# Big Data and Data Profiling
## Introduction

20.4.2017

Felix Naumann

# Overview

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Information Systems Team



Diana **Stephan**

Prof. Felix **Naumann**

Sebastian **Kruse**

Tim **Repke**

Thorsten **Papenbrock**

Dr. Ralf **Krestel**

project **DuDe**

**Duplicate Detection**

Data Change

Data Fusion

project **Stratosphere**

Entity Search

**Data Profiling**

**Information Integration**

**Web Science**

Data Scrubbing

project **DataChEx**

Data as a Service

**Information Quality**

**Data Cleansing**

**Text Mining**

Tobias **Bleifuß**

Web Data

Linked Open Data

RDF Data Mining

Dependency Detection

Hazar **Harmouch**

ETL Management

**Service-Oriented Systems**

**Entity Recognition**

**Opinion Mining**

project **Metanome**

John **Koumarelas**

Toni **Grütze**

Michael **Loster**

Ahmad **Samiei**

Zhe **Zuo**

Konstantina **Lazariduo**

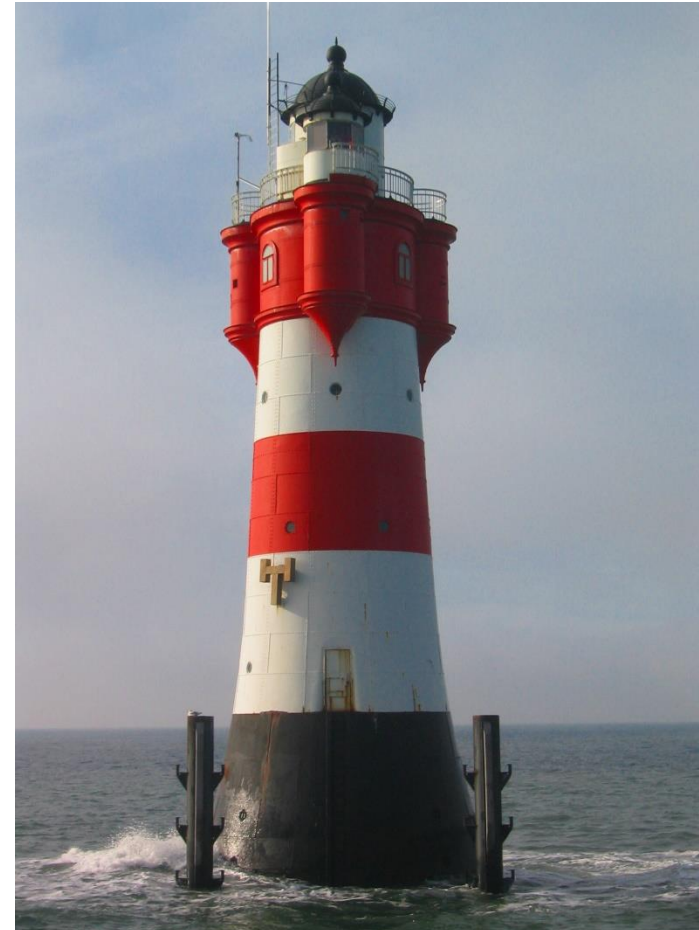Lan **Jiang**

# Other courses in this semester

- Lectures
  - DBS I (Bachelor)
  - Data Profiling
- Seminars
  - Bachelor: Text Mining
  - Master: Recommender Systems
- Bachelorproject
  - Ingestion – Commerzbank
- Masterproject
  - Hate Speech Detection



Extending the Database Relational Model to Capture More Meaning

E. F. CODD
IBM Research Laboratory

Felix Naumann
Data Profiling
Summer 2017

# Overview

1. Introduction to research group
2. **Lecture organisation**
3. (Big) data
   - Data sources
   - Profiling
4. Overview of semester

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Dates and exercises

- Lectures
  - Mondays 15:15 – 16:45
  - Thursdays 13:30 – 15:00
- Exercises
  - In parallel
- First lecture
  - 20.4.2017
- Last lecture
  - 27.7.2017
- See Web for timetable updates!

- Exam
  - Oral or written exam
  - Probably first week after lectures
- Prerequisites
  - To participate
    - Background in databases and their implementation (e.g. DBS I and II)
  - For exam
    - Attend lectures
    - Active participation in exercises
    - "Successfully" complete exercise tasks

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Feedback

- Evaluation at end of semester

- Question any time please!
  - During lectures
  - During consultation: Tuesdays 13-15
  - Email: naumann@hpi.de

- Also: Give feedback about
  - improving lectures
  - informational material
  - organization

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Literature

- No single textbook
- References to various papers during lecture
- All papers are available either via email from me or (preferred) from
  - Google Scholar: http://scholar.google.com/
  - DBLP: http://www.informatik.uni-trier.de/~ley/db/index.html
  - CiteSeer: http://citeseer.ist.psu.edu/
  - ACM Digital Library: www.acm.org/dl/
  - Homepages of authors

- Profiling relational data: a survey. Ziawasch Abedjan, Lukasz Golab, Felix Naumann, VLDB Journal, vol. 24(4):557-581 2015
  - https://hpi.de/fileadmin/user_upload/fachgebiete/naumann/publications/2015/dataprofiling_main.pdf

# Exercise

- Algorithm design and programming exercises
  - Data profiling (emphasis on efficiency and scalability)
  - Unique column combinations
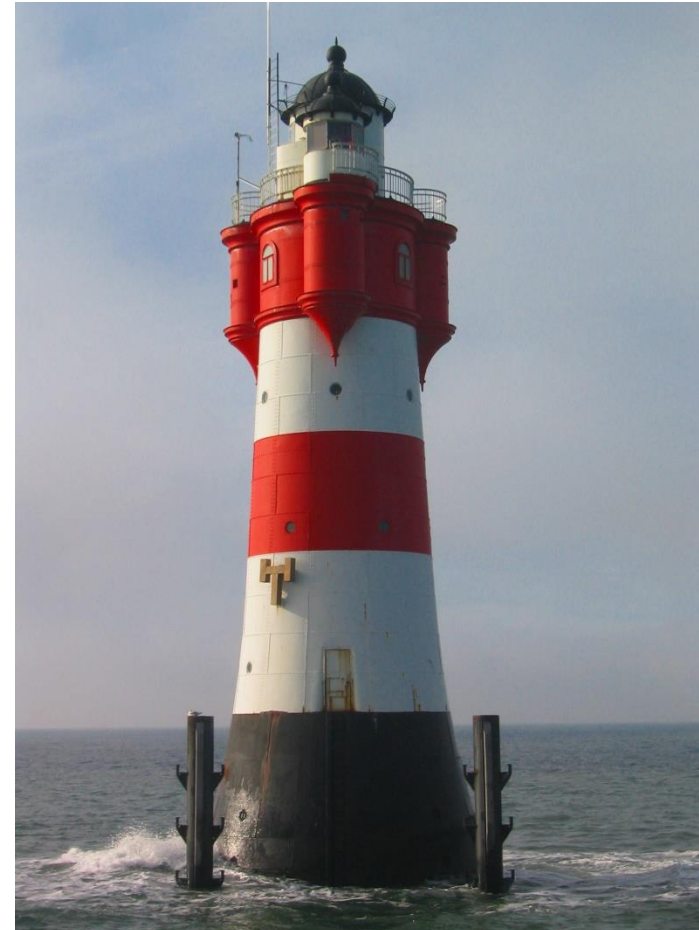  - Inclusion dependencies
  - Functional dependencies

- Self-motivation wrt. good solutions!

# Introduction: Audience

- Which semester?
- HPI or UP?
- Erasmus o.ä.?
  - English?
- Database knowledge?
  - Which other related lectures?
- Your motivation?

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Overview

1. Introduction to research group
2. Lecture organisation
3. **(Big) data**
   - Data sources
   - Profiling
4. Overview of semester

**Felix Naumann**
**Data Profiling**
**Summer 2017**
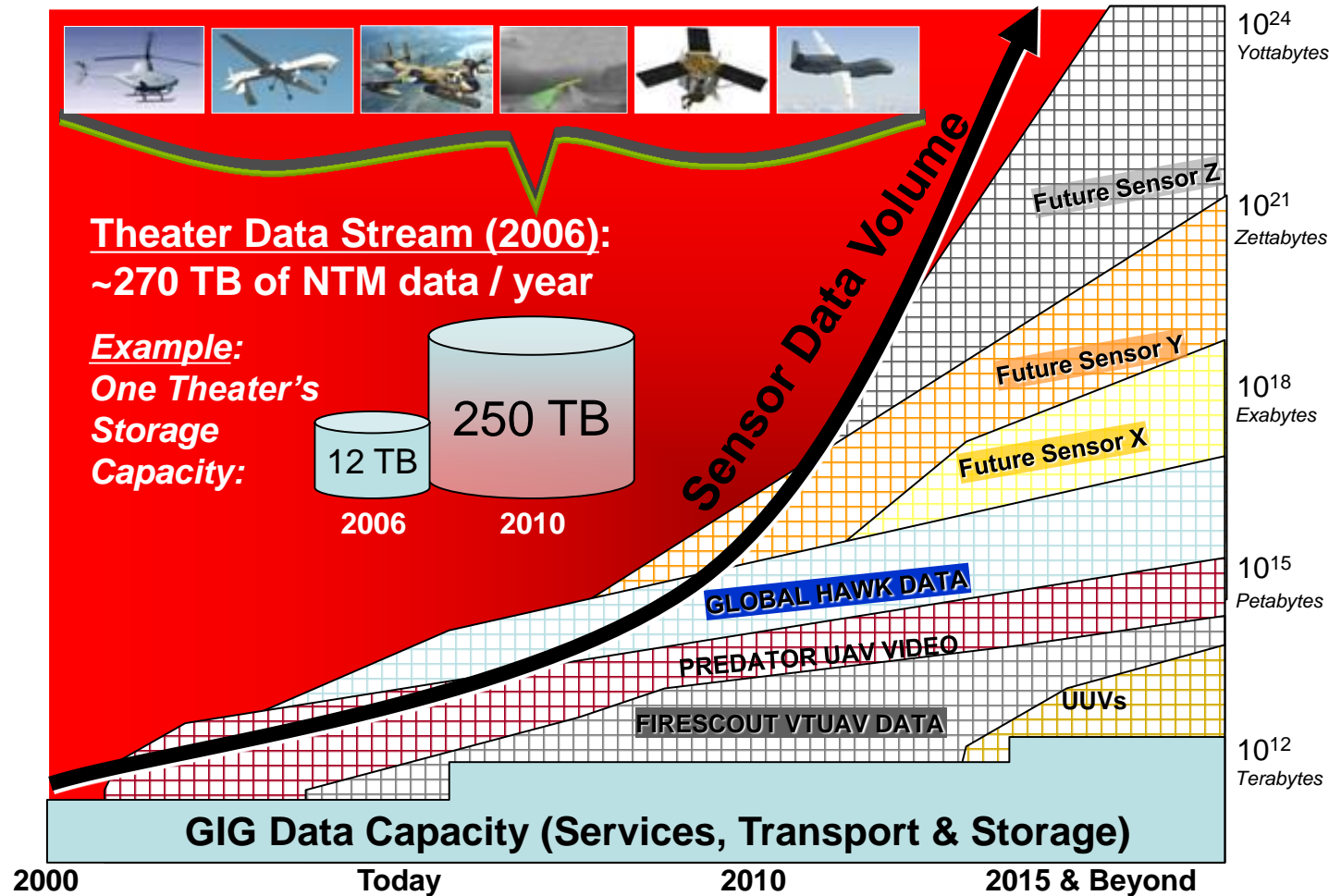
# Big Data Motivation

We're now entering what I call the "Industrial Revolution of Data," where the majority of data will be stamped out by machines: software logs, cameras, microphones, RFID readers, wireless sensor networks and so on.
These machines generate data a lot faster than people can, and their production rates will grow exponentially with Moore's Law. Storing this data is cheap, and it can be mined for valuable information.
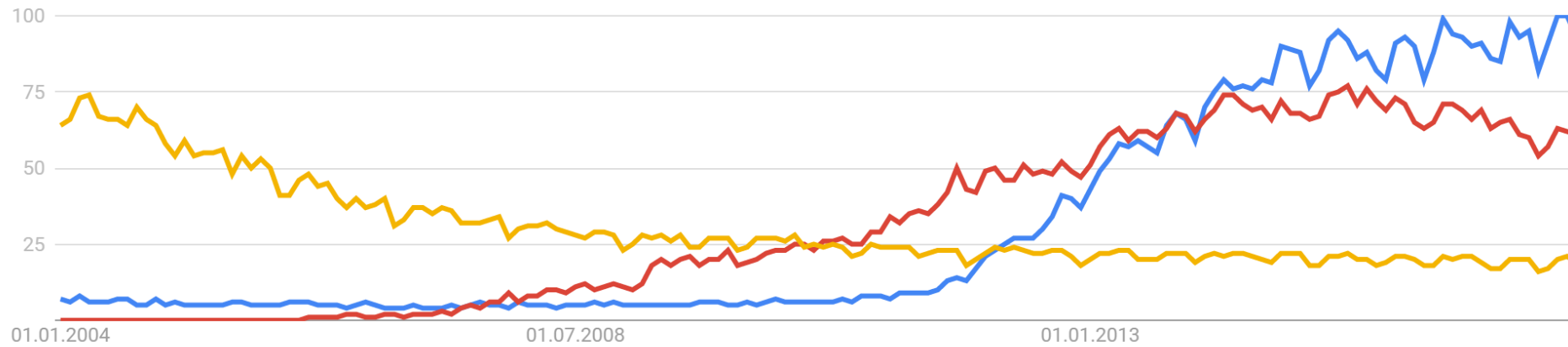
- Joe Hellerstein
  http://gigaom.com/2008/11/09/mapreduce-leads-the-way-for-parallel-programming/

**Felix Naumann**
**Data Profiling**
**Summer 2017**

**Theater Data Stream (2006):**
**~270 TB of NTM data / year**

*Example:*
*One Theater's Storage Capacity:*

12 TB — 2006
250 TB — 2010

Sensor Data Volume

$10^{24}$ *Yottabytes*

Future Sensor Z

$10^{21}$ *Zettabytes*

Future Sensor Y

Future Sensor X

$10^{18}$ *Exabytes*

GLOBAL HAWK DATA

PREDATOR UAV VIDEO

FIRESCOUT VTUAV DATA

UUVs

$10^{15}$ *Petabytes*

$10^{12}$ *Terabytes*

**GIG Data Capacity (Services, Transport & Storage)**

2000 — Today — 2010 — 2015 & Beyond

Felix Naumann
Data Profiling
Summer 2017

Using 1TB drives, this would require 1 trillion ($10^{12}$) drives!

13

Bob Gourley: Thoughts on the future of Information Sharing Technology

# Big Data trends



**data warehouse**   **big data**   **hadoop**

**Felix Naumann
Data Profiling
Summer 2017**

Quelle: Google trends, Apr. 2017

# Defining Big Data

Big data is a collection of data sets so large and complex that it becomes difficult to process using on-hand database management tools or traditional data processing applications.

| Capture | Sharing | Querying |
|---------|---------|----------|
| Extraction | Search | Analysis |
| Curation | Storage | Visualization |

If data is **too big**, **too fast**, or **too hard** for existing tools to process, it is Big Data.

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Gartner's 3 (+ 1) V's – Properties of Big Data

- **Volume**
  - 12 terabytes of Tweets: product sentiment analysis
  - 350 billion annual meter readings: predict power consumption
- **Velocity**
  - 5 million daily trade events: identify potential fraud
  - 500 million daily call detail records: predict customer churn faster
- **Variety**
  - 100's of live video feeds from surveillance cameras
  - 80% data growth in images, video and documents to improve customer satisfaction
- **Veracity** (Wahrhaftigkeit)
  - 1 in 3 business leaders don't trust the information they use to make decisions.

**Felix Naumann
Data Profiling
Summer 2017**

http://www.ibm.com/software/data/bigdata/

# More V's

- **Viscosity**
  - Integration and dataflow friction
- **Venue**
  - Different locations that require different access & extraction methods
- **Vocabulary**
  - Different language and vocabulary
- **Value**
  - Added-value of data to organization and use-case
- **Virality**
  - Speed of dispersal among community
- **Variability**
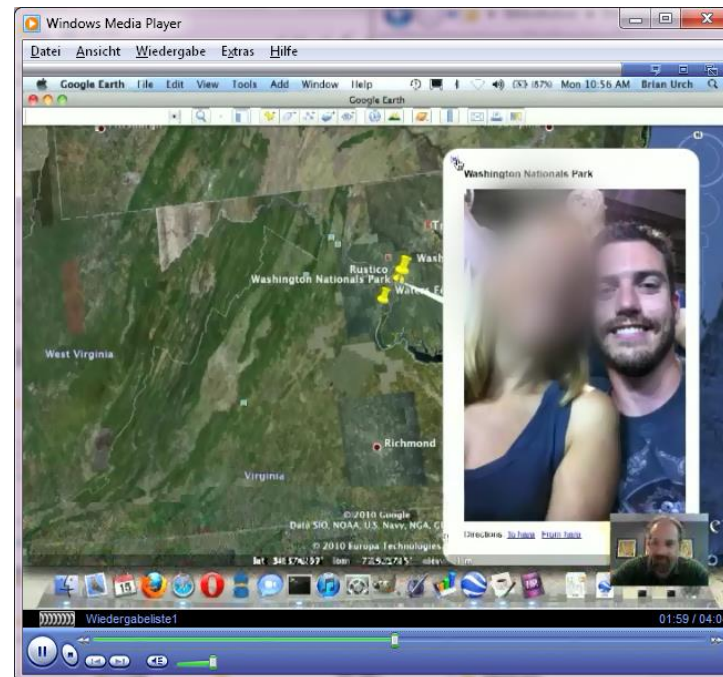  - Data, formats, schema, semantics change

# Big and Small

- Big Data can be very small
  - Streaming data from aircraft sensors
  - Hundred thousand sensors on an aircraft is "big data"
  - Each producing an eight byte reading every second
  - Less than 3GB of data in an hour of flying
    - (100,000 sensors x 60 minutes x 60 seconds x 8 bytes).
- Not all large datasets are "big".
  - Video streams plus metadata
  - Telco calls and internet connections
  - Can be parsed extremely quickly if content is well structured.
  - From http://mike2.openmethodology.org/wiki/Big_Data_Definition
- The task at hand makes data "big".

# „Big data" in business

- Has been used to sell more hardware and software.
- Has become a shallow buzzword.

- But: The actual big data is there, has added-value, and can be used effectively.
  - Data mining
  - Marketing / advertising
  - Collaborative filtering
  - Raytheon's RIOT software
  - NSA, etc.
  - Kreditech, Lenddo, Klout, …
  - …



Felix Naumann
Data Profiling
Summer 2017

19

# „Big data" in business

- Amazon.com
  - Millions of back-end operations every day
  - Catalog, searches, clicks, wish lists, shopping carts, third-party sellers, …
- Walmart
  - \> 1 million customer transactions per hour
  - 2.5 petabytes (2560 terabytes)
- Facebook
  - 250 PB, 600TB added daily (2013)
  - 1 billion photos on one day (Halloween)
- FICO Credit Card Fraud Detection
  - Protects 2.1 billion active accounts

**Felix Naumann
Data Profiling
Summer 2017**

# Big Government Data (USA)

- **Big Data Research and Development Initiative**
  - Explored how big data addresses important problems facing the government.
  - 84 different big data programs spread across six departments
- **Data.gov**
  - > 104.000 datasets
- **Government owns six of the ten most powerful supercomputers in the world.**
- **NASA Center for Climate Simulation**
  - 32 petabytes of climate observations and simulations

**Topics**

A-Z    1-9                                    Clear All

Manufacturing (70)

Ecosystems (75)

Climate (108)

Law (120)

World Wide Human Ge... (145)

Education (147)

BusinessUSA (218)

Agriculture (223)

Research (227)

Finance (235)

Safety (327)

Consumer (329)

Ocean (376)

Energy (623)

**Felix Naumann
Data Profiling
Summer 2017**

# Examples from Wikipedia – Big Science

- Large Hadron Collider
  - □ 150 million sensors; 40 million deliveries per second
  - □ 600 million collisions per second
  - □ Theoretically: 500 exabytes per day (500 quintillion bytes)
  - □ Filtering: 100 collisions of interest per second
    - – Reduction rate of 99.999% of these streams
  - □ 25 petabytes annual rate before replication (2012)
  - □ 200 petabytes after replication

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Examples from Wikipedia - Science

- Sloan Digital Sky Survey (SDSS)
  - Began collecting astronomical data in 2000
  - Amassed more data in first few weeks than all data collected in the history of astronomy.
  - 200 GB per night
  - Stores 140 terabytes of information
  - Large Synoptic Survey Telescope, successor to SDSS
    - Online in 2016
    - Will acquire that amount of data every five days.
- Human genome
  - Originally took 10 years to process;
  - Now it can be achieved in one day.

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Big Data = Science?

- The End of Theory: The Data Deluge Makes the Scientific Method Obsolete (Chris Anderson, Wired, 2008)
  - □ All models are wrong, but some are useful.  (George Box)
  - □ All models are wrong, and increasingly you can succeed without them. (Peter Norvig, Google)

- Before Big Data: Correlation is not causation!
- With Big Data: Who cares?
  - □ Traditional approach to science — hypothesize, model, test — is becoming obsolete.

**Felix Naumann**
**Data Profiling**
**Summer 2017**

- Petabytes allow us to say: "**Correlation is enough.**"

http://www.wired.com/science/discoveries/magazine/16-07/pb_theory

# Correlation vs. Causation



Quelle: Spurious correlations (www.tylervigen.com)

# Correlation vs. Causation



Internet Explorer vs Murder Rate

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Addressing Big Data: Parallelization

- Long tradition in databases
- Vertical and horizontal partitioning
- Shared nothing
- Each machine runs same single-machine program

- Other trends
  - Map/Reduce / Hadoop
  - Multicore CPUs
  - GPGPUs

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Levels of Parallelism on Hardware

- Instruction-level Parallelism
  - Single instructions are automatically processed in parallel
  - Example: Modern CPUs with multiple pipelines and instruction units.
- Data Parallelism
  - Different data can be processed independently
  - Each processor executes the same operations on its share of the input data.
  - Example: Distributing loop iterations over multiple processors
  - Example: GPU processing
- Task Parallelism
  - Different tasks are distributed among the processors/nodes
  - Each processor executes a different thread/process.
  - Example: Threaded programs.

Data profiling and data cleansing are prerequisites for all of these!

- Industry keynote speakers on credit ratings using big data
  - "If the data is out there, we will find it."
  - "… and that is why I closed my Twitter account."
  - "… and that is why I had my son close his Twitter account."

# Overview

1. Introduction to research group
2. Lecture organisation
3. (Big) data
   - **Data sources**
   - Profiling
4. Overview of semester

**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Open vs. closed source

## Open

- Linked data
  - http://linkeddata.org/
- Government data
  - data.gov, data.gov.uk
  - Eurostat
- Scientific data
  - Genes, proteins, chemicals
  - Scientific articles
  - Climate
  - Astronomy
- Published data
  - Tweet (limited)
  - Crawls
- Historical data
  - Stock prices

## Closed

- Transactional data
  - Music purchases
  - Retail-data
- Social networks
  - Tweets, Facebook data
  - Likes, ratings
- E-Mails
- Web logs
  - Per person
  - Per site
- Sensor data
- Military data

# The Linking Open Data cloud diagram

Linking Open Data cloud diagram 2017, by
Andrejs Abele, John P. McCrae, Paul Buitelaar,
Anja Jentzsch and Richard Cyganiak.
http://lod-cloud.net/

33

# Wikipedia Infoboxes

```
{{Infobox company
|name            = International Business Machines Corporation
|logo            = <br />[[File:IBM logo.svg|200px]]<br />
|caption         = Logo since 1972, designed by [[Paul Rand]]
|type            = [[Public company|Public]]
|traded_as       = {{New York Stock Exchange|IBM}}<br />[[Dow Jones Industrial Average|Dow Jones Compon
Component]]
|industry        = [[Personal computer hardware|Computer hardware]], [[Software|Computer software]], [[
services]], [[Information technology consulting|IT consulting]]
|products        = [[List of IBM products|See IBM products]]
|founder         = [[Charles Ranlett Flint]]
|foundation      = [[Endicott, New York|Endicott]], New York, U.S.<br />({{Start date|1911|06|16}})
|location_city   = [[Armonk, New York|Armonk]], New York
|location_country = U.S.
|area_served     = Worldwide
|key_people      = [[Ginni Rometty]]<br />({{small|Chairman, President, and CEO}})
|revenue         = {{Increase}} US$ 106.91 [[1000000000 (number)|billion]] <small>(2011)</small><ref na
|url=http://rcpmag.com/articles/2012/01/20/intel-ibm-exceed-earnings-estimates-google-falls-short.aspx|
International Business Machines Corporation |work=United States Securities and Exchange Commission}}</re
|operating_income = {{Increase}} US$ {{0|0}}20.28 billion <small>(2011)</small><ref name=10K/>
|net_income      = {{Increase}} US$ {{0|0}}15.85 billion <small>(2011)</small><ref name=10K/>
|assets          = {{Increase}} US$ 116.43 billion <small>(2011)</small><ref name=10K/>
|equity          = {{Decrease}} US$ {{0|0}}20.13 billion <small>(2011)</small><ref name=10K/>
|num_employees   = 433,362 <small>(2012)</small><ref name="Fortune 500: IBM employees"/>
```

**Felix Naumann**
**Data Profiling**
**Summer 2017**

34

# DBpedia statistics

- From 125 languages of Wikipedia
- 3 billion triples
  - 580 million English
- English DBpedia
  - 4.6 million things
  - 1,445,000 persons
  - 735,000 places
  - 411,000 creative works
  - 241,000 organizations
  - 251,000 species
  - ...



**Felix Naumann**
**Data Profiling**
**Summer 2017**

- http://wiki.dbpedia.org/about/facts-figures

# And more sources

- Government data
  - www.data.gov
    (380k data sets)
  - data.gov.uk (9k)
  - ec.europa.eu/eurostat
- Finance / business data
- Scientific databases
  - www.uniprot.org
  - skyserver.sdss.org
- The Web
  - HTML tables and lists: billions
  - General sources: Dbpedia (3.7m), freebase (23m), microformats…
  - Domain-specific sources: IMDB,
    Gracenote, isbndb, …



Browse Raw Datasets

| | Name | Popularity | Type |
|---|---|---|---|
| 1. | Worldwide M1+ Earthquakes, Past 7 Days   Geography and Environment   ANSS, geologist, plate, real time, environment, …  Real-time, worldwide earthquake list for the past 7 days | 167,711 views | |
| 2. | U.S. Overseas Loans and Grants (Greenbook)   Foreign Commerce and Aid   foreign assistance, economic assistance, Greenbook, …  These data are U.S economic and military assistance by country from 1946 to 2010. | 62,348 views | |
| 3. | CMS Medicare and Medicaid EHR Incentive Program, electronic health record products used for attestation   Science and Technology   electronic health record, …  Data set merges information about the Centers for Medicare and Medicaid Services, | 34,285 views | |
| 4. | Federal Data Center Consolidation Initiative (FDCCI) Data Center Closings 2010-2013   Federal Government Finances and Employment   fddci, …  Federal Data Center Consolidation Initiative (FDCCI) Data Center Closings 2010-2013 | 32,648 views | |
| 5. | TSCA Inventory   Geography and Environment   new chemicals, manufactured chemicals, …  This dataset consists of the non confidential identities of chemical substances | 27,007 views | |
| 6. | Data.gov Catalog   Other   dataset, metadata, catalog, data extraction tool, …  An interactive dataset containing the metadata for the Data.gov raw datasets and tools | 23,117 views | |
| 7. | US DOE/NNSA Response to 2011 Fukushima Incident: Radiological Air Samples   Geography and Environment   radiation, Japan, nuclear, Tohoku, …  Field Samples are physical media collected during the response which are | 22,458 views | |
| 8. | US DOE/NNSA Response to 2011 Fukushima Incident: Field Team Radiological Measurements   Geography and Environment   Japan, nuclear, Tohoku, radiation, …  Field Measurements describe &alpha; and &beta; activity and &gamma; exposure rate. | 20,940 views | |
| 9. | Federal Executive Branch Internet Domains   Federal Government Finances and Employment   .gov, domains, agencies, federal, registered  Listing of Federal Agency Internet Domains (This list is updated bi-weekly to reflect the | 17,267 views | |

# Getting the data

- Download
  - Data volumes make this increasingly infeasible
  - Fedex HDDs
  - Fedex tissue samples instead of sequence data
- Generating big (but synthetic) data
  1. Automatically insert interesting features and properties
  2. Then „magically" detect them
- Sharing data
  - Repeatability of experiments
  - Not possible for commercial organizations

# Pathologies of Big Data

- Store basic demographic information about each person
  - □ **age, sex, income, ethnicity, language, religion, housing status, location**
  - □ Packed in a 128-bit record
- World population: 6.75 billion rows, 10 columns, 128 bit each
  - □ About 150 GB
- What is the median age by sex for each country?
  - □ Algorithmic solution
    - – 500$ Desktop: I/O-bound; 15min reading the table
    - – 15,000$ Server with RAM: CPUI-bound; <1min
  - □ Database solution
    - – Aborted bulk load to PostgreSQL – disk full
      (bits vs. integer and DBMS inflation)
  - □ Small database solution (3 countries, 2% of data)
    - – **SELECT country,age,sex,count(*)
      FROM people GROUP BY country,age,sex;**
    - – *> 24h, because of poor analysis: Sorting instead of hashing*
    - – "PostgreSQL's difficulty here was in **analyzing** [=profiling] the stored data, not in storing it."
      - From http://queue.acm.org/detail.cfm?id=1563874

**Felix Naumann
Data Profiling
Summer 2017**

Visualization of edits by user „Pearle"

http://en.wikipedia.org/wiki/File:Viegas-UserActivityonWikipedia.gif

39

# Overview

1. Introduction to research group
2. Lecture organisation
3. (Big) data
   - Data sources
   - **Profiling**
4. Overview of semester



**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Definition Data Profiling

- Data profiling is the process of examining the data available in an existing data source [...] and collecting statistics and information about that data.

  - Wikipedia 03/2013

- Data profiling refers to the activity of creating small but informative summaries of a database.

  - Ted Johnson, Encyclopedia of Database Systems

- Data profiling vs. data mining
  - Data profiling gathers technical metadata to support data management
  - Data mining and data analytics discovers non-obvious results to support business management
  - Data profiling results: information about columns and column sets
  - Data mining results: information about rows or row sets (clustering, summarization, association rules, etc.)

- Define as a set of data profiling tasks / results

# Classification of Profiling Tasks

42

# Use Cases for Profiling

- Query optimization
  - Counts and histograms
- Data cleansing
  - Patterns and violations
- Data integration
  - Cross-DB inclusion dependencies
- Scientific data management
  - Handle new datasets
- Data analytics and mining
  - Profiling as preparation to decide on models and questions
- Database reverse engineering

- Data profiling as preparation for any other data management task

**Felix Naumann**
**Data Profiling**
**Summer 2017**

43

# Challenges of (Big) Data Profiling
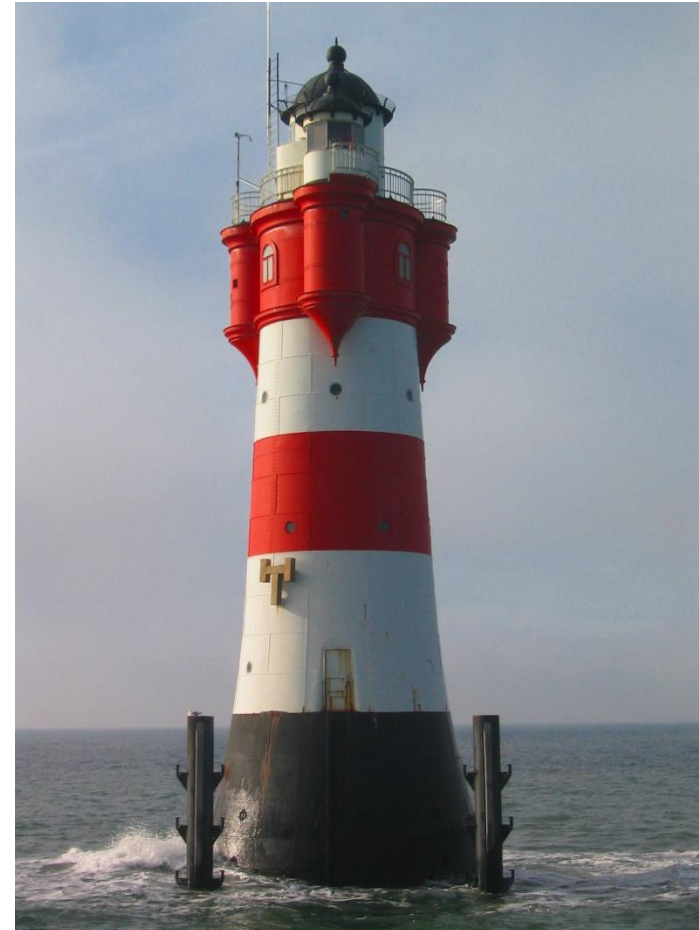
- Computational complexity
  - Number of rows
    - Sorting, hashing
  - Number of columns
    - Number of column combinations
- Large solution space
- I/O-bound due to large data sets and distribution

- New data types (beyond strings and numbers)
- New data models (beyond relational): RDF, XML, etc.
- New requirements
  - User-oriented
  - Streaming
  - Etc. – see next slide set

# Large solution space



- Size of lattice: $2^n - 1$ (empty set not considered)
- Nodes at level 1: n
- Nodes at level n: 1
- Nodes at level k: $\binom{n}{k} = \frac{n!}{(n-k)!k!}$
- Largest level at n/2: $\binom{n}{n/2} = \frac{n!}{\left(\frac{n}{2}!\right)^2}$

$$
\begin{array}{ccccccccccccccc}
& & & & & & & 1 & & & & & & & \\
& & & & & & 1 & & 1 & & & & & & \\
& & & & & 1 & & 2 & & 1 & & & & & \\
& & & & 1 & & 3 & & 3 & & 1 & & & & \\
& & & 1 & & 4 & & 6 & & 4 & & 1 & & & \\
& & 1 & & 5 & & 10 & & 10 & & 5 & & 1 & & \\
& 1 & & 6 & & 15 & & 20 & & 15 & & 6 & & 1 & \\
1 & & 7 & & 21 & & 35 & & 35 & & 21 & & 7 & & 1 \\
\end{array}
$$

- $\binom{n}{k} \in \Theta(n^k)$ , i.e., bounded from above and below

1. Introduction to research group
2. Lecture organisation
3. (Big) data
   ■ Data sources
   ■ Profiling
4. **Overview of semester**



**Felix Naumann**
**Data Profiling**
**Summer 2017**

# Schedule

- Big Data (today) and Data Profiling Introduction
- Data Stuctures
  - Lattices, Apriori traversal, complexity, agree-sets/evidence sets, PLIs, Bloom filters
- Unique Column Combinations (UCCs)
  - A-Priori, DUCC, HCA?, Gordian?, Swan?
- Scientific experiments
- Functional Dependencies (FDs)
  - TANE, FD-Mine, FDep, HyFD + approximate TANE
  - FD-measures for ranking (g1-3, support & confidence)
- Inclusion Dependencies (INDs)
  - Spider, Binder, Find2, zigzag, SINDY?, MANY, cINDs?
- Semantics
  - Key & FK detection, normalization, interpretation

**Felix Naumann**
**Data Profiling**
**Summer 2017**

47

# Schedule

- Order Dependencies (ODs)
  - ORDER, Szlichta-paper
- Denial Constraints
  - FastDC, Hydra
- Data Synopses
- Column Uniqueness: Approximately counting number of unique values
- Approximation
  - partial, conditional, approximate as concepts, selected approaches
- RDF-Profiling (optional)
- Outlook

- Guest lectures
  - Thomas Bläsius
  - Giuseppe Polese

**Felix Naumann**
**Data Profiling**
**Summer 2017**

48