# DATA COLLECTION

Slides by WESLEY WILLETT

VISUAL ANALYTICS

# WHERE DOES DATA COME FROM?

We tend to think of data as a thing...
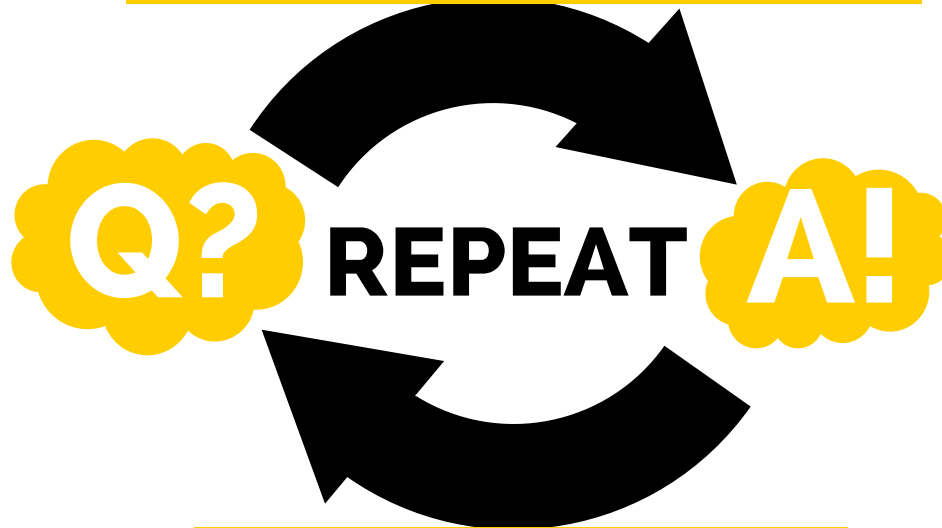
in a database...

somewhere...

# WHY DO YOU NEED DATA?

(HINT: Usually, because you have a question you need to answer!)

DATA ANSWERS

# ANALYSIS IS A CYCLE

GATHERING DATA,
APPLYING STATISTICAL TOOLS,
AND CONSTRUCTING GRAPHICS
TO ADDRESS QUESTIONS

Q?  REPEAT  A!

INSPECT "ANSWERS" AND
ASSESS NEW QUESTIONS

(...BUT OFTEN YOU <u>START</u> WITH A QUESTION AND NEED TO COLLECT  DATA TO FIT IT)

# CHOOSING A QUESTION

"How has language evolved over time?"

"What will the weather be like next month?"

"Are the right people seeing my advertisements?"

"What is the current temperature?"

# A PROBLEM OF SCALE

CHALLENGING
TO FIND DATA

"How has language evolved over time?"

"What will the weather be like next month?"

"Are the right people seeing my advertisements?"

"What is the current temperature?"

NOT AS
INTERESTING

# HOW TO OBTAIN DATA?

## COLLECT IT
– OBSERVATION
– SURVEYS
– LOGGING
– SENSORS
– CROWDSOURCING

## FIND OR EXTRACT IT
– OPEN CORPUSES
– DATA RETAILERS
– APIS
– SCRAPING THE WEB

## GENERATE IT
– SIMULATIONS

ALL OF THESE HAVE
PROS/CONS

# THIS LIST IS NOT EXHAUSTIVE

This lecture is intended to expose you to just a few useful data sources and collection methods.

# COLLECTING DATA

**Choosing the best way to capture information you need.**

# SURVEYS

Paper surveys / In person interviews

STILL ONE OF THE BEST WAYS TO GET DETAILED DATA OR DATA ABOUT SENSITIVE SUBJECTS

# SURVEYS ONLINE

# CROWDSOURCING DATA COLLECTION

# WEB LOGGING

Tracking Visits, Click-Throughs, and Traffic Patterns and other measures of User Activity.

- Google Analytics
- Open Web Analytics
- and many others...

# EDITS & ACCESSS LOGS ON

WIKIPEDIA

# SENSORS

- Weather stations
- Personal activity trackers
- Cameras
- Mobile phones

# HOW TO OBTAIN DATA?

## COLLECT IT
- OBSERVATION
- SURVEYS
- LOGGING
- SENSORS
- CROWDSOURCING

## FIND OR EXTRACT IT
- OPEN CORPUSES
- DATA RETAILERS
- APIS
- SCRAPING THE WEB

## GENERATE IT
- SIMULATIONS

# GENERATING DATA

**SIMULATIONS**

http://www.nasa.gov/content/a-portrait-of-global-winds/

SHARE

# Is It Better to Rent or Buy?

By **MIKE BOSTOCK, SHAN CARTER** and **ARCHIE TSE**

The choice between buying a home and renting one is among the biggest financial decisions that many adults make. But the costs of buying are more varied and complicated than for renting, making it hard to tell which is a better deal. To help you answer this question, our calculator takes the most important costs associated with buying a house and computes the equivalent monthly rent. **RELATED ARTICLE**

## Home Price

A very important factor, but not

EQUIVALENT RENT

– $8K

**If you can rent a similar home for less than ...**

# HOW TO OBTAIN DATA?

## COLLECT IT
- OBSERVATION
- SURVEYS
- LOGGING
- SENSORS
- CROWDSOURCING

## FIND OR EXTRACT IT
- OPEN CORPUSES
- DATA RETAILERS
- APIS
- SCRAPING THE WEB

## GENERATE IT
- SIMULATIONS

# FINDING AND EXTRACTING EXISTING DATA

## LARGE OPEN CORPUSES

# DBPEDIA

# About: Iceland

An Entity of Type : place, from Named Graph : http://dbpedia.org, within Data Space : dbpedia.org

Iceland /ˈaɪslənd/ (Icelandic: Ísland [ˈistlant]), sometimes referred to in full as the Republic of Iceland (Lýðveldið Ísland), is a Nordic island country marking the juncture between the North Atlantic and the Arctic Ocean, on the Mid-Atlantic Ridge. The country has a population of 325,671 and a total area of 103,000 km2 (40,000 sq mi), which makes it the most sparsely populated country in Europe.

| Property | Value |
|----------|-------|

# QUERYING DBPEDIA

D Virtuoso SPARQL Query Edi ×

dbpedia.org/sparql

## Virtuoso SPARQL Query Editor

About | Namespace Prefixes | Inference rules | iSPARQL

### Default Data Set Name (Graph IRI)

http://dbpedia.org

### Query Text

```
select distinct ?Concept where {[] a ?Concept} LIMIT 100
```

# WIKIDATA

WIKIDATA

Main Page | Discussion

Main page
Community portal
Project chat
Create a new item
Item by title
Recent changes
Random item
Help
Donate

Print/export
Create a book
Download as PDF
Printable version

Tools
What links here
Related changes
Special pages
Permanent link
Page information

open
collaborative
multilingual
free

## Welcome to Wikidata

the free knowledge base with 14,820,745 data items that anyone

Introduction • Project Chat • Community Portal • Help

**Welcome!**

Wikidata is a free linked database that can be read and edited by both humans and machines.

Wikidata acts as central storage for the **structured data** of its Wikimedia sister projects including Wikipedia

**Learn about data**

New to the wonderful world o
up to speed and feeling com

# PROJECT GUTENBERG

www.gutenberg.org/ebooks/search/?sort_order=downloads

en|fr

Project Gutenberg offers 46,845 free ebooks to download.

Donate

Flattr this!

**Search**   **Latest**   **Terms of Use**   **Bookmarks**   **Donate?**   **Mobile**

Search Project Gutenberg. <s:   Help

# All Books (sorted by popularity)

A   **Sort Alphabetically**

   **Sort by Release Date**

**The Kama Sutra of Vatsyayana**
Vatsyayana
13285 downloads

# GOOGLE N-GRAMS

Google Ngram Viewer  ×

https://books.google.com/ngrams/graph?content=Chocolate%2CVanilla%2CStrawber...

## Google books Ngram Viewer

Graph these comma-separated phrases: Chocolate,Vanilla,Strawberry   ☐ case-insensitive

between 1800 and 2000 from the corpus English with smoothing of 3 . **Search lots of books**

0.000350%

0.000300%

0.000250%

0.000200%                                                                              Chocolate

# FINDING AND EXTRACTING EXISTING DATA

## GOVERNMENT AND INTERNATIONAL DATA INITIATIVES

# GOVERNMENT INITIATIVES

## WWW.DATA.GOV (US)

## DATA.GOV.UK

## DATA.GOV.BE



**Data.gov**

https://www.data.gov

🇺🇸 DATA.GOV   DATA   TOPIC

The home of the data

Here you will find data, tools web and mobile applications.

SEAR

Health Care Provider Charge

BROWSE TOPICS

---

**Data Search | data.gov.uk**

data.gov.uk/data/search

DATA.GOV.UK Beta
Opening up Government

🏠 / Datasets

Search for data...

📍 or conduct map based search

Show Search Facets »

## 19422 Results

### Live traffic information from the Highway

Highways Agency

Live traffic information data showing traffic information on the road network in England, maintained by the Highways Agency. August 2013 Following a change of...

### Learning Aim Reference Service

Skills Funding Agency

Learning Aim Reference Service (LARS) service will offer a 'Quic facility, allowing users to search by most commonly used fields full set of search fields will still...

---

**| data.gov.be**

data.gov.be/fr/datasets

nl  fr  de  en

Beta

Data.gov.be

HOME  CONDITIONS D'UTILISATION  **DATA**  APPS  IDÉES  FORUM

Liste de sets de données disponibles comme "open data".

| Catégorie | Type | Granularité | |
|---|---|---|---|
| – Tout – | – Tout – | – Tout – | Appliquer |

| Titre | Catégorie | Type |
|---|---|---|
| Zones de stationnement voirie 2013 | Mobilité | Télécharg |
| Usages TIC des ménages wallons | TIC | Télécharg Service w |
| Usages TIC des citoyens wallons | Population, Economie, TIC | Télécharg Service w |
| UDP Mars 2013 par commune | Energie, Pouvoirs publics | Télécharg |
| UDP Mai 2013 par commune | Energie, Pouvoirs publics | Télécharg |

# NEW DATA INITIATIVES JUST TO TRACK ALL THE DATA INITIATIVES

Browse by country | re3dat ✕

www.re3data.org/browse/by-country/

## re3data.org
### REGISTRY OF RESEARCH DATA REPOSITORIES

Home      Search      **Browse**      Suggest      FAQ

About      Schema      Contact      Imprint

# INITIATIVES IN FRANCE

## HTTP://DATA.GOUV.FR



## HTTP://OPENDATA.PARIS.FR/EXPLORE/

# FINDING AND EXTRACTING EXISTING DATA

## OTHER PUBLIC DATA REPOSITORIES

# MORE REPOSITORIES

VISUALIZING.ORG

http://visualizing.org/data/browse

AMAZON PUBLIC DATA HOSTING

http://aws.amazon.com/publicdatasets/

GOOGLE PUBLIC DATA

http://www.google.com/publicdata/directory

# FINDING AND EXTRACTING EXISTING DATA

DATA RETAILERS

# DATA RETAILERS

AZURE DATA MARKET

https://datamarket.azure.com/browse/data

FACTUAL

http://www.factual.com/

AND AGAIN, THERE ARE MANY, MANY MORE...

# FINDING AND EXTRACTING EXISTING DATA

APIS

# TWITTER

Streaming APIs (live data by users and by topics)

The "Firehose" (all of live twitter)

Complete Archives via "Gnip" and eventually (maybe) the US Library of Congress

HTTPS://DEV.TWITTER.COM

# Tottenham Riots

402 sources sharing 551 tweets matching "tottenhamriots" or "tottenham"

## Search

(enter search terms here) [Search]

## Sort

# times retweeted

## Show Sources (showing 8 of 10 sources loaded)

All | Ordinary People | Journalists / Bloggers | Organizations | Uncategorized | Eyewitnesses

## Show Tweets

All | **Exclude RTs** | Images & Videos

---

👁 **Daniel Carr, @daniel_carr** (2 years, 3 months old)

Myself in 160 characters: Schizophrenic. Also a criminologist

**London, United Kingdom**

Journo/Blogger   **41** RTed   **56** Klout

NETWORK SKETCH
**213** Followers
**163** Following

FRIENDS' LOCATIONS
London, GB 34.48%
Glasgow, GB 5.17%
Manchester, GB 3.45%

TOP ENTITIES MENTIONED HISTORICALLY
Bruce Grove, Tottenham Hale, London, BBC, Haha,

**31 Tweets**

#tottenham #tottenhamriots Fire near Bruce Grove Station, larger one towards Lordship Lane
Aug. 6, 2011, 11:27 p.m.

#tottenham #tottenhamriots @MrsCheddies by Bruce Grove I mean north of previous fires, on High Rd towards Lordship Lane
Aug. 6, 2011, 11:24 p.m.

@hackneyhive yeah around that area there are 2 fires, one small now, one very large #tottenham #tottenhamriots

---

**Aidan Rowe, @Aidan_Rowe** (1 year, 2 months old)

Post-punk, proto-utopian, anarchist, activist, musician, blogger, student, failed comedian.
http://redwriters1.blogspot.com

Ordinary Person   **23** RTed   **49** Klout

NETWORK SKETCH
**215** Followers
**395** Following

FRIENDS' LOCATIONS
Dublin, IE 43.48%
London, GB 4.35%
Cork, IE 1.74%

TOP ENTITIES MENTIONED HISTORICALLY
Oslo, BBC, Dublin, Dermot Mulqueen, Johann Hari,

**5 Tweets**

"Why couldn't the people in #Tottenham just have held a nice dignified protest for us to ignore?" - Liberals #tottenhamriots
Aug. 7, 2011, 12:49 a.m.

Any reports of arrests? #tottenham #tottenhamriots Hope everyone is safe. #acab

# SRSR
[DIAKOPOULOS ET AL. 2012]

GOOGLE EARTH ENGINE

HTTPS://EARTHENGINE.GOOGLE.ORG/

1984

2012

# MORE APIS
## (APPLICATION PROGRAMMING INTERFACES)

### NEW YORK TIMES APIS

http://developer.nytimes.com/

 (Archival news articles from 1851, books, movies,
  geographical, and political data)

### ECHONEST APIS

http://developer.echonest.com/

(Incredibly detailed Song, Album, Artist data for millions of musicians)

### OPEN STREET MAP

http://wiki.openstreetmap.org/wiki/API

(Detailed location and map data for the whole world)

# AND THE LIST GOES ON!

Medical, Education, Health ×

www.programmableweb.com/category/all/apis?order=field_popu...

**ProgrammableWeb**

Search Over 12,008 APIs | **Search APIs**

Filter APIs

By Category | By Protocols/Formats | ☐ Include Deprecated APIs

**(PROGRAMMABLEWEB.COM IS A GREAT REFERENCE)**

API Name

Google Maps | Mapping | 12.05.2005

# FINDING AND EXTRACTING EXISTING DATA

## SCRAPING THE WEB

## WHY SCRAPE?

**No API exists** for the data you want
(can't access the right data, wrong format, etc.)

**Simplicity** – Usually don't need to authenticate, no rate-limiting, etc.

Want to capture **context of pages** or relationship between them.

# FOR EXAMPLE...

le TOUR de france

07/05 > 07/27/2014

EN

Sunday July 27th, 2014

Stage 21
Évry / Paris Champs-Élysées

STAGE FINISHED

**19:16** Top 5

**19:14** The winner is... Marcel Kittel

**19:10** All together with 3km to go

THE RACE | ROUTE | CLASSIFICATIONS | TEAMS | VIDEOS & PHOTOS | HISTORY | STORE | Search

PARIS TOURS
12/10/2014

PREVIOUS

SUNDAY, JULY 27TH - STAGE 21    137.5km

Évry / Paris Champs-Élysées

NEXT

# Overall individual time classification
Total distance covered: **3660.5 KM**

**LCL**

| RANK | RIDER | RIDER NO. | TEAM | TIMES | GAP |
|------|-------|-----------|------|-------|-----|
| 1. | 🇮🇹 NIBALI Vincenzo | 41 | ASTANA PRO TEAM | **89h 59' 06"** | |
| 2. | 🇫🇷 PÉRAUD Jean-Christophe | 81 | AG2R LA MONDIALE | **90h 06' 43"** | + 07' 37" |
| 3. | 🇫🇷 PINOT Thibaut | 127 | FDJ.FR | **90h 07' 21"** | + 08' 15" |
| 4. | 🇪🇸 VALVERDE BELMONTE Alejandro | 11 | MOVISTAR TEAM | **90h 08' 46"** | + 09' 40" |
| 5. | 🇺🇸 VAN GARDEREN Tejay | 141 | BMC RACING TEAM | **90h 10' 30"** | + 11' 24" |
| 6. | 🇫🇷 BARDET Romain | 82 | AG2R LA MONDIALE | **90h 10' 32"** | + 11' 26" |
| 7. | 🇨🇿 KONIG Leopold | 201 | TEAM NETAPP-ENDURA | **90h 13' 38"** | + 14' 32" |
| 8. | 🇪🇸 ZUBELDIA AGIRRE Haimar | 169 | TREK FACTORY RACING | **90h 17' 03"** | + 17' 57" |
| 9. | 🇳🇱 TEN DAM Laurens | 67 | BELKIN PRO CYCLING | **90h 17' 17"** | + 18' 11" |
| 10. | 🇳🇱 MOLLEMA Bauke | 61 | BELKIN PRO CYCLING | **90h 20' 21"** | + 21' 15" |
| 11. | 🇫🇷 ROLLAND Pierre | 151 | TEAM EUROPCAR | **90h 22' 13"** | + 23' 07" |
| 12. | 🇱🇺 SCHLECK Frank | 161 | TREK FACTORY RACING | **90h 24' 54"** | + 25' 48" |
| 13. | 🇧🇪 VAN DEN BROECK Jurgen | 131 | LOTTO-BELISOL | **90h 33' 07"** | + 34' 01" |
| 14. | 🇷🇺 TROFIMOV Yury | 29 | TEAM KATUSHA | **90h 35' 47"** | + 36' 41" |
| 15. | 🇳🇱 KRUIJSWIJK Steven | 64 | BELKIN PRO CYCLING | **90h 37' 21"** | + 38' 15" |
| 16. | 🇫🇷 FEILLU Brice | 211 | BRETAGNE - SECHE ENVIRONNEMENT | **90h 43' 05"** | + 43' 59" |
| 17. | 🇺🇸 HORNER Christopher | 114 | LAMPRE - MERIDA | **90h 43' 37"** | + 44' 31" |
| 18. | 🇪🇸 NIEVE ITURRALDE Mikel | 5 | TEAM SKY | **90h 45' 37"** | + 46' 31" |
| 19. | 🇫🇷 GADRET John | 13 | MOVISTAR TEAM | **90h 46' 36"** | + 47' 30" |

# SOMETIMES YOU DON'T NEED A SCRAPER!

A few tips and tricks...

# PULLING DATA TABLES FROM THE WEB

**Google Sheets**

## IMPORTHTML

Imports data from a table or list within an HTML page.

# Demographics of India

*This article is about the people from India. For other uses, see Indian (disambiguation).*

The **demographics of India** are inclusive of the second most populous country in the world, with over 1.21 billion people (2011 census), more than a sixth of the world's population. Already containing 17.5% of the world's population, India is projected to be the world's most populous country by 2025, surpassing China, its population reaching 1.6 billion by 2050.[4][5] Its population growth rate is 1.41%, ranking 102nd in the world in 2010.[6] Indian population reached the billion mark in 2000.

| Demographics of India | |
|---|---|
| **Population** | 1,236,344,631 (July 2014 est.)[1] (2nd) |
| **Growth rate** | 1.51% (2009 est.) (93rd) |
| **Birth rate** | 20.22 births/1,000 population (2013 est.) |
| **Death rate** | 7.4 deaths/1,000 population (2013 est.) |
| **Life expectancy** | 68.89 years (2009 est.) |
| • **male** | 67.46 years (2009 est.) |
| • **female** | 72.61 years (2009 est.) |
| **Fertility rate** | 2.44 children born/woman (SRS 2011) |
| **Infant mortality rate** | 44 deaths/1,000 live births (2011 est.) |
| **Age structure** | |

**Population distribution in India by states**

| Rank | State / Union Territory | Type | Population | % [18] | Area [19] (km²) | Density (/km²) | Males | Females | Sex Ratio [20] | Literacy | Rural [21] Population | Urban [21] Population |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Uttar Pradesh | State | 199,812,341 | 16.50 | 240,928 | 828 | 104,480,510 | 95,331,831 | 912 | 67.68 | 131,658,339 | 34,539,582 |
| 2 | Maharashtra | State | 121,455,333 | 9.28 | 307,713 | 365 | 58,243,056 | 54,131,277 | 929 | 82.34 | 55,777,647 | 41,100,980 |
| 3 | Bihar | State | 103,804,637 | 8.60 | 94,163 | 1,102 | 54,278,157 | 49,821,295 | 918 | 61.80 | 74,316,709 | 8,681,800 |
| 4 | West Bengal | State | 91,276,115 | 7.54 | 88,752 | 1,030 | 46,809,027 | 44,467,088 | 950 | 76.26 | 57,748,946 | 22,427,251 |
| 5 | Madhya Pradesh | State | 72,626,809 | 6.00 | 308,245 | 236 | 37,612,306 | 35,014,503 | 931 | 69.32 | 44,380,878 | 15,967,145 |
| 6 | Tamil Nadu | State | 72,147,030 | 5.96 | 130,058 | 555 | 36,137,975 | 36,009,055 | 996 | 80.09 | 34,921,681 | 27,483,998 |
| 7 | Rajasthan | State | 68,548,437 | 5.66 | 342,239 | 201 | 35,550,997 | 32,997,440 | 928 | 66.11 | 43,292,813 | 13,214,375 |
| 8 | Karnataka | State | 61,095,297 | 5.05 | 191,791 | 319 | 30,966,657 | 30,128,640 | 973 | 75.36 | 34,889,033 | 17,961,529 |
| 9 | Gujarat | State | 60,439,692 | 4.99 | 196,024 | 308 | 31,491,260 | 28,948,432 | 919 | 78.03 | 31,740,767 | 18,930,250 |

`=ImportHtml("http://en.wikipedia.org/wiki/Demographics_of_India", "table",4)`

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | Rank | State / Union Territory | Type | Population | % [18] | Area [19] (km²) |
| 2 | 1 | Uttar Pradesh | State | 199,812,341 | 16.5 | 240,928 |
| 3 | 2 | Maharashtra | State | 121,455,333 | 9.28 | 307,713 |
| 4 | 3 | Bihar | State | 103,804,637 | 8.6 | 94,163 |
| 5 | 4 | West Bengal | State | 91,276,115 | 7.54 | 88,752 |
| 6 | 5 | Madhya Pradesh | State | 72,626,809 | 6 | 308,245 |
| 7 | 6 | Tamil Nadu | State | 72,147,030 | 5.96 | 130,058 |

# PARSING

## Tabula



Tabula is a tool
locked inside P

### Extracted tabular data ×

| 2 | | |
|---|---|---|
| All Students | 79,858 | 99% |
| Gender | | |
| Male | 40,492 | 98% |
| Female | 39,134 | 99% |
| Ethnicity | | |
| White | 10,665 | 99% |
| Black | 49,379 | 99% |
| Latino/Hispanic | 13,717 | 98% |
| Asian | 4,746 | 100% |
| Native American | 132 | 99% |
| Multiracial | 941 | 98% |
| Other Groups | | |
| IEP | 11,471 | 98% |

☐ Use row/columns separators ❓

Close | Copy to clipboard as CSV | Download data ▾

Page 3

year of enrollment in a U.S. school.

# BUILDING A WEB SCRAPER

FETCHING DATA  +  PARSING DATA

YOU SHOULD **SEPARATE** THESE PROCESSES **WHENEVER POSSIBLE**!

# FETCHING

## DON'T DO EVERYTHING AT ONCE

Download complete pages and save them locally <u>before</u> you process them.

## DEALING WITH PAGINATION

If results or records are spread across multiple pages, you may need to parse the page to find the link to the next page.

# PARSING

## SERIOUSLY, DON'T DO EVERYTHING AT ONCE!

Processing data from local files means
you don't have to get it right the first time.

## USE YOUR BROWSER'S DEVELOPER TOOLS

All modern web browsers have built-in tools
that let you inspect web pages.

# BE CAREFUL - YOU CAN GET YOURSELF BLOCKED

Many sites will try to slow or block heavy access (both to prevent scraping and DoS attacks)

To get around this...You can introduce delays in your scraper or scrape from multiple locations.

# A FEW MORE NOTES ABOUT DATA MANAGEMENT

## FORMATS AND BEST-PRACTICES

# DATA FORMATS

STRUCTURED DATA is more like what you'd find in a traditional spreadsheet or database.

UNSTRUCTURED DATA can include raw text, streaming data, even images or video.

SEMI-STRUCTURED DATA is more organized, but doesn't follow a fixed schema (e.g. DBPEDIA data)

# CSV

## (Comma-Separated Value)

```
1  firstName,lastName,age,streetAddress,city,state
2  John,Smith,25,21 2nd Street,New York,NY,10021,2
```

| firstName | lastName | age | streetAddress | city | state | postalCode | homePhoneNumber | faxPhoneNumber | gender |
|---|---|---|---|---|---|---|---|---|---|
| John | Smith | 25 | 21 2nd Street | New York | NY | 10021 | 212 555-1239 | 646 555-4567 | male |

# XML

## (eXtensible Markup Language)

```xml
<person>
  <firstName>John</firstName>
  <lastName>Smith</lastName>
  <age>25</age>
  <address>
    <streetAddress>21 2nd Street</streetAddress>
    <city>New York</city>
    <state>NY</state>
    <postalCode>10021</postalCode>
  </address>
  <phoneNumbers>
    <phoneNumber type="home">212 555-1234</phoneNumber>
    <phoneNumber type="fax">646 555-4567</phoneNumber>
  </phoneNumbers>
  <gender>
```

| firstName | lastName | age | streetAddress | city | state | postalCode | homePhoneNumber | faxPhoneNumber | gender |
|-----------|----------|-----|---------------|------|-------|------------|-----------------|----------------|--------|
| John | Smith | 25 | 21 2nd Street | New York | NY | 10021 | 212 555-1239 | 646 555-4567 | male |

# JSON

## (JavaScript Object Notation)

```json
{
  "firstName": "John",
  "lastName": "Smith",
  "age": 25,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021"
  },
  "phoneNumber": [
    {
      "type": "home",
      "number": "212 555-1239"
    },
    {
      "type": "fax",
      "number": "646 555-4567"
    }
  ]
}
```

| firstName | lastName | age | streetAddress | city | state | postalCode | homePhoneNumber | faxPhoneNumber | gender |
|-----------|----------|-----|---------------|------|-------|------------|-----------------|----------------|--------|
| John | Smith | 25 | 21 2nd Street | New York | NY | 10021 | 212 555-1239 | 646 555-4567 | male |

# YAML

## (YAML Ain't Markup Language)

```yaml
---
  firstName: John
  lastName: Smith
  age: 25
  address:
      streetAddress: 21 2nd Street
      city: New York
      state: NY
      postalCode: 10021

  phoneNumber:
      -
        type: home
        number: 212 555-1234
      -
        type: fax
```

| firstName | lastName | age | streetAddress | city | state | postalCode | homePhoneNumber | faxPhoneNumber | gender |
|-----------|----------|-----|---------------|------|-------|------------|-----------------|----------------|--------|
| John | Smith | 25 | 21 2nd Street | New York | NY | 10021 | 212 555-1239 | 646 555-4567 | male |

# HANDLING DATA

## STORING DATA

– Always keep <u>backups</u>

– <u>Password protect</u> or <u>encrypt</u> any data with personal or sensitive information

## PROVENANCE

– Keep track of <u>where/when</u> data was collected

– Record any data processing steps so you (or others) can repeat them if necessary

# IP, COPYRIGHT, AND (RE)SHARING DATA

- Be sure you know who <u>owns</u> the data.

- Think early on about whether or not you'll need to <u>publish</u> or <u>(re)share</u> data.

- Be careful you <u>aren't violating copyright</u>, especially when scraping.

# PRIVACY AND ANONYMIZING DATA

- Any information that could be used to <u>identify individuals</u> is sensitive!

- There may be <u>legal reprecussions</u> for releasing it.

- In some cases you might need to <u>anonymize</u> data before sharing.

**JUST REMOVING NAMES IS OFTEN NOT ENOUGH!**

# OTHER INFORMATION CAN STILL BE UNIQUE



Medical Data — Ethnicity, Visit date, Diagnosis, Procedure, Medication, Total charge

ZIP, Birth date, Sex

Voter List — Name, Address, Date registered, Party affiliation, Date last voted

[L. Sweeney. 2002]
k-ANONYMITY: A MODEL FOR PROTECTING PRIVACY

# LOCATION DATA IS ESPECIALLY SENSITIVE



[de Montjoye et al. 2013]
Unique in the Crowd: The privacy bounds of human mobility

# REGULATIONS (ACADEMIA AND RESEARCH)

Institutional Review and Ethics Boards may need to approve experiments or data collection <u>before it happens.</u>

Studies involving people may need <u>informed consent.</u>

# REGULATIONS (INDUSTRY)

Some governments have placed limits on how long user data can be kept.

Some kinds of tracking (e.g., cookies) may now require opt-in or notifications. (However this varies by country).

# SOCIAL EXPERIMENTS

## Experimental evidence of massive-scale emotional contagion through social networks

Adam D. [...] [...]ock[c,d]

[a]Core Da[...] [...] [...]esearch and Education, University of California, San Francisco, CA 9414[...] [...] NY 14853

Edited [...] [...]3, 2013)

Emot[...] [...]rs via text-based cont[...] [...] of psy-without the[...] in laboratory exper[...] negative emotions to others. network, collected over a 20-y period [...] moods (e.g., depression, happiness) can be tr[...] networks [Fowler JH, Christakis NA (2008) *BMJ* 337:a[...] though the results are controversial. In an experiment with people who use Facebook, we test whether emotional contagion occurs outside of in-person interaction between individuals by reducing the amount of emotional content in the News Feed. When positive expressions were reduced, people produced fewer positive posts and more negative posts; when negative expressions were re-duced, the opposite pattern occurred. These results indicate that emotions expressed by others on Facebook influence our own emotions, constituting experimental evidence for massive-scale contagion via social networks. This work also suggests that, in contrast to prevailing assumptions, in-person interaction and non-verbal cues are not strictly necessary for emotional contagion, and

[...] [...] later seen by [...] (8). Because people's [...] content than one person can [...] stories, and activities undertaken by [...] primary manner by which people see content tha[...] Which content is shown or omitted in the News Feed [...] termined via a ranking algorithm that Facebook continually develops and tests in the interest of showing viewers the content they will find most relevant and engaging. One such test is [...]

## Senator asks FTC to investigate Facebook's mood study

After the social network altered the news feeds of nearly 700,000 users without telling them, Sen. Mark R. Warner wants to know if there should be oversight on these types of experiments.

# "EXPERIMENTING ON HUMAN BEINGS"

Dating Research from OkCupid

## We Experiment On Human Beings!

July 28th, 2014 by Christian Rudder

Tweet 2,760      f Share 10k

I'm the first to admit it: we might be popular, we might create a lot of great relationships, we might blah blah blah. But OkCupid doesn't

# IN SUMMARY: THERE ARE LOTS OF TOOLS AT YOUR DISPOSAL!

## COLLECT IT
– OBSERVATION
– SURVEYS
– LOGGING
– SENSORS
– CROWDSOURCING

## FIND OR EXTRACT IT
– OPEN CORPUSES
– DATA RETAILERS
– APIS
– SCRAPING THE WEB

## GENERATE IT
– SIMULATIONS

(...AND WE'LL
BE HAPPY TO DISCUSS
OR SUGGEST MORE)

**NEXT UP**

TUTORIAL 2 – BUILDING A WEB SCRAPER

NEXT WEEK

DATA CLEANING

# BEFORE NEXT WEEK'S CLASS

INSTALL :

**OpenRefine** (formerly Google Refine)
http://openrefine.org/

# DATA COLLECTION

Slides by WESLEY WILLETT

VISUAL ANALYTICS