


An End to Boring Data

With Visualizations in Python

HEATHER SHAPIRO
TECHNICAL EVANGELIST, MICROSOFT
[@microheather](#)



What will we cover?

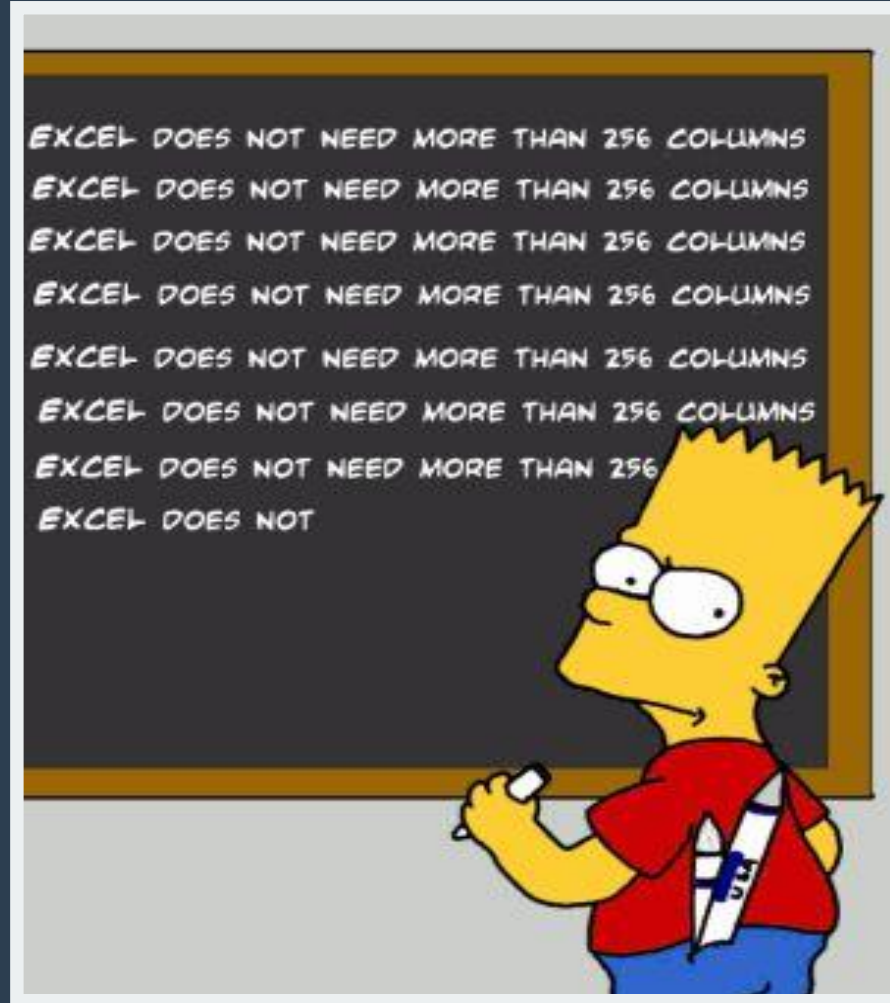
- Why data visualizations are important
- Case study on NYC Restaurant Ratings
- What libraries in python work best for different types of graphs

Why we need data viz

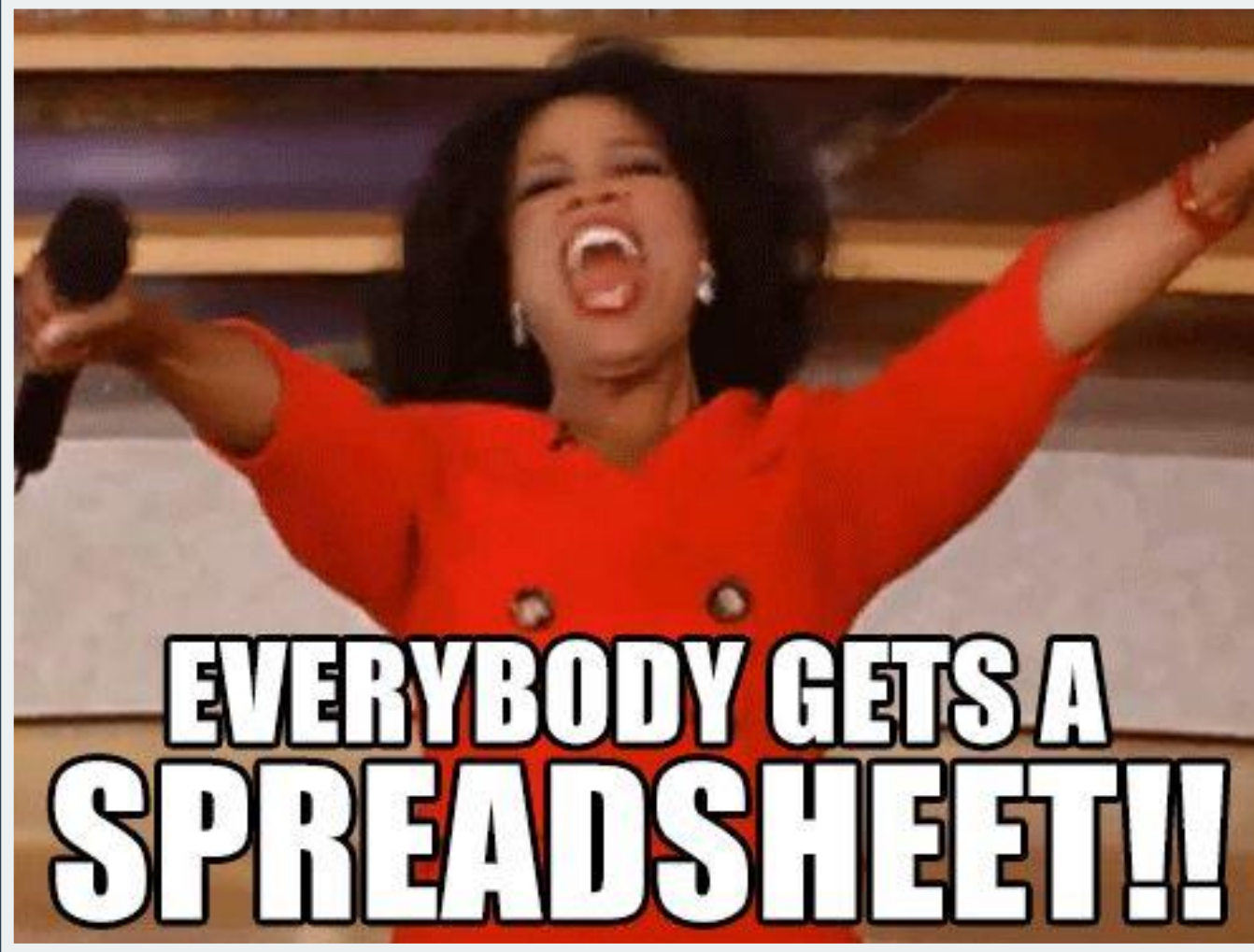
HARD TO UNDERSTAND



TOO MANY NUMBERS



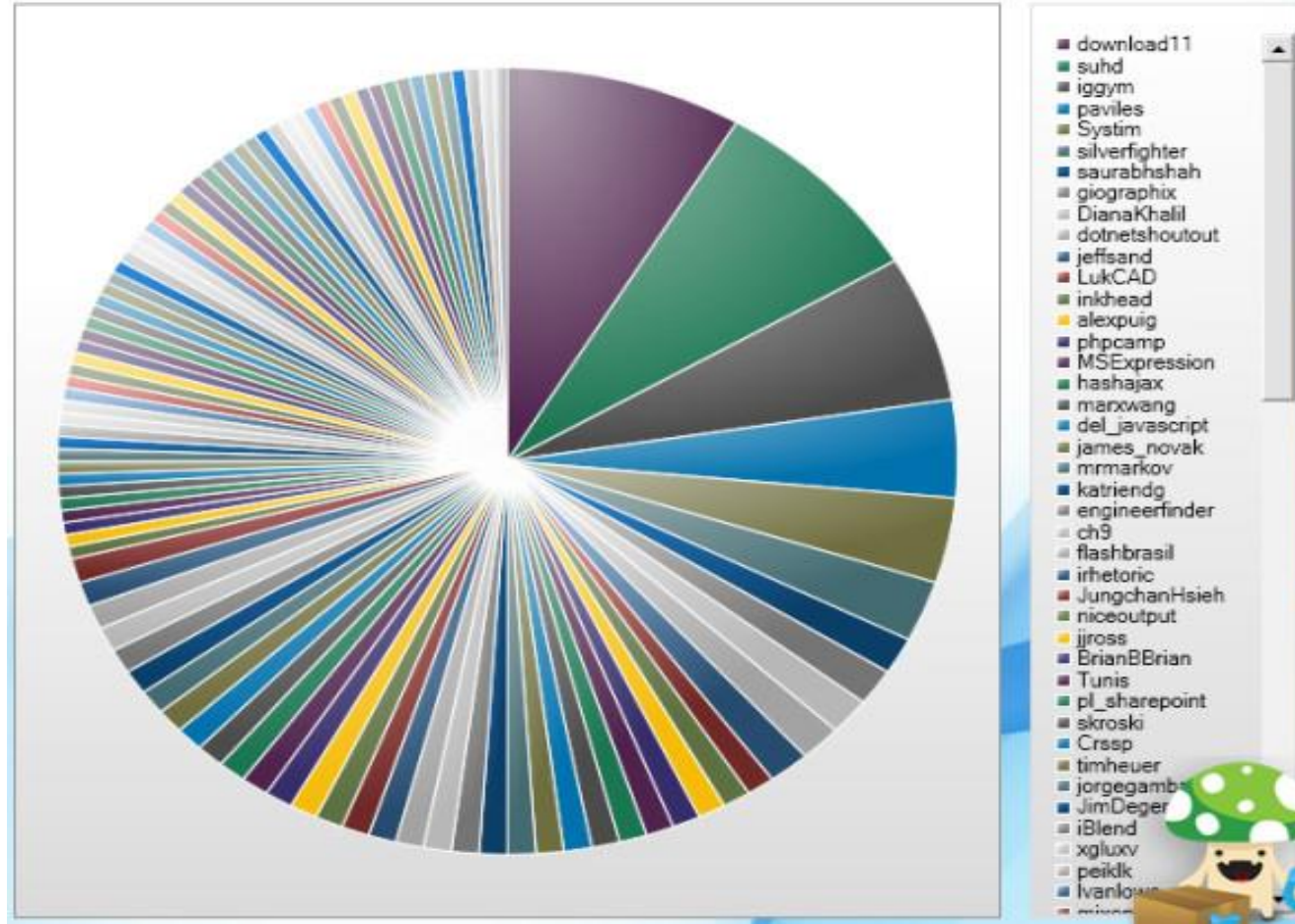
BORING MEETINGS



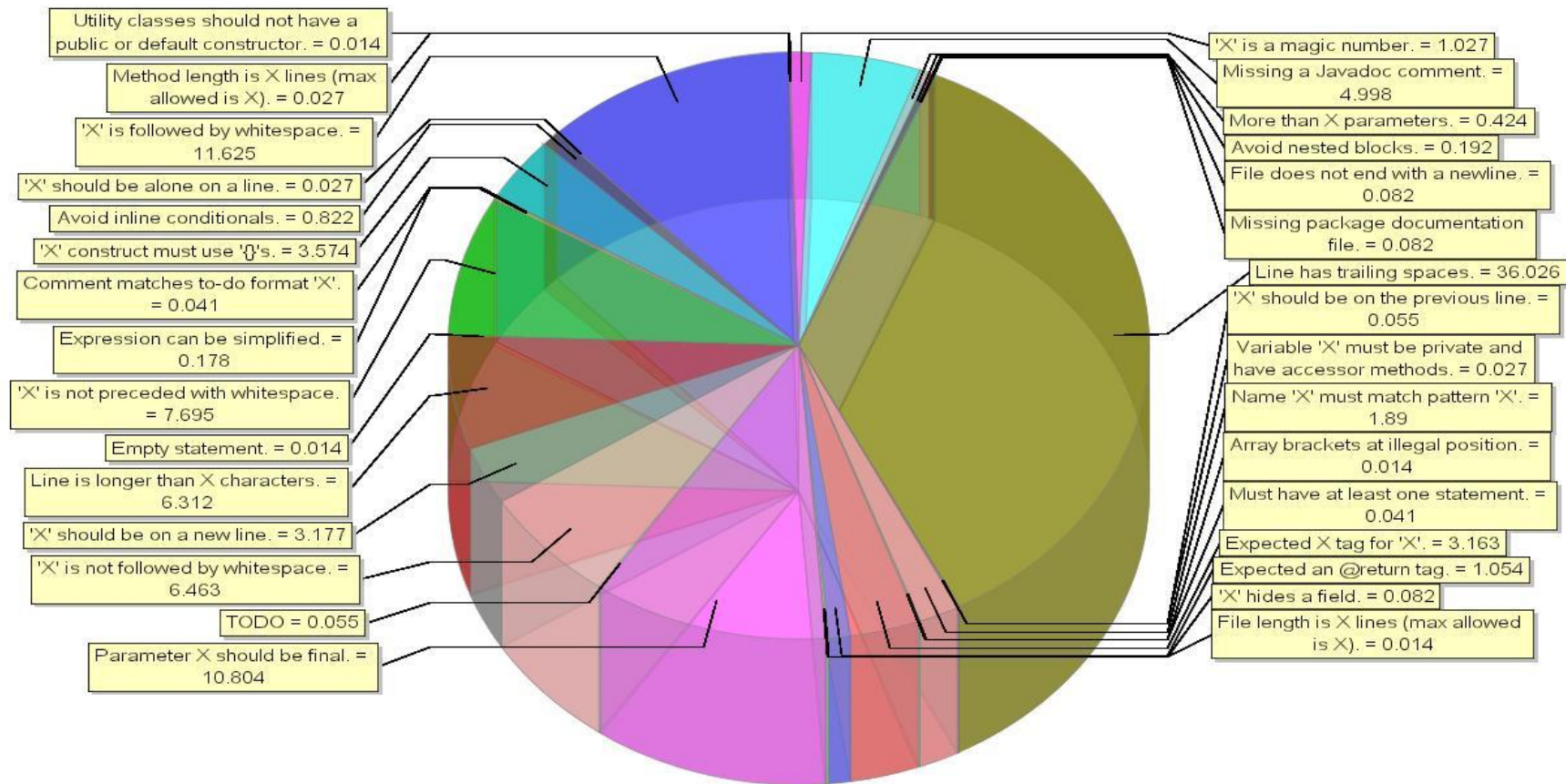
Visualizations gone wrong

TOO MANY VARIABLES

100 Most Active Tweeters



...?/



What data viz provides

- Helps the visual learner
- Makes sense of tremendous amounts of data
- Helps walk through a problem
- Tells a story in seconds



CASE STUDY

NYC Restaurant Ratings





GIFBG.com

HEATHER SHAPIRO | TECHNICAL EVANGELIST, MICROSOFT
@microheather



NYC OPEN DATA

The screenshot shows the NYC OpenData website. At the top, the logo 'NYC OpenData' is on the left, followed by a yellow badge that says '1500+ Data Sets Available'. To the right are social media icons for GitHub, a wrench, a question mark, NYC, Tumblr, and Twitter, along with 'Sign Up' and 'Sign In' links. The main banner features a blurred image of a crowd and the title '2015 Open Data Plan Update'. Below the title is a paragraph: 'On July 15, Mayor Bill de Blasio released Open Data for All, the annual update to the NYC Open Data Plan. As required by Local Law 11 of 2012, each City entity must identify and ultimately publish all of its digital public data for citywide aggregation and publication by 2018. [Click here to read Open Data for All.](#)' Below this are 'View' and 'More Stories' links. A search bar with the placeholder 'Search' is centered below the banner. Underneath the search bar is the text 'Click here to view the NYC OpenData dashboard'. At the bottom, there are eight icons arranged in two rows of four, each with a label: Business (briefcase), City Government (capitol dome), Education (graduation cap), Environment (leaf), Health (heart with pulse line), Housing (house), Public Safety (police cap), and Transportation (goggles).

NYC OpenData 1500+ Data Sets Available

2015 Open Data Plan Update

On July 15, Mayor Bill de Blasio released Open Data for All, the annual update to the NYC Open Data Plan. As required by Local Law 11 of 2012, each City entity must identify and ultimately publish all of its digital public data for citywide aggregation and publication by 2018. [Click here to read Open Data for All.](#)

View More Stories

Search

Click here to view the NYC OpenData dashboard

Business City Government Education Environment

Health Housing Public Safety Transportation

Dataset (File)	Data Field Name	Data Type	Length	Expected Values	Description
WEBEXTRACT	CAMIS	Varchar	10		This is a unique identifier for the entity (restaurant)
WEBEXTRACT	DBA	varchar	255		This field represents the name (doing business as) of the entity (restaurant)
				<ul style="list-style-type: none"> • 1 = MANHATTAN • 2 = BRONX • 3 = BROOKLYN • 4 = QUEENS • 5 = STATEN ISLAND • Missing 	
WEBEXTRACT	BORO	Varchar	1		Borough in which the entity (restaurant) is located. NOTE: There may be discrepancies between zip code and listed boro due to differences in an establishment's mailing address and physical location
WEBEXTRACT	BUILDING	Varchar	10		This field represents the building number for the entity (restaurant)
WEBEXTRACT	STREET	Varchar	100		This field represents the street name at which the entity (restaurant) is located.
WEBEXTRACT	ZIPCODE	Varchar	5		Zip code as per the address of the entity (restaurant)
WEBEXTRACT	PHONE	Varchar	20		Phone number
	CUISINE				
WEBEXTRACT	DESCRIPTION	Varchar	200		This field describes the entity (restaurant) cuisine.
					This field represents the date of inspection. NOTE: Inspection dates of 1/1/1900 mean an establishment has not yet had an inspection
WEBEXTRACT	INSPECTION DATE	Datetime	N/A		
				<ul style="list-style-type: none"> • Violations were cited in the following area(s). • No violations were recorded at the time of this inspection. • Establishment re-opened by DOHMH • Establishment re-closed by DOHMH • Establishment Closed by DOHMH. Violations were cited in the following area(s) and those requiring immediate action were addressed. • "Missing" = not yet inspected 	
WEBEXTRACT	ACTION	Varchar	150		This field represents the action that is associated with each restaurant inspection.
WEBEXTRACT	VIOLATION CODE	Varchar	3		This field represents each violation associated with a restaurant inspection.
	VIOLATION				
WEBEXTRACT	DESCRIPTION	Varchar	600		This field describes the violation codes
				<ul style="list-style-type: none"> • Critical • Not Critical • Not Applicable 	
WEBEXTRACT	CRITICAL FLAG	Varchar	1		Critical violations are those most likely to contribute to foodborne illness.
WEBEXTRACT	SCORE	Varchar	3		Total score for a particular inspection; updated based on adjudication results.
				<ul style="list-style-type: none"> • Not Yet Graded • A = Grade A • B = Grade B • C = Grade C • Z = Grade Pending • P=Grade Pending issued on re-opening following an initial inspection that resulted in a closure 	
WEBEXTRACT	GRADE	Varchar	1		This field represents the grade associated with this inspection. Grades given during a reopening inspection derived from the previous re-inspection.
WEBEXTRACT	GRADE DATE	Datetime	N/A		The date when the grade was issued to the entity (restaurant)
WEBEXTRACT	RECORD DATE	Datetime	N/A		The date when the webextract was run to produce this data set
				•Calorie Posting/ Compliance Inspection	

Steps Taken

- Load modules
- Load Restaurant Rating Data
- Understand the data
- Visualize 😊



TOOLS USED

- Pandas
- Matplotlib
- Basemap
- Folium
- Seaborn
- Bokeh
- Plot.ly



GETTING STARTED

HEATHER SHAPIRO | TECHNICAL EVANGELIST, MICROSOFT
@microheather

PANDAS



PANDAS

Import pandas as pd

- Python library to provide data analysis features
- Built on *NumPy*, *SciPy*, and *matplotlib*
- Key components
 - Series
 - DataFrames

LIBRARIES FOR STATISTICAL GRAPHS

HEATHER SHAPIRO | TECHNICAL EVANGELIST, MICROSOFT
[@microheather](#)



MATPLOTLIB

import matplotlib.pyplot as plt

- MATLAB-like plotting framework

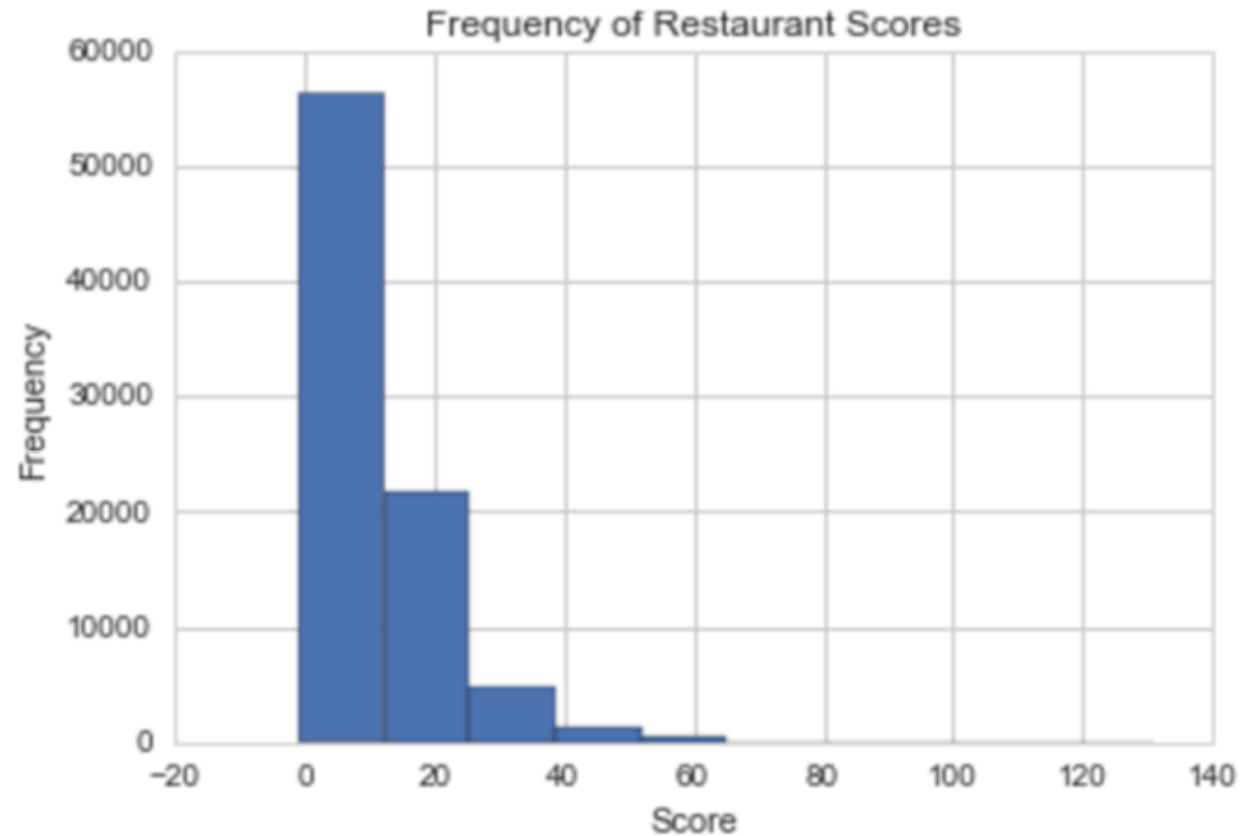
MATPLOTLIB

```
f, ax = plt.subplots() ## creates figure area with axes  
# histogram our data with numpy  
data = mRests['SCORE']
```

```
plt.hist(data)  
plt.xlabel('Score')  
plt.ylabel('Frequency')  
plt.title("Frequency of Restaurant Scores")
```

```
plt.show()
```

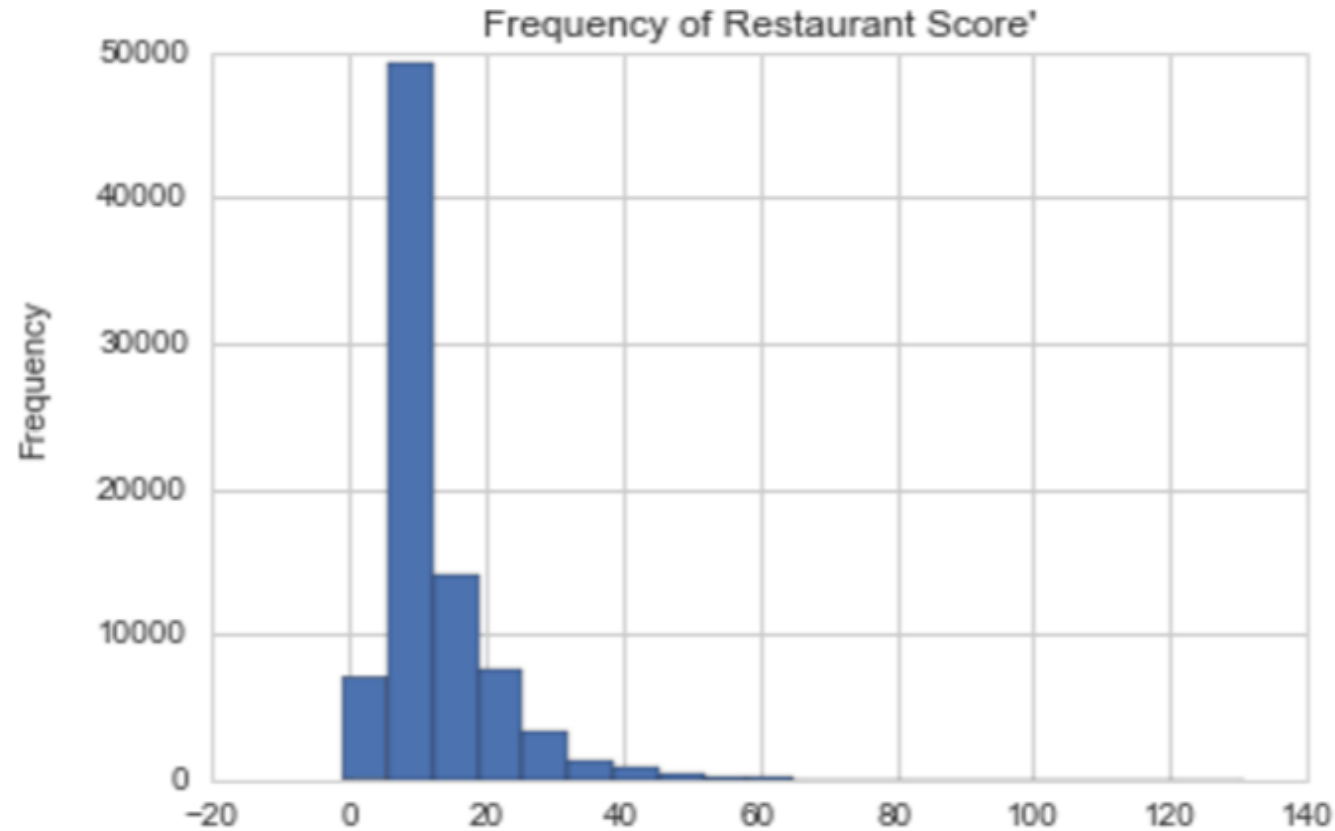
MATPLOTLIB



PANDAS SHORTHAND

```
mRests["SCORE"].hist(bins=20)  
plt.title("Frequency of Restaurant Score")
```

PANDAS SHORTHAND



seaborn

SEABORN

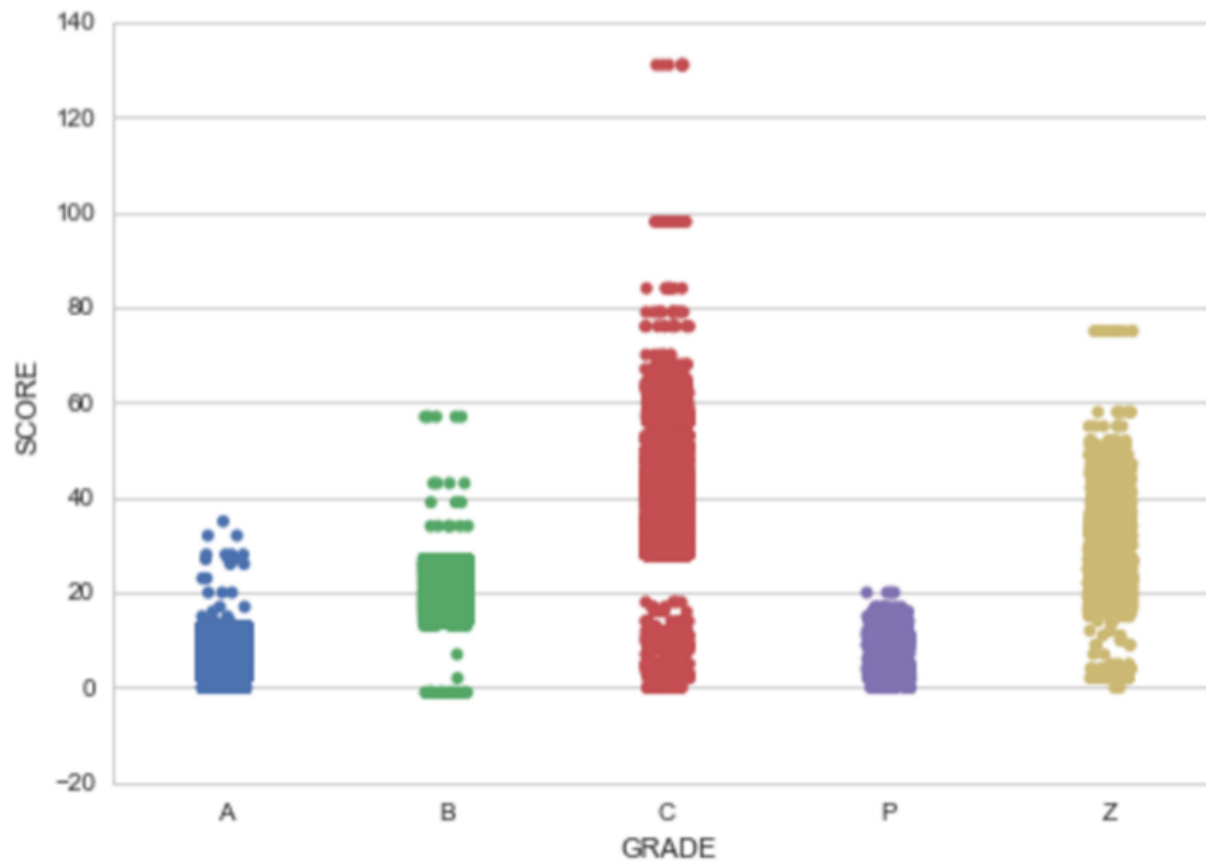
Import seaborn

- Built on top of matplotlib
- Creates more sophisticated graphs
- Look more professional

SEABORN

```
sns.stripplot(x="GRADE", y = "SCORE", data =  
mRests, jitter = True)
```

SEABORN



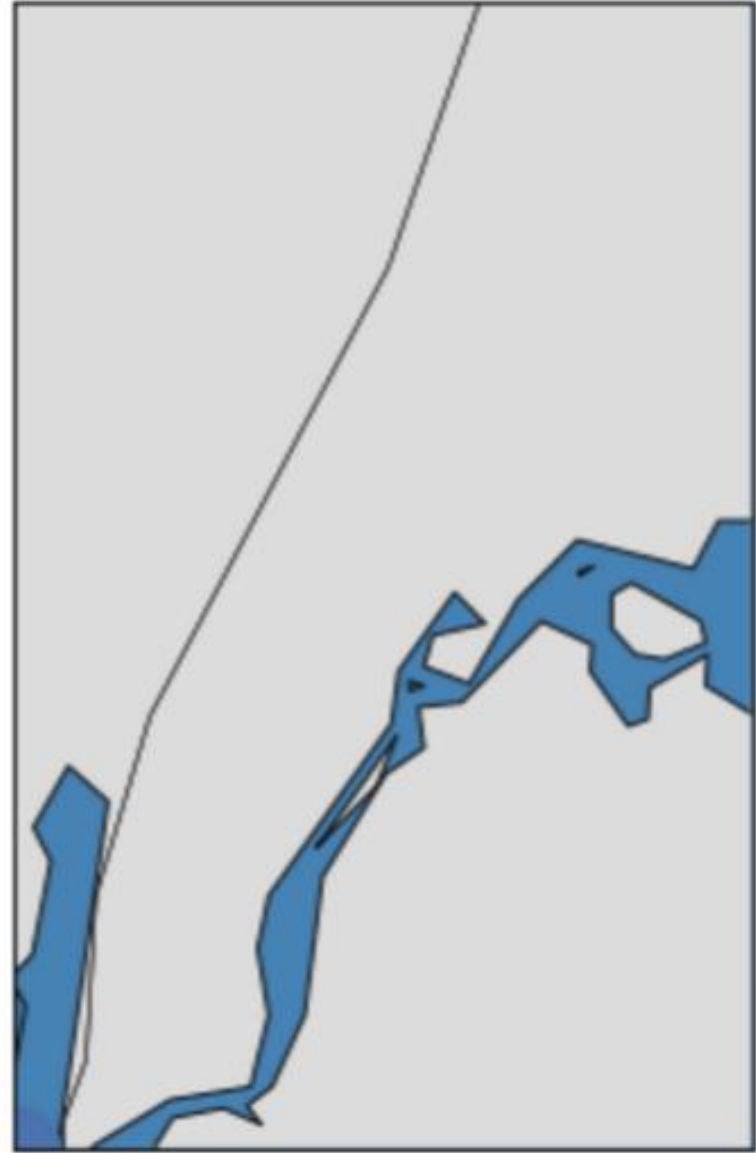
LIBRARIES FOR MAPPING

HEATHER SHAPIRO | TECHNICAL EVANGELIST, MICROSOFT
@microheather

BASEMAP

- Hard to install. There are a lot of prereqs and the documentation isn't there for windows

BASEMAP



FOLIUM

Import folium

- Visualize data on a Leaflet map
- Built-in tilesets from:
 - OpenStreetMap, MapQuest Open, MapQuest Open Aerial, Mapbox, and Stamen, and supports custom tilesets with Mapbox or Cloudmade API keys.

LIBRARIES FOR INTERACTIVE PLOTS

HEATHER SHAPIRO | TECHNICAL EVANGELIST, MICROSOFT
[@microheather](#)



bokeh



**GO PUT THEM
IN YOUR BLOGS 😊**

CLOSING THOUGHTS

- **Pandas** → handy for simple plots but you need to be willing to learn matplotlib to customize.
- **Seaborn** → supports more complex visualization approaches but still requires matplotlib. The color schemes are a nice bonus.

CLOSING THOUGHTS

- **Basemap** → Hard to install. Not very robust and there is not high granularity for the maps.
- **Folium** → Great documentation for mapping. Wish you could add more interactive widgets.

CLOSING THOUGHTS

- **bokeh** → Overkill for simple scenarios and documentation was not great.
- **Plotly** → most interactive graphs. You can save them offline and create rich web-based visualizations for your blog. Not good with city level data for maps.

What did we cover?

- Introduction to data visualizations in python
- How to walk through a data problem
- Which libraries are useful and for what
- Great way to update that blog

RESOURCES

- My blog 😊 www.Microheather.com
- My github: www.github.com/heatherbshapiro
- NYC Open Data: nycopendata.socrata.com
- Data Sets: data.gov
- Data Science VM in Azure : aka.ms/datasciencevm
- Azure Machine Learning: studio.azureml.net
- Channel9 and MVA

CONTACT ME

- Email: hshapiro@Microsoft.com
- Twitter [@microheather](https://twitter.com/microheather)

THANKS!

A young girl with dark hair, wearing a grey and white striped long-sleeved shirt and pink pants, is riding a bicycle on a paved path. The bicycle has a green frame and black handlebars. The background is slightly blurred, showing a paved path and some greenery. The word "THANKS!" is overlaid in large, bold, blue capital letters in the top left corner.

HEATHER SHAPIRO | TECHNICAL EVANGELIST, MICROSOFT
@microheather