# International Program in Survey and Data Science

Frauke Kreuter

JPSM – Uni Mannheim – IAB                    ASI 28.06.2017

The National Academies of
SCIENCES · ENGINEERING · MEDICINE

**REPORT**

# INNOVATIONS IN FEDERAL STATISTICS

## Combining Data Sources While Protecting Privacy

Chapman & Hall/CRC
Statistics in the Social and Behavioral Sciences Series

# BIG DATA AND SOCIAL SCIENCE

## A Practical Guide to Methods and Tools

Edited by
**Ian Foster, Rayid Ghani,
Ron S. Jarmin, Frauke Kreuter,
and Julia Lane**

**CRC Press**
Taylor & Francis Group
A CHAPMAN & HALL BOOK

Source: Roberto Rigobon
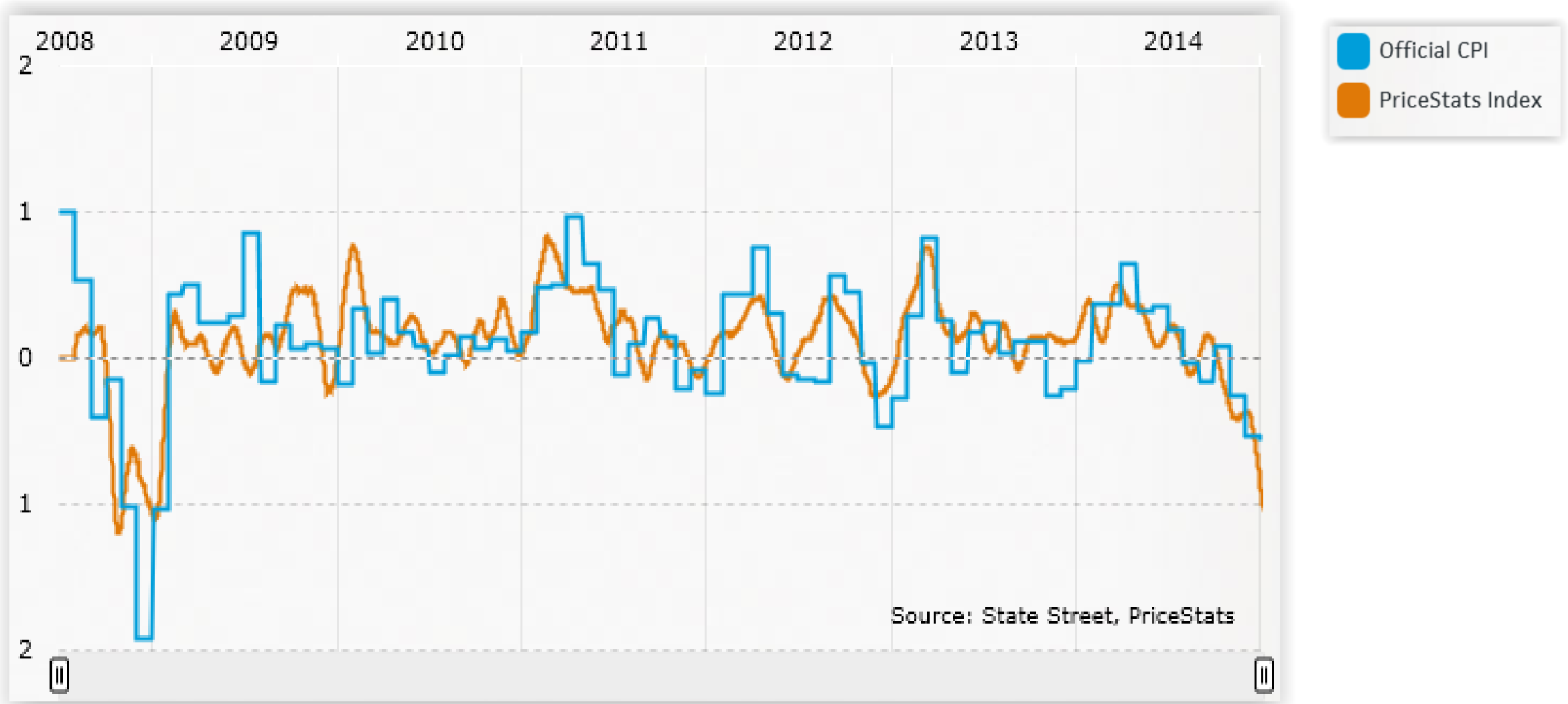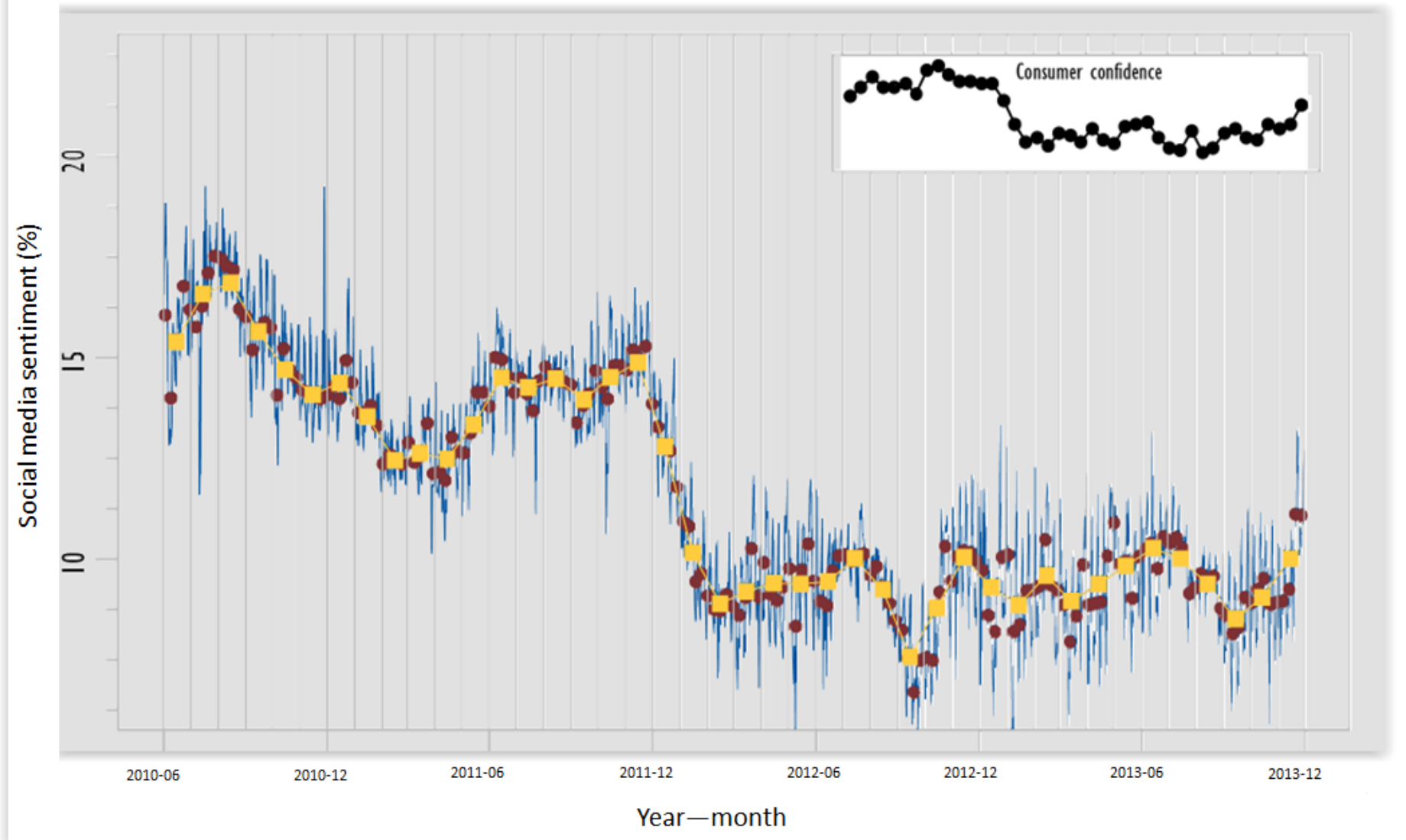
# The Excitement

4

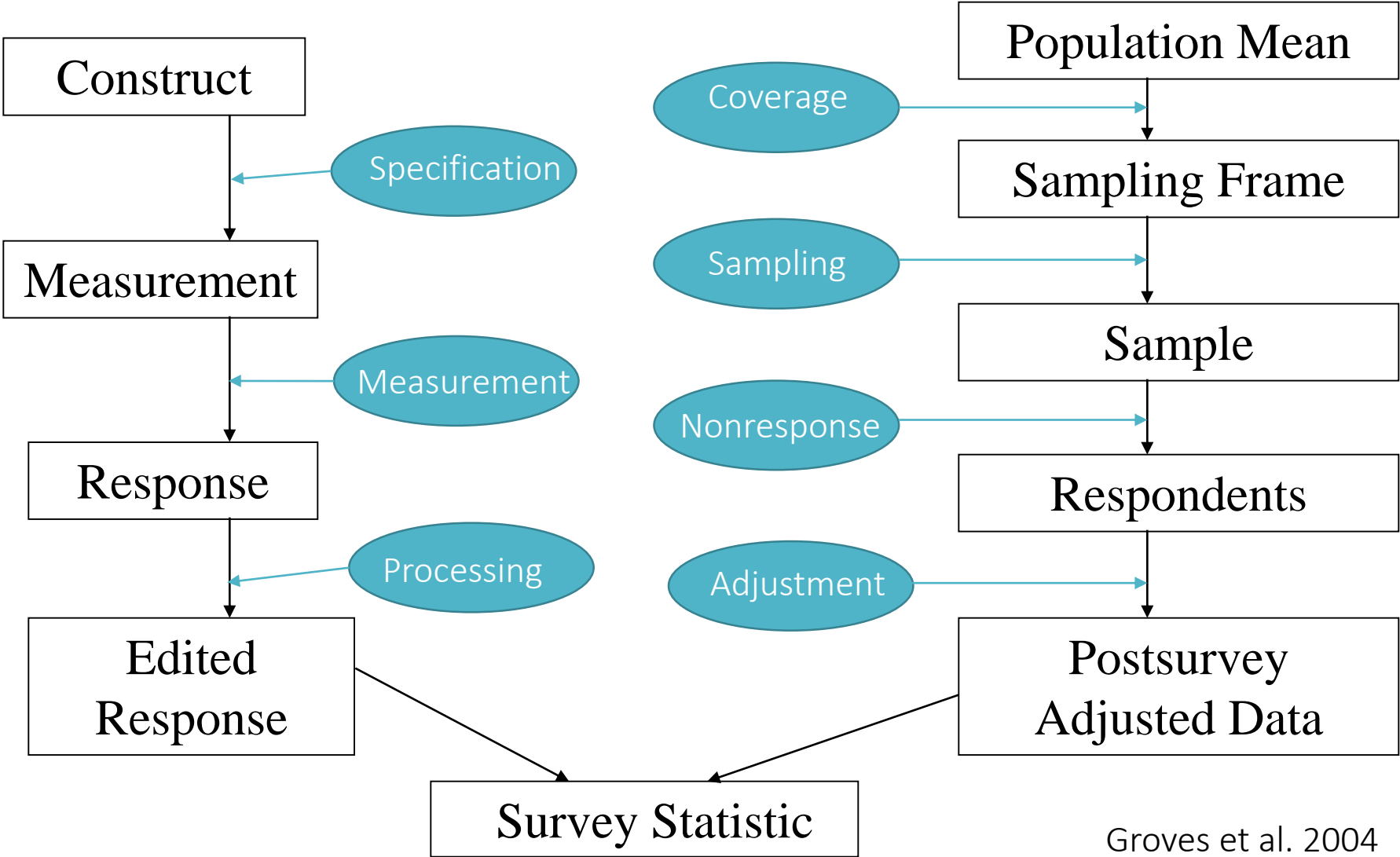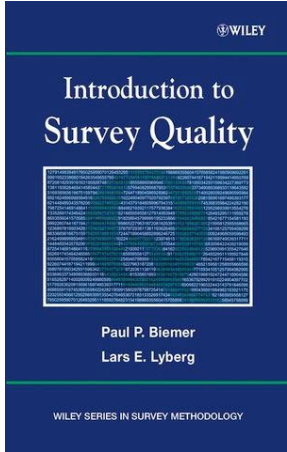US Aggregated Inflation Series, Monthly Rate,
PriceStats Index vs. Official CPI. Accessed January 18, 2015 from the PriceStats website.

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).

# The Doubt

# Data Generating Process



Groves et al. 2004

# Big Data Process Map

**Generate**

**ETL**

**Analyze**

Source 1

Source 2

•
•
•
•

Source M

Extract

Transform (Cleanse)

Load (Store)

Filter/Reduction (Sampling)

Computation/ Analysis (Visualization)

BIEMER, P. (2017) ERRORS AND INFERENCE, CHAPTER 10 IN BD AND SOCIAL SCIENCE

# Big Data Process Map

**Generation**

**Analyze**

Source 1

Source 2

•
•
•
•

Source M

Ähnlich wie in Umfragen:

Fehlende Werte;

Selbstselektion;

fehlende Meta-Daten

Transform (Cleanse)

Load (Store)

Filter/Reduction (Sampling)

Computation/ Analysis (Visualization)

BIEMER, P. (2017) ERRORS AND INFERENCE, CHAPTER 10 IN BD AND SOCIAL SCIENCE

Big Data Process Map

**Generation**

**ETL**

Source 1

Source 2

.
.
.

Source M

Extract

Transform (Cleanse)

Load (Store)

Ähnlich wie Datenaufbereitung in Erhebungen: Kodierung, Editierung, Säubern, Imputation, Integration, Zusammenführung von Datensätzen

Computation/ Analysis (Visualization)

# Big Data Process Map

**Generate**

**ETL**

**Analyze**

Source 1

Source 2

·
·
·

Source M

Extract

Transform (Cleanse)

Load (Store)

Filter/Reduction (Sampling)

Computation/ Analysis (Visualization)

# The Skills

**Content key words**

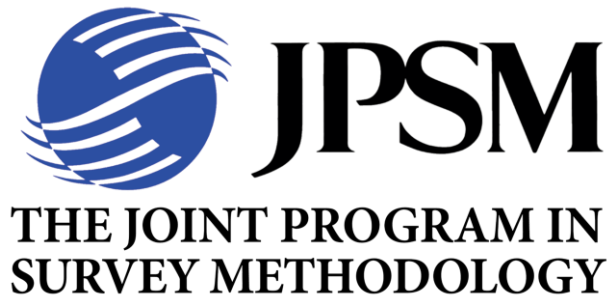| Stage | Keywords |
|---|---|
| Data Output/Access | Visualization, disclosure control, ethics, privacy |
| Data Analysis | Statistical methods, machine learning, Bayesian, hierarchical, small area estimation |
| Data Curation/Storage | Practical training in data base management, SQL, editing, coding, imputation, etc. |
| Data Generating Process | Designed (survey and admin) and organic data (transaction and aspirational), linkage, matching |
| Research Questions | Economics, public policy, criminology, journalism, public health, sociology, etc. |

# The Program

16

# Project coordinators and funding

# New program characteristics – In brief

- Multidisciplinary and modularized curriculum
- Relevant methods and tools
- Faculty from world-leading institutions

- Flexible web-based learning environment
- Live (video) interaction with faculty and students
- Face-to-face networking meetings

# INTERNATIONAL PROGRAM IN SURVEY AND DATA SCIENCE

offered through the University of Mannheim and the Joint Program in Survey Methodology (Universities of Maryland and Michigan, Westat)

**BE PART OF IT**

survey-data-science.net

We are pleased to announce the launch of the International Program in Survey and Data Science (IPSDS). Fundamental changes in the nature of data, their availability, the way in which they are collected, integrated, and disseminated are a big challenge for all those working with designed data from surveys as well as organic data. IPSDS was developed in response to the increasing demand from researchers and practitioners for the appropriate methods and right tools to face these changes. We offer a multidisciplinary curriculum, world-class faculty, and a web-based learning environment that allows you to take courses from anywhere in the world.

# Cooperation

## University Partners

- University of Maryland
- University of Michigan

- Catholic University of Santiago de Chile
- Australian National Unversity
- Beijing University
- Ashoka University (expressed interest)
- U. of Capetown (planned)

## Other Partners

- SRO - Michigan
- PEW
- German Record Linkage Center
- GESIS
- Bureau of Labour Statistics
- U.S. Census Bureau
- Statistics Netherlands

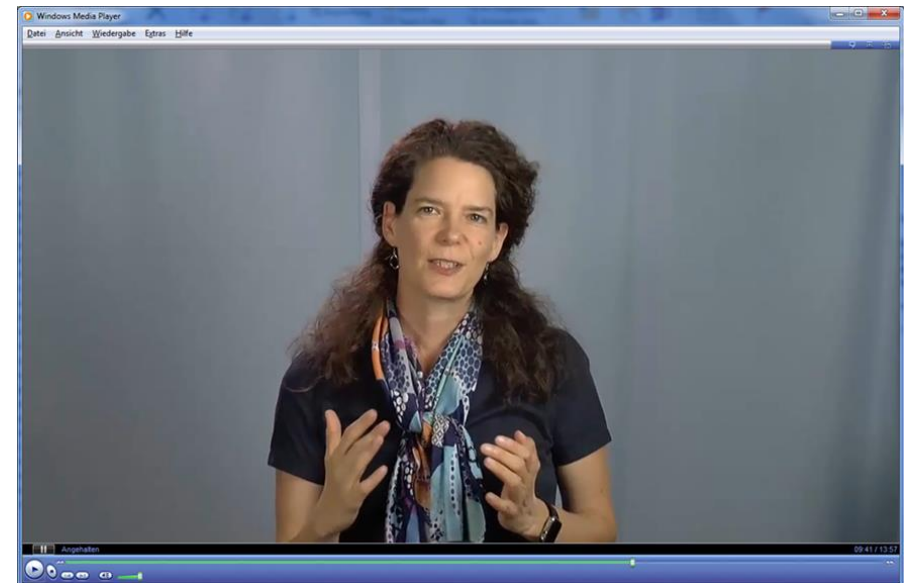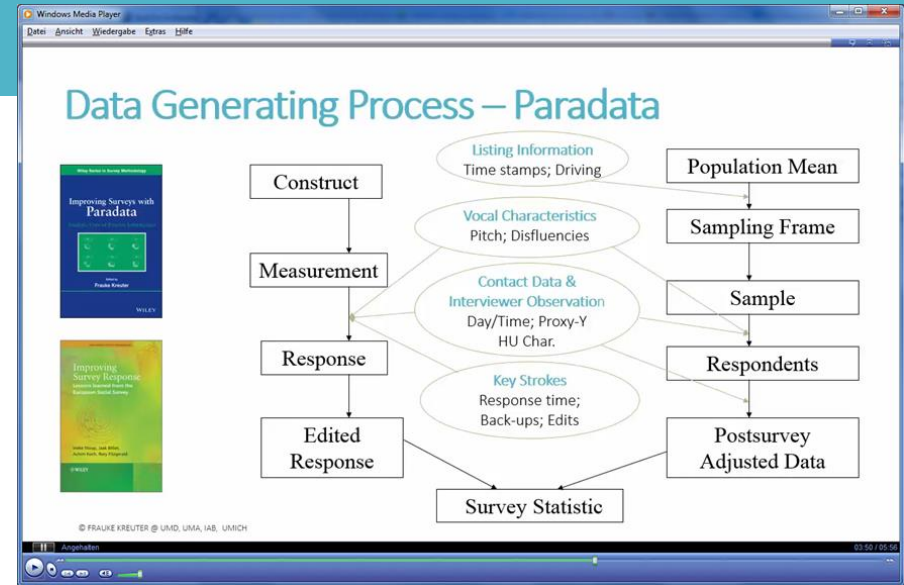| Data Output/Access | min. 3 credits/ 6 ECTS | Ethics 1 credit/2 ECTS | Data Confidentiality and Statistical Disclosure Control 2 credits/4 ECTS | Visualization 2 credits/4 ECTS | |
|---|---|---|---|---|---|
| Data Analysis | min. 6 credits/ 12 ECTS | GLM 3 credits/6 ECTS | Analysis of Complex Data 3 credits/6 ECTS | Propensity Score/Statistical Matching 3 credits/6 ECTS | Machine Learning I-III 1 credit/2 ECTS each |
| Data Curation/Storage | min. 3 credits/ 6 ECTS | Database Management 3 credits/6 ECTS | Data Munging I-III 1 credit/2 ECTS each | | |
| Data Generating Process | min. 4 credits/ 8 ECTS | Data Collection 3 credits/6 ECTS | Record Linkage 1 credit/2 ECTS | Practical Tools for Sampling and Weighting 3 credits/6 ECTS | Applied Sampling 3 credits/6 ECTS | Experimental Design 3 credits/6 ECTS |
| Research Question | min. 3 credits/ 6 ECTS | Fundamentals of Survey and Data Science 3 credits/6 ECTS | | | |

# Format

Each week set of videos (pre-recorded)

Lectures are broken into easily digestible sessions to help students to better focus on the material

Engage with the material at their own pace

**Frauke**

## NAVIGATION

- 📁 Home
  - 📁 Current course
    - 📁 SURV751
      - 📁 Participants

## ⚙ ADMINISTRATION

- 📁 Course administration
  - ✏ Turn editing on
  - ⚙ Edit settings
  - 📁 Users
  - ▼ Filters
  - 📁 Reports
  - ▦ Grades
  - ☁ Backup
  - 📁 Question bank
- 📁 Switch role to...
- 📁 My profile settings

## Introduction

To join the weekly online meeting, go to www.bluejeans.com and enter the meeting ID (611682210) under the join meeting tab.
NOTE: Blue Jeans is not currently compatible with Google Chrome. Users should use Safari, Internet Explorer, or Firefox as your browser when using Blue Jeans.

- 💬 News forum
- 💬 Discussion Forum
- 📕 Course Notes
- 📁 Data sets included in the course Notes
- 📄 Introduction and Syllabus
- 📕 Intro to R for SPSS Users

This file contains notes from a previous shortcourse introducing R to SPSS users. All of the homework assignments in this course will require R so if you aren't familiar with R here is some supplementary materials for you to use to help familiarize yourself with this software (which is free!) including downloading R and using a powerful package called Rcommander to read in and manage data files within the R environment.

## Week 1

Bluejeans Join Meeting [Tuesday, 02/02/2016, 06:00 p.m.-07:00 p.m.]

- 📁 Readings Week 1
  - 📕 Kreuter-Peng 01 26 14_manuscript.pdf
  - 📕 Public Opin Q-2015-Japec-839-80.pdf
  - 📕 Public Opin Q-2016-Schober-poq_nfv048.pdf
- 🎯 0. Introduction Big Data_1
  *January 29, 2016  Mediasite Presenter*

## LATEST NEWS

Add a new topic...

Next week's Class (Feb 23)
10:36 PM, Feb 15  Trent Buskirk

Older topics ...

## 📅 UPCOMING EVENTS

- ☑ Quiz 4 (Quiz opens)
  Tomorrow, 2:46 AM
- ☑ Quiz 3 (Quiz closes)
  Wednesday, February 24, 11:05 AM
- 📋 HW 3 Assignment
  Sunday, February 28, 10:00 PM
- ☑ Quiz 4 (Quiz closes)
  Monday, February 29, 2:46 AM

Go to calendar...

New event...

## 🕐 RECENT ACTIVITY

Activity since Saturday, February 20, 2016, 10:51 AM
Full report of recent activity...

**COURSE UPDATES:**

Added File
HW 2 Solutions

**ASSIGNMENTS SUBMITTED:**

9:04 PM, Feb 20
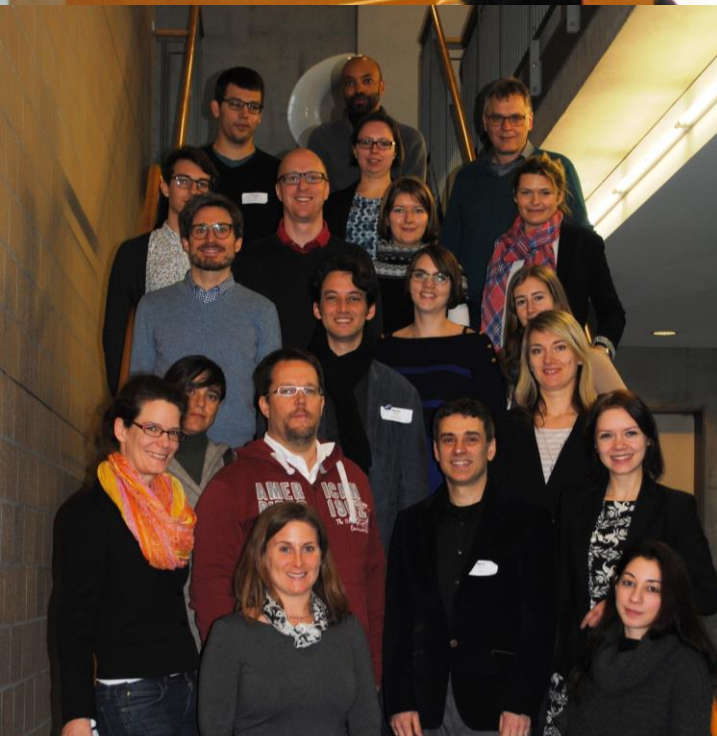
# Machine Learning Methods/Techniques

- There are many different machine learning methods available
- Many are non-parametric in nature and while a functional form can be specified, it is generally not a natural byproduct of the method
- Wu et al. (2008) provide an overview of ten of the top machine learning algorithms including (see http://bit.ly/1liWTir) :
  - ★ 🗎 K-means Clustering
  - 🗎 PageRank
  - ★ 🗎 K-nearest neighbors
  - 🗎 Support Vector Machines
  - 🗎 Decision Trees and Classification and Regression Trees
  - 🗎 Apriori Algorithm
  - 🗎 The EM Algorithm (Expectation-Maximization)
  - 🗎 Naïve Bayes
  - 🗎 Ensemble Methods (like AdaBoost and Random Forests).

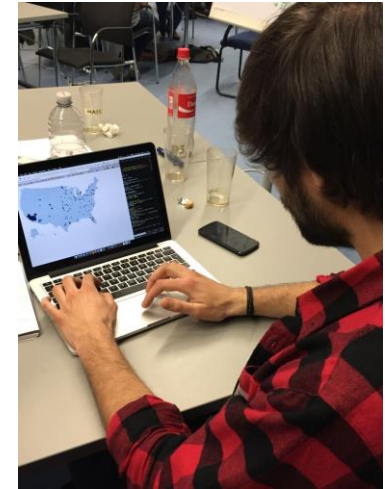0100000101010000010101000001001111101010010 Small Course **Big**

Kick-off 2/20/2016

# … recruitment and team work

# Interest

# Demand for our students

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

fkreuter@umd.edu