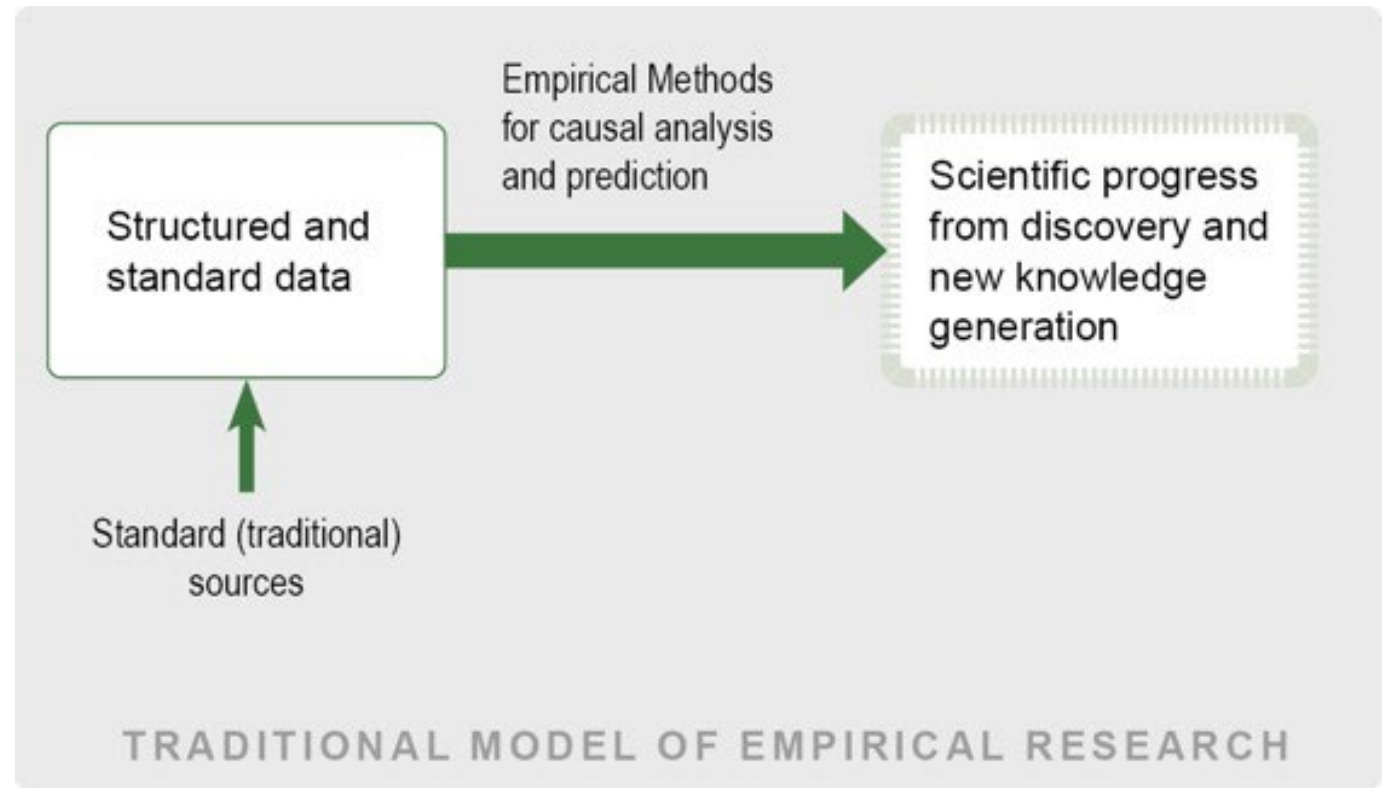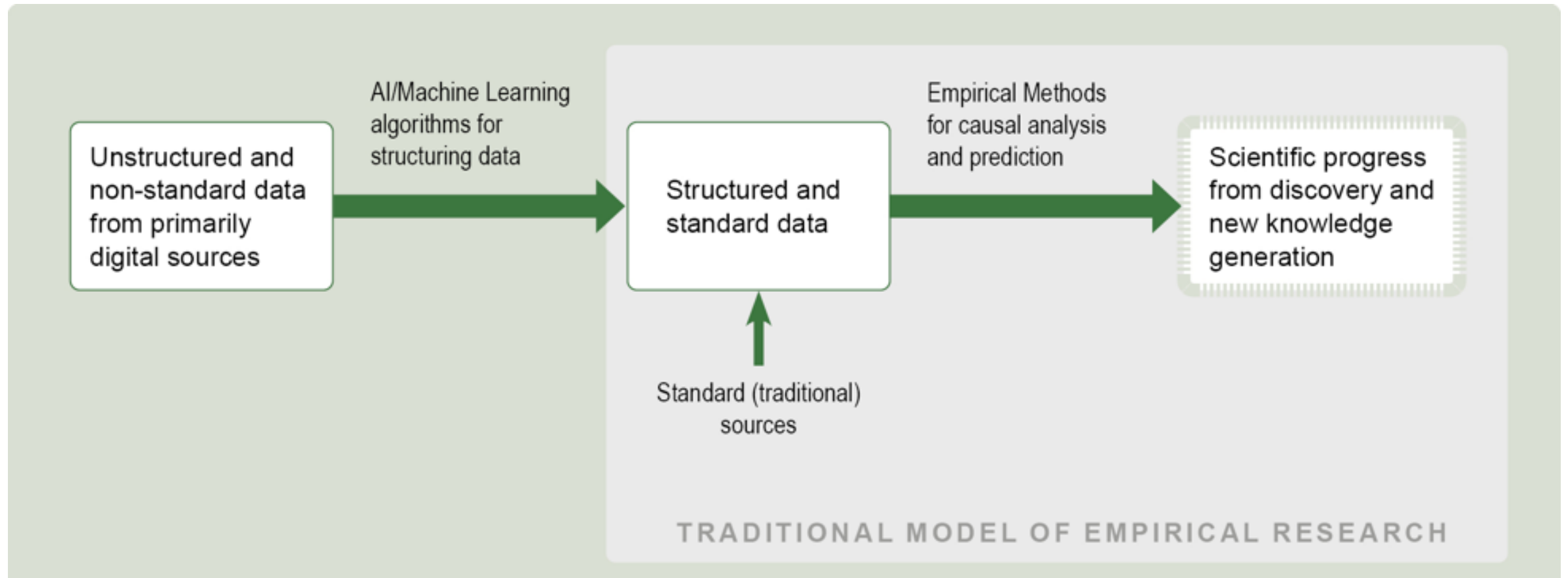# New Measurements in New Data Infrastructures.

## The long road from vision to practice.

**Frauke Kreuter**

Joint Program in Survey Methodology, University of Maryland
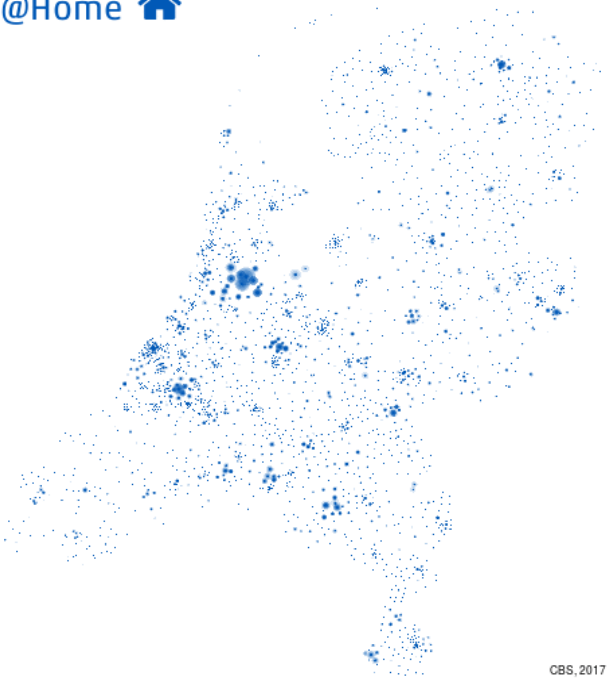LMU Munich Statistics and Data Science

Structured and standard data
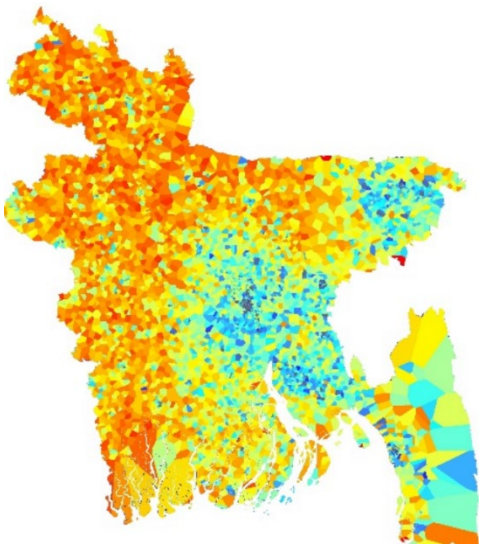
Empirical Methods for causal analysis and prediction

Scientific progress from discovery and new knowledge generation

Standard (traditional) sources

TRADITIONAL MODEL OF EMPIRICAL RESEARCH

Unstructured and non-standard data from primarily digital sources

AI/Machine Learning algorithms for structuring data

Structured and standard data

Empirical Methods for causal analysis and prediction

Scientific progress from discovery and new knowledge generation

Standard (traditional) sources

TRADITIONAL MODEL OF EMPIRICAL RESEARCH

# New Linked Measurements

| Sustainable Communities | No Poverty | Decent Work and Economic Growth |
|---|---|---|

@Home 🏠



CBS, 2017

https://doi.org/10.6084/m9.figshare.c.3662800.v1

See also Steele, J. et al. (2017): „Mapping poverty using mobile phone and satellite data.", *Journal of the Royal Society Interface*, 14.
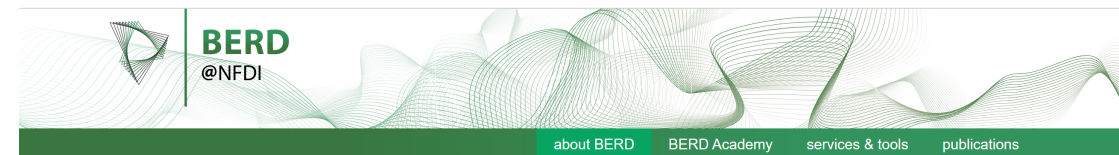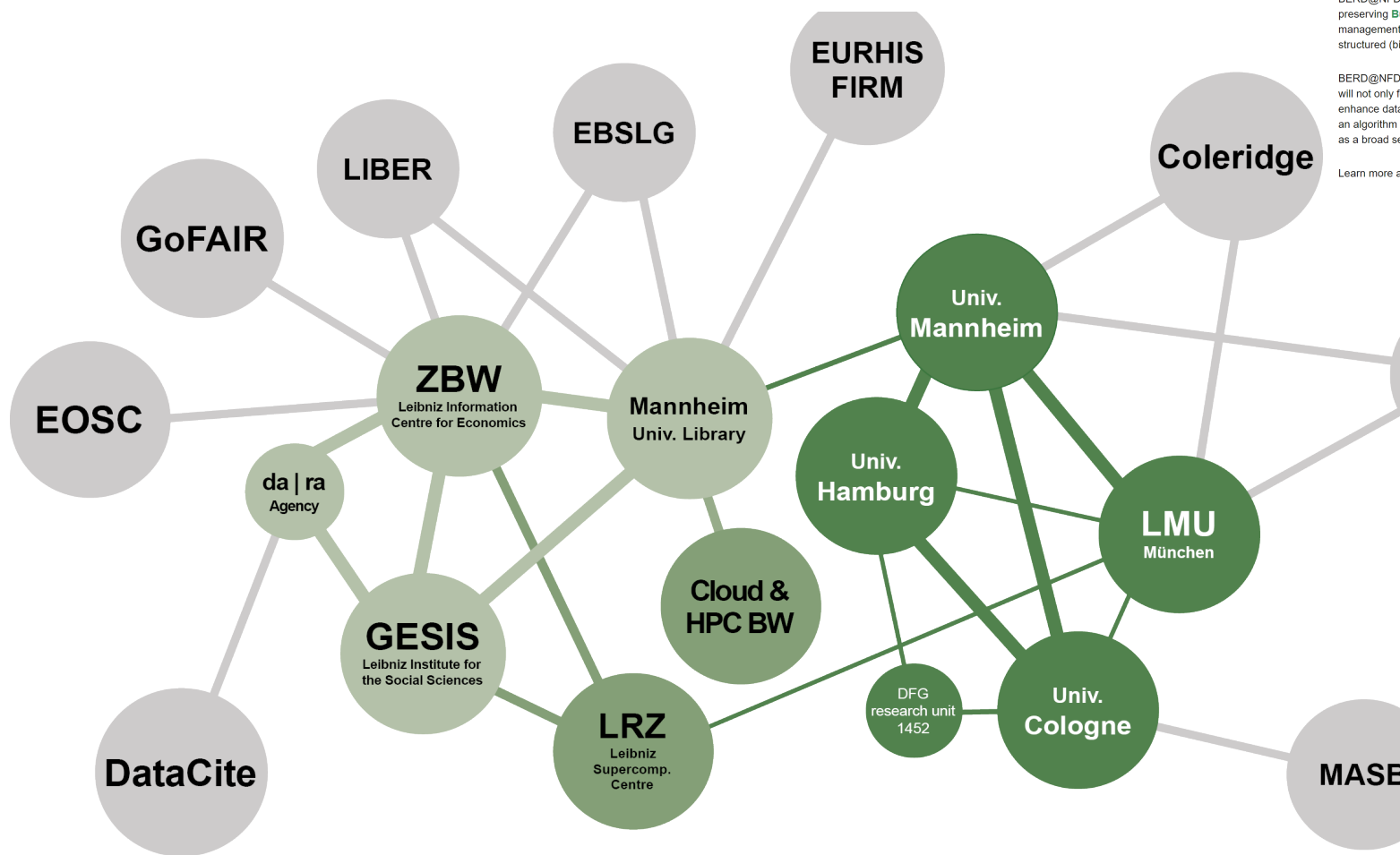
See also Jean, N. et al. (2016): „Use of satellite imagery and machine learning to predict poverty.", *Science*, 353(6301), 790-794.



See also Bansak, K. et al. (2018) Improving refugee integration through data-driven algorithmic assignment. Science, 359, 6373, 325-329.

# Context

# Lessons Learned

1. Creating new measures out of linked data is not enough. The effort has to be tied to a product.
2. Research questions need to guide decisions on measurements and data sources.
3. Data Science is a "Team Sport" … and needs to be treated as such.

# Lessons Learned

1. Creating new measures out of linked data is not enough. The effort has to be tied to a product.

# Democratizing our Data: A Challenge to Invest in Data and Evidence-based Policy

Kreuter, F., Ghani, R., & Lane, J. (2019). Change Through Data: A Data Analytics Training Program for Government Employees. *Harvard Data Science Review, 1*(2). https://doi.org/10.1162/99608f92.ed353ae3

https://coleridgeinitiative.org/



Training Module → Data producer     Data steward

Data user → Metadata          Metadata Data / Usage Feedback          Access Workflows Monitoring / Reporting

Data analysis Code Collaboration

**Documentation Module**
Explorer links metadata, codes, tools, publications

**Stewardship Module**
Approval workflow, monitoring, reporting

**Collaboration Module**
Interactive chat and code sharing

**Workspace and tools**
python    jupyter    R

Data in cloud
Alternative: local servers

**Security Module**
FedRAMP security certified

**Result**

- Trained Staff
- New Products
- New Network
- New Metrics
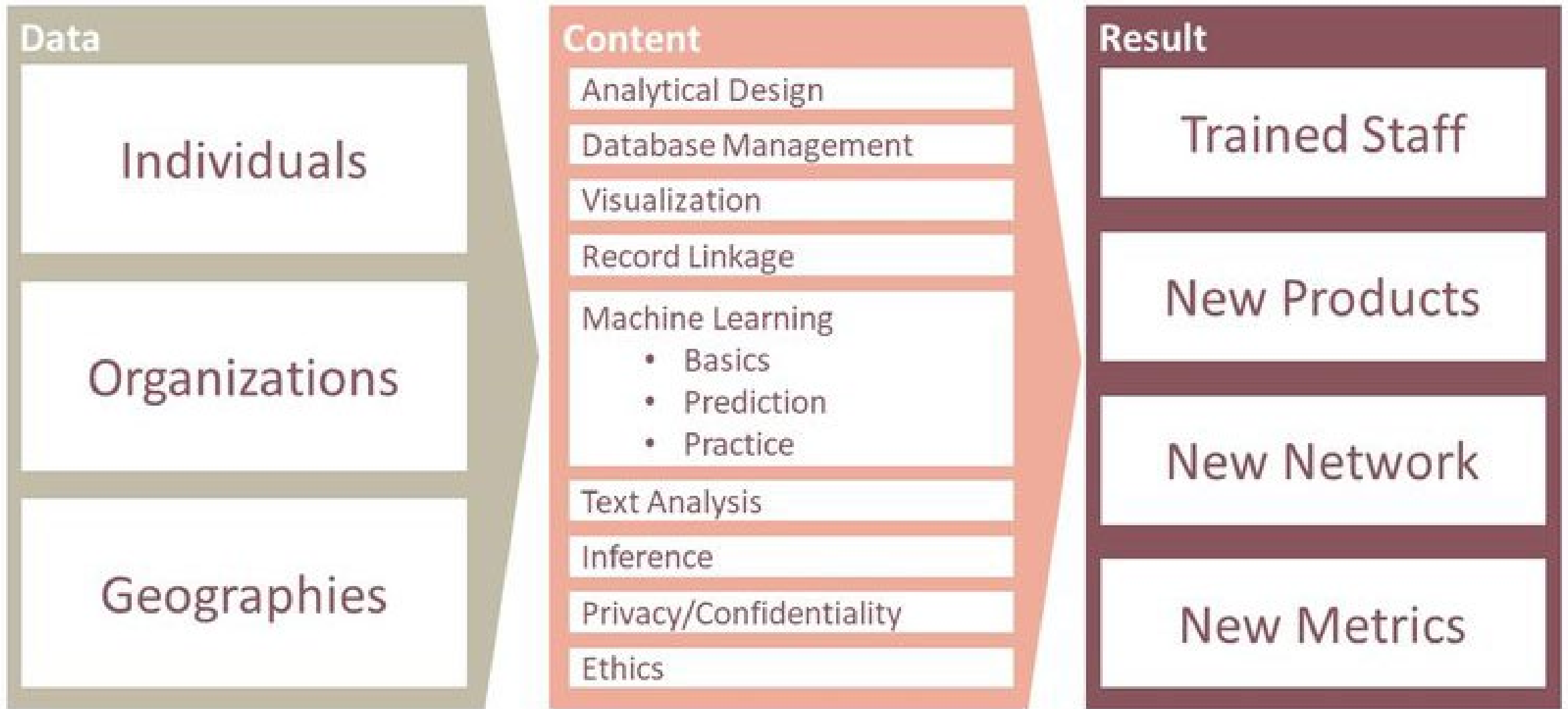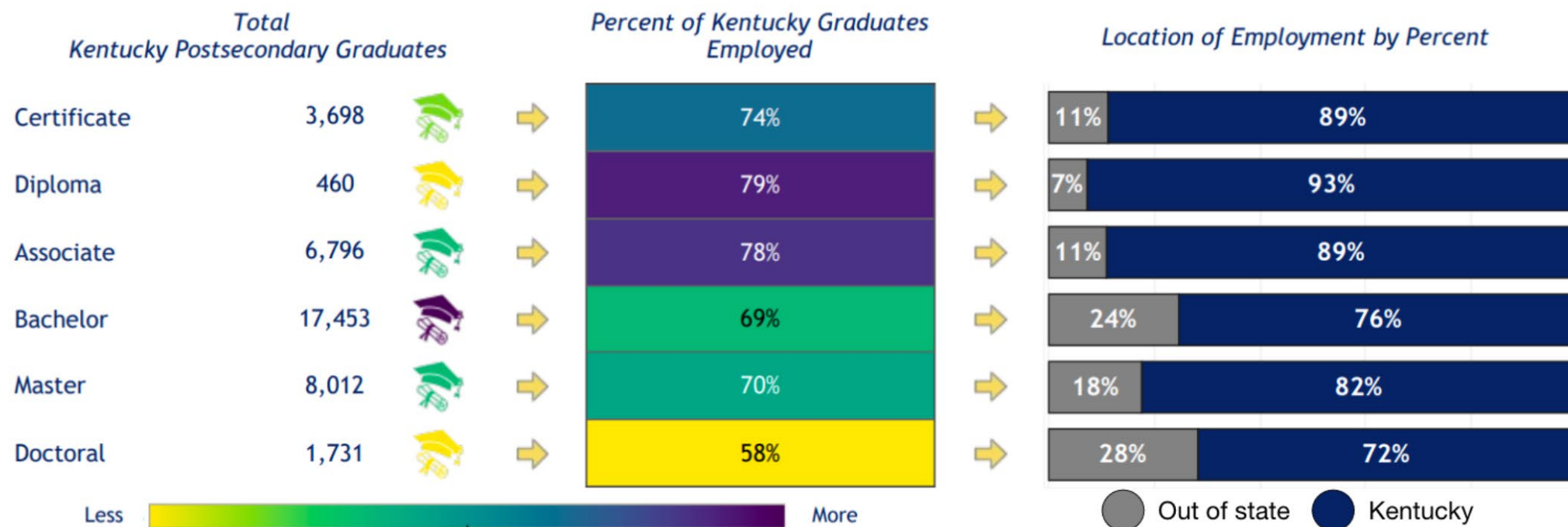
Kreuter, Ghani, Lane (2019): Change through data. A Data Analytics Training Program for Government Employees.
Harvard Data Science Review,1.2

Kreuter, Ghani, Lane (2019): Change through data. A Data Analytics Training Program for Government Employees. Harvard Data Science Review,1.2
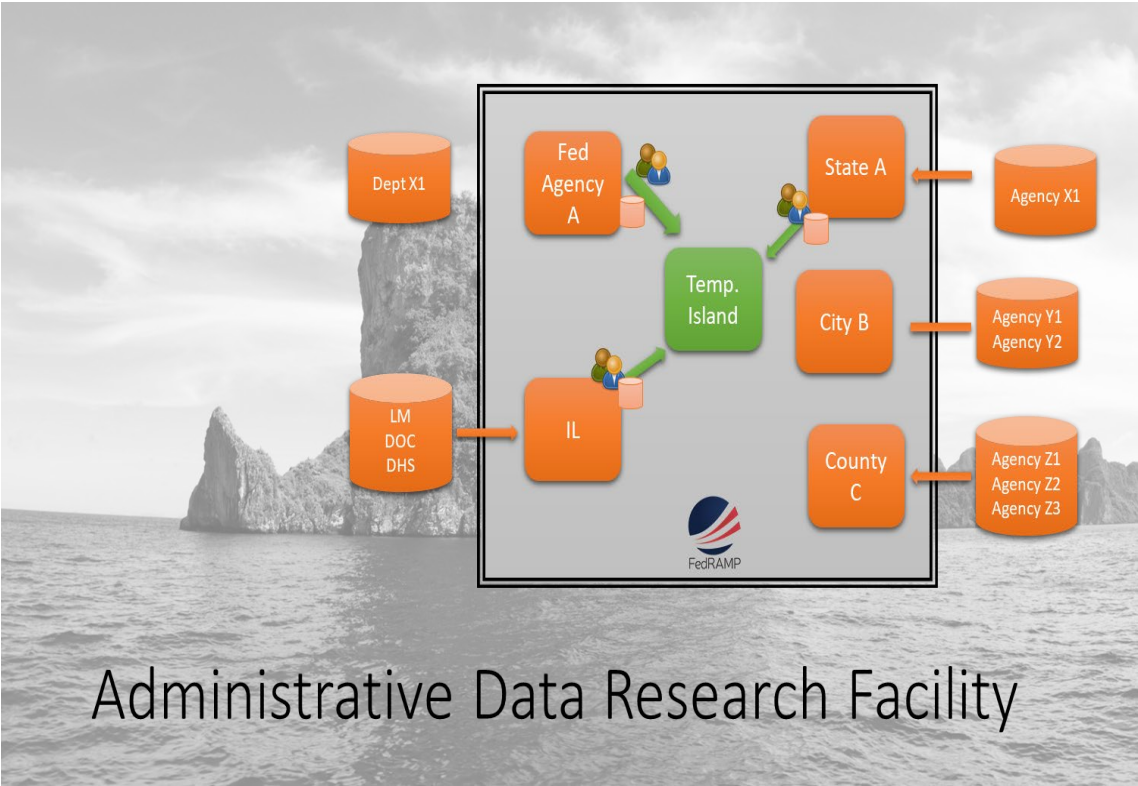
| Total Kentucky Postsecondary Graduates | | Percent of Kentucky Graduates Employed | Location of Employment by Percent | |
| --- | --- | --- | --- | --- |
| Certificate | 3,698 | 74% | 11% | 89% |
| Diploma | 460 | 79% | 7% | 93% |
| Associate | 6,796 | 78% | 11% | 89% |
| Bachelor | 17,453 | 69% | 24% | 76% |
| Master | 8,012 | 70% | 18% | 82% |
| Doctoral | 1,731 | 58% | 28% | 72% |

Less ──────────────────────── More

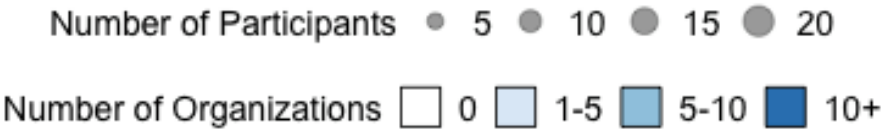Out of state    Kentucky

Enlarge Image

*Kentucky Postsecondary Graduate Outcomes by Credential Level: Five years post-graduation, includes all Kentucky graduates that were employed within Indiana, Kentucky, Ohio, or Tennessee from academic year 2013*

# Expansion



Administrative Data Research Facility

Quarter: 2016-Q4
Total Organizations: 0
Total Participants: 0



Number of Participants   ● 5  ● 10  ● 15  ● 20

Number of Organizations   ☐ 0  ☐ 1-5  ☐ 5-10  ☐ 10+

# Lessons Learned

2. Research questions need to guide decisions on measurements and data sources.

# Effects of Unemployment?



PSYCHOLOGISCHE MONOGRAPHIEN

DIE ARBEITSLOSEN VON
MARIENTHAL

EIN SOZIOGRAPHISCHER VERSUCH ÜBER DIE
WIRKUNGEN LANGDAUERNDER ARBEITSLOSIGKEIT

MIT EINEM ANHANG
ZUR GESCHICHTE DER SOZIOGRAPHIE

BEARBEITET UND HERAUSGEGEBEN VON DER
ÖSTERREICHISCHEN WIRTSCHAFTSPSYCHOLOGISCHEN
FORSCHUNGSSTELLE

VERLAG VON S. HIRZEL IN LEIPZIG 1933



MARIENTHAL

The Sociography of an
Unemployed Community

Marie Jahoda, Paul F. Lazarsfeld,
and Hans Zeisel



Source: Archives for the History of Sociology
in Austria (Graz), »Marienthal« Virtual Archives





Kreuter, Haas, Keusch, Bähr, Trappmann. 2018. "Collecting survey and smartphone sensor data with an app: Opportunities and challenges around privacy and informed consent." Social Science Computer Review

# Panel + Administrative Data as Frame

Sample of households with at least one welfare benefit recipient (at reference date)

Refreshed annually

Surveyed annually

Random household sample of resident population

Refreshed annually

Surveyed annually

**+**



Trappmann M., Christoph B., Achatz J., Wenzig C. (2009) PASS: a new panel study for labour market research, Int. J. of Manpower , 30, 7, pp.765-770

# Data from Smart Phone Sensors



- Network quality and location information (every half hour)

- Interaction history

- Characteristics of the social network

- Activity data (every two minutes)

- Smartphone usage

# Great Uptake

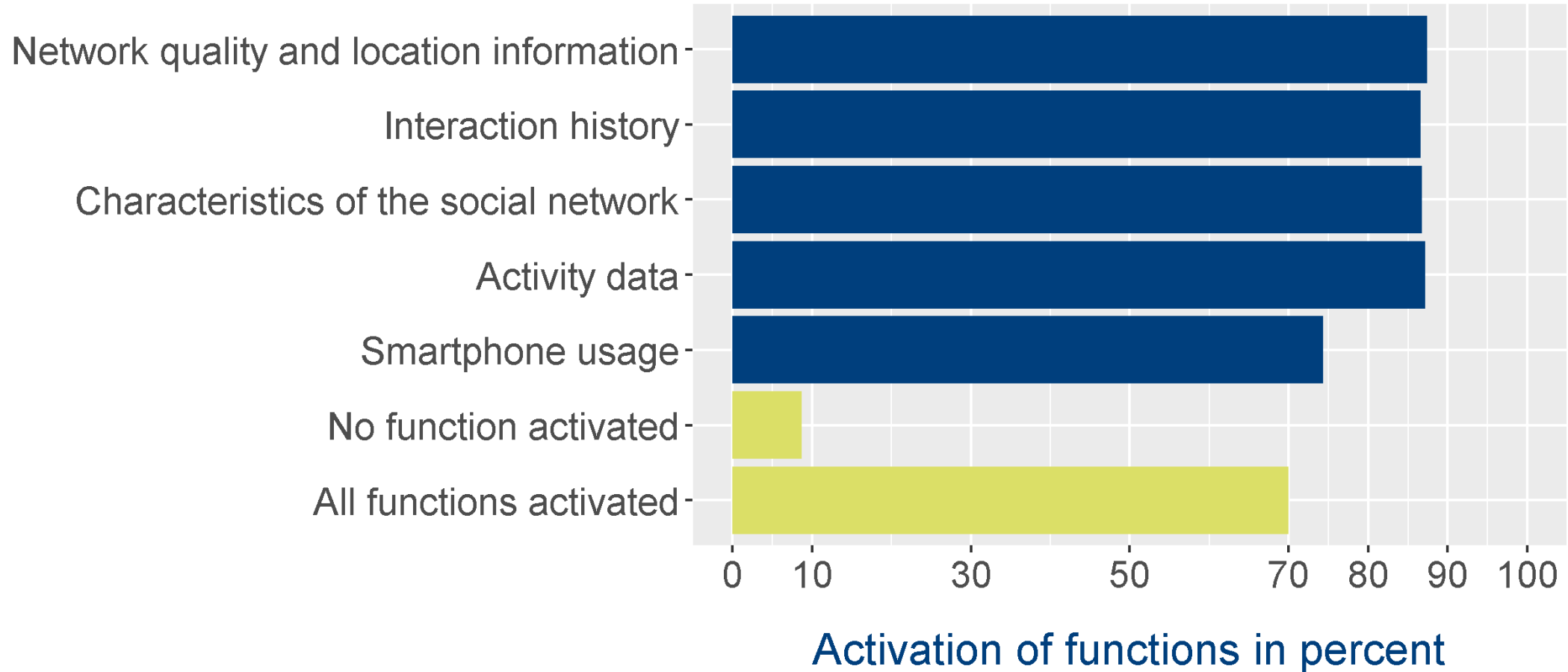# STANDARDIZED | NON-STANDARDIZED SOURCES

↓↓

Data for Social Sciences

## STRUCTURED DATA

Characteristics:

- e.g. survey data, administrative data
- high standardization
- homogenity of sources and formats
- standardized collecting, managing and analyzing
- standardized tools and methods
- sufficient computing and storage capacity
- no interconnected data infrastructure
- open legal and ethical issues

## UNSTRUCTURED DATA

Characteristics:

- e.g. text, video, audio
- low standardization
- heterogeneity of sources and formats
- no standardized collecting, managing and analyzing
- no standardized tools and methods
- no sufficient computing and storage capacity
- no developed data infrastructure
- open legal and ethical issues

# Lessons Learned

**3. Data Science is a "Team Sport" ... and needs to be treated as such.**

**METHODOLOGIST**

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 in survey research: AAPOR Task Force report. Japec, L.; Kreuter, F. ; Berg, M. ; Biemer, P. ; Decker, P.; Lampe, C.; Lane, J. ; O'Neil, C. Usher, A.. DOI: 10.1093/poq/nfv039. URL: http://poq.oxfordjournals.org/content/79/4/839.

# Work in Teams

# Lessons Learned

1. Creating new measures out of linked data is not enough. The effort has to be tied to a product.
2. Research questions need to guide decisions on measurements and data sources.
3. Data Science is a "Team Sport" ... and needs to be treated as such.