# Summary

## Big Data – Web Scraping and Machine Learning

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

# The Excitement

US Aggregated Inflation Series, Monthly Rate, PriceStats Index vs.
Official CPI. Accessed January 18, 2015 from the PriceStats website.

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

Social media sentiment (daily, weekly and monthly) in the Netherlands, June 2010 - November 2013. The development of consumer confidence for the same period is shown in the insert (Daas and Puts 2014).

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

Source: Roberto Rigobon

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

# Excitement over new data sources:

1. New research questions can be asked
   - spatial and temporal granularity
   - small parts of the populations
   - other form of data (text, visuals)

2. Reduced data collection costs

3. 'Instant' more timely availability

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

# Vs as defining characteristic



http://www.rosebt.com/blog/data-veracity

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

# ...one more V



| Volume | Velocity | Variety | Veracity* |
|---|---|---|---|
| **Data at Rest** | **Data in Motion** | **Data in Many Forms** | **Data in Doubt** |
| Terabytes to exabytes of existing data to process | Streaming data, milliseconds to seconds to respond | Structured, unstructured, text, multimedia | Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations |

http://www.rosebt.com/blog/data-veracity

# The Process

# Big Data Process Map

**Generate**

**ETL**

**Analyze**

Source 1

Source 2

Source M

Extract

Transform
(Cleanse)

Load (Store)

Filter/Reduction
(Sampling)

Computation/
Analysis
(Visualization)

# Big Data Process Map

**Generation**



Source 1

Source 2

Source M

Similar to data collection errors in surveys; data may be erroneous or missing; data generating units may be self-selected; meta-data may be lacking or absent

Transform (Cleanse)

Load (Store)

**Analyze**

Filter/Reduction (Sampling)

Computation/ Analysis (Visualization)

# Big Data Process Map

**Generation**

**ETL**

Source 1

Source 2

•
•
•

Source M

Extract

Transform (Cleanse)

Load (Store)

Similar to data processing stages in surveys; includes creating or enhancing meta-data; record matching; variable coding, editing, data munging or scrubbing, and data integration

Computation/ Analysis (Visualization)

# Big Data Process Map

**Generation**

**ETL**

**Analyze**

Source 1

Source 2

Source M

Extract

Similar to estimation and analysis error in surveys; includes weighting, modeling, estimation, graphing…

Load (Store)

Filter/Reduction (Sampling)

Computation/ Analysis (Visualization)
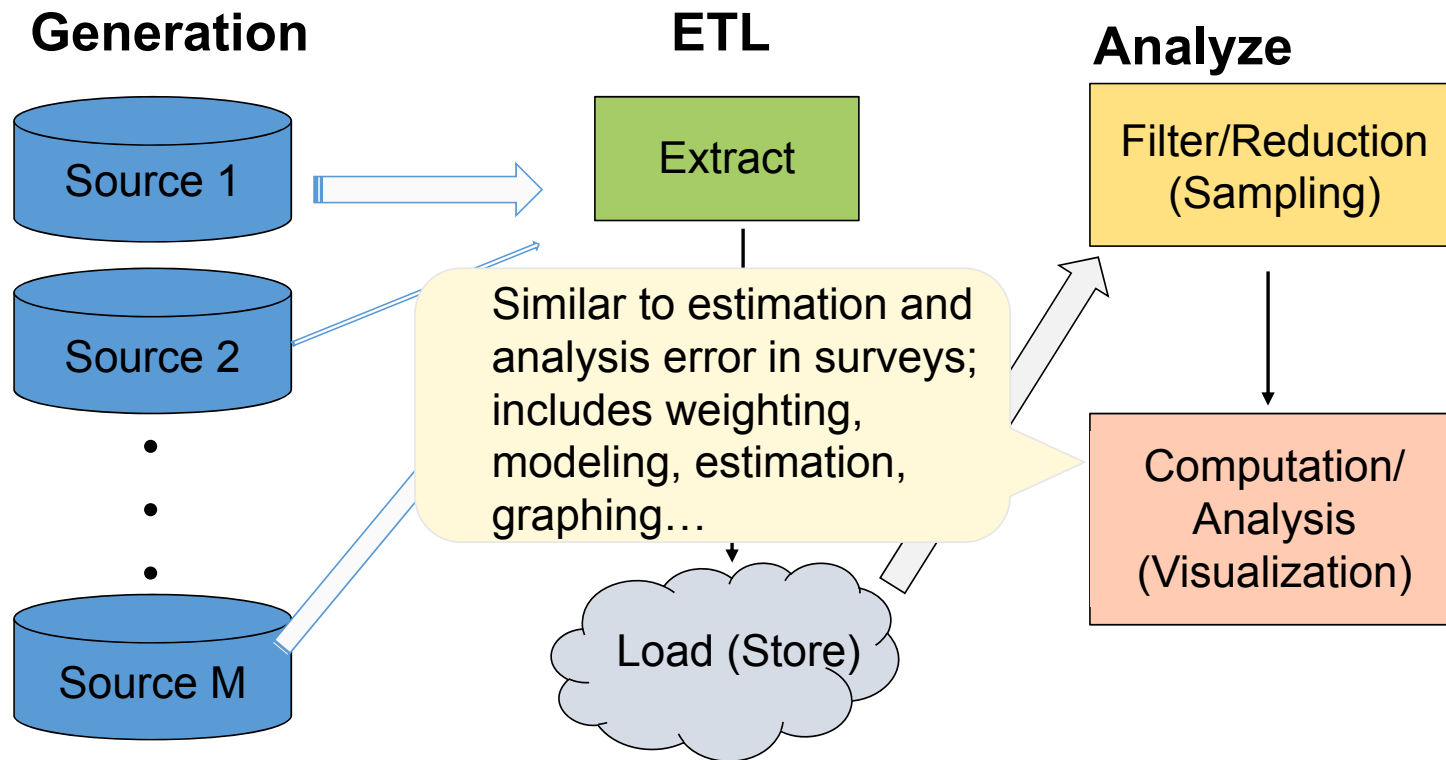
# Key Ingredients for Valid Inference

1. **Data generating process needs to be known**

2. **Framework as tool to identify errors**

3. **Model or break confounders**

4. **Know your inferential goal**

# The Skills

Content key words

| | |
|---|---|
| **Data Output/Access** | Visualization, disclosure control, ethics, privacy |
| **Data Analysis** | Statistical methods, machine learning, network Bayesian, hierarchical, small area, spatial |
| **Data Curation/Storage** | Data munging, database management, SQL, editing, coding, imputation, etc. |
| **Data Generating Process** | Web-, Mobile-, Phone-, F2F-Surveys, APIs, Web scraping, linkage, matching, sampling, weighting |
| **Research Questions** | Economics, public policy, criminology, journalism, public health, sociology, etc. |

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

Classical Statistical Approaches versus Statistical Machine Learning

Model Evaluation/Validation

**Database Management**

Programming with Big Data

# When to use different data management and analysis technologies

**Text files and scripting language**
- Your data is small
- Your analysis is simple
- You do not expect to repeat analyses over time

**Statistical packages**
- Your data is modest in size
- Your analysis maps well to your chosen statistical package

**Relational database**
- Your data is structured
- You will be analyzing data repeatedly over time

**NoSQL database**
- Your data is unstructured
- Your data is extremely large

**Classical Statistical Approaches versus Statistical Machine Learning**

**Model Evaluation/Validation**
**Database Management**
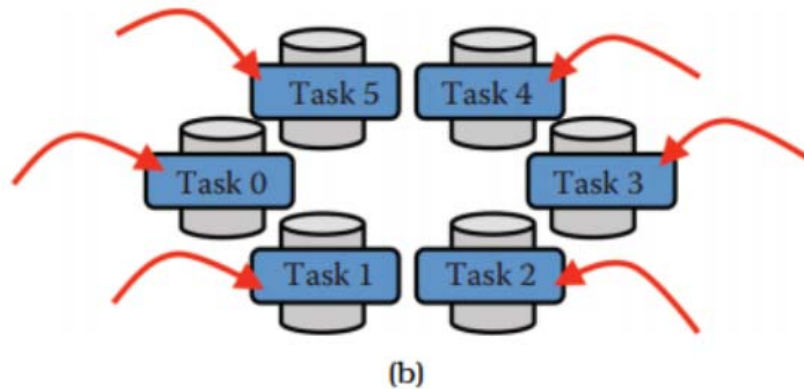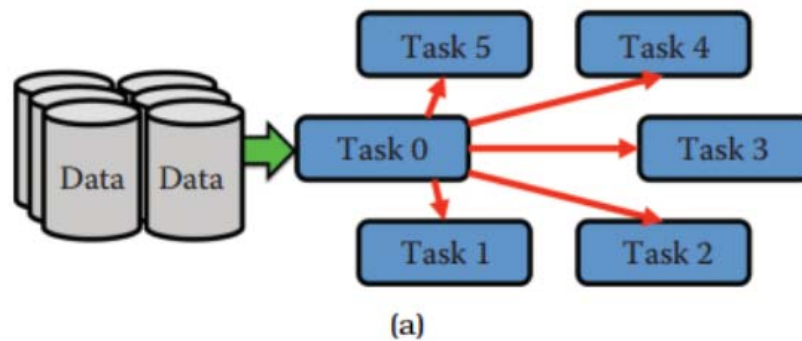
**Programming with Big Data**

**Figure 5.1.** (a) The traditional parallel computing model where data is brought to the computing nodes. (b) Hadoop's parallel computing model: bringing compute to the data [241]
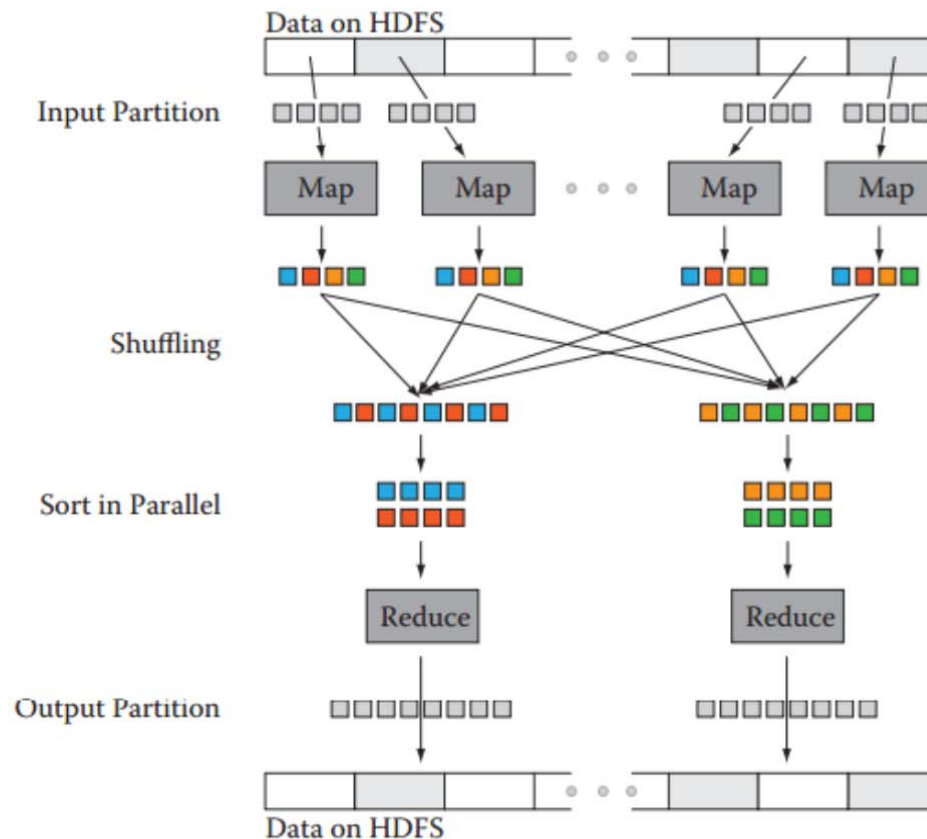
**Figure 5.2.** Data transfer and communication of a MapReduce job in Hadoop. Data blocks are assigned to several maps, which emit key–value pairs that are shuffled and sorted in parallel. The reduce step emits one or more pairs, with results stored on the HDFS

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

**RESEARCHER**

Team member with experience applying formal research methods, including survey methodology and statistics

**DOMAIN EXPERT**

User, analyst, or leaders with deep subject matter expertise related to the data, its appropriate use, and its limitations

**SYS ADMIN**

Team member responsible for defining and maintaining a computation infrastructure that enalbes large scale computation

**COMPUTER SCIENTIST**

Technically skilled team member with education in computer programming and data processing technology

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg          Usher 2015

# Learn More & Engage

http://coleridgeinitiative.org/

http://survey-data-science.net/

http://datafest.de

SPONSORED BY THE

Federal Ministry of Education and Research

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg

# Data Mining/Machine Learning Resources

- http://www.dataminingconsultant.com/resources.htm
- Data Mining Algorithms Explained Using R (2015)
  - http://bit.ly/1yZYHjK
- Data Mining for the Social Sciences (2015)
  - http://bit.ly/1DpPFC2
- An Introduction to Statistical Learning with Applications in R (2013)
  - Free PDF Version: http://bit.ly/1iUJso0
  - Online Resources for FREE lecture videos and labs in R
    - http://bit.ly/1snBMk5
- An overview of Machine Learning Functions available in R
  - http://cran.r-project.org/web/views/MachineLearning.html

Big Data Bundesbank - Kreuter, Kern, Schierholz, Sternberg