

---

Daniel Estrada Ocaña A01369854

---

## 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

```
1 # Import the packages that we will be using
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import pandas as pd
5
6
7
```

### ✓ Importing data

```
1 # Define where you are running the code: colab or local
2 RunInColab = True # (False: no | True: yes)
3
4 # If running in colab:
5 if RunInColab:
6     # Mount your google drive in google colab
7     from google.colab import drive
8     drive.mount('/content/drive')
9
10    # Find location
11    #!pwd
12    #!ls
13    #!ls "/content/drive/My Drive/Colab Notebooks/MachineLearningWithPython/"
14
15    # Define path del proyecto
16    Ruta = "/content/drive/MyDrive/Sistemas/4to_semestre/semanaTec/TC1002S"
17
18 else:
19     # Define path del proyecto
20     Ruta = ""
21
22    Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
23
24
25 1 # url string that hosts our .csv file
26 2 url = Ruta + "/NotebooksProfessor/datasets/iris/iris.csv"
27 3
28 4
29 5 # Read the .csv file and store it as a pandas Data Frame
30 6
31 7 df = pd.read_csv(url )
32 8
```

### ✓ Exploring the content of the data set

Get a general 'feel' of the data

1 df

	Ms_1	Ms_2	Ms_3	Ms_4	Type
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
...	...	...	...	...	...
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

150 rows × 5 columns

Next steps:

☒ View recommended plots

1 df.shape

(150, 5)

1 nrows = df.shape[0]

2 nrows

150

1 ncols = df.shape[1]

2 ncols

5

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.

1 # Number of times that each distinct value of a variable occurs in a data set

2 ms1 = df.value\_counts()

3 ms1

Ms_1	Ms_2	Ms_3	Ms_4	Type	
5.8	2.7	5.1	1.9	Iris-virginica	2
6.2	2.2	4.5	1.5	Iris-versicolor	1
	2.9	4.3	1.3	Iris-versicolor	1
	3.4	5.4	2.3	Iris-virginica	1
6.3	2.3	4.4	1.3	Iris-versicolor	1
				..	
5.4	3.9	1.3	0.4	Iris-setosa	1
		1.7	0.4	Iris-setosa	1
5.5	2.3	4.0	1.3	Iris-versicolor	1
	2.4	3.7	1.0	Iris-versicolor	1
7.9	3.8	6.4	2.0	Iris-virginica	1

Length: 149, dtype: int64

```

1 # Proportion of each distinct value of a variable occurs in a data set
2 propor = df.value_counts() / len(df)
3 propor

```

```

Ms_1  Ms_2  Ms_3  Ms_4  Type
5.8    2.7    5.1    1.9  Iris-virginica    0.013333
6.2    2.2    4.5    1.5  Iris-versicolor  0.006667
      2.9    4.3    1.3  Iris-versicolor  0.006667
      3.4    5.4    2.3  Iris-virginica    0.006667
6.3    2.3    4.4    1.3  Iris-versicolor  0.006667
      ...
5.4    3.9    1.3    0.4  Iris-setosa      0.006667
      1.7    0.4    0.4  Iris-setosa      0.006667
5.5    2.3    4.0    1.3  Iris-versicolor  0.006667
      2.4    3.7    1.0  Iris-versicolor  0.006667
7.9    3.8    6.4    2.0  Iris-virginica    0.006667
Length: 149, 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).

```

1 # Total number of observations
2
3 obser = len(df)
4 # Total number of null observations
5
6 nobser = df.isnull().sum()
7
8 # Total number of counts (excluding missing values)
9 tot = obser-nobser
10
11 tot

Ms_1    150
Ms_2    150
Ms_3    150
Ms_4    150
Type     150
dtype: int64

```

**There are no missing values**

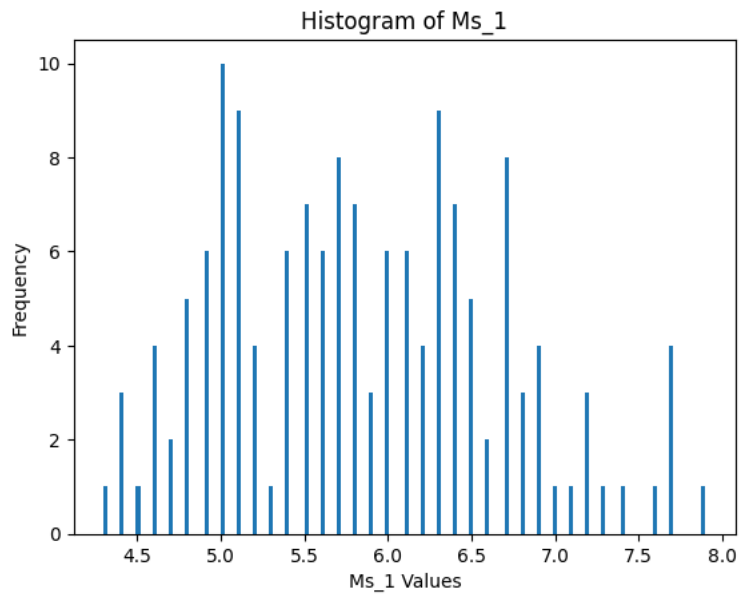
## ✓ Histogram

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

```

1 # Plot histogram of the Ms_1 only
2 plt.hist(df['Ms_1'], bins=len(df)) # Adjust the number of bins as needed
3 plt.title('Histogram of Ms_1')
4 plt.xlabel('Ms_1 Values')
5 plt.ylabel('Frequency')
6 plt.show()
7

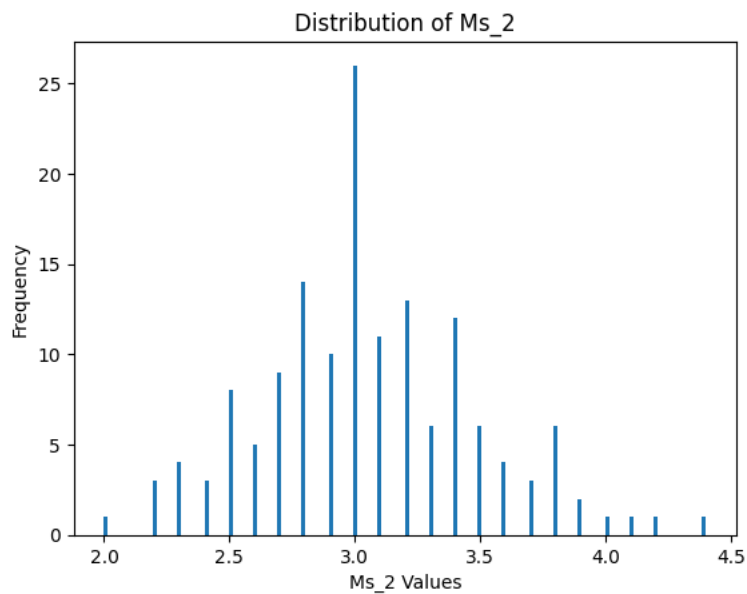
```



```

1 # Plot distribution of the Ms_2 only
2 plt.hist(df['Ms_2'], bins=len(df)) # Adjust the number of bins as needed
3 plt.title('Distribution of Ms_2')
4 plt.xlabel('Ms_2 Values')
5 plt.ylabel('Frequency')
6 plt.show()
7

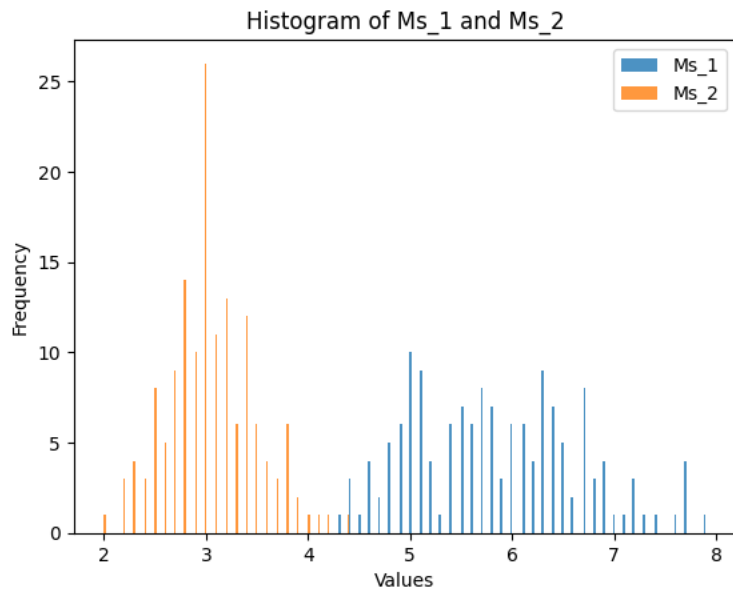
```



```

1 # Plot histogram of both the Ms_1 and the Ms_2
2 plt.hist(df['Ms_1'], bins=len(df), alpha=0.8, label='Ms_1') # Adjust the number of bins as needed
3 plt.hist(df['Ms_2'], bins=len(df), alpha=0.8, label='Ms_2') # Adjust the number of bins as needed
4 plt.title('Histogram of Ms_1 and Ms_2')
5 plt.xlabel('Values')
6 plt.ylabel('Frequency')
7 plt.legend()
8 plt.show()
9

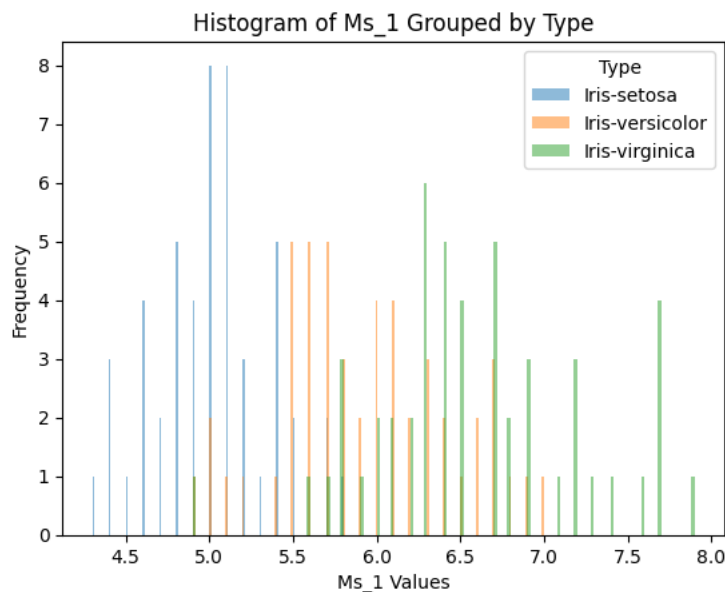
```



## ✓ 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.

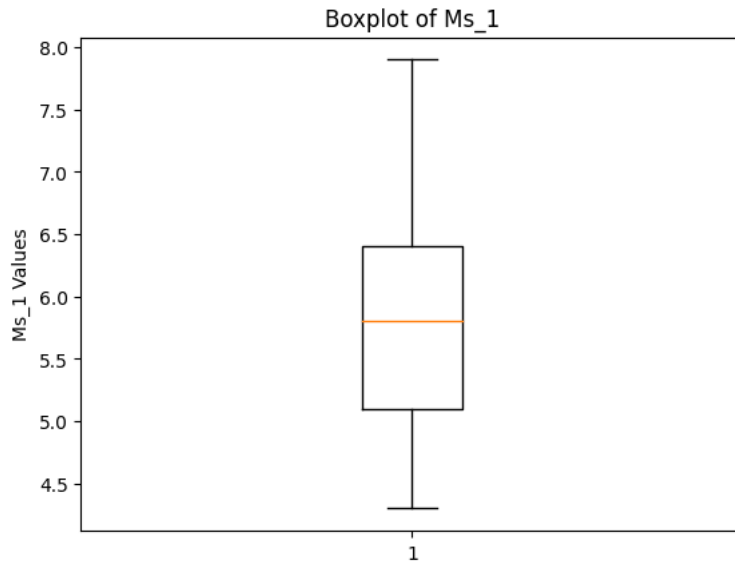
```
1 # Create histograms of the "Ms_1" grouped by "Type"
2
3 for type_name, group_data in df.groupby('Type'):
4     plt.hist(group_data['Ms_1'], bins=len(df), alpha=0.5, label=type_name)
5
6 # Add labels and legend
7 plt.title('Histogram of Ms_1 Grouped by Type')
8 plt.xlabel('Ms_1 Values')
9 plt.ylabel('Frequency')
10 plt.legend(title='Type')
11 plt.show()
12
```



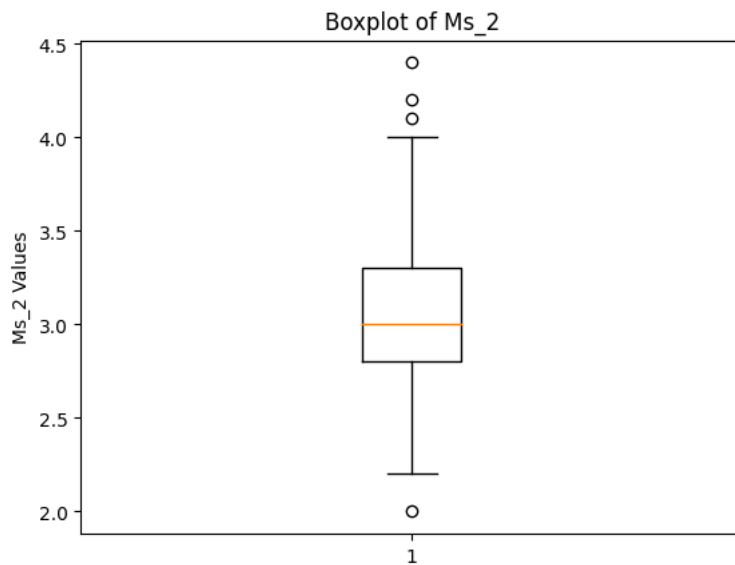
## ✓ Boxplots

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

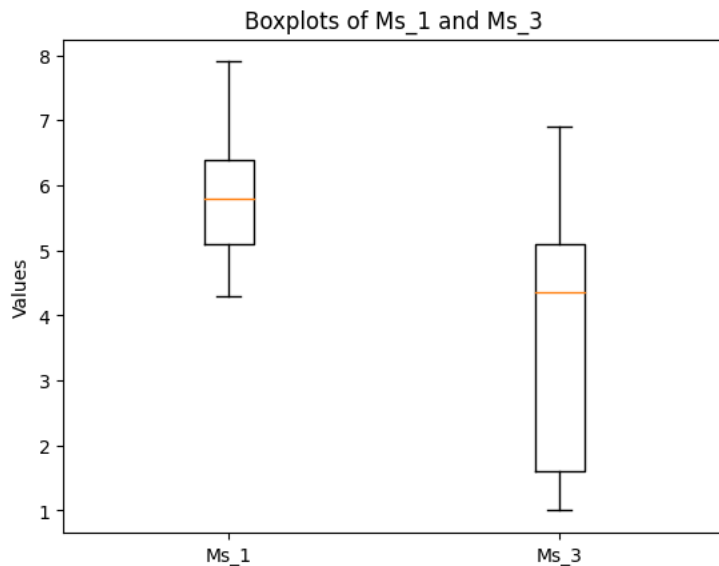
```
1 # Create the boxplot of the "Ms_1" amounts
2 plt.boxplot(df['Ms_1'])
3 plt.title('Boxplot of Ms_1')
4 plt.ylabel('Ms_1 Values')
5 plt.show()
6
```



```
1 # Create the boxplot of the "tips" amounts
2
3 plt.boxplot(df['Ms_2'])
4 plt.title('Boxplot of Ms_2')
5 plt.ylabel('Ms_2 Values')
6 plt.show()
```



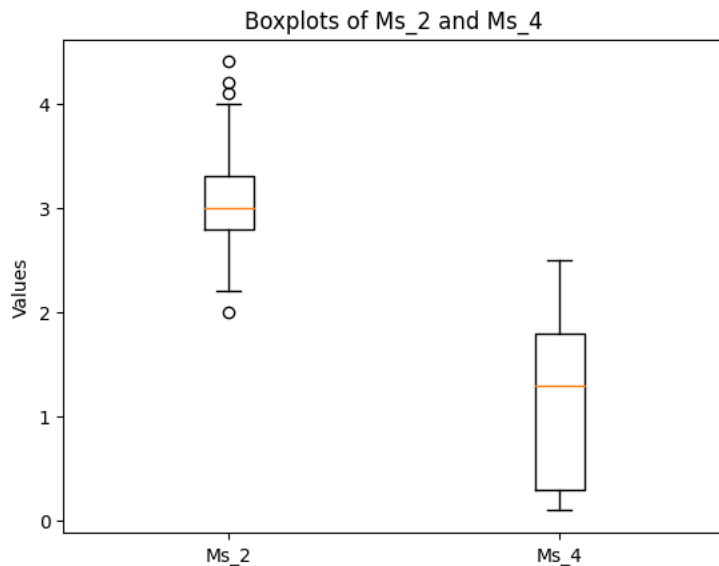
```
1 # Create the boxplots of the "Ms_1" and of the "Ms_3" amounts
2
3 plt.boxplot([df['Ms_1'], df['Ms_3']])
4 plt.title('Boxplots of Ms_1 and Ms_3')
5 plt.ylabel('Values')
6 plt.xticks([1, 2], ['Ms_1', 'Ms_3']) # Assign labels to x-axis ticks
7 plt.show()
```



```

1 # Create the boxplots of the "Ms_2" and of the "Ms_4" amounts
2 plt.boxplot([df['Ms_2'], df['Ms_4']])
3 plt.title('Boxplots of Ms_2 and Ms_4')
4 plt.ylabel('Values')
5 plt.xticks([1, 2], ['Ms_2', 'Ms_4']) # Assign labels to x-axis ticks
6 plt.show()
7

```



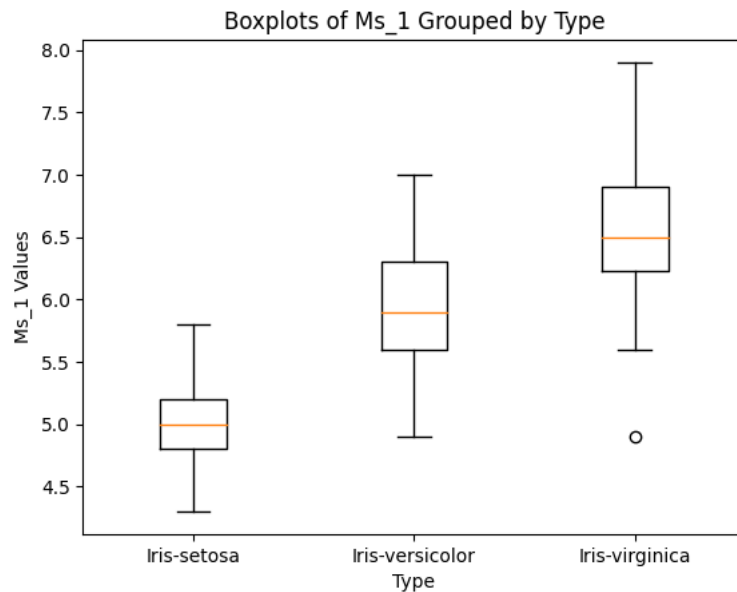
## ✓ Boxplots plotted by groups

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.

```

1 # Create side-by-side boxplots of the "Ms_1" grouped by "Type"
2
3 grouped_data = [group['Ms_1'] for name, group in df.groupby('Type')]
4
5 # Create side-by-side boxplots
6 plt.boxplot(grouped_data, labels=df['Type'].unique())
7 plt.title('Boxplots of Ms_1 Grouped by Type')
8 plt.xlabel('Type')
9 plt.ylabel('Ms_1 Values')
10 plt.show()

```



## ✓ Histograms and boxplots plotted by groups

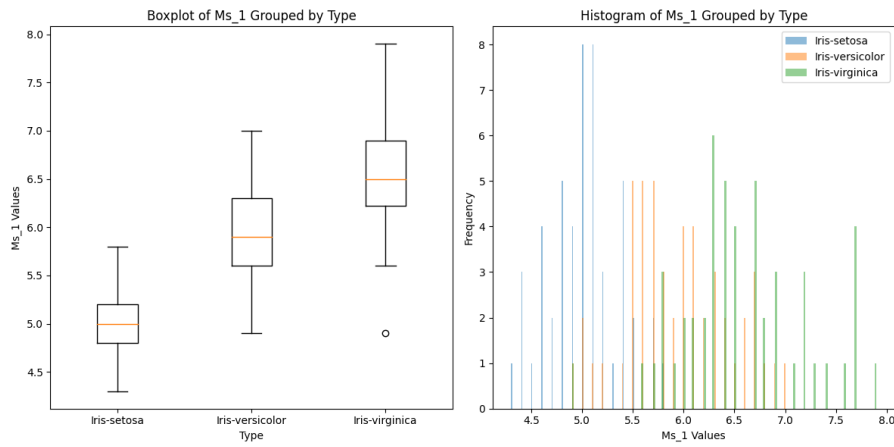
We can also create both boxplots and histograms of one quantitative variable grouped by another categorical variable.

```

1 # Create a boxplot and histogram of the "Ms_1" grouped by "Type"
2 # Create subplots with 1 row and 2 columns
3 fig, axs = plt.subplots(1, 2, figsize=(12, 6))
4
5 # Group the data by 'Type'
6 grouped_data = [group['Ms_1'] for name, group in df.groupby('Type')]
7
8 # Create boxplot
9 axs[0].boxplot(grouped_data, labels=df['Type'].unique())
10 axs[0].set_title('Boxplot of Ms_1 Grouped by Type')
11 axs[0].set_xlabel('Type')
12 axs[0].set_ylabel('Ms_1 Values')
13
14 # Create histogram
15 for name, group in df.groupby('Type'):
16     axs[1].hist(group['Ms_1'], bins=len(df), alpha=0.5, label=name)
17 axs[1].set_title('Histogram of Ms_1 Grouped by Type')
18 axs[1].set_xlabel('Ms_1 Values')
19 axs[1].set_ylabel('Frequency')
20 axs[1].legend()
21
22 plt.tight_layout()
23 plt.show()
24

```

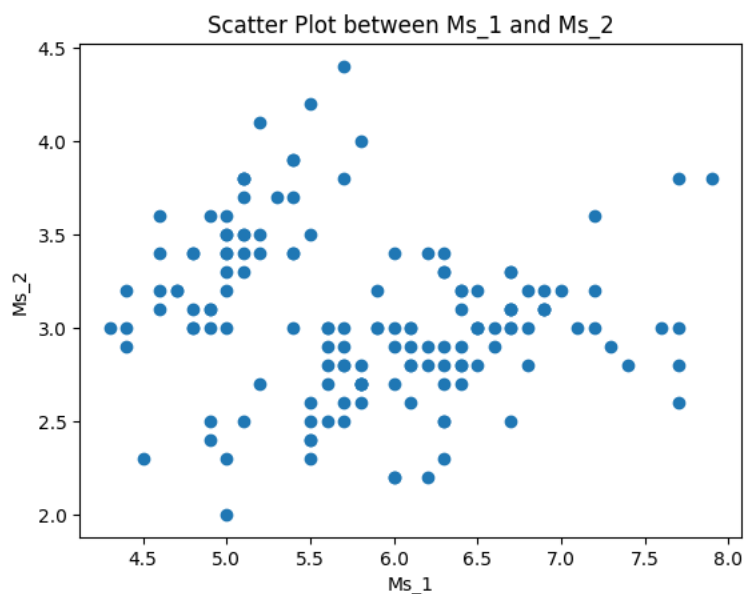




## ✓ Scatter plot

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

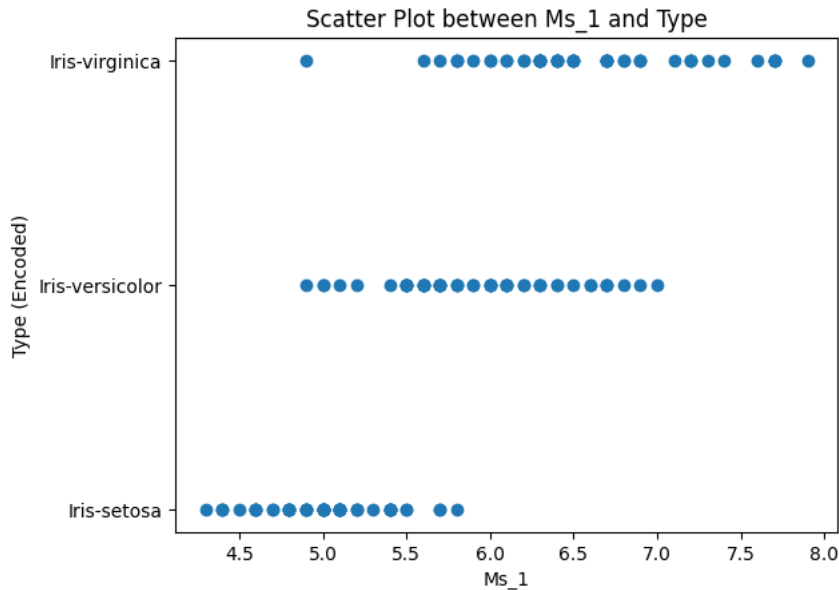
```
1 # scatter plot between two variables
2
3 plt.scatter(df['Ms_1'], df['Ms_2'])
4 plt.title('Scatter Plot between Ms_1 and Ms_2')
5 plt.xlabel('Ms_1')
6 plt.ylabel('Ms_2')
7 plt.show()
```



```

1 # scatter plot between two variables (one categorical)
2 df["Type encoded"] = df.Type.replace({"Iris-setosa": 1, "Iris-versicolor": 2, "Iris-virginica": 3})
3
4 plt.scatter(df['Ms_1'], df['Type encoded'])
5 plt.title('Scatter Plot between Ms_1 and Type')
6 plt.xlabel('Ms_1')
7 plt.ylabel('Type (Encoded)')
8 plt.yticks(df['Type encoded'].unique(), df['Type'].unique())
9 plt.show()
10

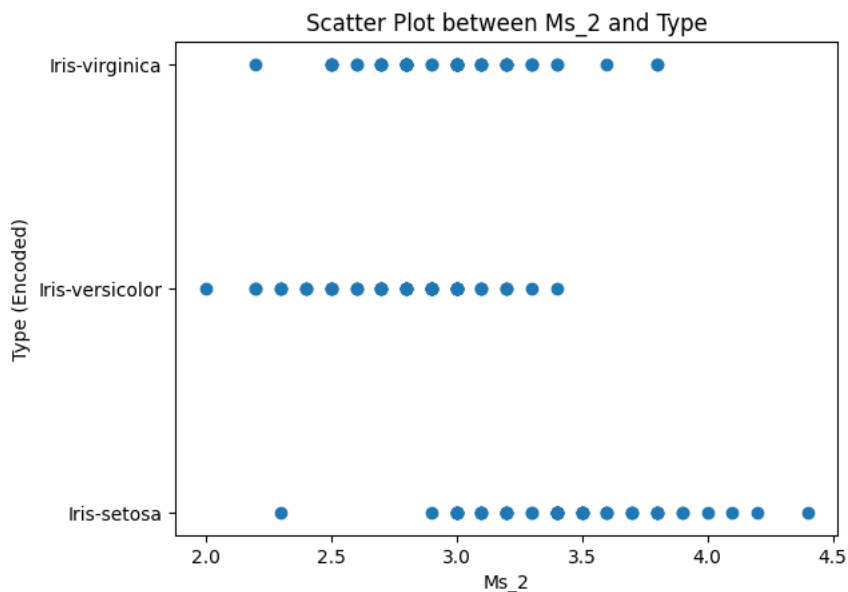
```



```

1 # scatter plot between two variables (one categorical)
2 df["Type encoded"] = df.Type.replace({"Iris-setosa": 1, "Iris-versicolor": 2, "Iris-virginica": 3})
3
4 plt.scatter(df['Ms_2'], df['Type encoded'])
5 plt.title('Scatter Plot between Ms_2 and Type')
6 plt.xlabel('Ms_2')
7 plt.ylabel('Type (Encoded)')
8 plt.yticks(df['Type encoded'].unique(), df['Type'].unique())
9 plt.show()
10

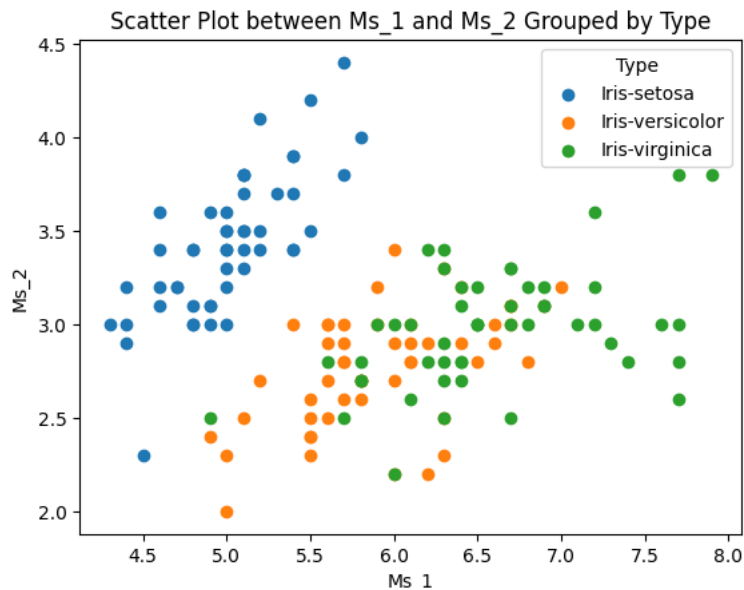
```



```

1 # scatter plot between two variables grouped according to a categorical variable
2
3 grouped_data = df.groupby('Type')
4
5
6 for name, group in grouped_data:
7     plt.scatter(group['Ms_1'], group['Ms_2'], label=name)
8
9
10 plt.title('Scatter Plot between Ms_1 and Ms_2 Grouped by Type')
11 plt.xlabel('Ms_1')
12 plt.ylabel('Ms_2')
13 plt.legend(title='Type')
14 plt.show()

```

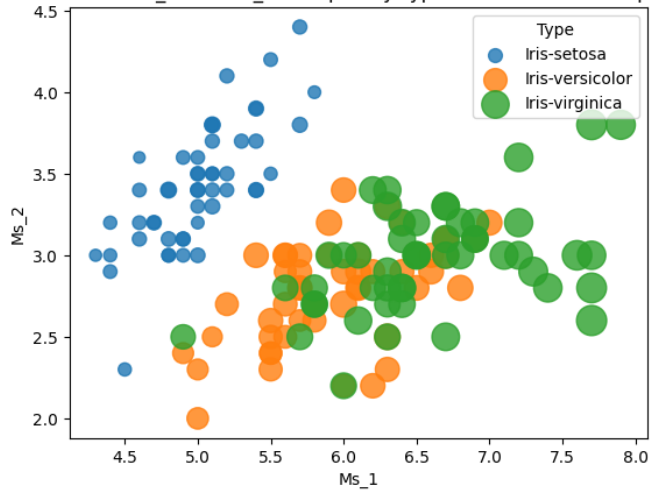


```

1 # scatter plot between two variables grouped according to a categorical variable and with size of markers
2
3
4 grouped_data = df.groupby('Type')
5
6
7 for name, group in grouped_data:
8     plt.scatter(group['Ms_1'], group['Ms_2'], s=group['Ms_3']*50, alpha=0.8, label=name)
9
10
11 plt.title('Scatter Plot between Ms_1 and Ms_2 Grouped by Type with Marker Size representing Ms_3')
12 plt.xlabel('Ms_1')
13 plt.ylabel('Ms_2')
14 plt.legend(title='Type')
15 plt.show()

```

Scatter Plot between Ms\_1 and Ms\_2 Grouped by Type with Marker Size representing Ms\_3



## Final remarks

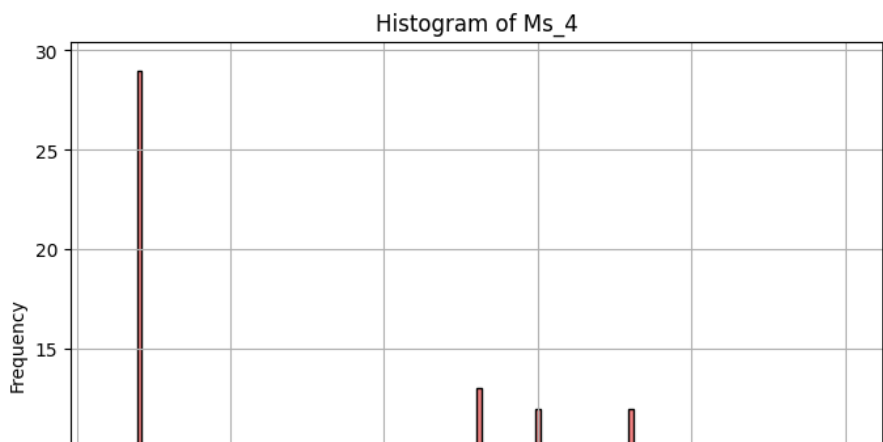
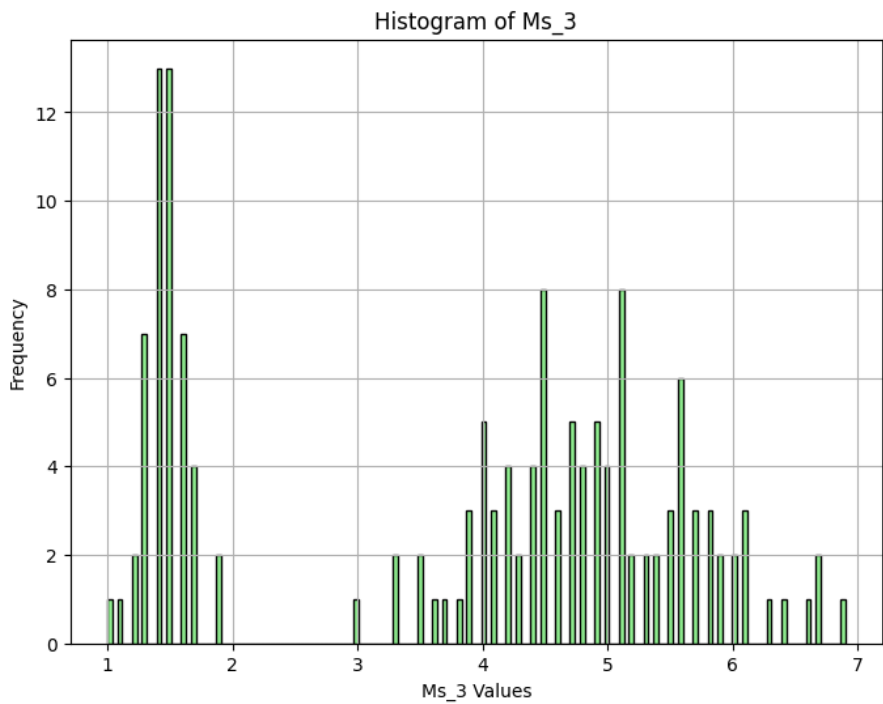
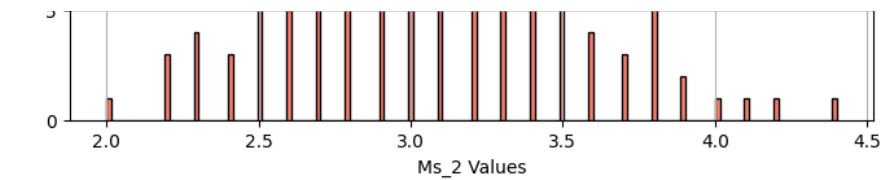
- Visualizing your data using **tables**, **histograms**, **boxplots**, **scatter plots** and other tools is essential to carry put analysis and extract conclusions
- There are several ways to do the same thing
- The **Seaborn** package provides visualisations tools that allow to explore data from a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

```
1 # Histogram of Ms_1
2 plt.figure(figsize=(8, 6))
3 plt.hist(df['Ms_1'], bins=len(df), color='skyblue', edgecolor='black')
4 plt.title('Histogram of Ms_1')
5 plt.xlabel('Ms_1 Values')
6 plt.ylabel('Frequency')
7 plt.grid(True)
8 plt.show()
9
10 # Histogram of Ms_2
11 plt.figure(figsize=(8, 6))
12 plt.hist(df['Ms_2'], bins=len(df), color='salmon', edgecolor='black')
13 plt.title('Histogram of Ms_2')
14 plt.xlabel('Ms_2 Values')
15 plt.ylabel('Frequency')
16 plt.grid(True)
17 plt.show()
18
19 # Histogram of Ms_3
20 plt.figure(figsize=(8, 6))
21 plt.hist(df['Ms_3'], bins=len(df), color='lightgreen', edgecolor='black')
22 plt.title('Histogram of Ms_3')
23 plt.xlabel('Ms_3 Values')
24 plt.ylabel('Frequency')
25 plt.grid(True)
26 plt.show()
27
28 # Histogram of Ms_4
29 plt.figure(figsize=(8, 6))
30 plt.hist(df['Ms_4'], bins=len(df), color='lightcoral', edgecolor='black')
31 plt.title('Histogram of Ms_4')
32 plt.xlabel('Ms_4 Values')
33 plt.ylabel('Frequency')
34 plt.grid(True)
35 plt.show()
```



## 2. Plot the histograms for each of the quantitative variables

```
1 fig, axs = plt.subplots(2, 2, figsize=(12, 8))
2
3 # Flatten the axs array for easier iteration
4 axs = axs.flatten()
5
6 # Plot histograms for each variable
7 for i, column in enumerate(['Ms_1', 'Ms_2', 'Ms_3', 'Ms_4']):
8     axs[i].hist(df[column], bins=len(df), color='skyblue', edgecolor='black')
9     axs[i].set_title(f'Histogram of {column}')
10    axs[i].set_xlabel(f'{column} Values')
11    axs[i].set_ylabel('Frequency')
12
13 # Adjust layout
14 plt.tight_layout()
15 plt.show()
```