

Introduction

Lecture BigData Analytics

Julian M. Kunkel

julian.kunkel@googlemail.com

University of Hamburg / German Climate Computing Center (DKRZ)

16-10-2015



Outline

1 Introduction

2 BigData Challenges

3 Analytical Workflow

4 Use Cases

5 Programming

6 Summary

About DKRZ

German Climate Computing Center (DKRZ)



Partner for Climate Research
Maximum Compute Performance.
Sophisticated Data Management.
Competent Service.

Scientific Computing

- Research Group of Prof. Ludwig at the University of Hamburg
- Embedded into DKRZ



Research

- Analysis of parallel I/O
- I/O & energy tracing tools
- Middleware optimization
- Alternative I/O interfaces
- Data reduction techniques
- Cost & energy efficiency

Lecture

Concept of the lecture

- The lecture is focussing on applying technology and some theory
- Theory
 - Data models and processing concepts
 - Algorithms and data structures
 - System architectures
 - Statistics and machine learning
- Applying technology
 - Learning about various state-of-the art technology
 - Hands-on for understanding the key concepts
 - Languages: Java, Python, R
- The domain of big data is overwhelming, especially in terms of technology
- It is a crash course for several topics such as statistics and databases
- ⇒ it is not the goal to learn and understand every aspect in this lecture

Lecture (2)

Slides

- Many openly accessible sources have been used
- Citation to them by a number
- The reference slide provides the link to the source
- For figures, a reference is indicated by *Source: [Author]¹ [title] [ref]*
- In the title, an [ref] means that this reference has been used for the slide, some text may be taken literally

Excercise

- Weekly delivery, processing time about 8 hours / per week estimated
- Teamwork of 2 or 3 people (groups are mandatory!)
- Supported by: Hans Ole Hatzel

¹If available

Idea of BigData

Methods of obtaining knowledge (Erkenntnissprozess)

Theory (model), hypothesis, experiment, analysis (repeat)

- Explorative: start theory with observations of phenomena
- Constructivism: starts with axioms and reason implications

The Fourth Paradigm

- (Big) Data + Analytics ⇒ Insight (prediction of the future)
 - For industry: insight = business advantage and money...
- Analytics: follow an explorative approach and study the data
 - To infer knowledge, use statistics / machine learning
- Construct a theory (model) and validate it with the data

Example Models

Similarity is a (very) simplistic model and predictor for the world

- Humans use this approach in their cognitive process
- Uses the advantage of BigData

Weather prediction

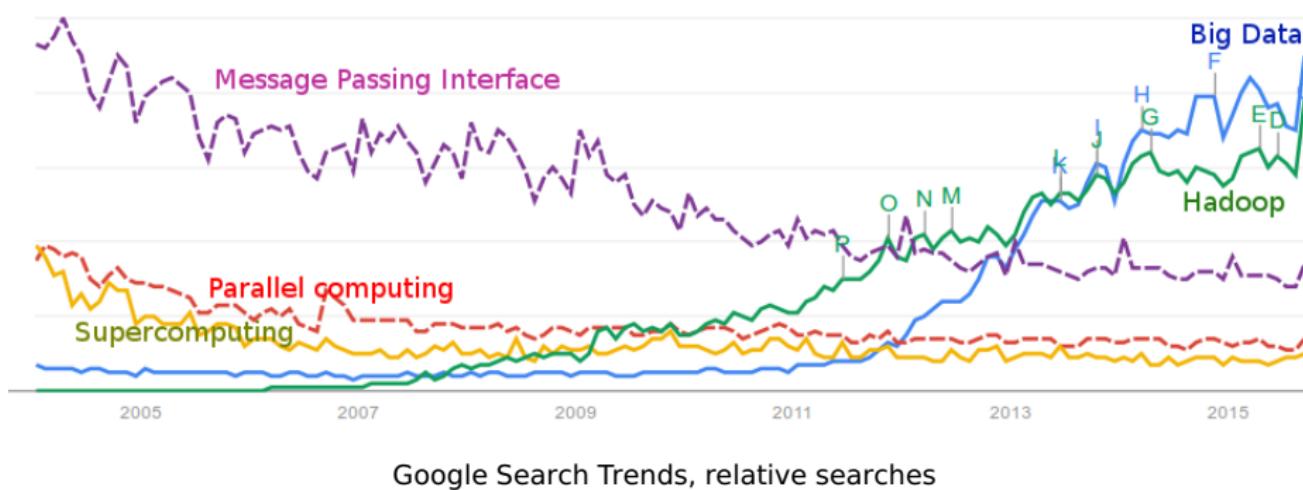
- You may develop and rely on complex models of physics
- Or use a simple model for a particular day; e.g. expect it to be similar to the weather of the day over the last X years
- Used by humans: rule of thumb for farmers

Preferences of Humans

- Identify a set of people which liked items you like
- Predict you like also the items those people like (items you haven't rated so far)

Relevance of Big Data

- Big Data Analytics is emerging
- Relevance increases compared to supercomputing



1 Introduction

2 BigData Challenges

- Volume
- Velocity
- Variety
- Veracity
- Value

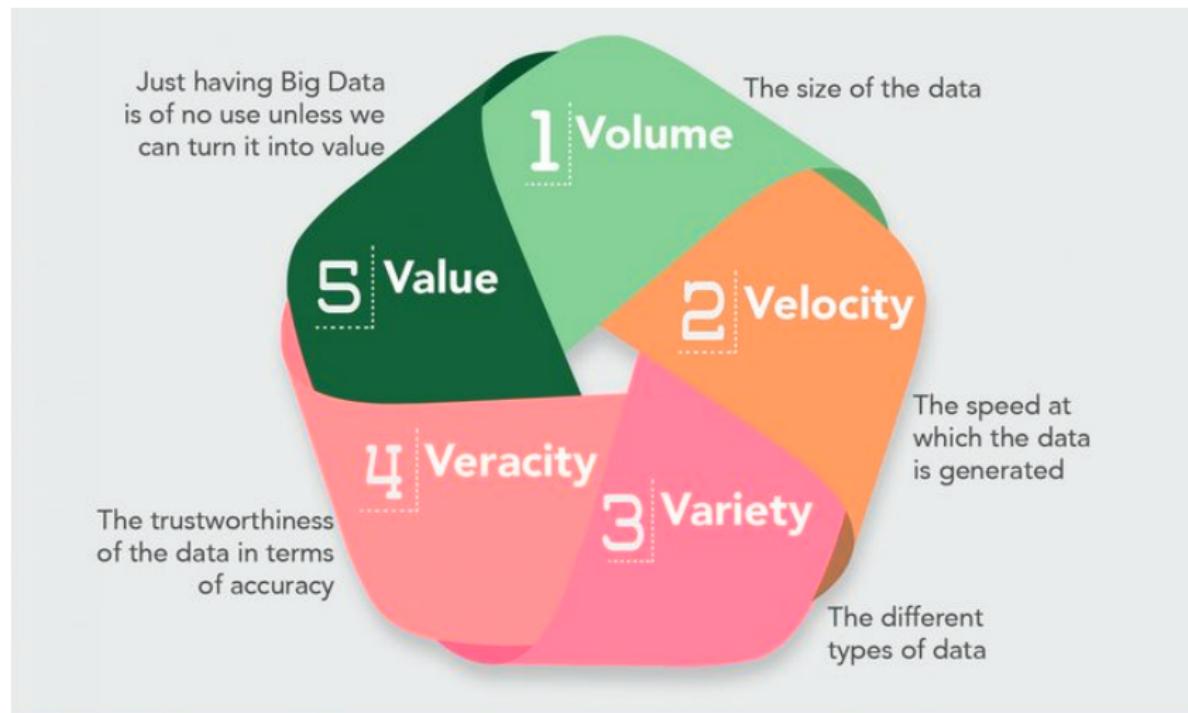
3 Analytical Workflow

4 Use Cases

5 Programming

6 Summary

BigData Challenges & Characteristics



Source: MarianVesper [4]

Volume: The size of the Data

What is Big Data

Terrabytes to 10s of petabytes

What is not Big Data

A few gigabytes

Examples

- Wikipedia corpus with history ca. 10 TByte
- Wikimedia commons ca. 23 TByte
- Google search index ca. 46 Gigawebpages²
- YouTube per year 76 PByte (2012³)

²<http://www.worldwidewebsize.com/>

³<https://sumanrs.wordpress.com/2012/04/14/youtube-yearly-costs-for-storagenetworking-estimate/>

Velocity: Data Volume per Time

What is Big Data

30 KiB to 30 GiB per second
(902 GiB/year to 902 PiB/year)

What is not Big Data

A never changing data set

Examples

- LHC (Cern) with all experiments about 25 GB/s ⁴
- Square Kilometre Array 700 TB/s (in 2018) ⁵
- 50k Google searches per s ⁶
- Facebook 30 Billion content pieces shared per month ⁷

⁴<http://home.web.cern.ch/about/computing/processing-what-record>

⁵<http://venturebeat.com/2014/10/05/how-big-data-is-fueling-a-new-age-in-space-exploration/>

⁶<http://www.internetlivestats.com/google-search-statistics/>

⁷<https://blog.kissmetrics.com/facebook-statistics/>

Data Sources

Enterprise data

- Serves business objectives, well defined
- Customer information
- Transactions, e.g. Purchases

Experimental/Observational data (EOD)

- Created by machines from sensors/devices
- Trading systems, satellites
- Microscopes, video streams, Smart meters

Social media

- Created by humans
- Messages, posts, blogs, Wikis

Variety: Types of Data

■ Structured data

- Like tables with fixed attributes
- Traditionally handled by relational databases

■ Unstructured data

- Usually generated by humans
- E.g. natural language, voice, Wikipedia, Twitter posts
- Must be processed into (semi-structured) data to gain value

■ Semi-structured data

- Has some structure in tags but it changes with documents
- E.g. HTML, XML, JSON files, server logs

What is Big Data

- Use data from multiple sources and in multiple forms
- Involve unstructured and semi-structured data

Veracity: Trustworthiness of Data

What is Big Data

- Data involves some uncertainty and ambiguities
- Mistakes can be introduced by humans and machines
 - People sharing accounts
 - Like sth. today, dislike it tomorrow
 - Wrong system timestamps

Data Quality is vital!

Analytics and conclusions rely on good data quality

- Garbage data + perfect model => garbage results
- Perfect data + garbage model => garbage results

GIGO paradigm: *Garbage In – Garbage Out*

Value of Data

What is Big Data

- Raw data of Big Data is of low value
 - For example, single observations
- Analytics and theory about the data increases the value

Analytics transform big data into smart data!

Types of Data Analytics and Value of Data

- 1 Descriptive analytics (Beschreiben)**
 - “What happened?”
- 2 Diagnostic analytics**
 - “Why did this happen, what went wrong?”
- 3 Predictive analytics (Vorhersagen)**
 - “What will happen?”
- 4 Prescriptive analytics (Empfehlen)**
 - “What should we do and why?”

The level of insight and value of data increases from step 1 to 4

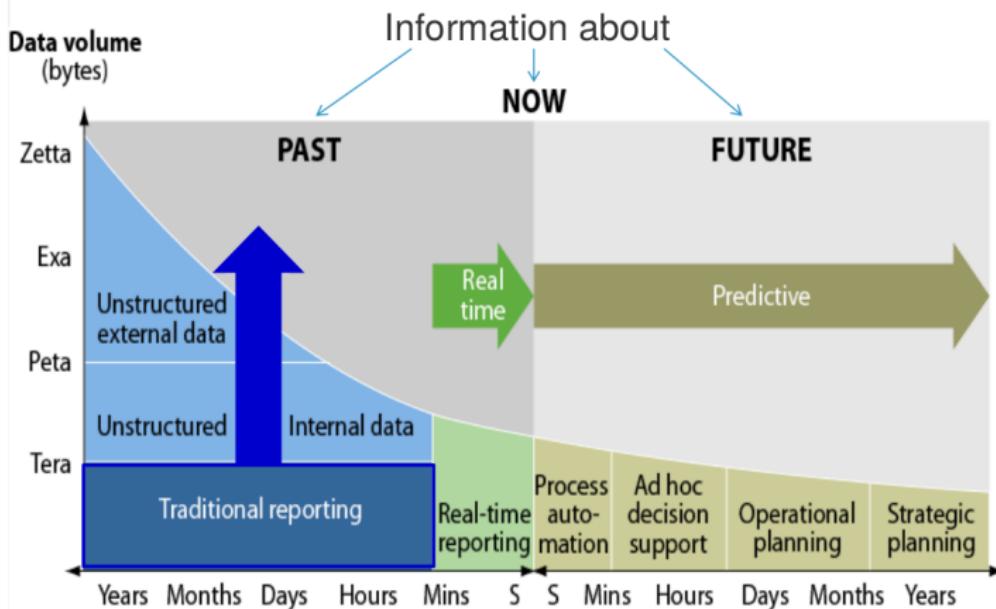
The Value of Data (alternative view)



Source: Dursun Delen, Haluk Demirkhan [9]

The Value of Data (alternative view 2)

Most BI remains backward-looking



Source: Forrester report. Understanding The Business Intelligence Growth Opportunity.
20-08-2011

1 Introduction

2 BigData Challenges

3 Analytical Workflow

- Value Chain
- Roles
- Privacy

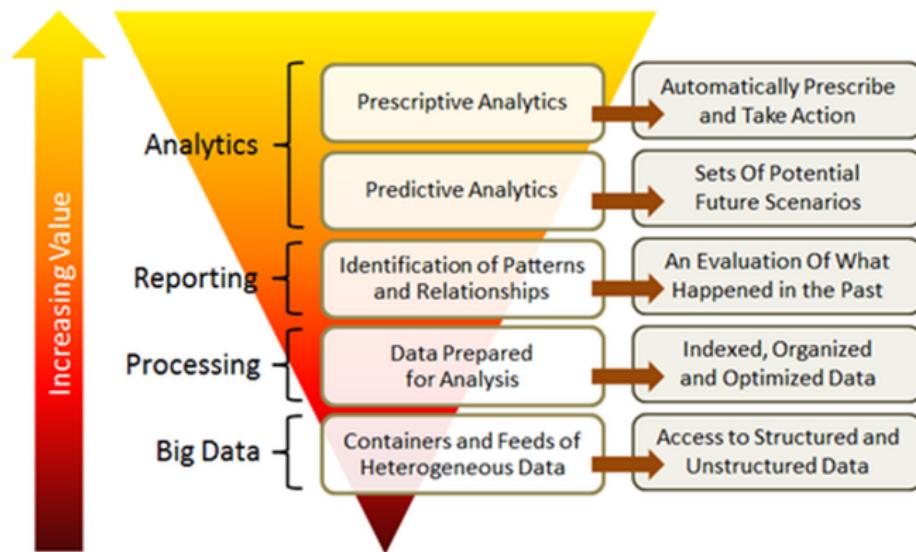
4 Use Cases

5 Programming

6 Summary

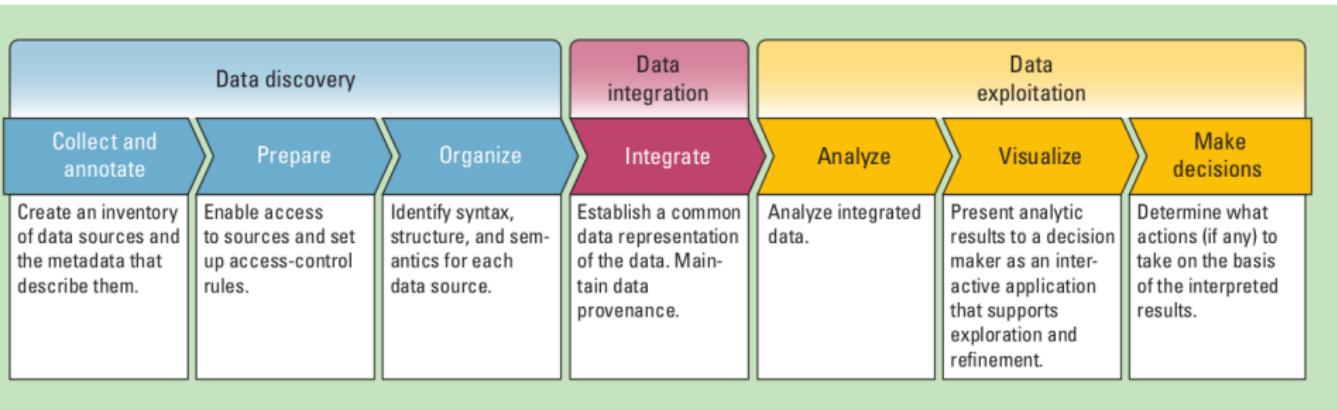
Big Data Analytics Value Chain

There are many visualizations of the processing and value chain [8]



Source: Andrew Stein [8]

Big Data Analytics Value Chain (2)



Source: Miller and Mork [7]

Roles in the Big Data Business

Data scientist

Data science is a systematic method dedicated to knowledge discovery via data analysis [1]

- In business, optimize organizational processes for efficiency
- In science, analyze experimental/observational data to derive results

Data engineer

Data engineering is the domain that develops and provides systems for managing and analyzing big data

- Build modular and scalable data platforms for data scientists
- Deploy big data solutions

Typical Skills

Data scientist

- Statistics + (Mathematics)
- Computer science
 - Programming e.g.: Java, Python, R, (SAS, ...)
 - Machine learning
- Some domain knowledge for the problem to solve

Data engineer

- Computer science
 - Databases
 - Software engineering
 - Massively parallel processing
 - Real-time processing
- Languages: C++, Java, Python
- Understand performance factors and limitations of systems

Data Science vs. Business Intelligence (BI)

Characteristics of BI

- Provides pre-created dashboards for management
 - Repeated visualization of well known analysis steps
- Deals with structured data
- Typically data is generated within the organization
- Central data storage (vs. multiple data silos)
- Handled well by specialized database techniques

Typical types of insight

- Customer service data: “what business causes the largest customer wait times”
- Sales and marketing data: “which marketing is most effective”
- Operational data: “efficiency of the help desk”
- Employee performance data: “who is most/least productive”

Privacy

Be aware of privacy issues if you deal with personal/private information.
German privacy laws are more strict than those of other countries

Ziel des Datenschutzes

Recht auf informationelle Selbstbestimmung

- Schutz des Einzelnen vor beeinträchtigung des Persönlichkeitsrechts durch den Umgang mit seinen personenbezogenen⁸ Daten
- Besonderer Schutz für Daten über Gesundheit, ethnische Herkunft, religiöse, gewerkschaftschliche oder sexuelle Orientierung

⁸§3 BDSG, Einzelangaben über persönliche oder sachliche Verhältnisse einer bestimmten oder bestimmbaren natürlichen Person

Wichtige Grundsätze des Gesetzes [10]

- Verbotsprinzip mit Erlaubnisvorbehalt
 - Erhebung, Verarbeitung, Nutzung und Weitergabe von personenbezogenen Daten sind verboten
 - Nutzung nur mit Rechtsgrundlage oder mit Zustimmung der Person
- Unternehmen mit 10 Personen benötigen Datenschutzbeauftragten
- Verfahren zur automatischen Verarbeitung sind vom Datenschutzbeauftragten zu prüfen und anzeigenpflichtig
- Sitz der verantwortlichen Stelle maßgeblich
 - Bei einer Niederlassung in D gilt BDSG
- Prinzipien: Datenvermeidung, -sparsamkeit
- Schutz vor Zugriffen, Änderungen und Weitergabe
- Betroffene haben Recht auf Auskunft, Löschung oder Sperrung
- Anonymisierung/Pseudonymisierung: Ist die Zuordnung zu Einzelpersonen (nahezu) ausgeschlossen, so können Daten verarbeitet werden

1 Introduction

2 BigData Challenges

3 Analytical Workflow

4 Use Cases

■ Overview

5 Programming

6 Summary

THE BIG PICTURE ON HADOOP

Apache Hadoop is an open source software framework created in 2005.
Engineered for Big Data and large-scale processing applications.



MOST COMMONLY USED HADOOP SERVICES



TOP APPLICATION TYPES THAT BENEFIT FROM HADOOP



PROBLEM OR OPPORTUNITY?



THE FUTURE OF HADOOP

61%

of organizations plan to deploy Hadoop or have already deployed it

\$50.2B

Worldwide sales based on Hadoop technology are forecasted to reach \$50.2 billion by 2020



This infographic is brought to you by StackIQ (www.stackiq.com), makers of stacki - the fastest open source bare metal installer. Download it at www.stackiq.com.

Apache Hadoop, Hadoop, the logo for Hadoop Apache Hive, Hive, Apache Hbase, Hbase, Apache Pig, and Pig are all trademarks of Apache Software Foundation. All other trademarks are the property of their respective owners.

Sources: TDWI (The Data Warehousing Institute), Solix Technologies (The Current State of Hadoop)
Product and company names mentioned herein may be trademarks and/or registered trademarks of their respective companies.

Source: [21]

Use Cases for BigData Analytics

Increase efficiency of processes and systems

- Advertisement: Optimize for target audience
- Product: Acceptance (like/dislike) of buyer, dynamic pricing
- Decrease financial risks: fraud detection, account takeover
- Insurance policies: Modeling of catastrophes
- Recommendation engine: Stimulate purchase/consume
- Systems: Fault prediction and anomaly detection
- Supply chain management

Science

- Epidemiology research: Google searches indicate Flu spread
- Personalized Healthcare: Recommend good treatment
- Physics: Finding the Higgs-Boson, analyze telescope data
- Enabler for social sciences: Analyze people's mood

Big Data in Industry

INDUSTRY	USE CASE	DATA TYPE							
		Sensor	Server Logs	Text	Social	Geographic	Machine	Clickstream	Structured
Financial Services	New Account Risk Screens	✓	✓						
	Trading Risk		✓						
	Insurance Underwriting	✓		✓	✓				
Telecom	Call Detail Records (CDR)					✓	✓		
	Infrastructure Investment		✓					✓	
	Real-time Bandwidth Allocation	✓	✓	✓	✓				
Retail	360° View of the Customer			✓				✓	
	Localized, Personalized Promotions					✓			
	Website Optimization						✓		
Manufacturing	Supply Chain and Logistics	✓							
	Assembly Line Quality Assurance	✓							
	Crowd-sourced Quality Assurance				✓				
Healthcare	Use Genomic Data in Medical Trials	✓					✓		
	Monitor Patient Vitals in Real-Time							✓	
Pharmaceuticals	Recruit and Retain Patients for Drug Trials				✓			✓	
	Improve Prescription Adherence			✓	✓				✓
Oil & Gas	Unify Exploration & Production Data	✓			✓			✓	
	Monitor Rig Safety in Real-Time	✓						✓	
Government	ETL Offloaded Response to Federal Budgetary Pressures						✓		
	Sentiment Analysis for Government Programs				✓				

Source: [20]

Example Use Case: Deutschland Card [2]

Goals

- Customer bonus card which tracks purchases
- Increase scalability and flexibility
- Previous solution based on OLAP

Big Data Characteristics

- Volume: $O(10)$ TB
- Variety: mostly structured data, schemes are extended steadily
- Velocity: data growth rate $O(100)$ GB / month

Results

- Much better scalability of the solution
- From dashboards to ad-hoc analysis within minutes

Example Use Case: DM [2]

Goals

- Predict required employees per day and store
- Prevent staff changes on short-notice

Big Data Characteristics

- Input data: Opening hours, incoming goods, empl. preferences, holidays, weather ...
- Model: NeuroBayes (Bayes + neuronal networks)
- Predictions: Sales, employee planning
- 450.000 predictions per week

Results

- Daily updated sales per store
- Reliable predictions for staff planning
- Customer and employee satisfaction

Example Use Case: OTTO [2]

Goals

Optimize inventory and prevent out-of-stock situations

Big Data Characteristics

- Input data: product characteristics, advertisement
- Volume/Velocity: 135 GB/week, 300 million records
- Model: NeuroBayes (Bayes + neuronal networks)
- 1 billion predictions per year

Results

- Better prognostics of product sales (up to 40%)
- Real time data analytics

Example Use Case: Smarter Cities (by KTH) [2]

Goals

- Improve traffic management in Stockholm
- Prediction of alternative routes

Big Data Characteristics

- Input data: Traffic videos/sensors, weather, GPS
- Volume/Velocity: 250k GPS-data/s + other data sources

Results

- 20% less traffic
- 50% reduction in travel time
- 20% less emissions

Example Facebook Studies

Insight from [11] by exploring posts

- Young narcissists tweet more likely.
Middle-aged narcissists update their status
- US students post more problematic information than German students
- US Government checks tweets/facebook messages for several reasons
- Human communication graph has an average diameter of 4.74

Manipulation of news feeds [13]

- News feeds have been changed to analyze people's behavior in subsequent posts
- Paper: "Experimental evidence of massive-scale emotional contagion through social networks"

From Big Data to the Data Lake [20]

- With cheap storage costs, people promote the concept of the data lake
- Combines data from many sources and of any type
- Allows for conducting future analysis and not miss any opportunity

Attributes of the data lake

- Collect everything: all data, both raw sources over extended periods of time as well as any processed data
 - Decide during analysis which data is important, e.g. no “schema” until read
- Dive in anywhere: enable users across multiple business units to refine, explore and enrich data on their terms
- Flexible access: enable multiple data access patterns across a shared infrastructure: batch, interactive, online, search, and others

1 Introduction

2 BigData Challenges

3 Analytical Workflow

4 Use Cases

5 Programming

- Java
- Python
- R

6 Summary

Programming BigData Analytics

High-level concepts

- SQL and derivatives
- Domain-specific languages (Cypher, PigLatin)

Programming languages

- Java interfaces are widely available but low-level
- Python and R have connectors to popular BigData solutions

In the exercises, we'll learn and use basics of those languages/interfaces

Introduction to Java

- Developed by Sun Microsystems in 1995
- Object oriented programming language
- OpenJDK implementation is open source
- Source code ⇒ byte code ⇒ just-in-time compiler
 - Byte code is portable & platform independent
 - Virtual machine abstracts from systems
- Strong and static type system
- Popular language for Enterprise & Big Data applications
 - Most popular programming language (Pos. 1 on the TIOBE index)
- Development tools: Eclipse

Specialties

- Good runtime and compile time error reporting
- Generic data types (vs. templates of C++)
- Introspection via. Reflection

Example Java Program

```
1 import java.util.Scanner;
2 import java.io.FileReader;
3 import java.io.FileNotFoundException;
4 // compile with javac program.java
5 // run with java program
6 public class program{
7     // the main method is part of a class
8     public static void main(String [ ] args) throws FileNotFoundException{
9         try{
10             // read from file "program.java" and create simple tokens
11             Scanner data = new Scanner(new FileReader("program.java"));
12             while(data.hasNext()){
13                 System.out.println(data.next());
14             }
15         }catch(Exception e){
16             // handle error here, we'll just rethrow the error
17             throw(e);
18         }
19     }
20 }
```

Example Java Classes

```
1 // Run: javac classes1.java and java Rabbit
2 // An abstract class is not completely implemented
3 abstract class Animal{
4     // instance member
5     private float weight;
6     // not-implemented instance function
7     public abstract String name();
8     // constructor
9     public Animal(float weight){ this.weight = weight; }
10    public String toString(){ return "I'm a " + name() + " with " +
11        weight + " kg"; }
12 }
13
14 class Rabbit extends Animal{
15     // invoke the constructor of the parent
16     public String name(){ return "Rabbit"; }
17     public Rabbit(){ super(2.5f); }
18
19     // the main method is part of a class
20     public static void main(String [ ] args){
21         Animal a = new Rabbit();
22         System.out.println(a); // I'm a Rabbit with 2.5 kg
23     }
24 }
25
26
```

Introduction to Python

- Open source
- Position 5 on TIOBE index
- Interpreted language
- Weak type system (errors at runtime)
- Development tools: any editor, interactive shell
- Note: Use and learn python3 explicitly
- Recommended plotting library: matplotlib⁹

Specialties

- Strong text processing
- Simple to use
- Support for object oriented programming
- Indentation is relevant for code blocks

⁹<http://matplotlib.org/gallery.html>

Example Python Program

```
1 #!/bin/env python
2 import re # use the module 're'
3
4 # function reading a file
5 def readFile(filename):
6     with open(filename, 'r') as f:
7         data = f.readlines()
8         f.close()
9     return data
10    return [] # return an empty array/list
11
12 # the main function
13 if __name__ == "__main__":
14     data = readFile('intro.py')
15     # iterate over the array
16     for x in data:
17         # extract imports from a python file using a regex
18         m = re.match("import[ \t]+(?P<WHAT>[^# ]*)", x)
19         if m:
20             print(m.group("WHAT"))
21             # dictionary (key value pair)
22             dic = m.groupdict()
23             dic.update( {"FILE" : 'intro.py'}) # append a new dict. with one key
24             # use format string with dictionary
25             print("Found import '%(WHAT)s' in file %(FILE)s" % dic )
26             # Prints: Found import 're' in file intro.py
```

Example Python Classes

```
1 from abc import abstractmethod
2
3 class Animal():
4     # constructor, self are instance methods, else class methods
5     def __init__(self, weight):
6         self.__weight = weight # private variables start with __
7
8     # decorator
9     @abstractmethod
10    def name(self):
11        return self.__class__.__name__ # reflection like
12
13    def __str__(self):
14        return "I'm a %s with weight %f" % (self.name(), self.__weight)
15
16 class Rabbit(Animal):
17     def __init__(self):
18         # super() is available with python 3
19         super().__init__(2.5)
20
21     def name(self):
22         return "Small Rabbit" # override name
23
24 if __name__ == "__main__":
25     r = Rabbit()
26     print(r) # print: I'm a Small Rabbit with weight 2.500000
```

Introduction to R

- Based on S language for statisticians
- Open source
- Position 19 on TIOBE index
- Interpreter with C modules (packages)
 - Easy installation of packages via CRAN¹⁰
- Popular language for data analytics
- Development tools: RStudio (or any editor), interactive shell
- Recommended plotting library: ggplot2¹¹

Specialties

- Vector/matrix operations. Note: Loops are slow, so avoid them
- Table data structure (data frames)

¹⁰Comprehensive R Archive Network

¹¹<http://docs.ggplot2.org/current/>

Course for Learning R Programming

```
1 # Run with "Rscript intro.R" or run "R" and copy&paste into interactive shell
2 # Installing a new package is as easy as:
3 install.packages("swirl")
4 # Note: sometimes packages are not available on all mirrors!
5 library(swirl) # load the package
6
7 help(swirl) # read help about the function swirl
8
9 swirl() # start an interactive course to learn R
10
11 # a simple for loop
12 for (x in 1:10){
13   if (x < 5){
14     print(x)
15   }else{
16     print(x * 2)
17   }
18 }
```

Example R Program

```
1 # create an array
2 x = c(1, 2, 10:12)
3
4 # apply an operator on the full vector and output it
5 print( x*2 ) # prints: 2 4 20 22 24
6
7 # slice arrays
8 print ( x[3:5] ) # prints: 10 11 12
9 print( x[c(1,4,8)] ) # prints: 1 11 NA
10
11 r = runif(100, min=0, max=100) # create array with random numbers
12 m = matrix(r, ncol=4, byrow = TRUE) # create a matrix
13
14 # slice matrix rows "m[row(s), column(s)]"
15 print( m[10:12, ] ) # Output:
16 #          [,1]      [,2]      [,3]      [,4]
17 #[1,] 85.46609 60.749703 10.5062183 7.449173
18 #[2,] 79.76042 52.199321 96.9699856 97.877946
19 #[3,] 37.34286 8.266282  0.3398741  1.957607
20
21 # slice rows & columns
22 print ( m[10, c(1,4)] ) # Output: [1] 85.466085 7.449173
23
24 # subset the table based on a mask
25 set = m[ (m[,1] < 20 & m[,2] > 2) , ]
```

Accessing CSV Files with R

```
1 # function to create a table (data frame) and fill it with random data
2 createTable = function (size){
3   tbl = read.table(text="", col.names = c("Type", "Time"))
4   tbl[1:size, ] = 0 # initialize size times a full rows
5   tbl$Time = runif(size, min=0, max=100) # address by column name
6   # create random types, factor() for nominal data and
7   # ordered() for ordinal data
8   tbl$type = factor(round(runif(size, min=0.5, max=3.49)),
9     levels=1:3, # three categories
10    labels=c("unknown", "good", "bad"))
11   tbl$type[size] = "bad" # assign last element to be bad
12   return (tbl)
13 }
14 # change columnnames
15 colnames(tbl) = c("Typ", "Duration")
16
17 d = createTable(5)
18 # Assign the column with the name
19 print( d )
20 print( summary (d) ) # some statistics about d
21 # Write CSV incl. header
22 write.table(d, file = "mydata.csv", sep=",", row.names=FALSE)
23 # reread table
24 d = read.table("mydata.csv", header = TRUE, sep = ",")
```

Summary

- Big data analytics
 - Explore data and model causalities to gain knowledge & value
- Challenges: 5 Vs – Volume, velocity, variety, veracity, value
- Data sources: Enterprise, humans, Exp./Observational data (EOD)
- Types of data: Structured, unstructured and semi-structured
- Levels of analytics: Descriptive, predictive and prescriptive
- Roles in big data business: Data scientist and engineer
- Data science != business “intelligence”

Bibliography

- 1 Book: Lillian Pierson. **Data Science for Dummies**. John Wiley & Sons
- 2 Report: Jürgen Urbanski et.al. **Big Data im Praxiseinsatz – Szenarien, Beispiele, Effekte**. BITKOM
- 3 <http://winfwiki.wi-fom.de/>
- 4 Forrester Big Data Webinar. Holger Kisker, Martha Bennet. Big Data: Gold Rush Or Illusion?
- 5 <http://blog.eoda.de/2013/10/10/veracity-sinnhaftigkeit-und-vertrauenswuerdigkeit-von-bigdata-als-kernherausforderung-im-informationszeitalter/>
- 6 <http://lehrerfortbildung-bw.de/kompetenzen/projektkompetenz/methoden/erkenntnis.htm>
- 7 Gilbert Miller, Peter Mork From Data to Decisions: A Value Chain for Big Data.
http://www.fh-schmalkalden.de/Englmeier-p-790/_ValueChainBigData.pdf
- 8 Andrew Stein. The Analytics Value Chain. <http://steinvox.com/blog/big-data-and-analytics-the-analytics-value-chain/>
- 9 Dursun Delen, Haluk Demirkiran,. Decision Support Systems, Data, information and analytics as services.<http://j.mp/11bl9b9>
- 10 Wikipedia
- 11 Kashmir Hill. 46 Things We've Learned From Facebook Studies. Forbe.
<http://www.forbes.com/sites/kashmirhill/2013/06/21/46-things-weve-learned-from-facebook-studies/>
- 12 Hortonworks <http://hortonworks.com/>
- 13 http://www.huffingtonpost.com/2014/12/10/facebook-most-popular-paper_n_6302034.html
- 20 <http://hortonworks.com/blog/enterprise-hadoop-journey-data-lake/>
- 21 http://www.stackit.com/hadoop/?utm_campaign=Stackit+Hadoop+Infographic