

Data sciences as essential job skill!

A social science, economics, and public policy training perspective

Frauke Kreuter

JPSM – Uni Mannheim – IAB

Paderborn 11/12/17



AAPOR Report on Big Data

AAPOR Big Data Task Force
February 12, 2015

Prepared for AAPOR Council by the Task Force, with Task Force members including:

Lilli Jasec, Co-Chair, Statistics Sweden
Frauke Kreuter, Co-Chair, JPSM at the U. of Maryland, U. of Mannheim & IAB
Marcus Berg, Stockholm University
Paul Biemer, RTI International
Paul Decker, Mathematica Policy Research
Cliff Lampe, School of Information at the University of Michigan
Julia Lane, American Institutes for Research
Cathy O'Neil, Johnson Research Labs
Abe Usher, HumanGeo Group

Acknowledgement: We are grateful for comments, feedback and editorial help from Eran Ben-Porath, Jason McMillan, and the AAPOR council members.

The National Academies of
SCIENCES • ENGINEERING • MEDICINE

REPORT

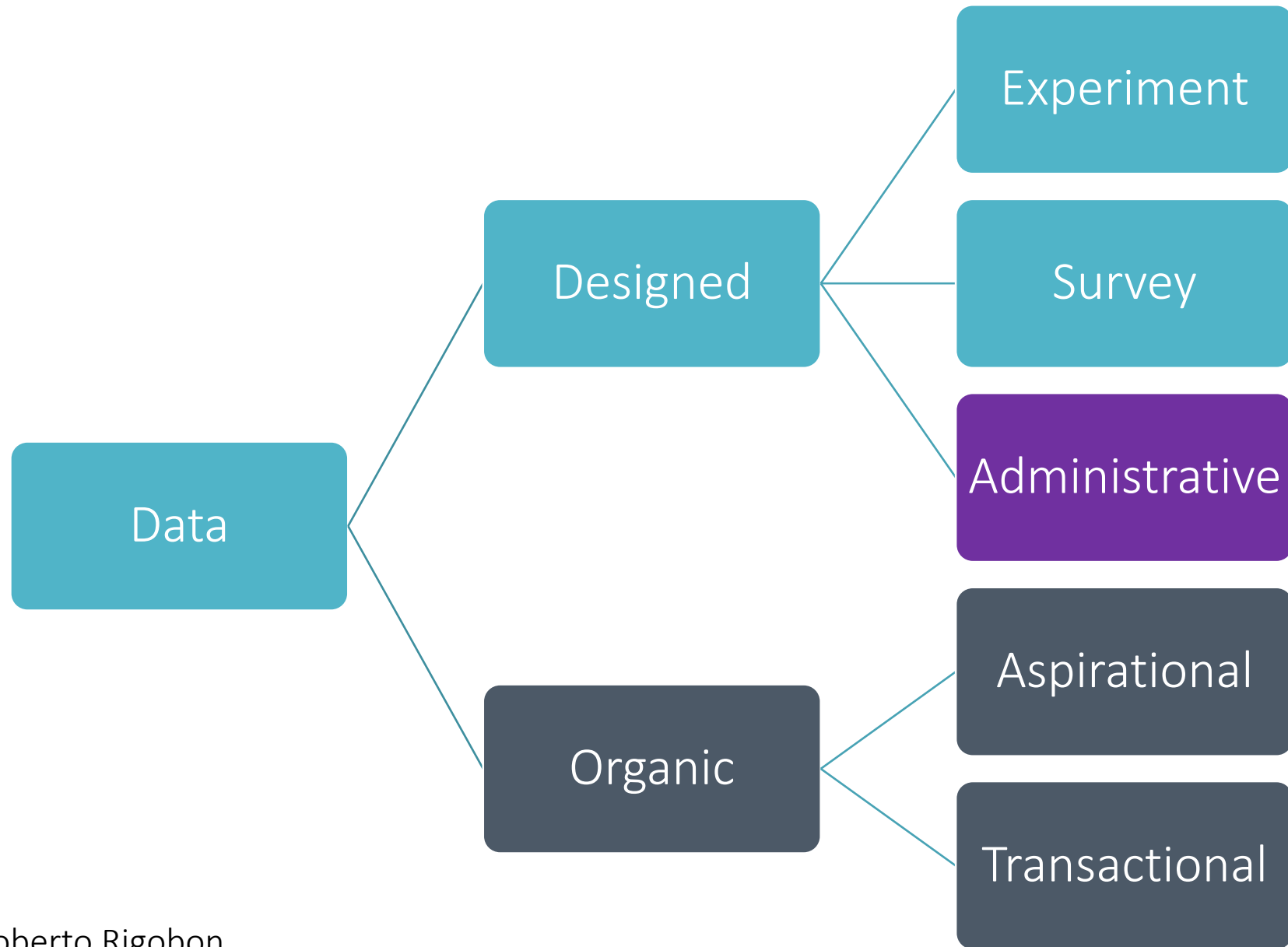
INNOVATIONS IN FEDERAL STATISTICS

Combining Data Sources While
Protecting Privacy

THE PROMISE OF EVIDENCE-BASED POLICYMAKING

Report of the Commission on Evidence-Based Policymaking





Source: Roberto Rigobon

Social Science

Marienthal study:

- social inclusion,
- coping behavior,
- day structure of the unemployed

MoDeM @ IAB

Sensor data replace time intensive observations in original study

- Radius of action
- Walking speed, sports activities
- Social networks
- Media use

...

Linkage with survey and admin data

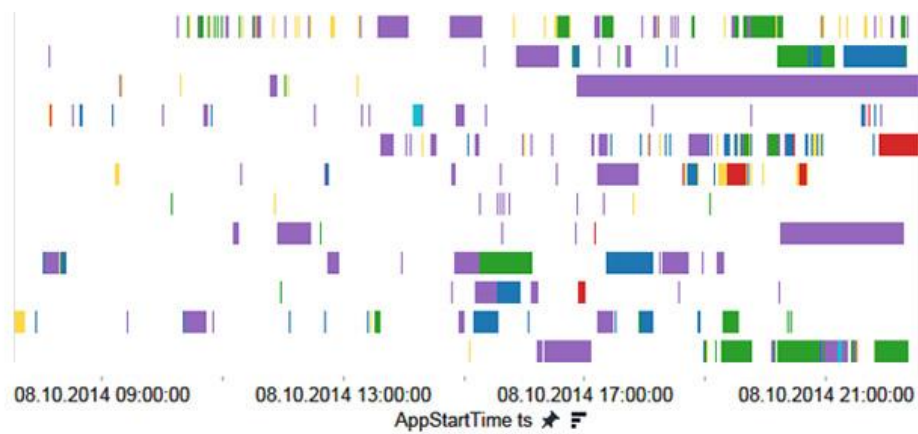
With a new introduction by **Christian Fleck**

MARIENTHAL

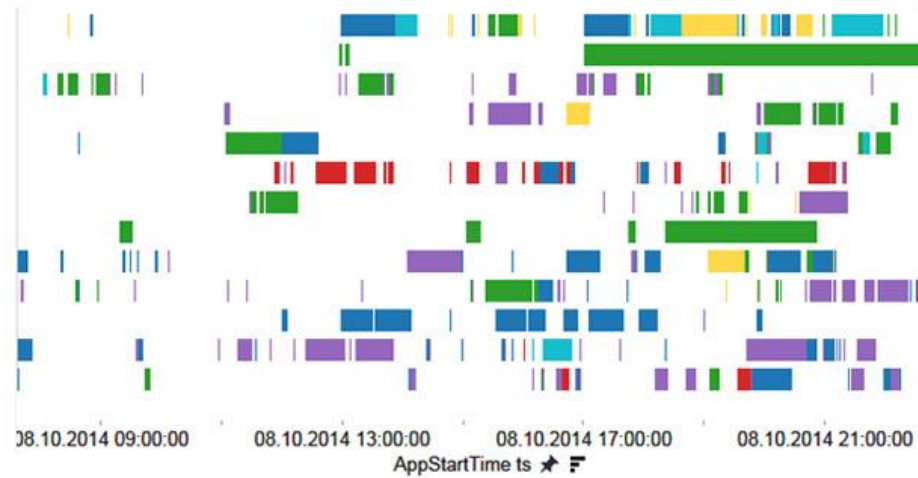


**The Sociography of an
Unemployed Community**

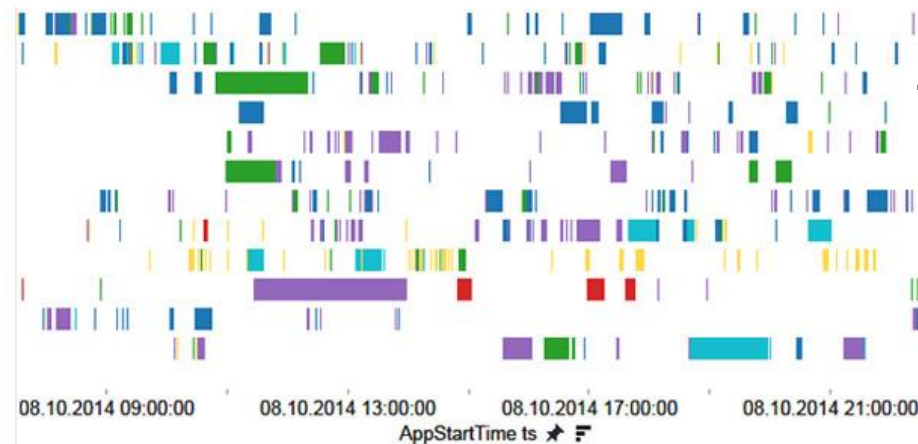
**Marie Jahoda, Paul F. Lazarsfeld,
and Hans Zeisel**



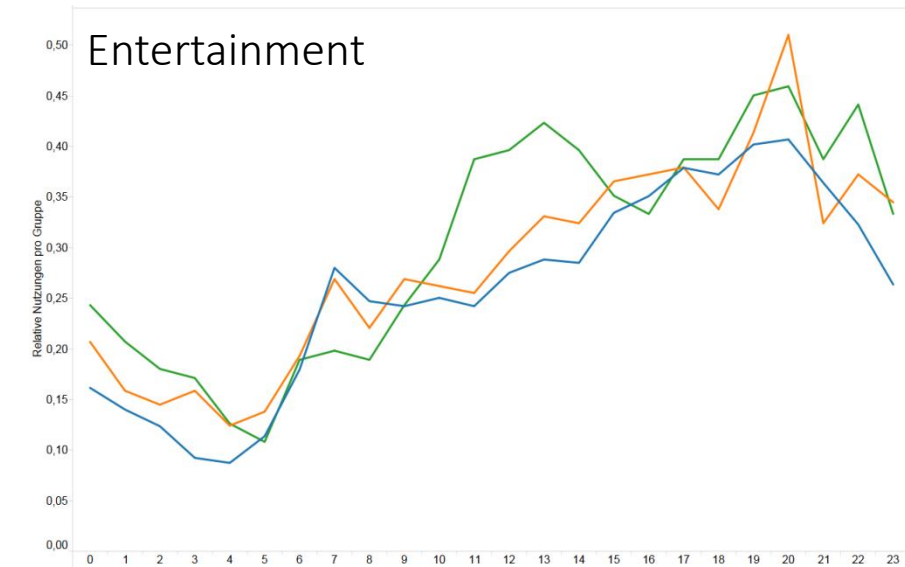
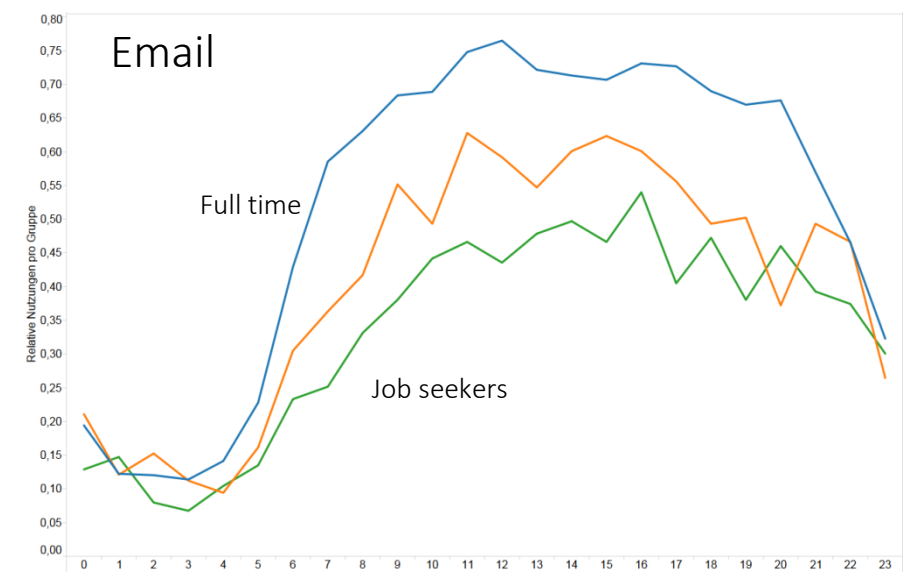
Full-time employed
→ App use past 5pm



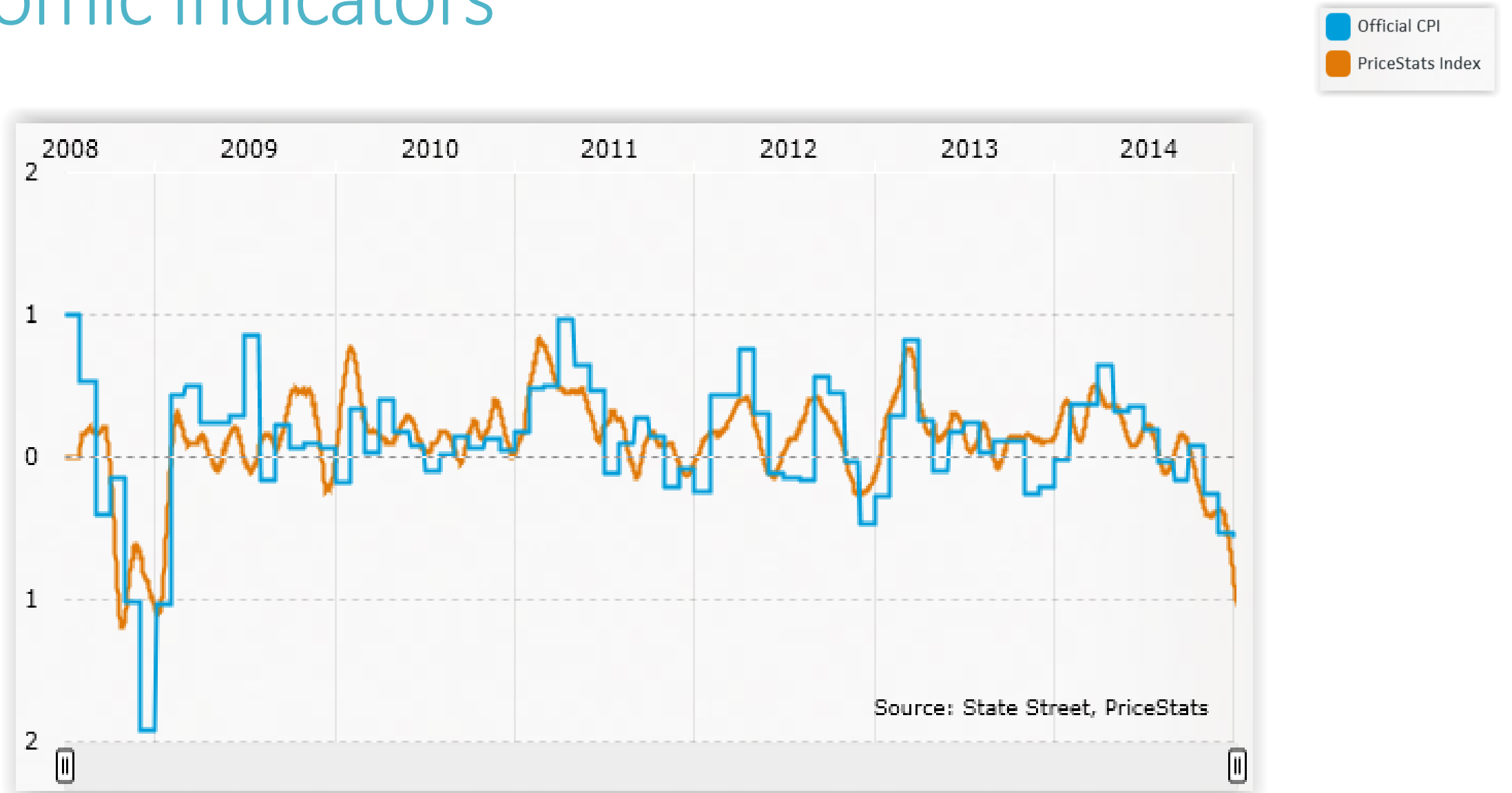
Part-time employed
→ App use at noon



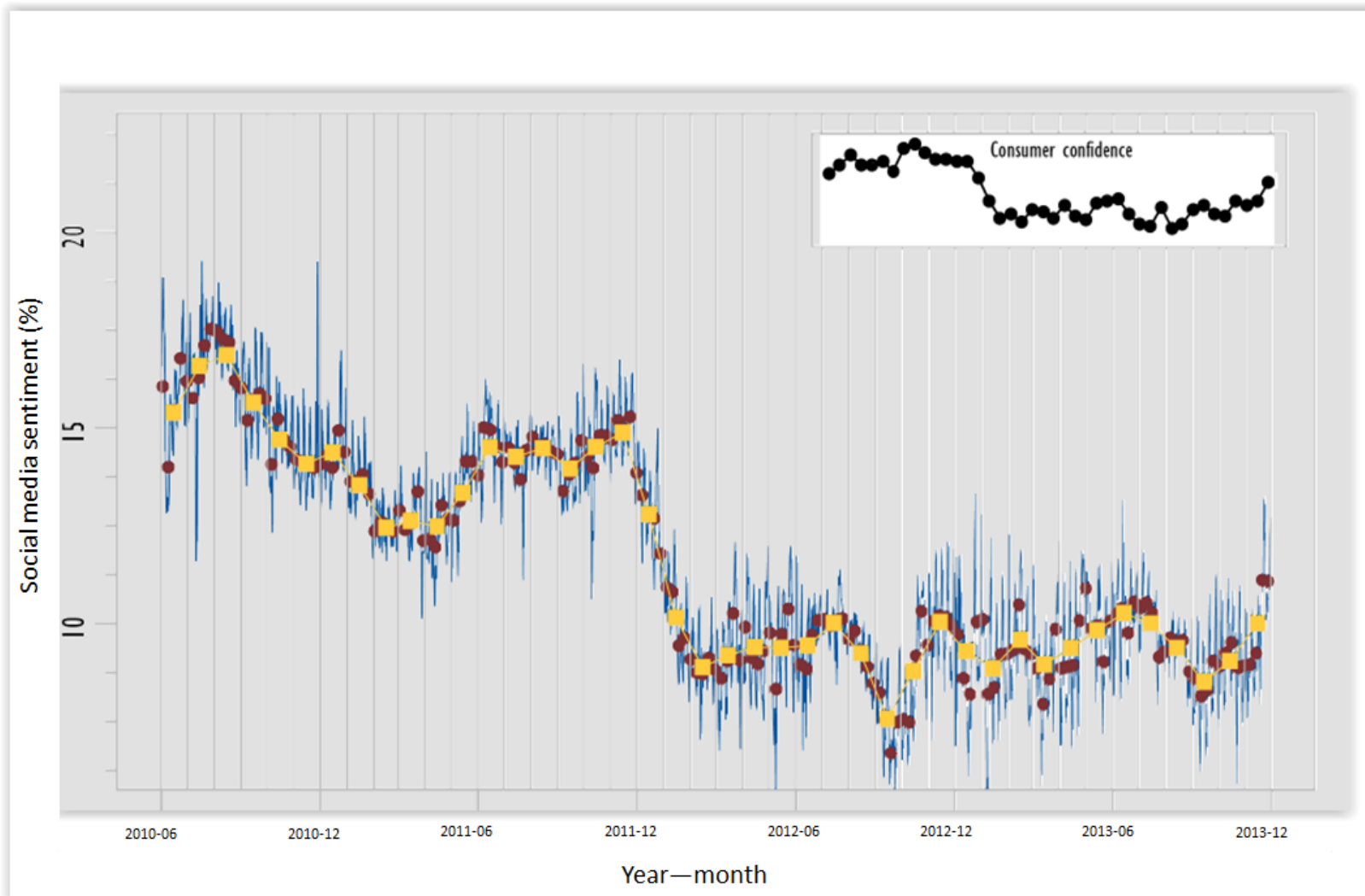
Job seekers
→ Continuous app use



Economic Indicators



US Aggregated Inflation Series, Monthly Rate, PriceStats Index vs. Official CPI the PriceStats website. . 1/1/2015



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).

Public Policy

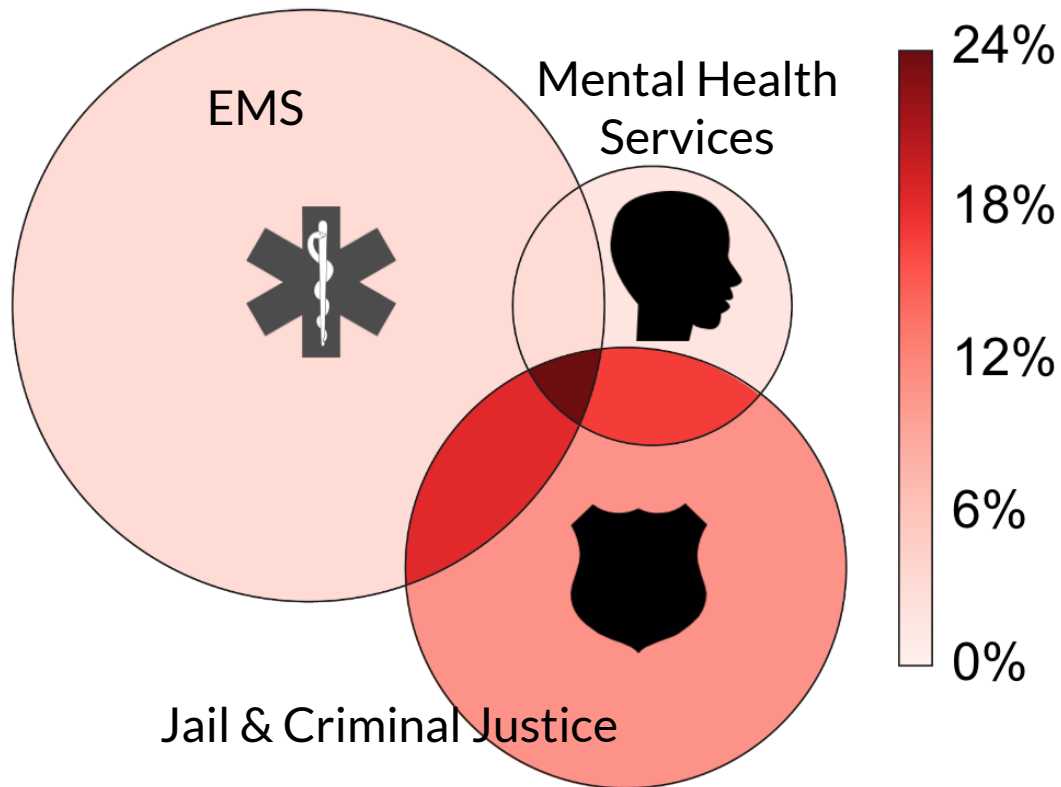
11 million people move through 3,100 Jails

\$22 Billion in costs

- 64 % suffer from mental illness,
- 68% have a substance abuse disorder
- 44 % suffer from chronic health problems



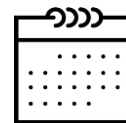
Hope: Combined data and predictive systems can support **targeted, preventative interventions** to help people at **risk of interactions** with the criminal justice system



Of the top 200 predicted individuals



104 went to jail within 1 year



19 years total jail time

Skills



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 enables large scale computation

RESEARCHER

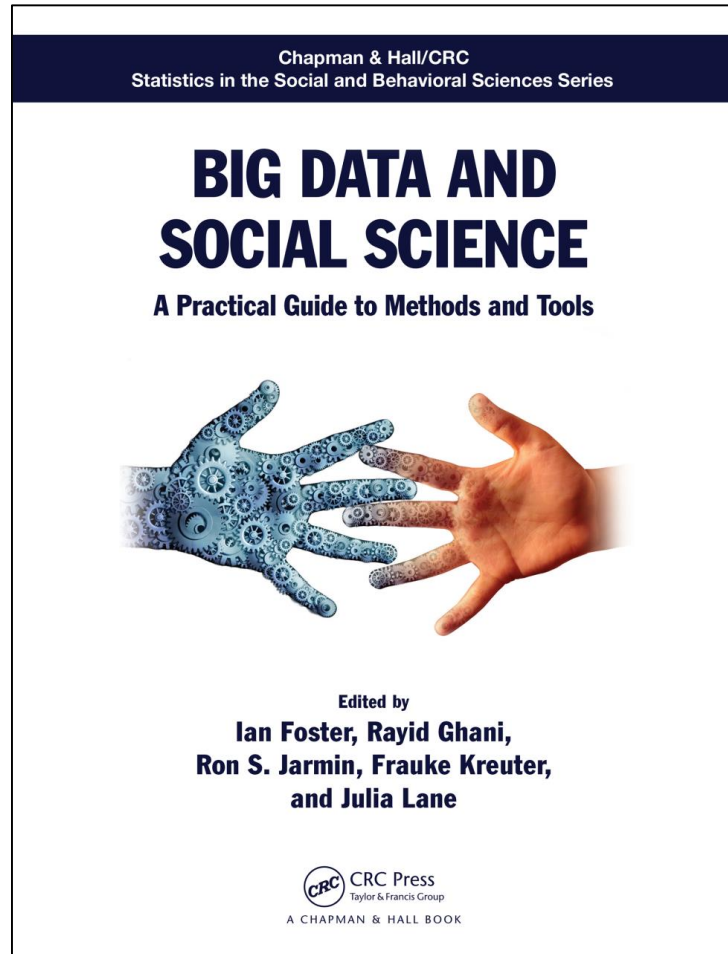
Team member with experience applying formal research methods, including survey methodology and statistics

COMPUTER SCIENTIST

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

1st Example – Coleridge Initiative

Professional Training Workshops



Three Classes

- Different cohorts (ex-offenders, welfare recipients and veterans)
- Joined with housing, transportation and jobs data

Class Format

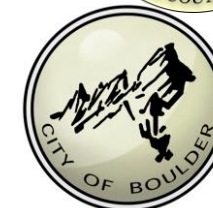
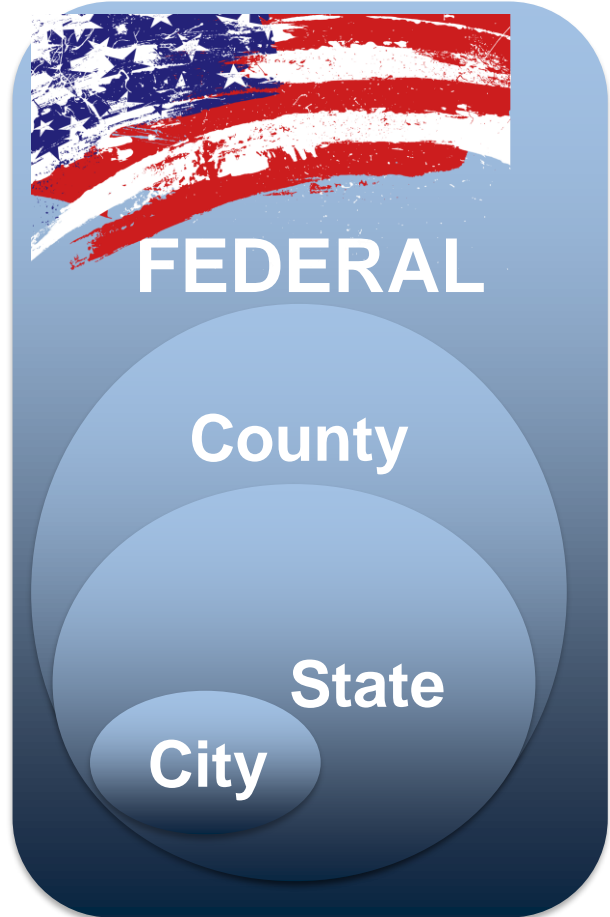
- Module 1: Foundations – Research Questions, Python, SQL
- Module 2: Data Acquisition – Web Scraping, API, Record Linkage
- Module 3: Data Analysis – Machine Learning, Networks, Text, Spatial
- Module 4: Visualization, Inference, Ethics, Privacy

Additional Information

- Final reports are all virtual
- Teaching Assistants and facilitators will be at each site for each module

Networks: The first two classes brought together

~40 agencies from city, state, county and federal agencies



Approach



Approach

Taught using **hands-on projects** with **real microdata** in a **secure environment** so that participants can learn the basics of how to:

- Code and collect new data
- Work with spatial data
- Manage complex data,
- Apply machine learning, text and network analysis
- Visualize relationships
- Address inference issues
- Manage privacy and confidentiality

Content Example

Problem

[Go back to Table of Contents](#)

First, turning some
can you take base

Four Main

- **Description:**
- **Prediction:**
- **Detection:**
- **Behavior Change:**

Model Evaluation

[Go back to Table of Contents](#)

In this phase, you take the predictors from your test set and apply your model to them, then assess the quality of the model by comparing the *predicted values* to the *actual values* for each record in your testing data set.

- **Performance Estimation:** How well will our model do once it is deployed and applied to new data?

Now let's use the model we just fit to make predictions on our test dataset, and see what our accuracy score is:

Python's [scikit-learn](#) is a commonly used, well documented Python library for machine learning. This library can help you split your data into training and test sets, fit models and use them to predict results on new data, and evaluate your results.

We will start with the simplest [LogisticRegression](#) model and see how well that does.

You can use any number of metrics to judge your models (see [model evaluation](#)), but we'll use [accuracy_score\(\)](#) (ratio of correct predictions to total number of predictions) as our measure.

```
# Let's fit a model
from sklearn import linear_model
model = linear_model.LogisticRegression(penalty='l1', C=1e5)
model.fit( X_train, y_train )
print(model)
```

t action

Content Example: API

Isochrone API

- Back to the [Table of Contents](#)

Gives the area (as a polygon) a traveler can reach from a location. The user can set if so desired, [description here](#)

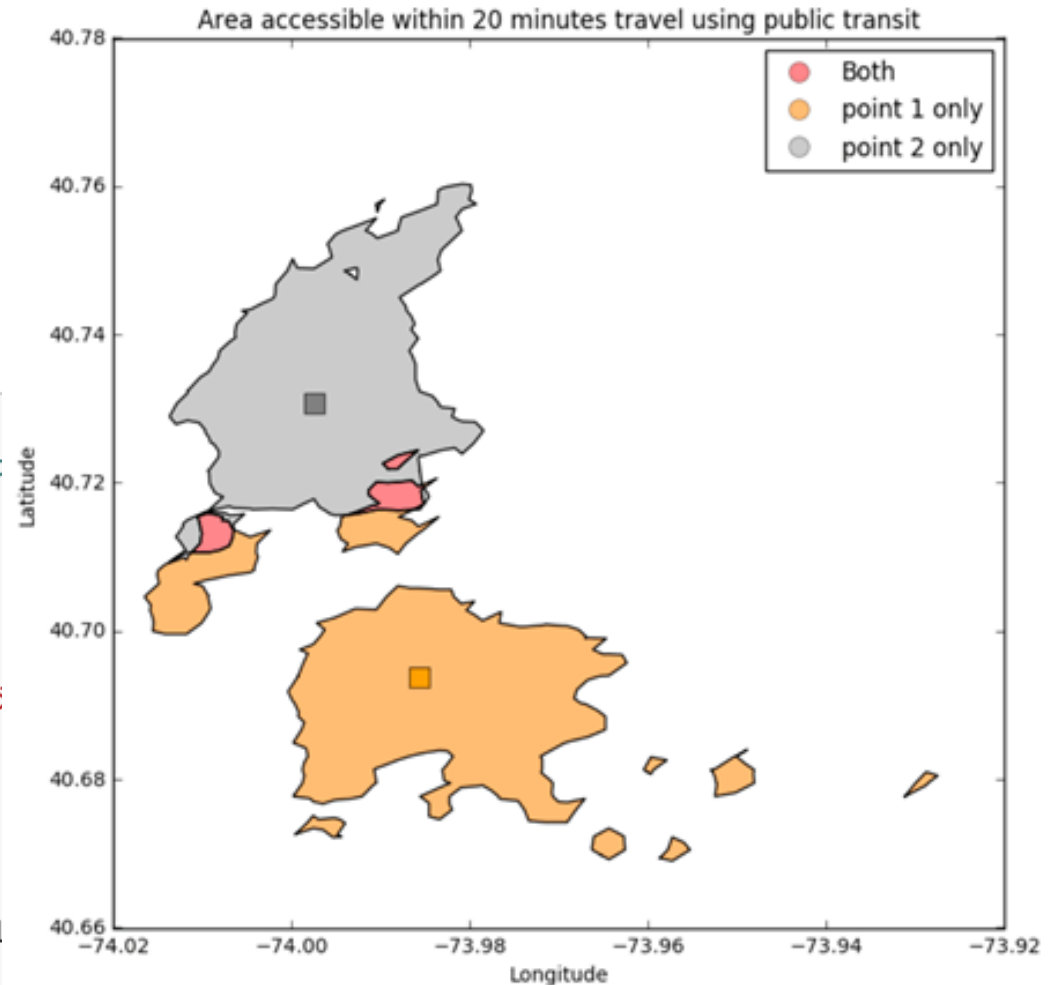
```
: # set start location
start_point = [40.693856, -73.985754] # Downtown Manhattan

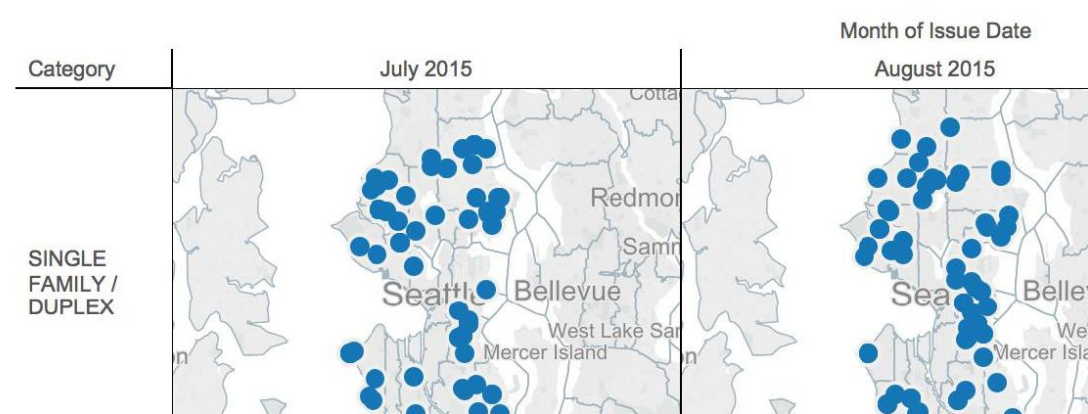
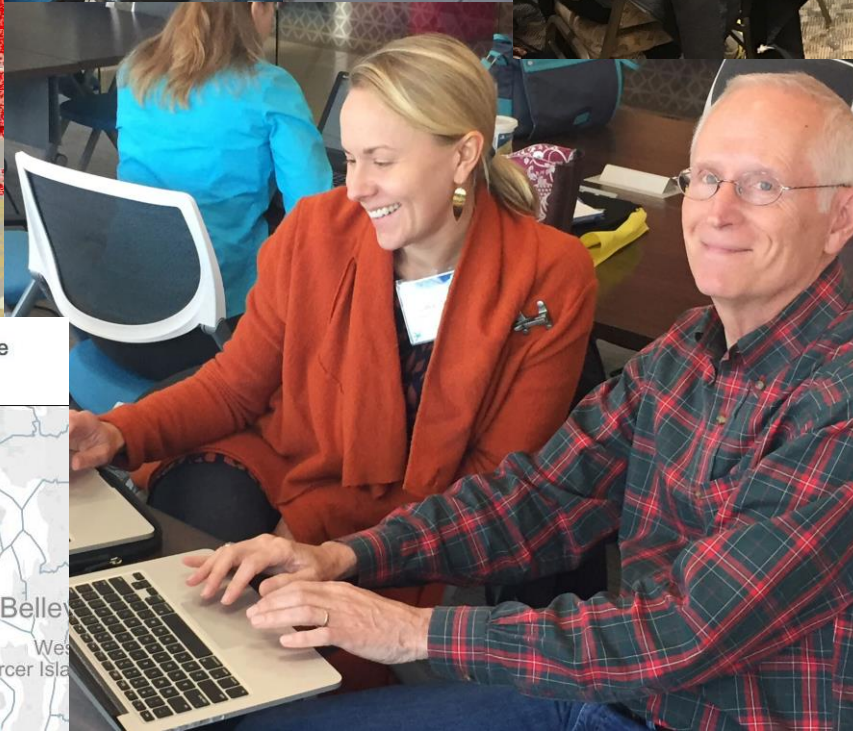
# set travel time - in seconds
travel_time = 60*20 # 20 minutes

# create query URL
qry_url = '{}isochrone?fromPlace={},{}&mode=transit&time={}'.format(
    base_url, start_point[0], start_point[1], travel_time)

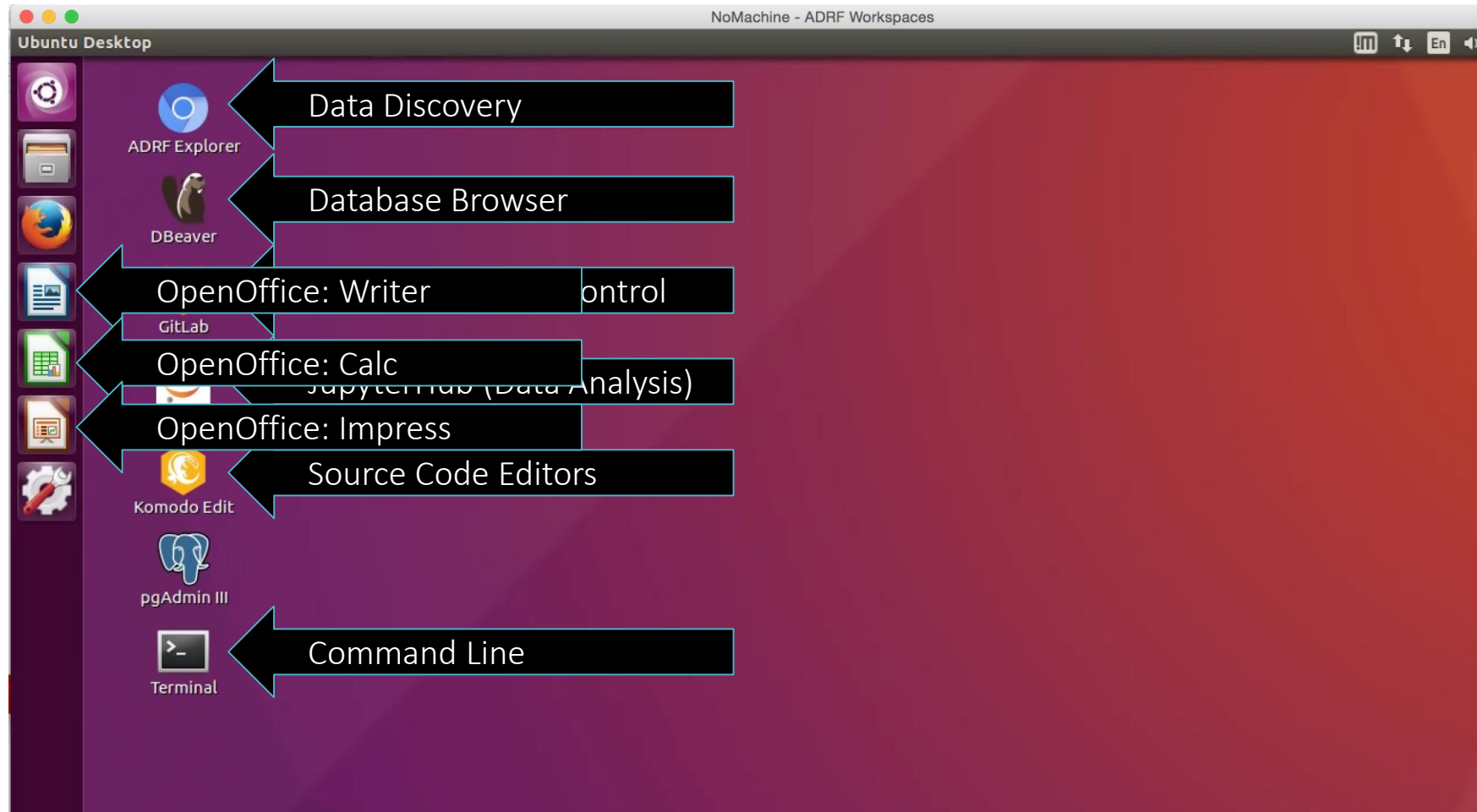
# get json request
iso_json = json.loads(requests.get(qry_url).text)

## load isochrone into a geopandas dataframe
iso_gdf = gpd.GeoDataFrame.from_features(iso_json['features'])
```





Collaborative secure environment



What our **participants** say about the program

"Love the Jupyter notebooks!! ... I love how the code snippets and explanations are set up in the Jupyter notebooks. The format of going through it individually and discussing questions/challenges in your group, with the experts available when needed, worked really well for my learning style."



Danielle Fulmer
Director of Business
Analytics



I could see our agency benefiting potentially from something like this in that, as the system builds out and collects additional resources/datasets that impact criminal justice system practices, this may be an option for a place for us to look for the results of studies using evidence based practices.



Katy Fitzgerald
Management Analyst



What they put on their *out of office*

"Thanks for your email. I will be away from the office starting on April 18th and returning Monday, April 24th."

Where did I go?

Why so long?

Why?

What can I do?

University of Maryland - College Park, Maryland

I am at the most awesome, intense, hands-on training for using data to improve public policy.

Noemi Reyes and I were accepted into the program – AND received a very generous scholarship from the Laura and John Arnold Foundation.

If you have data you would like to contribute, send me an email with the Subject: Dane Data – that will get the attention of the little robot I set up to prioritize that message! We can set up a time to chat when I get back."



Noyemi Reyes
Research Analyst



2nd Example – International Program in Survey and Data Science



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



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.

Problem we tried to solve – In brief

- Allow for multidisciplinary curriculum
- Modularized – adapt to prior skills and work needs
- Relevant methods and tools
- Mix of faculty from academia and industry

Key elements:

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

Partners and Funding

University Partners

- University of Maryland
- University of Mannheim
- Catholic University of Santiago de Chile
- Australian National University
- 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

SPONSORED BY THE



Federal Ministry
of Education
and Research



Modules

Data Output/Access

Learn how to communicate results, distribute and store your data; Ethics

Data Analysis

Learn a variety of analysis methods suited for different data types

Data Curation/Storage

Learn how to curate and manage data

Data Generating Process

Understand how to collect data yourself, and how data are generated through administrative and processes.

Research Question

Learn how to formulate your research goal and which data are best suited to achieve this goal.

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
participants to better focus on the
material

Engage with the material at their
own pace

The screenshot displays a course interface. On the left, a file explorer lists items for Week 2 and Week 3. Week 2 includes '2. Model Eval Validation', '3. K-Means Clustering', 'Homework Assignment 1', 'Quiz 2', and 'HW Number 1 Solutions'. Week 3 includes '4. K-Nearest Neighbors', '5. CARTS', 'HW 2 Assignment', 'Quiz 3', and 'HW 2 Solutions'. The central video player shows a presentation titled 'Machine Learning Methods/Techniques' with a list of topics: 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, and Ensemble Methods (like AdaBoost and Random Forests). The video is playing at 02:15 / 44:08. On the right, a sidebar shows 'HW 2 Assignment' at 8:58 PM, Feb 21, by Man Kaiwen.

2. Model Eval Validation
January 12, 2016 Mediasite Presenter

3. K-Means Clustering
January 12, 2016 Mediasite Presenter

Homework Assignment 1
data file for homework number 1
Tasks for Homework Number 1

Quiz 2
HW Number 1 Solutions

This is a .R file that can be opened using Notepad or other text editor (or Word) for the tasks of HW 1.

Week 3

Bluejeans Join Meeting [Tuesday, 02/16/2016]

4. K-Nearest Neighbors
January 12, 2016 Mediasite Presenter

5. CARTS
January 12, 2016 Mediasite Presenter

HW 2 Assignment
Tasks for HW Number 2
Datasets for HW 2

Quiz 3
HW 2 Solutions

Here is the R script file containing the solutions for HW 2.

3. K-Means Clustering
www.jpsmcclases.umd.edu/Mediasite/Play

Machine Learning Methods/Techniques

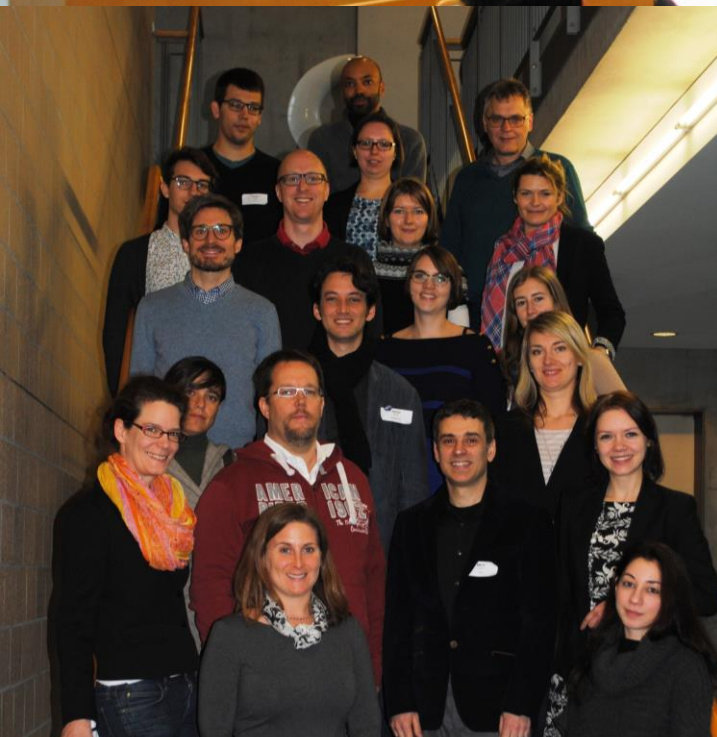
- There are many different machine learning methods available
- Many are non-parametric in nature and while a functional form can be specified, it is generally not a natural byproduct of the method
- Wu et al. (2008) provide an overview of ten of the top machine learning algorithms including (see <http://bit.ly/liWTir>):
 - 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).

01000001010000001010100000101011101010010 Small Course Big

Playing
02:15 / 44:08

HW 2 Assignment
8:58 PM, Feb 21
Man Kaiwen
HW 2 Assignment

Annual „Connect“ Event (next June 9, 2018)



Lessons Learning

- Modular approach much appreciated by working professionals
- Learning with application at hand is key
- Participants from all sectors and disciplines
- Very high demand on graduates
- Privacy and confidentiality very important
- Hardest to learn and hardest to teach:
Asking the right question!

<http://coleridgeinitiative.org>
<http://survey-data-science.net/>

fkreuter@umd.edu