



**matplotlib**  
Version 2.2.2

**Seaborn**



Bokeh

**CHULA ENGINEERING**  
Foundation toward Innovation

**COMPUTER**

**Chula Big Data and IoT**  
Center of Excellence  
(CUBIC)

**Data ind**  
Data Mining Group  
Machine Intelligence and Knowledge Discovery Lab  
Chulalongkorn University

# Data Visualizations in Python

## Python for Data Analytics

**Peerapon Vateekul, Ph.D.**

Department of Computer Engineering,  
Faculty of Engineering, Chulalongkorn University

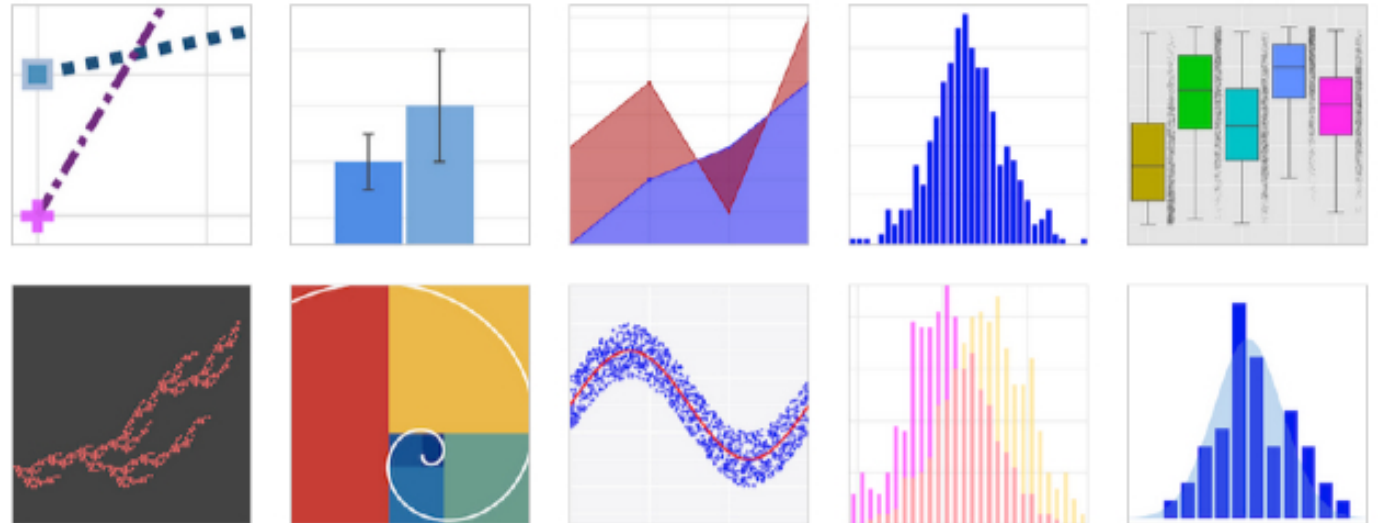
[Peerapon.v@chula.ac.th](mailto:Peerapon.v@chula.ac.th)



# Outlines

2

- Overview
- Exploratory Data Analysis (EDA)
- Scatter Plots
- Line Plots
- Histograms Plots
- Bar Plots
- Box Plots



## Reference:

1. <https://towardsdatascience.com/5-quick-and-easy-data-visualizations-in-python-with-code-a2284bae952f>



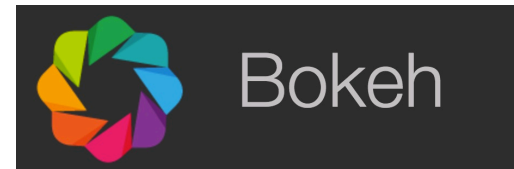
# Overview

3

- **Matplotlib:** Python 2D plotting library
- **Seaborn:** Statistical data visualization (On top; Matplotlib)
- **Bokeh:** an interactive visualization library that targets modern web browsers for presentation.

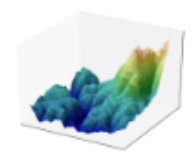
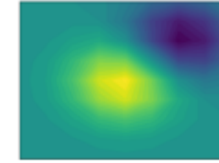
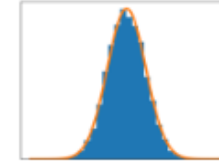
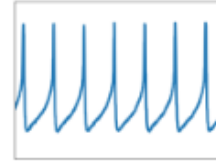


Seaborn





# Overview - Matplotlib



4

- Matplotlib is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms.
- Matplotlib can be used in Python scripts, the Python and [IPython](#) shells, the [Jupyter](#) notebook, web application servers, and four graphical user interface toolkits.
- Matplotlib tries to make easy things easy and hard things possible.
- You can generate plots, histograms, power spectra, bar charts, errorcharts, scatterplots, etc., with just a few lines of code.
  - For examples, see the [sample plots](#) and [thumbnail gallery](#).
- For simple plotting the pyplot module provides a MATLAB-like interface, particularly when combined with IPython.
- For the power user, you have full control of line styles, font properties, axes properties, etc, via an object oriented interface or via a set of functions familiar to MATLAB users.



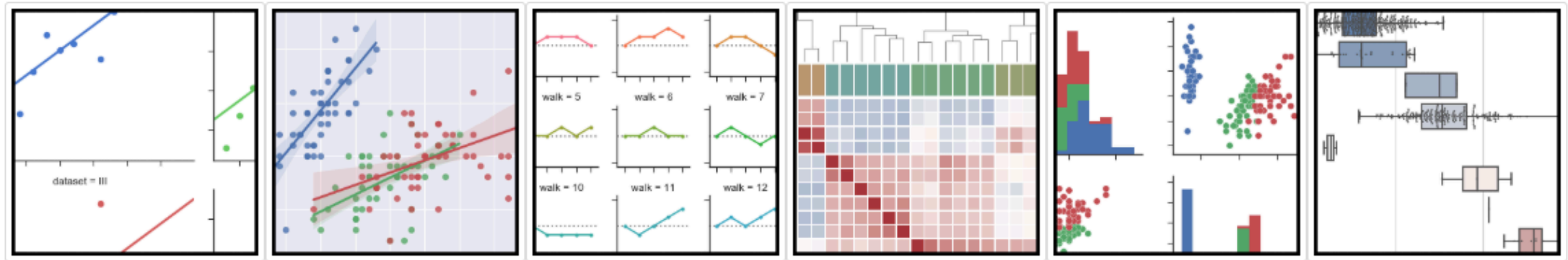


# Overview - Seaborn

# Seaborn

5

- Seaborn is a Python visualization library based on matplotlib.
- It provides a high-level interface for drawing attractive statistical graphics.
- On top; Matplotlib library.



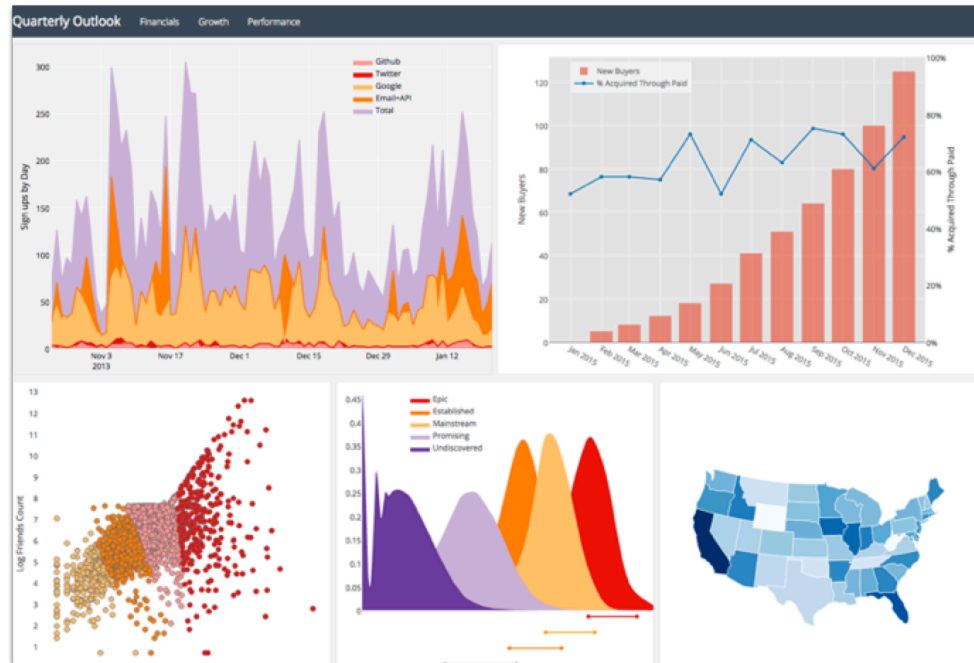


# Overview - Plotly



6

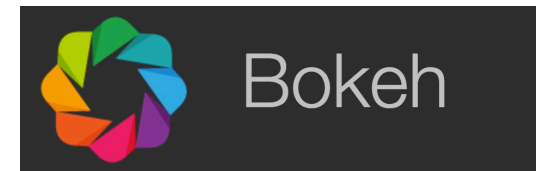
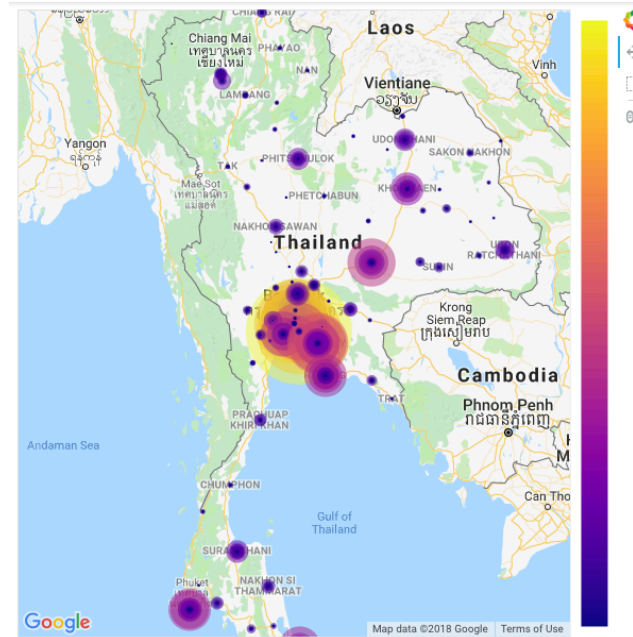
- Plotly creates leading open source tools for composing, editing, and sharing **interactive** data visualization via the Web.
- Our collaboration servers (available in cloud or on premises) allow data scientists to showcase their work, make graphs without coding, and collaborate with business analysts, designers, executives, and clients.





# Overview - Bokeh

- Bokeh is an interactive visualization library that targets modern web browsers for presentation.
- Its goal is to provide elegant, concise construction of versatile graphics, and to extend this capability with high-performance interactivity over very large or streaming datasets.
- Bokeh can help anyone who would like to quickly and easily create interactive plots, dashboards, and data applications.





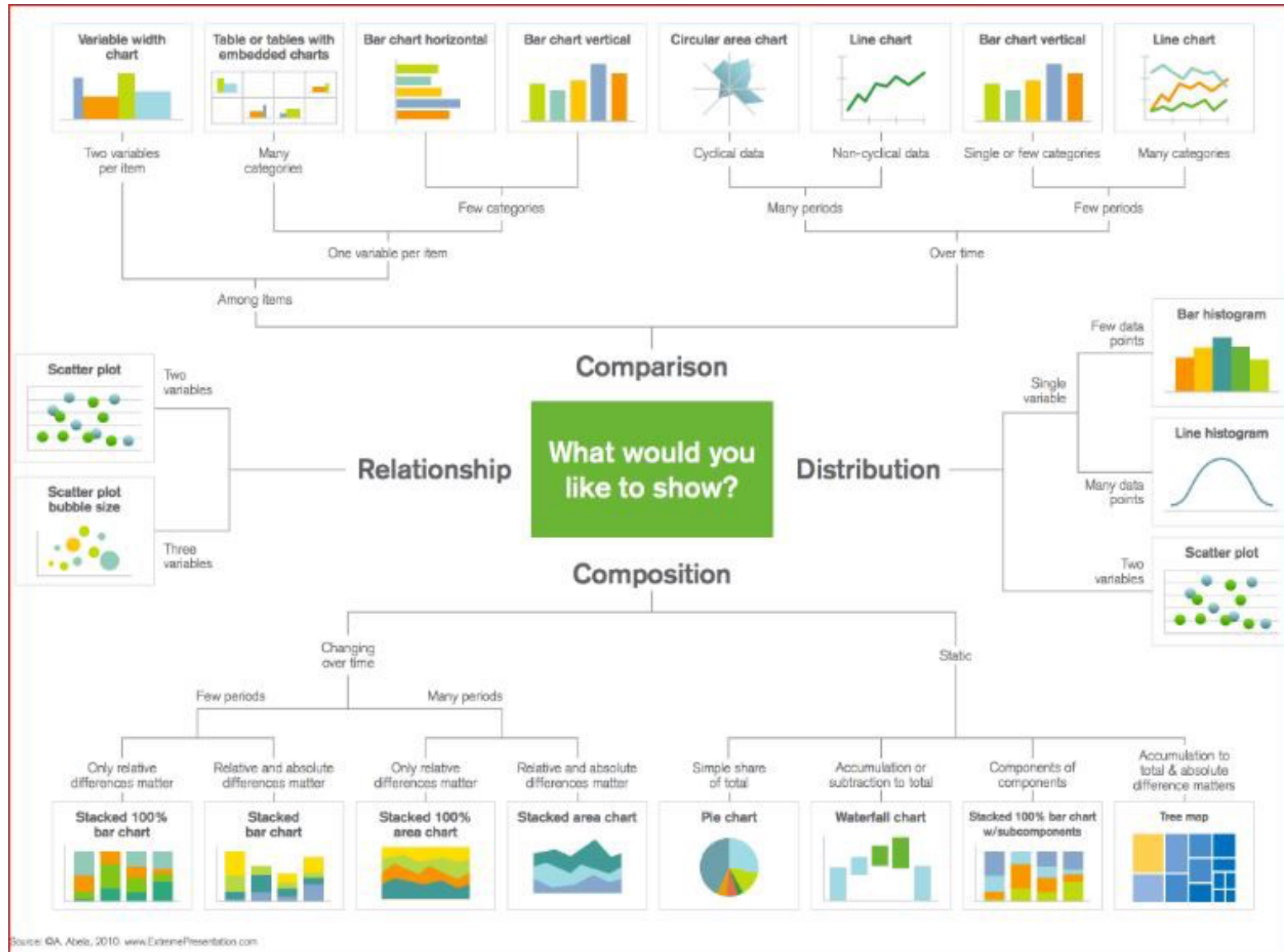
# Exploratory Data Analysis (EDA)

- Data Visualization is a big part of a data scientist's jobs.
- In the early stages of a project, you'll often be doing an **Exploratory Data Analysis (EDA)** to gain some insights into your data.
- Creating visualizations really helps make things clearer and easier to understand, especially with larger, high dimensional datasets.
- Towards the end of your project, it's important to be able to present your final results in a clear, concise, and compelling manner that your audience, whom are often non-technical clients, can understand.





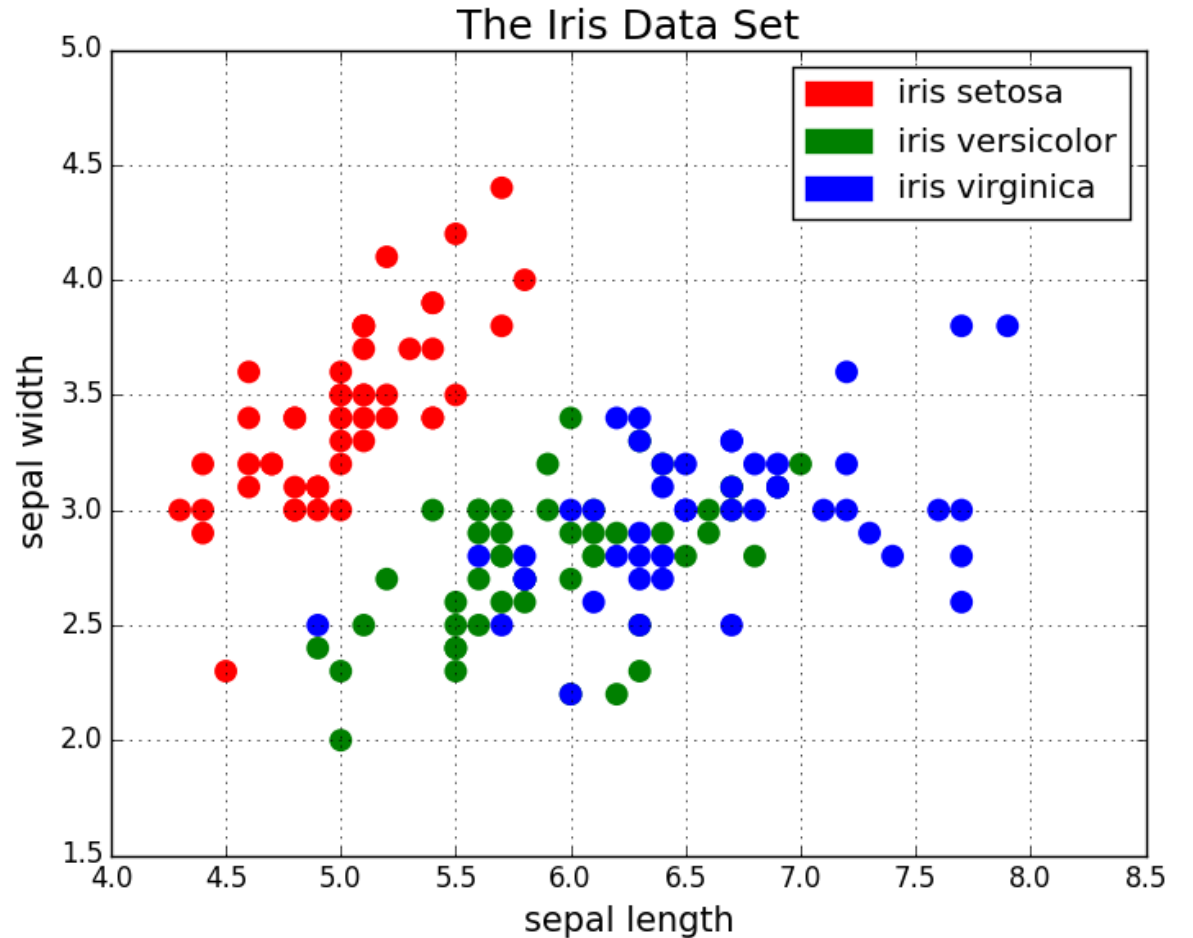
# EDA (cont.)





# Scatter Plots

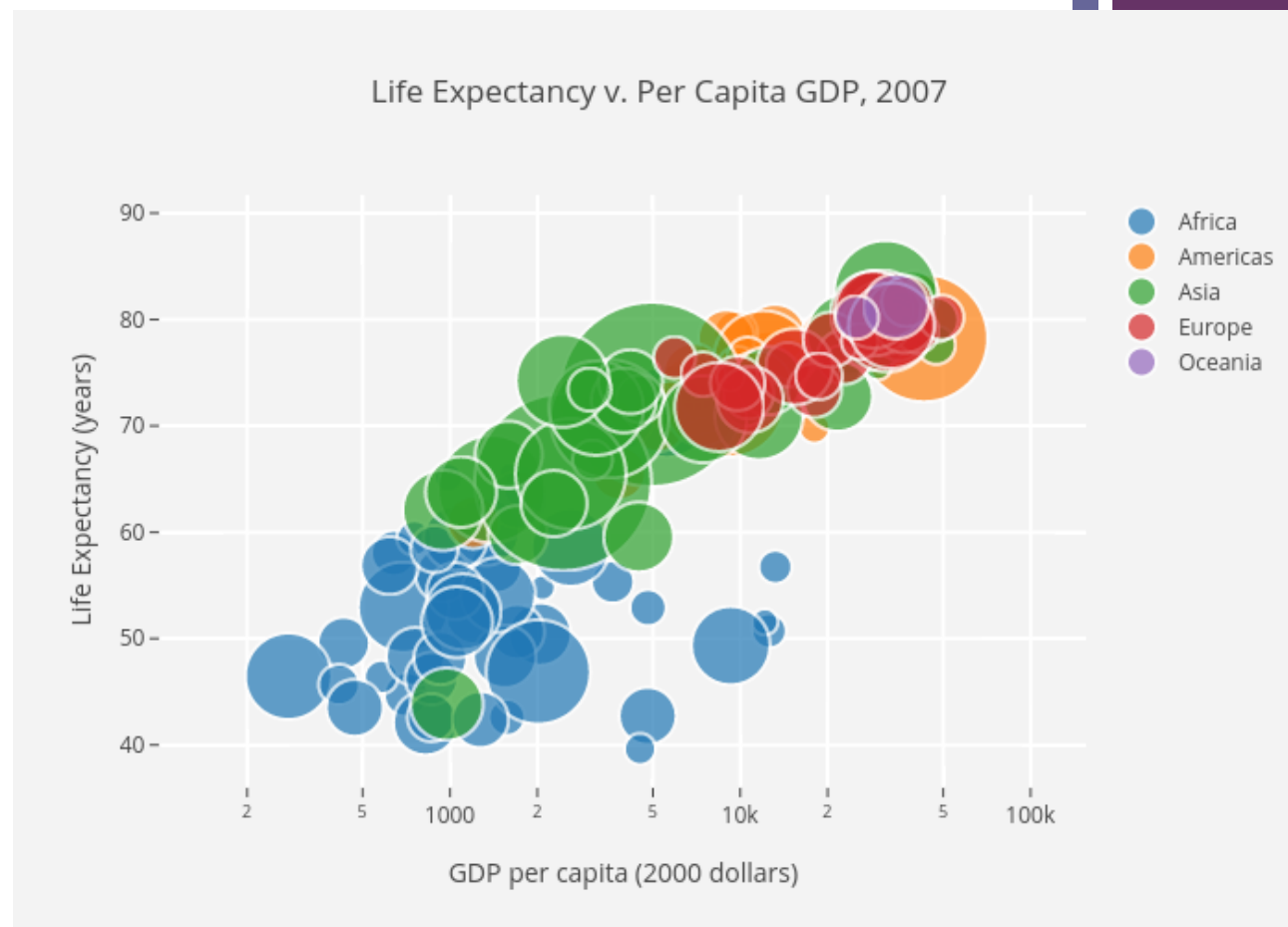
- Scatter plots are great for showing the relationship between two variables since you can directly see the raw distribution of the data.
- You can also view this relationship for different groups of data simple by colour coding the groups as seen in the first figure below.
- Want to visualize the relationship between three variables? No problemo! Just use another parameters, like point size, to encode that third variable as we can see in the second figure below. b





## Scatter Plots (cont.)

- Scatter plots are great for showing the relationship between two variables since you can directly see the raw distribution of the data.
- You can also view this relationship for different groups of data simple by colour coding the groups as seen in the first figure below.
- Want to visualize the relationship between three variables? No problemo! Just use another parameters, like point size, to encode that third variable as we can see in the second figure below. b



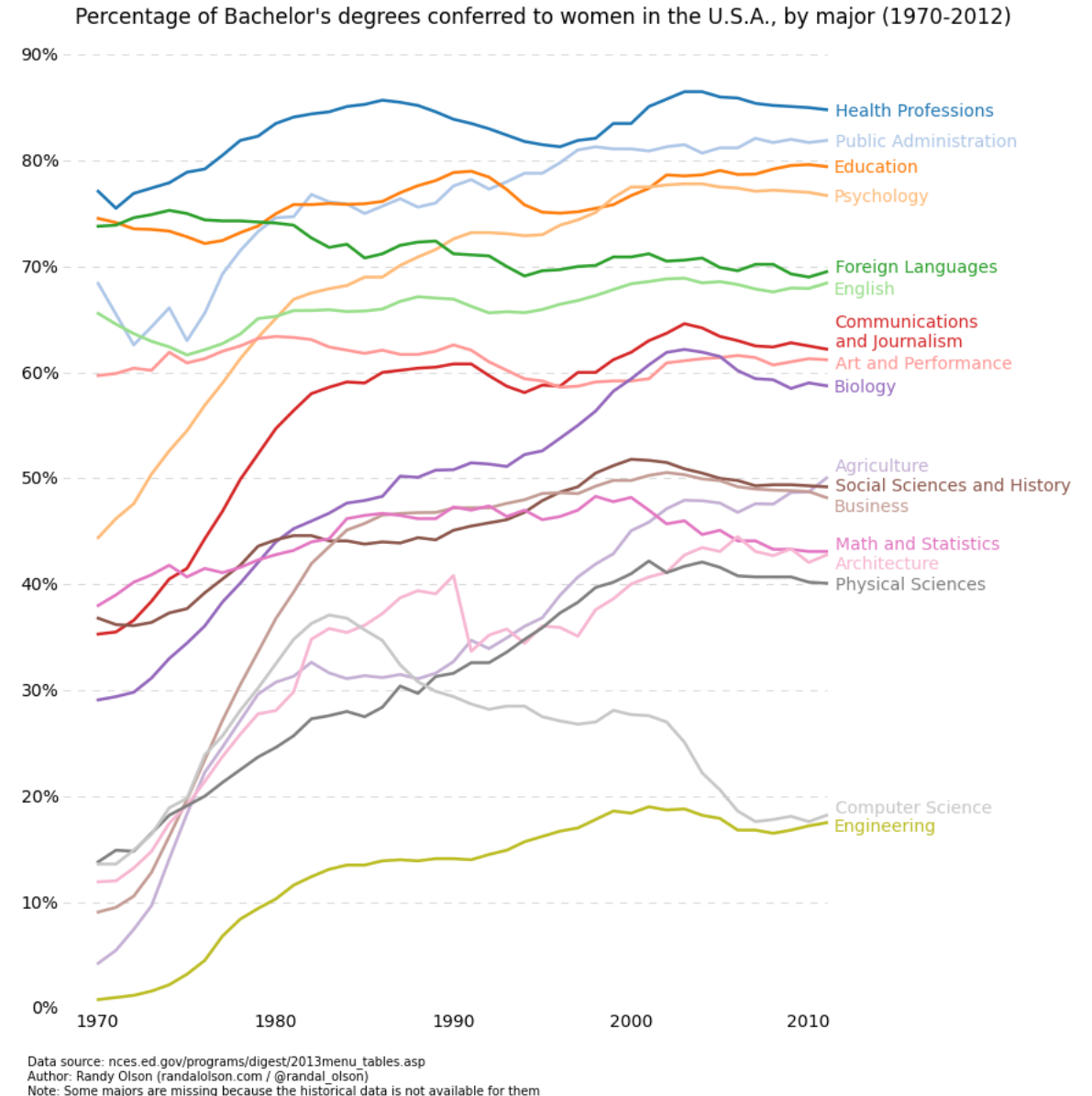
# Scatter Plots (cont.)

- Now for the code.
- We first import Matplotlib's pyplot with the alias "plt". To create a new plot figure we call `plt.subplots()`.
- We pass the x-axis and y-axis data to the function and then pass those to `ax.scatter()` to plot the scatter plot.
- We can also set the point size, point color, and alpha transparency.
- You can even set the y-axis to have a logarithmic scale. The title and axis labels are then set specifically for the figure.
- That's an easy to use function that creates a scatter plot end to end!

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 def scatterplot(x_data, y_data, x_label="", y_label="", title="", color = "r", yscale_log=False):
5
6     # Create the plot object
7     _, ax = plt.subplots()
8
9     # Plot the data, set the size (s), color and transparency (alpha)
10    # of the points
11    ax.scatter(x_data, y_data, s = 10, color = color, alpha = 0.75)
12
13    if yscale_log == True:
14        ax.set_yscale('log')
15
16    # Label the axes and provide a title
17    ax.set_title(title)
18    ax.set_xlabel(x_label)
19    ax.set_ylabel(y_label)
```

# Line Plots

- Line plots are best used when you can clearly see that one variable varies greatly with another i.e they have a high covariance. Lets take a look at the figure below to illustrate.
- We can clearly see that there is a large amount of variation in the percentages over time for all majors.
- Plotting these with a scatter plot would be extremely cluttered and quite messy, making it hard to really understand and see what's going on.
- Line plots are perfect for this situation because they basically give us a quick summary of the covariance of the two variables (percentage and time). Again, we can also use grouping by colour encoding.





## Line Plots (cont.)

- Here's the code for the line plot.
- It's quite similar to the scatter above. with just some minor variations in variables.

```
1  def lineplot(x_data, y_data, x_label="", y_label="", title=""):
2      # Create the plot object
3      _, ax = plt.subplots()
4
5      # Plot the best fit line, set the linewidth (lw), color and
6      # transparency (alpha) of the line
7      ax.plot(x_data, y_data, lw = 2, color = '#539caf', alpha = 1)
8
9      # Label the axes and provide a title
10     ax.set_title(title)
11     ax.set_xlabel(x_label)
12     ax.set_ylabel(y_label)
```

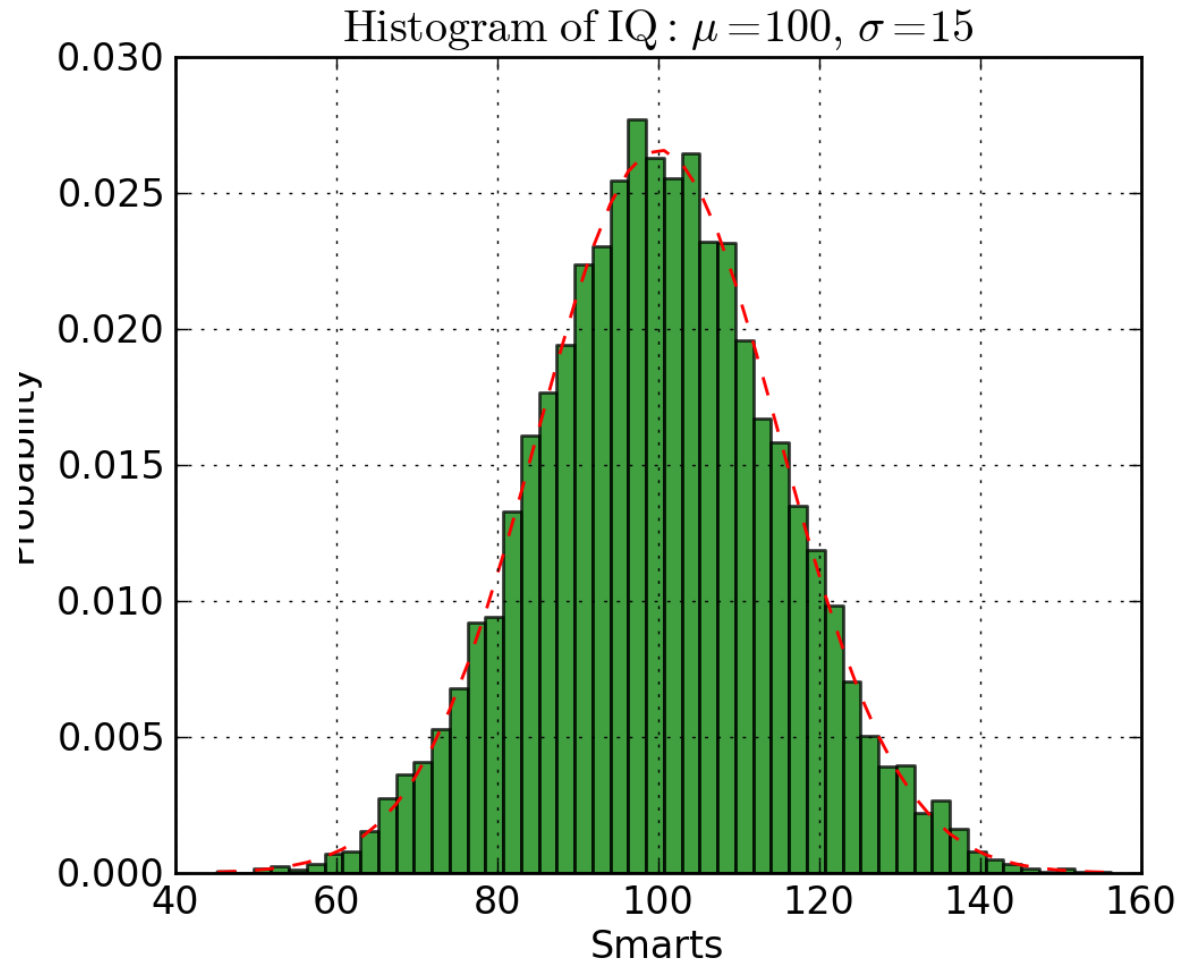
line\_plot.py hosted with ❤️ by GitHub

[view raw](#)



# Histograms Plots

- Histograms are useful for viewing (or really discovering) the distribution of data points.
- Check out the histogram below where we plot the frequency vs IQ histogram.
- We can clearly see the concentration towards the center and what the median is.
- We can also see that it follows a Gaussian distribution.
- Using the bars (rather than scatter points, for example) really gives us a clearly visualization of the relative difference between the frequency of each bin.
- The use of bins (discretization) really helps us see the “bigger picture” where as if we use all of the data points without discrete bins, there would probably be a lot of noise in the visualization, making it hard to see what is really going on.





# Histograms Plots (cont.)

- The code for the histogram in Matplotlib is shown below. There are two parameters to take note of.
- Firstly, the `n_bins` parameters controls how many discrete bins we want for our histogram.
- More bins will give us finer information but may also introduce noise and take us away from the bigger picture; on the other hand, less bins gives us a more “birds eye view” and a bigger picture of what’s going on without the finer details.
- Secondly, the `cumulative` parameter is a boolean which allows us to select whether our histogram is cumulative or not. This is basically selecting either the Probability Density Function (PDF) or the Cumulative Density Function (CDF).

```
1 def histogram(data, n_bins, cumulative=False, x_label = "", y_label = "", title = ""):
2     _, ax = plt.subplots()
3     ax.hist(data, n_bins = n_bins, cumulative = cumulative, color = '#539caf')
4     ax.set_ylabel(y_label)
5     ax.set_xlabel(x_label)
6     ax.set_title(title)
```

histogram.py hosted with ❤ by GitHub

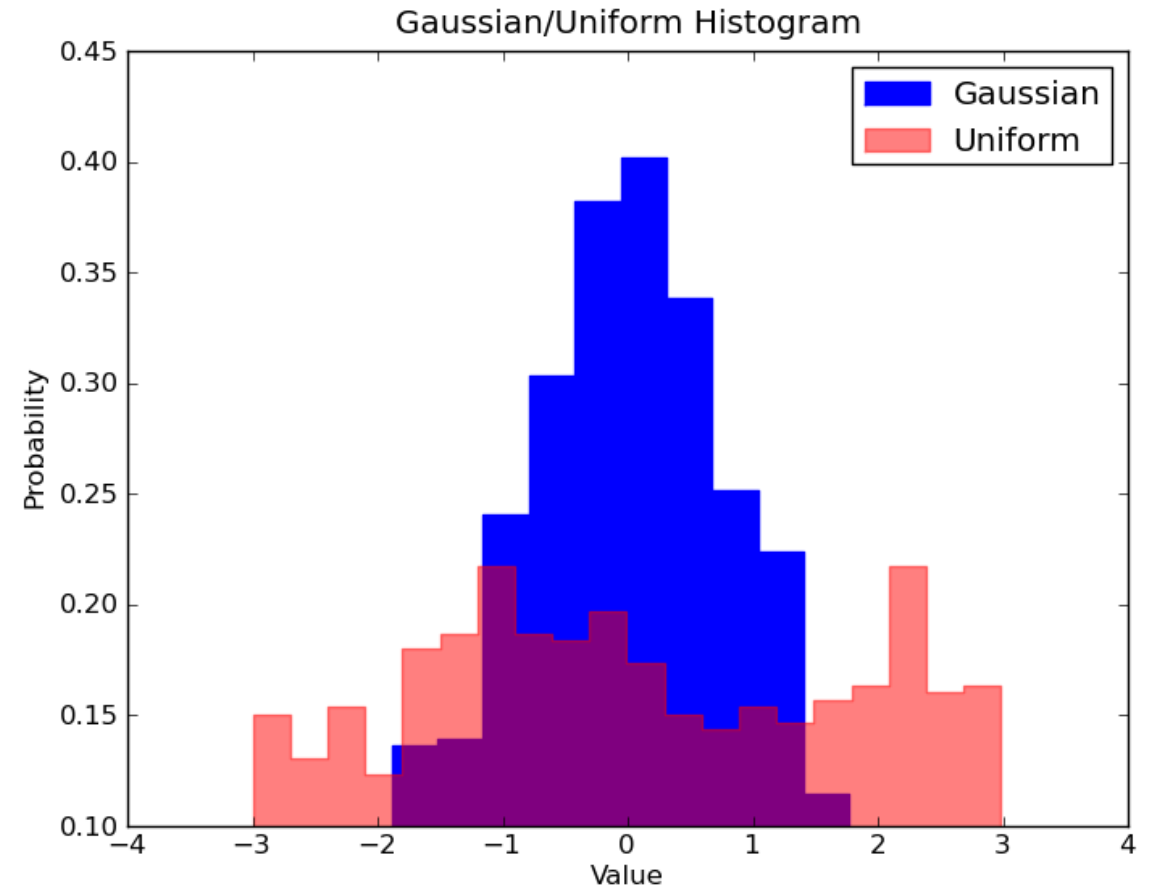
[view raw](#)





# Histograms Plots (cont.)

- Imagine we want to compare the distribution of two variables in our data.
- One might think that you'd have to make two separate histograms and put them side-by-side to compare them.
- But, there's actually a better way: we can overlay the histograms with varying transparency.
- Check out the figure below. The Uniform distribution is set to have a transparency of 0.5 so that we can see what's behind it. This allows use to directly view the two distributions on the same figure.



# Histograms Plots (cont.)

- There are a few things to set up in code for the overlaid histograms.
- First, we set the horizontal range to accommodate both variable distributions.
- According to this range and the desired number of bins we can actually compute the width of each bin.
- Finally, we plot the two histograms on the same plot, with one of them being slightly more transparent.

```
1 # Overlay 2 histograms to compare them
2 def overlaid_histogram(data1, data2, n_bins = 0, data1_name="", data1_color="#539caf", data2_color="#e377c2", data2_name=""):
3     # Set the bounds for the bins so that the two distributions are fairly compared
4     max_nbins = 10
5     data_range = [min(min(data1), min(data2)), max(max(data1), max(data2))]
6     binwidth = (data_range[1] - data_range[0]) / max_nbins
7
8
9     if n_bins == 0:
10         bins = np.arange(data_range[0], data_range[1] + binwidth, binwidth)
11     else:
12         bins = n_bins
13
14     # Create the plot
15     _, ax = plt.subplots()
16     ax.hist(data1, bins = bins, color = data1_color, alpha = 1, label = data1_name)
17     ax.hist(data2, bins = bins, color = data2_color, alpha = 0.75, label = data2_name)
18     ax.set_ylabel(y_label)
19     ax.set_xlabel(x_label)
20     ax.set_title(title)
21     ax.legend(loc = 'best')
```



# Bar Plots

- Bar plots are most effective when you are trying to visualize categorical data that has few (probably  $< 10$ ) categories.
- If we have too many categories then the bars will be very cluttered in the figure and hard to understand.
- They're nice for categorical data because you can easily see the difference between the categories based on the size of the bar (i.e magnitude); categories are also easily divided and colour coded too.
- There are 3 different types of bar plots we're going to look at: regular, grouped, and stacked. Check out the code below the figures as we go along.



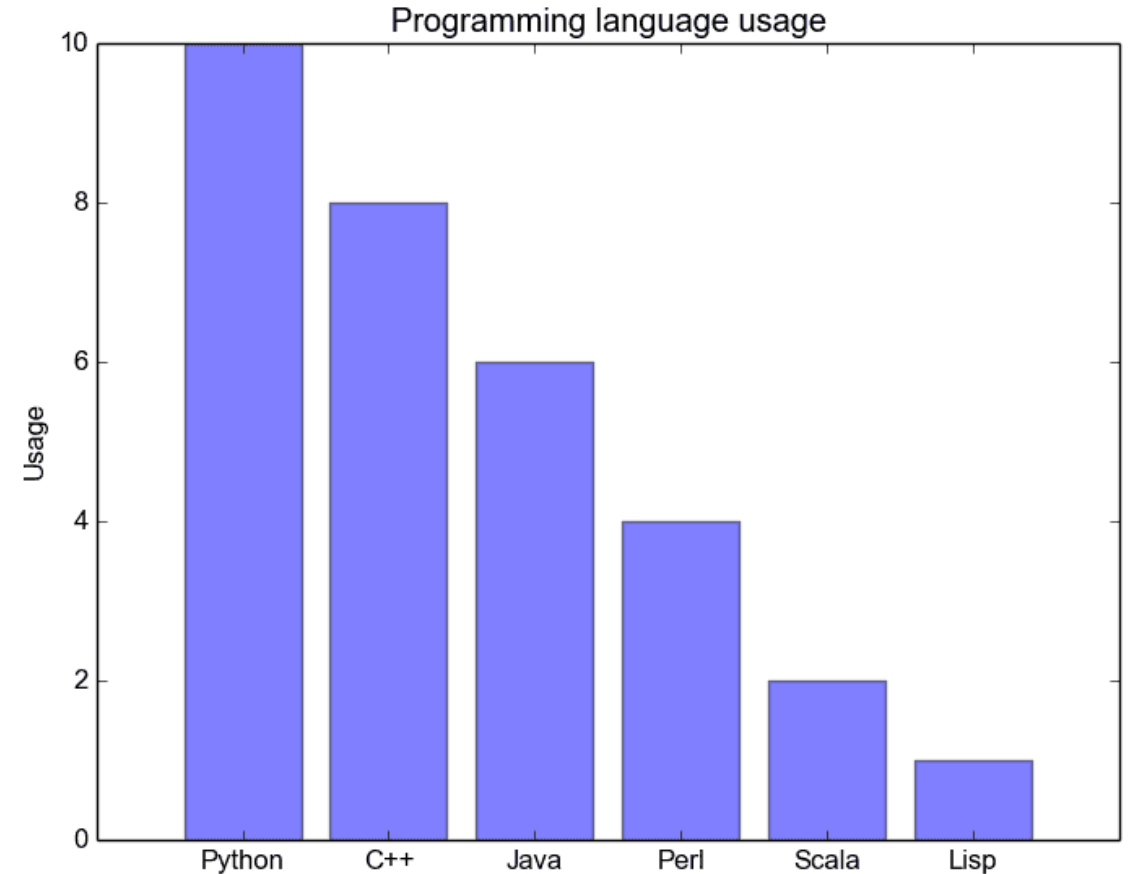
## Bar Plots (cont.)

- The regular barplot is in the next figure.
- In the `barplot()` function, `x_data` represents the tickers on the x-axis and `y_data` represents the bar height on the y-axis.
- The error bar is an extra line centered on each bar that can be drawn to show the standard deviation.
- Grouped bar plots allow us to compare multiple categorical variables. Check out the second bar plot below.
- The first variable we are comparing is how the scores vary by group (groups G1, G2, ... etc). We are also comparing the genders themselves with the colour codes.
- Taking a look at the code, the `y_data_list` variable is now actually a list of lists, where each sublist represents a different group. We then loop through each group, and for each group we draw the bar for each tick on the x-axis; each group is also colour coded.



## Bar Plots (cont.)

- Stacked bar plots are great for visualizing the categorical make-up of different variables.
- In the stacked bar plot figure below we are comparing the server load from day-to-day.
- With the colour coded stacks, we can easily see and understand which servers are worked the most on each day and how the loads compare to the other servers on all days.
- The code for this follows the same style as the grouped bar plot.
- We loop through each group, except this time we draw the new bars on top of the old ones rather than beside them.

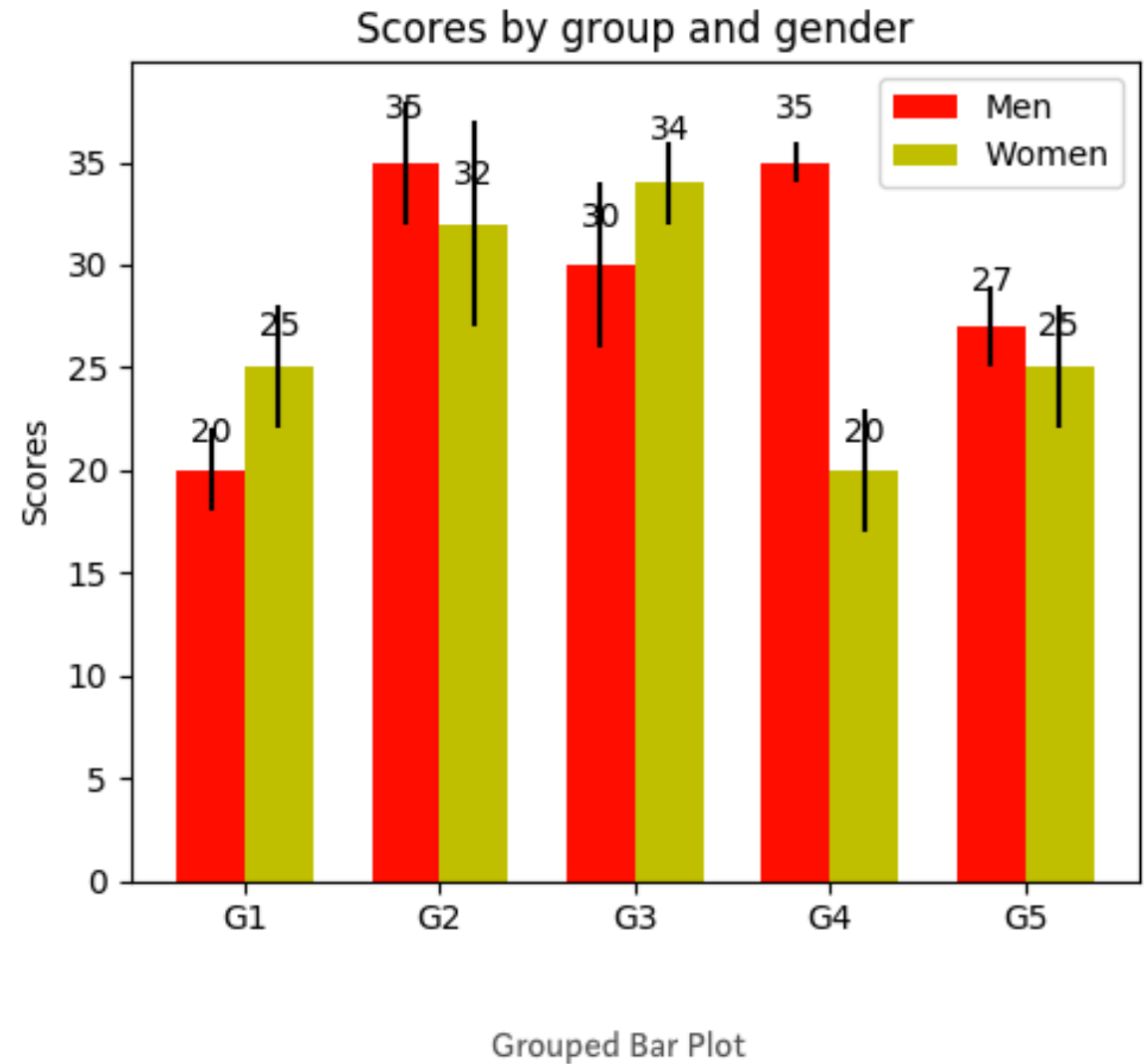




## Bar Plots (cont.)

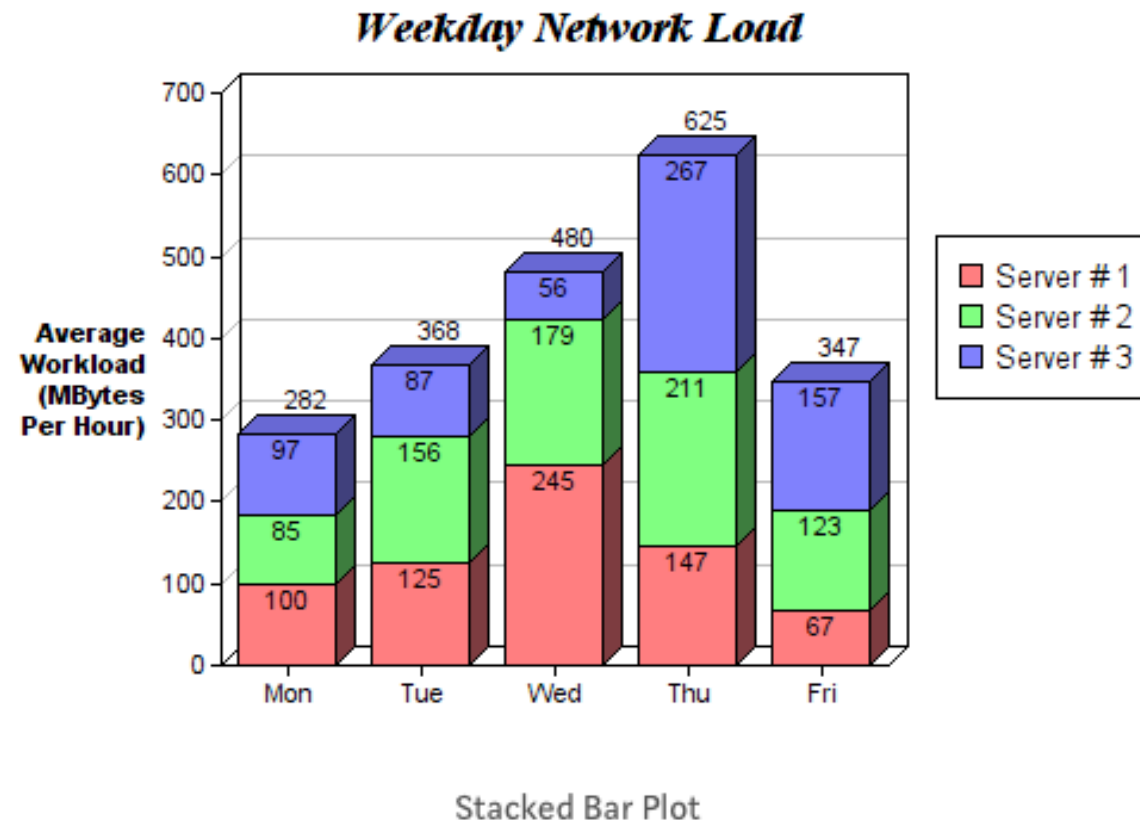
22

- Stacked bar plots are great for visualizing the categorical make-up of different variables.
- In the stacked bar plot figure below we are comparing the server load from day-to-day.
- With the colour coded stacks, we can easily see and understand which servers are worked the most on each day and how the loads compare to the other servers on all days.
- The code for this follows the same style as the grouped bar plot.
- We loop through each group, except this time we draw the new bars on top of the old ones rather than beside them.



# Bar Plots (cont.)

- Stacked bar plots are great for visualizing the categorical make-up of different variables.
- In the stacked bar plot figure below we are comparing the server load from day-to-day.
- With the colour coded stacks, we can easily see and understand which servers are worked the most on each day and how the loads compare to the other servers on all days.
- The code for this follows the same style as the grouped bar plot.
- We loop through each group, except this time we draw the new bars on top of the old ones rather than beside them.



```
1 def barplot(x_data, y_data, error_data, x_label="", y_label="", title=""):
2     _, ax = plt.subplots()
3     # Draw bars, position them in the center of the tick mark on the x-axis
4     ax.bar(x_data, y_data, color = '#539caf', align = 'center')
5     # Draw error bars to show standard deviation, set ls to 'none'
6     # to remove line between points
7     ax.errorbar(x_data, y_data, yerr = error_data, color = '#297083', ls = 'none', lw = 2,
8     ax.set_ylabel(y_label)
9     ax.set_xlabel(x_label)
10    ax.set_title(title)
11
```



```
14 def stackedbarplot(x_data, y_data_list, colors, y_data_names="", x_label="", y_label="", title=""):
15     _, ax = plt.subplots()
16     # Draw bars, one category at a time
17     for i in range(0, len(y_data_list)):
18         if i == 0:
19             ax.bar(x_data, y_data_list[i], color = colors[i], align = 'center', label = y_data_names[i])
20         else:
21             # For each category after the first, the bottom of the
22             # bar will be the top of the last category
23             ax.bar(x_data, y_data_list[i], color = colors[i], bottom = y_data_list[i - 1],
24                   label = y_data_names[i])
25     ax.set_ylabel(y_label)
26     ax.set_xlabel(x_label)
27     ax.set_title(title)
28     ax.legend(loc = 'upper right')
```

# Bar Plots (cont.)

```
31 def groupedbarplot(x_data, y_data_list, colors, y_data_names="", x_label="", y_label="", title=""):
32     _, ax = plt.subplots()
33     # Total width for all bars at one x location
34     total_width = 0.8
35     # Width of each individual bar
36     ind_width = total_width / len(y_data_list)
37     # This centers each cluster of bars about the x tick mark
38     alteration = np.arange(-(total_width/2), total_width/2, ind_width)
39
40     # Draw bars, one category at a time
41     for i in range(0, len(y_data_list)):
42         # Move the bar to the right on the x-axis so it doesn't
43         # overlap with previously drawn ones
44         ax.bar(x_data + alteration[i], y_data_list[i], color = colors[i], label = y_data_names[i])
45     ax.set_ylabel(y_label)
46     ax.set_xlabel(x_label)
47     ax.set_title(title)
48     ax.legend(loc = 'upper right')
```

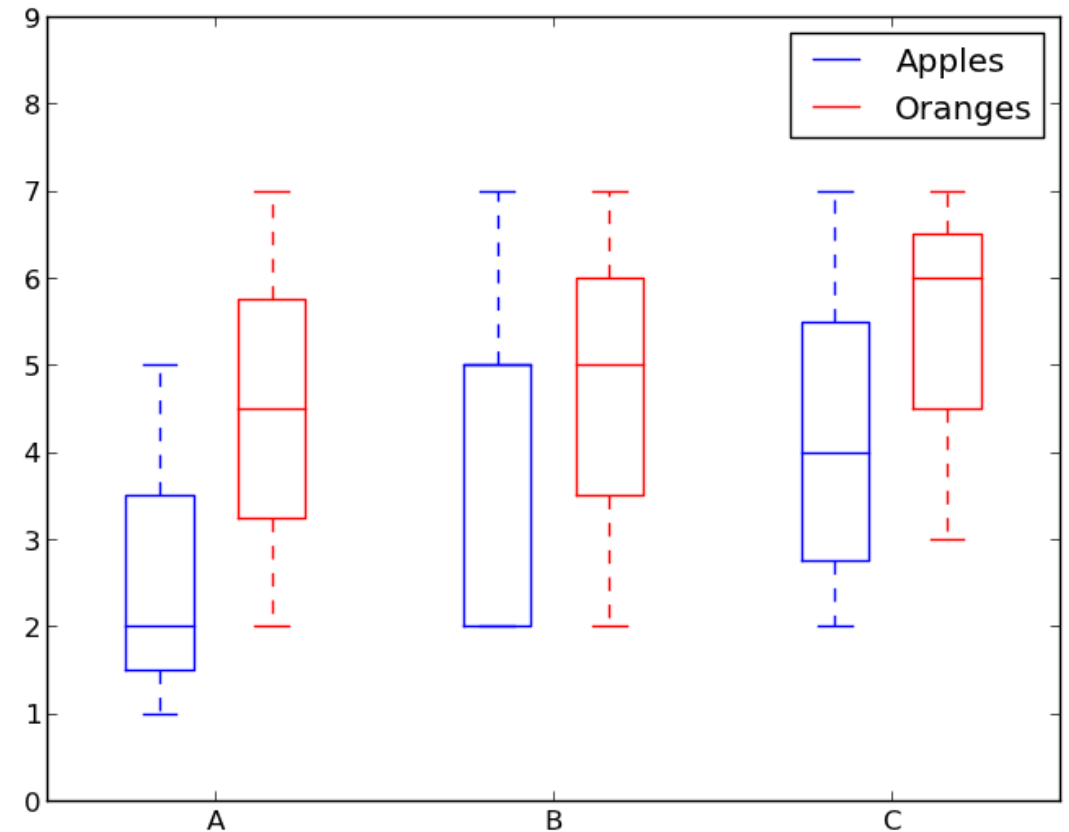
# Box Plots

- We previously looked at histograms which were great for visualizing the distribution of variables.
- But what if we need more information than that? Perhaps we want a clearer view of the standard deviation? Perhaps the median is quite different from the mean and thus we have many outliers? What if there is so skew and many of the values are concentrated to one side?
- That's where boxplots come in. Box plots give us all of the information above.
- The bottom and top of the solid-lined box are always the first and third quartiles (i.e 25% and 75% of the data), and the band inside the box is always the second quartile (the median).
- The whiskers (i.e the dashed lines with the bars on the end) extend from the box to show the range of the data.



## Box Plots (cont.)

- Since the box plot is drawn for each group/variable it's quite easy to set up. The `x_data` is a list of the groups/variables.
- The Matplotlib function `boxplot()` makes a box plot for each column of the `y_data` or each vector in sequence `y_data`; thus each value in `x_data` corresponds to a column/vector in `y_data`.
- All we have to set then are the aesthetics of the plot.



# Box Plots (cont.)

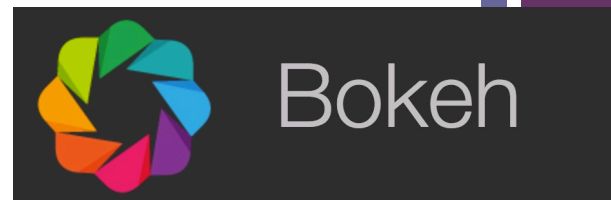
- Since the box plot is drawn for each group/variable it's quite easy to set up. The `x_data` is a list of the groups/variables.
- The Matplotlib function `boxplot()` makes a box plot for each column of the `y_data` or each vector in sequence `y_data`; thus each value in `x_data` corresponds to a column/vector in `y_data`.
- All we have to set then are the aesthetics of the plot.

```
1 def boxplot(x_data, y_data, base_color="#539caf", median_color="#297083", x_label="", y_label=""):  
2     _, ax = plt.subplots()  
3  
4     # Draw boxplots, specifying desired style  
5     ax.boxplot(y_data  
6                 # patch_artist must be True to control box fill  
7                 , patch_artist = True  
8                 # Properties of median line  
9                 , medianprops = {'color': median_color}  
10                # Properties of box  
11                , boxprops = {'color': base_color, 'facecolor': base_color}  
12                # Properties of whiskers  
13                , whiskerprops = {'color': base_color}  
14                # Properties of whisker caps  
15                , capprops = {'color': base_color})  
16  
17     # By default, the tick label starts at 1 and increments by 1 for  
18     # each box drawn. This sets the labels to the ones we want  
19     ax.set_xticklabels(x_data)  
20     ax.set_ylabel(y_label)  
21     ax.set_xlabel(x_label)  
22     ax.set_title(title)
```

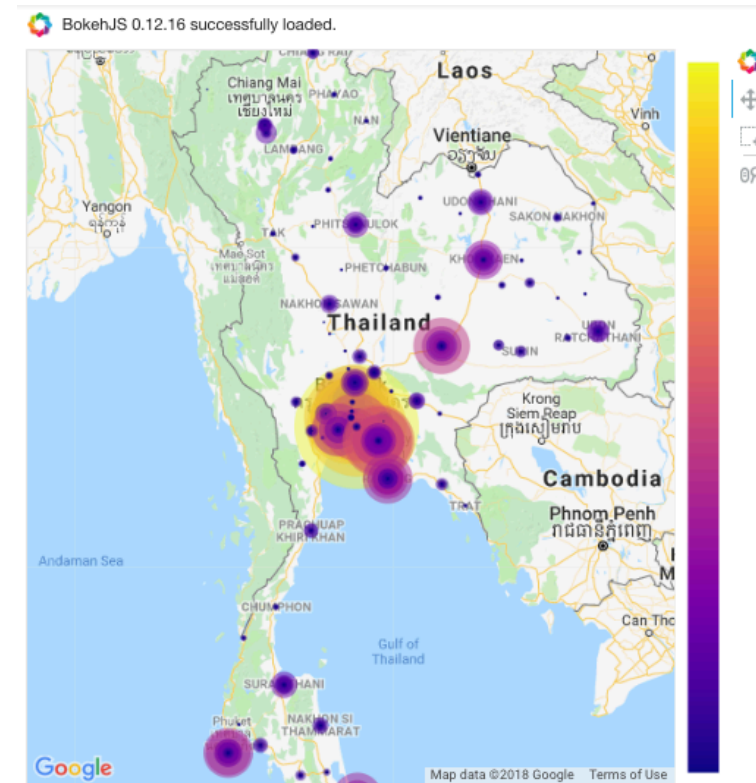


# Extra Lab: Geo Plot

30



	PRODUCT	Delinquent_Status	REGION	PROVINCE	AGE	ZIP_CODE	REGION_ENG	ACTIVE_ACCTS	ACTIVE_OS	TOTAL_WO_ACCTS	TOTAL_WO_OS
90221	Class C Loan	Z.WO	ภาคตะวันออกเฉียงเหนือ	อุบลราชธานี	07.>60	34320	4.North East	0	0.0	2	1946.89
90222	Class C Loan	Z.WO	ภาคตะวันออกเฉียงเหนือ	อุบลราชธานี	07.>60	34330	4.North East	0	0.0	2	70583.79
90223	Class C Loan	Z.WO	ภาคตะวันออกเฉียงเหนือ	อุบลราชธานี	07.>60	34340	4.North East	0	0.0	1	49854.61
90224	Class C Loan	Z.WO	ภาคตะวันออกเฉียงเหนือ	อุบลราชธานี	07.>60	34350	4.North East	0	0.0	2	288280.08
90225	Class C Loan	Z.WO	ภาคตะวันออกเฉียงเหนือ	อุบลราชธานี	07.>60	34360	4.North East	0	0.0	2	115262.40





Any Questions?

