

Chapter 01

Introduction

Dr. Steffen Herbold
herbold@cs.uni-goettingen.de

Outline

- Introduction to Big Data
- Data Science and Business Intelligence
- The Skillset of Data Scientists
- Summary

What is „Big Data“?!?

Is this really
about size?



Naive Definition

- Naive definition:
 - Big data only depends on the data size
 - 1 Gigabyte? 1 Terabyte? 1 Petabyte?
- Naive interpretation misses important aspects
 - Time:
 - Analyzing 1 Gigabyte of data per day is different from analyzing 1 Gigabyte of data per second
 - Diversity:
 - Analyzing spread sheets with numeric data is different from analyzing Web pages that contain a mixture of text and images
 - Distribution:
 - Analyzing data from a single source is different from analyzing data from multiple sources

Definition of Big Data

- Following Gartner's IT Glossary:
 - Big data is high-**volume**, high-**velocity** and/or high-**variety** information assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and process automation.

- The three Vs

- Volume
- Velocity
- Variety



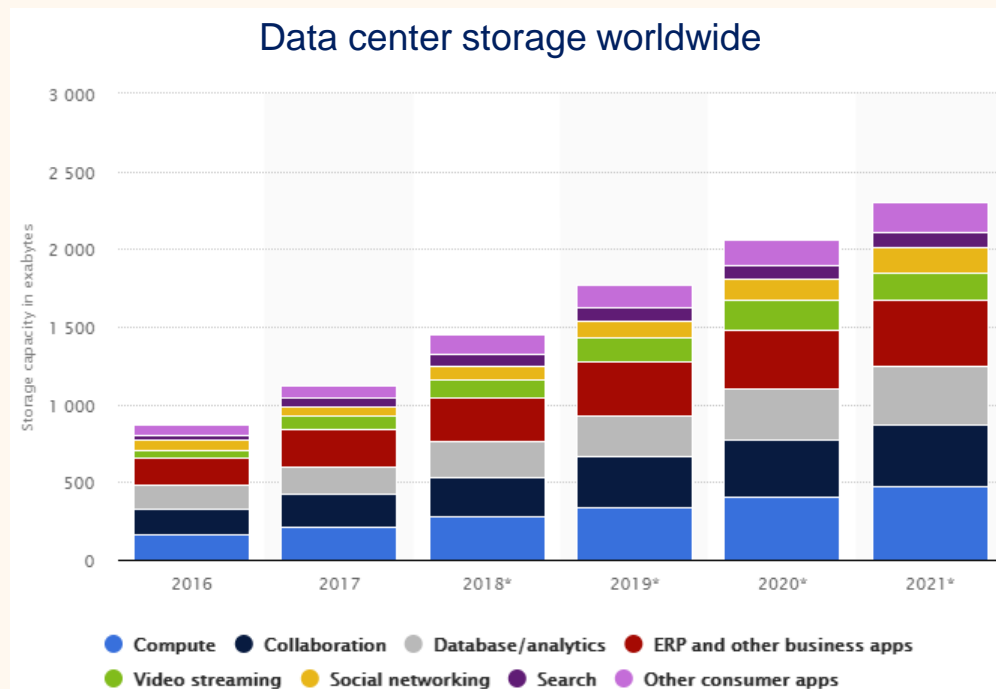
Some people actually use 10 Vs to define big data!

- Variability
- Veracity
- Validity
- Vulnerability
- Volatility
- Visualization
- Value



The 3 Vs: Volume

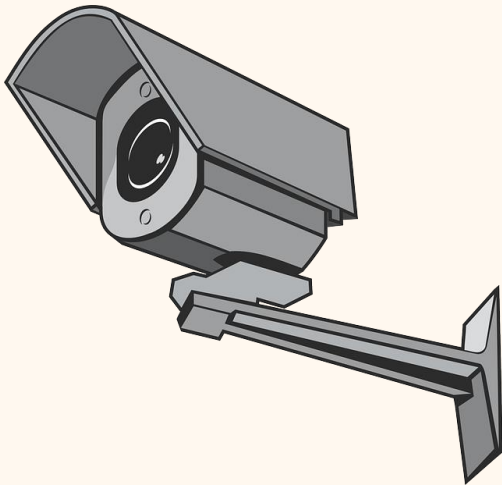
- Scale of the data must be „big“
 - No clear definition
 - „that demand [...] innovative forms of information processing“ (Gartner)



© Statista 2018

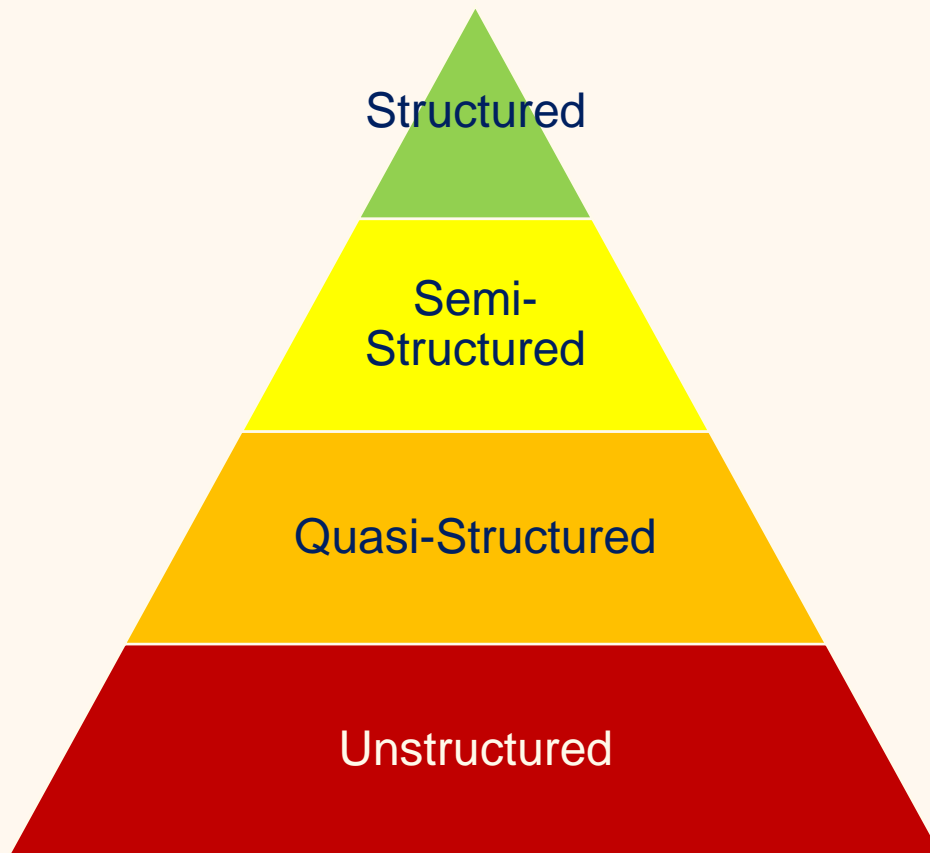
The 3 Vs: Velocity

- Speed at which new data is created
- Speed at which data must be processed and analyzed
 - Often close to real-time



The 3 Vs: Variety

- Diversity in data types and data sources



- Data with defined types and structure
- Example: comma separated values

- Textual data with parseable pattern
- Example: XML files with schema

- Textual data with erratic formats that can be formatted with effort
- Example: Clickstream data

- Data that has no inherent structure, often with multiple formats
- Example: Web site, videos

Outline

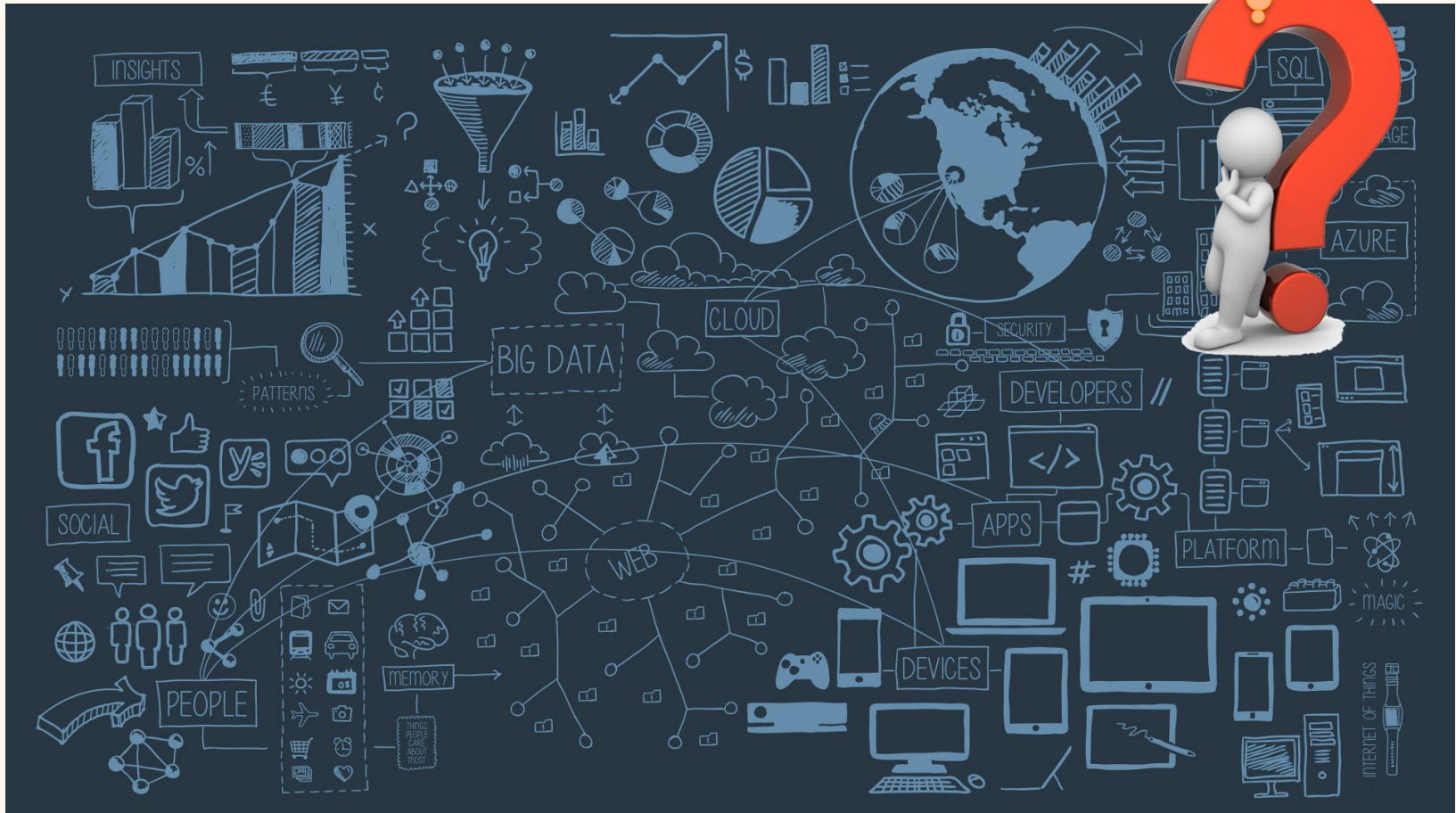
- Introduction to Big Data
- **Data Science and Business Intelligence**
- The Skillset of Data Scientists
- Summary

Defining Data Science

- Unfortunately, there is no clear definition (yet?)
- Goal is the extraction of knowledge from data
- Combination of techniques from different disciplines
- Scientific principles guide the data analysis

What is „Data Science“?!?

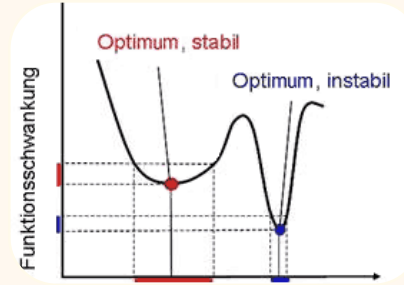
Tools? Big Data?
Machine Learning?



Mathematical Aspects



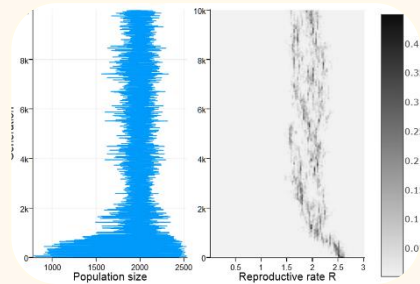
Computational
Geometry



Optimization



Stochastics

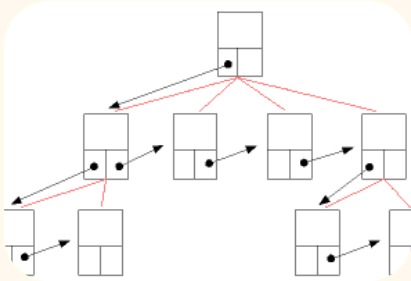


Scientific
Computing

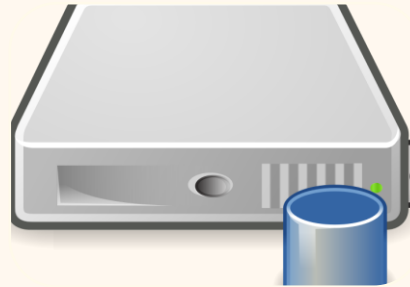


Machine
Learning

Computer Science Aspects



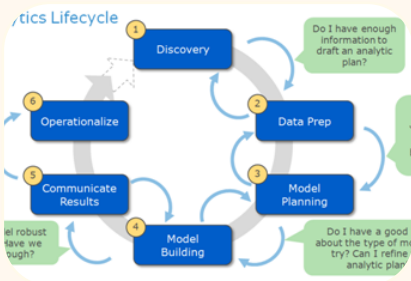
Data Structures and Algorithms



Databases



Distributed Computing



Software Engineering

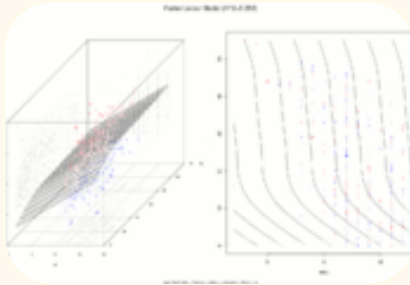


Artificial Intelligence

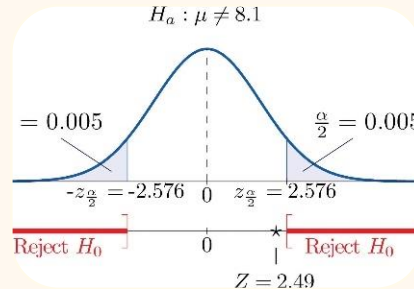


Machine Learning

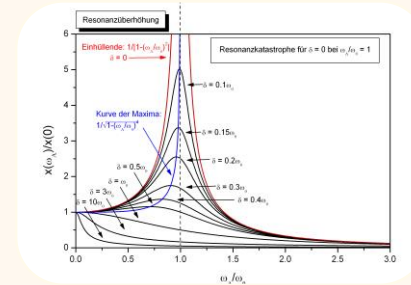
Statistical Aspects



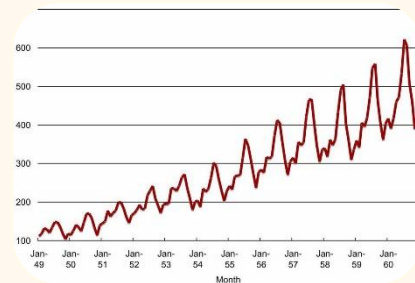
Linear Models



Statistical Tests



Inference



Time Series Analysis



Machine Learning

Applications



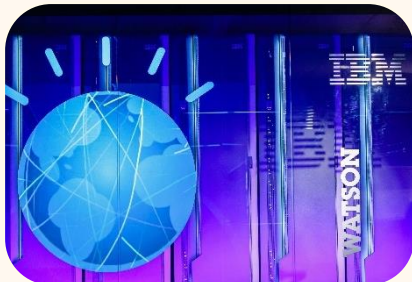
Intelligent Systems



Robotics



Marketing



Medicine



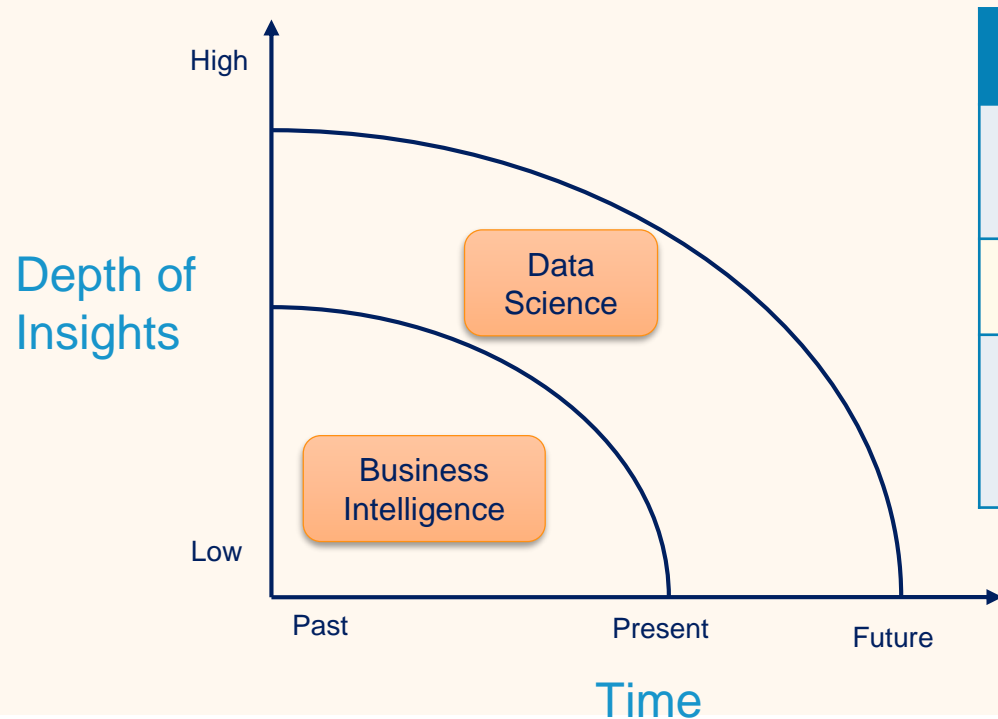
Autonomous Driving



Social Networks

Data Science vs. Business Intelligence

- Business Intelligence (Gartner IT Glossary)
 - [...] best practices that enable access to and analysis of information to improve and optimize decisions and performance.



	Business Intelligence	Data Science
Techniques	Dashboards, alerts, queries	Optimization, predictive modelling, forecasting
Data Types	Structured, data warehouses	Any kind, often unstructured
Common questions	What happened...? How much did...? When did...?	What if...? What will...? How can we...?

More Data → More Opportunities



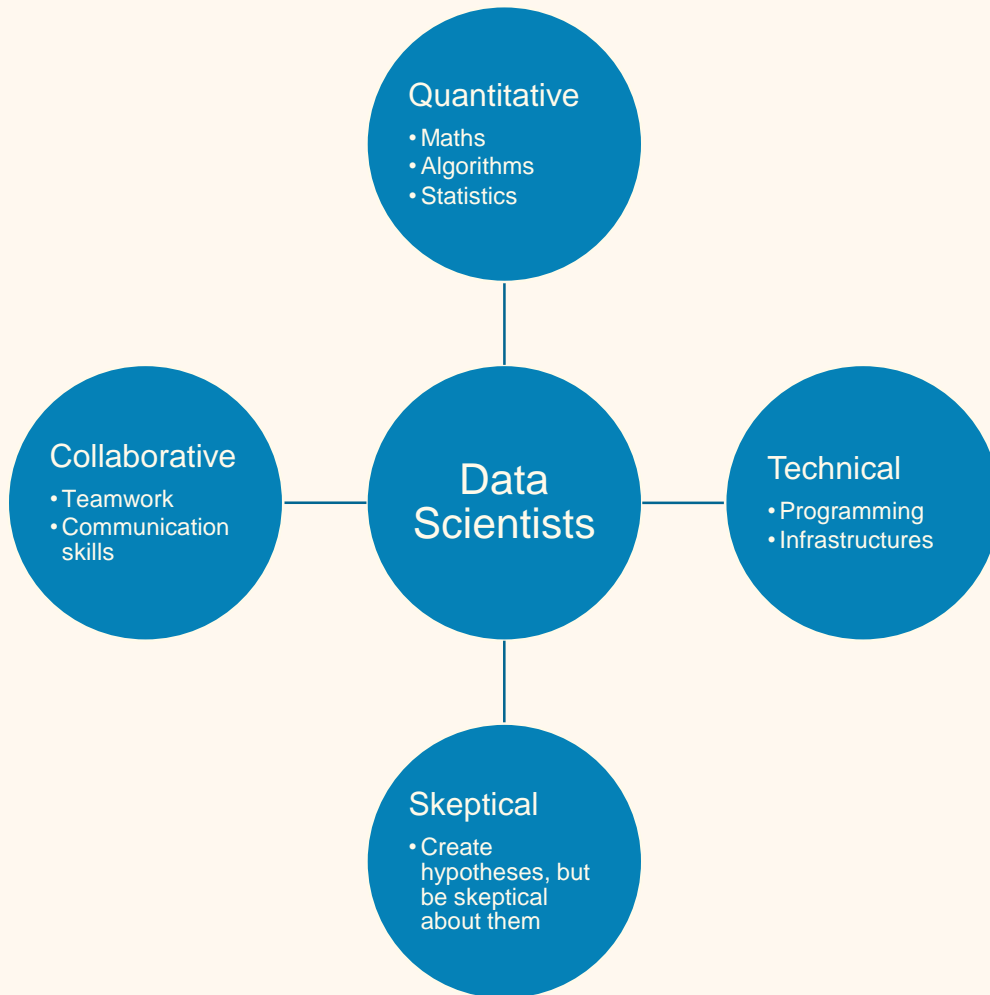
Outline

- Introduction to Big Data
- Data Science and Business Intelligence
- **The Skillset of Data Scientists**
- Summary

What are Data Scientists?

- Not computer scientists
 - But should know about databases, data structures, algorithms, etc.
- Not mathematicians
 - But should know about optimization, stochastics, etc.
- Not statisticians
 - But should know about regression, statistical tests, etc.
- Not domain experts
 - But must work together with them

Skills of Data Scientists



A bit of everything

... but actually as much as possible of everything

Different types of Data Scientists

- According to Microsoft Research:

- Polymath
 - „Do it all“
- Data Evangelist
 - Data analysis, disseminating and acting on insights
- Data Preparer
 - Querying existing data, preparing data for analysis
- Data Shapers
 - Analyzing and preparing data
- Data Analyzer
 - Analyzing data
- Platform Builder
 - Collect data and create infrastructures
- Moonlighters (50%/20%)
 - „Spare time“ data scientists
- Insight Actors
 - Use the outcome and act on insights.

Miyung Kim, Thomas Zimmermann, Robert DeLine, Andrew Begel: Data Scientists in Software Teams: State of the Art and Challenges, IEEE Transactions on Software Engineering (Online First)

Outline

- Introduction to Big Data
- Data Science and Business Intelligence
- The Skillset of Data Scientists
- **Summary**

Summary

- Big data has a high volume, velocity, and variety
 - Different data structures
 - Structured, semi-structured, quasi-structured, unstructured
 - Data science is a very diverse discipline
 - Maths, computer science, statistics, applications
- Data scientists require a diverse skillset