

# Data preparation

“Give me six hours to  
chop down a tree and I  
will spend the first four  
sharpening the axe.”

— Abraham Lincoln (?)

# Know Your Questions



# Gathering your data



# Manual



2010 ANNUAL REPORT — CATALOGUED ITEMS

# Self Generated or User Generated

Untitled form All changes saved in Drive

Name and attach values list

Which Sheet? Which data? Name the values list!

Zapbook IA

Select Sheet  
Live Integrations with Categories

Select column header  
Name

Column must have a header label in row 1.

Data preview...  
Name

---

Google Sheets  
Trello  
Google Calendar  
Gmail  
Slack  
MailChimp  
Evernote  
Twitter

Short answer text

Back Next

FORMRANGER

Assign Form questions to be populated from values lists in a Sheet or student rosters created in the [Doctopus Add-on](#). Only multiple choice, checkbox, list, and grid type questions are supported.

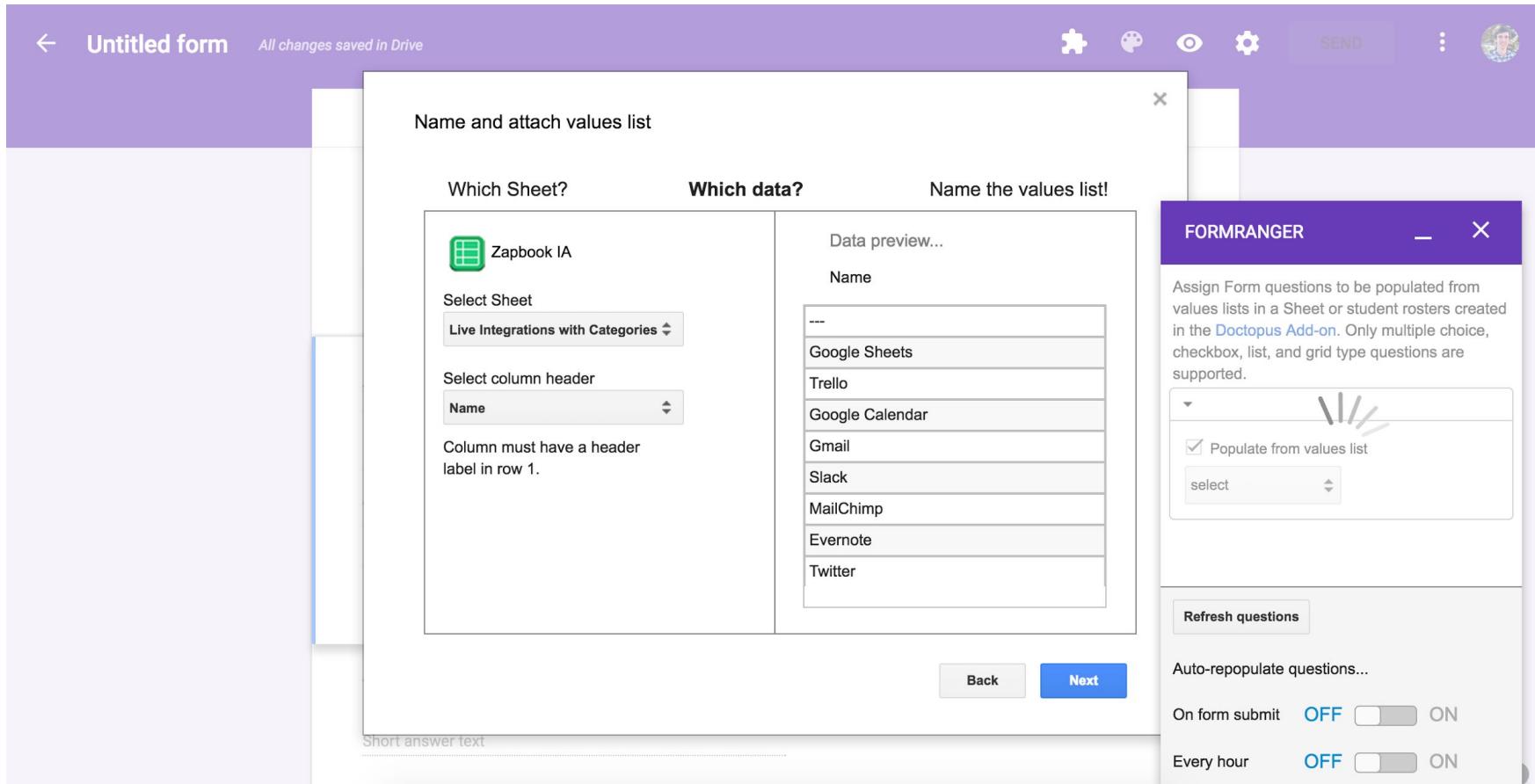
Populate from values list  
select

Refresh questions

Auto-repopulate questions...

On form submit OFF ON

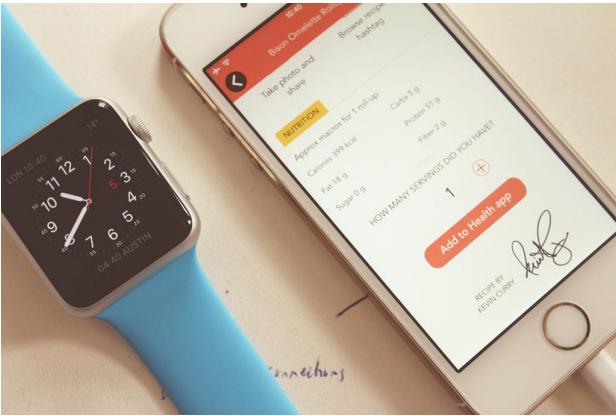
Every hour OFF ON



**Self Generated or User Generated**



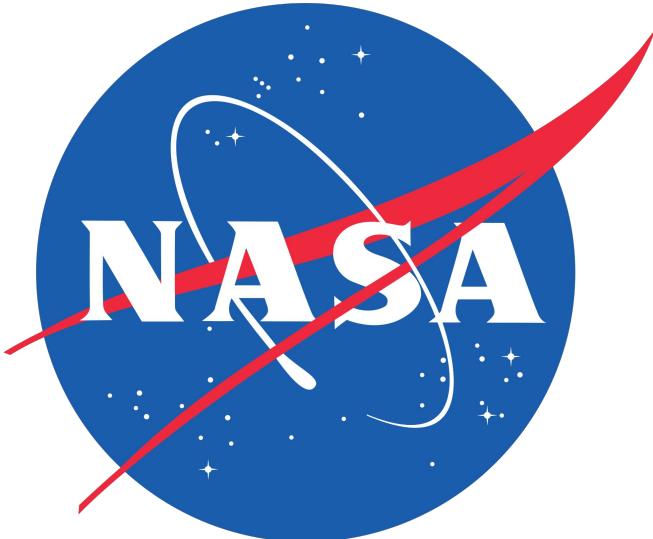
# Fit Bit, Nike+, iPhone Health App, Strava, Biometrics etc.



## Existing Data Sources



reddit



twitter



U B E R

# Tidy data / Don't mix different type of data in the same column

Summit Name	Country	Lat.
Mt. Everest	Nepal/Tibet	27°59' N
Nanga Parbat	Pakistan	35°14' N
Namcha Barwa	China	29°38' N
Jengish Chokusu	Kyrgyzstan/China	42°02' N
Aconcagua	Argentina	32°39' S
Chimborazo	Ecuador	1°28' S
Mt. McKinley (Denali)	US	63°06' N
Mt. Logan	Canada	60°34' N

Columns:  
attribute  
(or variable)

Elevation (ft.)	Prominence (ft.).	First ascent
29028	29028	1953
26657	15118	1953
25531	13471	1992
24406	13609	1956
22841	22841	1897
20561	13523	1880
20320	20138	1913
19550	17224	1925

Cristobal Colon	Colombia	10°50' N
Gora Elbrus	Russia	43°21' N
Pico de Orizaba (Citlaltepetl)	Mexico	19°02' N
Damavand	Iran	35°57' N
Bogda Shan	China	43°48' N
Vinson Massif	Antarctica	78°31' S
Puncak Jaya	Indonesia	4°04' S
Mont Blanc	France/Italy	45°50' N
Klyuchevskaya Volcano	Russia	56°03' N
Mauna Kea	US	19°49' N
Kinabalu	Malaysia	6°04' N

Rows: observations

## Tidy data is easier to work with

For most charting tools (Excel, Illustrator, etc.), omit commas, symbols, letters for quantitative data:

\$2,745        #        2745

+76%        #        76

65mph        #        65

212°F        #        212

March 1, 2017 # 20170301

73°41'W        #        -73.41

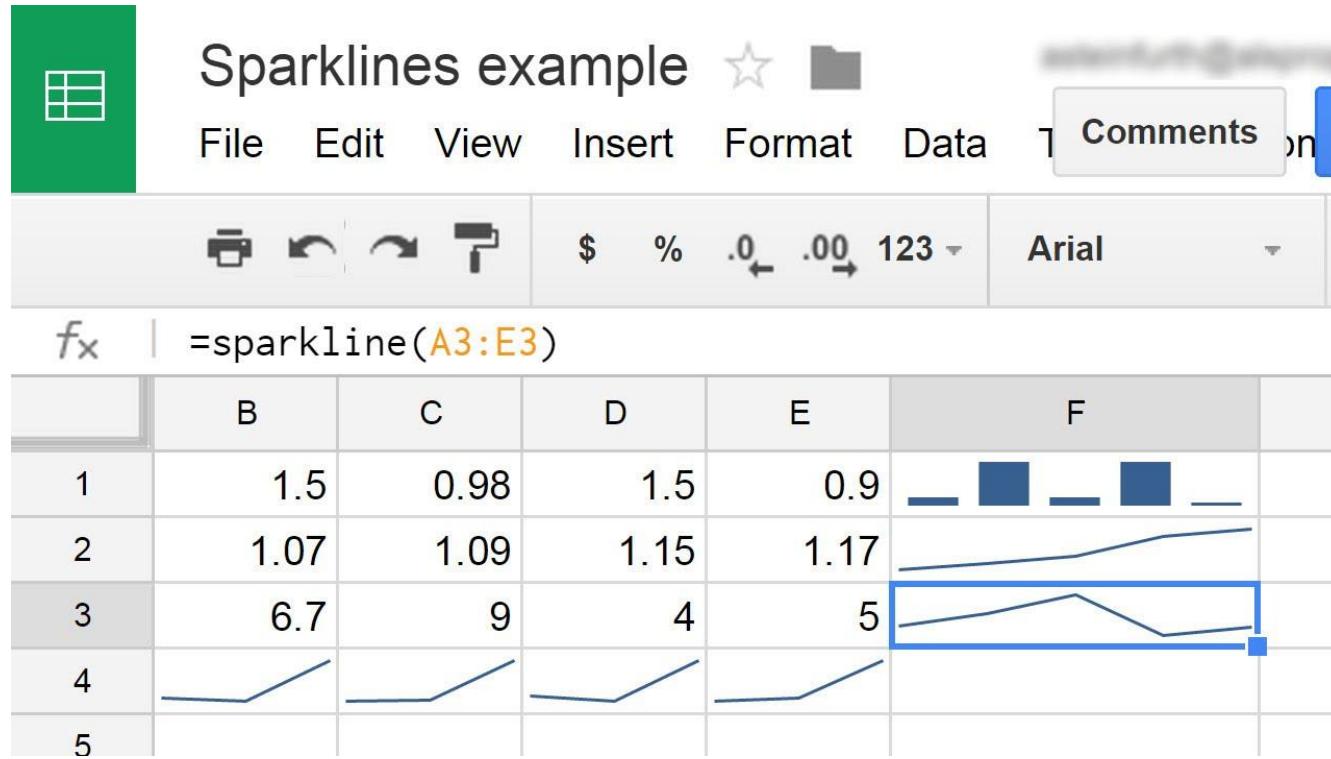
# CSV files / Comma-separated values

Rank	Summit Name	Country	Lat.	Long.	Elevation (ft.)	Elevation (m)	Prominence (ft.)	First ascent
1	Mt. Everest	Nepal/Tibet	27.59	86.55	29028	8848	29028	1953
2	Nanga Parbat	Pakistan	35.14	74.36	26657	8125	15118	1953
3	Namcha Barwa	China	29.38	95.03	25531	7782	13471	1992
4	Jengish Chokusu	Kyrgyzstan/China	42.02	80.07	24406	7439	13609	1956
5	Aconcagua	Argentina	-32.39	-70.01	22841	6962	22841	1897
6	Chimborazo	Ecuador	-1.28	-78.49	20561	6267	13523	1880
7	Mt. McKinley (Denali)	US	63.06	-151.03	20320	6194	20138	1913
8	Mt. Logan	Canada	60.34	-140.24	19550	5959	17224	1925
9	Kilimanjaro	Tanzania	-3.04	37.21	19340	5895	19308	1889
10	Cristobal Colon	Colombia	10.50	-73.41	18701	5700	18074	1939
11	Gora Elbrus	Russia	43.21	42.26	18510	5642	15554	1874
12	Pico de Orizaba (Citlaltepetl)	Mexico	19.02	-97.16	18491	5636	16148	1848
13	Damavand	Iran	35.57	52.07	18405	5610	15311	905
14	Bogda Shan	China	43.48	88.20	17864	5445	13523	1981
15	Vinson Massif	Antarctica	-78.31	-85.36	16050	4892	16050	1966
16	Puncak Jaya	Indonesia	-4.04	137.11	16023	4884	16023	1962
17	Mont Blanc	France/Italy	45.50	6.52	15777	4809	15406	1786
18	Klyuchevskaya Volcano	Russia	56.03	160.39	15584	4750	15252	1788
19	Mauna Kea	US	19.49	-155.28	13796	4205	13796	1823
20	Kinabalu	Malaysia	6.04	116.34	13435	4095	13435	1851

# Data analysis

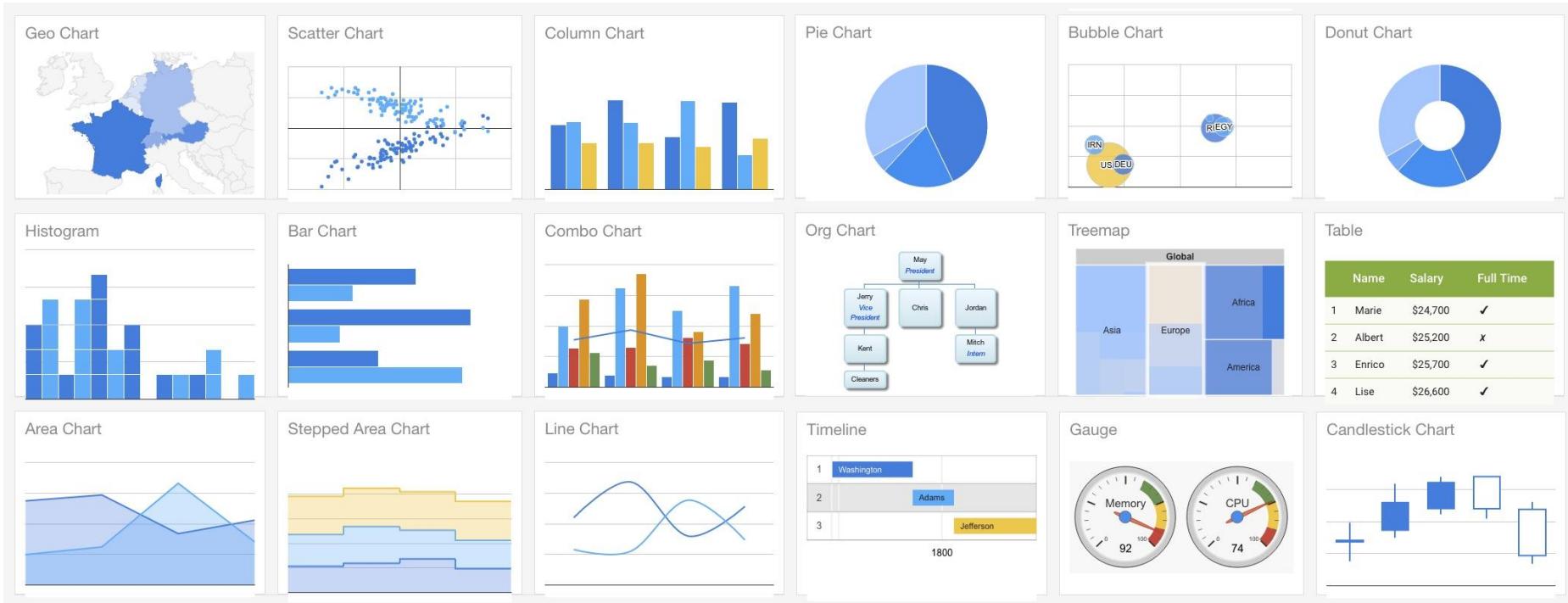
# Suggested tool

Sparklines - Google Sheets ([EXAMPLE](#))



# Suggested tool

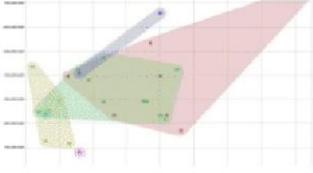
Google Charts - requires some coding knowledge



# Suggested tool

## RAWGraphs

**Convex Hull**  
Dispersion



In mathematics, the convex hull is the smallest convex shape containing a set of points. Applied to a scatterplot, it is useful to identify points belonging to the same category.

Based on  
<http://blocks.org/rmbostock/4341699>

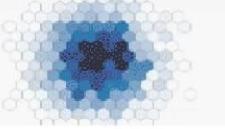
**Convex Hull**  
Dispersion



**Delaunay Triangulation**  
Dispersion



**Hexagonal Binning**  
Dispersion



**Scatter Plot**  
Dispersion



**Voronoi Tessellation**  
Dispersion



**Box plot**  
Distribution



**Circular Dendrogram**  
Hierarchy



**Cluster Dendrogram**  
Hierarchy



**Circle Packing**  
Hierarchy (weighted)



**Clustered Force Layout**  
Hierarchy (weighted)



**Sunburst**  
Hierarchy (weighted)

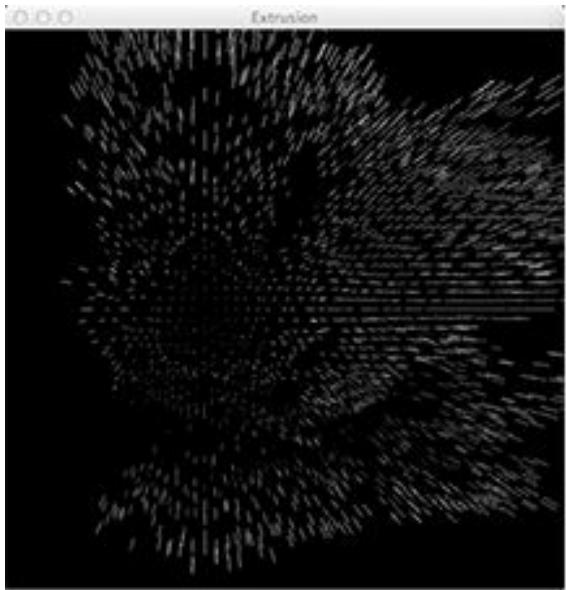


**Treemap**  
Hierarchy (weighted)



# Suggested tool

Processing



Display Window

The screenshot shows the Processing IDE interface. At the top, there's a toolbar with various icons and tabs labeled "Extrusion". Below the toolbar is a text editor containing Java-like pseudocode for a "Shape" class. The code includes methods for setup and draw, and it involves loading an image, extracting pixel values, and calculating depth based on color. The bottom part of the interface shows a message area with the text "Done Saving." and a text area below it.

```
// Extrusion.  
//  
// Converts a flat image into spatial data points and rotates the points  
// around the center.  
//  
#Shape {  
    float outline = 100;  
    int[][] offsets;  
    int[][] values;  
    float angle;  
  
    void setup() {  
        size(500, 500, P3D);  
  
        alpha1a = new int[values.length][height];  
        values = new int[values.length][height];  
        offset1a();  
  
        // Load the image onto a new array.  
        // Extract the values and store in an array.  
        a = loadImage("yellowH.jpg");  
        a.loadPixels();  
        for (int i = 0; i < a.height; i++) {  
            for (int j = 0; j < a.width; j++) {  
                offset1a[i][j] = a.pixels[i * a.width + j];  
                values[i][j] = int(blue(offset1a[i][j]));  
            }  
        }  
    }  
  
    void draw() {  
    }  
}  
  
Done Saving.  
15
```

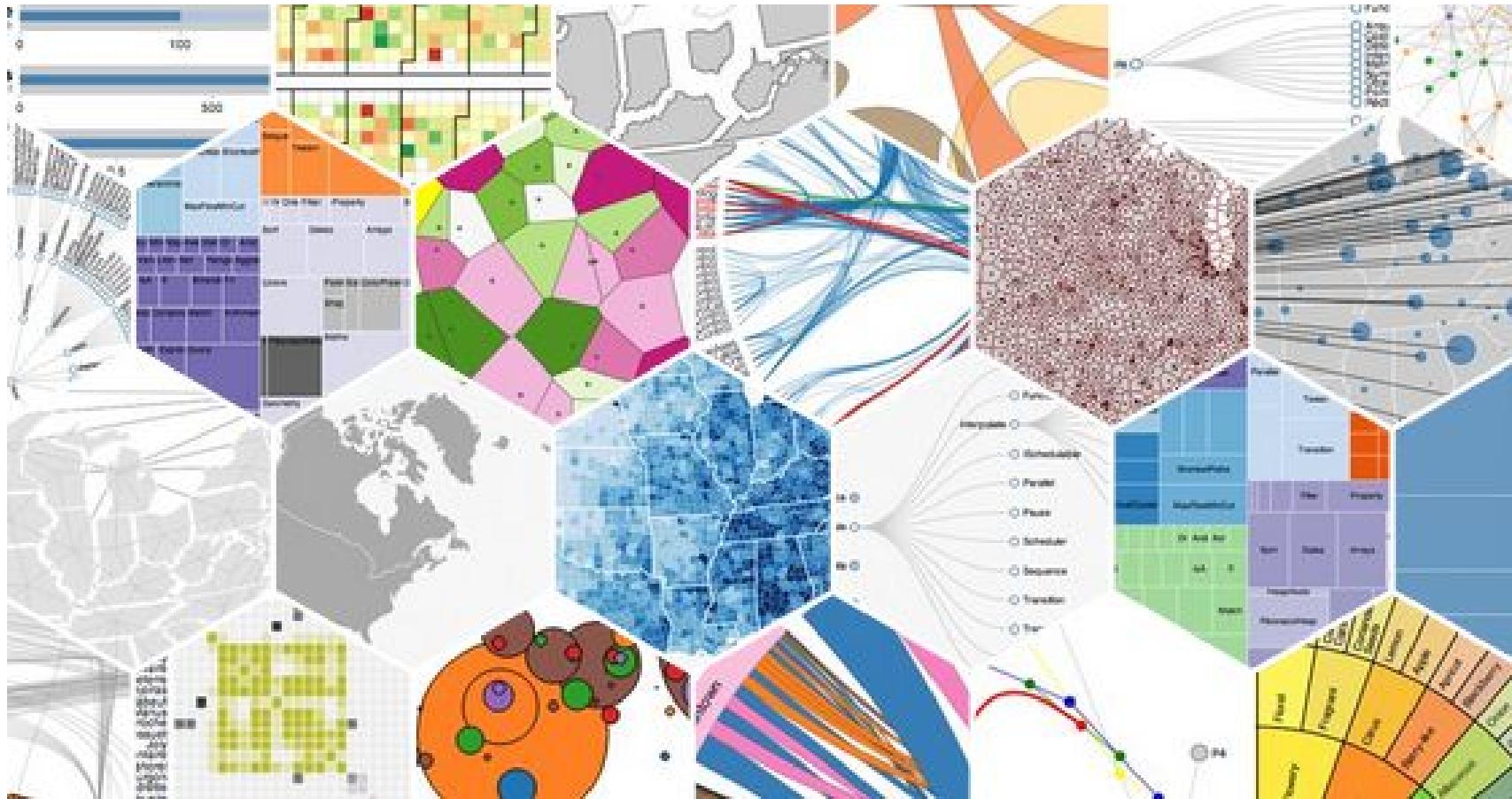
Toolbar  
Tabs

Text Editor

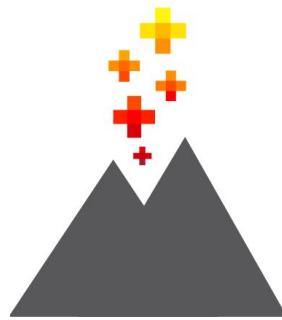
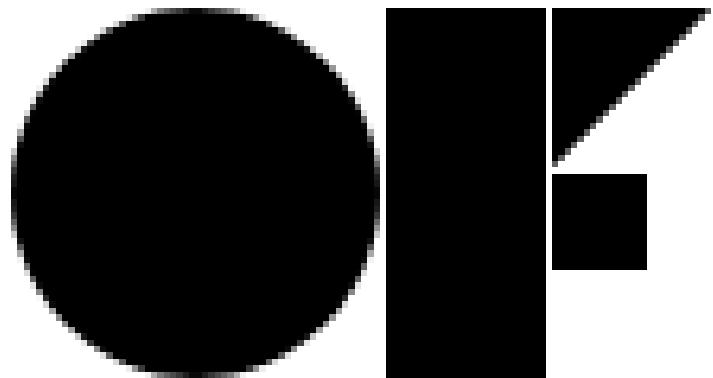
Message Area

Text Area

# Advanced: D3 d3js.org



## Advanced: C++ Libraries



## Advanced: Python

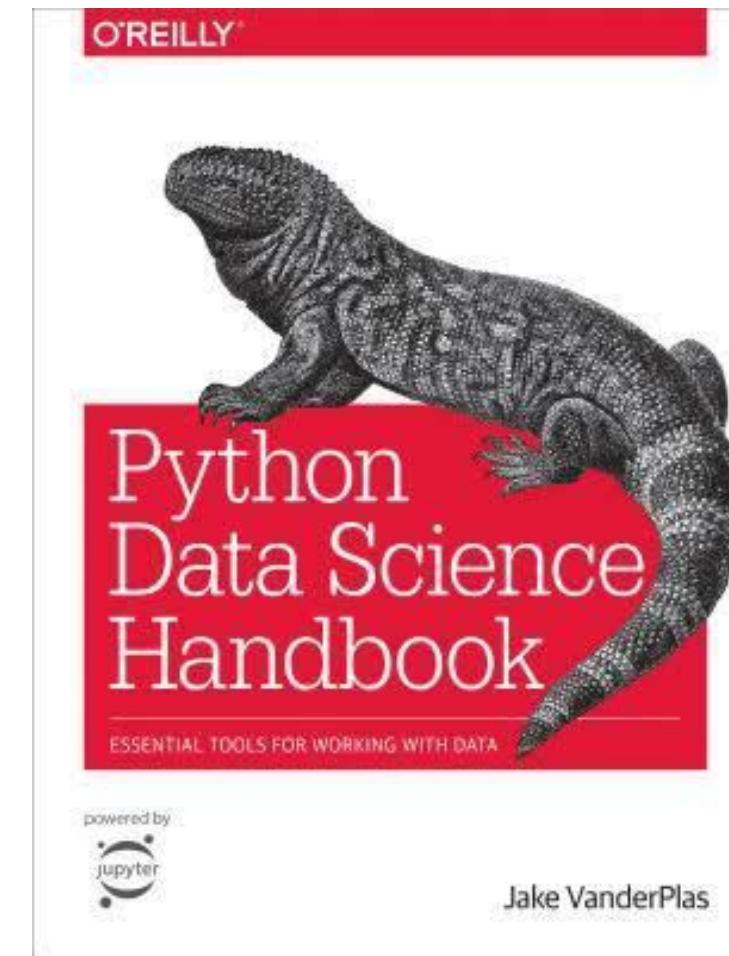
```
dict.py
Advanced: Python
{'a': 5, 'b': 7, 'A': 5, 't': 8, 'B': 7}
char_frequency = dict()
for key, value in char.items():
    if key.lower() in char_frequency:
        char_frequency[key.lower()] += value
    else:
        char_frequency[key.lower()] = value
print(char_frequency)
```

The Python logo, consisting of two interlocking snakes, is overlaid on the code. One snake is blue and the other is yellow, both with white eyes.

## Python resources

Python Data Science Handbook:  
Essential Tools for Working with Data

[Python Tutorial](#)





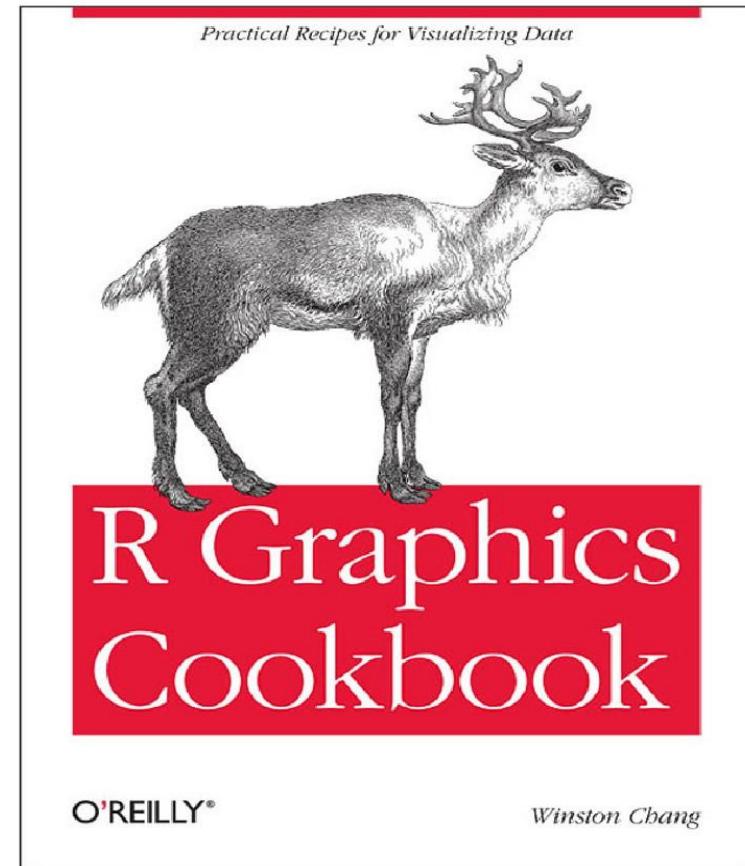
## R resources

First download R

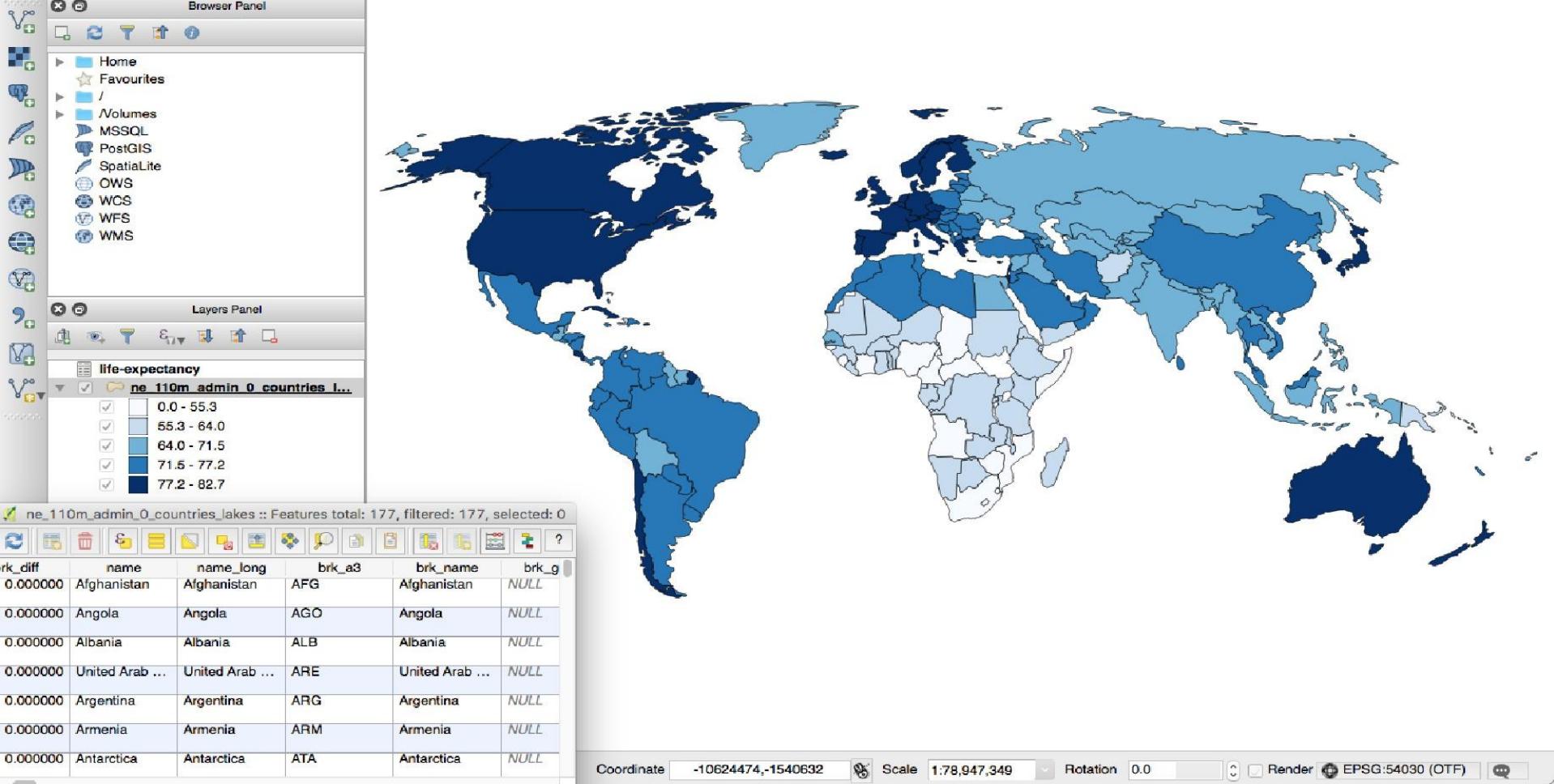
Then download R Studio

R Graphics Cookbook

[R Tutorial](#)



# Advanced: QGIS - Free and Open Source Geographic Information System



# QGIS resources

QGIS download

Tutorials

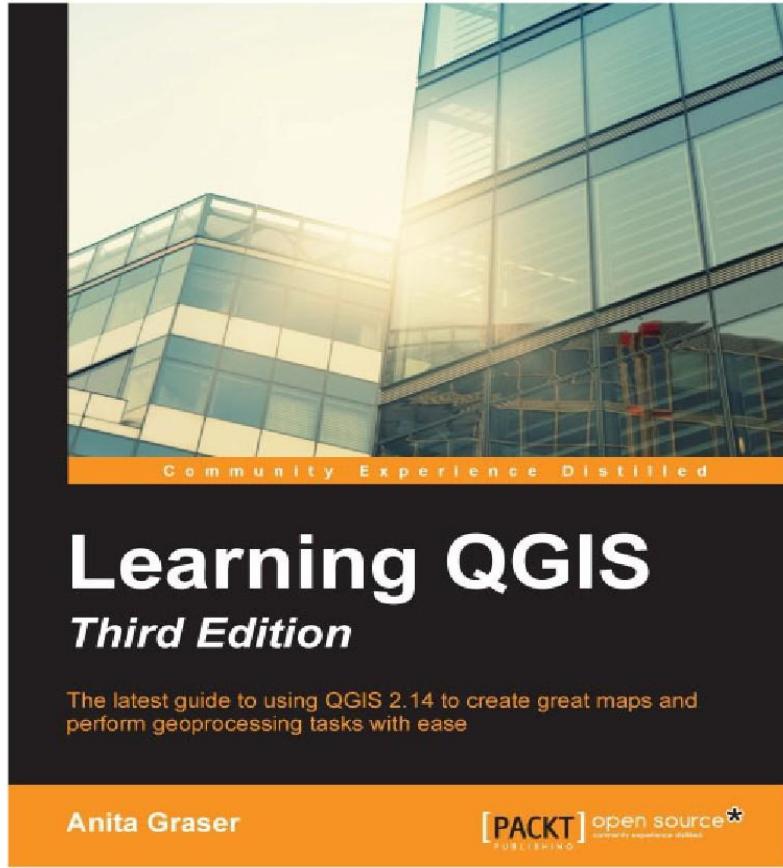
Learning QGIS book

OpenStreetMap

OpenStreetMap data extracts

Natural Earth

Trimble Data Marketplace



# Data Art vs. Data Visualization

# Thousands of Exhausted Things - The Office of Creative Research

[LINK](#)



# Serendipity - Kyle McDonald



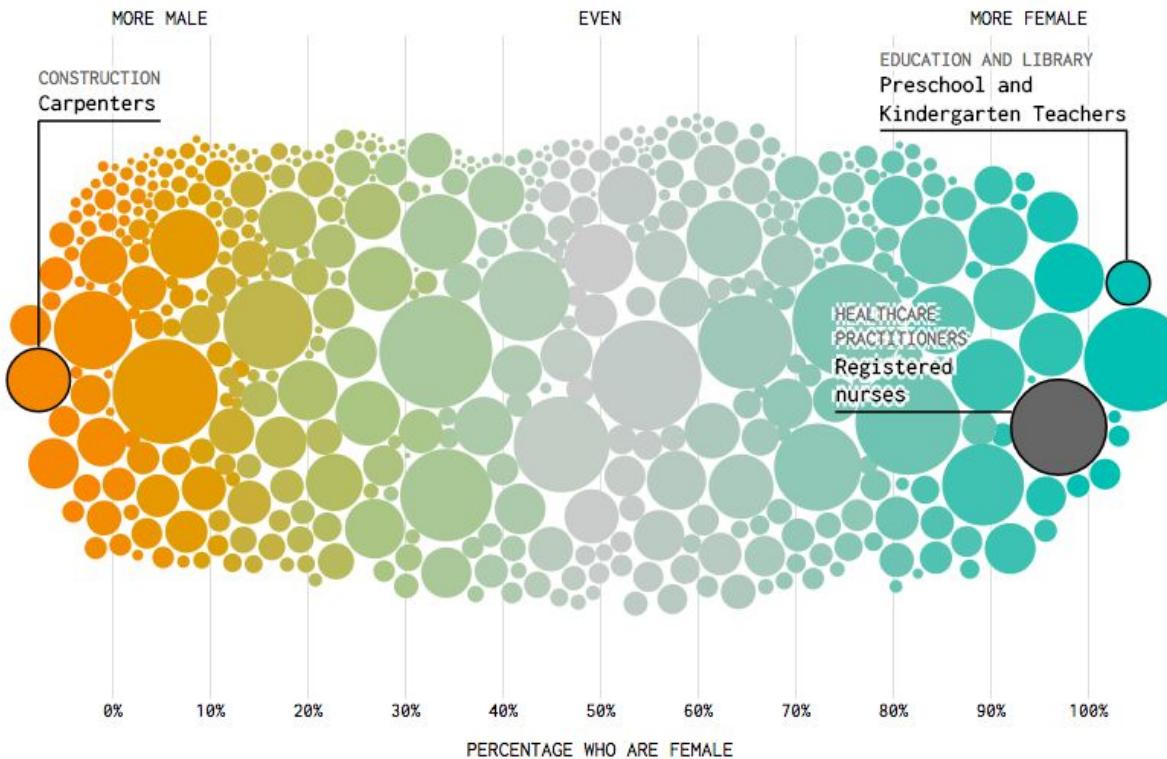
# I Want You to Want Me - Jonathan Harris



# Occupations - Nathan Yau

## MALE AND FEMALE OCCUPATIONS IN 2015

*Larger circles represent more common jobs.*



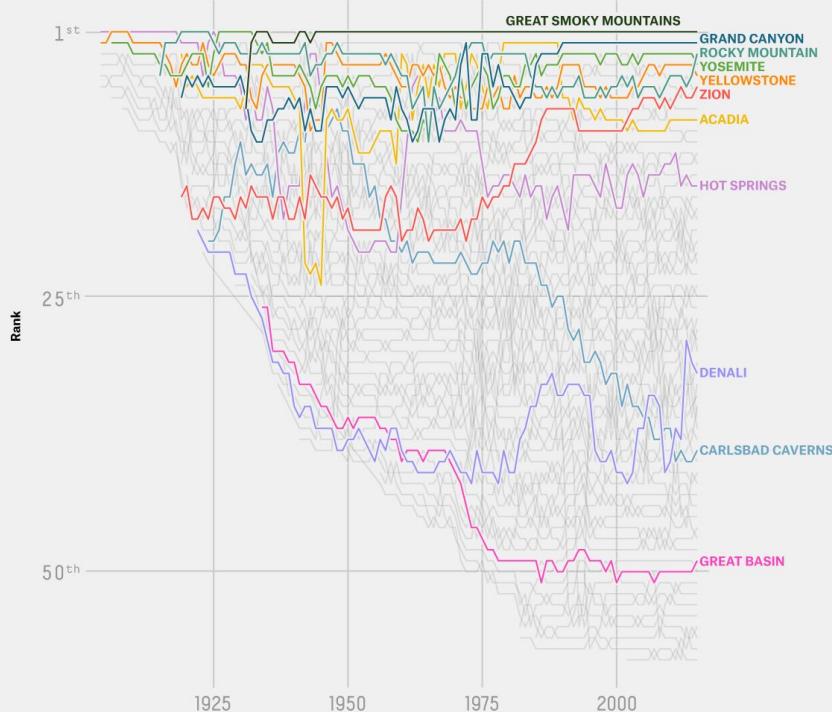
## Spies in the Skies - Buzzfeed



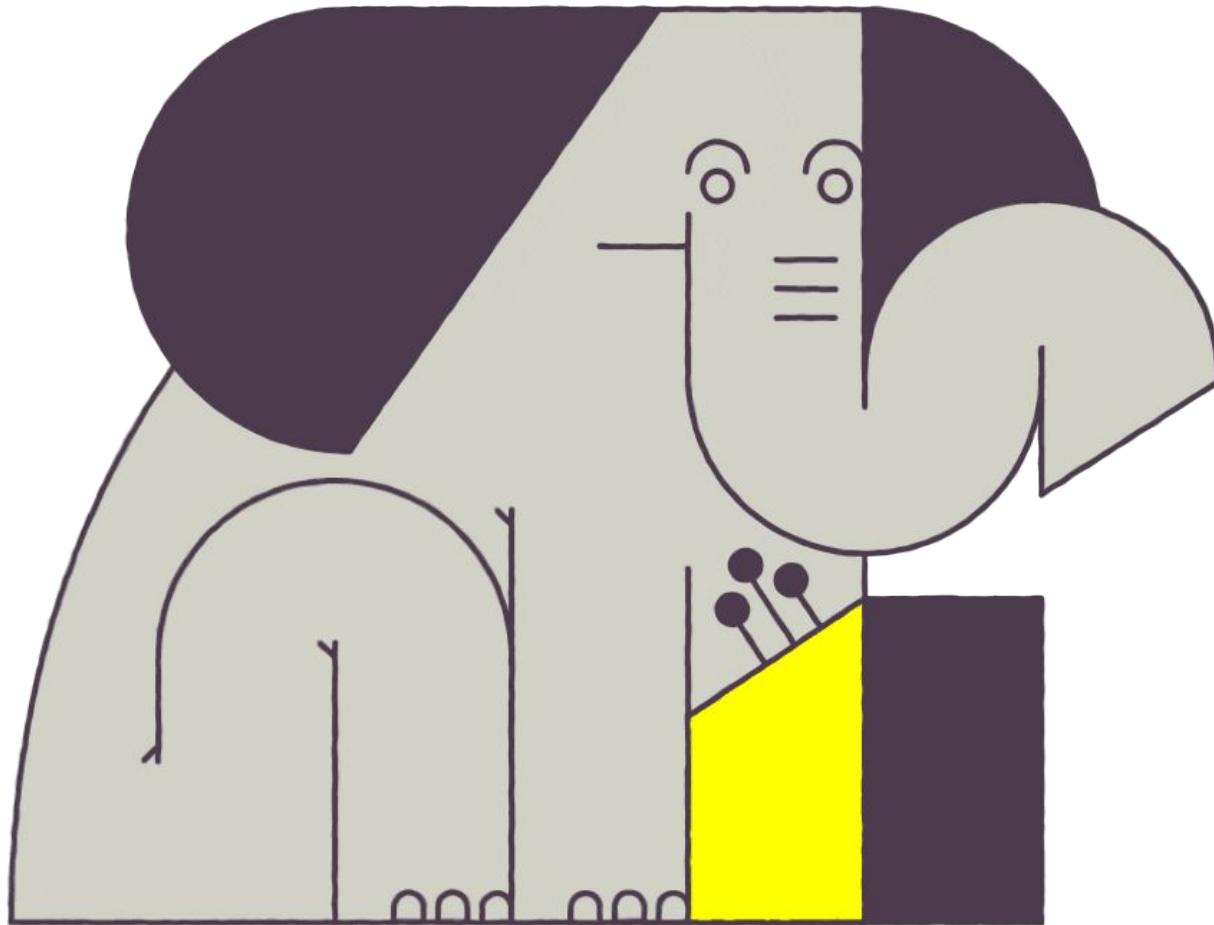
# National Parks-FiveThirtyEight

## The most popular national parks

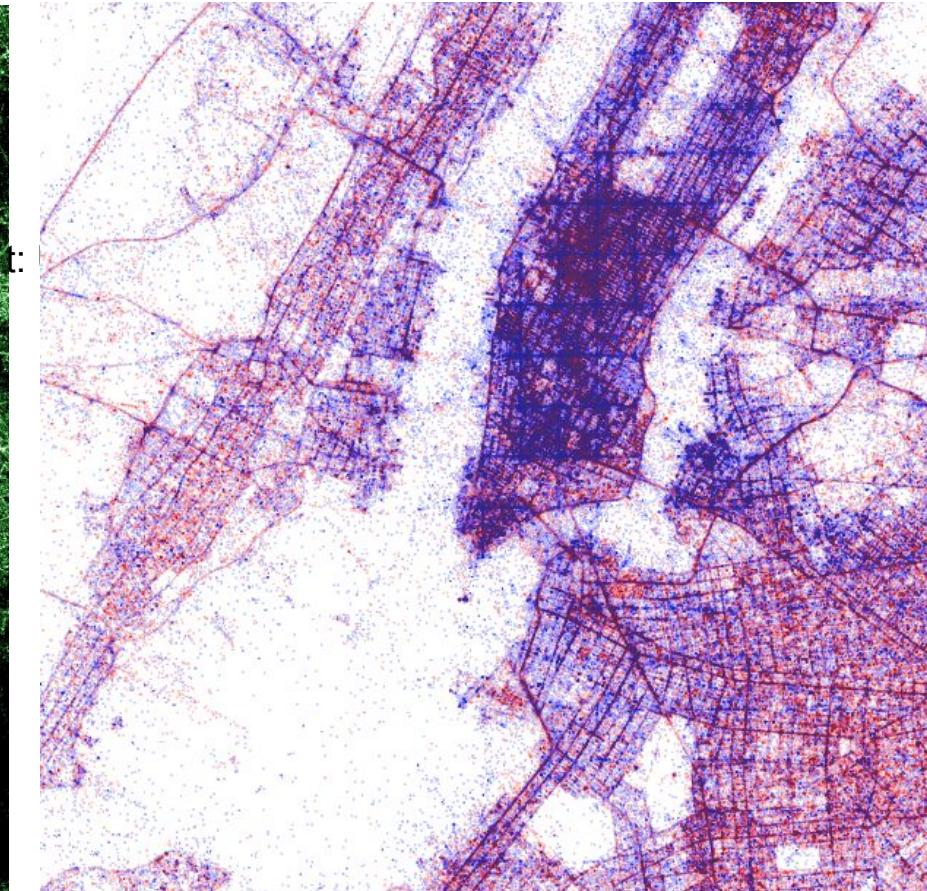
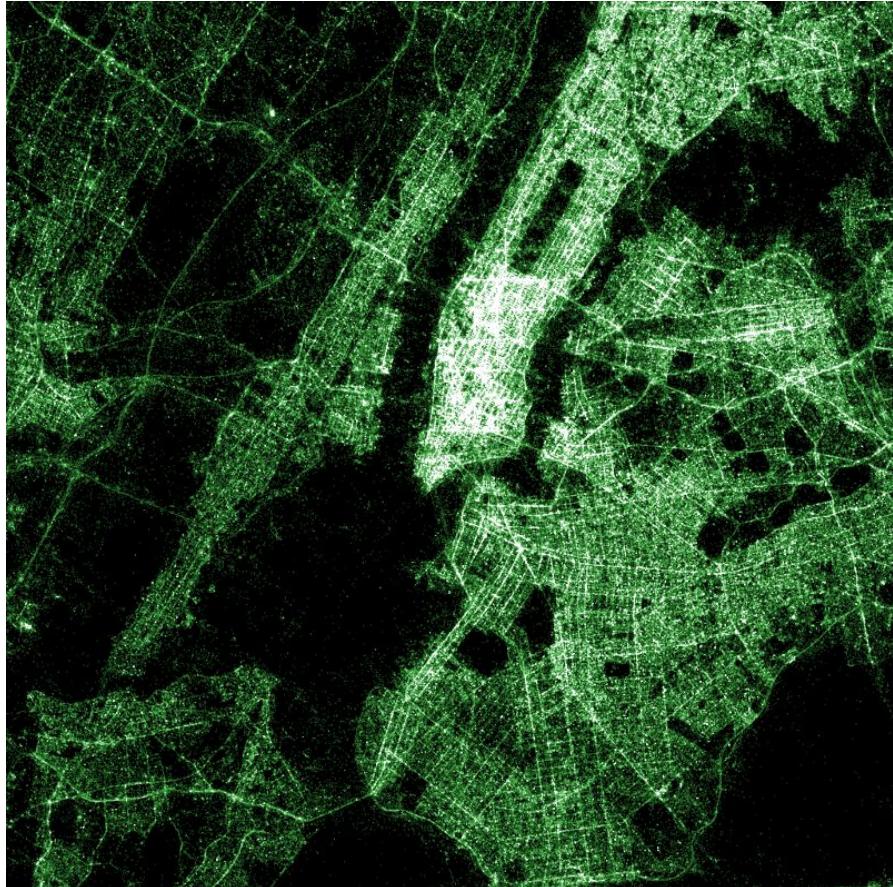
National parks ranked by number of visitors in a given year



## Republican Stump Speech-FiveThirtyEight



## Shazam and One Billion Recognitions-Umar Hansa



# Assignment

**Take your assigned data set and formulate the following:**

1. Before looking at the specifics of the data, develop a set of questions you want answered with this data set. You may find that the data doesn't quite provide enough of what you want to concern, in that case, modify your questions or gather extra data.
2. Your data should be clean, but in the case that it's not, tidy it according to your needs. Document this process.
3. Organize your data according to all that you know about organizational methods. Be prepared to defend this as they relate to the questions you want answered and the message to be relayed.
4. Use an easier charting software (Sheets with sparklines, google charts, Processing, or RawGraphs (or if you feel inclined, one of the advanced methods or languages) to visualize your data.

**We will present in class next Monday. BTW- your work can be more data art-but still must present all of the above.**