

Data Warehouses

Krzysztof Dembczyński

Intelligent Decision Support Systems Laboratory (IDSS)
Poznań University of Technology, Poland



Software Development Technologies
Master studies, first semester
Academic year 2017/18 (winter course)

Goal: understanding data ...



Goal: ... to make data analysis efficient.

- **Buzzwords:** Big Data, Data Science, Machine learning, NoSQL ...

- **Buzzwords:** Big Data, Data Science, Machine learning, NoSQL ...
- **How Big Data Changes Everything:**
 - ▶ Several books showing the impact of Big Data revolution (e.g., **Disruptive Possibilities: How Big Data Changes Everything** by Jeffrey Needham).

- **Buzzwords:** Big Data, Data Science, Machine learning, NoSQL ...
- **How Big Data Changes Everything:**
 - ▶ Several books showing the impact of Big Data revolution (e.g., **Disruptive Possibilities: How Big Data Changes Everything** by Jeffrey Needham).
- **Computerworld (Jul 11, 2007):**
 - ▶ *12 IT skills that employers can't say no to:*
 - 1) *Machine learning*...

- **Buzzwords:** Big Data, Data Science, Machine learning, NoSQL ...
- **How Big Data Changes Everything:**
 - ▶ Several books showing the impact of Big Data revolution (e.g., **Disruptive Possibilities: How Big Data Changes Everything** by Jeffrey Needham).
- **Computerworld (Jul 11, 2007):**
 - ▶ *12 IT skills that employers can't say no to:*
 - 1) *Machine learning*

...
- **Three priorities of Google announced at BoxDev 2015:**
 - ▶ Machine learning – speech recognition
 - ▶ Machine learning – image understanding
 - ▶ Machine learning – preference learning/personalization

- **Buzzwords:** Big Data, Data Science, Machine learning, NoSQL ...
- **How Big Data Changes Everything:**
 - ▶ Several books showing the impact of Big Data revolution (e.g., **Disruptive Possibilities: How Big Data Changes Everything** by Jeffrey Needham).
- **Computerworld (Jul 11, 2007):**
 - ▶ *12 IT skills that employers can't say no to:*
 - 1) Machine learning

...
- **Three priorities of Google announced at BoxDev 2015:**
 - ▶ Machine learning – speech recognition
 - ▶ Machine learning – image understanding
 - ▶ Machine learning – preference learning/personalization
- **OpenAI** founded in 2015 as a non-profit artificial intelligence research company.

Data mining

- Data mining is the discovery of **models** for data, ...
- But what is a model?

if all you have is a hammer, everything looks like a nail

How to characterize your data?

- **Database programmer** usually writes:

```
select avg(column), std(column) from data
```

How to characterize your data?

- **Database programmer** usually writes:

```
select avg(column), std(column) from data
```

- **Statistician** might decide that the data comes from a Gaussian distribution and use a formula to compute the most likely parameters of this Gaussian: the mean and standard deviation.

How to characterize your data?

- **Database programmer** usually writes:

```
select avg(column), std(column) from data
```

- **Statistician** might decide that the data comes from a Gaussian distribution and use a formula to compute the most likely parameters of this Gaussian: the mean and standard deviation.
- **Machine learner** will use the data as training examples and apply a learning algorithm to get a model that predicts future data.

How to characterize your data?

- **Database programmer** usually writes:

```
select avg(column), std(column) from data
```

- **Statistician** might decide that the data comes from a Gaussian distribution and use a formula to compute the most likely parameters of this Gaussian: the mean and standard deviation.
- **Machine learner** will use the data as training examples and apply a learning algorithm to get a model that predicts future data.
- **Data miner** will discover the most frequent patterns.

They all want to understand data and use this knowledge for making better decisions

Data+ideas vs. statistics+algorithms

- About the Amazon's recommender system:

It's often more important to creatively invent new data sources than to implement the latest academic variations on an algorithm.

Data+ideas vs. statistics+algorithms

- About the Amazon's recommender system:
It's often more important to creatively invent new data sources than to implement the latest academic variations on an algorithm.
- WhizBang! Labs tried to use machine learning to locate people's resumes on the Web: the algorithm was not able to do better than procedures designed by hand, since a resume has a quite standard shape and sentences.

Data+computational power

- Object recognition in computer vision:

Data+computational power

- Object recognition in computer vision:
 - ▶ Scanning large databases **can perform better** than the best computer vision algorithms!

Data+computational power

- Object recognition in computer vision:
 - ▶ Scanning large databases **can perform better** than the best computer vision algorithms!
- Automatic translation

Data+computational power

- Object recognition in computer vision:
 - ▶ Scanning large databases **can perform better** than the best computer vision algorithms!
- Automatic translation
 - ▶ Statistical translation based on large corpora **outperforms** linguistic models!

Human computation

- CAPTCHA and reCAPTCHA
- ESP game
- Check a lecture given by Luis von Ahn:
http://videlectures.net/iaai09_vonahn_hc/
- Amazon Mechanical Turk

Data+ideas vs. statistics+algorithms

Those who ignore Statistics are condemned to reinvent it.

Brad Efron

- In Statistics, a term **data mining** was originally referring to attempts to extract information that was not supported by the data.
- Bonferroni's Principle: "if you look in more places for interesting patterns than your amount of data will support, you are bound to find crap".
- Rhine paradox.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ XBox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ XBox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ Xbox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ Xbox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.
 - ▶ Google and PageRank.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ Xbox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.
 - ▶ Google and PageRank.
 - ▶ Clustering of Cholera cases in 1854.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ XBox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.
 - ▶ Google and PageRank.
 - ▶ Clustering of Cholera cases in 1854.
 - ▶ Win one of the Kaggle's competitions!!! <http://www.kaggle.com/>.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ XBox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.
 - ▶ Google and PageRank.
 - ▶ Clustering of Cholera cases in 1854.
 - ▶ Win one of the Kaggle's competitions!!! <http://www.kaggle.com/>.
 - ▶ Autonomous cars.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ XBox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.
 - ▶ Google and PageRank.
 - ▶ Clustering of Cholera cases in 1854.
 - ▶ Win one of the Kaggle's competitions!!! <http://www.kaggle.com/>.
 - ▶ Autonomous cars.
 - ▶ Deep learning.

Data+ideas vs. statistics+algorithms

- The data mining algorithms can perform quite well!!!
 - ▶ XBox Kinect: object tracking vs. pattern recognition (check: http://videlectures.net/ecmlpkdd2011_bishop_embracing/).
 - ▶ Pattern finding: association rules.
 - ▶ Netflix: recommender system.
 - ▶ Google and PageRank.
 - ▶ Clustering of Cholera cases in 1854.
 - ▶ Win one of the Kaggle's competitions!!! <http://www.kaggle.com/>.
 - ▶ Autonomous cars.
 - ▶ Deep learning.
 - ▶ And many others.

Data+ideas+computational power+statistics+algorithms

Processing of massive data sets

- To make the data analysis efficient, we need to organize data in a way that ensures efficient storage and access.

Processing of massive data sets

- To make the data analysis efficient, we need to organize data in a way that ensures efficient storage and access.
- Different data management technologies:
 - ▶ File management systems,
 - ▶ Database management systems (hierarchical, network-based, relational),
 - ▶ Data warehouses,
 - ▶ NoSQL.

Data warehouses

- Data warehouse is a first component of the **decision support/business intelligence system**.

Data warehouses

- Data warehouse is a first component of the **decision support/business intelligence system**.
- **Information processing**: querying, basic statistical analysis, reporting using cross-tabs, tables, charts, or graphs, low-cost Web-based accessing tools integrated with Web browsers.

Data warehouses

- Data warehouse is a first component of the **decision support/business intelligence system**.
- **Information processing**: querying, basic statistical analysis, reporting using cross-tabs, tables, charts, or graphs, low-cost Web-based accessing tools integrated with Web browsers.
- **Exploratory querying**: OLAP operations for multidimensional data view, finding unexpected facts in databases.

Data warehouses

- Data warehouse is a first component of the **decision support/business intelligence system**.
- **Information processing**: querying, basic statistical analysis, reporting using cross-tabs, tables, charts, or graphs, low-cost Web-based accessing tools integrated with Web browsers.
- **Exploratory querying**: OLAP operations for multidimensional data view, finding unexpected facts in databases.
- **Approximate queries**: response times are often impractical for large data warehouses: use fast, approximate answers.

Data warehouses

- Data warehouse is a first component of the **decision support/business intelligence system**.
- **Information processing**: querying, basic statistical analysis, reporting using cross-tabs, tables, charts, or graphs, low-cost Web-based accessing tools integrated with Web browsers.
- **Exploratory querying**: OLAP operations for multidimensional data view, finding unexpected facts in databases.
- **Approximate queries**: response times are often impractical for large data warehouses: use fast, approximate answers.
- **Knowledge discovery**: finding hidden patterns and associations, analytical models for prediction and clustering, visualization.

To be learned in the upcoming semester ...

The aim and the scope of the course

- **Aim:** To get to know how to design and construct data warehouses for efficient data processing.
- **Scope:** We will learn about:
 - ▶ Dimensional modeling,
 - ▶ ETL process,
 - ▶ OLAP systems,
 - ▶ MapReduce systems (Spark),
 - ▶ Processing of massive data.

Main information about the course

- Instructor:
 - ▶ dr inż. Krzysztof Dembczyński (kdembczynskicsputpoznanpl)
- Website:
 - ▶ www.cs.put.poznan.pl/kdembczynski/lectures/dw

Lectures

- Main topics of lectures:
 - ▶ Introduction
 - ▶ Evolution of database systems
 - ▶ Dimensional modeling
 - ▶ ETL and OLAP systems
 - ▶ MapReduce systems (Spark)
 - ▶ Processing of massive data.

Labs

- Strong connection between lectures and labs.
- Software: Spark.
- List of tasks and exercises for each meeting (also homeworks).
- Small programming projects and short exercises.
- Main topics:
 - ▶ Dimensional modeling
 - ▶ ETL process
 - ▶ MapReduce in Spark

Evaluation

- **Lecture:**

Test: 75 % (min. 50%)

Labs: 25 % (min. 50%)

- **Labs:**

Regular exercises and home works: 4x25 points (min. 50%)

- **Scale:**

90 % – 5.0 80 % – 4.5 70 % – 4.0

60 % – 3.5 50 % – 3.0 < 50 % – 2.0

- **Bonus points for all:** up to 10 percent points.

Bibliography

- H. Garcia-Molina, J. D. Ullman, and J. Widom. *Systemy baz danych. Kompletny podręcznik. Wydanie II.*
Helion, 2011
- Z. Królikowski. *Hurtownie danych: logiczne i fizyczne struktury danych.*
Wydawnictwo Politechniki Poznańskiej, 2007
- R. Kimball and M. Ross. *The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling, 3rd Edition.*
John Wiley & Sons, 2013
- A. Rajaraman and J. D. Ullman. *Mining of Massive Datasets.*
Cambridge University Press, 2011
<http://www.mmds.org>