# Demystifying Digital Scholarship

## An Introduction to Data Wrangling



Sherman Centre for Digital Scholarship

12-February, 2016

# Outline and Objectives

**Topics**

Let's talk about data

What is data 'wrangling'?
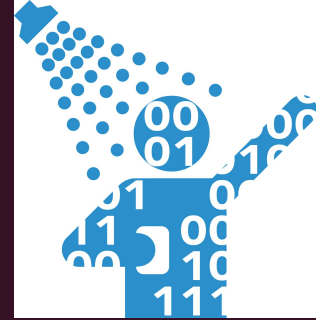
Why you should 'wrangle'?

How to 'wrangle'?

**Objectives**

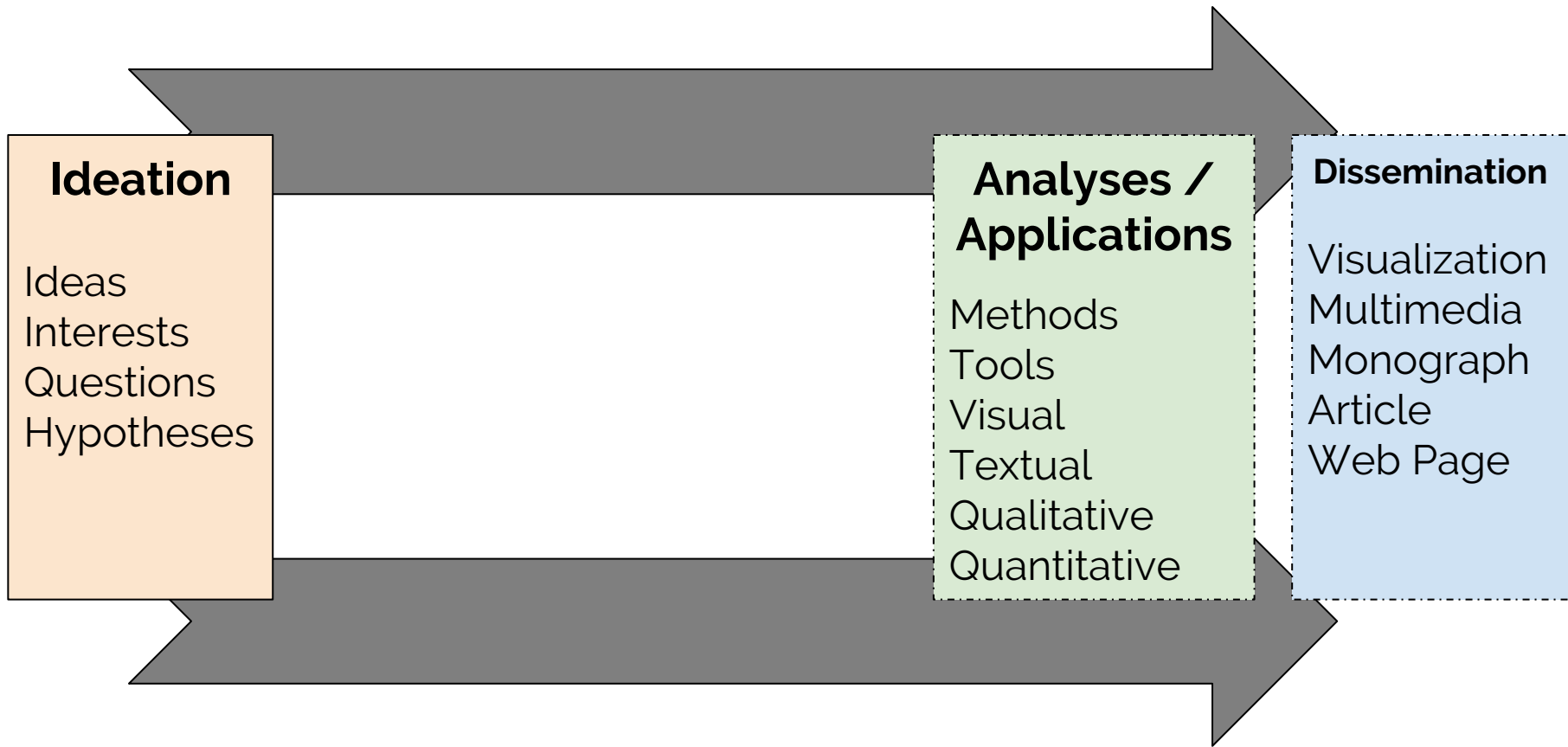Present and discuss strategies for structuring data

# Data & Data Wrangling

# Every Project Has Data!

➢ Chances are that your project contains at least one (and likely more) data types:
  → Text, images, tags, geographical coordinates, categorical items, records, metadata, multimedia, etc.

➢ Understanding your data and your intended actions is a critical part of developing a DS/DH project
  → It guides your data activities
  → It helps inform you of the ways in which your data can be used -- by you, your collaborators and others in your research community

**Ideation**

Ideas
Interests
Questions
Hypotheses

**Analyses /
Applications**

Methods
Tools
Visual
Textual
Qualitative
Quantitative

**Dissemination**

Visualization
Multimedia
Monograph
Article
Web Page

**Ideation**

Ideas
Interests
Questions
Hypotheses

**Data**

**Analyses / Applications**

Methods
Tools
Visual
Textual
Qualitative
Quantitative

**Dissemination**

Visualization
Multimedia
Monograph
Article
Web Page

**Ideation**

Ideas
Interests
Questions
Hypotheses

**Data Preparation**

Getting data

**Analyses / Applications**

Methods
Tools
Visual
Textual
Qualitative
Quantitative

**Dissemination**

Visualization
Multimedia
Monograph
Article
Web Page

# Ideation

Ideas
Interests
Questions
Hypotheses

## Data Preparation

Getting data

Assessing it

Cleaning it
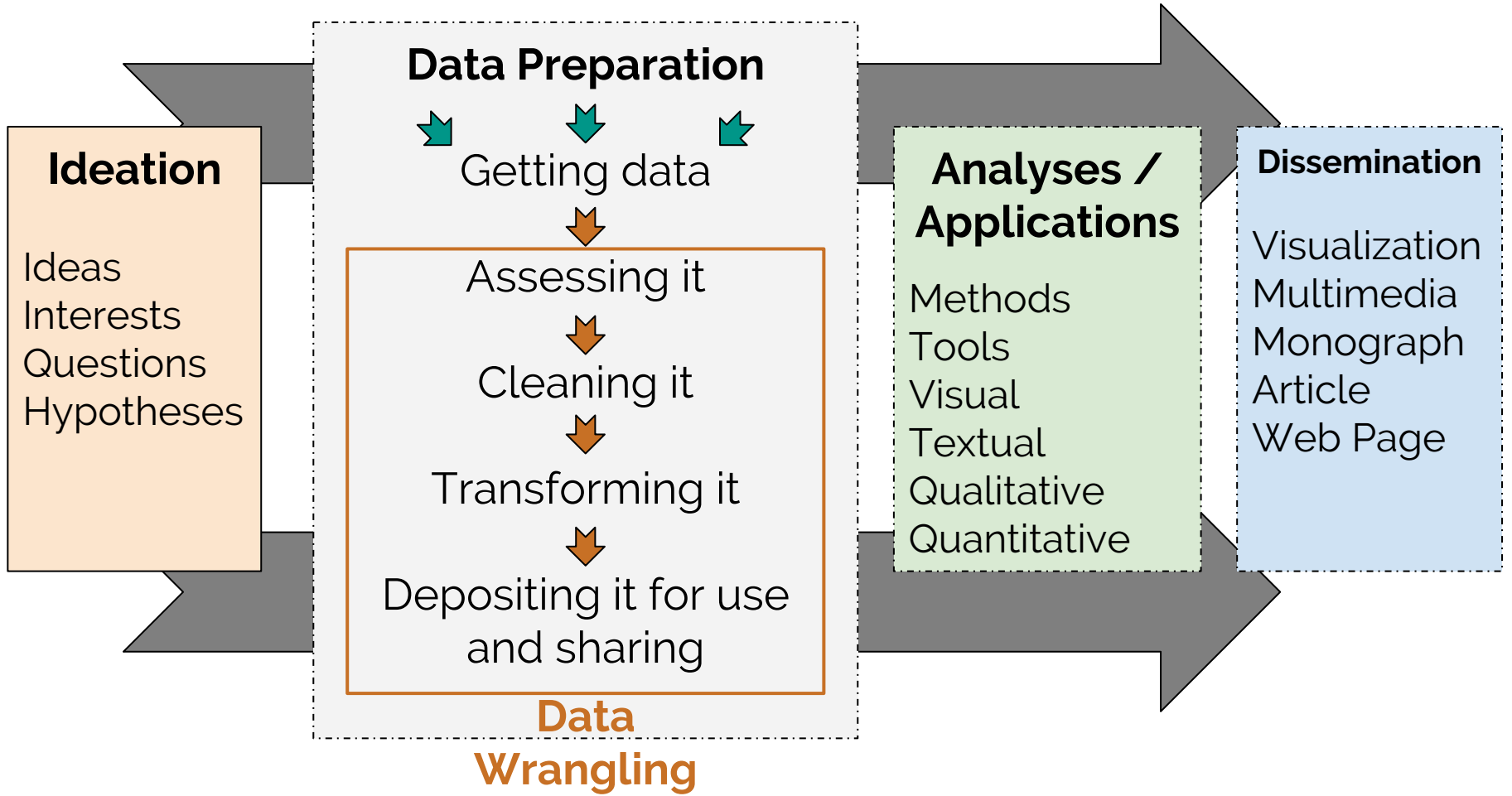
Transforming it

Depositing it for use and sharing

## Data Wrangling

## Analyses / Applications

Methods
Tools
Visual
Textual
Qualitative
Quantitative

## Dissemination

Visualization
Multimedia
Monograph
Article
Web Page

# Defining data wrangling

**wran·gle**  *v.tr*

➢ To manage or herd

➢ To manage or control

➢ To grasp and maneuver (something); wrestle

➢ To win or obtain by argument

# Defining data wrangling

**wran·gle**  *v.tr*

➢ To manage or herd

➢ To manage or control

➢ To grasp and maneuver (something); wrestle

➢ To win or obtain by argument

**wran·gle**  *v.intr*

➢ To attempt to deal with or understand something; contend or struggle

# Defining data wrangling

**wran·gle**  *v.tr*

➢ To manage or herd

➢ To manage or control

➢ To grasp and maneuver (something); wrestle

➢ To win or obtain by argument

**wran·gle**  *v.intr*

➢ To attempt to deal with or understand something; contend or struggle

**wran·gle**  *n.*

➢ An angry, noisy argument or dispute.

# Defining data wrangling

**In the context of data, DATA WRANGLING:**

➢ Is the process of cleaning and conditioning data into a usable format

➢ May be a manual, semi-automated or automated process

➢ Produces data that connects to tools, collaborators and communities

# Why wrangle?

➤ Even if you understand and work well with your data, it doesn't mean that a computer will be able to use it to the same extent.

➤ Computers (like people) are only as flexible/adaptable as far as they have been trained or instructed.

➤ Therefore, it often takes work to structure your information/data in a way that can be used in a computing environment.

# Why wrangle?

Because this happens →

| Data |
|------|
| Your data |
| Your  data |
| Your data |
| Your |
| Data |

, and this can be a problem

From School of Data's Data Cleaning Module:

"the Invisible Man is in your spreadsheet, messing with your data"

http://schoolofdata.org/handbook/courses/data-cleaning-invisible-man-in-spreadsheets/

# Activity: Why wrangle?

All workshop materials are available in the following Google Drive folder:

https://goo.gl/u53tkz

Check out **Sample Copy - Crowdsourcing Exercise**

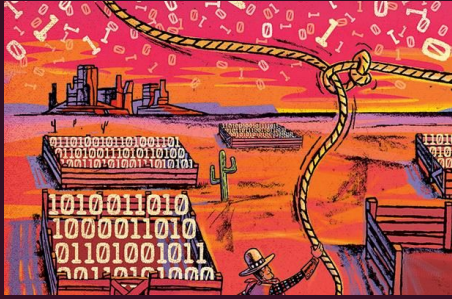('pop out' into a new window; Make a copy and play with it, if you want)

➢ What's not quite right?

➢ What could possibly go wrong?

➢ What needs to be cleaned?

➢ How best to clean this data?
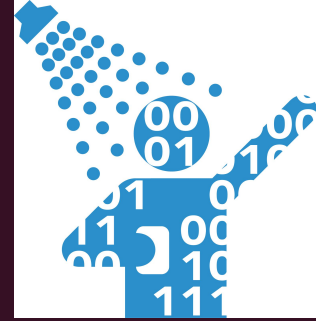    What if it had 30000 rows? What if it amassed 30000 rows / day?

Secret Location
Eva's Map

# Understanding your needs
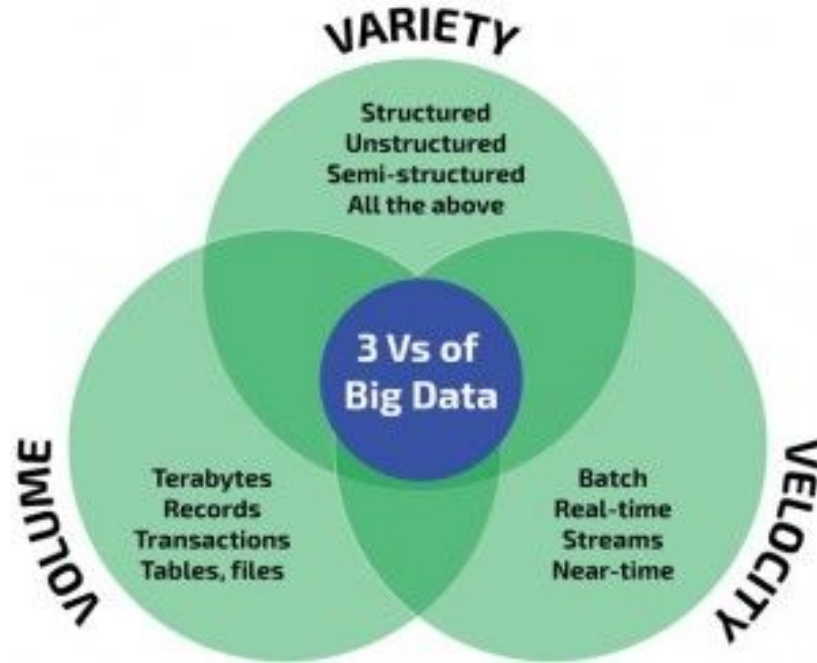
What is your 'data'?

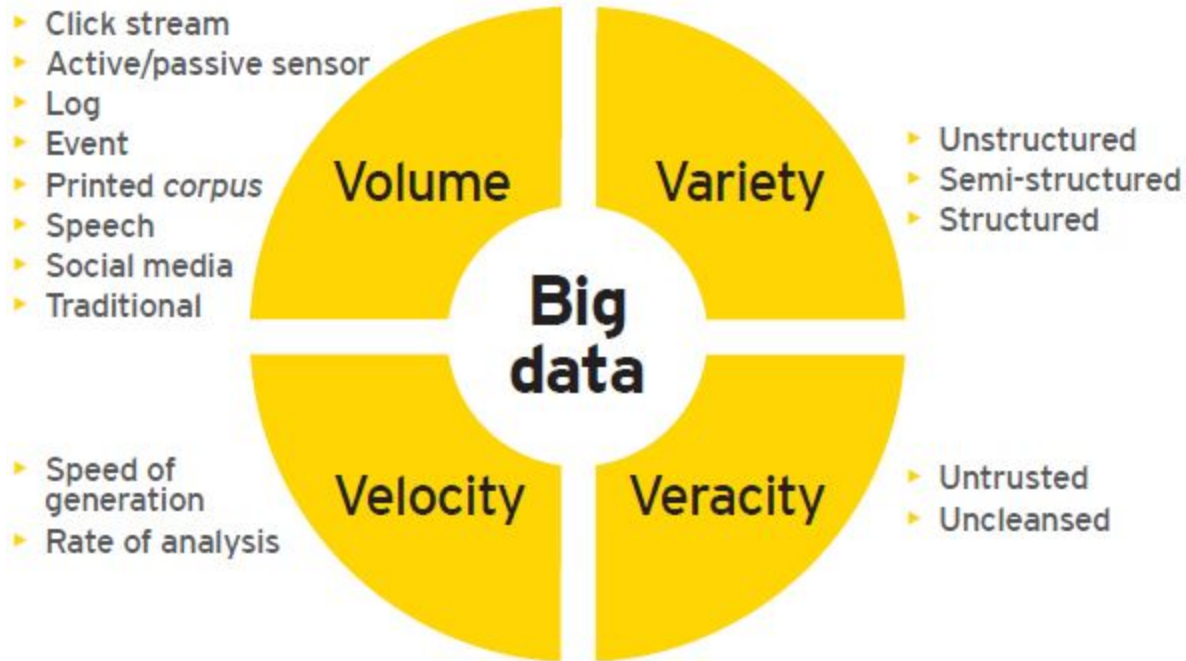What kind of 'wrangling' is required?

# Data wrangling: BIG impact, even for small data

Whether your data is small or "BIG", conditioning it is critical.

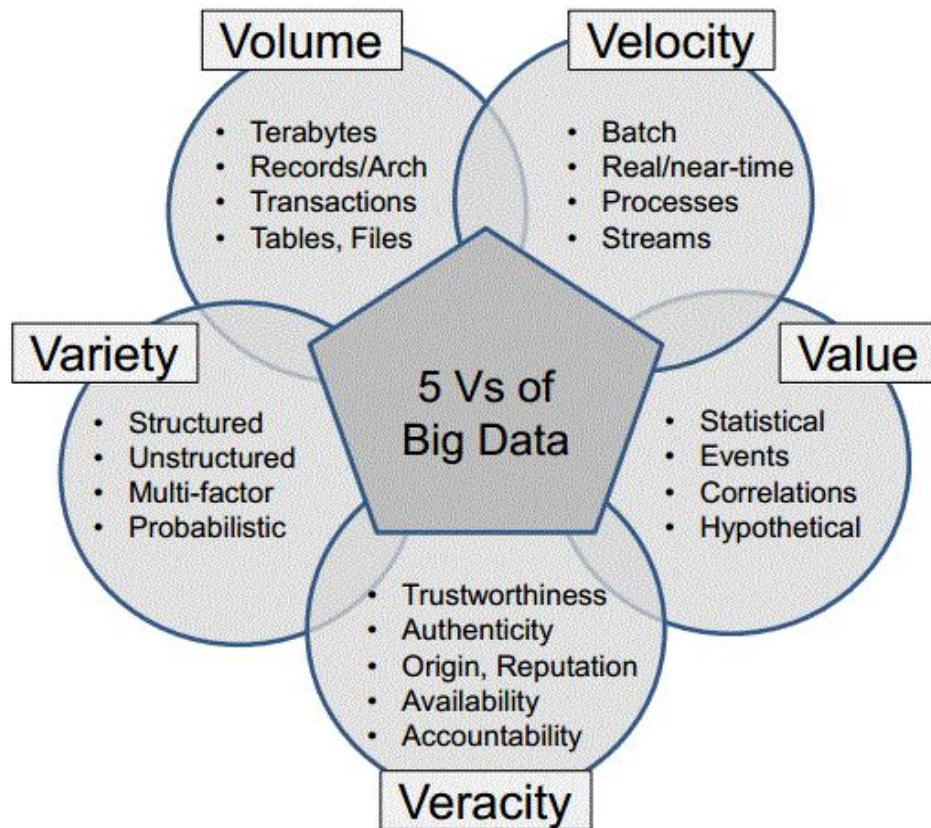# Big data is defined by the 3 "V"s

# ...or the 4 "V"s



- Click stream
- Active/passive sensor
- Log
- Event
- Printed *corpus*
- Speech
- Social media
- Traditional

**Volume**

**Variety**

- Unstructured
- Semi-structured
- Structured

**Big data**

- Speed of generation
- Rate of analysis

**Velocity**

**Veracity**

- Untrusted
- Uncleansed

# ...or 5?



**Volume**
- Terabytes
- Records/Arch
- Transactions
- Tables, Files

**Velocity**
- Batch
- Real/near-time
- Processes
- Streams

**Variety**
- Structured
- Unstructured
- Multi-factor
- Probabilistic

**Value**
- Statistical
- Events
- Correlations
- Hypothetical

**Veracity**
- Trustworthiness
- Authenticity
- Origin, Reputation
- Availability
- Accountability

5 Vs of Big Data

# 6? Really?



**Volume**
- Multi-domain data
- User/device data
- Geolocation data

**Value**
- Spectrum modeling
- Spectrum prediction
- Spectrum management

**Velocity**
- Data in motion
- Stream computing
- Batch algorithms
- Real-time algorithms

**Big Spectrum Data**

**Variety**
- Crowd sensing
- Geolocation database
- Heterogeneous sensors
- Different data types

**Viability**
- Variable selection
- Variable relevance
- Variable relationship

**Veracity**
- Spectrum data quality
- Data uncertainty
- Data security

# Factors Influencing Your Wrangling Needs

**Your data might...**

➢ be voluminous and dynamic

➢ come from very diverse sources in a variety of formats and structures

➢ vary considerably in quality and consistency (i.e. may require heavy scrutiny)
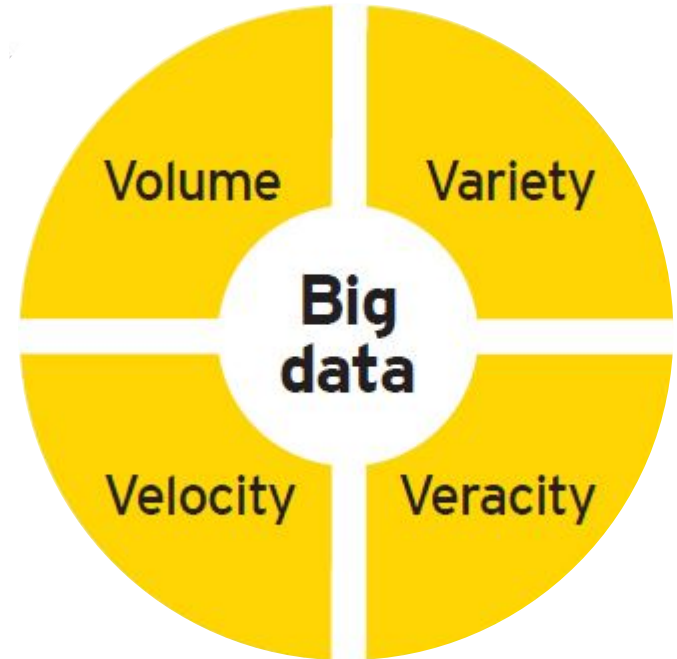
**Or it might not be...**
**All of these factors influence your data wrangling needs and potential solutions!**

# What else matters?

**Other factors that will affect the the nature of your wrangling activities:**

➢ What tools are you using?

➢ Are automated approaches available?

➢ What formats are required for later analyses?

➢ How tolerant are you to errors?

# Getting & Assessing Your Data

➢ Downloading Files

➢ API access

➢ Web-scraping

# Bulk Downloads

**Common data formats**

➢ ASCII text

➢ Delimited ASCII (.csv, .tsv)

➢ PDF

➢ JSON (https://goo.gl/632RGb)

➢ Markup languages (TEI, HTML, XML)

➢ Multimedia (.wav, .ogg, .mp4, .mkv)

➢ Other, program-specific and ad-hoc file formats (fixed width, SAS, xls)

```
{
  "name": "ballparks",
  "type": "FeatureCollection",
  "features": [{
    "type": "Feature",
    "geometry": {
      "type": "Point",
      "coordinates": [-112.066564, 33.445081]
    },
    "properties": {
      "Class": "Majors",
      "League": "Major League Baseball",
      "Team": "Arizona Diamondbacks",
      "Ballpark": "Chase Field",
      "Lat": "33.445081",
      "Long": "-112.066564"
    }
  },
```

```
<typeOfResource>cartographic</typeOfResource>
<genre authority="lctgm">Aerial photographs</genre>
<originInfo>
<publisher>Air Photo Division, Energy Mines + Resources</publisher>
<place>
<placeTerm type="text">[Place of publication unknown]</placeTerm>
</place>
<dateCreated>1966</dateCreated>
<dateOther>1966</dateOther>
</originInfo>
<language>
<languageTerm type="code" authority="iso639-2b">eng</languageTerm>
</language>
<physicalDescription>
<extent>[1:37,000 approximately]</extent>
```

# API Access

API = Application program interface
➢ set of protocols/tools for building software applications
➢ governs how software should interact with each other and user interfaces

Reddit API: https://www.reddit.com/dev/api
New York Times API: http://developer.nytimes.com/docs

# Web Scraping

Scraping vs. Parsing:

➢ Parsing: data being extracted is intended as input to another program

➢ Scraping: data being extracted is intended for display to an end user

e.g. > `wget https://www.reddit.com/r/sandersforpresident`

➢ Scrapes web page html to file →
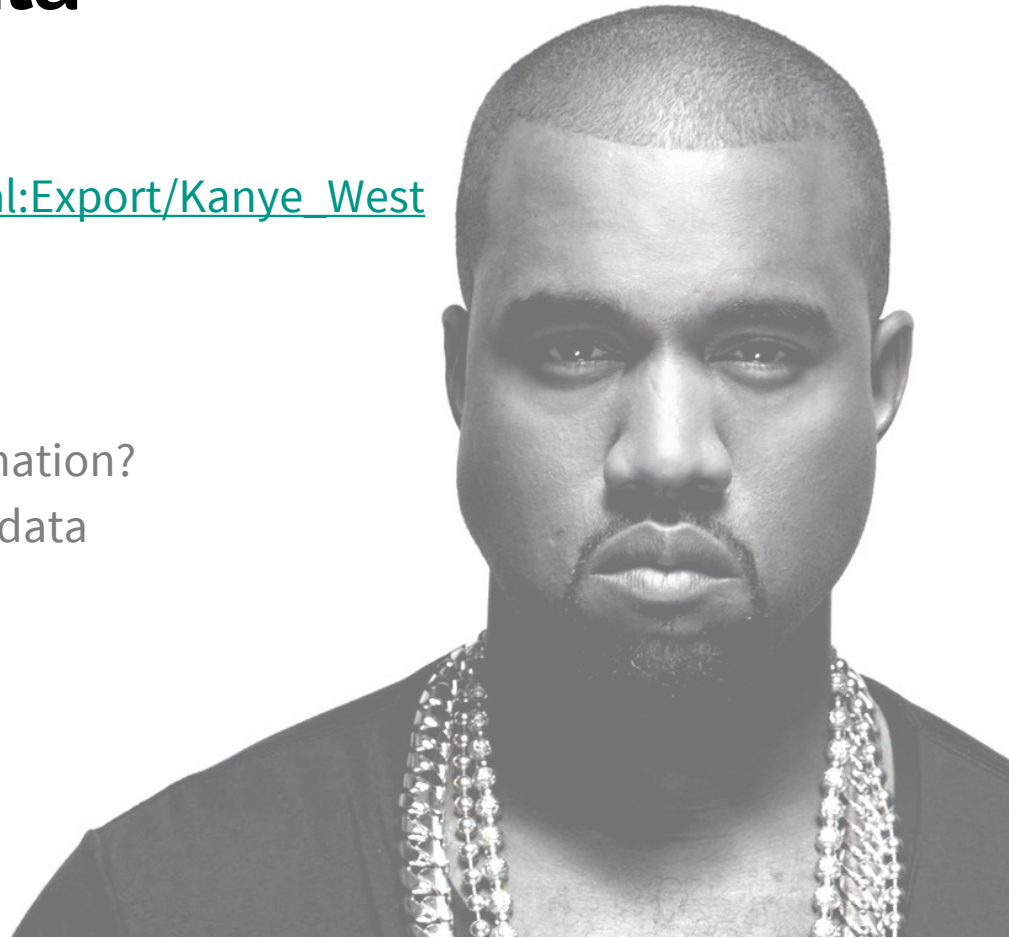
# Activity: Assessing Data

Wikipedia page exported to XML:

➢ https://en.wikipedia.org/wiki/Special:Export/Kanye_West

➢ https://goo.gl/pCDg30

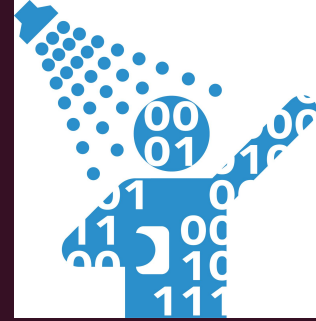Goal: Explore which wikipedia articles

are linked to Kanye?

➢ Can you extract the necessary information?

➢ Are there patterns/landmarks in the data
    that you can take advantage of?

➢ What is the consistency of the
    information?

# Cleaning & Transforming Your Data



| | A | B | C |
|---|---|---|---|
| 1 | Data | Results | Formula |
| 2 | drapes | 1 | =countif(A2:A6, "drapes") |
| 3 | grapes | 1 | =countif(A2:A6, A2) |
| 4 | grapeshot | 2 | =countif(A2:A6, "?rapes") |
| 5 | grapefruit | 3 | =countif(A2:A6, "?rapes*") |
| 6 | grapevine | 4 | =countif(A2:A6,"grape*") |
| 7 | 100 | 1 | =countif(A7:A10, "100") |
| 8 | 1,000 | 1 | =countif(A7:A10, A7) |
| 9 | 10,000 | 2 | =countif(A7:A10, "<=1000") |
| 10 | 100,000 | 3 | =countif(B7:B10, "<>"&C12) |
| 11 | | 4 | =countif(B7:B10, "<="&D12) |
| 12 | More Data: | 1,000 | 100,000 |

# Data Cleaning: Goal & Activities

Goal: Create data sets that are **consistent** and **interoperable** with other data of interest

Data cleaning (scrubbing) may include:
➢ Detecting and remediating corrupt/inaccurate records
➢ Removing typographical errors
➢ Validating against a known list of possibilities (e.g. verify a string as a postal code)
➢ Eliminating duplicate entries
➢ Harmonizing and standardizing (e.g. represent 'St', 'St.', 'Street' as 'street')

Cleaning may be carried out
➢ Manually (interactive)
➢ Semi-automatically (guided)
➢ Automatically (scripted)

# Approaches For Cleaning & Transformation

- ➤ Spreadsheets are often a good place to start
  - → http://schoolofdata.org/courses/#IntroDataCleaning
- ➤ Text Editors - Find and replace
- ➤ Command line - regular expressions (http://pythex.org/; http://www.regular-expressions.info/)
- ➤ Customized applications
  - → Stanford data wrangler: http://vis.stanford.edu/wrangler/
  - → Trifacta: https://www.trifacta.com/
  - → Open (Google) Refine: http://openrefine.org/
- ➤ Fully-automated scripts

ACCOU NTA NT,
AUDITOR AND ADJUSTER.

CITY OF HAMILTON

q6 Jas Ramage, laborer
96 Sand Garrity
98 Geo 1-ake, granite cutter
too Alex Martin, laborer
tog Godfrey Auld, laborer
to6 Jas Flynn, laborer
110 Jas Furlong, laborer
112 11 A Broadbent, brakinn
ta Henry N'ernon, directory
pulilisher
16 Matthew Hunter, carprar
118 Herbert Dixon

Yoik 11 intersect(

36 A H McKeown, tnerchant
40 Thos 11 Jermvn. bookkpr
42 Jas Mclaugrain, machist

Main st Intersects
68 Mrs Mary Liven

# So, when to let the computer take over?

Whenever it works and will save you time!



**Time / Effort to Complete** (vertical axis, left)

**Potential for Calamity** (vertical axis, right)

a file;
page of
text, etc.

many, diverse
datasets;
corpora

**# of files/items to wrangle**

# Spreadsheets: The frenemy of research



**Doing it manually**

Time / Effort to Complete

Potential for Calamity

a file; page of text, etc.

many, diverse datasets; corpora

**# of files/items to wrangle**

# Spreadsheets: The frenemy of research



Time / Effort to Complete

# of files/items to wrangle

Potential for Calamity

a file; page of text, etc.

many, diverse datasets; corpora

**Doing it manually**

**Doing it automatically**

**Also see this: https://xkcd. com/1205/ (thanks Chris!)**

Russia Submits Syria Cease-Fire Plan to U.S.

NATO Deploys Ships to Help in Migrant Crisis

EUROPE

# Flaws Found in Study Favored by Backers of Austerity

By **MATTHEW DALTON** and **GEOFFREY T. SMITH**

April 18, 2013 4:29 p.m. ET

But a study published this week by three economists at the University of Massachusetts says the Reinhart-Rogoff findings contain basic errors, including one involving their spreadsheets that omitted five countries from the result.

Correcting for these errors largely causes the Reinhart-Rogoff finding to disappear, according to the new study by the economists Thomas Herndon, Michael Ash and Robert Pollin.

**EuSpRiG**
European Spreadsheet
Risks Interest Group

| | |
|---|---|
| Identifier: | FH1217 |
| Title: | 1,791 voters inflated to 4,870 |
| Source: | http://www.chieftain.com/news/local/ortiz-scales-down-inactive-voter-count/article_abcb94c0-06a4-11e1-8668-001cc4c03286.html |
| Organization: | Pueblo County |
| Region: | USA |
| Release Date: | 04 November 2011 |
| Risk: | False Elections |
| Tags: | Government |

Thursday, Ortiz ultimately agreed with Gessler that the state database showed 1,791 inactive ballots were returned in Pueblo County. Ortiz blamed simple counting errors in his office for his inaccurate number. The state computer data is cumbersome to read, so Ortiz's office converted its display to a different form of spreadsheet. That's where the counting error occurred, the clerk said. Anytime there is a dispute over tallying votes, it always raises suspicions about possible voter fraud.

| | |
|---|---|
| Identifier: | FH1210 |
| Title: | $1M went missing as staff managed "monstrous spreadsheets." |
| Source: | http://www.metrowestdailynews.com/top_stories/x1876834739/Framingham-discovers-1-5-million-error |
| Organization: | Town of Framingham (municipality) |
| Region: | USA |
| Release Date: | 15 October 2011 |
| Discrepancies: | 12,000,000 |
| Risk: | Money Loss |
| Tags: | Government |

The town mistakenly reckoned it had $1.5 million more in this year's budget than it actually has and must now use $600,000 in unexpected state aid to help fill that gap, officials said yesterday.

Chief Financial Officer Mary Ellen Kelley said she takes responsibility for the mistake, which she found Wednesday night on the debt services line item in the $208.6 million fiscal 2012 operating budget. She said a figure went missing as staff managed "monstrous spreadsheets."

"It's frustrating," Kelley said yesterday. "I hate when we make mistakes. People are human and they do make mistakes, but I hate it."

http://www.eusprig.org/horror-stories.htm

# Activity: Why wrangle?

All workshop materials are available in the following Google Drive folder: https://goo.gl/u53tkz

Open one of the folders in the Google Drive; explore the contents
➢ Scraped reddit subreddit pages in JSON and HTML
  → Google sheet, transformed from html page (import-reddit-xxxxxx) - created with www.import.io/
➢ Google sheet of structured twitter data (#xxxxxxxx) - created with https://tags.hawksey. info/

Devise a plan of how you might use these data to explore a research question.
➢ How will you analyze the data?
➢ In what ways will you need to clean | transform the data?
➢ What tools might work for this?

**Trump:**

#Trump2016

https://www.reddit.com/r/the_donald

**Bernie:**

#FeelTheBern

https://www.reddit.com/r/sandersforpresident

**Hillary:**

#ImWithHer

https://www.reddit.com/r/hillaryclinton/

# Depositing for Use & Sharing

# Depositing for Use & Sharing

Your data can (and should!?) outlive your research project
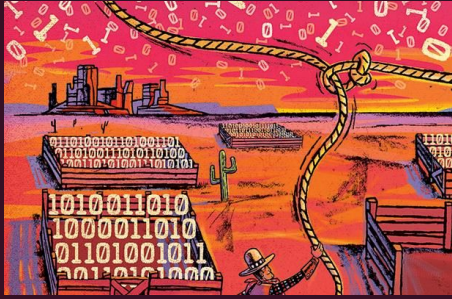
Consider how you can make your data more:
➢ adaptable, flexible, extensible to other uses
➢ secure and sustainable
➢ collaborative / social
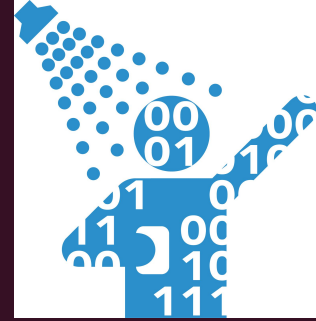
# Depositing for Use & Sharing

How might you need to package/adjust your data in order to:
➢    be better able to show (demonstrate) it to people?
➢    make it work with different tools?
  →    Keep data is translatable formats (or formats that work well with visualization)
  →    Explore formatting requirements/compatibility for other community tools

➢    have it be discoverable and understandable to others?
➢    maximize its lifespan?
  →    Document it!
  →    Organize it.
  →     Archive it to an appropriate repository (IR, data repository)

# Final Thoughts and Strategies

# Experimentation and Documentation

Often, experimentation and iteration are important in establishing the best way to get your data into 'shape'.

Starting with a sample of data is a good approach.

It's important to document your outcomes **and** your process

# How to save time in the long run...

➢ Look for tools that exist to help you wrangle your data
  → automated or semi-automated (guided) cleaning
  → data transformation
  → converting between data formats

➢ Seek out tutorials / instruction for the tools you're using

➢ Control your data at the point of collection - refine your process to reduce 'garbage in'
  → e.g. when crowdsourcing data -- use controlled fields and vocabularies; insert data validation processes