- In 1993: 100.000G transferred over the Internet per year
- In 2008: 100.000G transferred over the Internet in a *second*

# How much information is being produced?

- Data, data everywhere:
  http://www.economist.com/node/15557443

# How much information is being produced?

- We are producing more data than we are able to store
- We should extract and describe useful data
- We can store that data
- We should also try to predict future data

# Knowledge discovery

- Discover patterns and models in data. These patterns need to be:
  - (i) Valid: hold for new data with high probability
  - (ii) Useful: we can base further actions on them
  - (iii) Unexpected: non-obvious
  - (iv) Understandable: humans can interpret them

# Examples

# Examples

- PageRank
  - (i) Valid: holds for all new data
  - (ii) Useful: we base rankings of search results on PageRank
  - (iii) Unexpected: non-obvious and non-trivial to calculate
  - (iv) Understandable: popularity, importance of Web pages, etc.

# Examples

# Examples

- Knowledge Graph

  (i) Valid: holds for new data with high probability
  (ii) Useful: users can explore connections between concepts
  (iii) Unexpected: non-obvious and non-trivial
  (iv) Understandable: related, similar, etc

# Examples

- https://www.google.at/#q=albert+einstein
- http://www.youtube.com/watch?v=mmQl6VGvX-c

# Examples

# Examples

- TechMeme
    - (i) Valid: holds for new data with high probability
    - (ii) Useful: users can find and explore the news about technology
    - (iii) Unexpected: non-obvious and non-trivial
    - (iv) Understandable: summaries of tech news, etc.

# Examples

# Examples

- Amazon product recommendations
  - (i) Valid: holds for new data with high probability
  - (ii) Useful: users can find and explore new products
  - (iii) Unexpected: non-obvious and non-trivial
  - (iv) Understandable: related articles, etc.

# Examples



Figure 9: Earthquake location estimation based on tweets. Balloons show the tweets on the earthquake. The cross shows the earthquake center. Red represents early tweets; blue represents later tweets.

# Examples



Figure 10: Typhoon trajectory estimation based on tweets.

# Examples

- Twitter earthquake and typhoon prediction
    - (i) Valid: holds for new data with high probability
    - (ii) Useful: can save lives
    - (iii) Unexpected: non-obvious and non-trivial
    - (iv) Understandable: trajectories of typhoons, positions of earthquakes

# Examples

# Examples

- Graph of ideas http://zoom.it/l3dq
    - (i) Valid: holds for new data with high probability
    - (ii) Useful: development of ideas
    - (iii) Unexpected: non-obvious and non-trivial
    - (iv) Understandable: who influenced whom, etc

# Knowledge discovery vs. data mining

- Knowledge discovery refers to the entire process, of which knowledge is the end-product
- It is iterative and interactive
- Data mining refers to a specific step in this process
- It is the step consisting of applying data analysis and discovery algorithms that produce a particular enumeration of patterns over data
- Additional steps are necessary to ensure that the process produces useful knowledge

# Steps in the knowledge discovery process

1. Developing an **understanding of the application domain** and the relevant prior knowledge and identifying the goal of the KDD process from the customers viewpoint

2. Creating a target data set: **selecting a data set** or focusing on a subset of variables or data samples on which discovery is to be performed

3. **Data cleaning and preprocessing**: basic operations such as the removal of noise. If appropriate collecting the necessary information to model or account for noise, deciding on strategies for handling missing data fields, accounting for time sequence information and known changes

# Steps in the knowledge discovery process

4. Data reduction and projection: **finding useful features to represent the data** depending on the goal of the task. Using dimensionality reduction or transformation methods to reduce the effective number of variables under consideration or to find invariant representations for the data

5. Matching the goals of the KDD process step to a particular **data mining method** e.g. summarization, classification, regression, clustering, etc

6. Choosing the **data mining algorithms**: selecting methods to be used for searching for patterns in the data. This includes deciding which models and parameters may be appropriate e.g. models for categorical data are different than models on vectors over the reals. Matching a particular data mining method with the overall criteria of the KDD process e.g. the enduser may be more interested in understanding the model than its predictive capabilities

# Steps in the knowledge discovery process

7. **Data mining searching for patterns** of interest in a particular representational form or a set of such representations, classification rules or trees, regression, clustering and so forth. The user can significantly aid the data mining method by correctly performing the preceding steps

8. **Interpreting mined patterns**: possibly return to any of the steps for further iteration. This step can also involve visualization of the extracted patterns, models or visualization of the data given the extracted models

9. **Consolidating discovered knowledge**: incorporating this knowledge into another system for further action or simply documenting it and reporting it to interested parties. This also includes checking for and resolving potential conflicts with previously believed or extracted knowledge
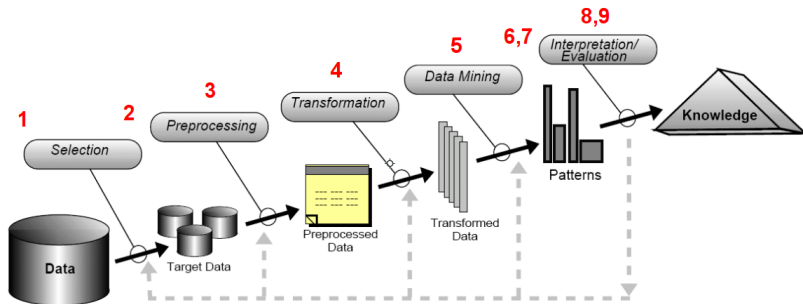
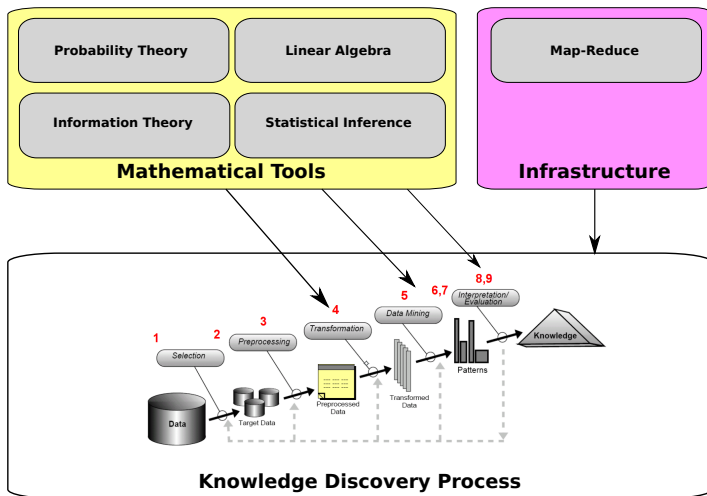# Steps in the knowledge discovery process

**Reading!**

Knowledge Discovery and Data Mining: Towards a Unifying Framework (1996) Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth

# Steps in the knowledge discovery process

# Big picture: KDDM



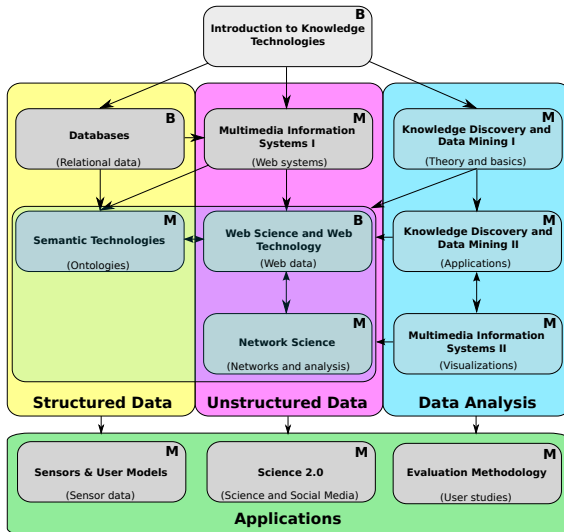**Knowledge Discovery Process**

# Teaching @ KTI and KDDM1

- KDDM1: Basics, theory and KDD process until step 8 (no visualization, no interpretation). We will only shortly mention graph-based data mining
- KDDM2: Implementation and practice of the theory from KDDM1
- MMIS2: KDD process steps 8 and 9 (interpretation and visualization)
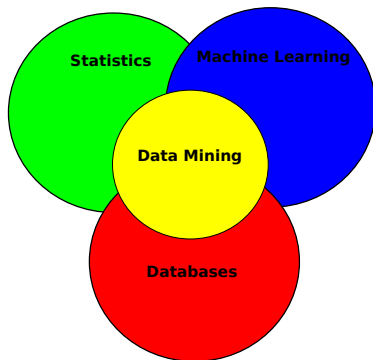- Network Science: graph and networks mining

# Teaching @ KTI



+ Projects, Bachelor Thesis, Master Projects, Master Thesis, PhD Thesis

# Data mining and other fields

- Data mining overlaps with

  (i) Databases: Large-scale data, simple queries
  (ii) Machine learning: Small data, complex models, model parameters
  (iii) Statistics: Theory, predictive models, no algorithms

- Data mining: **Algorithms**, simple and predictive models, large-scale data

# Data mining and other fields

# KDDM1: some thoughts

- Big data, data science, . . .
- Many buzzwords
- But in the end:
  1. You need to know how to program and how to develop software systems
  2. You need to understand the math behind it: linear algebra, probability, statistics
- If you know both of these you are in the top 10% of developers in the field ;)

# KDDM1: some highlights

- Bayesian inference: update your prior knowledge with new evidence
- Recommender systems: decompose the matrix and find out what your users really like :)
- PCA: reduce the dimensionality of the data
- LSA: analyze relationships between set of documents with linear algebra