

Visualizing Data in Python

When working with a new dataset, one of the most useful things to do is to begin to visualize the data. By using **tables**, **histograms**, **boxplots**, **scatter plots** and other visual tools, we can get a better idea of what the data may be trying to tell us, and we can gain insights into the data that we may have not discovered otherwise.

In this notebook will use the [Seaborn](https://seaborn.pydata.org/) data processing library, which is a higher-level interface to **Matplotlib** that can be used to simplify many visualization tasks

The **Seaborn** provides visualisations tools that will allow to explore data from a graphical perspective.

Acknowledgments

- Data from <https://www.coursera.org/> from the course "Understanding and Visualizing Data with Python" by University of Michigan

▼ Importing libraries

```
# Import the packages that we will be using
```

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
```

▼ Importing data

```
# Define where you are running the code: colab or local
RunInColab = True # (False: no | True: yes)
```

```
# If running in colab:
```

```
if RunInColab:
```

```
    # Mount your google drive in google colab
    from google.colab import drive
    drive.mount('/content/drive')
```

```
    # Find location
```

```
    #!pwd
```

```
    #!ls
```

```
    #!ls "/content/drive/My Drive/Colab Notebooks/MachineLearningWithPython/"
```

```
    # Define path del proyecto
```

```
    Ruta = "/content/drive/MyDrive/Colab Notebooks/MachineLearningWithPython/"
```

```
else:
```

```
    # Define path del proyecto
```

```
    Ruta = ""
```

```
    Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", fo
```

```
# url string that hosts our .csv file
```

```
url = Ruta + "cartwheel.csv"
```

```
# Read the .csv file and store it as a pandas Data Frame
```

```
df = pd.read_csv(url)
```

▼ Exploring the content of the data set

Get a general 'feel' of the data

```
vars = ["Age", "Gender", "Glasses", "Height", "Wingspan", "CwDistance", "Complete", "Score"]  
df = df[vars].dropna()  
  
df
```

17	27.0	M	Y	66.00	66.0	74	Y	5
18	23.0	M	Y	70.00	69.0	64	Y	3
19	24.0	F	Y	68.00	66.0	85	Y	8
20	23.0	M	Y	69.00	67.0	66	N	2
21	29.0	M	N	71.00	70.0	101	Y	8
22	25.0	M	N	70.00	68.0	82	Y	4
23	26.0	M	N	69.00	71.0	63	Y	5
24	23.0	F	Y	65.00	63.0	67	N	3
25	28.0	M	N	75.00	76.0	111	Y	10
26	24.0	M	N	78.40	71.0	92	Y	7
27	25.0	M	Y	76.00	73.0	107	Y	8
28	32.0	F	Y	63.00	60.0	75	Y	8
29	38.0	F	Y	61.50	61.0	78	Y	7
30	27.0	F	Y	62.00	60.0	72	Y	8
31	33.0	F	Y	65.30	64.0	91	Y	7
32	38.0	F	N	64.00	63.0	86	Y	10
33	27.0	M	N	77.00	75.0	100	Y	8
34	24.0	F	N	67.80	62.0	98	Y	9
35	27.0	M	N	68.00	66.0	74	Y	5
36	25.0	F	Y	65.00	64.5	92	Y	6

▼ Frequency tables

The `value_counts()` method can be used to determine the number of times that each distinct value of a variable occurs in a data set. In statistical terms, this is the "frequency distribution" of the variable. The `value_counts()` method produces a table with two columns. The first column contains all distinct observed values for the variable. The second column contains the number of times each of these values occurs. Note that the table returned by `value_counts()` is actually a **Pandas** data frame, so can be further processed using any Pandas methods for working with data frames.

```
# Number of times that each distinct value of a variable occurs in a data set
df.Complete.value_counts()
```

```
Y    41
N     8
Name: Complete, dtype: int64
```

```
# Proportion of each distinct value of a variable occurs in a data set
x =df.Complete.value_counts()
x/x.sum()
```

```
Y    0.836735
N    0.163265
Name: Complete, dtype: float64
```

Note that the `value_counts()` method excludes missing values. We confirm this below by adding up observations to your data frame with some missing values and then computing `value_counts()` and comparing this to the total number of rows in the data set, which is 28. This tells us that there are $28 - (21+6) = 1$ missing values for this variable (other variables may have different numbers of missing values).

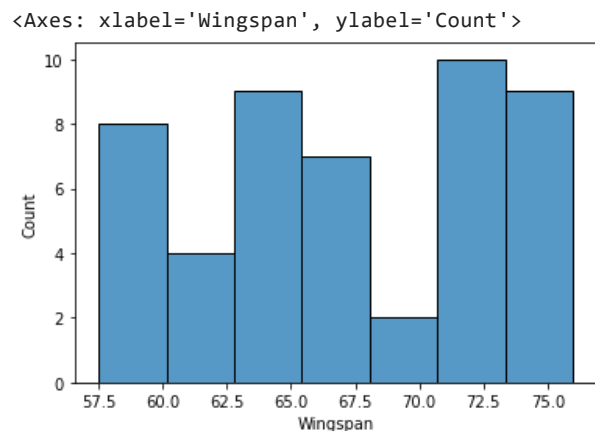
```
# Total number of observations
p= df.Complete.shape
print('Number of observations:', p )
# Total number of null observations
q= df.Complete.isnull
print('Number of null observations:', q )
# Total number of counts (excluding missing values)
r= df.Complete.value_counts
print('Number of counts:',r)

42    Y
43    Y
44    Y
45    Y
46    N
47    N
48    Y
Name: Complete, dtype: object>
Number of counts: <bound method IndexOpsMixin.value_counts of 0    Y
1     Y
2     Y
3     Y
4     N
5     N
6     Y
7     Y
8     N
9     Y
10    Y
11    Y
12    Y
13    Y
14    Y
15    Y
16    N
17    Y
18    Y
19    Y
20    N
21    Y
22    Y
23    Y
24    N
25    Y
26    Y
27    Y
28    Y
29    Y
30    Y
31    Y
32    Y
33    Y
34    Y
35    Y
36    Y
37    Y
38    Y
39    Y
40    Y
41    Y
42    Y
43    Y
44    Y
45    Y
46    N
47    N
48    Y
Name: Complete, dtype: object>
```

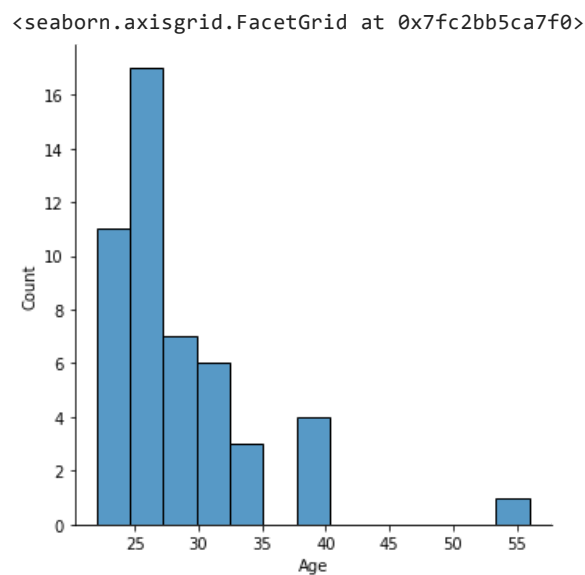
▼ Histogram

It is often good to get a feel for the shape of the distribution of the data.

```
# Plot histogram of the total bill only
sns.histplot(data=df["Wingspan"])
```



```
# Plot distribution of the tips only
sns.displot(df["Age"],kde= False,)
```



```
# Plot histogram of both the Age and the Wingspan
df1 = df.loc[:,["Age","Wingspan"]]
sns.histplot(data=df1)
```

<Axes: ylabel='Count'>



▼ Histograms plotted by groups

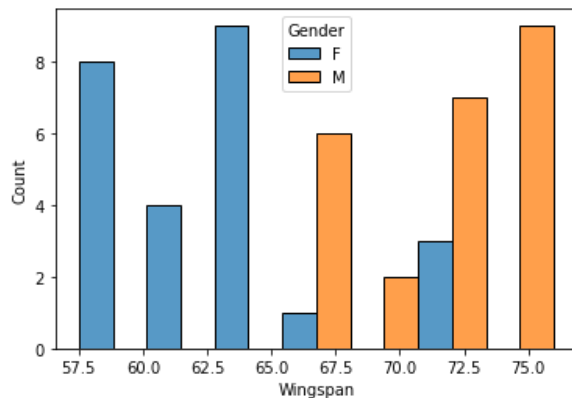
While looking at a single variable is interesting, it is often useful to see how a variable changes in response to another. Thus, we can create a histograms of one quantitative variable grouped by another categorical variables.

```
~ | [blue bar] [orange bar] |
```

```
# Create histograms of the "Wingspan" grouped by "Gender"
```

```
sns.histplot(data=df,x="Wingspan", hue="Gender", multiple="dodge")
```

<Axes: xlabel='Wingspan', ylabel='Count'>



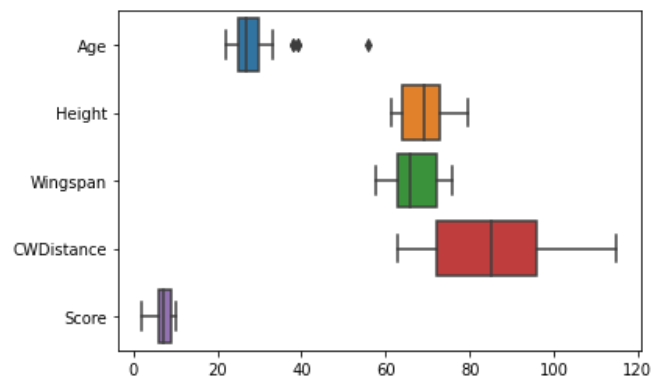
▼ Boxplots

Boxplots do not show the shape of the distribution, but they can give us a better idea about the center and spread of the distribution as well as any potential outliers that may exist. Boxplots and Histograms often complement each other and help an analyst get more information about the data

```
# Create the boxplot of the "total bill" amounts
```

```
sns.boxplot(data=df, orient = 'h')
```

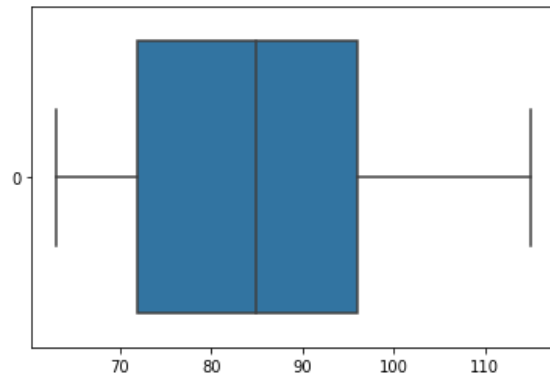
<Axes: >



```
# Create the boxplot of the "tips" amounts
```

```
sns.boxplot(data=df["CWDistance"], orient = 'h')
```

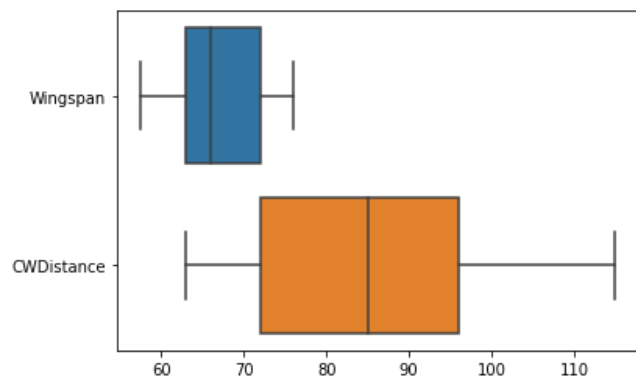
<Axes: >



```
# Create the boxplots of the "Wingspan" and of the "Height" amounts
```

```
sns.boxplot(data=df[["Wingspan","CWDistance"]], orient = 'h')
```

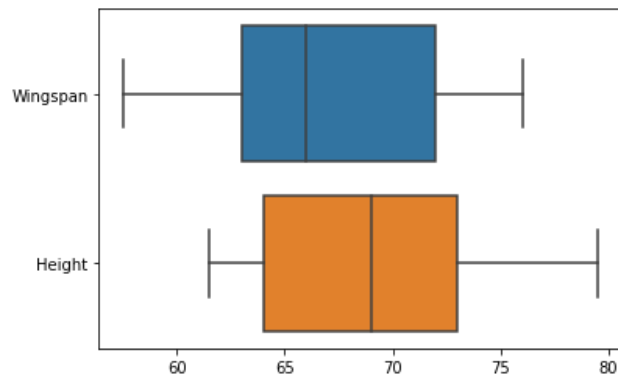
<Axes: >



```
# Create the boxplots of the "Wingspan" and of the "tips" amounts
```

```
sns.boxplot(data=df[["Wingspan","Height"]], orient = 'h')
```

<Axes: >



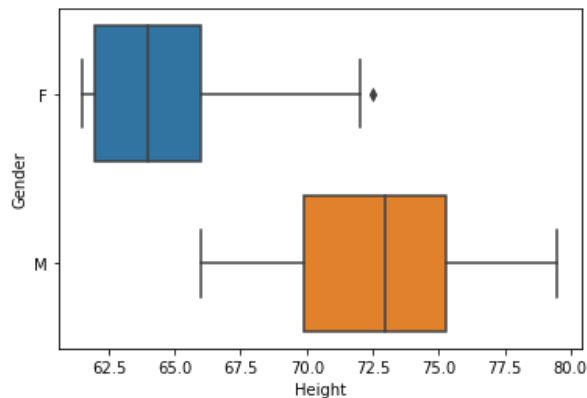
▼ Boxplots plotted by groups

While looking at a single variable is interesting, it is often useful to see how a variable changes in response to another. Thus, we can create a side-by-side boxplots of one quantitative variable grouped by another categorical variables.

```
# Create side-by-side boxplots of the "Height" grouped by "Gender"
```

```
sns.boxplot(data=df, x="Height", y="Gender")
```

<Axes: xlabel='Height', ylabel='Gender'>

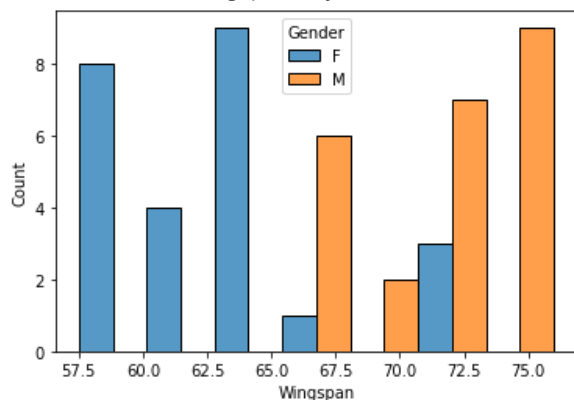


▼ Histograms and boxplots plotted by groups

We can also create both boxplots and histograms of one quantitative variable grouped by another categorical variable

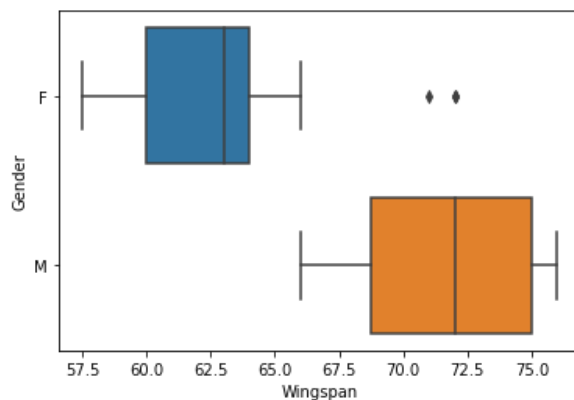
```
# Create a boxplot and histogram of the "tips" grouped by "Gender"
sns.histplot(data=df, x="Wingspan", hue="Gender", multiple="dodge")
```

<Axes: xlabel='Wingspan', ylabel='Count'>



```
sns.boxplot(data=df, x="Wingspan", y="Gender")
```

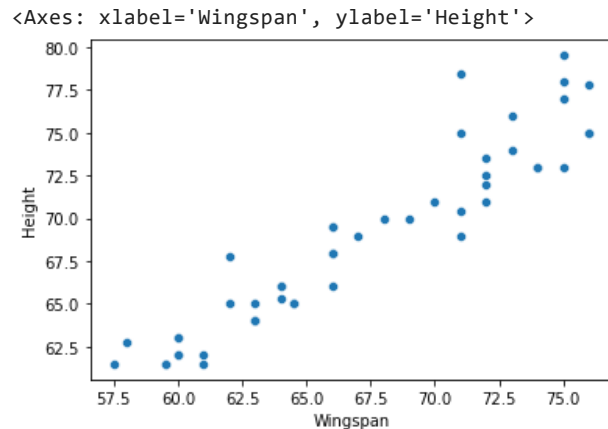
<Axes: xlabel='Wingspan', ylabel='Gender'>



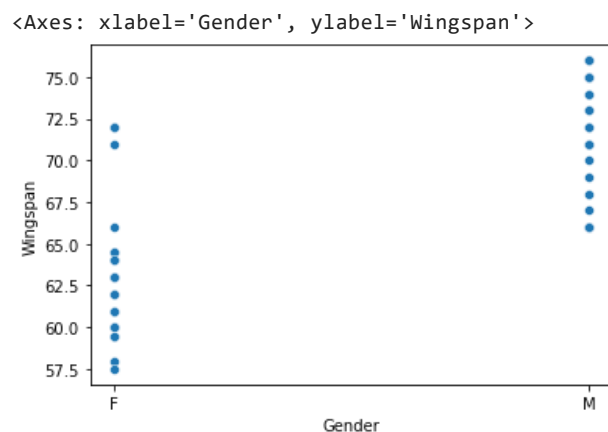
▼ Scatter plot

Plot values of one variable versus another variable to see how they are correlated


```
# scatter plot between two variables
sns.scatterplot(data=df, x="Wingspan", y="Height")
```

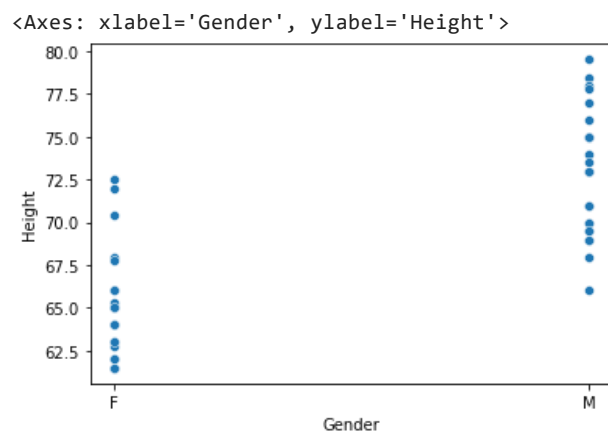


```
# scatter plot between two variables (one categorical)
sns.scatterplot(data=df, x="Gender", y="Wingspan")
```



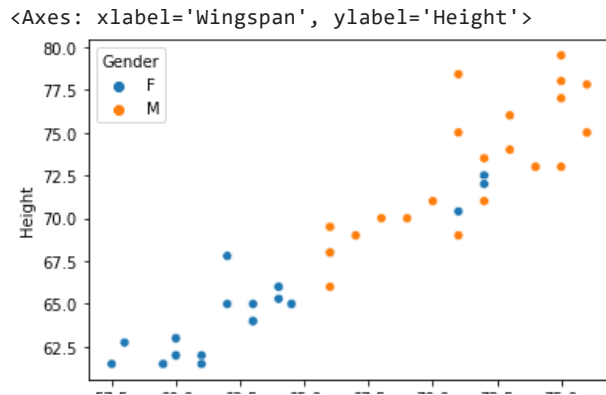
```
# scatter plot between two variables (one categorical)
```

```
sns.scatterplot(data=df, y="Height", x="Gender")
```



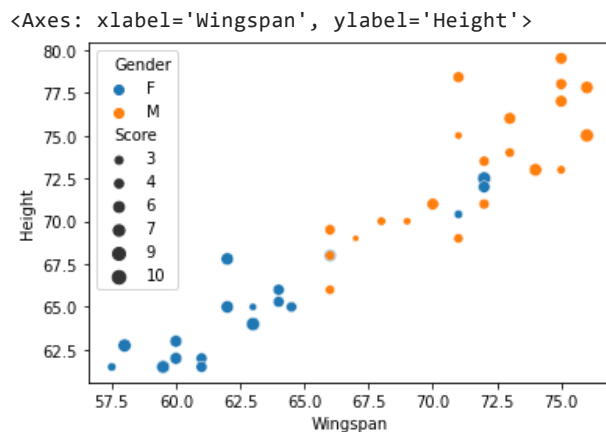
```
# scatter plot between two variables grouped according to a categorical variable
```

```
sns.scatterplot(data=df, x="Wingspan", y="Height", hue="Gender")
```



scatter plot between two variables grouped according to a categorical variable and with size of markers

```
sns.scatterplot(data=df, x="Wingspan", y="Height", hue="Gender", size="Score")
```



Final remarks

- Visualizing your data using **tables, histograms, boxplots, scatter plots** and other tools is essential to carry put analysis and extract conclusions
- There are several ways to do the same thing
- The **Seaborn** package provides visualisations tools that allow to explore data from a graphical perspective

▼ Activity: work with the iris dataset

Repeat this tutorial with the iris data set and respond to the following inquiries

1. Plot the histograms for each of the four quantitative variables
2. Plot the histograms for each of the quantitative variables
3. Plot the boxplots for each of the quantitative variables
4. Plot the boxplots of the petal width grouped by type of flower
5. Plot the boxplots of the setal length grouped by type of flower
6. Provide a description (explanation from your observations) of each of the quantitative variables