

Introduction

- Making plots and static or interactive visualizations is one of the most important tasks in data analysis. It may be a part of the exploratory process; for example, helping identify outliers, needed data transformations, or coming up with ideas for models.
- Matplotlib is the most extensively used library of python for data visualization due to its high flexibility and extensive functionality that it provides.

Table of Contents

1. Setting up
 - Importing matplotlib
 - Matplotlib for Jupyter notebook
 - Dataset
 - Documentation
2. Matplotlib basics
 - Make a simple plot
 - Labels, and Legends
 - Size, Colors, Markers, and Line Styles
 - Figures and subplots
3. Line Chart
4. Bar Chart
5. Histogram
6. Box plot
7. Violin plot
8. Scatter plot
9. Bubble plot

1. Setting up

Importing matplotlib

Just as we use the `np` shorthand for NumPy and the `pd` shorthand for Pandas, we will use standard shorthands for Matplotlib import:

```
import matplotlib.pyplot as plt
```

We import the `pyplot` interface of matplotlib with a shorthand of `plt` and we will be using it like this in the entire notebook.

Matplotlib for Jupyter notebook

You can directly use matplotlib with this notebook to create different visualizations in the notebook itself. In order to do that, the following command is used:

```
%matplotlib inline
```

Documentation

All the functions covered in this notebook and their detail description can be found in the [official matplotlib documentation](#).

```
In [1]: 1 # importing required libraries
2 import numpy as np
3 import pandas as pd
4
5 # importing matplotlib
6 import matplotlib.pyplot as plt
7
8 # display plots in the notebook itself
9 %matplotlib inline
```

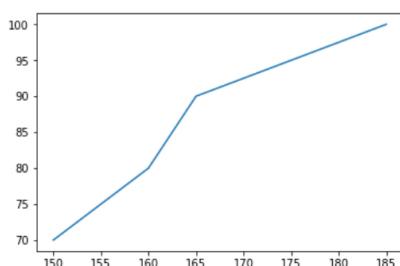
2. Matplotlib basics

Make a simple plot

Let's create a basic plot to start working with!

```
In [14]: 1 height = [150,160,165,185]
2 weight = [70, 80, 90, 100]
3
4 # draw the plot
5 plt.plot(height, weight)
```

```
Out[14]: [<matplotlib.lines.Line2D at 0x7f6caeecb080>]
```



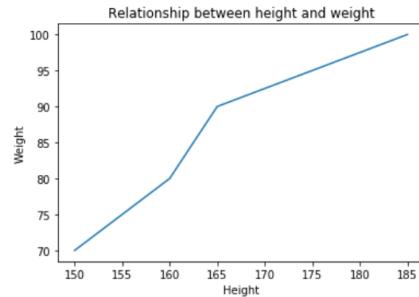
We pass two arrays as our input arguments to `plot()` method and invoke the required plot. Here note that the first array appears on the x-axis and second array appears on the y-axis of the plot.

Title, Labels, and Legends

- Now that our first plot is ready, let us add the title, and name x-axis and y-axis using methods `title()`, `xlabel()` and `ylabel()` respectively.

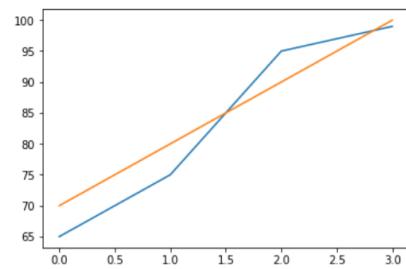
```
In [15]: 1 # draw the plot
2 plt.plot(height,weight)
3 # add title
4 plt.title("Relationship between height and weight")
5 # label x axis
6 plt.xlabel("Height")
7 # label y axis
8 plt.ylabel("Weight")
```

Out[15]: `Text(0, 0.5, 'Weight')`



```
In [16]: 1 calories_burnt = [65, 75, 95, 99]
2
3 # draw the plot for calories burnt
4 plt.plot(calories_burnt)
5 # draw the plot for weight
6 plt.plot(weight)
```

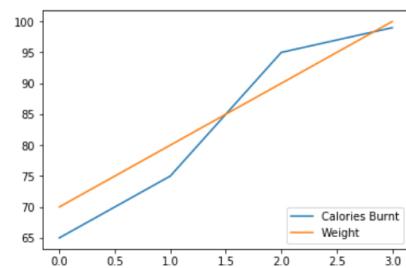
Out[16]: `[<matplotlib.lines.Line2D at 0x7f6cae0de48>]`



- Adding **legends** is also simple in matplotlib, you can use the `legend()` which takes **labels** and **loc** as label names and location of legend in the figure as parameters.

```
In [17]: 1 # draw the plot for calories burnt
2 plt.plot(calories_burnt)
3 # draw the plot for weight
4 plt.plot(weight)
5
6 # add legend in the lower right part of the figure
7 plt.legend(labels=['Calories Burnt', 'Weight'], loc='lower right')
```

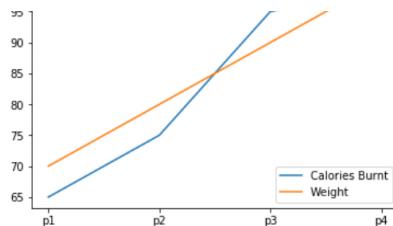
Out[17]: `<matplotlib.legend.Legend at 0x7f6caed72da0>`



- Notice that in the previous plot, we are not able to understand that each of these values belong to different persons.
- Look at the X axis, can we add labels to show that each belong to different persons?
- The labeled values on any axis is known as a **tick**.
- You can use the `xticks` to change both the location of each tick and it's label. Let's see this in an example

```
In [33]: 1 # draw the plot
2 plt.plot(calories_burnt)
3 plt.plot(weight)
4
5 # add legend in the lower right part of the figure
6 plt.legend(labels=['Calories Burnt', 'Weight'], loc='lower right')
7
8 # set labels for each of these persons
9 plt.xticks(ticks=[0,1,2,3], labels=['p1', 'p2', 'p3', 'p4']);
```

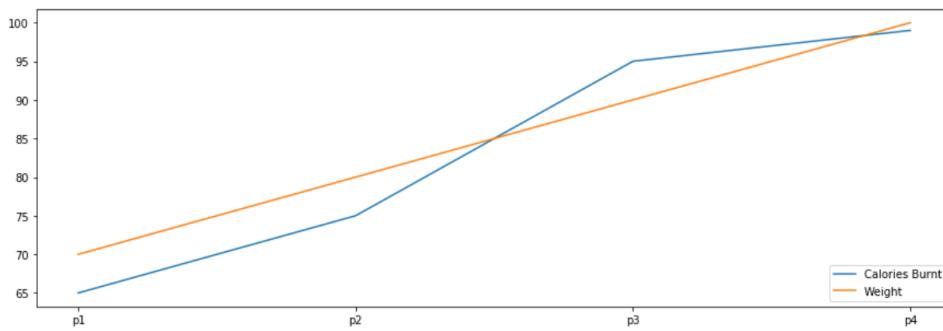




Size, Colors, Markers and Line styles

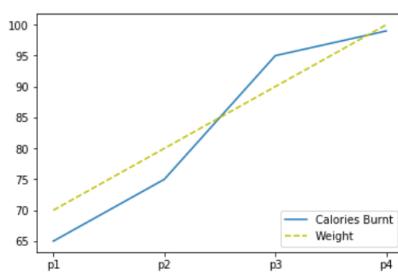
- You can also specify the size of the figure using method `figure()` and passing the values as a tuple of the length of rows and columns to the argument `figsize`.
- The values of length are considered to be in **inches**.

```
In [34]: 1 # figure size in inches
2 plt.figure(figsize=(15,5))
3
4 # draw the plot
5 plt.plot(calories_burnt)
6 plt.plot(weight)
7
8 # add Legend in the lower right part of the figure
9 plt.legend(labels=['Calories Burnt', 'Weight'], loc='lower right')
10
11 # set Labels for each of these persons
12 plt.xticks(ticks=[0,1,2,3], labels=['p1', 'p2', 'p3', 'p4']);
```



- With every X and Y argument, you can also pass an optional third argument in the form of a string which indicates the colour and line type of the plot.
- The default format is `b-` which means a **solid blue line**. In the figure below we use `go` which means **green circles**. Likewise, we can make many such combinations to format our plot.

```
In [31]: 1 # draw the plot
2 plt.plot(calories_burnt)
3 plt.plot(weight, 'y--')
4
5 # add Legend in the lower right part of the figure
6 plt.legend(labels=['Calories Burnt', 'Weight'], loc='lower right')
7
8 # set labels for each of these persons
9 plt.xticks(ticks=[0,1,2,3], labels=['p1', 'p2', 'p3', 'p4']);
```



- We can also plot multiple sets of data by passing in multiple sets of X and Y axis in the `plot()` method as shown.

Figure and subplots

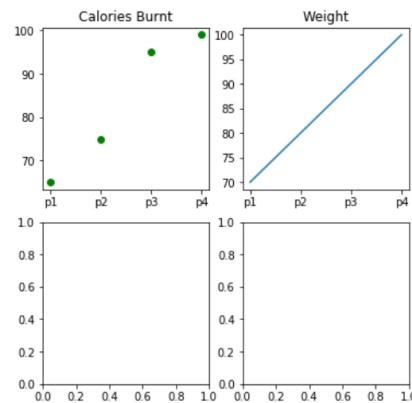
- We can use `subplots()` method to add more than one plots in one figure.
- The `subplots()` method takes two arguments: they are `nrows`, `ncols`. They indicate the number of rows, number of columns respectively.
- This method creates two objects: **figure** and **axes** which we store in variables `fig` and `ax`.
- You plot each figure by specifying its position using row index and column index. Let's have a look at the below example:

```
In [60]: 1 # create 2 plots
2 fig, ax = plt.subplots(nrows=2, ncols=2, figsize=(6,6))
3
4 # plot on 0 row and 0 column
5 ax[0,0].plot(calories_burnt,'go')
6
7 # plot on 0 row and 1 column
8 ax[0,1].plot(weight)
9
10 # set titles for subplots
11 ax[0,0].set_title("Calories Burnt")
```

```

12 ax[0,1].set_title("Weight")
13
14 # set ticks for each of these persons
15 ax[0,0].set_xticks(ticks=[0,1,2,3]);
16 ax[0,1].set_xticks(ticks=[0,1,2,3]);
17
18 # set labels for each of these persons
19 ax[0,0].set_xticklabels(labels=['p1', 'p2', 'p3', 'p4']);
20 ax[0,1].set_xticklabels(labels=['p1', 'p2', 'p3', 'p4']);

```

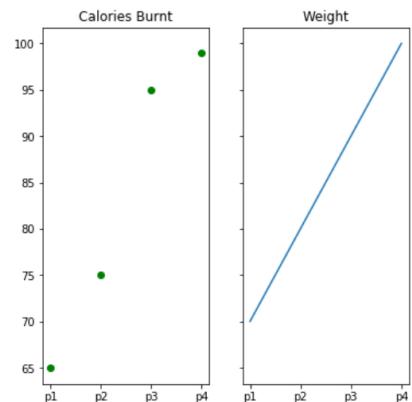


- Notice that in the above figure we have two empty plots, that is because we created 4 subplots (2 rows and 2 columns).
- As a data scientist, there will be times when you need to have a common axis for all your subplots. You can do this by using the `sharex` and `sharey` parameters of `subplot()`.

```

In [62]: 1 # create 2 plots
2 fig, ax = plt.subplots(nrows=1, ncols=2, figsize=(6,6), sharex=True, sharey=True)
3
4 # plot on 0 row and 0 column
5 ax[0].plot(calories_burnt,'go')
6
7 # plot on 0 row and 1 column
8 ax[1].plot(weight)
9
10 # set titles for subplots
11 ax[0].set_title("Calories Burnt")
12 ax[1].set_title("Weight")
13
14 # set ticks for each of these persons
15 ax[0].set_xticks(ticks=[0,1,2,3]);
16 ax[1].set_xticks(ticks=[0,1,2,3]);
17
18 # set labels for each of these persons
19 ax[0].set_xticklabels(labels=['p1', 'p2', 'p3', 'p4']);
20 ax[1].set_xticklabels(labels=['p1', 'p2', 'p3', 'p4']);

```



- Notice in the above plot, now both x and y axes are only labelled once for each of the outer plots. This is because the inner plots "share" both the axes.
- Also, there are only **two plots** since we decreased the number of rows to 1 and columns to 2 in the `subplot()`.
- You can learn more about [subplots here](#).

Load dataset

Let's load a dataset and have a look at first 5 rows.

```

In [63]: 1 # read the dataset
2 data_BM = pd.read_csv('bigmart_data.csv')
3 # drop the null values
4 data_BM = data_BM.dropna(how="any")
5 # view the top results
6 data_BM.head()

```

	Item_Identifier	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location
0	FDA15	9.300	Low Fat	0.016047	Dairy	249.8092	OUT049	1999	Medium	
1	DRC01	5.920	Regular	0.019278	Soft Drinks	48.2692	OUT018	2009	Medium	
2	FDN15	17.500	Low Fat	0.016760	Meat	141.6180	OUT049	1999	Medium	

4	NCD19	8.930	Low Fat	0.000000	Household	53.8614	OUT013	1987	High
5	FDP36	10.395	Regular	0.000000	Baking Goods	51.4008	OUT018	2009	Medium

3. Line Chart

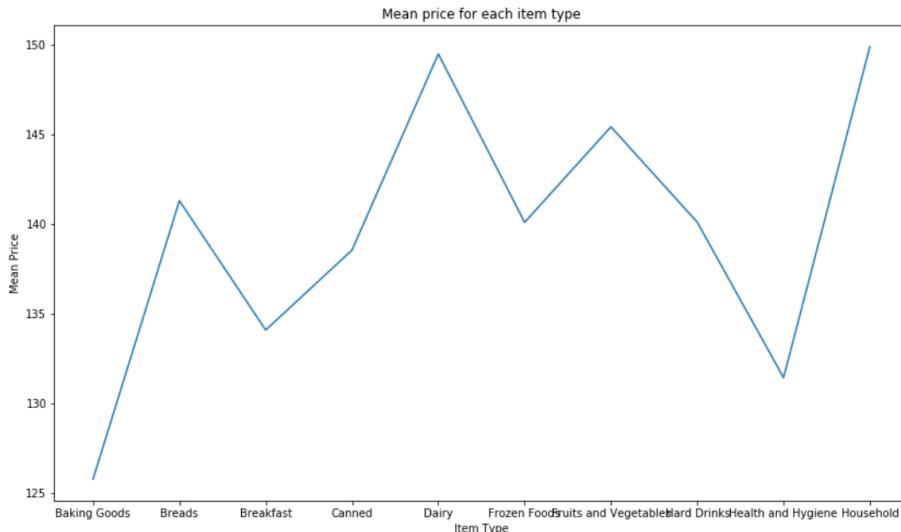
- We will create a line chart to denote the **mean price per item**. Let's have a look at the code.
- With some datasets, you may want to understand changes in one variable as a function of time, or a similarly continuous variable.
- In matplotlib, **line chart** is the default plot when using the `plot()`.

```
In [41]: 1 price_by_item = data_BM.groupby('Item_Type').Item_MRP.mean()[:10]
2 price_by_item
```

```
Out[41]: Item_Type
Baking Goods      125.795653
Breads           141.300639
Breakfast         134.090683
Canned            138.551179
Dairy             149.481471
Frozen Foods     140.095830
Fruits and Vegetables 145.418257
Hard Drinks       140.102908
Health and Hygiene 131.437324
Household          149.884244
Name: Item_MRP, dtype: float64
```

```
In [42]: 1 # mean price based on item type
2 price_by_item = data_BM.groupby('Item_Type').Item_MRP.mean()[:10]
3
4 x = price_by_item.index.tolist()
5 y = price_by_item.values.tolist()
6
7 # set figure size
8 plt.figure(figsize=(14, 8))
9
10 # set title
11 plt.title('Mean price for each item type')
12
13 # set axis labels
14 plt.xlabel('Item Type')
15 plt.ylabel('Mean Price')
16
17 # set xticks
18 plt.xticks(labels=x, ticks=np.arange(len(x)))
19
20 plt.plot(x, y)
```

```
Out[42]: [<matplotlib.lines.Line2D at 0x7f6caec49dd8>]
```



4. Bar Chart

- Suppose we want to have a look at **what is the mean sales for each outlet type?**
- A bar chart is another simple type of visualization that is used for categorical variables.
- You can use `plt.bar()` instead of `plt.plot()` to create a bar chart.

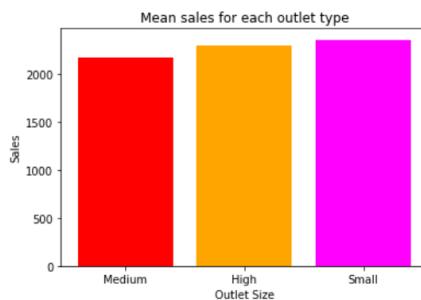
```
In [43]: 1 # sales by outlet size
2 sales_by_outlet_size = data_BM.groupby('Outlet_Size').Item_Outlet_Sales.mean()
3
4 # sort by sales
5 sales_by_outlet_size.sort_values(inplace=True)
6
7 x = sales_by_outlet_size.index.tolist()
8 y = sales_by_outlet_size.values.tolist()
9
10 # set axis labels
11 plt.xlabel('Outlet Size')
12 plt.ylabel('Sales')
13
14 # set title
15 plt.title('Mean sales for each outlet type')
16
17 # set xticks
```

```

18 plt.xticks(labels=x, ticks=np.arange(len(x)))
19
20 plt.bar(x, y, color=['red', 'orange', 'magenta'])

```

Out[43]: <BarContainer object of 3 artists>



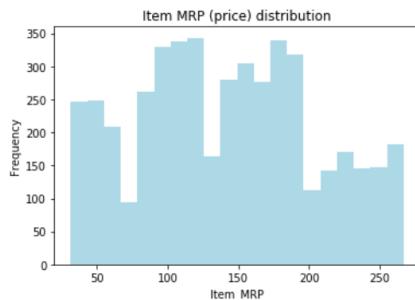
5. Histogram

- Distribution of item price
- Histograms are a very common type of plots when we are looking at data like height and weight, stock prices, waiting time for a customer, etc which are continuous in nature.
- Histogram's data is plotted within a range against its frequency.
- Histograms are very commonly occurring graphs in probability and statistics and form the basis for various distributions like the normal -distribution, t-distribution, etc.
- You can use `plt.hist()` to draw a histogram. It provides many parameters to adjust the plot, you can [explore more here](#).

```

In [44]: 1 # title
2 plt.title('Item MRP (price) distribution')
3
4 # xlabel
5 plt.xlabel('Item_MRP')
6
7 # ylabel
8 plt.ylabel('Frequency')
9
10 # plot histogram
11 plt.hist(data_BM['Item_MRP'], bins=20, color='lightblue');

```



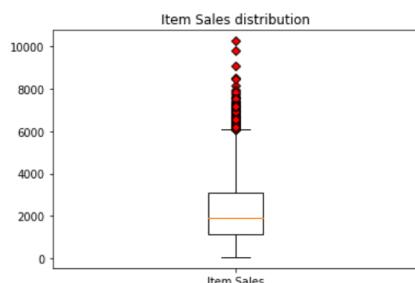
6. Box Plots

- Distribution of sales
- Box plot shows the three quartile values of the distribution along with extreme values.
- The "whiskers" extend to points that lie within 1.5 IQRs of the lower and upper quartile, and then observations that fall outside this range are displayed independently.
- This means that each value in the boxplot corresponds to an actual observation in the data.
- Let's try to visualize the distribution of `Item_Outlet_Sales` of items.

```

In [46]: 1 data = data_BM[['Item_Outlet_Sales']]
2
3 # create outlier point shape
4 red_diamond = dict(markerfacecolor='r', marker='D')
5
6 # set title
7 plt.title('Item Sales distribution')
8
9 # make the boxplot
10 plt.boxplot(data.values, labels=['Item Sales'], flierprops=red_diamond);

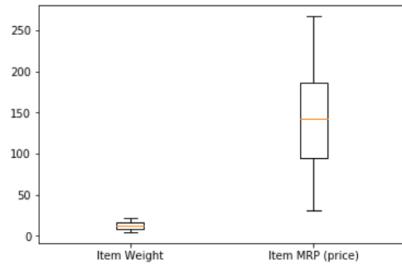
```



- You can also create multiple boxplots for different columns of your dataset.
- In order to plot multiple boxplots, you can use the same `subplots()` that we saw earlier.
- Let's see `Item_Weight` `Item_MRP` distribution together

```
- zero_out Item_Weight, Item_MRP distribution regular
```

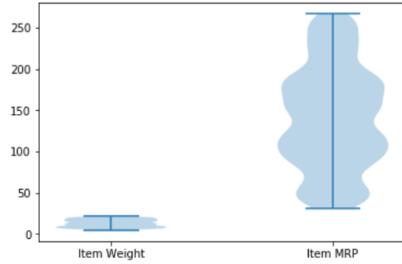
```
In [47]: 1 data = data_BM[['Item_Weight', 'Item_MRP']]  
2  
3 # create outlier point shape  
4 red_diamond = dict(markerfacecolor='r', marker='D')  
5  
6 # generate subplots  
7 fig, ax = plt.subplots()  
8  
9 # make the boxplot  
10 plt.boxplot(data.values, labels=['Item Weight', 'Item MRP (price)'], flierprops=red_diamond);
```



7. Violin Plots

- Density distribution of item weights and item price

```
In [48]: 1 data = data_BM[['Item_Weight', 'Item_MRP']]  
2  
3 # generate subplots  
4 fig, ax = plt.subplots()  
5  
6 # add labels to x axis  
7 plt.xticks(ticks=[1,2], labels=['Item Weight', 'Item MRP'])  
8  
9 # make the violinplot  
10 plt.violinplot(data.values);
```



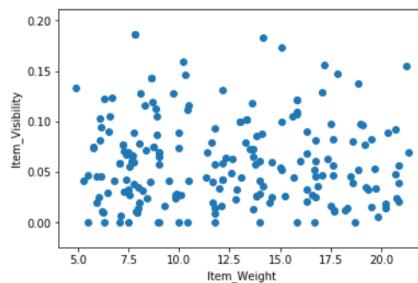
8. Scatter Plots

- Relative distribution of item weight and it's visibility
- It depicts the distribution of two variables using a cloud of points, where each point represents an observation in the dataset.
- This depiction allows the eye to infer a substantial amount of information about whether there is any meaningful relationship between them.

NOTE : Here, we are going to use only a subset of the data for the plots.

```
In [49]: 1 # set label of axes  
2 plt.xlabel('Item_Weight')  
3 plt.ylabel('Item_Visibility')  
4  
5 # plot  
6 plt.scatter(data_BM["Item_Weight"][:200], data_BM["Item_Visibility"][:200])
```

```
Out[49]: <matplotlib.collections.PathCollection at 0x7f6cae3edd30>
```



9. Bubble Plots

- Relative distribution of sales, item price and item visibility
- Let's make a scatter plot of Item_Outlet_Sales and Item_MRP and make the size of bubbles by the column Item_Visibility.
- Bubble plots let you understand the interdependent relations among 3 variables.

Note that we are only using a subset of data for the plots.

```
In [50]: 1 # set label of axes
```

```
2 plt.xlabel('Item_MRP')
3 plt.ylabel('Item_Outlet_Sales')
4
5 # set title
6 plt.title('Item Outlet Sales vs Item MRP (price)')
7
8 # plot
9 plt.scatter(data_BM["Item_MRP"][:100], data_BM["Item_Outlet_Sales"][:100], s=data_BM["Item_Visibility"][:100]*1000, c='red')
```

Out[50]: <matplotlib.collections.PathCollection at 0x7f6cae3cbf28>

