After importing different libraries and reading the dataset, we tried to see first few rows of the data using method head().



Using describe(), we get to see different statistical factors of columns.

Using info(), we get to see count of non-null values in each column; also, we can see the datatype for each column.



```
[13]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   PassengerId  891 non-null    int64
 1   Survived     891 non-null    int64
 2   Pclass       891 non-null    int64
 3   Name         891 non-null    object
 4   Sex          891 non-null    object
 5   Age          714 non-null    float64
 6   SibSp        891 non-null    int64
 7   Parch        891 non-null    int64
 8   Ticket       891 non-null    object
 9   Fare         891 non-null    float64
 10  Cabin        204 non-null    object
 11  Embarked     889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

Using value_counts() in a column here gave us the counts of each category type in the column (like 0 and 1 are the categories here).

Using shape, we can check the number of rows and column of a dataframe.

Using column, we can see the column names.

Then we tried to drop 2 columns (Passenger Id and Cabin) as they were not giving any intelligence to our further analysis.

Below screenshot showing mean() for Fare and Age using methods. The values were rounded off using round method.

```
[45]: print("Average Fare:", round(data["Fare"].mean(),3))
      Average Fare: 32.204

[49]: print("Average Age:", round(data["Age"].mean(),2))
      Average Age: 29.7
```

This is a simple countplot which shows the plot for count of survived passengers and the count of the passengers who were unable to survive.

This is a pie-plot where it shows the distribution of passengers as per the class (1, 2 &3)



Again a countplot which also shows the survival with classes as well.

A histplot which shows the distribution of age. It shows that most of the people aboard were of the age range of 20 and 40.



The scatterplot which shows that most of the people of all ages bought the tickets of price below 100 Dollars.

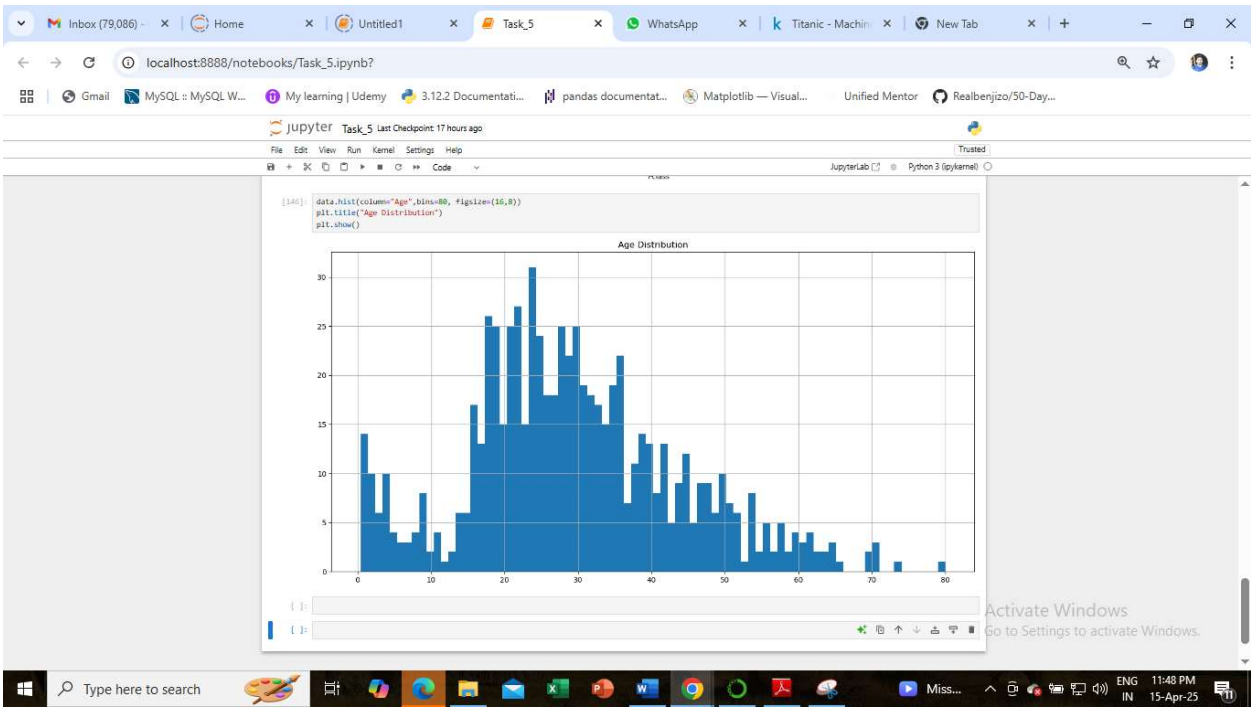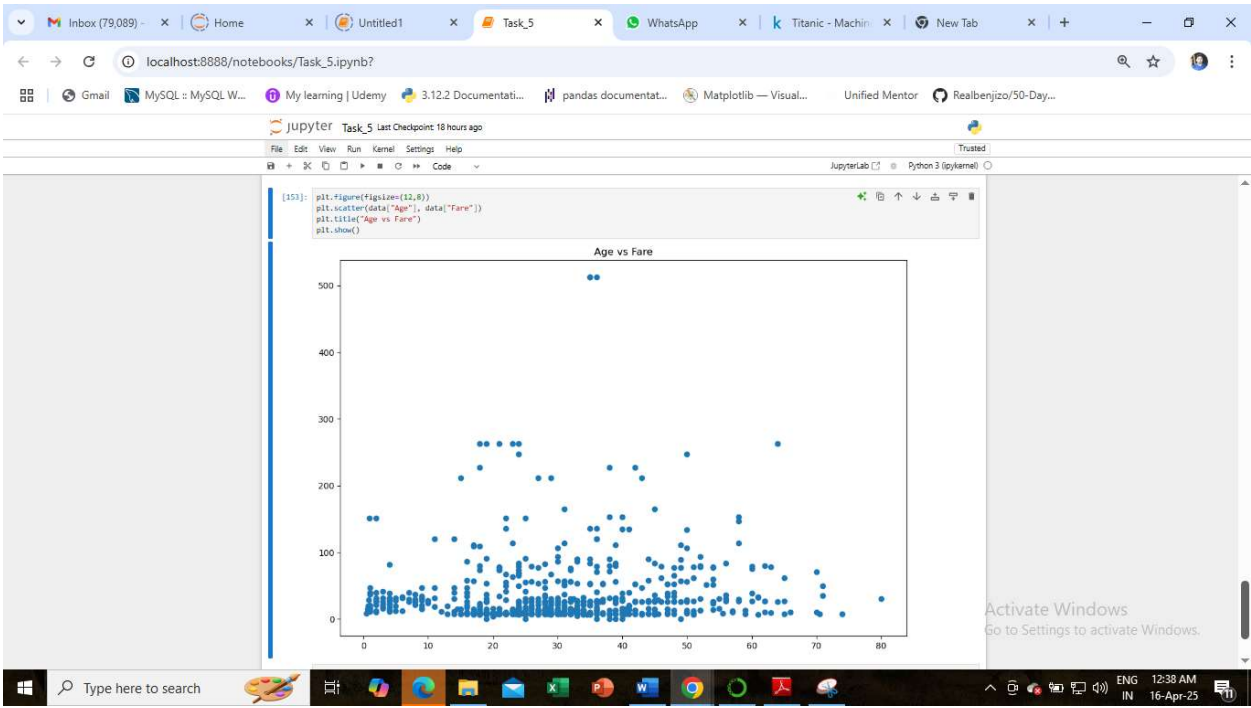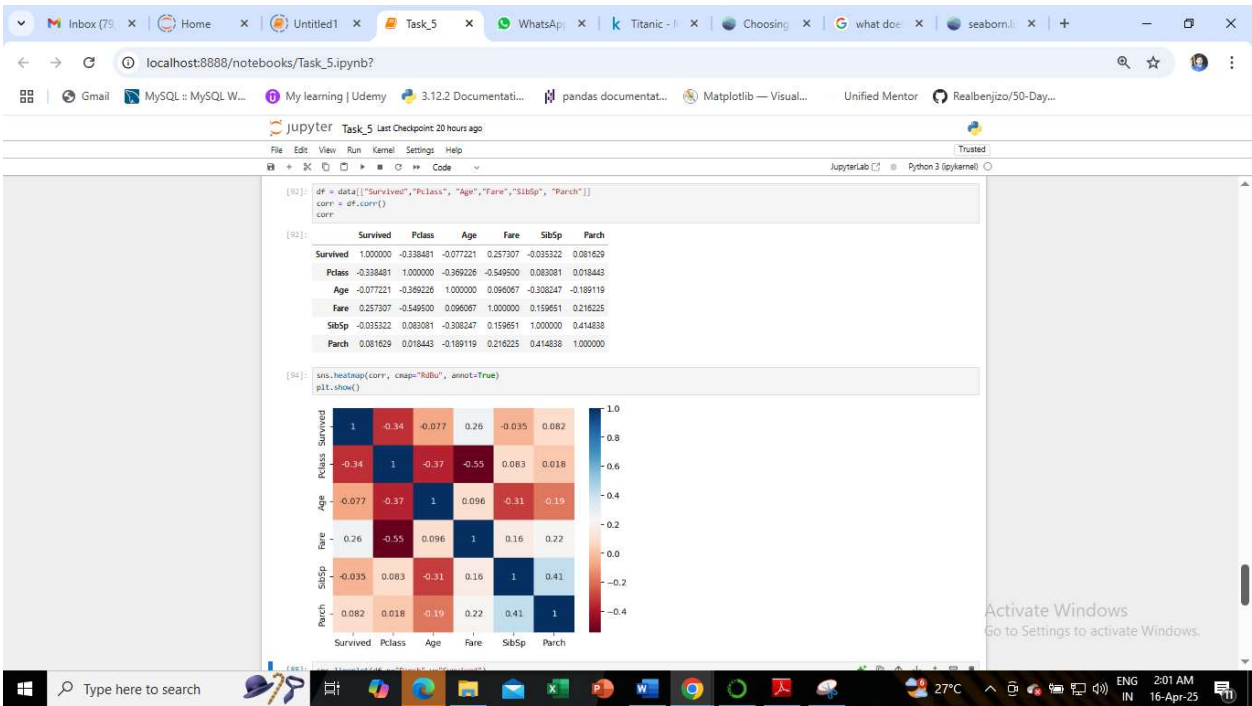The heatmap below shows that some of the column values have positive correlation with few other column values while others have a negative correlation. Like, for example, Pclass have a negative correlation with the Survived.



The pairplot() gave us plots of one column values against all other column values. Like for example, we can see here that with the increase of fare the number diminished for people irrespective of the fact whether they survived or not.

The boxplot again showing us that most of the people were in the age range of 20 and 40. The median lies in between that age range and there are some outliers in the range of 65-80.



Summary of findings-

➢ Most of the people are in the age range of 24-40.
➢ There were quite enough data columns where numerical analysis could be done.
➢ Passenger Id and Cabin columns were removed as they were not helping any inferences and plots.
➢ There were more people who could not survive.
➢ Most of the people made through journey in class 3.
➢ Comparatively, higher number of passengers died from class 3.
➢ Most of the passengers paid for the fare below 100 Dollars.