



# Datenanalyse & Big Data

Harald Schilly, Wifi Wien, 20160616



# Geschichte der Technologien

# Mathematik/Statistik

- Wahrscheinlichkeitstheorie: neue Ansätze und Methoden, philosophische Überlegungen, ...
- Wissenschafts&Erkenntnistheorie (Falsifikationismus)
- Bayes Statistik (Rechenaufwändiger, subjektivistischer Wahrscheinlichkeitsbegriff, Vorwissen, Verteilungen, Laplace 1812)
- Bootstrapping (1979)
- Optimierung: LP, MILP, konvex, NLP
- Graphentheorie

# ICT

- Verteiltes Rechnen ([Amdahl's Law](#))
- Funktionales Programmieren / immutable Datenstrukturen
- Approximative Algorithmen
- Relationale Datenbanken und NoSQL
  - guaranteed vs. eventual consistency
- SQL und DataFrames (R)
- [Moore's Law](#) / [Kryder's Law](#)
- NUMA Architektur
- Container & Cluster Management ([Kubernetes](#), [Mesos](#), ...)

# Visualisierung



Kommunikation der Ergebnisse, Überblick, sinnliche Erfassung

- [John Tukey](#)
- [Edward Tufte](#)
- [Leland Wilkinson](#)'s [Grammar of Graphics](#)



# Multidisciplinary: Data Science

Extraktion von Information aus Daten (information retrieval)

## Wissenschaften

- Mathematik
- Statistik/Wahrscheinlichkeitstheorie
- Informationstheorie
- Signalverarbeitung
- Mustererkennung
- Optimierung
- Visualisierung/Design

## ICT

- Programmierung
- Datenbanken
- Machine Learning
- Storage
- Cluster (HPC/HTC)
- Streaming



### Computer Science

• Leibniz – Binary Logic.

- Turing machines
- Information Theory
- Weiner & Cybernetics
- Von Neumann Architecture.

- Babbage, Lovelace
- Boolean Algebra
- Punch cards.

- First IBM Computers
- DBMS.

- Sort & Search Algorithms – Dijkstra, Kruskal, Shell Sort, ...
- Heuristics – Simulated Annealing, ...

- Graph Algorithms
- Multigrid methods
- Tree based methods.

- Text/ string search
- 1974 Peter Naur "Concise Survey of Computer Methods". **Data Science, Datalogy**
- Knuth – Art of Computer Programming.

- 1989 First KDD Workshop
- Gregory Piatsky-Shapiro.

- Database Marketing
- Data Mining, Knowledge Discovery
- "Data science, classification, and related methods."

### Data Technology

- Cartography
- Astronomical Charts.

- William Playfair
- Charles Minard
- Florence Nightingale.

- Removable Disk drives
- Relational DBMS.

- Desktop, floppy
- SQL, OOP
- High level languages.

- William Cleveland: Data Science
- Leo Breimann: Statistical Modeling: 2 Cultures.

### Visualization

- Calculus
- Logarithms
- Newton-Raphson.

- Optimization Methods
- Fourier and other transforms
- Matrix & Generalizations
- Non-euclidean geometries.

- Applications to Military, manufacturing, Communications.

- Networks
- Assignment Problems
- Automation
- Scheduling

- Edward Tufte.

- Grammar of Graphics
- Word Cloud, Tag Cloud.

### Mathematics/ OR

- Probability
- Correlation
- Bayes Theorem.

- Regression, Least Squares
- Time Series.

- Theoretical Foundations of Modern Statistics
- Hypothesis, DOE
- Mathematical Statistics.

- Bayesian Methods
- Time Series Methods (Box Cox, Survival, etc.)
- Stochastic Methods.

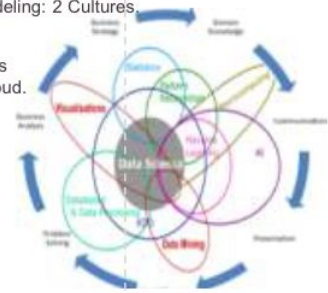
- 1962 John W. Tukey, Future of Data Analysis

- 1976 – SAS Institute
- 1977 The International Association for Statistical Computing (IASC).

- Decision Science
- Pattern recognition
- Machine learning.

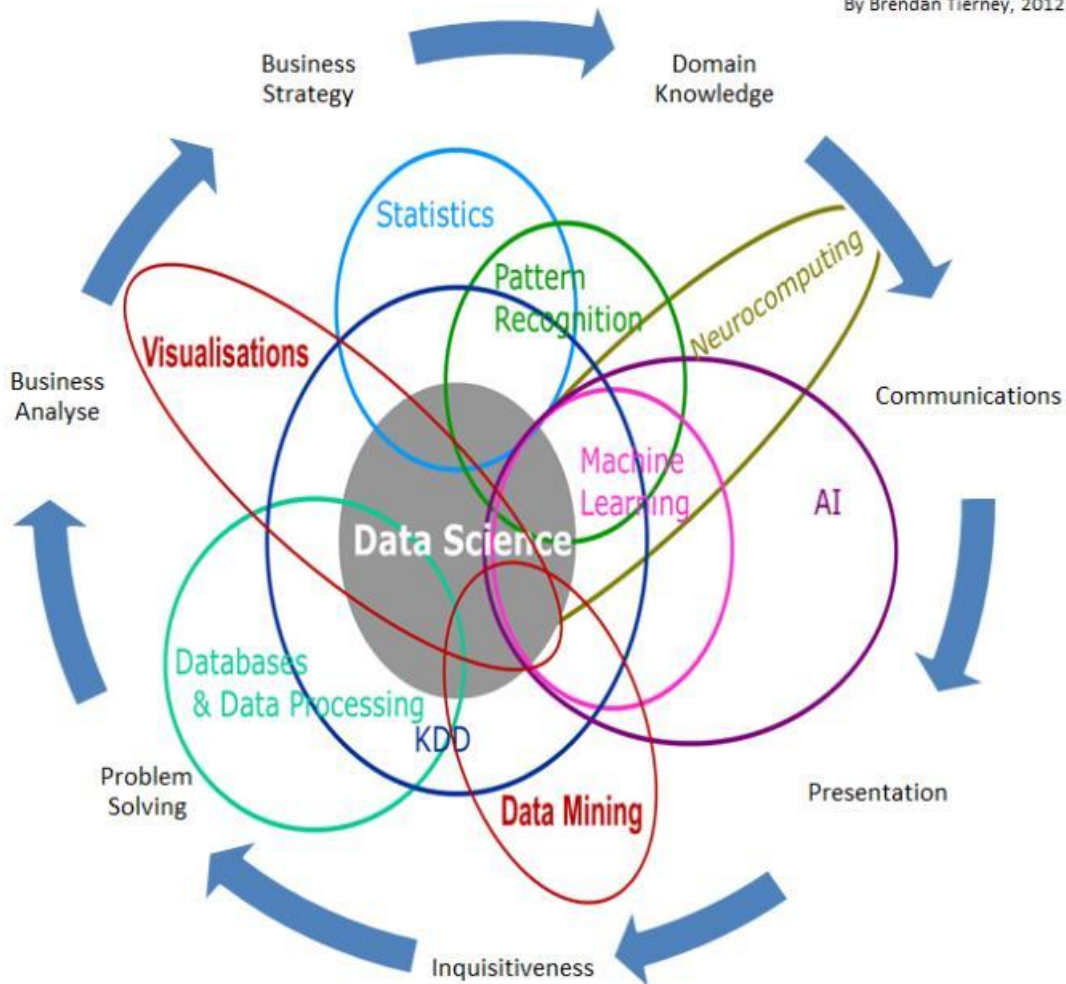
- Simulation, Markov
- Computational Statistics.

### Statistics



# Data Science Is Multidisciplinary

By Brendan Tierney, 2012







# Industrielle Forschung

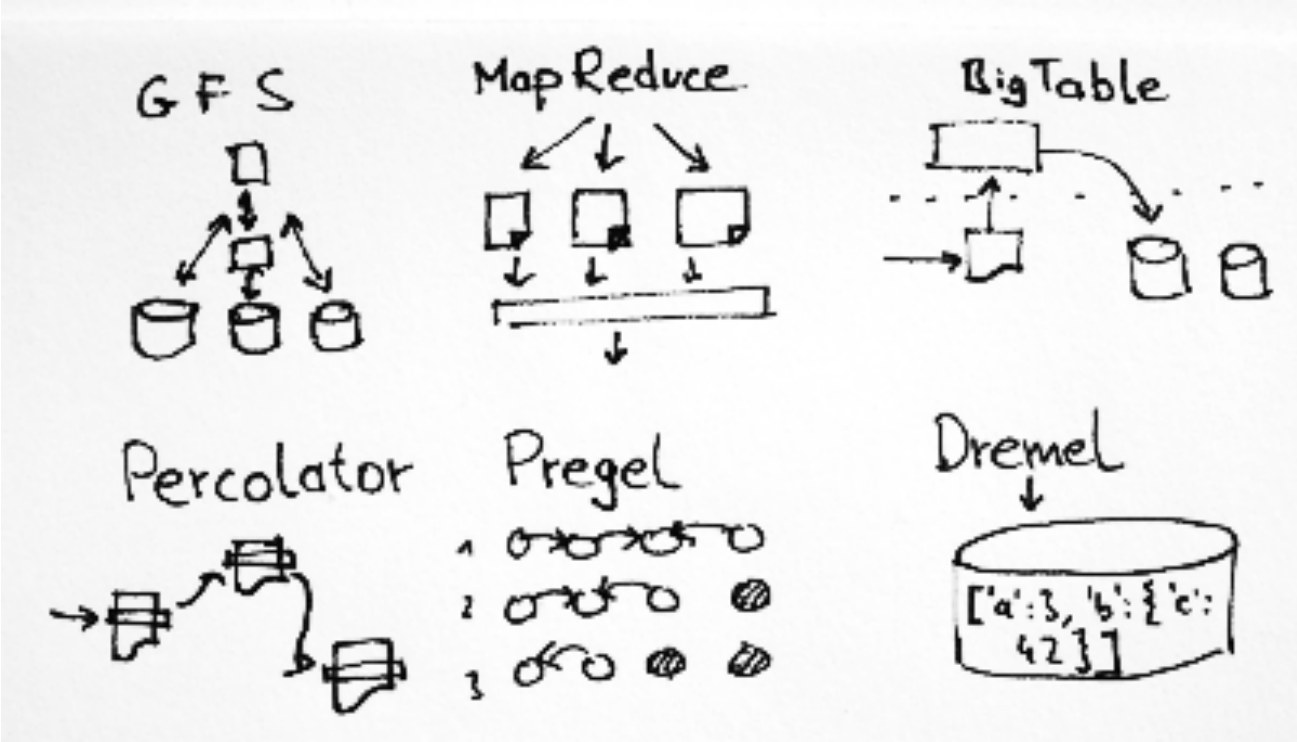
(zuvor militärische)



# Industry Research / Google & Co.

- Google Filesystem 2003: [dx.doi.org/10.1145%2F945445.945450](https://dx.doi.org/10.1145%2F945445.945450)
- MapReduce 2004: [research.google.com/archive/mapreduce.html](https://research.google.com/archive/mapreduce.html)
  - $\text{Map}(k_1, v_1) \rightarrow \text{list}(k_2, v_2)$  und  $\text{Reduce}(k_2, \text{list}(v_2)) \rightarrow \text{list}(v_3)$
- BigTable 2006: [research.google.com/archive/bigtable.html](https://research.google.com/archive/bigtable.html) (u.a. [columnar storage](#))
- Pregel 2010: [dl.acm.org/citation.cfm?id=1807184](https://dl.acm.org/citation.cfm?id=1807184)
- Dremel 2010: [Interactive Analysis of Web-Scale Datasets](#) (BigQuery)
- Stream Processing Frameworks:  
[Perculator](#) (Google Index), S4 (Yahoo!), [Storm](#) (Twitter)
- TensorFlow 2015: [tensorflow.org](https://tensorflow.org) ([whitepaper](#))

# Industry Research / Google & co





# Software

unvollständiger Überblick

# Software Lösungen / Distributed Computing



- MPI - Klassisch, Rechenintensiv
  - [mpich](#)
- [Cassandra](#), [HBase](#), [Hive](#), ...
  - inspiriert von BigTable
- Hadoop: [hadoop.apache.org](#)
  - HDFS (GFS, ...)
  - (mit Oozie) von Yahoo!
- **Spark**: [spark.apache.org](#)
  - Verallgemeinertes Modell
  - umfasst ETL, Computation, GraphX, Machine-Learning und Stream Processing.
- Giraph: [giraph.apache.org](#)
  - Pregel, bei Facebook für Graph der Freundschaften

# Software Lösungen: Spark (2011)



- **Idee:** Hadoop ist langsam wegen Festplatte → nützen wir den Speicher!
- Handling von Datenverarbeitungsschritten und Fehlern? → RDDs!
- High-Level API zur Beschreibung der Schritte, Caching, und Realisierung.
- Framework zum verteilten Ausführen, mit Management, etc. (aufgesetzt auf Zookeeper, Mesos)
- Gewinn Daytona GraySort contest  
(3x schneller, 10x weniger Nodes als Hadoop)
- Machine Learning, Graph Analyse, Streaming, etc.



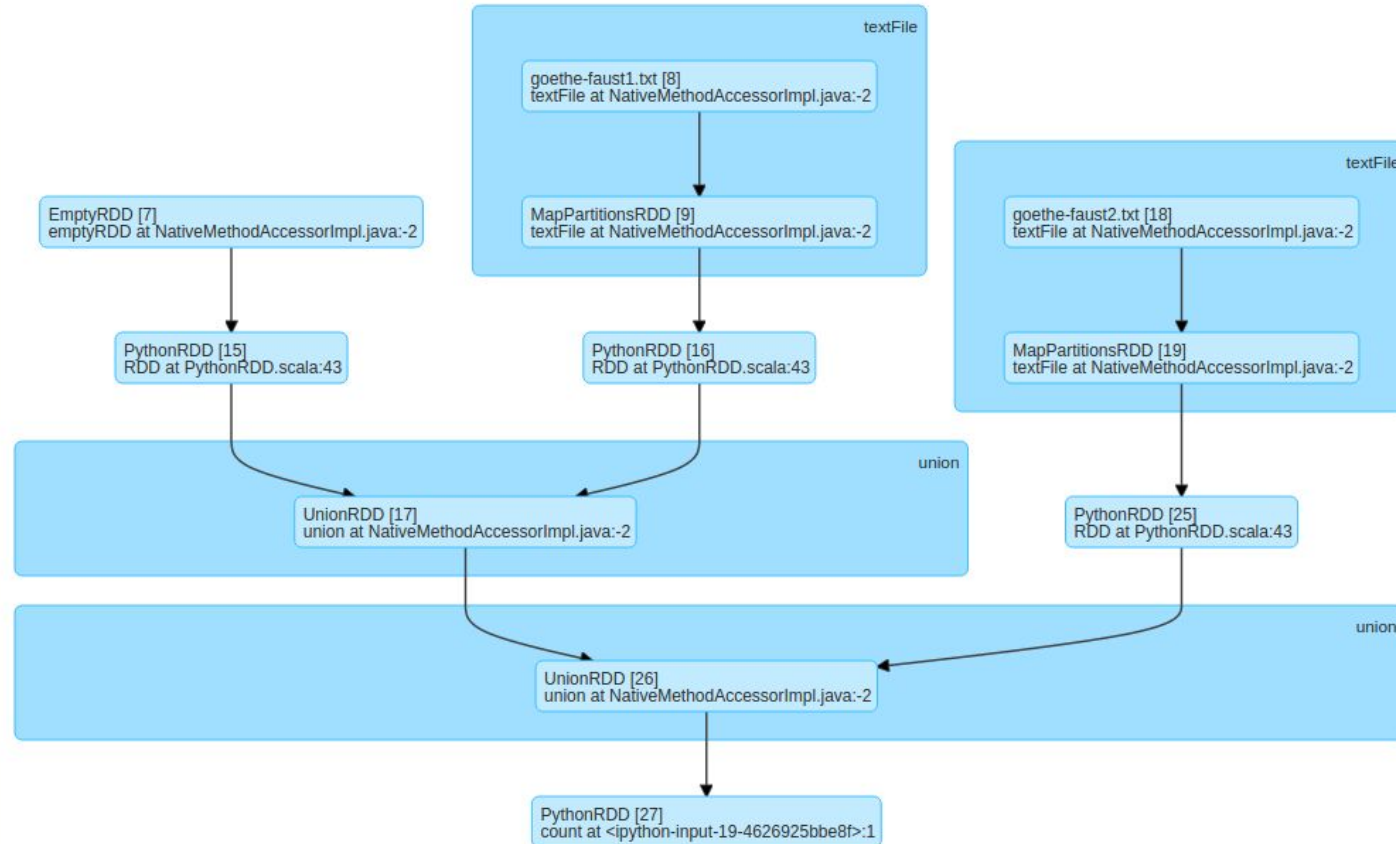
# Software Lösungen:



- Start 2011
- RDDs Initialisierung: `range`, `parallelize`, oder `textFile` (auch `hdfs://`)
- Zwei funktionale Grundbausteine:
  - **Transformationen:** RDD → RDD  
`map`, `flatMap`, `filter`, `sample`, `union`, `intersection`, `distinct`, `groupByKey`/  
`reduceByKey`/ `aggregateByKey`/ `sortByKey`, `join`, `cartesian`, `pipe`, ...
  - **Aktionen:** RDD → Value  
`collect`, `count`, `reduce`, `first` / `take` / `takeOrdered`, `foreach`, ...
- **Extras:** `cache`, `saveAs...`, `sc.accumulator`, `sc.broadcast`, ...

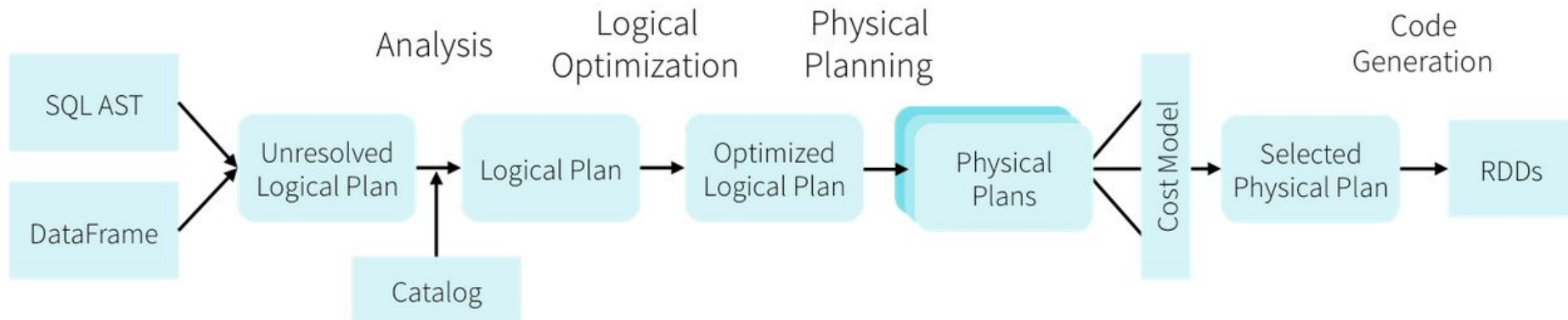
# Software Lösungen: Spark - Lineage Graph

Stage 12





# Software Lösungen: Spark - Query Execution



<https://databricks.com/blog/>

# Details for Stage 30 (Attempt 0)



Total Time Across All Tasks: 6 ms

Locality Level Summary: Process local: 2

Input Size / Records: 2.7 KB / 40

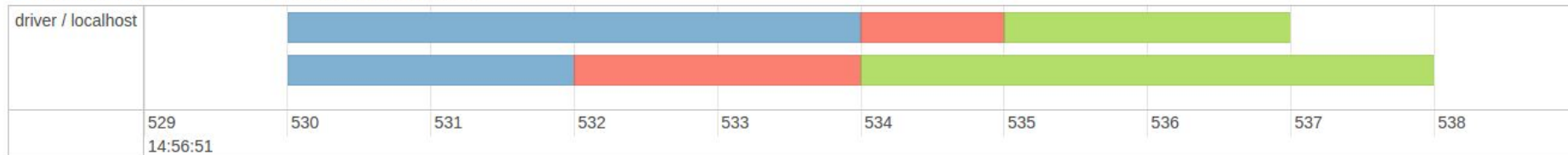
[▶ DAG Visualization](#)

[▶ Show Additional Metrics](#)

[▼ Event Timeline](#)

Enable zooming

- Scheduler Delay
- Task Deserialization Time
- Shuffle Read Time
- Executor Computing Time
- Shuffle Write Time
- Result Serialization Time
- Getting Result Time



## Summary Metrics for 2 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	2 ms	2 ms	4 ms	4 ms	4 ms
GC Time	0 ms	0 ms	0 ms	0 ms	0 ms
Input Size / Records	1376.0 B / 20	1376.0 B / 20	1376.0 B / 20	1376.0 B / 20	1376.0 B / 20

## Aggregated Metrics by Executor

Executor ID ▲	Address	Task Time	Total Tasks	Failed Tasks	Succeeded Tasks	Input Size / Records
driver	localhost:55233	15 ms	2	0	2	2.7 KB / 40

## Tasks

Index ▲	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time	Duration	GC Time	Input Size / Records	Errors
0	51	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2016/01/13 14:56:51	2 ms		1376.0 B (memory) / 20	

# Machine Learning in



Alle Funktionen in Scala; Schnittstellen für Python hat großen Umfang.

Overview: DataFrame→Transformation (Feature Extraction, ...) → Estimation → Model&Parameter

**MLLib**: <http://spark.apache.org/docs/latest/mllib-guide.html>

- Klassisch: Statistiken, Sampling, Optimierung (L-BFGS, ...), ...
- Clustering: KMeans, gauss mixture, LDA, ...
- Regression/Classification: Linear (auch SVMs, Logistisch, ...), NaiveBayes, DecisionTree, Gradient-Boosted Trees, RandomForest, Isotonic,
- Recommendation: ALS
- Accelerated failure time model: AFT
- Python API: <http://spark.apache.org/docs/latest/api/python/index.html>
- oft auch support für “streaming data”...

# Stream Processing in



- Quellen: Kafka, Flume, Twitter, ZeroMQ, Kinesis, or TCP sockets
- Datenpakete in Zeitfenstern als “DStream” (=RDD für Streaming)
- Windowing (sliding interval), Verarbeitung in Transformationen und ML-Modellen, ...
- Ausgabe
- <http://spark.apache.org/docs/latest/streaming-programming-guide.html>



# Software Lösungen: Drill



Drill: [drill.apache.org](http://drill.apache.org) - schema-free SQL engine

- “Treat your data like a table even when it's not”:  
mischen von Daten über Grenzen der Datenbank/Files hinweg, etc.
- Ideen von Dremel
- Durch SQL kompatibel zu BI Toolchain ([Tableau](#), etc.)

# Software Lösungen / weitere

- Parquet: [parquet.apache.org](https://parquet.apache.org) - columnar storage
- Hive: [hive.apache.org](https://hive.apache.org) - DataWarehouse, query large datasets
- HBase: [hbase.apache.org](https://hbase.apache.org) - Facebook, Yahoo! , ...
- Accumulo: [accumulo.apache.org](https://accumulo.apache.org) - BigTable + sorting
- Presto: [prestodb.io](https://prestodb.io) - Facebook (300PB warehouse), Airbnb, Dropbox
- Impala: [impala.io](https://impala.io) - SQL query engine für HDFS, HBase
- **BlinkDB**: [blinkdb.org](https://blinkdb.org) - schneller als in-Memory? → Ja!  
approximative Abfragen (stratified/poissonized sampling)



# Chancen

# Wissenschaftliche Methodik

- Hypothese, Empirie, ...
- Datenbasierte Entscheidungen
- Nachvollziehbarkeit
- Automatisierung
- Szenarienanalyse
- bessere Validierung
- Statistik
- Unsupervised Learning
- Predictive Modelling





# Datenverarbeitung

- ETL Pipeline
- weniger Mehrgleisigkeit
- schnellere Resultate
- Flexibilität
- weniger vendor lock-in
- ad-hoc
- Skalierbarkeit
- Mashup



# Visualisierung

- bessere Kommunikation von Effekten/Resultaten
- interaktive Webseiten statt statischer Dokumente  
(Bokeh, Shiny, Spyre, ...)



2012

HBR.ORG

# Harvard Business Review

OCTOBER 2012

46 **The Big Idea**  
The True Measures  
Of Success  
Michael J. Mauboussin

54 **International Business**  
10 Rules for Managing  
Global Innovation  
Kerley Wilson and Voss L. Ditz

55 **Leadership**  
What Ever Happened  
To Accountability?  
Thomas C. Ricks

GETTING  
CONTROL  
OF

# BIG DATA



How vast new streams of  
information are changing  
the art of management  
**PAGE 59**





May 2011

# Big data: The next frontier for innovation, competition, and productivity



*FIN*

© 2016, Harald Schilly <harald@schil.ly>

Lizenz: CC BY-SA 4.0