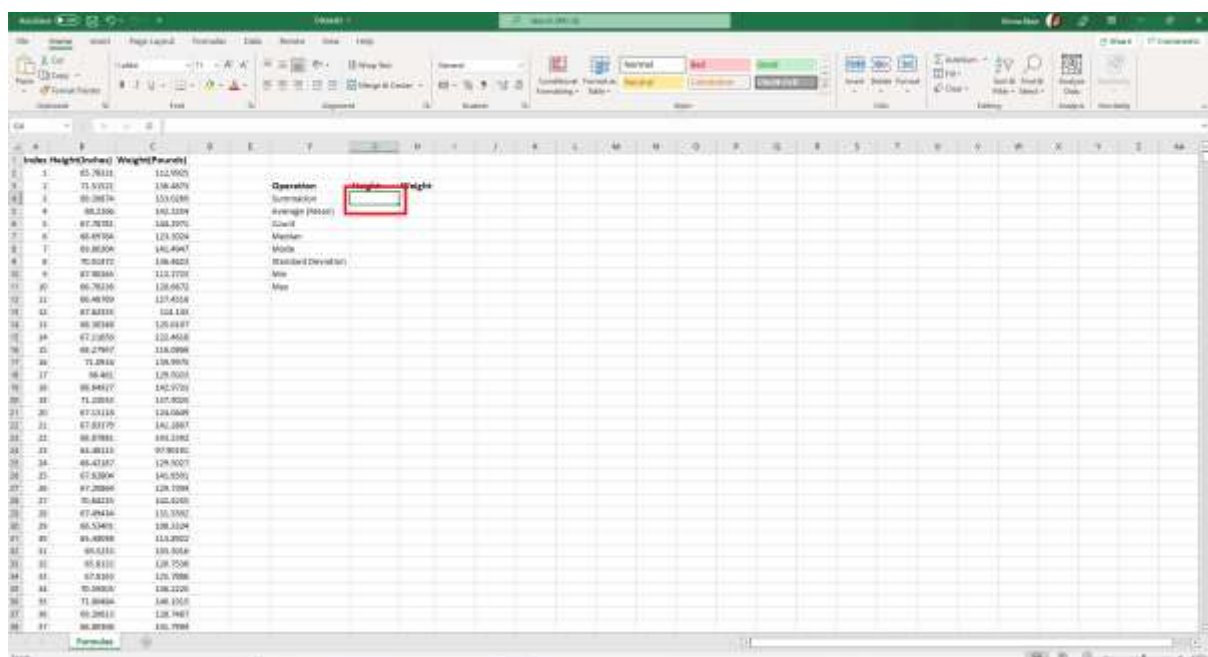## Week-5
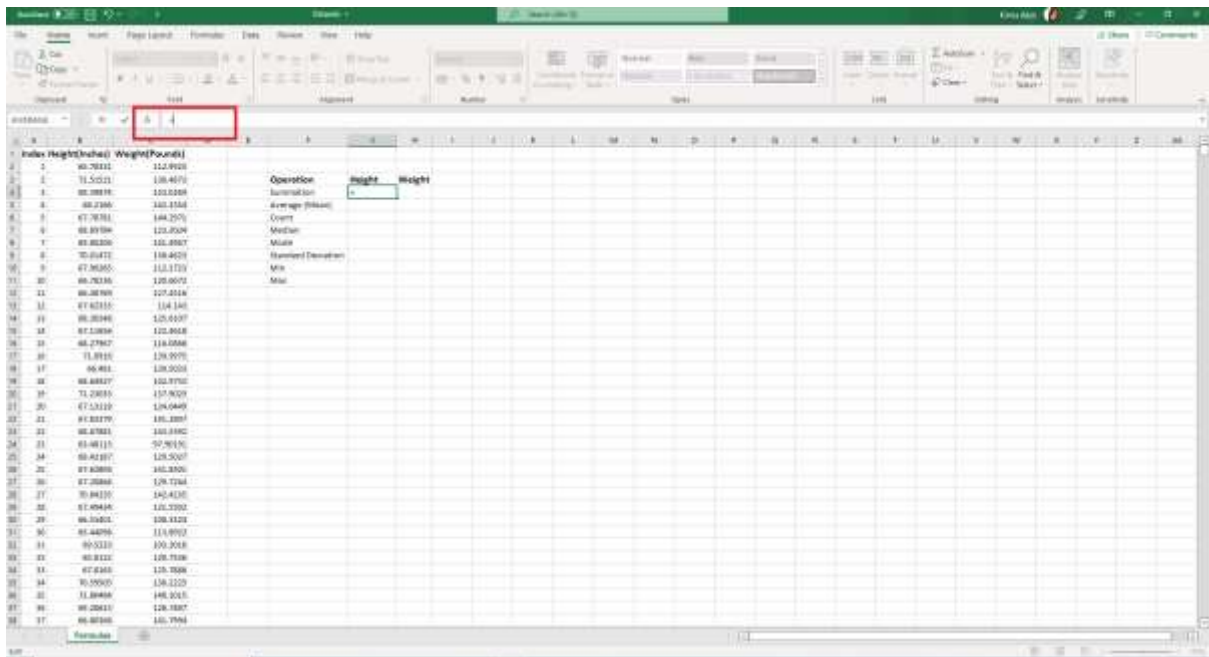## Data Processing and Visualization using Off-the Shelf Tools
## <u>Lab Guideline</u>

**In today's Lab you will use Microsoft Excel as for performing Data preprocessing operations concerning "Data Cleaning", "Basic Statistical Analysis" as well as "Visualization". The dataset(s) you will be using are available in a excel (.csv) file.**
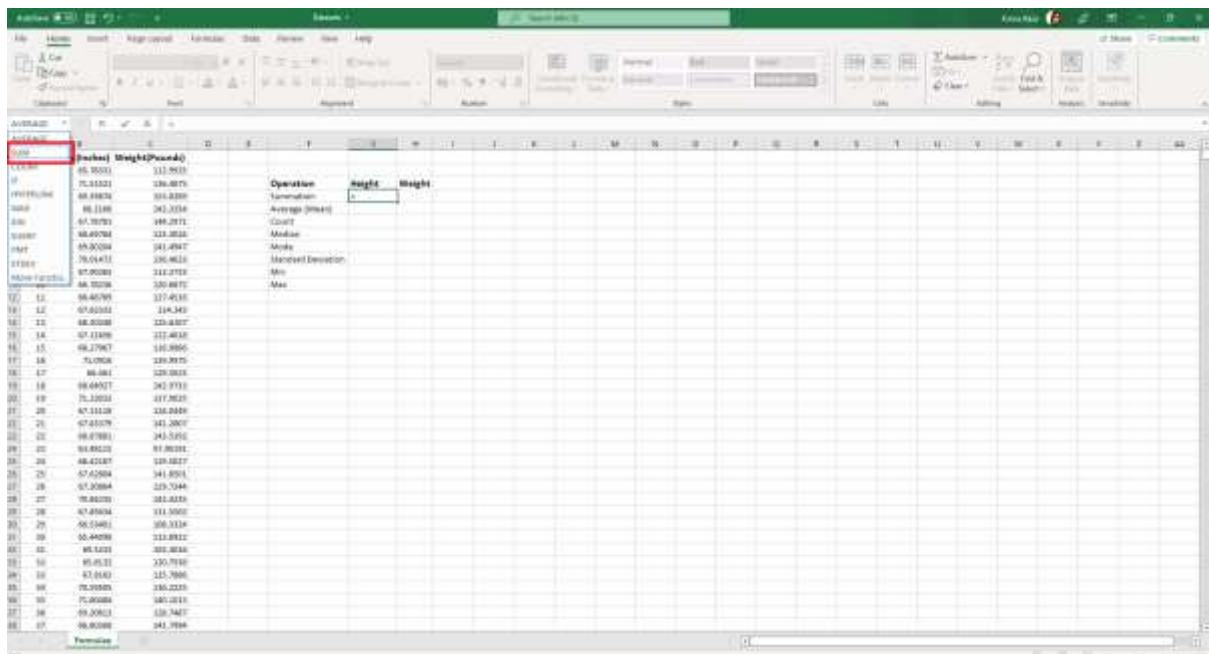
**Applying Useful Formula and Basic Statistical Analysis:**

1. Download the Datasets Excel (.csv) file from Moodle. Open the "Formulae" spread sheet. You will find Height(inches) and weight(pounds) data.
2. At this step you will use readily available formulas available in MS Excel for performing some basic statistical operations.
3. To do so, following are the steps that you would continue:
   a) Assume that, you want to calculate the summation of "Height" data available to you. Before you start the operation, you must select a cell where your desired "Summation" output will be held/stored. After clicking/selecting the cell you type equal "=" in the formula box
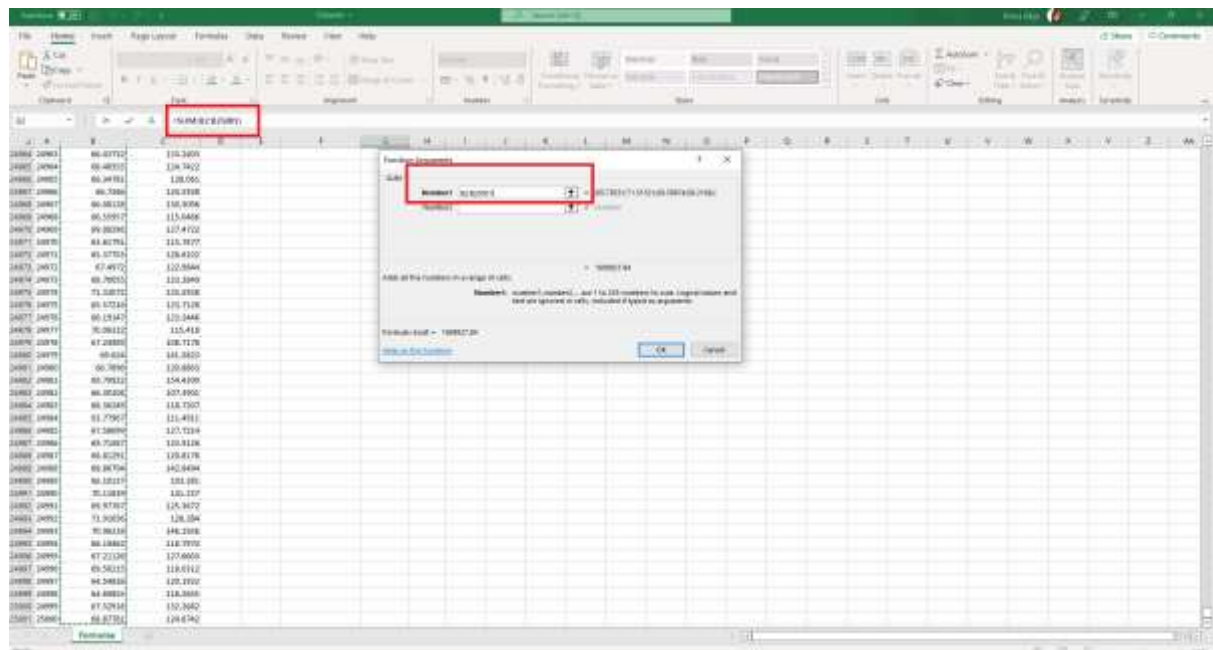   b) Then your formula operation mode is activated.

c) Keeping everything consistent till the previous step, click on the top left corner. You will find a formula list available. From there, you will now select "summation" operation.
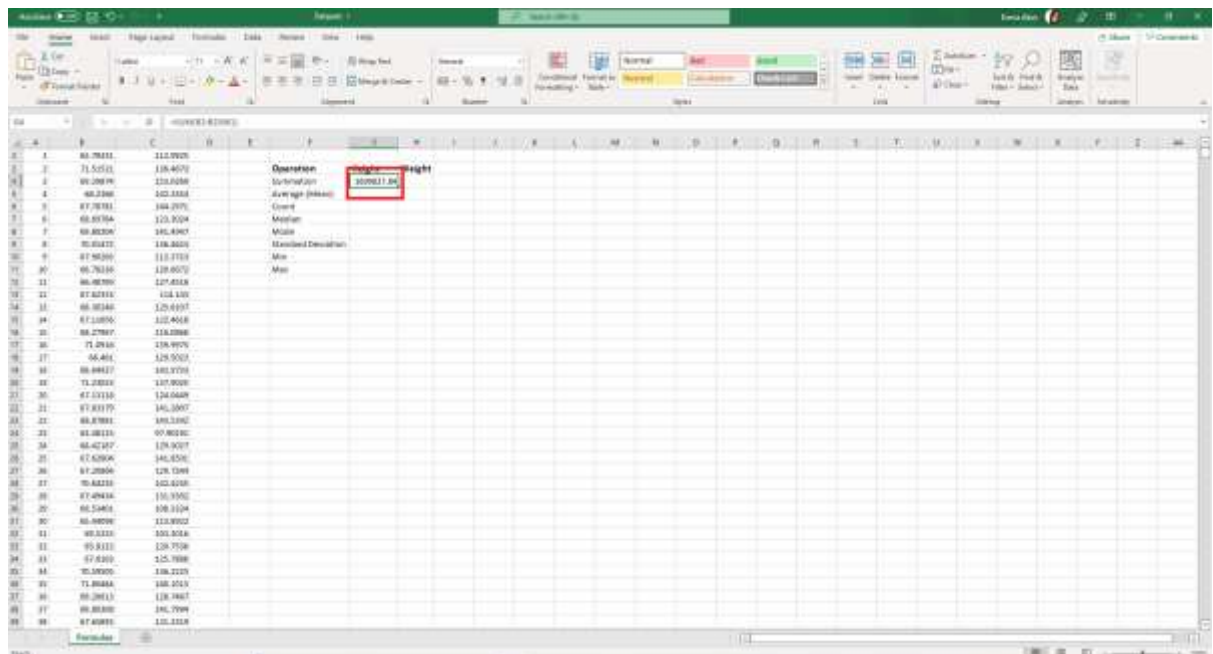


d) Once you select a function from the list a new window will pop up. You have to write "From" and "To" cell number(s) in here. We are taking cell "B2 to cell B25001" As we

are up for summing up the total height from the data present here. We write it like,
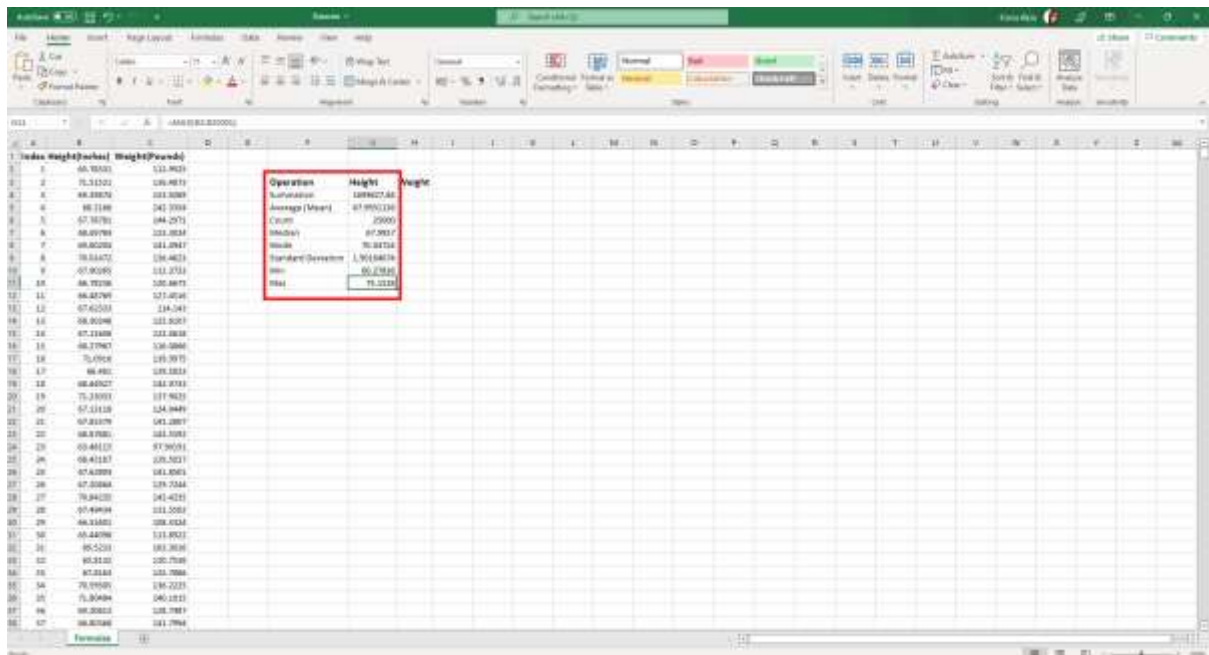


B2:B25001 in the text box.

e) After giving the input and clicking "OK", you should now get the result of the summation of the "Height" in the cell you selected.



Following the same steps, you now should be able to perform other statistical operations from the dataset.
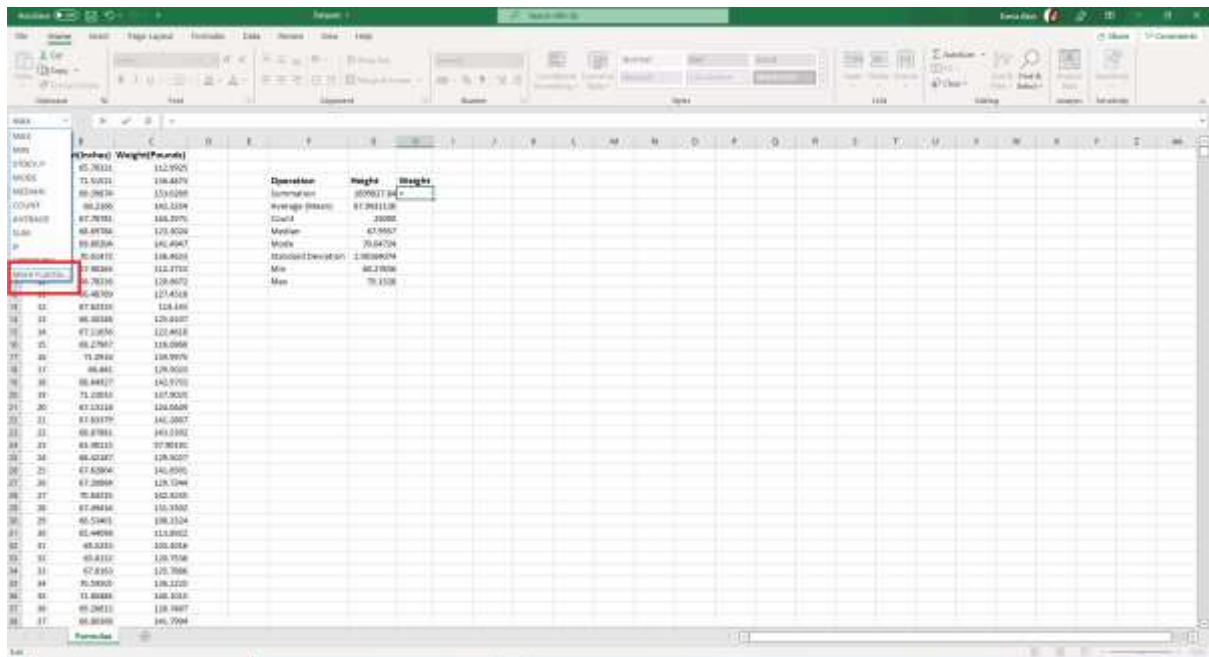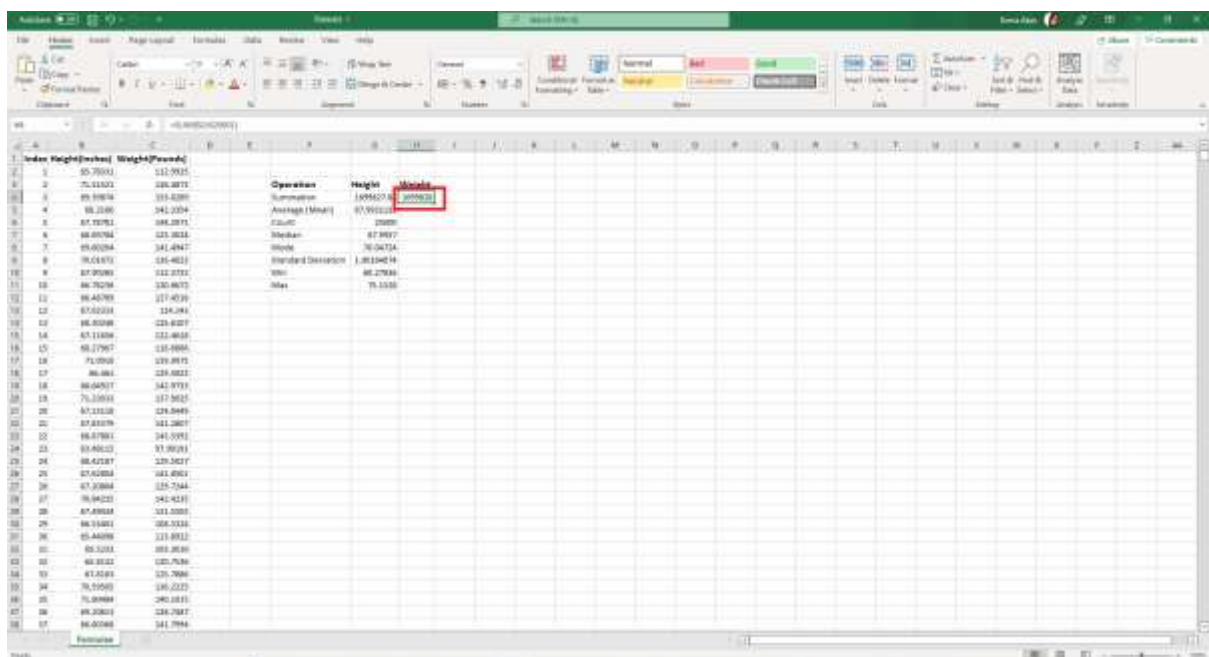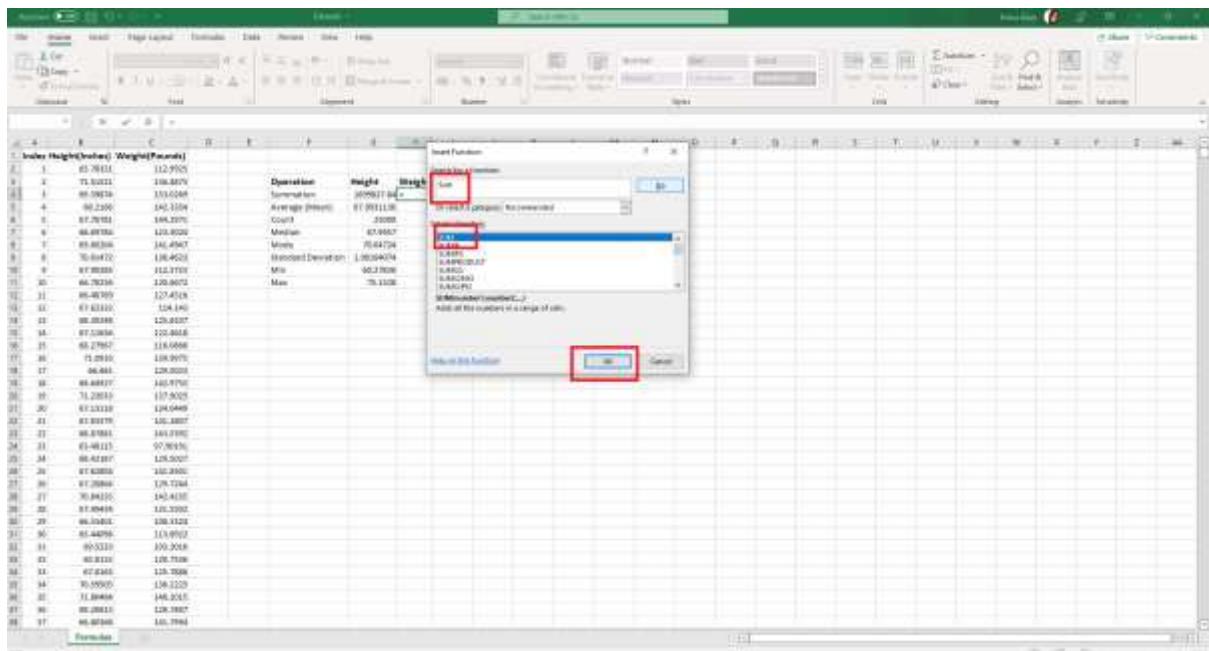
Kimia Aksir

## Useful Tips:

Instead of scrolling down and searching for the operation type. You can manually search for the operation you want to perform by using following steps:

Assume, you want to calculate the summation of the "Weight". So, you select cell. Go to choose function and go for "More Options".



Kimia Aksir

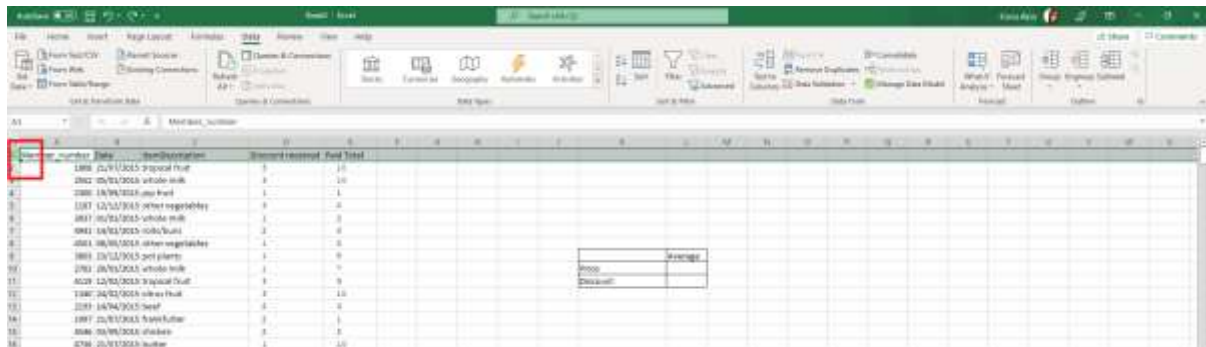You can manually type/input the operation "Sum", and it will give you the result in the selected cell.





## Task1:

Similarly, calculate the Average(Mean), Median, Mode and Standard Deviation of the "Weight" data given in the same dataset.

Kimia Aksir

## Data Cleaning:

There is another dataset given to you, called "Data Cleaning". In this Dataset you will find data of a grocery store.
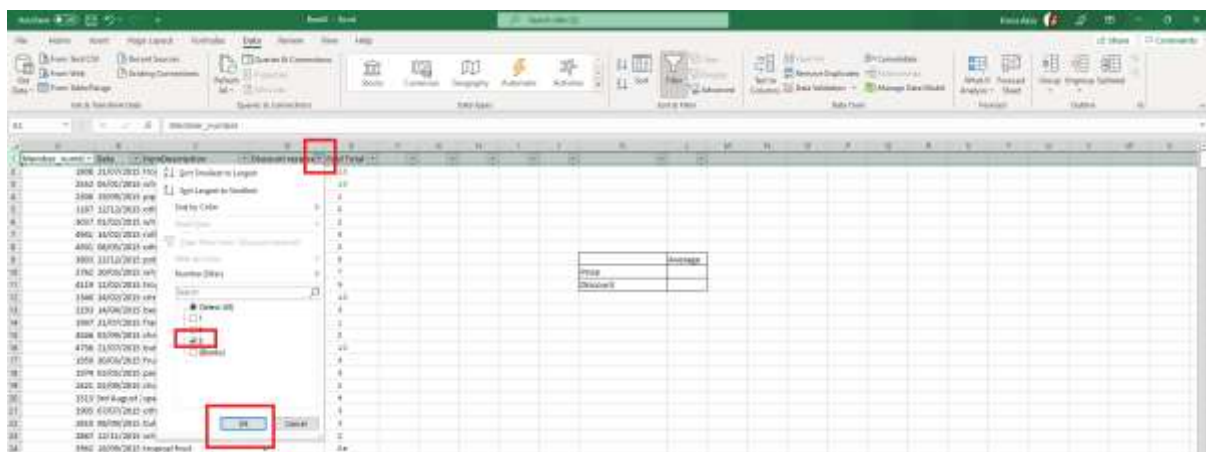
1. First observe if you can see any inconsistent formatting or, missing value problem in the data there.
2. In this step you will see how you can add filters and so, it should be much easier for you to see the inconsistencies in the dataset.
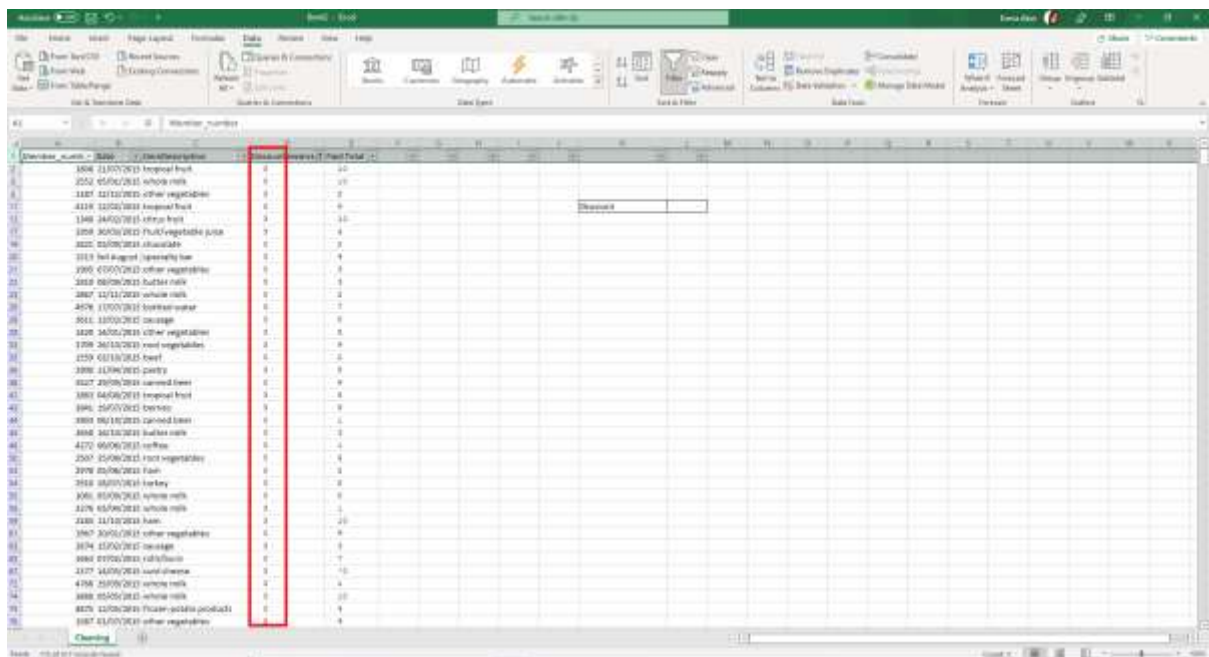   a) Select the row that contains the Title of your attributes:



   b) Keeping this row selected, go to option "Data" and select "Filter". This is how you should be able to add "Filter" on every attribute present in this dataset.



   c) Let's say, you want to see only the customers got discount of GBP 3. You can see them by applying filter list.



Kimia Aksir

Now, if you use filter on "Paid total" column, you will see all the unique values in the filter list. In the list you will see there are some values which are wrongly recorded (Some special characters and/or, character inputs have come along).



After they are selected. You can now easily see only these values and do the cleaning as they are already spotted from the filter.



Kimia Aksir

You can now correct them and make the data consistent, or, even discard the tuples if there's no hint left to correct them. This attribute contains total paid amount and only removing the characters from the values should work appropriately.

## Task:2
Use filter on "Discount" attribute select the records where "Discount" value is missing and discard those tuples.
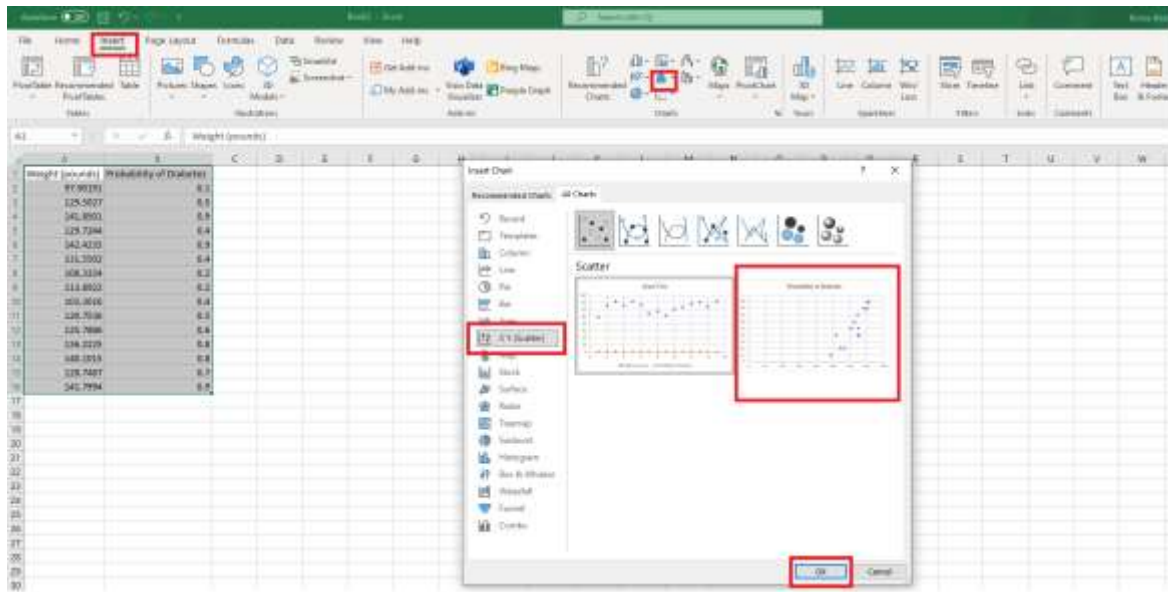
## Task:3

Have a look at the date column make the date format consistent, if you observe any inconsistency exists.
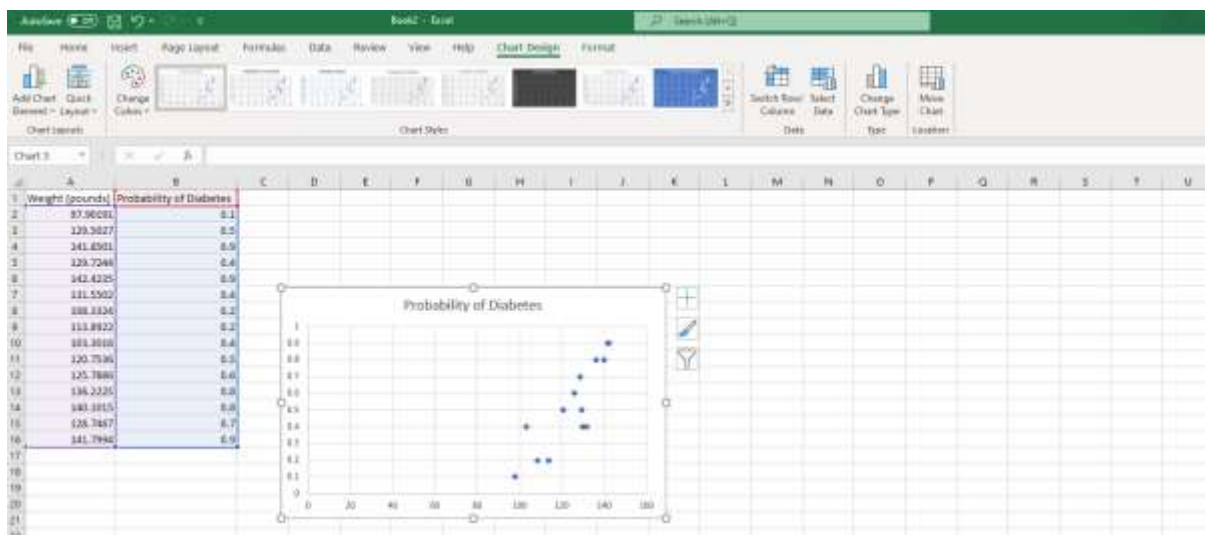
## Data Visualization:

As we know already, we can represent useful and meaningful insights from data through Data Visualization. We will now, work on Data representation. For that, you will use "Visualization" spread sheet from the Datasets excel (.csv) file. The Dataset contains data of weight and probability of diabetes.

You first select the data in the Dataset, Go to option "Insert" from above, go to chart and graph option. Pick X. Y. scatter and select "OK".



You now get a plotting of probability of diabetes against the weight data.



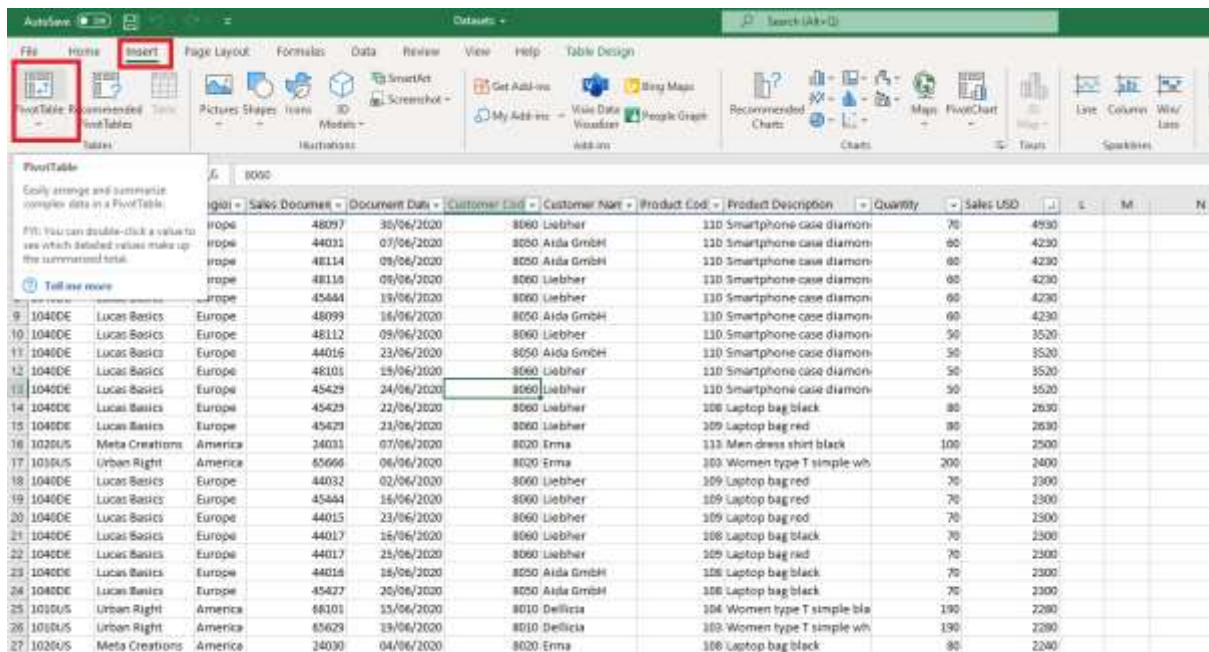You con use the same data and find a different type of visualization output of it.

Kimia Aksir

## Task:4

Now, represent the sales data, available in "Sales" spreadsheet using a Bar Chart (Column).
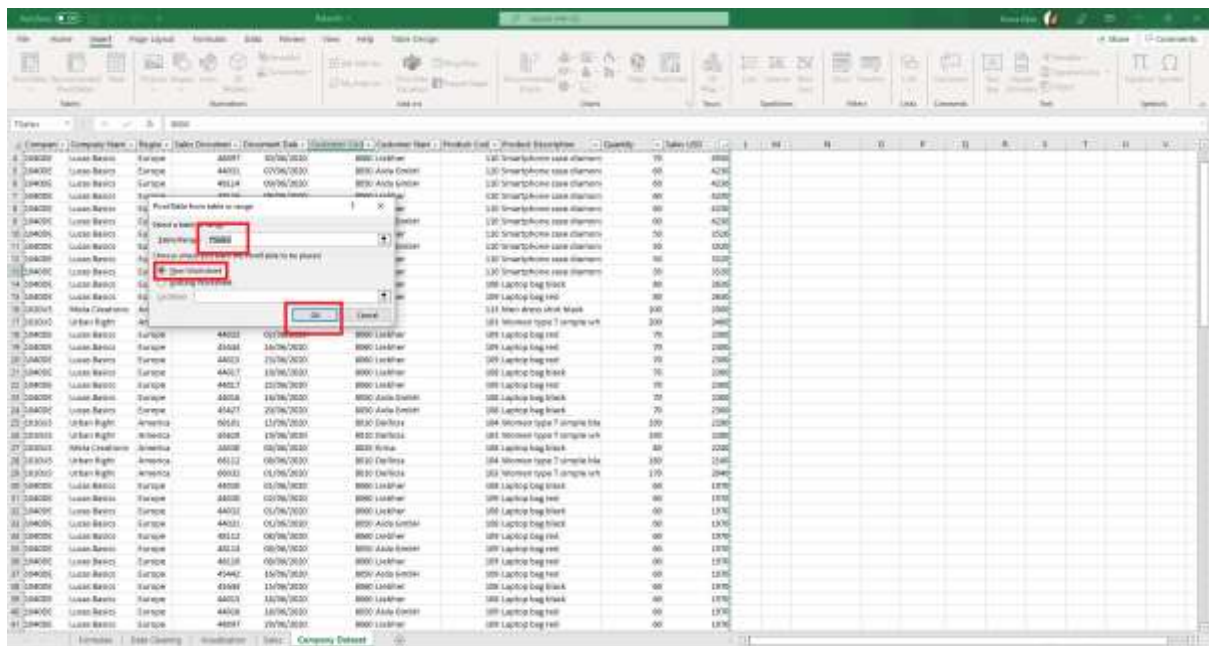
## Pivot Table:

We will now work on Pivot table. For that, you will use "Company Dataset" spreadsheet from the same dataset file excel file.
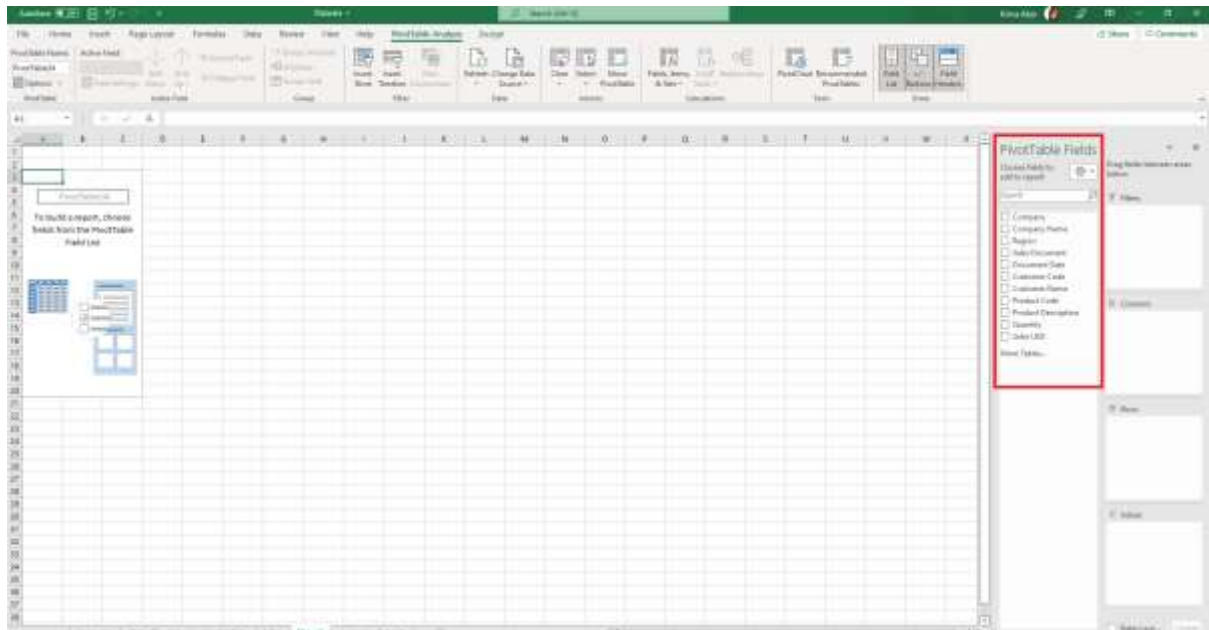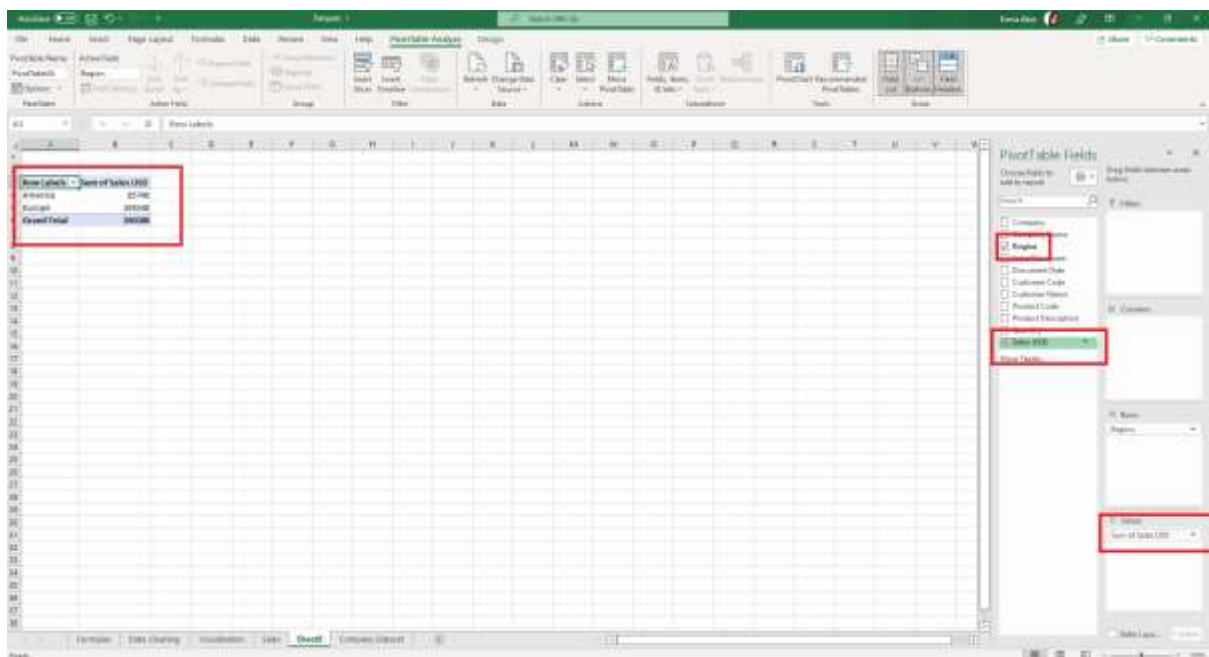
Open the spread sheet. Go to "Insert", Pivot Table.



You will select ok for Tsale table already created for you.



Kimia Aksir

Once the pivot gets opened in a new worksheet. You will see by manually selecting the parameters from the table you can generate any report you want.



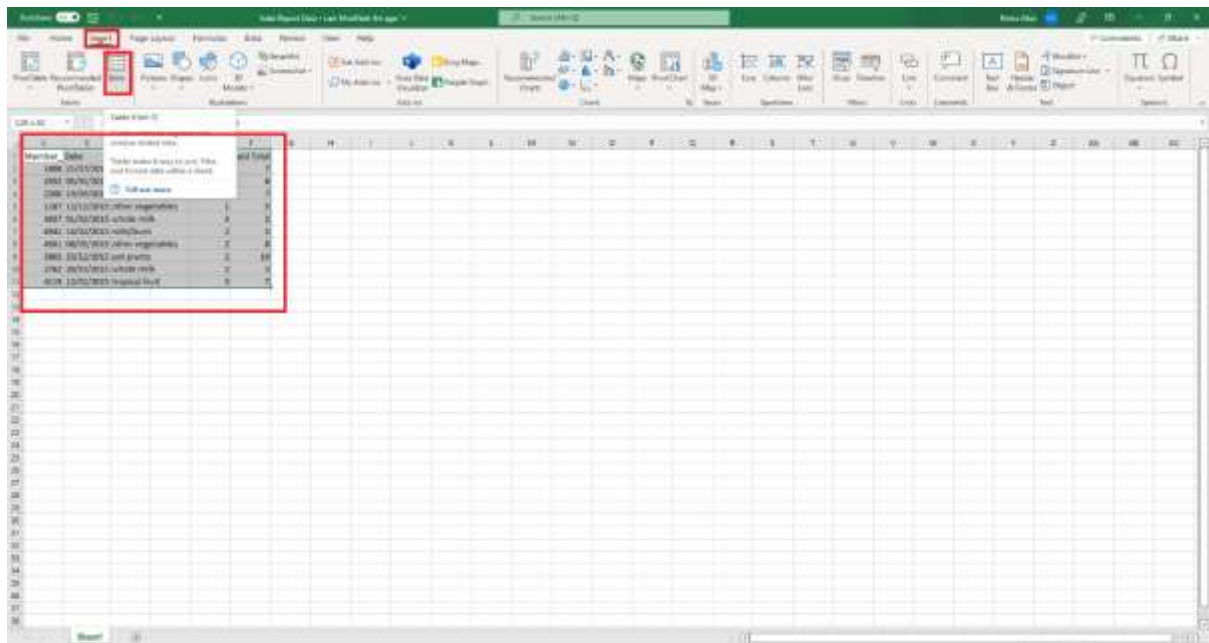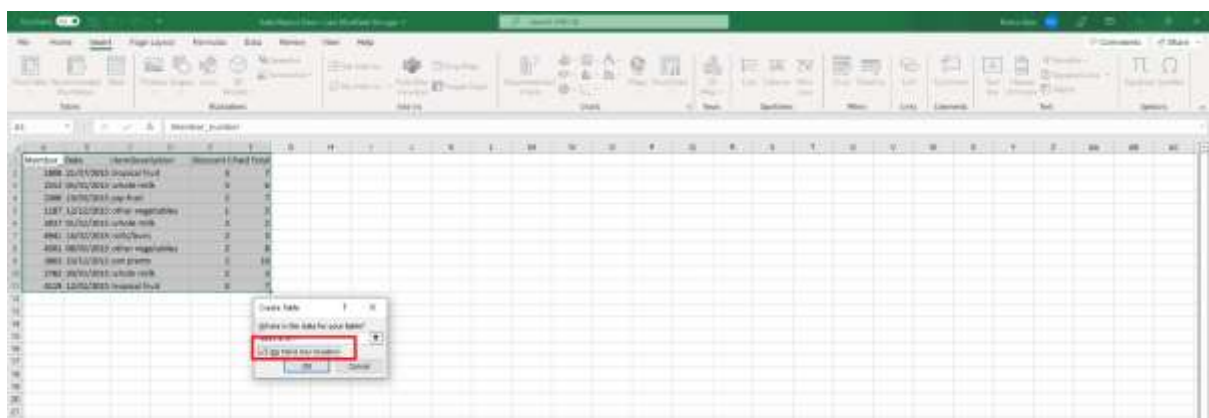Let's find the total sales (USD) in different regions.



## Task:5

Find the Region wise Average of Sales (USD)

Kimia Aksir

You can now create and manage a Pivot Table on your own. Use "Sales Report Data" dataset and try to follow the steps given below.
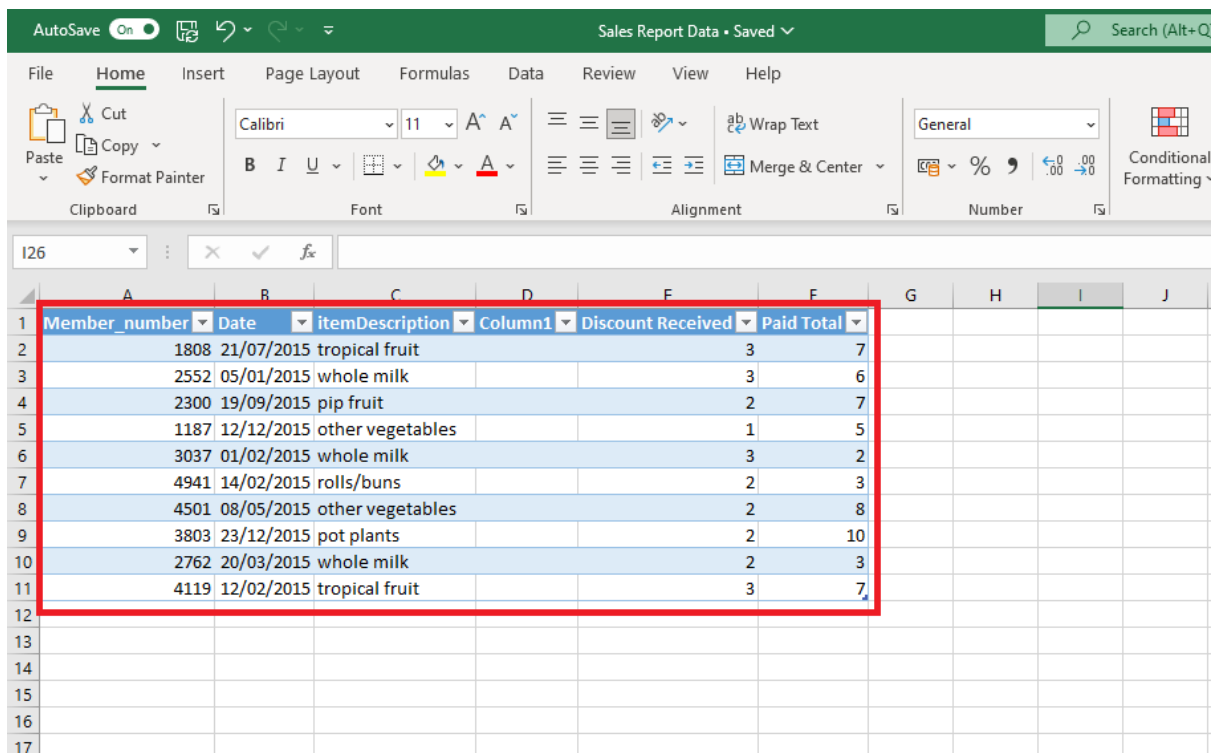
Selecting the Dataset Go to "Insert" menu tab from the top and select "Table":



Remember to tick "My table has header" manually if it's not automatically selected. Because without Header Pivot Table won't work!
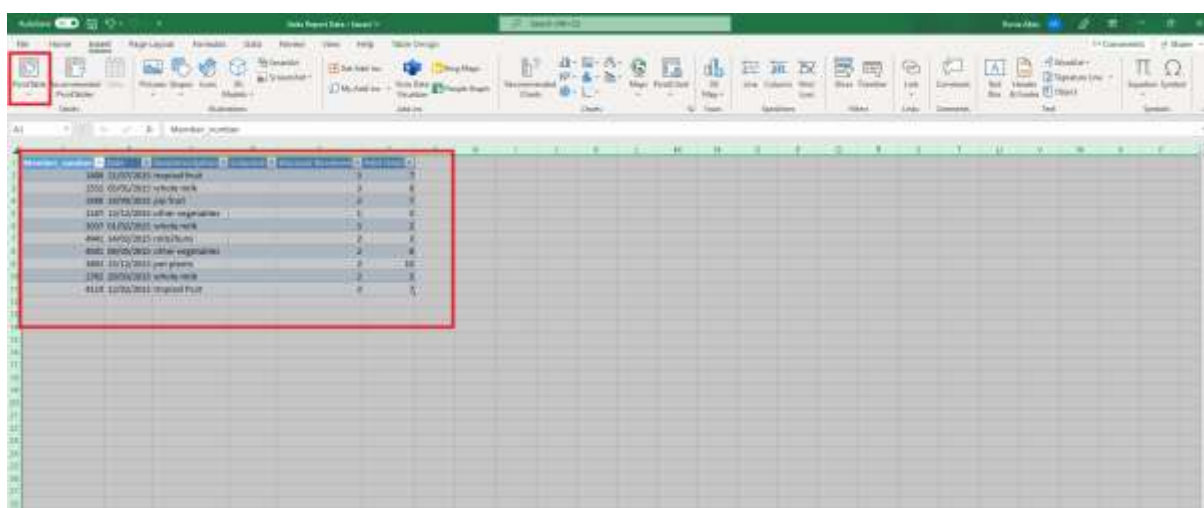


Kimia Aksir

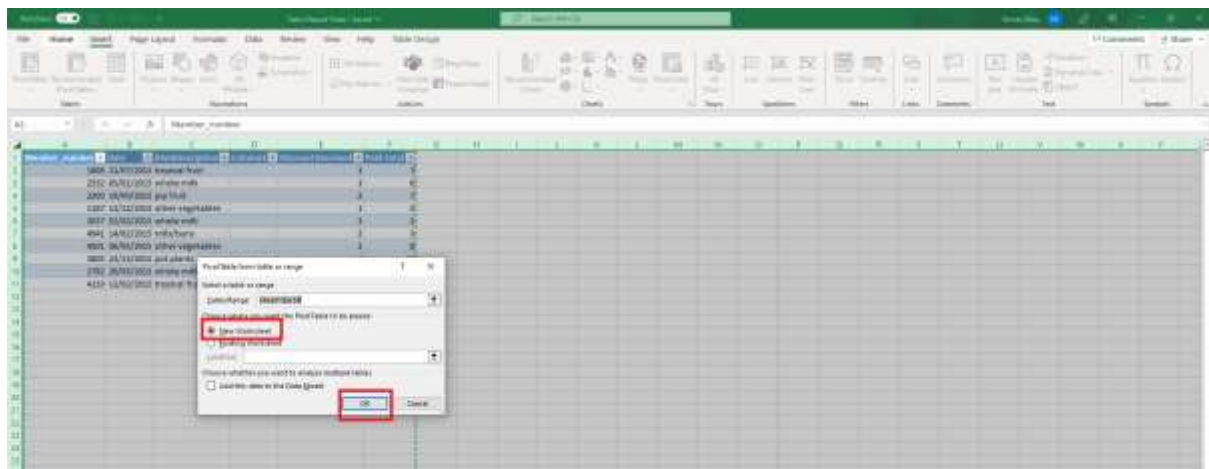Finally, You get your table as shown below with all the headers present in it:
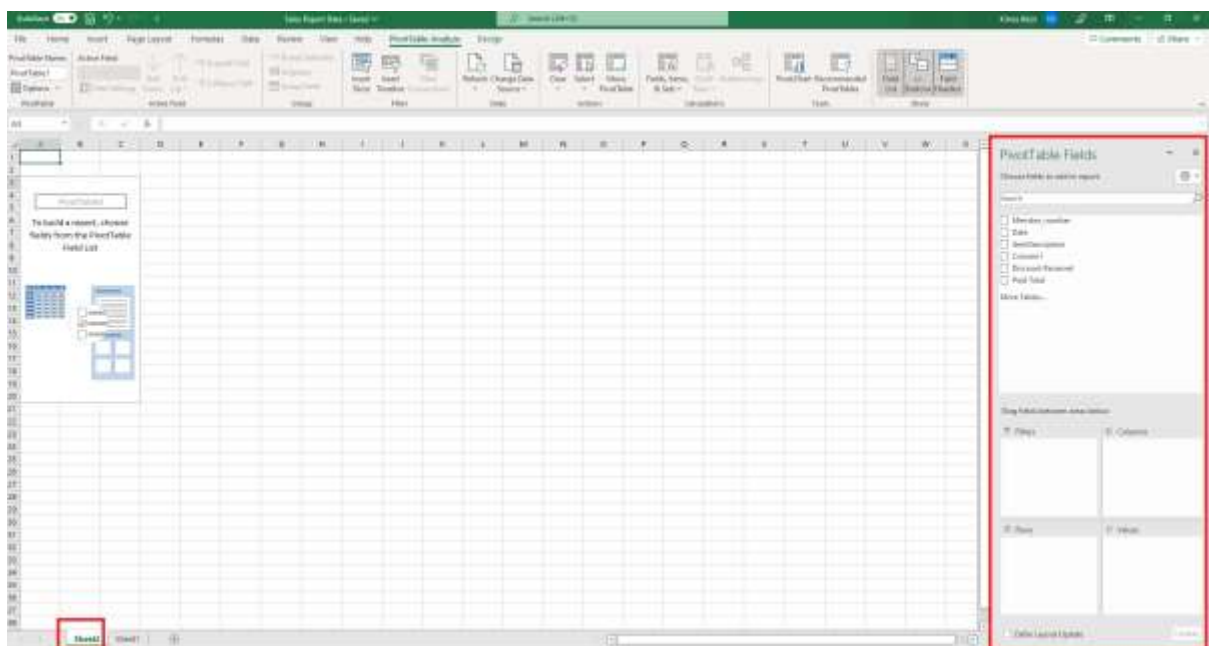


Keeping the table selected, go to the tab "Insert" and then "Pivot table". We may also want to choose "Recommended Pivot", however, choosing the former option will always allow us to customize the features on our own.
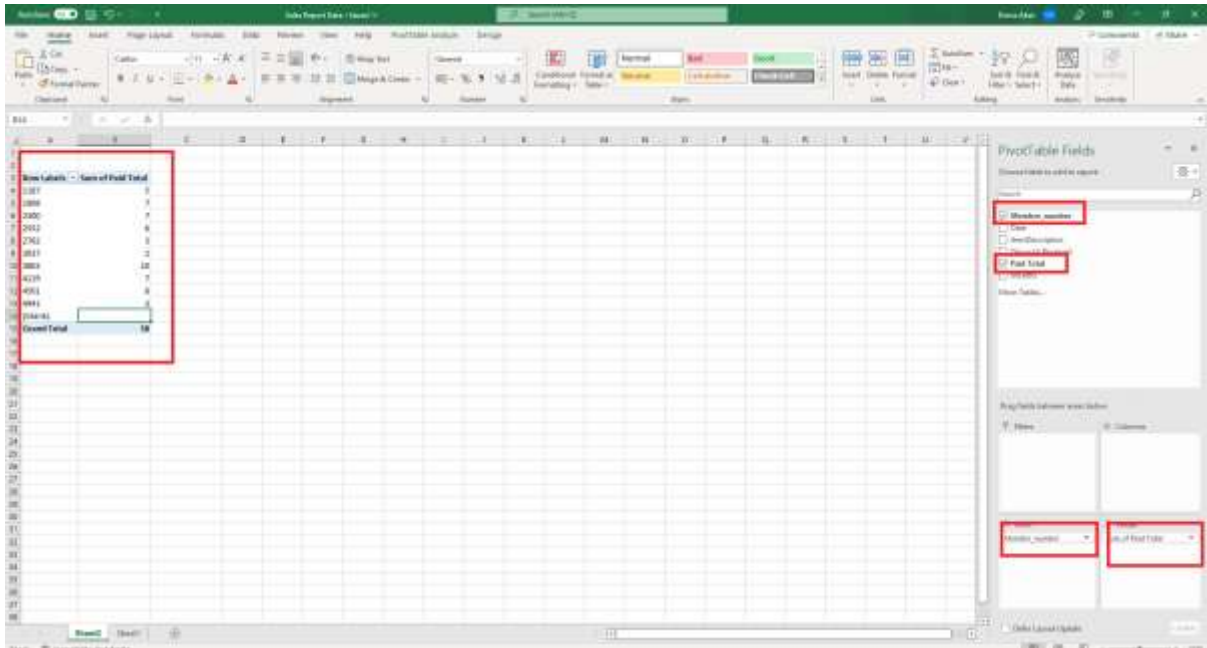


Kimia Aksir

Choosing "New Worksheet" for presenting the Pivot table is recommended here. However, you can also keep the Pivot on the same worksheet along with the table populated with data.



As you can see, in a separate worksheet the Pivot table features appeared. Now, you can generate any report or, check the pattern of the data in the dataset using this pivot table.



Let's generate a report on each member and they paid total.

Kimia Aksir

You can also try to add a few for rows in the table and check how it reflects in the Pivot table.

Kimia Aksir