


**DATABASE
SYSTEMS
GROUP**

Ludwig-Maximilians-Universität München
Institut für Informatik
Lehr- und Forschungseinheit für Datenbanksysteme



LMU

Knowledge Discovery in Databases II


Summer Semester 2017

Lecture 1: Introduction and outlook

Lectures : Prof. Dr. Peer Kröger, Yifeng Lu
Tutorials: Yifeng Lu


Script © 2015, 2017 Eirini Ntoutsis, Matthias Schubert, Arthur Zimek, Peer Kröger, Yifeng Lu

[http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_\(KDD_II\)](http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_(KDD_II))



**DATABASE
SYSTEMS
GROUP**

Course organization




LMU


- **Time and location**
 - Lectures: Thursday, **09:00-11:30**, room B 101 (Oettingenstr. 67)
 - Tutorial: Monday, 14:00-16:00, room A U115 (HGB)
 - Tutorial: Monday, 16:00-18:00, room A U115 (HGB)
 - All information and news can be found at:
[http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_\(KDD_II\)](http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_(KDD_II))
- **Exam**
 - Written exam, 90 min
 - 6 ECTS points
 - Registration for the written exam through UniWorX

Knowledge Discovery in Databases II: Introduction and overview

2




Chapter overview




- Knowledge Discovery in Databases, Big Data and Data Science
- Data Mining with Vectorized Data (Recap KDD I)
- Topics of KDD II
- Literature and supplementary materials

Knowledge Discovery in Databases II: Introduction and overview


3



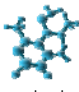
Motivation



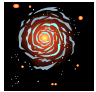
- Large amounts of data in multiple applications




connection data




molecule
process data



telescope data



transaction data



Web data/
click streams

...

• Manual analysis is infeasible

➔ **Knowledge Discovery in Databases and Data Mining**


Goals

- Descriptive modeling: Explains the characteristics and behavior of observed data
- Predictive modeling: Predicts the behavior of new data based on some model


Important: The extracted models/patterns don't have to apply to 100 % of the cases.

Knowledge Discovery in Databases II: Introduction and overview

4



What is KDD?




*Knowledge Discovery in Databases (KDD) is the **nontrivial process** of identifying **valid, novel, potentially useful, and ultimately understandable patterns in data.***

[Fayyad, Piatetsky-Shapiro, and Smyth 1996]


Remarks:

- *nontrivial*: it is not just the avg
- *valid*: to a certain degree the discovered patterns should also hold for new, previously unseen problem instances
- *novel*: at least to the system and preferable to the user
- *potentially useful*: they should lead to some benefit to the user or task
- *ultimately understandable*: the end user should be able to interpret the patterns either immediately or after some postprocessing

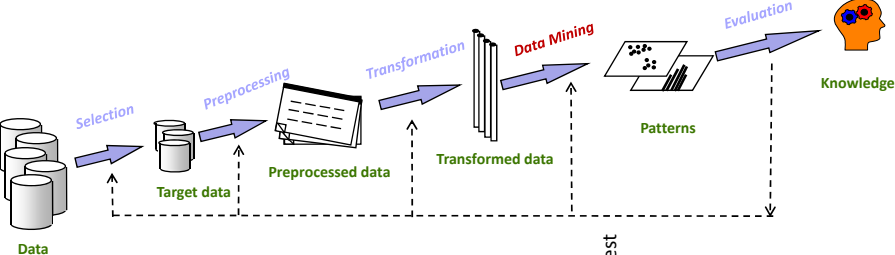
Knowledge Discovery in Databases II: Introduction and overview
5



The KDD process



[Fayyad, Piatetsky-Shapiro & Smyth, 1996]



Selection:

- Select a relevant dataset or focus on a subset of a dataset
- File / DB

Preprocessing/Cleaning:

- Integration of data from different data sources
- Noise removal
- Missing values

Transformation:

- Select useful features
- Feature transformation/discretization
- Dimensionality reduction


Data Mining:

- Search for patterns of interest

Evaluation:


- Evaluate patterns based on interestingness measures
- Statistical validation of the models

Knowledge Discovery in Databases II: Introduction and overview
6



**DATABASE
SYSTEMS
GROUP**

KDD landscape today




LMU

- Internet
- Internet of things
- Data intensive science / eScience
- Big data
- Data science
- ...


Knowledge Discovery in Databases II: Introduction and overview

7



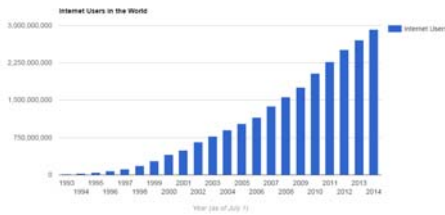
**DATABASE
SYSTEMS
GROUP**

Internet




LMU

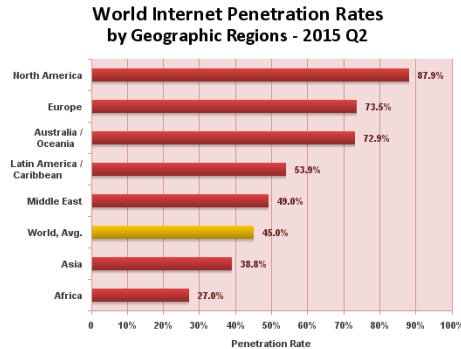
- Internet users (Source: <http://www.internetlivestats.com/internet-users/>)



Internet Users in the World



Web 2.0: A world of opinions




**World Internet Penetration Rates
by Geographic Regions - 2015 Q2**

Region	Penetration Rate
North America	87.9%
Europe	73.5%
Australia / Oceania	72.9%
Latin America / Caribbean	53.9%
Middle East	49.0%
World, Avg.	45.0%
Asia	38.8%
Africa	27.6%

Source: Internet World Stats - www.internetworldstats.com/stats.htm
 Penetration Rates are based on a world population of 7,260,621,118 and 3,270,490,584 estimated Internet users on June 30, 2015.
 Copyright © 2015, Miniwatts Marketing Group


Knowledge Discovery in Databases II: Introduction and overview

8



DATABASE
SYSTEMS
GROUP

Internet of Things



LMU

- The Internet of Things (IoT) is the network of physical objects or "things" embedded with electronics, software, sensors, and network connectivity, which enables these objects to collect and exchange data.

Source: https://en.wikipedia.org/wiki/Internet_of_Things




Image source: <http://tinyurl.com/prtfqxf>


During 2008, the number of things connected to the internet surpassed the number of people on earth... By 2020 there will be 50 billion ... vs 7.3 billion people (2015).

These things are everything, smartphones, tablets, refrigerators cattle.

Source: <http://blogs.cisco.com/diversity/the-internet-of-things-infographic>


Knowledge Discovery in Databases II: Introduction and overview

9



DATABASE
SYSTEMS
GROUP


The Fourth Paradigm: Data Intensive Science 1/2




LMU


Science Paradigms


- Thousand years ago:
science was **empirical**
describing natural phenomena
- Last few hundred years:
theoretical branch
using models, generalizations
- Last few decades:
a **computational** branch
simulating complex phenomena
- Today: **data exploration** (eScience)
unify theory, experiment, and simulation
 - Data captured by instruments or generated by simulator
 - Processed by software
 - Information/knowledge stored in computer
 - Scientist analyzes database/files using data management and statistics



$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G \rho}{3} - K \frac{c^2}{a^2}$$








Slide from: http://research.microsoft.com/en-us/um/people/gray/talks/nrc-cstb_escience.ppt


Knowledge Discovery in Databases II: Introduction and overview

10



DATABASE
SYSTEMS
GROUP

The Fourth Paradigm: Data Intensive Science 2/2



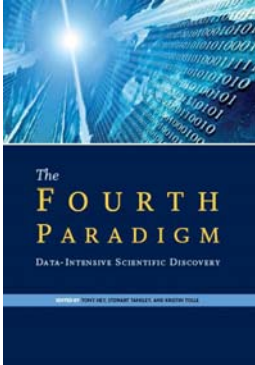
LMU

“Increasingly, scientific breakthroughs will be powered by advanced computing capabilities that help researchers manipulate and explore massive datasets.”

-The Fourth Paradigm – Microsoft


Examples of e-science applications:

- Earth and environment
- Health and wellbeing
 - E.g., The Human Genome Project (HGP)
- Citizen science
- Scholarly communication
- Basic science
 - E.g., CERN




Knowledge Discovery in Databases II: Introduction and overview

11



DATABASE
SYSTEMS
GROUP

Big Data



LMU

“Big data is a broad term for datasets so large or complex that traditional data processing applications are inadequate. Challenges include analysis, capture, data curation, search, sharing, storage, transfer, visualization, and information privacy.”

Source: https://en.wikipedia.org/wiki/Big_data

Capturing the value of big data:

- 300 billion USD potential value for the north American health system per year
- 250 billion Euro potential value for the public sector in Europe per year
- 600 billion USD potential value through the use for location based services

Source: McKinsey Report “Big data: The next frontier for innovation, competition, and productivity”, June 2011:


Data Scientist: The sexiest job of the 21st century:

“The United States alone faces a shortage of 140,000 to 190,000 people with deep analytical skills as well as 1.5 million managers and analysts to analyze big data and make decisions based on their findings.”


Source: <http://tinyurl.com/cplxu6p>

Knowledge Discovery in Databases II: Introduction and overview

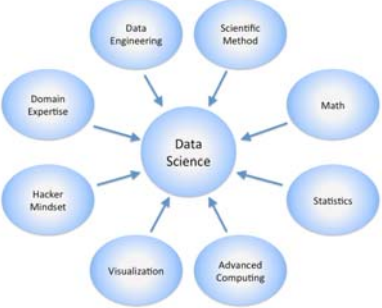
12



Data Science




- Science of managing and analyzing data to generate knowledge
- Very similar to KDD, but
 - Data Science is broader in its topics. (result representation, actions..)
 - Integrates all scientific directions being concerned with data analyses and knowledge representation.
 - New computational paradigms and hardware systems.




Wrap up: Many sciences worked on the topics for last decades. Data Science can be seen as an umbrella comprising all of these areas.

Knowledge Discovery in Databases II: Introduction and overview

13



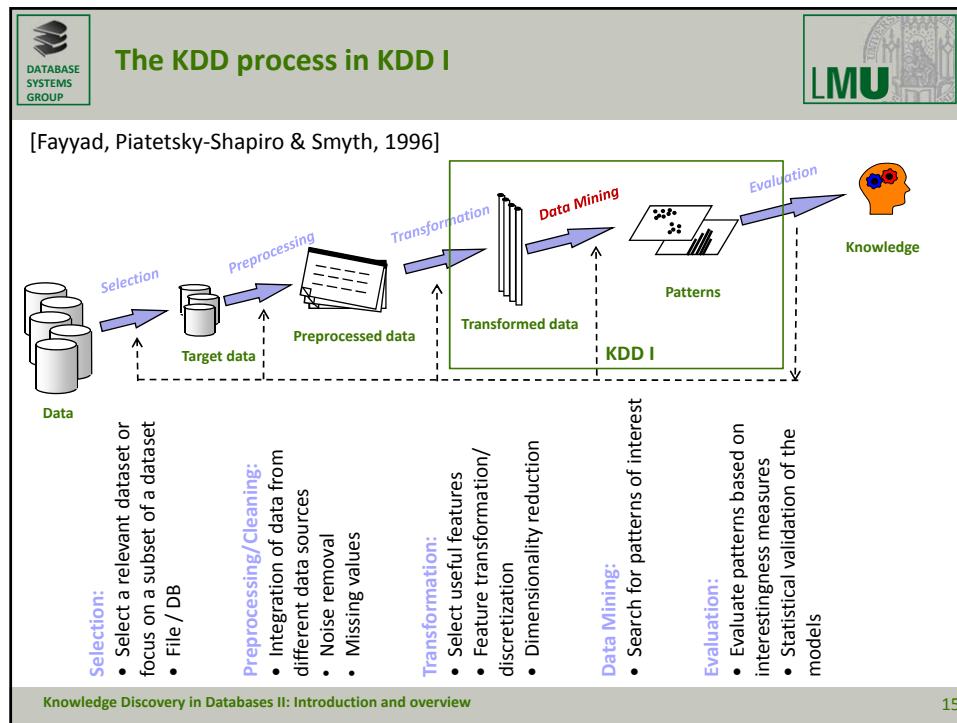
Chapter overview



- Knowledge Discovery in Databases, Big Data and Data Science
- Data Mining with Vectorized Data (Recap KDD I)
- Topics of KDD II
- Literature and supplementary materials

Knowledge Discovery in Databases II: Introduction and overview

14




KDD I topics

- Clustering
 - partitioning, agglomerative, density-based, grid-based
- Classification
 - NN-classification, Bayesian classifiers, SVMs, decision trees
- Association rule mining and frequent pattern mining
 - Apriori, FP-growth, FI, MFI, CFI
- Regression
- Outlier Detection


Most of the methods covered by KDD I assume the data to be a set of *feature vectors*

Knowledge Discovery in Databases II: Introduction and overview 16



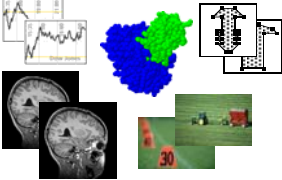
DATABASE
SYSTEMS
GROUP

Feature Vectors/Feature Transformation




LMU

- Isn't this assumption to work with feature vectors extremely limiting?
 - Well ...
- The concept of „Feature Transformation“ (Similarity modelling)
 - Extract characteristic (**numeric**) features from each object
 - Each object is represented as a high-dimensional (feature) vector
 - Characteristic features: similar vectors indicate similar objects



Data Space

Feature Transformation



Histogramms

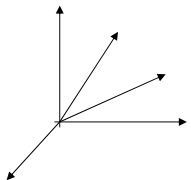
Moment Invariants

Covering

Sectoring

Fourier Transformation


...



Feature Space


Knowledge Discovery in Databases II: Introduction and overview

17



DATABASE
SYSTEMS
GROUP

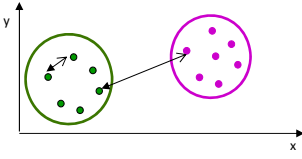
Clustering 1/3

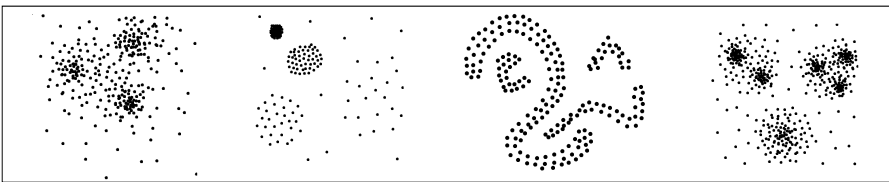


LMU

- **Goal:**


Group objects into groups so that the objects belonging in the same group are similar (high intra-cluster similarity), whereas objects in different groups are different (low inter-cluster similarity)
- Similarity/ distance function
- Unsupervised learning
- What is a good clustering ???






Knowledge Discovery in Databases II: Introduction and overview

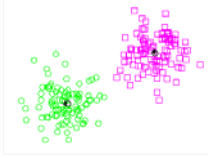
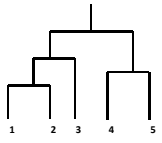

18



Clustering 2/3




- Partitioning clustering:
 - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
 - Typical methods: k-means, k-medoids, CLARANS
- Hierarchical clustering:
 - Create a hierarchical decomposition of the set of data (or objects) using some criterion
 - Typical methods: Diana, Agnes, BIRCH, ROCK, CHAMELEON
- Density-based clustering:
 - Based on connectivity and density functions
 - Typical methods: DBSCAN, OPTICS






Knowledge Discovery in Databases II: Introduction and overview

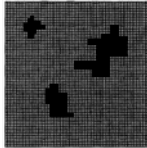
19



Clustering 3/3



- Grid-based clustering:
 - based on a multiple-level granularity structure
 - Typical methods: STING, CLIQUE
- Model-based clustering:
 - A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
 - Typical methods: EM, SOM, COBWEB
- User-guided or constraint-based clustering:
 - Clustering by considering user-specified or application-specific constraints
 - Typical methods: COD (obstacles), constrained clustering



Knowledge Discovery in Databases II: Introduction and overview

20



Classification 1/3



Given:

- a dataset of instances $D=\{t_1, t_2, \dots, t_n\}$ and
- a set of classes $C=\{c_1, \dots, c_k\}$

the classification problem is to define a mapping $f:D \rightarrow C$ where each instance t_i in D is assigned to one class c_j .

ID	Alter	Autotyp	Risk
1	23	Familie	high
2	17	Sport	high
3	43	Sport	high
4	68	Familie	low
5	32	LKW	low

Training set

A simple classifier:

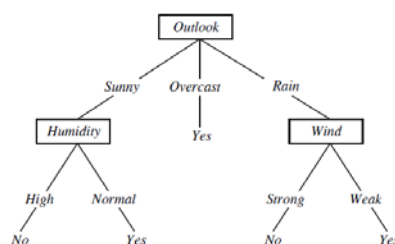
- if Alter > 50 then Risk= low;
- if Alter ≤ 50 and Autotyp=LKW then Risk=low;
- if Alter ≤ 50 and Autotyp ≠ LKW then Risk = high.



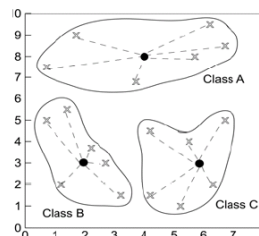
Classification 2/3




- Decision trees/ Partitioning




- Nearest Neighbors/ Lazy learners



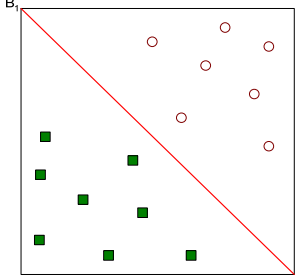


DATABASE
SYSTEMS
GROUP

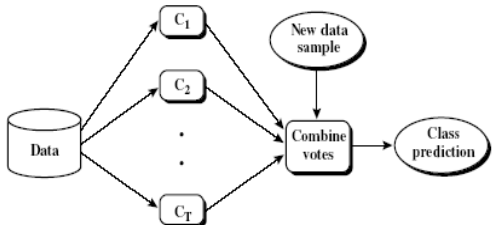
Classification 3/3



- SVM




- Ensembles




Knowledge Discovery in Databases II: Introduction and overview

23



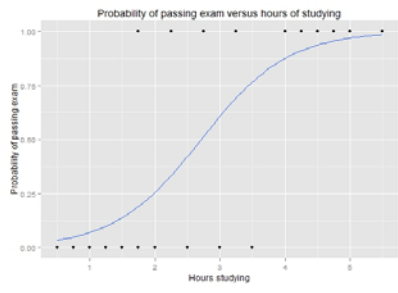
DATABASE
SYSTEMS
GROUP

Regression

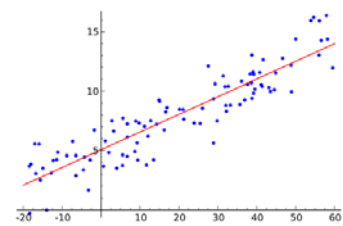


- Mapping objects to real values:
 - ⇒ determine the value for a new object
 - ⇒ describe the connection between description space and prediction space
- Supervised learning task

Logistic regression




Linear regression




Knowledge Discovery in Databases II: Introduction and overview

24





DATABASE
SYSTEMS
GROUP

Association rules/ frequent patterns 1/3





LMU

- Frequent patterns are patterns that appear frequently in a dataset.
 - Patterns: items, substructures, subsequences ...
- Typical example: Market basket analysis


Tid	Transaction items
1	Butter, Bread, Milk, Sugar
2	Butter, Flour, Milk, Sugar
3	Butter, Eggs, Milk, Salt
4	Eggs
5	Butter, Flour, Milk, Salt, Sugar

- We want to know: What products were often purchased together?
 - e.g.: beer and diapers?  
- Applications:
 - Improving store layout
 - Sales campaigns
 - Cross-marketing
 - Advertising

The parable of the beer and diapers:
http://www.theregister.co.uk/2006/08/15/beer_diapers/


Knowledge Discovery in Databases II: Introduction and overview

25



DATABASE
SYSTEMS
GROUP

Association rules/ frequent patterns 2/3



LMU

- Problem 1:** Frequent Itemsets Mining (FIM)
- Given:
 - A set of items I
 - A transactions database DB over I
 - A *minSupport* threshold s
- Goal: Find all frequent itemsets in DB , i.e.:

$$\{X \subseteq I \mid \text{support}(X) \geq s\}$$

TransaktionsID	Items
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F

Support of 1-Itemsets:

(A): 75%, (B), (C): 50%, (D), (E), (F): 25%,

Support of 2-Itemsets:


(A, C): 50%,

(A, B), (A, D), (B, C), (B, E), (B, F), (E, F): 25%

- Popular methods: Apriori, FPGrowth


Knowledge Discovery in Databases II: Introduction and overview

26



DATABASE
SYSTEMS
GROUP

Association rules/ frequent patterns 3/3



LMU

- **Problem 2: Association Rules Mining**
- Given:
 - A set of items I
 - A transactions database DB over I
 - A *minSupport* threshold s and a *minConfidence* threshold c
- Goal: Find all association rules $X \rightarrow Y$ in DB w.r.t. minimum support s and minimum confidence c , i.e.:
 - $\{X \rightarrow Y \mid \text{support}(X \cup Y) \geq s, \text{confidence}(X \rightarrow Y) \geq c\}$
 - These rules are called strong.

TransaktionsID	Items
2000	A,B,C
1000	A,C
4000	A,D
5000	B,E,F


Association rules:

$A \Rightarrow C$ (Support = 50%, Confidence= 66.6%)

$C \Rightarrow A$ (Support = 50%, Confidence= 100%)


Knowledge Discovery in Databases II: Introduction and overview

27



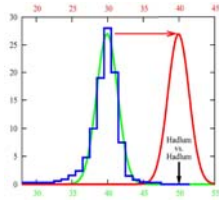
DATABASE
SYSTEMS
GROUP

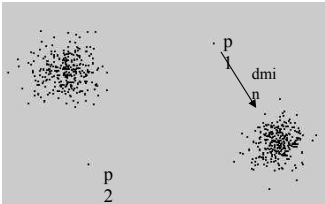
Outlier detection 1/2



LMU


- Goal: find objects that are considerably different from most other objects or unusual or in some way inconsistent with other objects
- Statistical approaches
 - Keys:
 - Probabilistic models
 - Deviation from models
- Distance-based approaches






Knowledge Discovery in Databases II: Introduction and overview

28



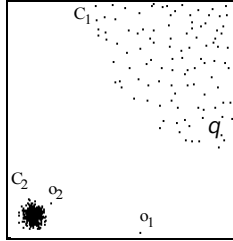
DATABASE
SYSTEMS
GROUP

Outlier detection 2/2

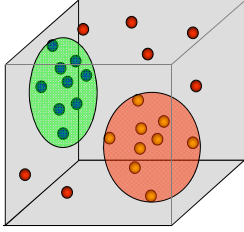


LMU

- Density-based approaches




- Clustering-based approaches




Knowledge Discovery in Databases II: Introduction and overview

29



DATABASE
SYSTEMS
GROUP

KDD I Recap




LMU


- In KDD I, we focus on how to solve specific data mining tasks
- Observations:
 - Almost all methods work on feature vectors (only)
 - Similarity / Distance measures play a key role in various data mining tasks
 - Clustering, Classification, Prediction, etc.
 - However, only simple distance functions were introduced
- In real world, useful information hidden in data with different forms
 - Suitable Feature Transformation not easy to find
 - Feature Transformation is a simple model that might lose object semantics (compare: relational vs. object model, table vs. graphs, ...)
- How to handle different types of data?
 - KDD II

Knowledge Discovery in Databases II: Introduction and overview

30




Chapter overview




- Knowledge Discovery in Databases, Big Data and Data Science
- Data Mining with Vectorized Data (Recap KDD I)
- Topics of KDD II
- Literature and supplementary materials

Knowledge Discovery in Databases II: Introduction and overview

31




KDD I vs. KDD II



- Simple data types in KDD I
 - Vector Data
- KDD II: How to deal with different complex objects.
 - Graph
 - Text
 - High-dimensional
 - Time serious
 - Shapes
 - Spatial-temporal data
 - Multi-media data
 - Heterogeneous
 -


Knowledge Discovery in Databases II: Introduction and overview

32



DATABASE
SYSTEMS
GROUP

But Before We Start: Data Cleaning




LMU

- “Dirty” in Data:
 - Dummy Values, Absence of Data, Multipurpose Fields, Contradicting Data, etc.
- Steps in Data Cleaning
 - Parsing: locates and identifies individual data elements in raw data
 - Correcting: corrects parsed individual data components using sophisticated data algorithms
 - Standardizing: applies conversion routines to transform data into standard formats
 - Matching: Searching and matching records within and across data based on predefined rules
 - Consolidating: Merges data into one representation


Knowledge Discovery in Databases II: Introduction and overview

33



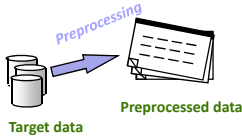
DATABASE
SYSTEMS
GROUP

Data Cleaning



LMU

- ...may take >60% of effort
- Integration of data from different sources
 - Mapping of attribute names (e.g. C_Nr → O_Id)
 - Joining different tables
 (e.g. Table1 = [C_Nr, Info1]
 and Table2 = [O_Id, Info2] ⇒
 JoinedTable = [O_Id, Info1, Info2])
- Elimination of inconsistencies
- Elimination of noise
- Computation of Missing Values (if necessary and possible)
 - Fill in missing values by some strategy (e.g. default value, average value, or application specific computations)
 - Uncertainty: Model each missing value by a (discrete) sample of possible values or a (continuous) distribution of possible values




Preprocessing

Target data Preprocessed data


Knowledge Discovery in Databases II: Introduction and overview

34



DATABASE
SYSTEMS
GROUP

Data Cleaning (Example)




LMU

- Data Quality Mining with Association Rules
 - Association rule mining generates rules for all transactions with confidence level
 - For each transaction:
 - Determine transaction type
 - Generate all related association rules
 - Summing the confidence values of the rules it violates
 - Based on the score, user can decide whether to accept or reject the data

Association Rule	Confidence
Model: S-Class → Engine: Petrol	90%
Model: S-Class → Equip: AirCondTypeC	75%
Model: S-Class → Equip: AutoWindshWiper	75%
Model: S-Class → Equip: NavigSystemD	75%
⋮	⋮


Knowledge Discovery in Databases II: Introduction and overview

35



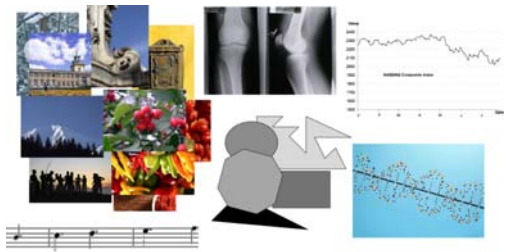
DATABASE
SYSTEMS
GROUP

Complex Object - High-dimensional data




LMU

- New applications deal with high-dimensional data (business intelligence: customers, sensors; multimedia: images, videos; biology: genes, molecules)
- High-dimensional points are abstracted to feature vectors




Knowledge Discovery in Databases II: Introduction and overview

36




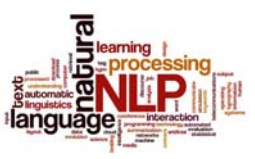
DATABASE
SYSTEMS
GROUP


Complex Object - Text

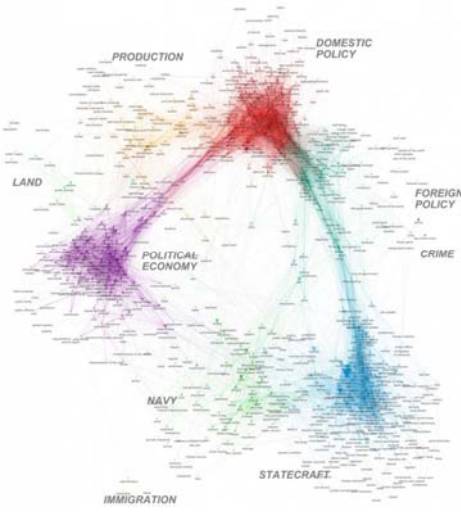


- Text: Sequence of Characters
 - Sentiment analysis
 - NLP
 - Books, static text corpi
 - Streams: Twitter, ...










The global network structure of the SoU address, 1790–2014 [from: sciencenode.org]


Knowledge Discovery in Databases II: Introduction and overview

37





DATABASE
SYSTEMS
GROUP

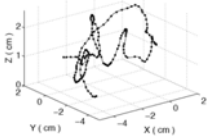
Complex Object – Sequence and Time Series Data




- Sequence: log of events happened in order
- Time series are a special type of sequences
 - Typically, values that are recorded over time
 - Index set I_n represents specific points in time
- Examples for **univariate time series**:
 - stock prices
 - audio data
 - temperature curves
 - ECG
 - amount of precipitation
- Examples for **multivariate time series**:
 - trajectories (spatial positions)
 - video data (e.g., color histograms)
 - combinations of sensor readings
- Similarity models of time series are often based on sequence similarity models










```
ATGAATTAGCTAAGGTTGTAGCTTATTTCCATAGG
GTTTIGCTCCGGACCATCCGGTCGTAGCGCGATT
GACTTGCCGGGTTGTGCCCGTATCCAGGTCACGA
CCTCATGGGAACTAGTGGCTGCCGGCAGTATCCT
GGTACGCACCTCATGTGGTATGCGTGGCTGTTGGTC
CGTATATGGACCTATATATGGATCGAAGC
```


Knowledge Discovery in Databases II: Introduction and overview

38

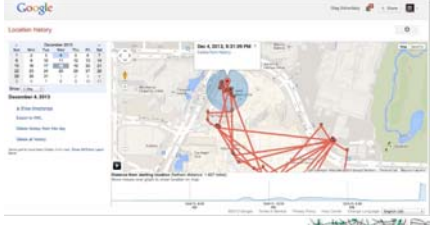





**DATABASE
SYSTEMS
GROUP**

Complex Object - Spatial-temporal data




- Objects moving in space and time
- Location-based services
- Gestures
- ...


Knowledge Discovery in Databases II: Introduction and overview

39

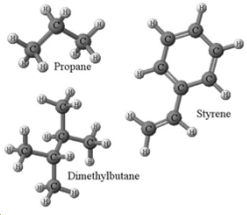


**DATABASE
SYSTEMS
GROUP**


Complex Object - Graph



- Graphs, graphs everywhere!
 - Chemical data analysis, proteins
 - Biological pathways/networks
 - Program control flow, traffic flow, work flow analysis
 - XML, Web, social network analysis
- Graphs form a complex and expressive data type
 - Trees, lattices, sequences, and items are degenerated graphs
 - Different applications result in different kinds of graphs and tasks
 - Diversity of graphs and tasks → diversity of challenges
 - Complexity of algorithms: many problems are of high complexity (NP-complete or even P-SPACE!)




Propane
Dimethylbutane
Styrene




Knowledge Discovery in Databases II: Introduction and overview

40

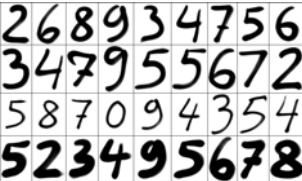
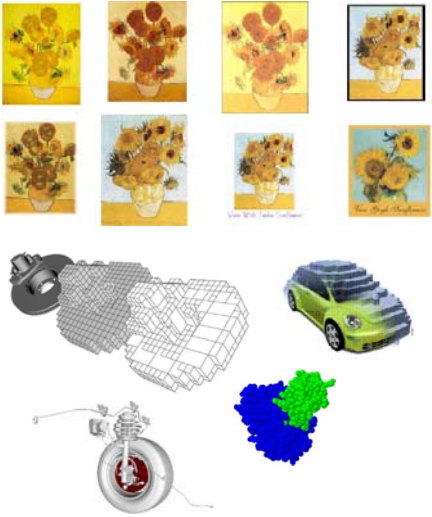


DATABASE
SYSTEMS
GROUP

Complex Object - Shapes




- (Objects in) Images
- 2D/3D objects


Knowledge Discovery in Databases II: Introduction and overview

41




DATABASE
SYSTEMS
GROUP

Complex Object - Multi-media data





- Rapid spread of multi-media data
- Nearly all device can generate and share multi-media data


<http://www.bing.com/>




<http://www.google.com/>










images








videos




Knowledge Discovery in Databases II: Introduction and overview

42



DATABASE
SYSTEMS
GROUP

Chapter overview



- Knowledge Discovery in Databases, Big Data and Data Science
- Data Mining with Vectorized Data (Recap KDD I)
- Topics of KDD II
- Literature and supplementary materials

Knowledge Discovery in Databases II: Introduction and overview

43



DATABASE
SYSTEMS
GROUP


Literature




- Han J., Kamber M., Pei J. (English)
Data Mining: Concepts and Techniques
3rd ed., Morgan Kaufmann, 2011 
- Tan P.-N., Steinbach M., Kumar V. (English)
Introduction to Data Mining
Addison-Wesley, 2006 
- Mitchell T. M. (English)
Machine Learning
McGraw-Hill, 1997 
- Lescovec J, Rajaraman A., Ulman J.
Mining of Massive Datasets
Cambridge University Press, 2014 
- Ester M., Sander J. (German)
Knowledge Discovery in Databases: Techniken und Anwendungen
Springer Verlag, September 2000 

Knowledge Discovery in Databases II: Introduction and overview

44




Further book titles




- C. M. Bishop, „*Pattern Recognition and Machine Learning*“, Springer 2007.
- S. Chakrabarti, „*Mining the Web: Statistical Analysis of Hypertext and Semi-Structured Data*“, Morgan Kaufmann, 2002.
- R. O. Duda, P. E. Hart, and D. G. Stork, „*Pattern Classification*“, 2ed., Wiley-Inter-science, 2001.
- D. J. Hand, H. Mannila, and P. Smyth, „*Principles of Data Mining*“, MIT Press, 2001.
- U. Fayyad, G. Piatetsky-Shapiro, P. Smyth: „*Knowledge discovery and data mining: Towards a unifying framework*“, in: Proc. 2nd ACM Int. Conf. on Knowledge Discovery and Data Mining (KDD), Portland, OR, 1996

Knowledge Discovery in Databases II: Introduction and overview

45



Online Resources



- *Mining Massive Datasets* class by Jure Lescovec, Anand Rajaraman and Jeffrey D. Ullman
 - <https://www.coursera.org/course/mmds>
- *Machine Learning* class by Andrew Ng, Stanford
 - <http://ml-class.org/>
- *Introduction to Databases* class by Jennifer Widom, Stanford
 - <http://www.db-class.org/course/auth/welcome>
- Kdnuggets: Data Mining and Analytics resources
 - <http://www.kdnuggets.com/>

Knowledge Discovery in Databases II: Introduction and overview

46