# Geospatial Data Acquisition and Evaluation

## Center for Geographic Analysis

Harvard University

**Nicole Alexander, Ph. D.**
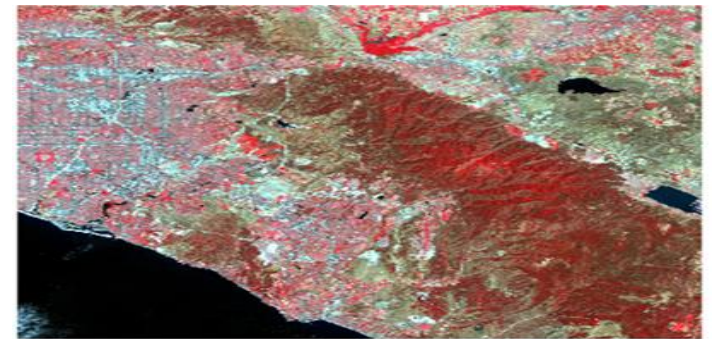
nalexander@cga.harvard.edu

Harvard DataFest 2018

# Geospatial Data Acquisition and Evaluation

- Geospatial Data Sources
- Data Transfer
- Metadata
- Geocoding Data Science Datasets in R using Google's API

# Geospatial Data Acquisition and Evaluation

- How data are captured determines the quality of decisions that can be made from analyzing the data
  - *Primary* sources: obtained through direct measurement
  - *Secondary* sources: derived from other sources
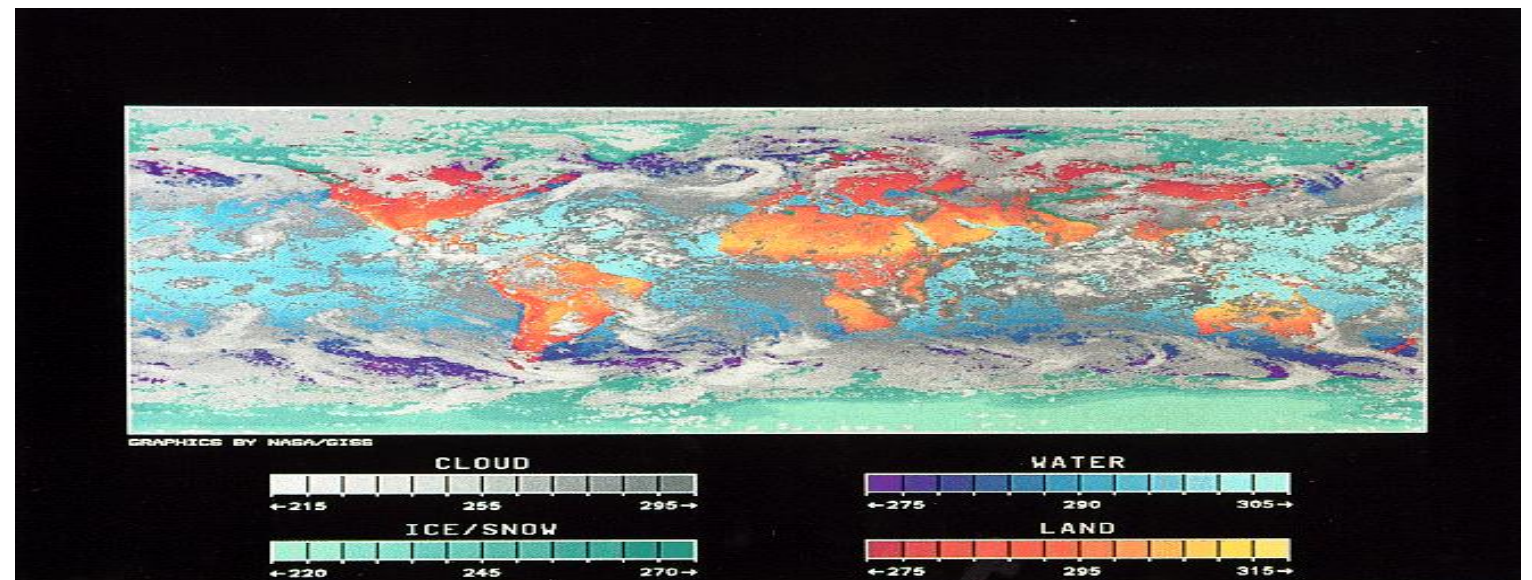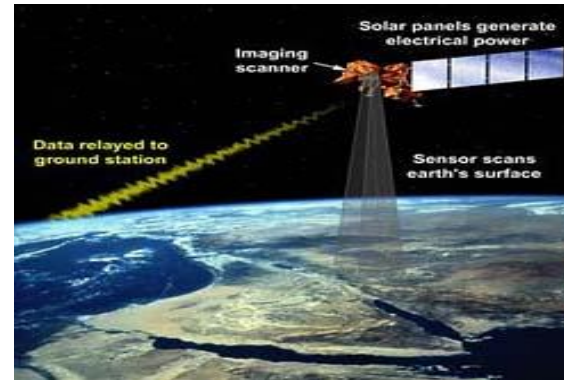- Data accuracy can more reliably be determined from primary sources

# Geospatial Data Sources

|  | RASTER | VECTOR |
|---|---|---|
| **Primary** | • Digital satellite remote-sensing images<br>• Digital aerial photographs | • GPS measurements<br>• Field survey measurements<br>• LiDAR |
| **Secondary** | • Scanned maps and photographs<br>• Digital elevation models from topographic map contours | • Topographic maps<br>• Toponymy (place-name) databases<br>• Geocoding |

# Primary Raster Data Capture

- Remote sensing
  - Satellite
  - Aircraft

- Image Resolution
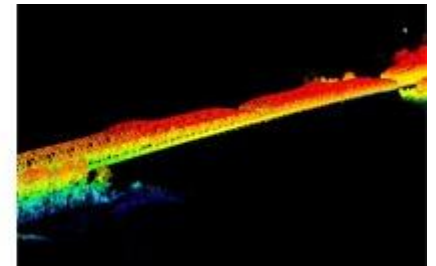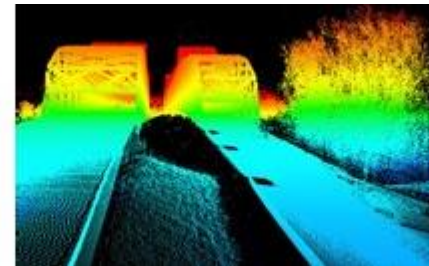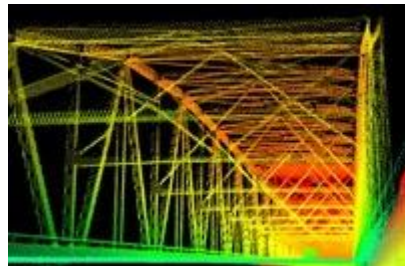  - Spatial
  - Spectral
  - Temporal

# Remote Sensing Data Capture

- Captures data over a large geographic areas
  - Total ground coverage range from 9 x 9 km – 200 x 200 km
- Pixel size determines spatial resolution of an image
  - Spatial accuracy of features increases as pixel size decreases
- Satellite systems capture data in the range of 0.5 m – 1 km pixel size
- Camera systems capture data in the range of 0.01 m – 5 m pixel size
- Costly compared with other methods of data capture
- Data volumes can be very large

# Primary Vector Data Capture



- Main sources
  - GPS
  - Surveying

- Remote Sensing
  - LiDAR (Light Detection And Ranging)
    - a "cloud" of points that reflects the surface

# Primary Vector Data Capture

- GPS
  - Recreational: low precision 6 – 12m
  - Mapping and GIS: medium precision 30cm – 5m
  - Surveying: high precision 5mm – 1cm
- Surveying
  - Used for large scale mapping of small areas and property boundaries
  - Capable of 1 mm accuracy
  - Equipment and crews are expensive
- LiDAR
  - 30,000 points per second at an accuracy of around 15cm
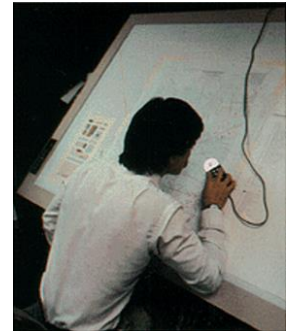  - Often rasterized to create DEMs

# Secondary Raster Data Capture

- Scanning of hardcopy media
  - Building plans, CAD drawings, property deeds, film, paper maps, aerial photographs, images, etc.
  - Spatial resolution of scanners in the range of 400 – 900 dpi (16 – 40 dots per mm)
- DEM generation from topographic map contours or LiDAR

# Secondary Vector Data Capture

- Vectorizing raster data
- Digitizing
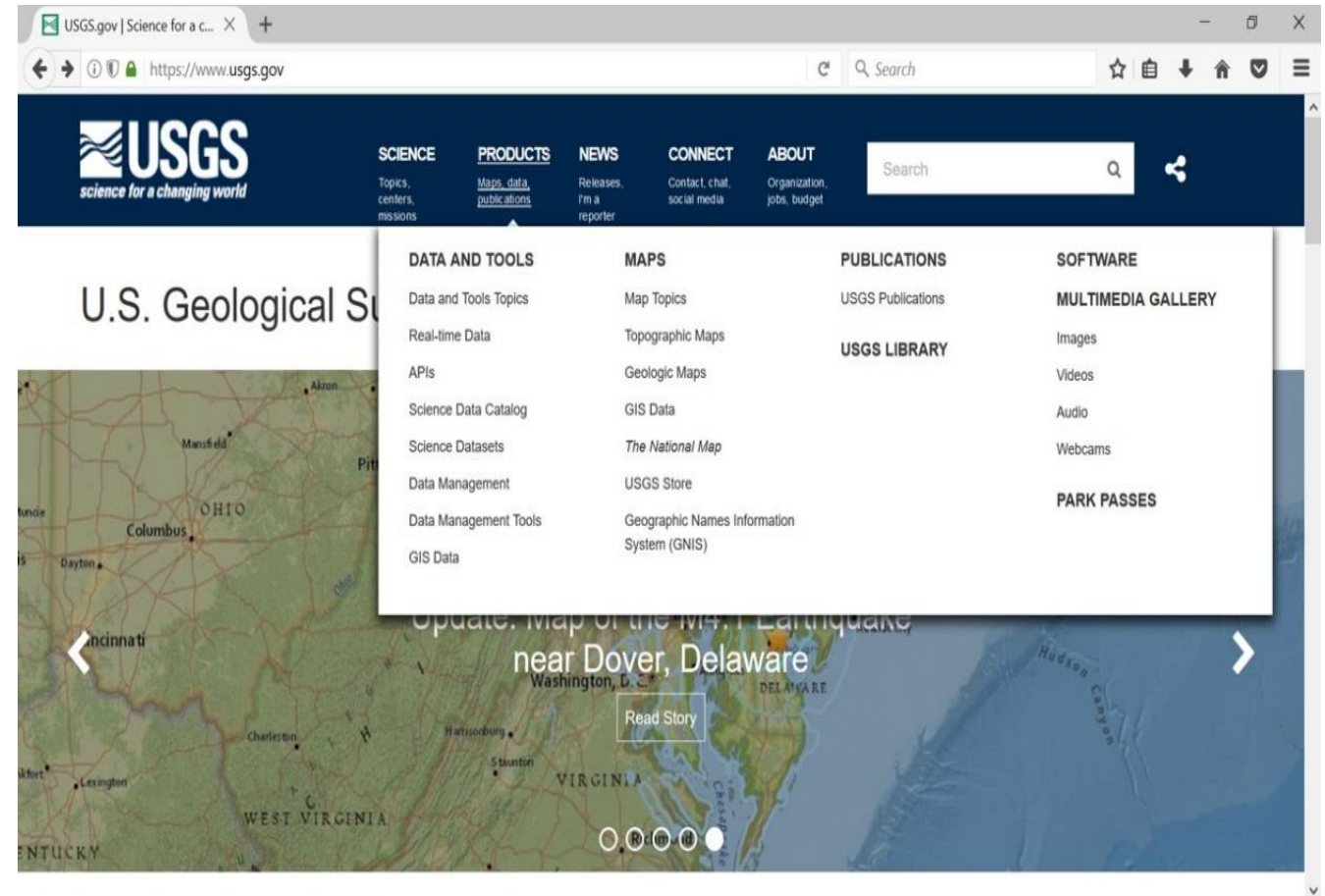- Geocoding
- Photogrammetry
- COGO – Coordinate Geometry

# Data Transfer: Obtaining Data from External Sources

- U.S. Geological Survey
- U.S. Census Bureau
- OpenStreetMap
- GeoNames
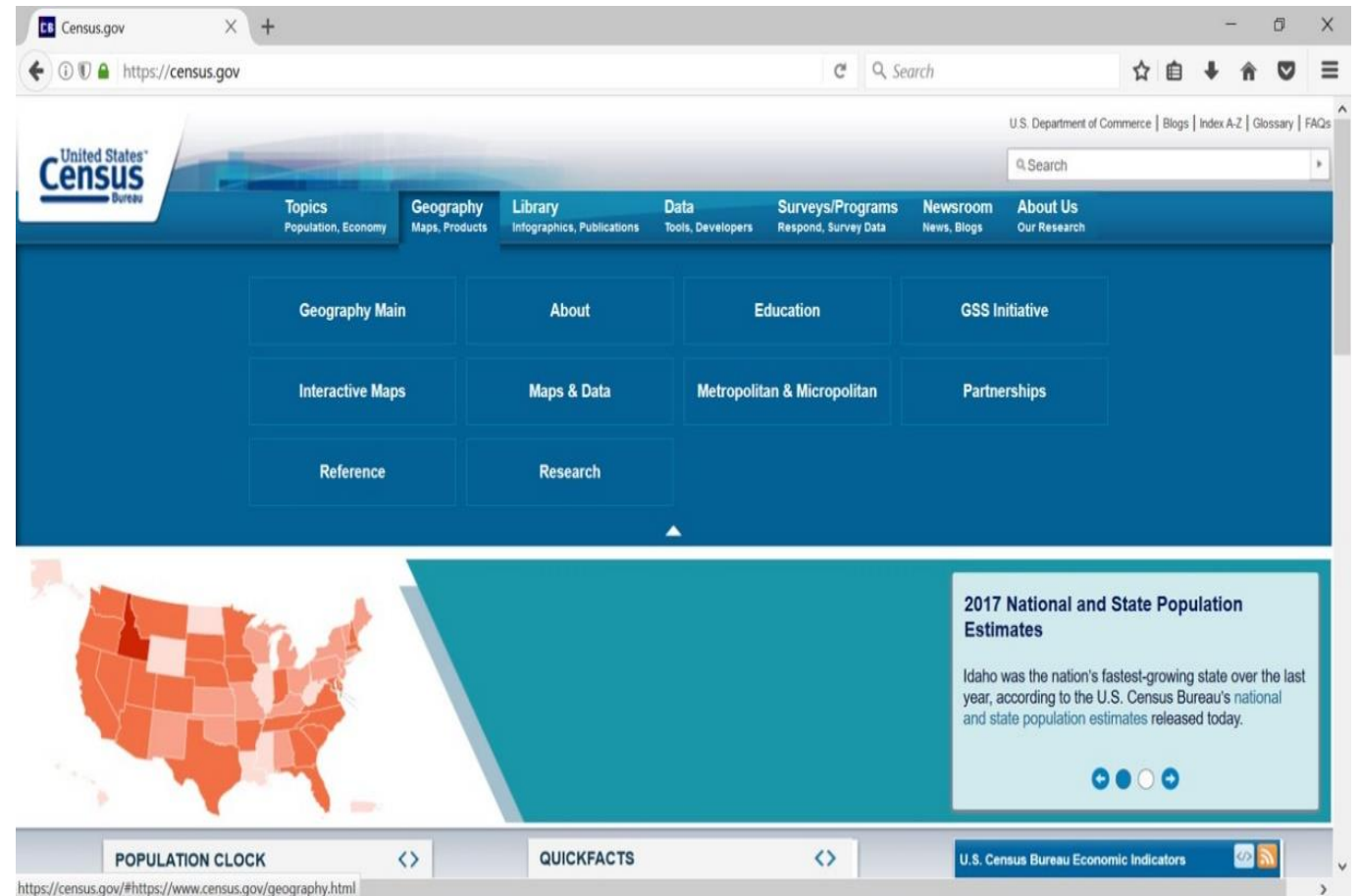- Other Geospatial Data Sites

# U.S. Geological Survey (usgs.gov)

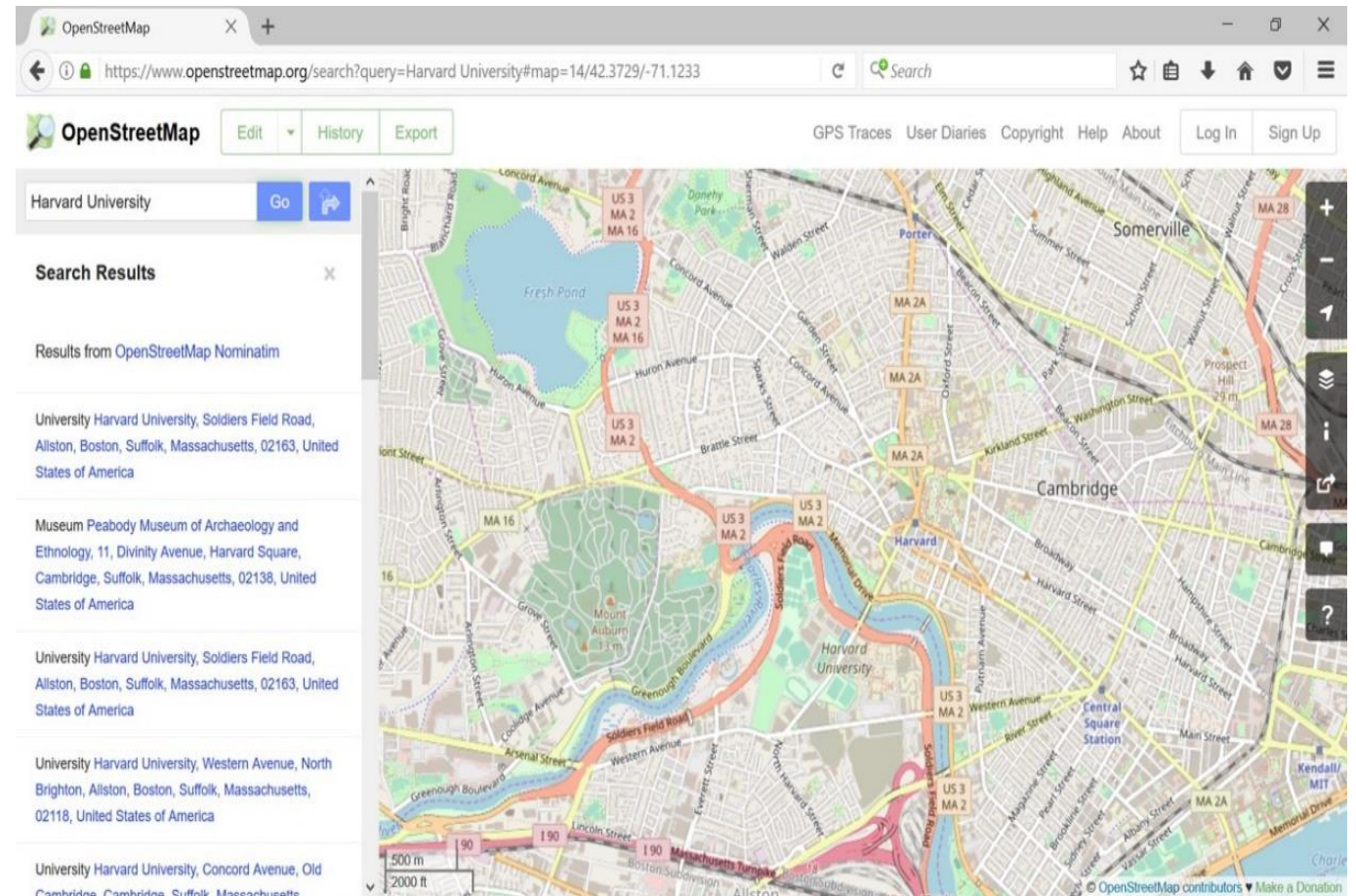- The major provider of geospatial data in the US

# U.S. Census Bureau (census.gov)
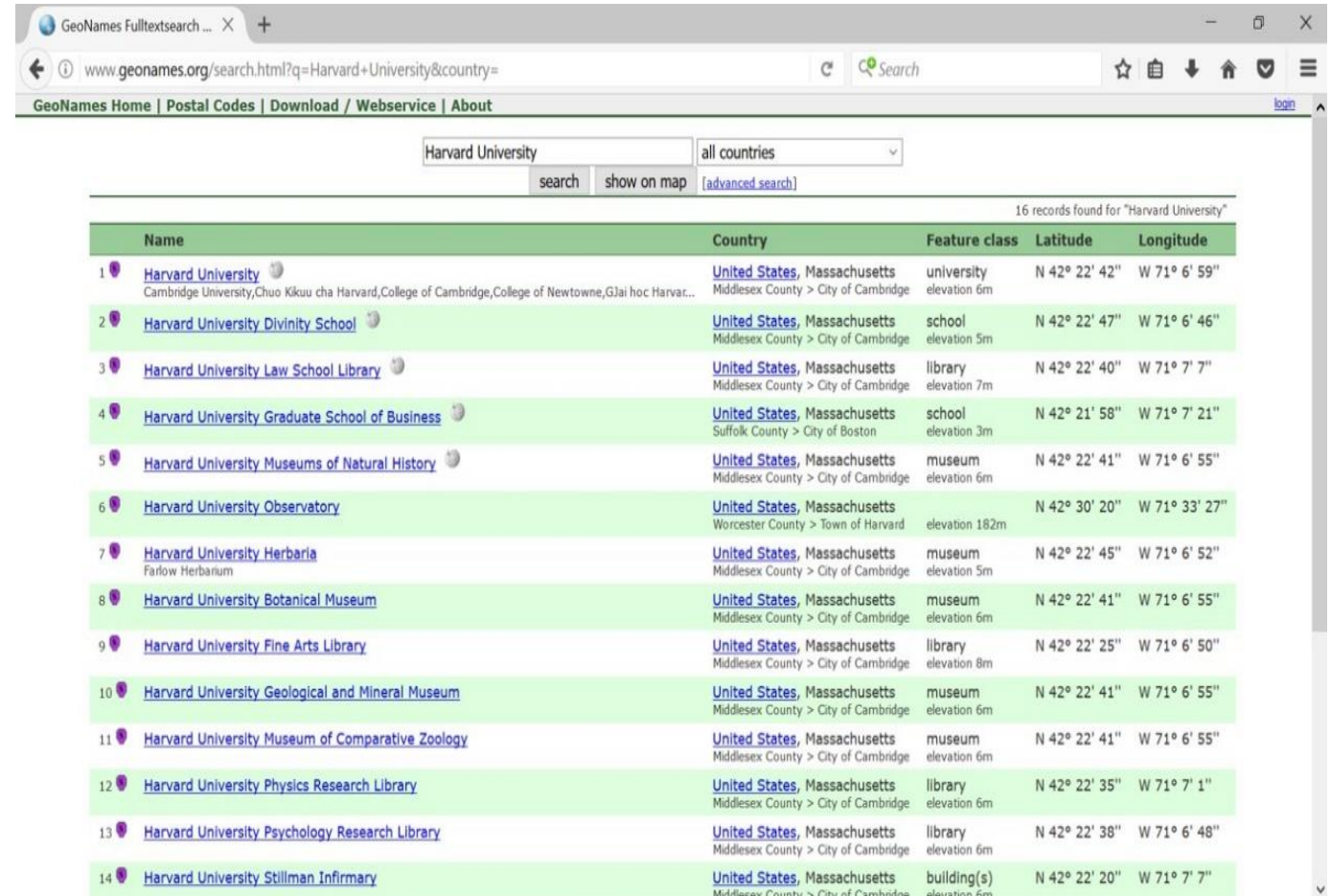
- Provides data to support the US decennial census

# OpenStreetMap (OSM) (openstreetmap.org)

- Map data of the world
- Created and maintained by a community of mappers

# GeoNames (geonames.org)

- Global geographical database of place names

# Other Geospatial Data Sites

- Harvard University
  - CGA: http://gis.harvard.edu/resources/data
  - Harvard Geospatial Library: http://hgl.harvard.edu
  - Harvard WorldMap: http://worldmap.harvard.edu
  - Harvard Map Collection: http://hcl.harvard.edu/libraries/maps/collections/digital.html#overview
- Local
  - MassGIS: http://www.mass.gov/mgis
  - City of Boston: https://data.boston.gov/dataset?groups=geospatial
  - Metro Boston Data Common: http://www.metrobostondatacommon.org/
- National
  - US Federal Government: http://data.gov
  - US Geological Survey: http://viewer.nationalmap.gov/viewer/
- Global
  - The ESRI Data and Maps: http://bit.ly/NBoQzQ
  - ArcGIS Online Services from ESRI: http://www.arcgis.com/home/

# Metadata

- Data about the geospatial data:
  - Identification
  - Data quality
  - Coordinate system
  - Attributes, etc.
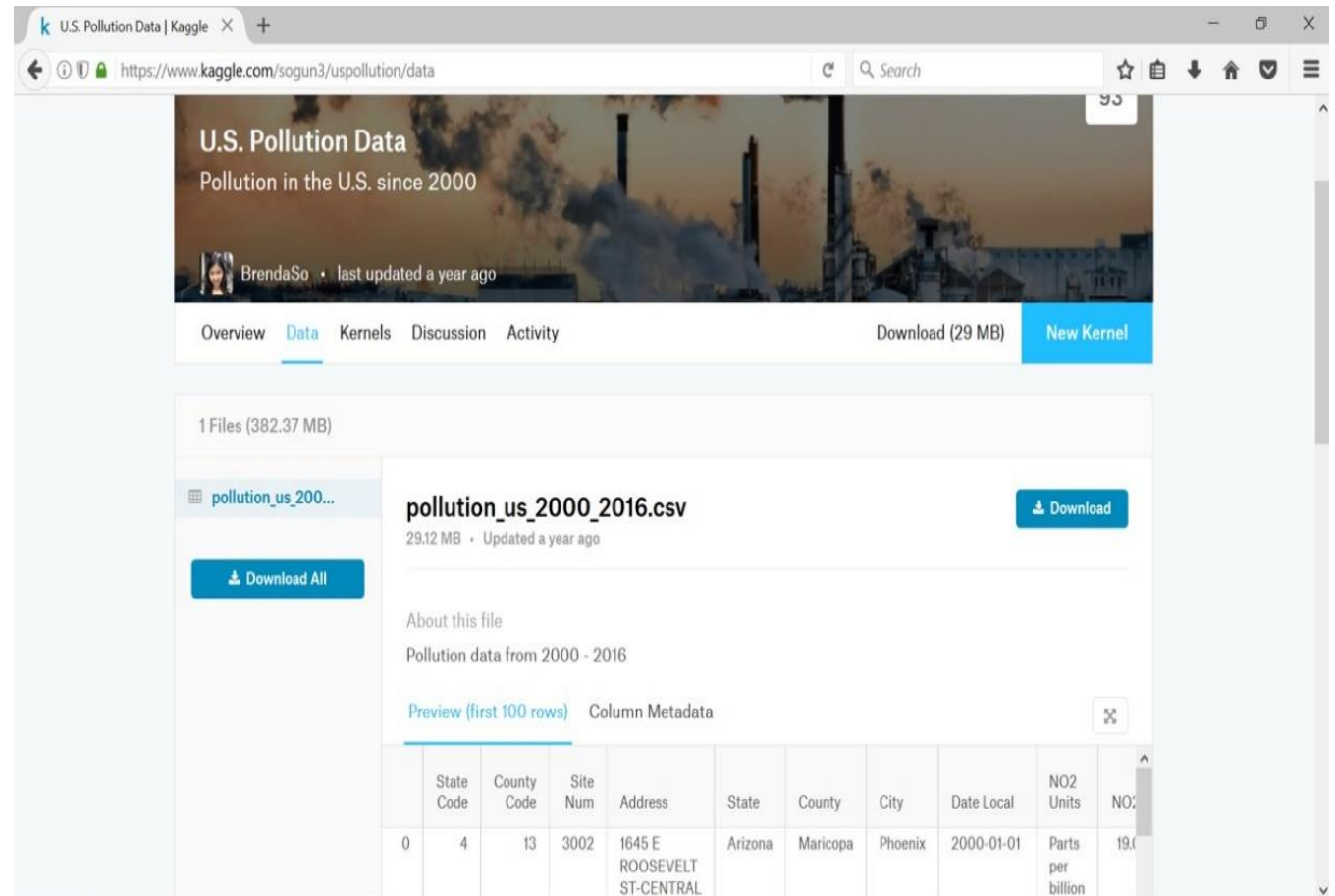- Especially important when using public data

# Geocoding Data Science Datasets in R using Google's API

- Kaggle
- Geocoding a CSV of Addresses using Google's API

# Kaggle (kaggle.com)

- Data science and machine learning site

- Datasets containing place-name or address information can be geocoded to perform spatial analysis

# Demo: Geocoding a CSV of Addresses in R using Google's Geocoding API

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Go to file/function          Addins ▾

googleRgeocode.R ×

☐ Source on Save

```r
1  #1. Install ggmap package (and dependencies) from the Repository in RStudio: Packages -> install -> ggmap
2
3  #2. Open a new R script: File -> New File -> R Script
4
5  #3. Load ggmap
6  library(ggmap)
7
8  #4. Select the CSV file from the file chooser
9  fileToLoad <- file.choose(new = TRUE)
10
11 #5. Read the CSV data and store it in a variable (origAddress)
12 origAddress <- read.csv(fileToLoad, stringsAsFactors = FALSE)
13
14 #6. Initialize the data frame
15 geocoded <- data.frame(stringsAsFactors = FALSE)
16
17 #7. Loop through the addresses in the CSV file to get the latitude and longitude
18 # of each address and add it to the
19 # origAddress data frame in new columns lat and lon
20 for(i in 1:nrow(origAddress))
21 {
22   result <- geocode(origAddress$addresses[i], output = "latlona", source = "google")
23   origAddress$lon[i] <- as.numeric(result[1])
24   origAddress$lat[i] <- as.numeric(result[2])
25   origAddress$geoAddress[i] <- as.character(result[3])
26 }
27
28 #8. Write a CSV file containing origAddress to the working directory
29 write.csv(origAddress, "geocoded.csv", row.names=FALSE)
30
```

22:85     (Top Level) ▾

Console   Terminal ×

~/ ⇨

>

# Data Acquisition and Evaluation Summary

- Data collection can be expensive and time-consuming
  - Main techniques
    - Primary
      - Raster – e.g. remote sensing
      - Vector – e.g. GPS, field survey and LiDAR
    - Secondary
      - Raster – e.g. scanning
      - Vector – e.g. digitizing and geocoding

- Conversion of existing data and online data options available

- Always ask first: **to buy or to build?**

# References

- Longley, P. A., M. F. Goodchild, D. J. Maguire, D. W. Rhind (2010). Geographic Information Systems & Science (3rd Ed). John Wiley & Sons, Inc.

- Clarke, K.C. (2010). Getting Started with GIS (5th Ed). Prentice-Hall, Inc., London.

- Chang, K-T. (2010). Introduction to Geographic Information Systems (5th Ed). McGraw-Hill.

- Center for Geographic Analysis (CGA) Harvard: GIS Training Materials.

- Storybench: http://www.storybench.org/geocode-csv-addresses-r/

- Google: https://developers.google.com/maps/documentation/geocoding/start

# Q&A