# What are Gathered Data? By Trent D. Buskirk

# What are Gathered Data?

- **Gathered data:**

  - Data obtained from an existing source most likely not designed for research purposes or distribution as such.

  - Data generated from some external procedure, tool, meter or sensor

  - Often referred to as "organic data" or "Big Data"

  - Social Media Data (e.g. Twitter); Sensor Data (e.g. Bluetooth wearable); Web-scraped or Open Data (e.g. Zillow.com)

# Big Data as a Process

| Big Data Generation | Big Data Management | Big Data Analytics |
|---|---|---|
| **Data Origin:**<br>• IoT Sensors<br>• Social Media<br>• Administrative<br>• Open Data<br>• Other | **Data Acquisition:**<br>• Data Collection<br>• Recording | **Modeling:**<br>• Feature Engineering<br>• Feature Reduction<br>• Model Construction<br>• Cross Validation |
| **Data Structure:**<br>• Structured<br>• Unstructured<br>• Semi-Structured | **Data Processing:**<br>• Extraction<br>• Cleaning<br>• Transformation<br>• Annotation<br>• Curation | **Analysis:**<br>• Visualization<br>• Hypothesis Testing<br>• Feature Importance<br>• Evaluation |
| **Data Attributes:**<br>• Volume<br>• Velocity<br>• Variety<br>• Veracity<br>• And others<br><br>**Data Elements:**<br>• Metadata<br>• User Tags<br>• Log files<br>• Other | **Data Access & Storage:**<br>• Integration<br>• Record Linkage<br>• Aggregation<br>• Representation<br>• Load<br>• Privacy | **Interpretation:**<br>• Inference<br>• Insight Extraction |

# Gathered Data… or Not? (1)

- **Public Use Survey Data Files**

  - These records are publicly available and often housed and accessed online

    - For example, Current Population monthly supplement survey data can be accessed and downloaded from census.gov

    - These data were designed for research or estimation purposes and while possibly a secondary data source for your research these data were designed and specific variables were intentionally included and processed for release in the public data files.

    - These data would not fall under the Gathered Data category.

# Gathered Data… or Not? (2)

- **Administrative Records**

  - While data designs can leverage information housed within administrative records, more often than not these types of data files, while structured, are broadly considered a type of "Big Data."

    - Housing Sales data for a given county are often available from a county assessor's office either in tabular form that can be scraped or as a tabular file that can be downloaded or through an Application Program Interface.

    - While these data are structured they were likely not created or generated for research purposes *a priori*.

    - There is considerable variety in not only the types of administrative records that are available but also in how they can be accessed.

# Gathered Data… or Not? (3)

- Application Programming Interfaces (APIs) allow users to more easily access web data stored, curated or created from online venues.

- Generally speaking data accessed through an API call would be considered gathered data

  - Unless the data source itself was presenting designed data (such as a call to generate tables created from the current population survey's online survey tabulation tool).

- APIs and Web-scraping from the same website may not produce equivalent data and we explore this in an upcoming lecture.

# Gathered Data: Summary

- Gathered data

  - Generated for a purpose other than analysis

  - Sourced from sensors, meters, web pages, social media, transactions, digital footprints and the like.

  - Data are opportunistically used for research or estimation purposes and are likely not the intended purpose of their generation.

  - While gathered data are not designed for research they can present viable auxiliary sources of information for investigating issues of interest to researchers and practitioners

# What's Next?

- In the next lecture, you will have a chance to **learn how to scrape data from a simple table presented online.**

# Readings

- [Ward and Barker (2013) on Big Data Definitions](#)

**M UNIVERSITY OF MICHIGAN**

© Faculty Presenter