# What is data?

*GEOG 30323*

*August 25, 2015*
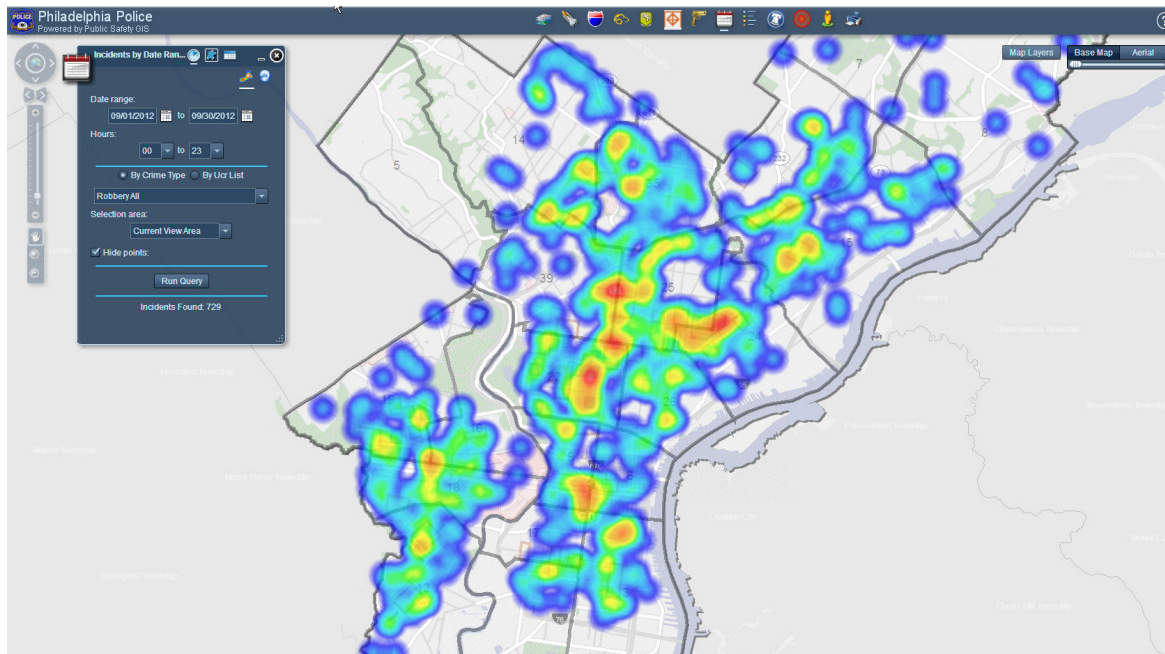
## Introductions

## My data journey



## Course goals

## Syllabus

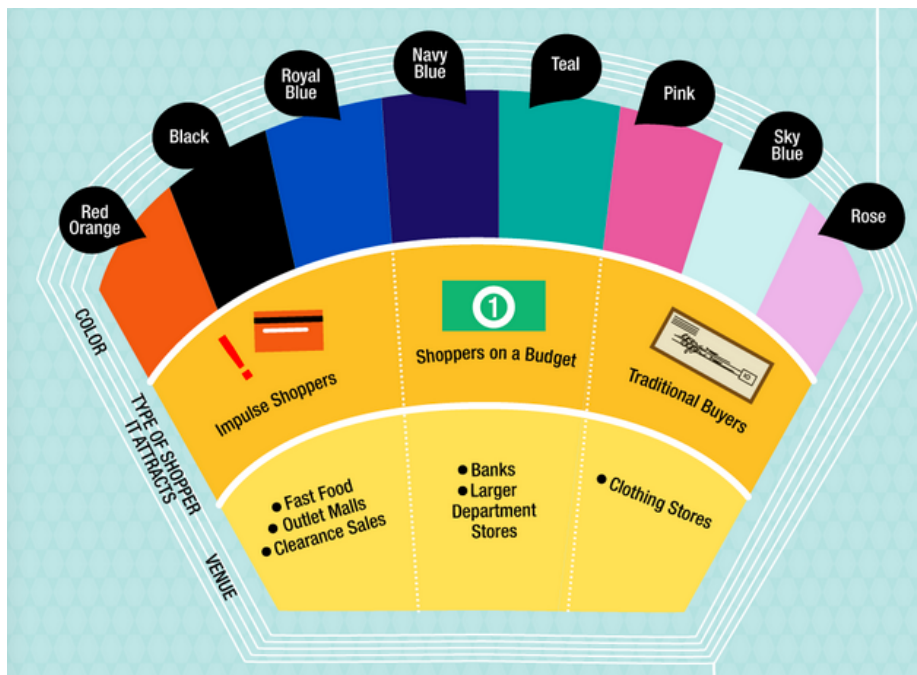## What is data?

Brainstorm: how might you use data if you were…

- …a librarian?
- …a police officer?
- …a small shoe store owner?
- …a general manager of a professional basketball franchise?
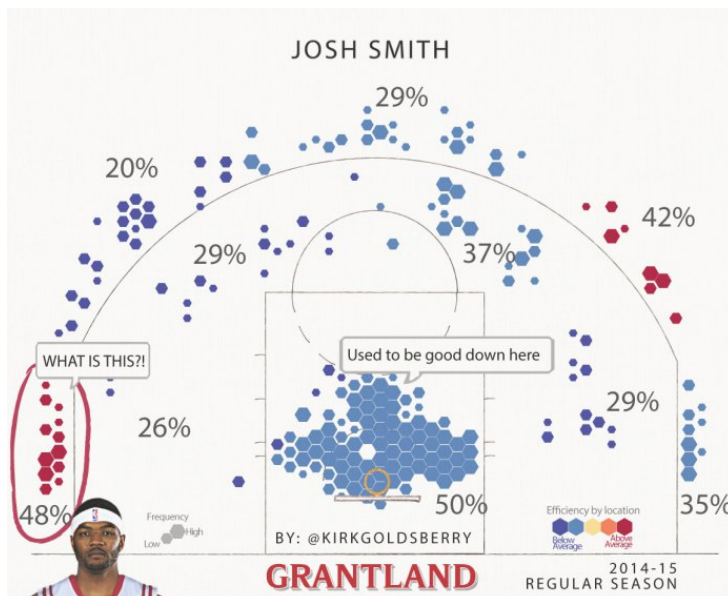
## Policing

Source: technical.ly (https://technical.ly/philly/2012/11/20/philadelphia-police-gis-new-system-crime-tracking/)

# Retail



Source: Humayun Khan/Shopify (http://www.shopify.com/blog/14254569-why-all-sale-signs-are-red-the-science-of-color-in-retail)

# Sports

Source: Kirk Goldsberry/*Grantland* (http://grantland.com/the-triangle/nba-least-efficient-shooters-josh-smith-michael-carter-williams-trey-burke-nerlens-noel-dion-waiters/)

# Data science



Source: *The Onion* (http://www.theonion.com/article/amazoncom-recommendations-understand-area-woman-be-2121)

## What you will learn in this course

This course covers the process of **exploratory data analysis**, which includes how to:

- Find data
- Load data
- Clean data
- Summarize data
- Visualize data
- Present data
- Example question: what are the most popular female baby names in Texas, and how has this changed over time?

# Finding data



Source: Social Security Administration (http://www.ssa.gov/oact/babynames/limits.html)

# Loading and cleaning data

```
TX - Notepad

File  Edit  Format  View  Help

TX,F,1910,Mary,895
TX,F,1910,Ruby,314
TX,F,1910,Annie,277
TX,F,1910,Willie,260
TX,F,1910,Ruth,252
TX,F,1910,Gladys,240
TX,F,1910,Maria,223
TX,F,1910,Frances,197
TX,F,1910,Margaret,194
TX,F,1910,Helen,189
TX,F,1910,Thelma,188
TX,F,1910,Mildred,186
TX,F,1910,Bessie,181
TX,F,1910,Lillian,180
TX,F,1910,Edna,178
TX,F,1910,Ethel,176
TX,F,1910,Lillie,170
TX,F,1910,Dorothy,167
TX,F,1910,Lucille,166
TX,F,1910,Minnie,164
TX,F,1910,Elizabeth,161
TX,F,1910,Hazel,151
TX,F,1910,Alice,149
TX,F,1910,Myrtle,149
TX,F,1910,Bertha,146
TX,F,1910,Opal,146
TX,F,1910,Irene,144
TX,F,1910,Emma,142
TX,F,1910,Marie,142
TX,F,1910,Mattie,142
TX,F,1910,Lois,140
TX,F,1910,Pauline,139
TX,F,1910,Juanita,133
TX,F,1910,Velma,133
TX,F,1910,Clara,129
```

# Summarizing data

|         | state | gender | year | name     | count | per1000    | rank |
|---------|-------|--------|------|----------|-------|------------|------|
| 177374  | TX    | F      | 2013 | Sophia   | 2381  | 14.905098  | 1    |
| 177375  | TX    | F      | 2013 | Emma     | 2043  | 12.789213  | 2    |
| 177376  | TX    | F      | 2013 | Isabella | 1941  | 12.150691  | 3    |
| 177377  | TX    | F      | 2013 | Mia      | 1830  | 11.455829  | 4    |
| 177378  | TX    | F      | 2013 | Olivia   | 1598  | 10.003506  | 5    |
| 177379  | TX    | F      | 2013 | Emily    | 1449  | 9.070763   | 6    |
| 177380  | TX    | F      | 2013 | Sofia    | 1380  | 8.638822   | 7    |
| 177381  | TX    | F      | 2013 | Abigail  | 1250  | 7.825020   | 8    |
| 177382  | TX    | F      | 2013 | Ava      | 1227  | 7.681040   | 9    |
| 177383  | TX    | F      | 2013 | Victoria | 1067  | 6.679437   | 10   |

# Visualizing data

## Female baby names per 1000 in the SSA data (TX)



# Presenting data

# Tools of the trade

- Your main tool: the Python programming language
- Guest appearances from: CartoDB, Tableau

- More to come on Python next week!

---

## Assignment 1: Data journalism

- **Data journalism**: movement in journalism toward interactive, data-driven content
- Popular publications: The Upshot (*New York Times*) (http://www.nytimes.com/upshot/?_r=0), FiveThirtyEight (http://fivethirtyeight.com/), *The Guardian* (http://www.theguardian.com/data), ProPublica (https://www.propublica.org/), many more

Your task:

- Find an example of data journalism from the past two years in a recognized publication
- Write a response that 1) summarizes the article, 2) discusses its data sources, and 3) discusses its methods of analysis and visualization.
- At least one paragraph, no longer than one page