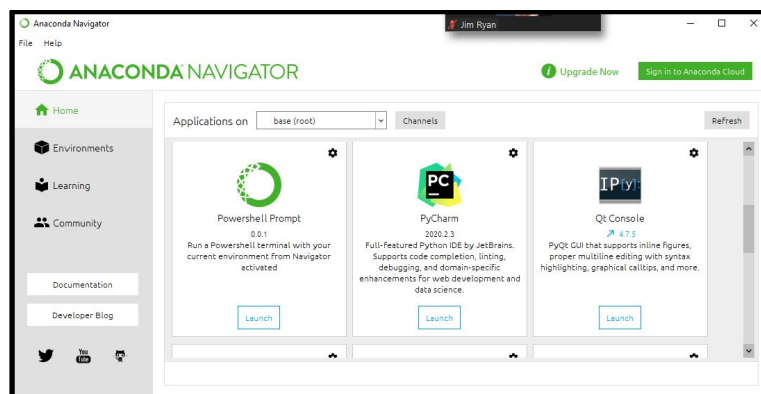# Data Visualization Fundamentals

## ⬥ Data Visualization in Python:

As we have covered in earlier weekly modules, data visualization is the discipline of trying to understand data by placing it in a visual context so that patterns, trends and correlations that might not otherwise be detected can be exposed.  Python offers multiple great graphing libraries that come packed with lots of different features. In prior tutorials, we have looked at boxplots, heatmaps, and histograms.  These are just a few of many charting options available in Python.  In this tutorial, we will be using the pandas, seaborn, and matplotlib libraries.  The following URLs list (1) data visualization packages we discuss in this tutorial, (2) types of charts available via the Seaborn package library, and (3) eight other data visualization package libraries outside of this tutorial:
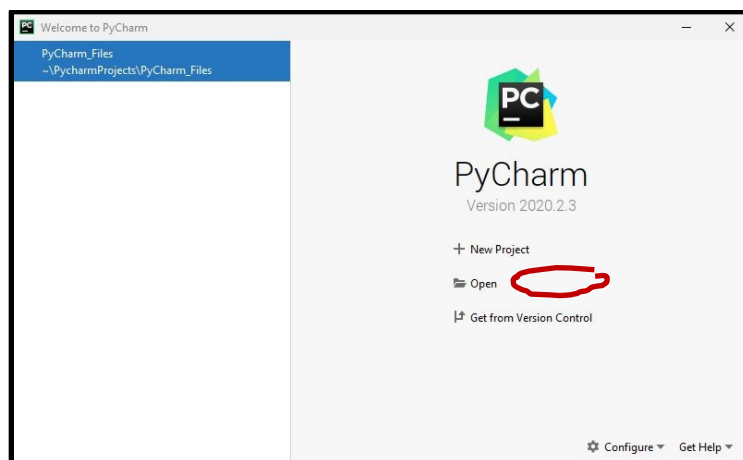
(1)  - https://towardsdatascience.com/introduction-to-data-visualization-in-python-89a54c97fbed

(2)  - https://python-graph-gallery.com/all-charts/

(3)  - https://mode.com/blog/python-data-visualization-libraries/
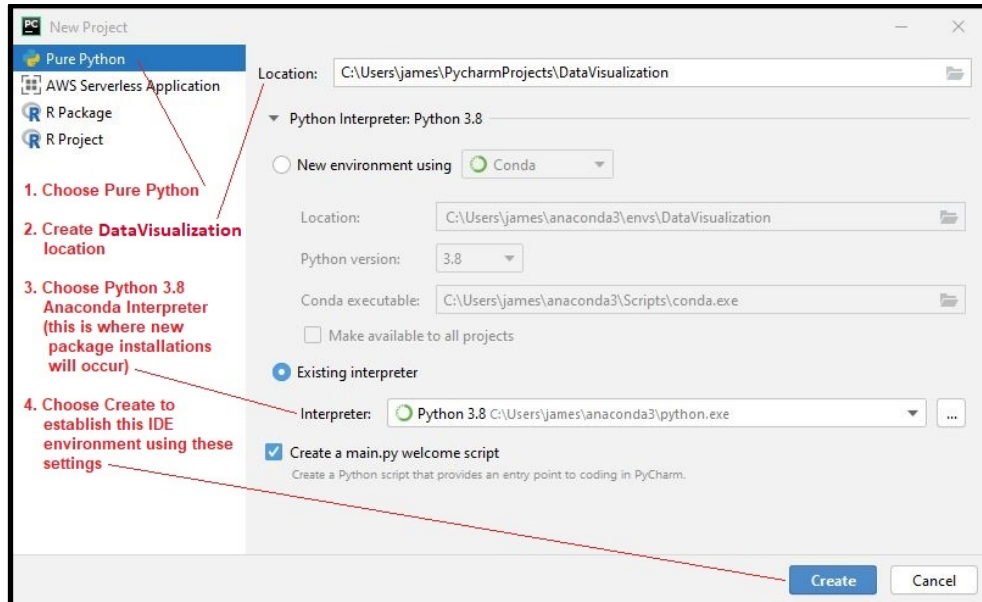
## ⬥ Create a DataVisualization Project in PyCharm

0) To ensure your various Python package libraries are available to import, be sure to use Anaconda Navigator to launch PyCharm.



1) Launch PyCharm and choose a New Project.

1) Choose a Pure Python project,
2) Change your project location to create a DataVisualization folder
3) Be sure to choose the Python interpreter as Python 3.8 associated with Anaconda.
4) Choose Create to open PyCharm with these settings.



*These images can be viewed at 400% magnification for details of the screen shots.*

5) *Download* the single_family_home_values.csv file from the Week 11 Canvas Module to the DataVisualization directory created in the last step.

## 🞣 Reading a CSV file into a Pandas DataFrame:

Pandas is an open-source Python Library used for high-performance data manipulation and data analysis using its powerful data structures.  Pandas use Series and DataFrame on top of Numpy arrary.  In this tutorial, we will look at DataFrame or df.  Python with pandas is in use in a variety of academic and commercial domains, including Finance, Economics, Statistics, Advertising, Web Analytics, and more. Using Pandas, we can accomplish five typical steps in the processing and analysis of data, regardless of the origin of data — load, organize, manipulate, model, and analyze the data.  For more specific details about Pandas not covered in this tutorial, please refer to the following URL:
https://www.tutorialspoint.com/python_data_science/python_pandas.htm

1) Anaconda will load the libraries in PyCharm for you.  To access the three libraries we are using for Tutorial Nine (e.g., pandas, seaborn, and matplotlib.pyplot), enter the following code into Main.py in the DataVisualization project directory of PyCharm:

```
import pandas as pd
import seaborn as sn
import matplotlib.pyplot as plt
```

2) To format the run window to handle the size of the csv data file we will be wrangling in Tutorial Seven, enter the following Python code:

```
desired_width = 400
pd.set_option('display.width', desired_width)
```

```
pd.set_option('display.max_columns', 20)
```

3)  To read the single_family_home_values.csv file into a pandas dataframe, type the following
    Python code into main.py via PyCharm in the line below the previous code above.

```
df = pd.read_csv(r'single_family_home_values.csv')  # reads Zillow file
```

4)  Once the csv file is loaded into a dataframe (df), we can begin data wrangling.  Any of the
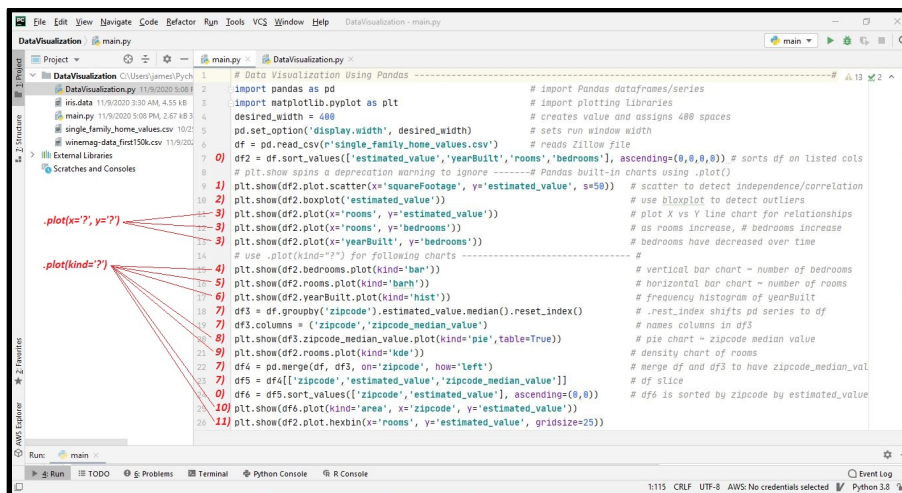    following Python functions may be used with a df, once it has been assigned.

```
print(df.head(2))      # returns top two rows
print(df.tail(5))      # returns bottom five rows
print(type(df))        # returns type of df
print(df.shape)        # returns row & column count
print(df.info())       # returns column attributes
print(df.describe())   # returns profiling stats
```

5)  In addition to 'single_family_home_values.csv', we will also use data files 'iris.data' and
    'winemag-data_first150k.csv'.  All three of these data files are available for download via
    the Week 11 Module in Canvas.  Feel free to use either of these data sets to practice data
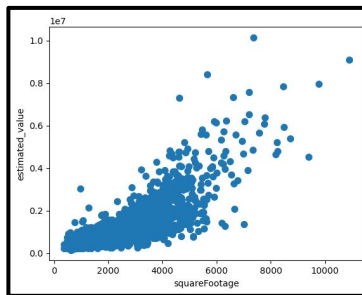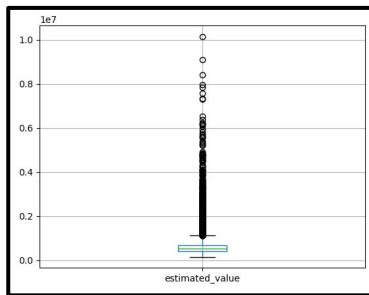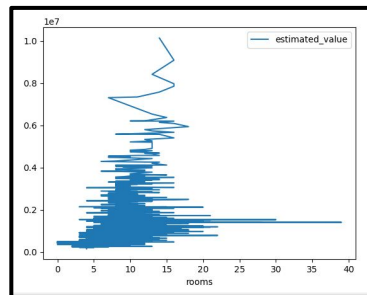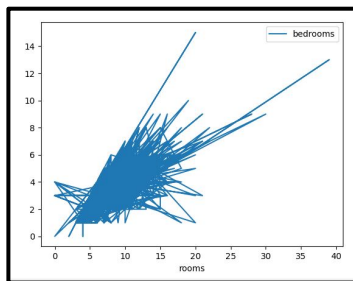    visualization tools.

## 🞣 Data Visualization using Pandas:

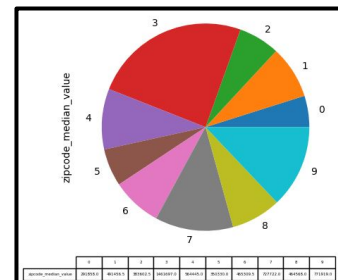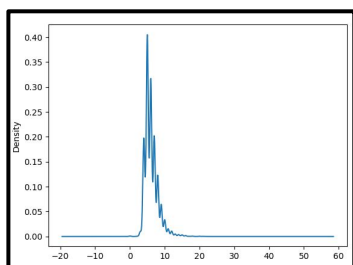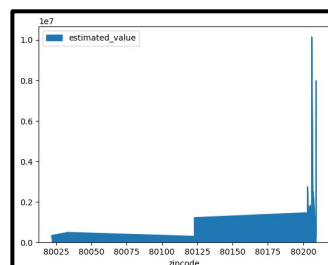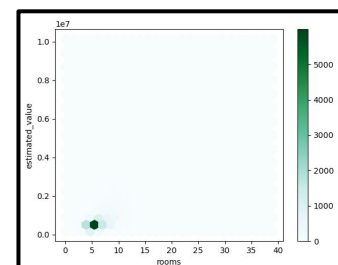The Pandas package libraries have built-in data visualization tools for many different types of
charts and plots. For more information outside this Tutorial on data visualization using Pandas
refer to the URL: https://pandas.pydata.org/pandas-docs/stable/user_guide/visualization.html
The following PyCharm screenshot shows the Python code for (0) sorting data frame values,  (1)
graphing a scatter plot, (2) boxplot, (3) three different x-y plots, (4) vertical bar chart, (5)
horizontal bar chart, (6) frequency histogram, (7)  group by, merge, or slice of data frames, (8)
graphing a pie chart with data table, (9) density distribution chart, (10) area chart, and (11)
hexoganol bin chart .  All these data visualization examples use the plt.show() function from the
matplotlib.pyplot package library to display the graphic.

A screen scraping tool like the Windows Snipping Tool can be used to save a plot.show() graphic or as we learned in Tutorial Eight-B under Network Mining, the plt.savefig() function from matplotlib.pyplot will save a copy of the graphic to disk in the Python project directory.  From the Python code in the previous PyCharm screenshot, the following graphics are chart screenshots.


*(1) Scatterplot*


*(2) boxplot*


*(3) plot(x,y)*


*(3) plot(x,y)*


*(3) plot(x,y)*


*(4) vertical bar*


*(5) horizontal bar*


*(6) frequency histogram*


*(8) pie chart with data table*


*(9) density distribution chart*


*(10) area chart*
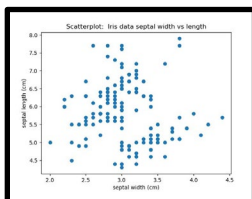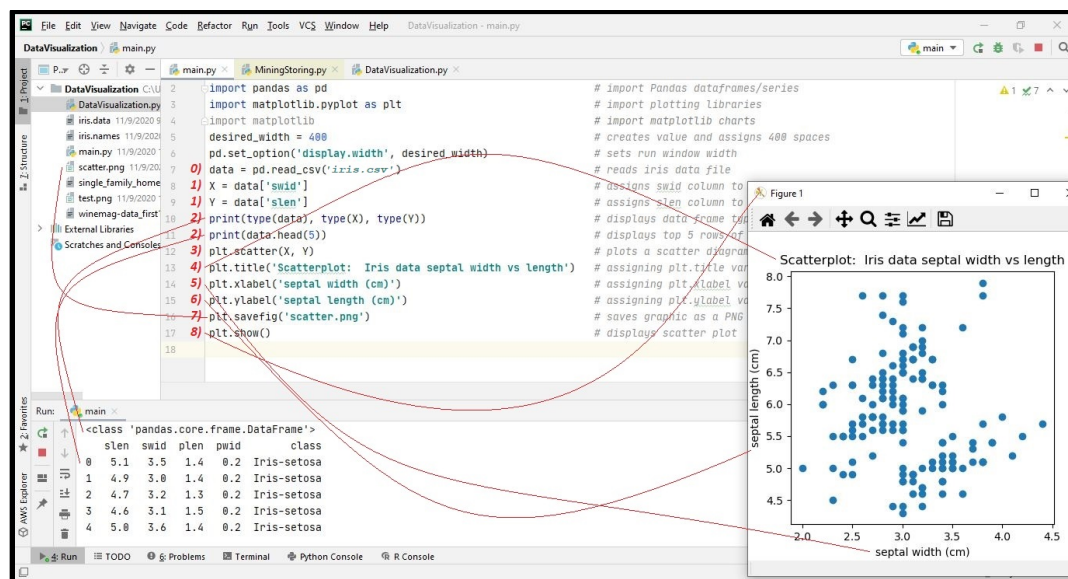

*(11) hexagonal bin chart*

**Exercise deliverable:** Construct an example screenshot of using Pandas built-in data visualization functions in Python code in either of the Week 11 data files or your data project csv. If you are using the csv above, then choose different column attributes for the charts. (1) Choose four different charts. (2) Capture a PyCharm screenshot of the Python code creating the four charts. (3) Either use a screen scraper or plt.savfig() function to save charts to disk as PNG files. Save the PyCharm screenshot and chart PNG files for use later in this tutorial.

## Data Visualization Using Matplotlib:

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. Matplotlib makes easy things easy and hard things possible. In this and prior Tutorials, we have used the plt.show() and plt.savefig() functions to display or save charts. Consider all of the charts available via the matplotlib libraries via the URL: https://matplotlib.org/gallery/index.html . In the PyCharm screenshot below, we use matplotlib.pyplot to create a scatterplot from the 'iris.csv' file. The Python code assigns plt.title, plt.Xlabel, and plt.Ylabel variables to appear on the scatterplot diagram. Specifically, the Python code demonstrates (0) reading the CSV file, (1) create a data frame slice for X and Y, (2) print the data frame type and the top five rows of the data frame, (3) create the scatterplot using X and Y data frames, (4) assign the plot title, (5) assign the X label title, (6) assign the Y label title, (7) save the scatterplot to disk as file scatter.png, and (8) display the scatterplot chart. Most matplotlib functions are used with NumPy and this example is an exception.





The figure to the left is the scatterplot saved to disk in the PyCharmProjects subdirectory. Saving a chart to disk allows you to access the graphic at a later date, rather than recreating the graphic. Note the plt.title, plt.Xlabel, and plt.Ylabel variables on the scatter.png graphic. All illustrations in this Tutorial may be viewed at 400 times.

**Exercise deliverable:** Construct an example scatterplot from your project data set.  Be sure to assign titles for the chart as well as the X and Y axis.  Remember the utility of scatter diagrams is the review the independence between variables X and Y.  Variable Y is not independent of variable X if any correlation (positive or negative) can be construed from the scatter diagram.  Save the PyCharm screenshot showing the Python code and the scatter diagram PNG file (i.e., use a screen scraper or plt.savefig() for later in this tutorial.  (If you have issues using your project data set, then use one of the data sets in Week 11 module).

## Using Seaborn for Data Visualization:

Seaborn is a Python data visualization library based on matplotlib and integrates closely with Pandas data structures. It provides a high-level interface for drawing attractive and informative statistical graphics.  In prior Tutorials, we have used the seaborn boxplot as well as heat map.  More graphics are available to review via the URL:
https://seaborn.pydata.org/examples/index.html

In the PyCharm screenshot below, the we combine Pandas data frames and Seaborn plotting functions to create four different statistical charts (e.g., box, pair, strip, and violin).  Specifically, the Python code (0) reads the Zillow home sales CSV into a Pandas df, (1) slices the df to remove outliers, (2) displays a box chart with the inter-quartile range, (3) displays a pair chart comparing scatter diagrams and area charts for estimated value and lastSaleAmount, (4) displays a strip chart of estimated home value distributions by zip code, and (5) a violin chart of estimated home value box plots by zip code.  All for of these charts help describe distributions and the plt.show() function from matplotlib.pyplot enables the display of the seaborn graphic.  All Python code from the PyCharm screenshots in Tutorial Nine may be found in the file DataVisualization.py, which can be downloaded from the Week 11 Module in Canvas.

The four charts created by the Python code in the previous PyCharm screenshot are viewable as



***Box Chart***



***Pair Chart***



***Strip Chart***



***Violin Chart***

**Exercise deliverable:** Construct two graphic examples of a box chart, pair chart, strip chart, and/or violin chart for practice in data visualization tools.  Capture the chart graphics using the plt.savefig() function or a screen scraper. Use your project data set or either of the CSV files in the Week 11 Module.  Save the PyCharm screenshot displaying the Python code and the two chart graphics for use later in this tutorial.

## ✚ Deliverables from Exercises Above:

Take the noted deliverables (i.e., PyCharm screenshots) from each exercise above (i.e., total of 4 screenshots), insert the screenshots into a MS Word document as you label each screenshot deliverable from the corresponding exercise, save the document as one PDF, and submit the PDF to the Tutorial Seven Canvas Assignment uplink having the following qualities:

a.  Your screenshot for each of the two exercises above is included.

b.  Each screenshot is labeled appropriately for each exercise.

      i.  Pandas Built-in Data Visualization …………….………….. 5 screenshots

     ii.  Matplotlib Data Visualization  ………………….…………. 2 screenshots

    iii. Seaborn Data Visualization ………………………………. 3 screenshot

    Total screenshots to submit ………………………………… 10 screenshots

c. All screenshots are legible output for each PyCharm exercise deliverable.