

Data-X Spring 2019: Homework 04

Name: Shrey Samdani

SID: 3032000414

In this homework, you will do some exercises with plotting.

REMEMBER TO DISPLAY ALL OUTPUTS. If the question asks you to do something, make sure to print your results.

1.

Data:

Data Source: Data file is uploaded to bCourses and is named: **Energy.csv**

The dataset was created by Angeliki Xifara (Civil/Structural Engineer) and was processed by Athanasios Tsanas, Oxford Centre for Industrial and Applied Mathematics, University of Oxford, UK).

Data Description:

The dataset contains eight attributes of a building (or features, denoted by X1...X8) and response being the heating load on the building, y1.

- X1 Relative Compactness
- X2 Surface Area
- X3 Wall Area
- X4 Roof Area
- X5 Overall Height
- X6 Orientation
- X7 Glazing Area
- X8 Glazing Area Distribution
- y1 Heating Load

Q1.1

Read the data file in python. Check if there are any NaN values, and print the results.

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

energy = pd.read_csv("Energy.csv")
print("# of NaN values = ",sum(pd.isnull(energy).any(axis=1)))
```

of NaN values = 0

Q 1.2

Describe (using python function) data features in terms of type, distribution range (max and min), and mean values.

```
In [2]: print(energy.dtypes)
energy.describe()
```

```
X1    float64
X2    float64
X3    float64
X4    float64
X5    float64
X6     int64
X7    float64
X8     int64
Y1    float64
dtype: object
```

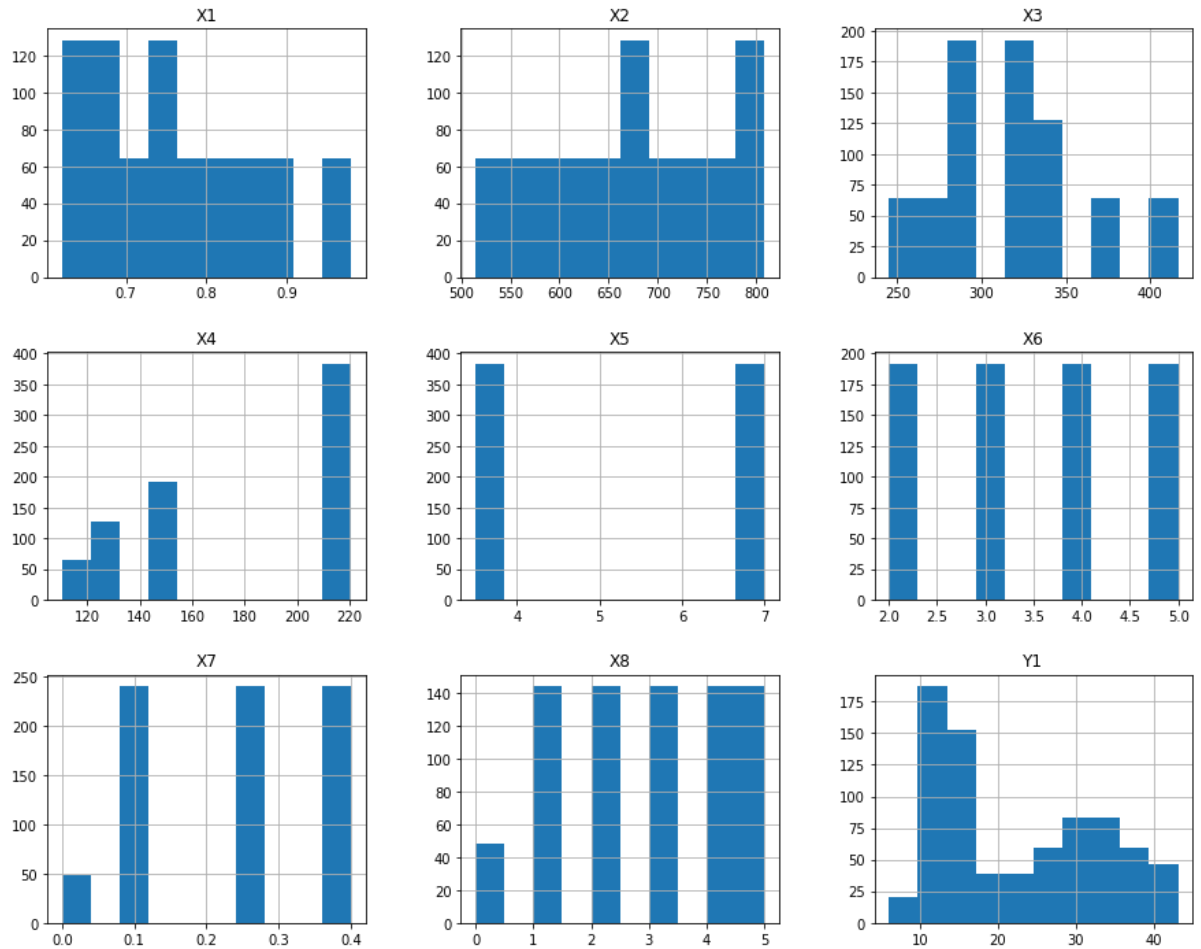
Out[2]:

	X1	X2	X3	X4	X5	X6	X7	
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.0
mean	0.764167	671.708333	318.500000	176.604167	5.250000	3.500000	0.234375	2.8
std	0.105777	88.086116	43.626481	45.165950	1.75114	1.118763	0.133221	1.5
min	0.620000	514.500000	245.000000	110.250000	3.50000	2.000000	0.000000	0.0
25%	0.682500	606.375000	294.000000	140.875000	3.50000	2.750000	0.100000	1.7
50%	0.750000	673.750000	318.500000	183.750000	5.25000	3.500000	0.250000	3.0
75%	0.830000	741.125000	343.000000	220.500000	7.00000	4.250000	0.400000	4.0
max	0.980000	808.500000	416.500000	220.500000	7.00000	5.000000	0.400000	5.0

Q 1.3

Plot feature distributions for all the attributes in the dataset (Hint - Histograms are one way to plot data distributions). This step should give you clues about data sufficiency.

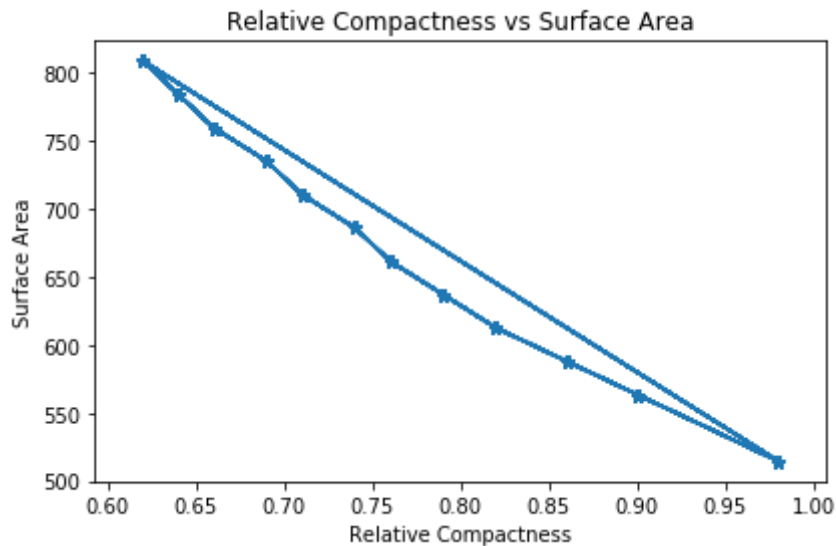
```
In [3]: energy.hist(figsize = (15, 12))
plt.show()
```



Q1.4

Create a combined line and scatter plot for attributes 'X1' and 'X2' with a marker (*). You can choose either of the attributes as x & y. Label your axes and give a title to your plot.

```
In [4]: # your code
plt.scatter(energy['X1'],energy['X2'], marker = "*")
plt.plot(energy['X1'],energy['X2'], marker = "*")
plt.xlabel("Relative Compactness")
plt.ylabel("Surface Area")
plt.title("Relative Compactness vs Surface Area")
plt.rcParams['figure.dpi'] = 150
plt.tight_layout()
plt.show()
```

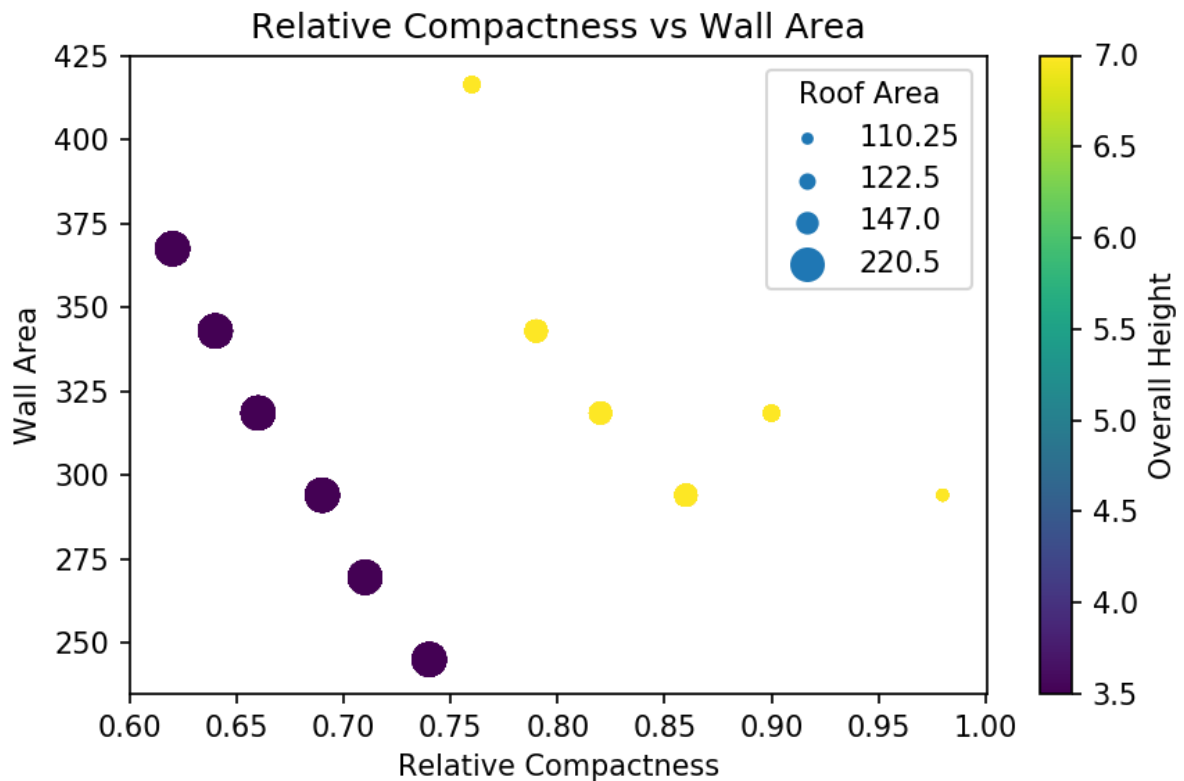


Q1.5

Create a scatter plot for how 'Wall Area' changes with 'Relative Compactness'. Give different colors for different 'Overall Height' and different bubble sizes by 'Roof Area'. Label the axes and give a title. Add a legend to your plot.

```
In [5]: for size,group in energy.groupby("X4"):
        plt.scatter(group["X1"], group["X3"], c=group["X5"], s=group["X4"]-100, label = size,vmin=3.5, vmax=7)
plt.xlabel("Relative Compactness")
plt.colorbar(label = "Overall Height")
plt.ylabel("Wall Area")
plt.title("Relative Compactness vs Wall Area")

plt.tight_layout()
plt.legend(title = "Roof Area")
plt.show()
```



2.

Q 2.1a.

Create a dataframe called `icecream` that has column `Flavor` with entries `Strawberry`, `Vanilla`, and `Chocolate` and another column with `Price` with entries `3.50`, `3.00`, and `4.25`. Print the dataframe.

```
In [6]: icecream = pd.DataFrame({"Flavor": ["Strawberry", "Vanilla", "Chocolate"],  
                                "Price": [3.5, 3, 4.25]})  
icecream
```

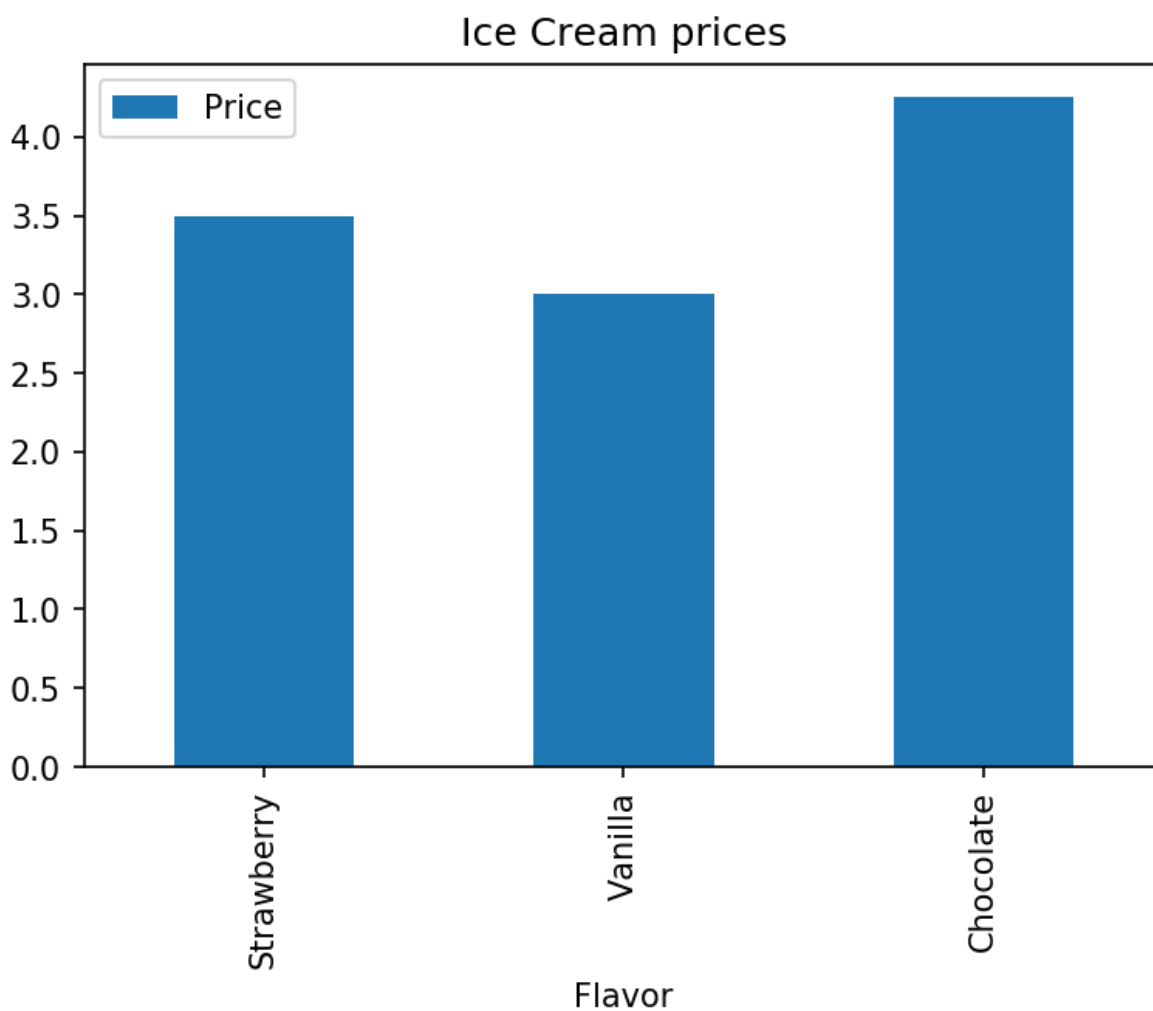
Out[6]:

	Flavor	Price
0	Strawberry	3.50
1	Vanilla	3.00
2	Chocolate	4.25

Q 2.1b

Create a bar chart representing the three flavors and their associated prices. Label the axes and give a title.

```
In [7]: icecream.plot.bar(x = "Flavor")  
plt.title("Ice Cream prices")  
plt.show()
```



Q 2.2

Create 9 random plots in a figure (Hint: There is a numpy function for generating random data).

The top three should be scatter plots (one with green dots, one with purple crosses, and one with blue triangles). The middle three graphs should be a line graph, a horizontal bar chart, and a histogram. The bottom three graphs should be trigonometric functions (one sin, one cosine, one tangent). Keep in mind the range and conditions for the trigonometric functions.

All these plots should be on the same figure and not 9 independent figures.

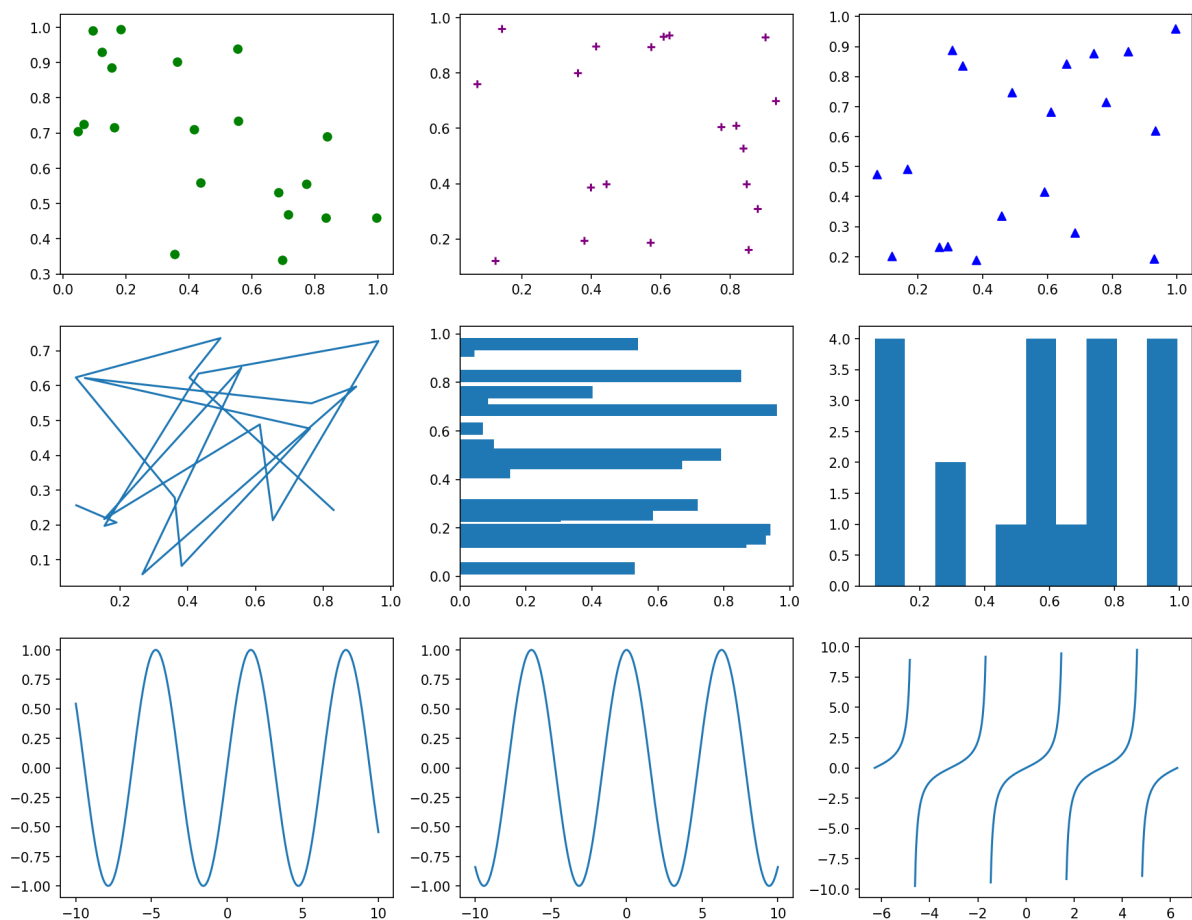
```

In [8]: fig, ax = plt.subplots(3,3, figsize = (15,12))
data = np.random.rand(6,2,20)
trigX = np.linspace(-10,10,200)
ax[0,0].scatter(data[0][0],data[0][1], c='g')
ax[0,1].scatter(data[1][0],data[1][1], c='purple', marker="+")
ax[0,2].scatter(data[2][0],data[2][1], c='b', marker="^")
ax[1,0].plot(data[3][0],data[3][1])
ax[1,1].barh(y=data[4][0],width=data[4][1], height = 0.05)
ax[1,2].hist(data[5][0])
ax[2,0].plot(trigX,np.sin(trigX))
ax[2,1].plot(trigX,np.cos(trigX))

#tangent
x = np.linspace(-1.0, 1.0, 1000)
x *= 2 * np.pi
y = np.tan(x)
y[y > 10] = np.nan
y[y < -10] = np.nan
ax[2,2].plot(x,y)
plt.show()

```

/Users/shreysamdani/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:18: RuntimeWarning: invalid value encountered in less



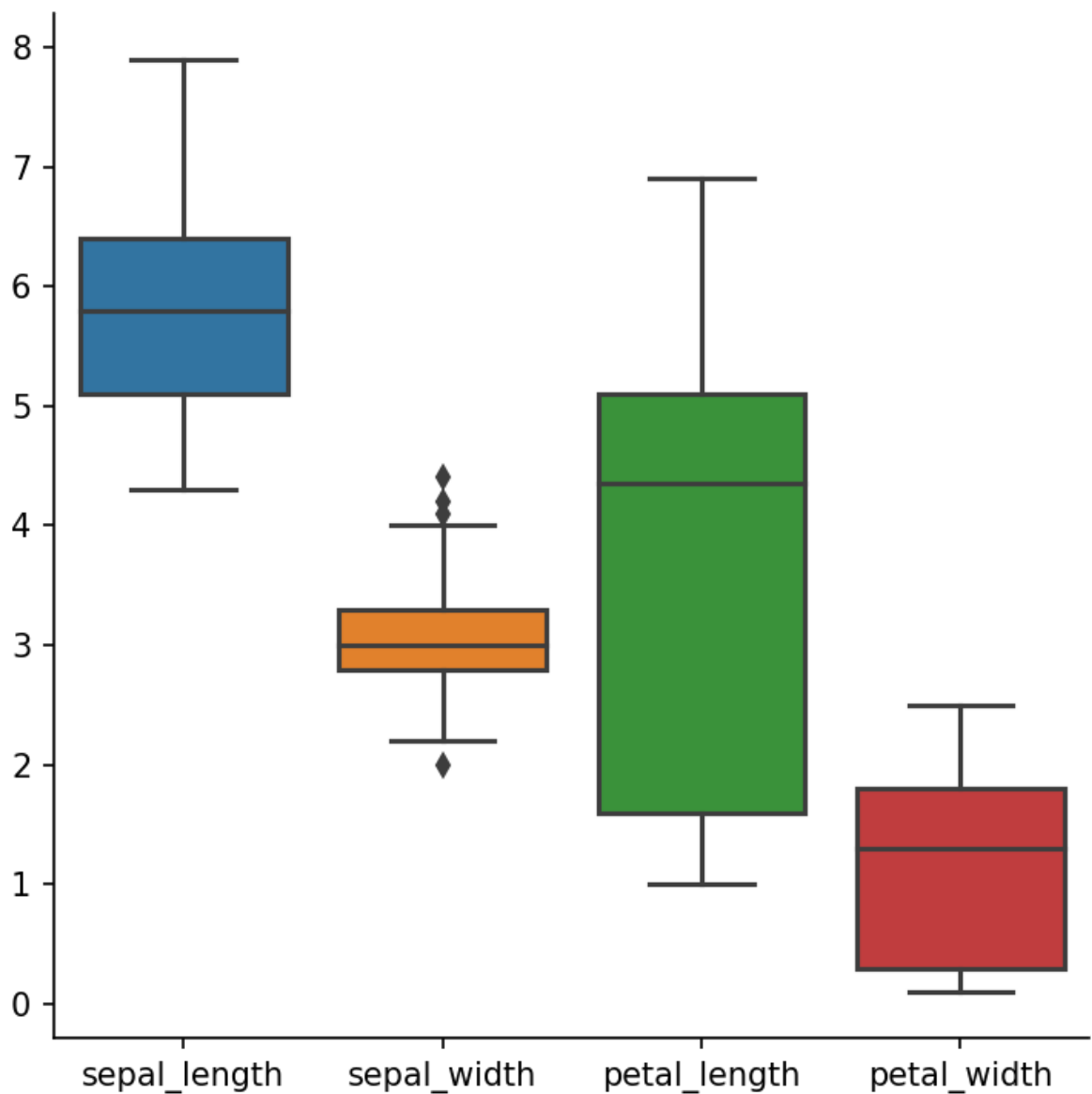
3.

Q 3.1

Load the 'Iris' dataset using seaborn. Create a box plot for the attributes 'sepal_length', 'sepal_width', 'petal_length' and 'petal_width' in the Iris dataset.

```
In [9]: import seaborn as sns
df = sns.load_dataset("iris")

sns.catplot(data=df, kind='box');
```



Q 3.2

In a few sentences explain what can you interpret from the above box plot.

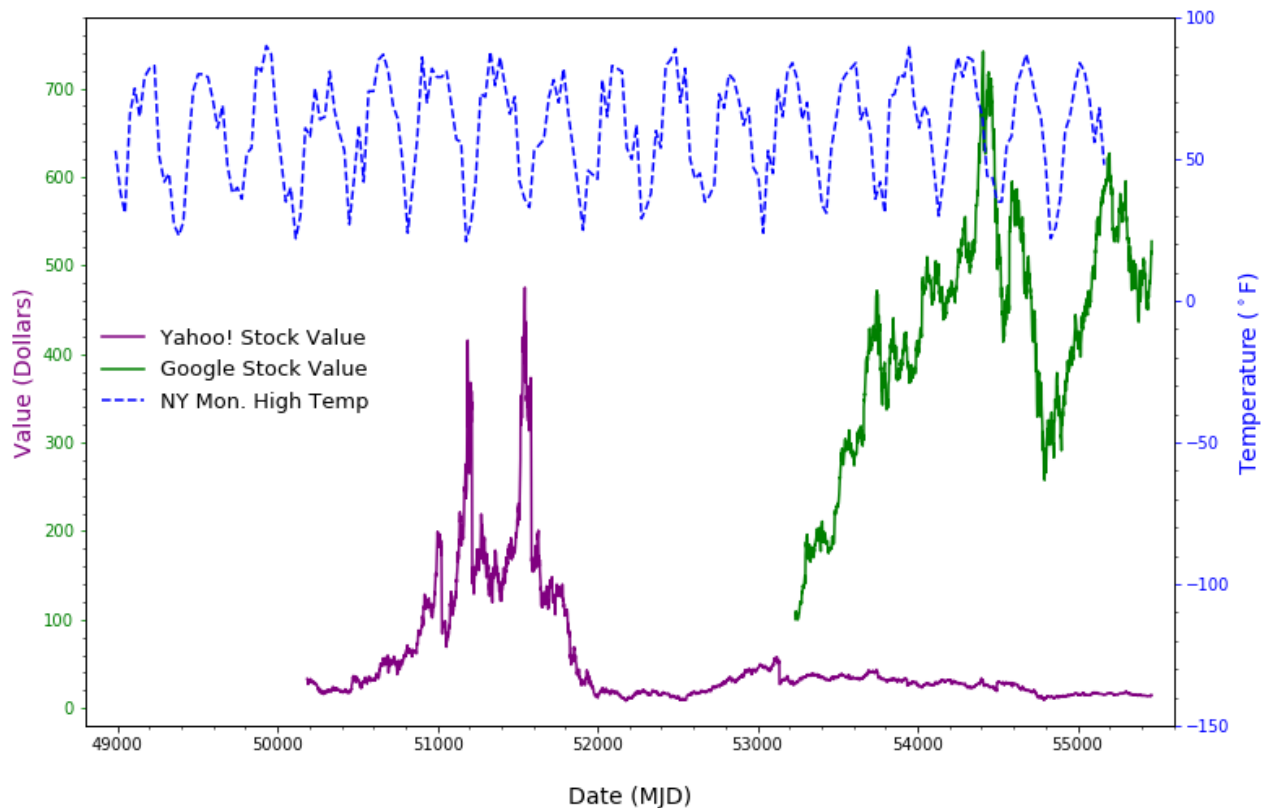
All the data seems relatively symmetrical, there is only a little bit of left skew in the petal length. There are a few outliers in the petal width. The petal length also has a large range in comparison to the other features.

Q 4.

The data files needed:

`google_data.txt`, `ny_temps.txt` & `yahoo_data.txt`

Use your knowledge with `Python`, `NumPy`, `pandas` and `matplotlib` to reproduce the plot below:

New York Temperature, Google, and Yahoo!

```
In [10]: google = pd.read_table("google_data.txt")
ny = pd.read_table("ny_temps.txt")
yahoo = pd.read_table("yahoo_data.txt")
fig, ax1 = plt.subplots()

ax1.plot(google["Modified Julian Date"], google["Stock Value"], c="g", label="Google Stock Value")
ax1.plot(yahoo["Modified Julian Date"], yahoo["Stock Value"], c="purple", label="Yahoo! Stock Value")
ax1.tick_params(axis='y', labelcolor="g", labelsizes=7)

plt.grid(False)
plt.ylabel("Value (Dollars)", color="purple")
plt.xlabel("Date (MJD)")
plt.tick_params(axis='x', labelsizes=7, color="black", length = 5)
plt.minorticks_on()

ax2 = ax1.twinx()
ax2.set_ylim([-150,100])

ax2.plot(ny["Modified Julian Date"], ny["Max Temperature"], "--", c="b", label="NY Mon. High Temp")
ax2.tick_params(axis='y', labelcolor="blue")
fig.legend(loc = "center left", fontsize="x-small", frameon = False, bbox_to_anchor=(0.1, 0.5))
ax2.grid(False)

plt.title("New York Temperature, Google, and Yahoo!", fontweight='bold')
plt.tight_layout()
plt.minorticks_on()
plt.ylabel("Temperature °F", color="blue")

plt.show()
```

