

Ludwig-Maximilians-Universität München
Lehrstuhl für Datenbanksysteme und Data Mining
Prof. Dr. Thomas Seidl

Knowledge Discovery and Data Mining 1

(Data Mining Algorithms 1)

Winter Semester 2022/23

Further Topics



Agenda

1. Introduction
2. Preliminaries: Data
3. Supervised Learning
4. Unsupervised Learning
5. Process Mining
6. Further Topics
 - 6.1 Privacy in Data Mining
 - 6.2 Outlook

Data Privacy

Situation

- ▶ Huge volume of data is collected
- ▶ From a variety of devices and platforms (e.g. Smartphones, Wearables, Social Networks, Medical systems)
- ▶ Capturing human behaviors, locations, routines, activities and affiliations
- ▶ Providing an opportunity to perform data analytics

Data Abuse is inevitable

- ▶ It compromises individual's privacy
- ▶ Or breaches the security of an institution

Data Privacy

- ▶ These privacy concerns need to be mitigated
- ▶ They have prompted huge research interest to *protect data*
- ▶ But,

Strong Privacy Protection \implies Poor Data Utility
Good Data Utility \implies Weak Privacy Protection



Challenge

Find a good trade-off between Data Utility and Privacy

Objectives of Privacy Preserving Data Mining

- ▶ Ensure data privacy
- ▶ Maintain a good trade-off between data utility and privacy

Paradigms

- ▶ k -Anonymity
- ▶ l -Diversity
- ▶ Differential Privacy

Linkage Attack

Method

Different public records can be linked to it to breach privacy

Hospital Records

<i>Private</i>		<i>Public</i>		
Name	Sex	Age	Zip	Disease
Alice	F	29	52062	Breast Cancer
Janes	F	27	52064	Breast Cancer
Jones	M	21	52066	Lung Cancer
Frank	M	35	52072	Heart Disease
Ben	M	33	52078	Fever
Betty	F	37	52080	Nose Pains

Public Records from Sport Club

<i>Public</i>				
Name	Sex	Age	Zip	Sport
Alice	F	29	52062	Tennis
Theo	M	41	52066	Golf
John	M	24	52062	Soccer
Betty	F	37	52080	Tennis
James	M	34	82066	Soccer

k -Anonymity

Privacy paradigm for protecting data records before *publication*

Three kinds of attributes:

1. *Key Attributes*: Uniquely identifiable attributes (e.g., Name, Social Security Number, Telephone Number)
2. *Quasi-identifier*: Groups of attributes that can be combined with external data to uniquely re-identify an individual (e.g. (Date of Birth, Zip Code, Gender))
3. *Sensitive Attributes*: An attacker should not be able to combine these with the key attributes. (e.g. Disease, Salary, Habit, Location etc.)

k -Anonymity

Attention

Hiding key attributes alone does **not** guarantee privacy.

An attacker may be able to break the privacy by combining the quasi-identifiers from the data with those from publicly available information.

Definition: k -Anonymity

Given a set of quasi-identifiers in a database table, the database table is said to be *k -Anonymous*, if the sequence of records in each quasi-identifier exists at least k times.

Ensure privacy by *Suppression* or *Generalization* of quasi-identifiers.

k-Anonymity: Suppression

Suppression

Accomplished by replacing a part or the entire attribute value by placeholder, e.g. “?”
(= generalization)

Example

- ▶ Suppress Postal Code: 52062 \mapsto 52???
- ▶ Suppress Gender: Male \mapsto ?; Female \mapsto ?

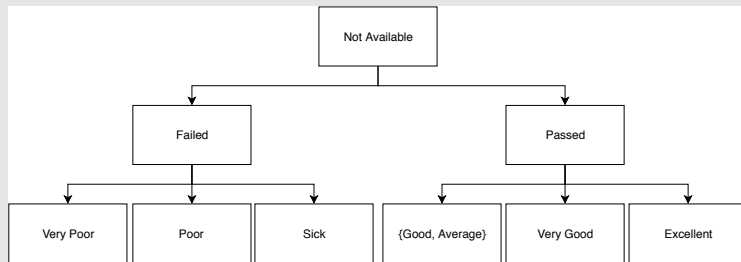
k-Anonymity: Generalization

Generalization

Accomplished by aggregating values from fine levels to coarser resolution using generalisation hierarchy.

Example

Generalize exam grades:



Shortcomings: Background Knowledge Attack

Background Knowledge Attack

Lack of diversity of the sensitive attribute values (homogeneity)

Example

- *Background Knowledge*: Alice is (i) 29 years old and (ii) female
- *Homogeneity*: All 2*-aged females have Breast Cancer.
⇒ Alice has BC!

Release			
Quasi Identifier			Sensitive
Sex	Age	Zip	Disease
F	2?	520??	Breast Cancer
F	2?	520??	Breast Cancer
M	2?	520??	Lung Cancer
M	3?	520??	Heart Disease
M	3?	520??	Fever
F	3?	520??	Nose Pains

This led to the creation of a new privacy model called *l*-diversity

Distinct l -Diversity

An quasi-identifier is l -diverse, if there are at least l different values. A dataset is l -diverse, if all QIs are l -diverse.

Example

Not "diverse"

Quasi Identifier	Sensitive
QI 1	Headache
QI 1	Headache
QI 1	Headache
QI 2	Cancer
QI 2	Cancer

2-diverse

Quasi Identifier	Sensitive
QI 1	Headache
QI 1	Cancer
QI 1	Headache
QI 2	Headache
QI 2	Cancer

Other Variants

- ▶ *Entropy I-Diversity*: For each equivalent class, the entropy of the distribution of its sensitive values must be at least I
- ▶ *Probabilistic I-Diversity*: The most frequent sensitive value of an equivalent class must be at most $1/I$

Limitations

- ▶ Not necessary at times
- ▶ Difficult to achieve: For large record size, many equivalent classes will be needed to satisfy *I-Diversity*
- ▶ Does not consider the distribution of sensitive attributes

Background Attack Assumption

- ▶ k -Anonymity and l -Diversity make assumptions about the adversary
- ▶ They at times fall short of their goal to prevent data disclosure
- ▶ There is another privacy paradigm which does not rely on background knowledge, called *Differential Privacy*

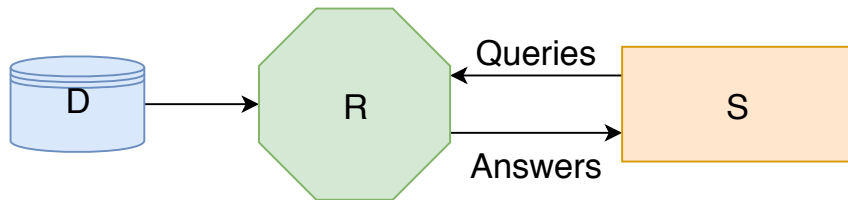
Differential Privacy

Core Idea

Privacy through data perturbation.

- ▶ The addition or removal of one record from a database should not reveal any information to an adversary, i.e. your *presence* or *absence* does not reveal or leak any information.
- ▶ Use a randomization mechanism to perturb query results of count, sum, mean functions, as well as other statistical query functions.

Differential Privacy



Definition

A randomized mechanism $R(x)$ provides ϵ -differential privacy if for any two databases D_1 and D_2 that differ on at most one element, and all outputs $S \subseteq \text{Range}(R)$

$$\frac{\Pr[R(D_1) \in S]}{\Pr[R(D_2) \in S]} \leq \exp(\epsilon)$$

ϵ is a parameter called *privacy budget/level*.

Data Perturbation

Data perturbation is achieved by noise addition.

Some Kinds of Noise

- ▶ Laplacian noise
- ▶ Gaussian noise
- ▶ Exponential Mechanism

Agenda

1. Introduction
2. Preliminaries: Data
3. Supervised Learning
4. Unsupervised Learning
5. Process Mining
6. Further Topics
 - 6.1 Privacy in Data Mining
 - 6.2 Outlook

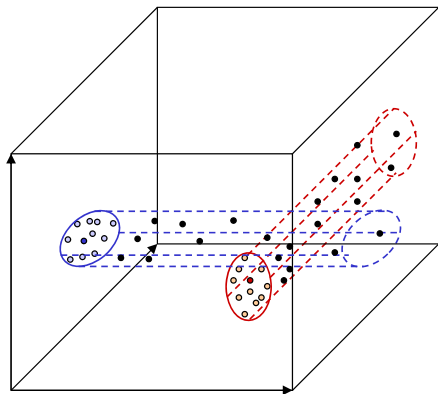
High-Dimensional Data

► Challenges:

- *Curse of dimensionality*: distances become more and more similar
- Datasets become sparse.
- Expensive distance measures
- Degeneration of index structures
- Unintuitive properties in high dimensions.

► Tasks

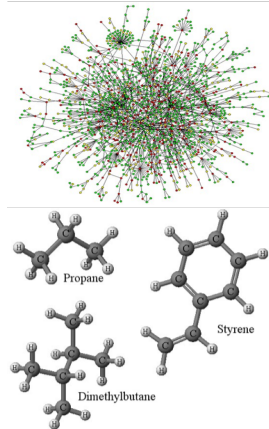
- Feature Selection
- Feature Reduction / Metric Learning
- Clustering in High-Dimensional Spaces



↪ Knowledge Discovery in Databases 2 (summer)

Graph Data

- ▶ Graphs are everywhere!
 - ▶ Chemical data analysis, proteins
 - ▶ Biological pathways/networks
 - ▶ Program control flow, traffic flow
 - ▶ Web graph, social network analysis
- ▶ Typical tasks
 - ▶ Measure similarity between graphs
 - ▶ Find frequent patterns in graphs
 - ▶ Generate "realistic" synthetic graphs
 - ▶ Identify groups in social networks
 - ▶ Integrate additional information



~> Knowledge Discovery in Databases 2 (summer)

Further Machine Learning Methods

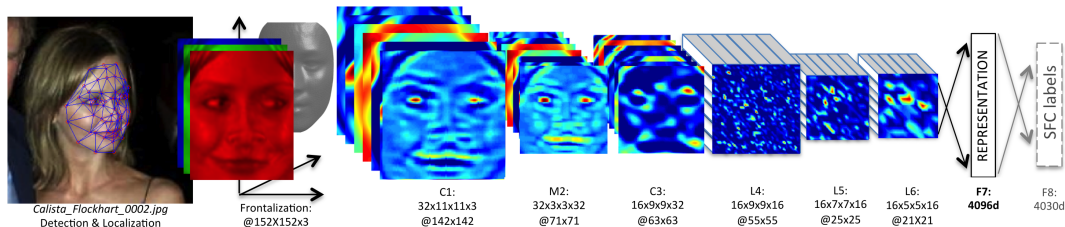


Image Source: Taigman, et al. "Deepface: Closing the gap to human-level performance in face verification." CVPR'14.

- ▶ Graphical Models
- ▶ Generative Models

- ▶ Neural Networks
- ▶ Deep Learning

~> Machine Learning (summer), Deep Learning and Artificial Intelligence (winter)

Decision Making / Planning

► Setting:

- Agents are in some environment, observe, and have to take actions that influence the environment.

► Methods:

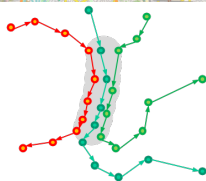
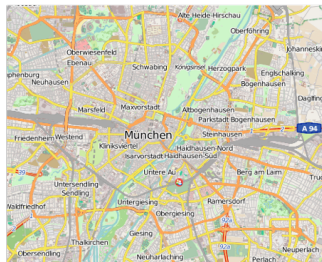
- Deterministic/Stochastic Planning
- A*-Search
- Model-Free Reinforcement Learning
- Q-Learning
- Adversarial Search (e.g. Alpha-Beta Pruning)



↪ Deep Learning (winter), AI for Games (summer)

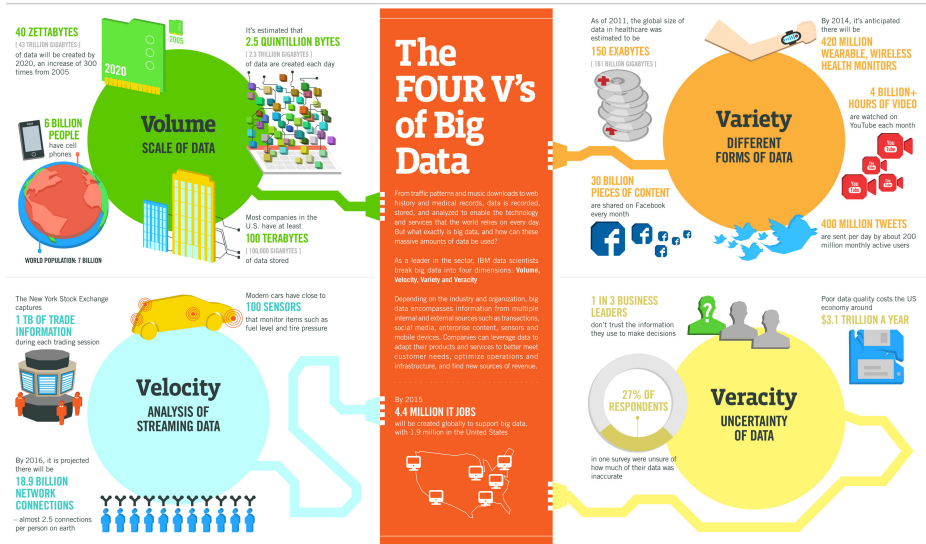
Spatial Data

- ▶ Mining spatial data
 - ▶ Spatial clustering, outlier detection, prediction, rule mining, ...
- ▶ Spatial data management
 - ▶ Process spatial queries without scanning the whole database
 - ▶ Spatial index structures: BSP-tree, R-tree, Quad-tree, ...
- ▶ Mining trajectory data
 - ▶ Similarity models for trajectories
 - ▶ Trajectory compression
 - ▶ Mining patterns in trajectories (encounters, flocks, ...)



~> AI for Games (summer)

Big Data



Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, GAS

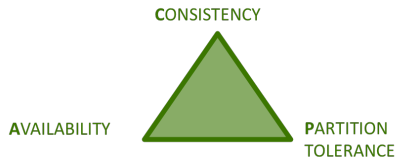
IBM

Big Data Management

- ▶ Vertical scaling limited and expensive
 ~> Distributed storage
- ▶ NoSQL databases
 - ▶ Redis
 - ▶ MongoDB
 - ▶ Cassandra
 - ▶ Neo4J
- ▶ Distributed file systems
 - ▶ GFS (Google)
 - ▶ HDFS (Hadoop)
 - ▶ S3 (Amazon)



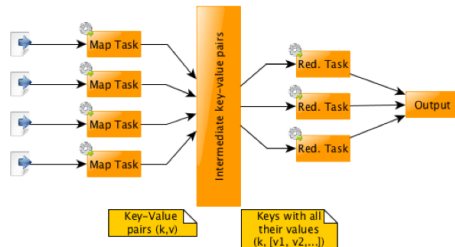
<https://www.greentree.com/latest-news/avoiding-cumulus-congestus>



~> Big Data Management and Analytics (winter)

Distributed Data Processing

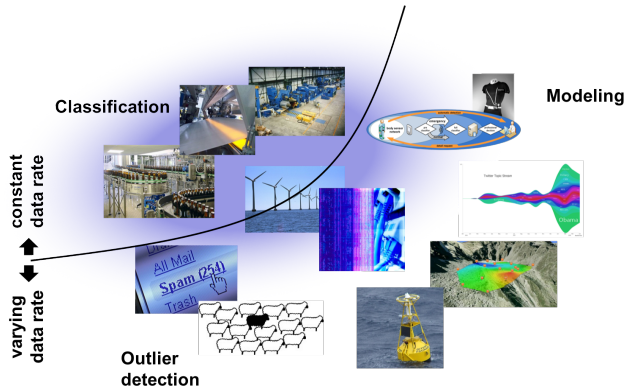
- ▶ Processing and analyzing big data
- ▶ Map-Reduce: Programming model for distributed processing of large datasets
 - ▶ Algorithms are specified as sequences of map and reduce functions
 - ▶ Programs are automatically parallelized and executed on a cluster
 - ▶ System is tolerant to hardware faults
- ▶ Frameworks
 - ▶ Apache Spark (batch processing)
 - ▶ Apache Flink (stream processing)



↪ Big Data Management and Analytics (winter)

Stream Data

- ▶ Data objects arrive over time in a continuous data stream
- ▶ Challenges
 - ▶ Infinite stream
 - ▶ Limited time and memory
 - ▶ Evolving distribution
 - ▶ Varying data rates
 - ▶ Concept drift
- ▶ Typical tasks
 - ▶ Sampling and buffering
 - ▶ Stream statistics
 - ▶ Aging mechanisms



~> Big Data Management and Analytics (winter)

Seminars, Practicals, Theses

Dive deeper into specific topics and get hands-on experience:

- ▶ Master Seminar "Recent Developments in Data Science" (summer and winter)
- ▶ Master Practical "Big Data Science" (summer)
- ▶ Master Practical "Applied Reinforcement Learning" (summer)
- ▶ Individual Bachelor and Master Theses

Conclusion

Best wishes for your exams, and see you at further opportunities!