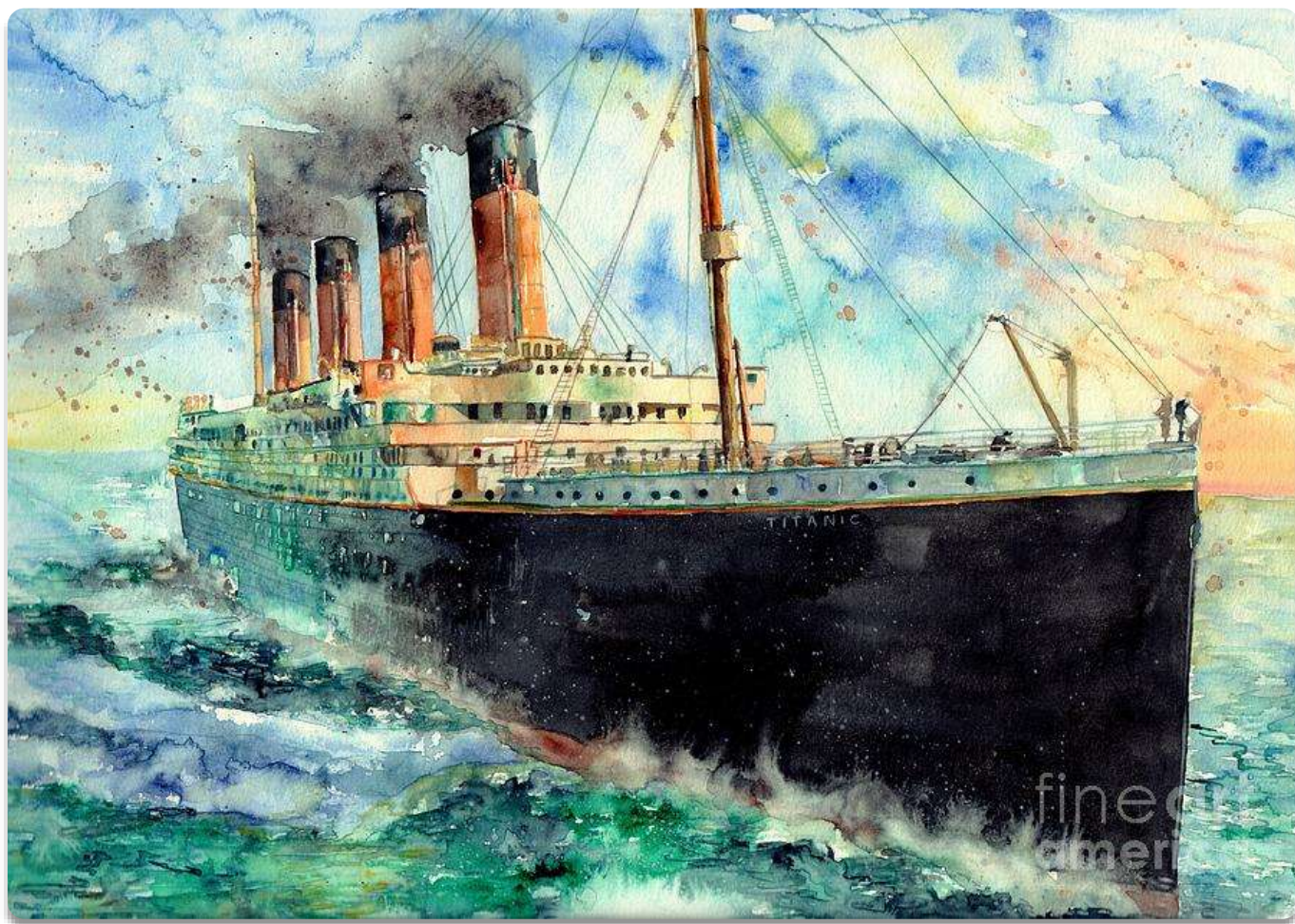


# DS 540 (STATISTICAL PROGRAMMING) FINAL PROJECT

PREDICTING SURVIVAL OF TITANIC PASSENGERS  
(EDA AND MACHINE LEARNING)



# TITANIC SHIP





# Details of Data set



Comprises of 891 rows and 8 columns



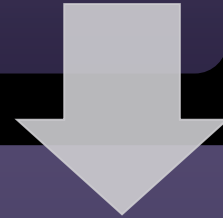
Consists of Categorical and Numerical Data



Consists of 2 missing values

# AIM OF PROJECT

To find the basic statistics of passengers travelled in the ship.



To predict the survival rate of passengers in the ship.

```
import pandas as pd # Import pandas
import numpy as np # Import numpy
import seaborn as sns # Import seaborn
import plotly.express as px # Import plotly express for Interactive Chart
import matplotlib.pyplot as plt # Import matplotlib
from sklearn.metrics import r2_score # r2score
from sklearn.linear_model import LinearRegression # Linear regression model
from matplotlib.animation import FuncAnimation # Import Animation Function
from sklearn.model_selection import train_test_split # train test split
from sklearn.tree import DecisionTreeRegressor # Import Decision Tree Regr
from sklearn.ensemble import RandomForestRegressor # random forest model
from math import sqrt # For squareroot operation
from sklearn.cluster import KMeans # kmeans clustering
from plotly.offline import iplot, init_notebook_mode # Standard plotly impo
init_notebook_mode()
init_notebook_mode(connected=True)
import cufflinks as cf
cf.go_offline()

%matplotlib inline
```

# LIBRARIES USED

# GRAPHS USED FOR FINDINGS

Violin Plot

Bar Plot

Histogram

Count Plot

Heat Map

Cat Plot

Swarm Plot

Interactive  
3D Plot

Dist plot

Scatter  
plot

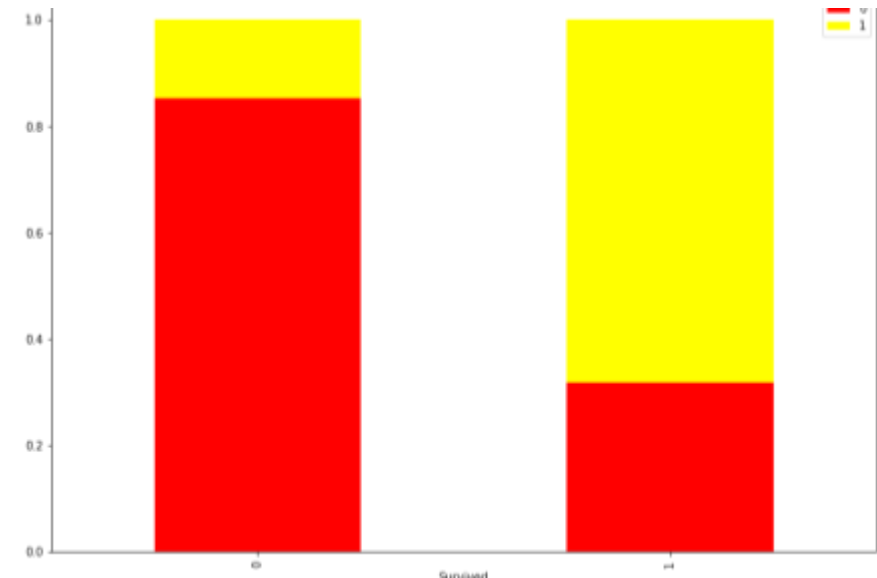
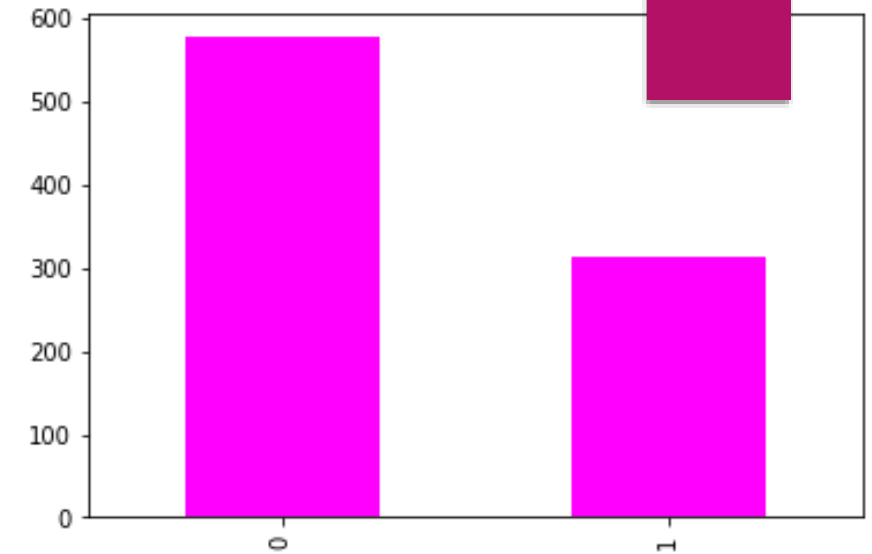
Pie Chart

KMeans  
Clustering

Logistic  
Regression

# Find out total no of male/female passengers

- ▶ The first fig above shows that males are more (577) than females (314) travelling in the ship
- ▶ The sec fig above shows that females(1) survival ratio is high than males (0). Yellow color is survived and red is represented by not survived.

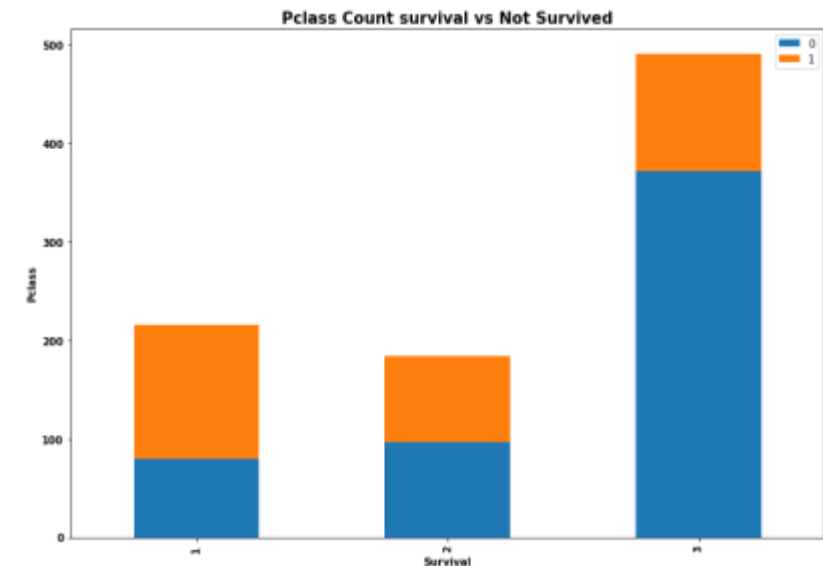
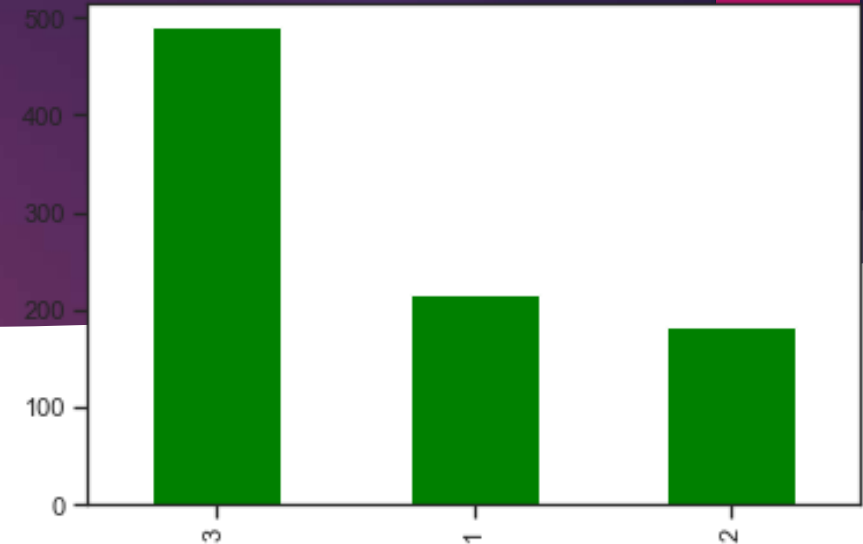




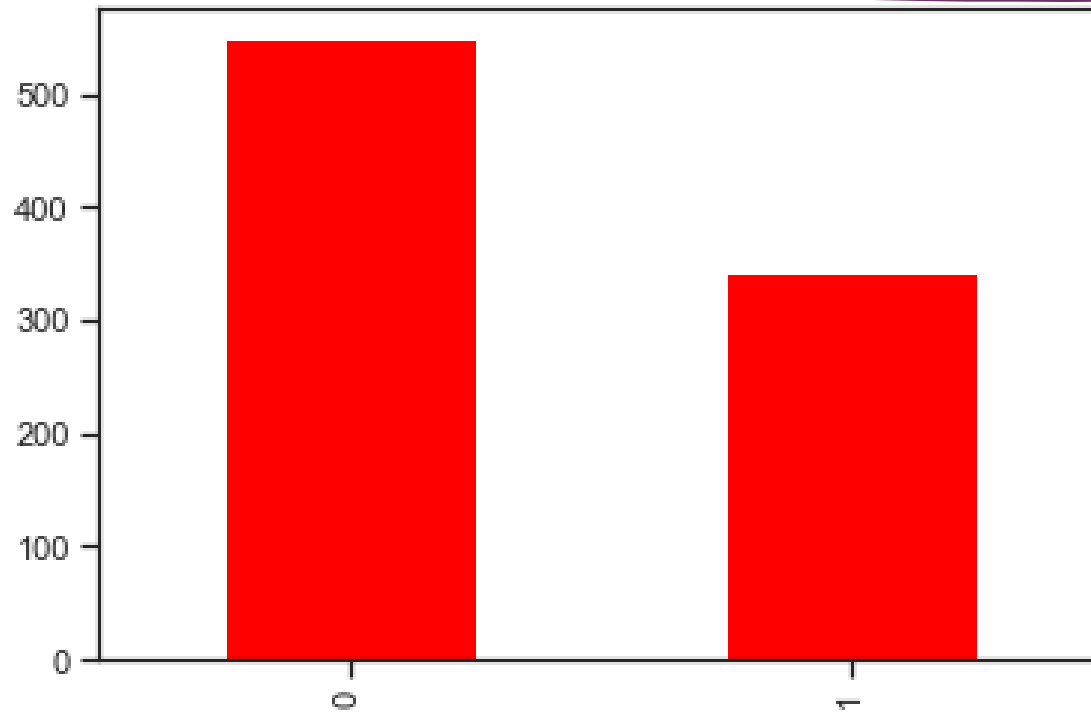
Find total number of passengers in each passenger class

The first fig above shows that more number of passengers (491) are travelling in 3rd class, 216 are travelling in 1st class and least number (184) passengers are travelling in 2nd class

The second fig above shows that survival rate is high in Pclass 1 and survival rate is almost equal in Pclass2 and Pclass1. Orange color is survived and blue color represented by not survived

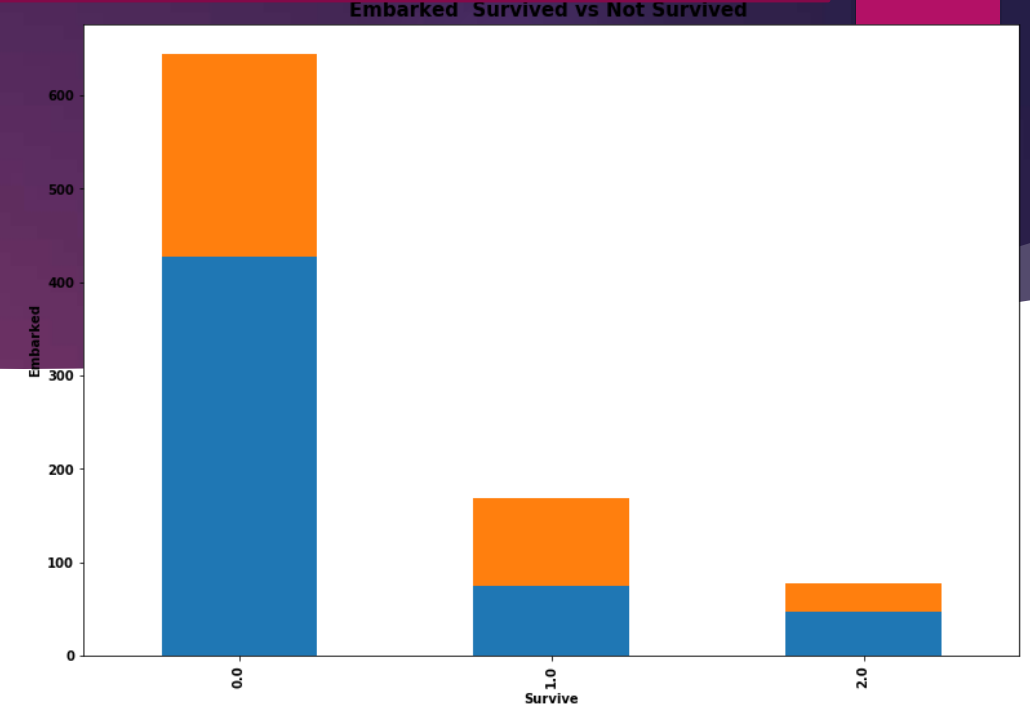
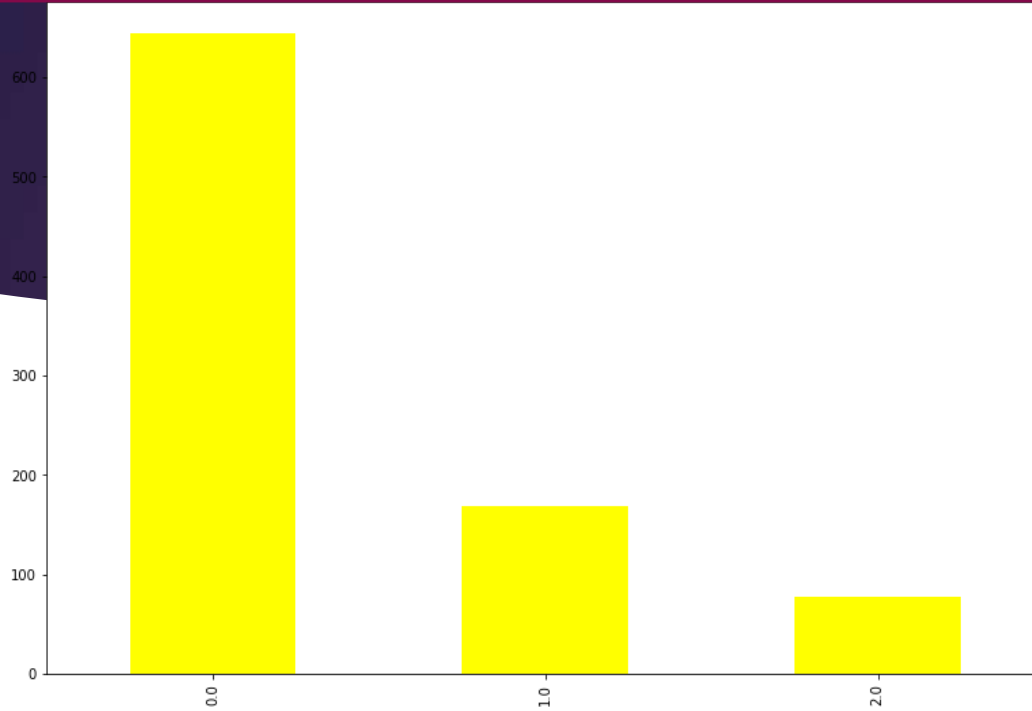


# Find out total number of survived/not survived passengers



The fig above shows that survived ratio is less (342) when compared to not survived (549)

## Find total number of passengers by point of Embarkation



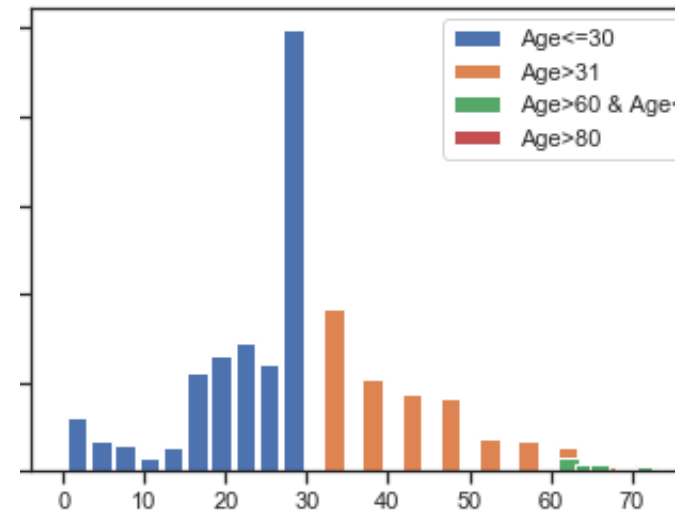
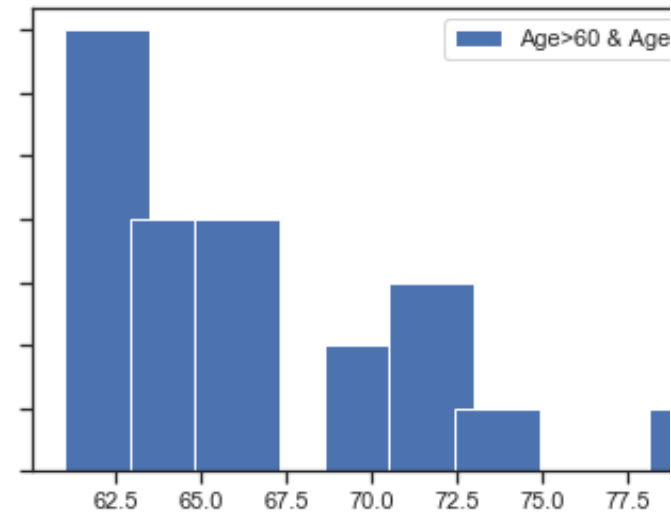
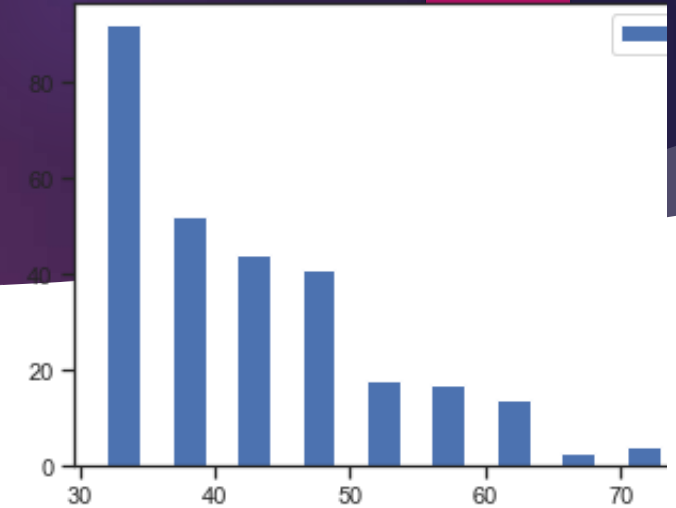
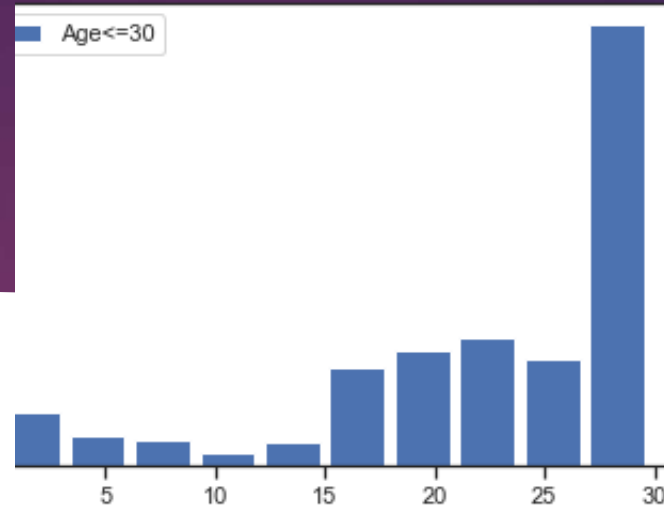
The first fig above shows that More number of passengers are from southhampston(644),the count of passengers from cherbourg are 168 and least number of passengers are from Queenstown(77)

The sec fig above shows that survived ratio is more in passengers from southhampston,next is cherbourg and least is queenstown.

## Histogram of Number of passengers in each age group

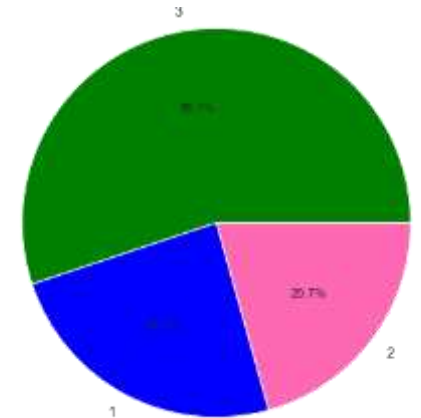
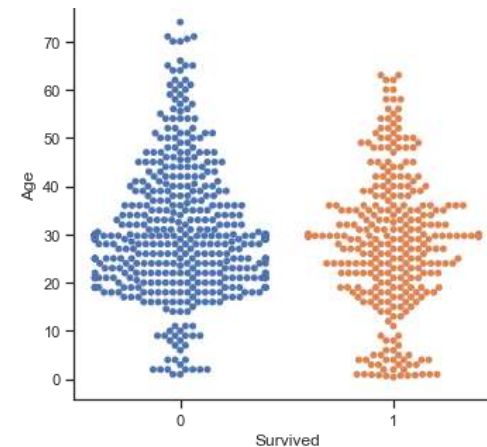
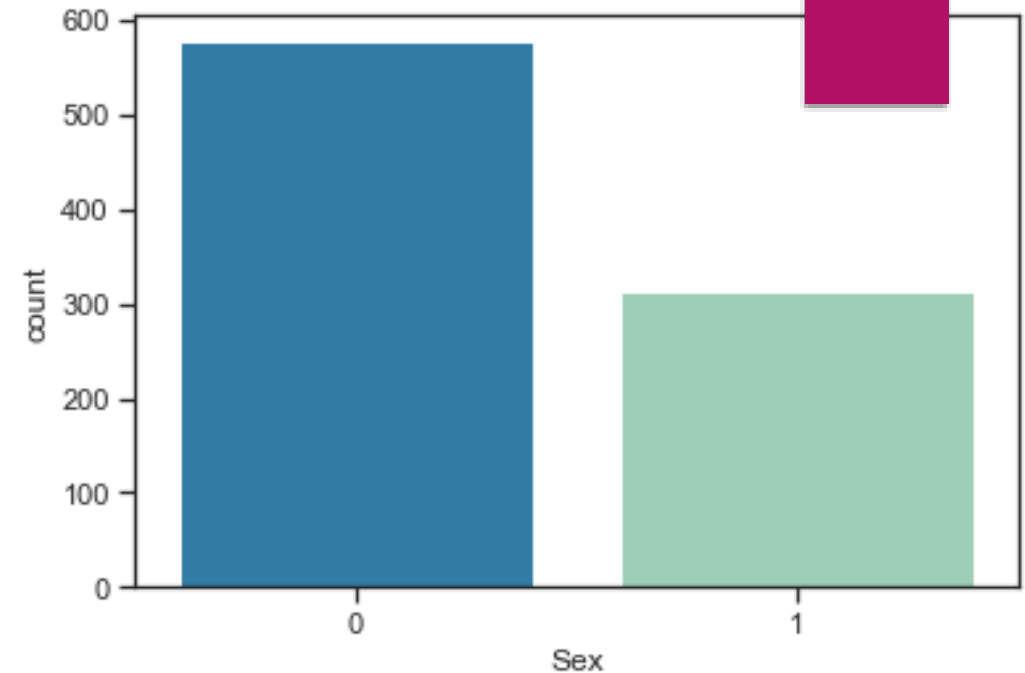
The above figs shows that more number of passengers are in the age group 25-30 and least number of passenger in the age group from 70 and above

Infants below 1 yr are the youngest passengers and passengers above 70 are the oldest passengers travelling in the ship



# Histogram, Pie Chart and swarm plot

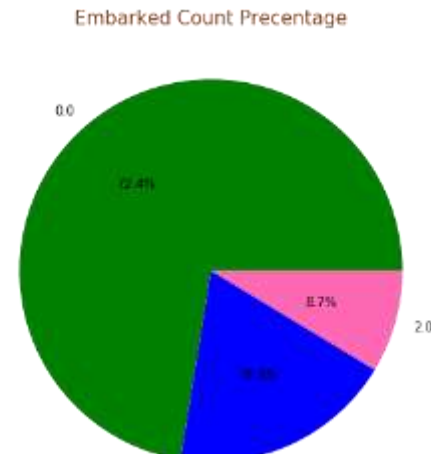
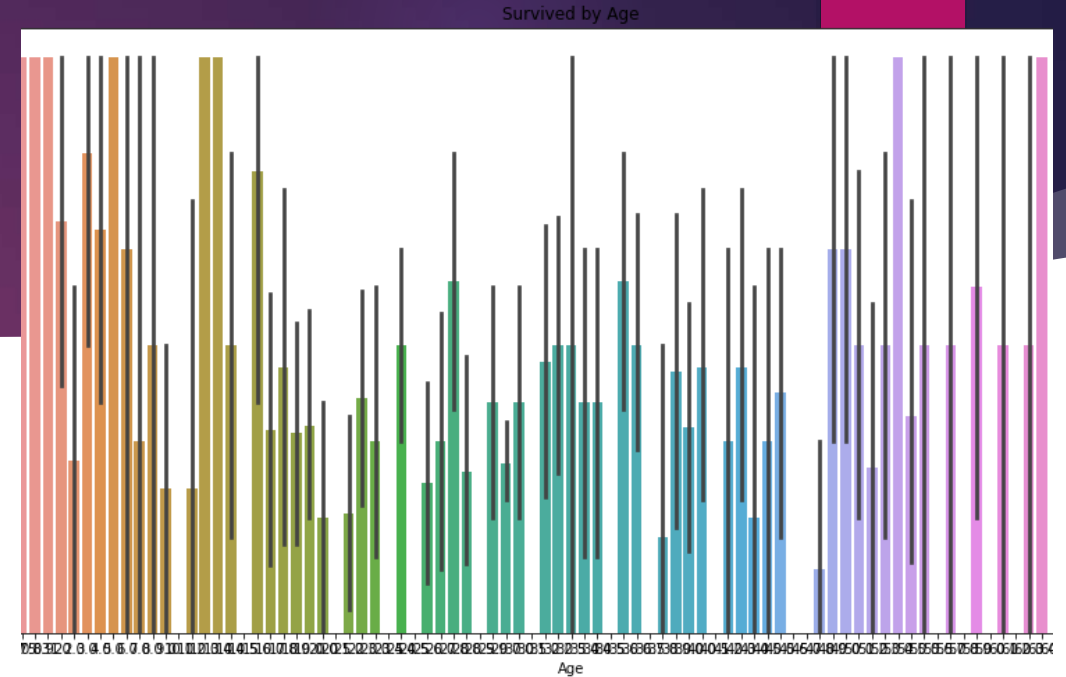
- ▶ The first fig above shows that males are more than females
- ▶ The sec pie chart shows that 55% are travelling in 3rd class, 24% are travelling in 1st class and 20% are travelling in 2nd class
- ▶ The third swarm plot above shows that in both male and females survived passengers are more in the age group between 20-30



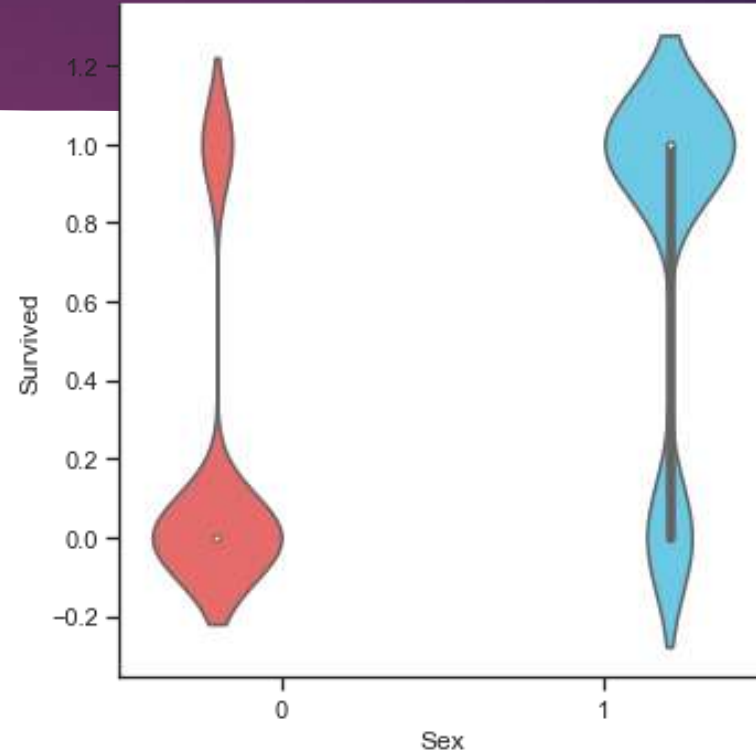
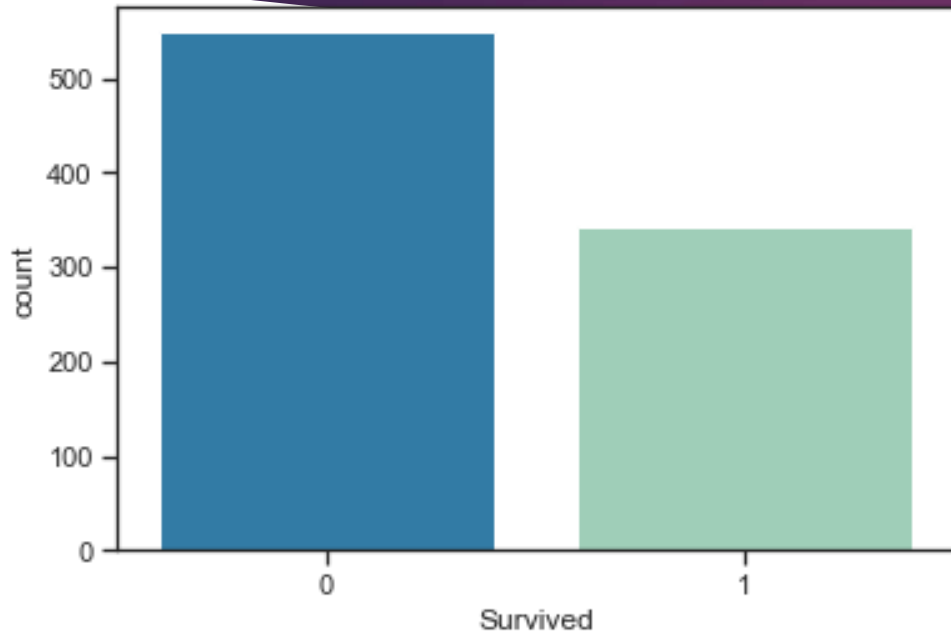


# Interactive Histogram Pie chart and bar chart,countplot

- ▶ The first fig above shows that count of passengers between 28-29 is 224,1 passenger is travelling in age of 80,5 passengers travelling in the age of 70-71 14 are travelling in age group of 0-2 yrs old.
- ▶ The sec pie chart shows that 72% passengers point of embarkation is southhampston,cherbourg is 18% and 8.7% are from queenstown
- ▶ The third fig above shows survival rate in each age group

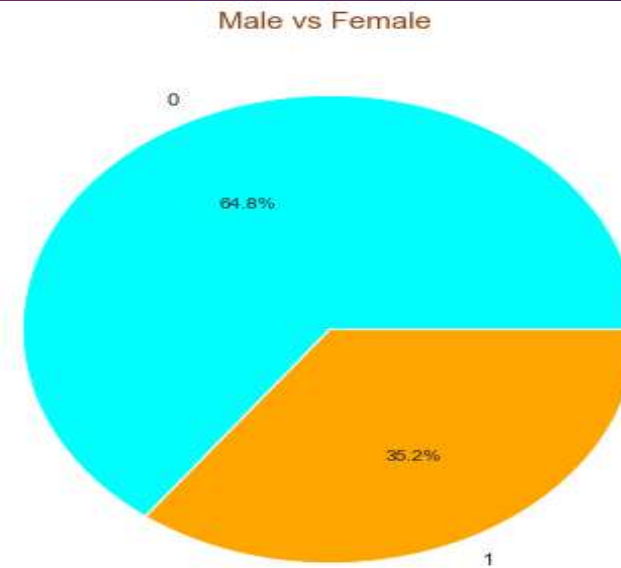
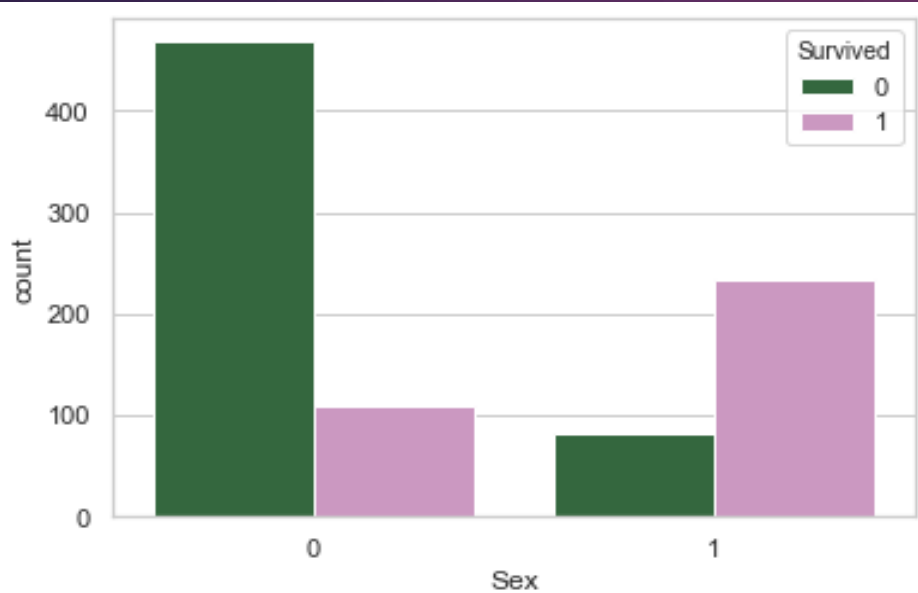


# Histogram & Violin plot



The first fig above shows that not survived (549) are higher than the survived (342)  
The sec violin plot above shows that survived ratio is high in females than males

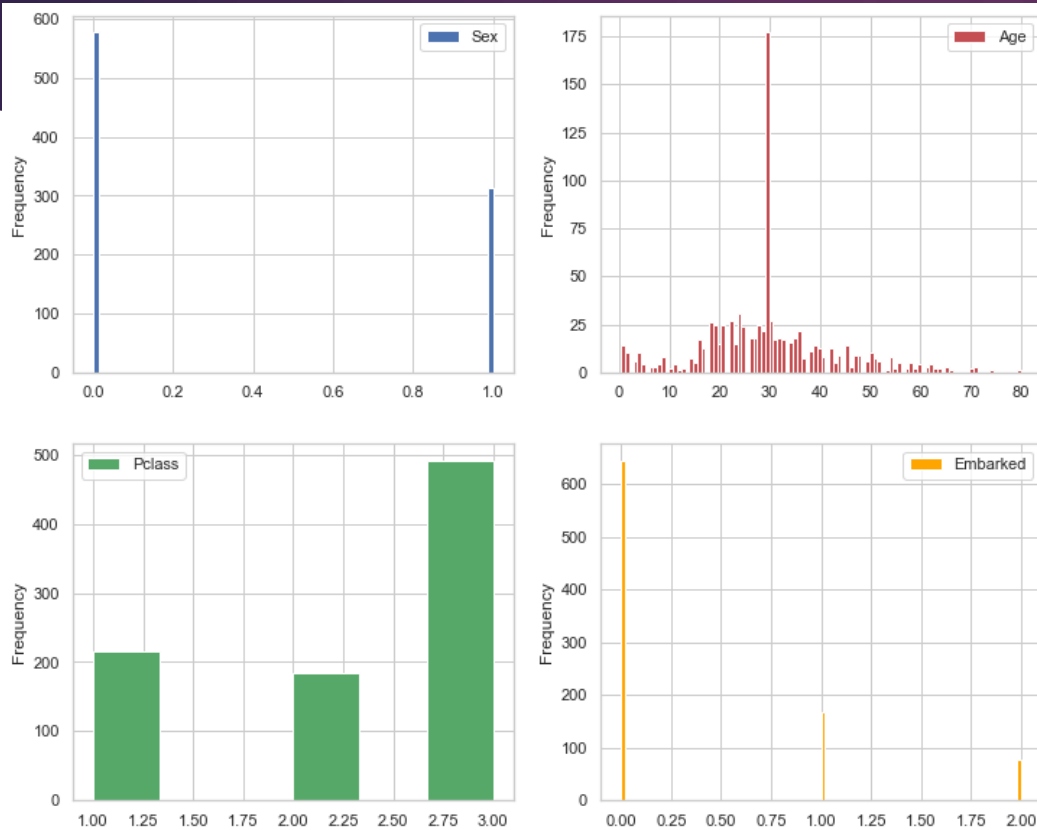
# Count plot & Pie chart



The first fig shows that survived ratio is high in females and not survived ratio is high in males

The sec fig shows that 64.8% are males and 36.2% are females travelling in the ship

# Histogram of count of passengers in each variable

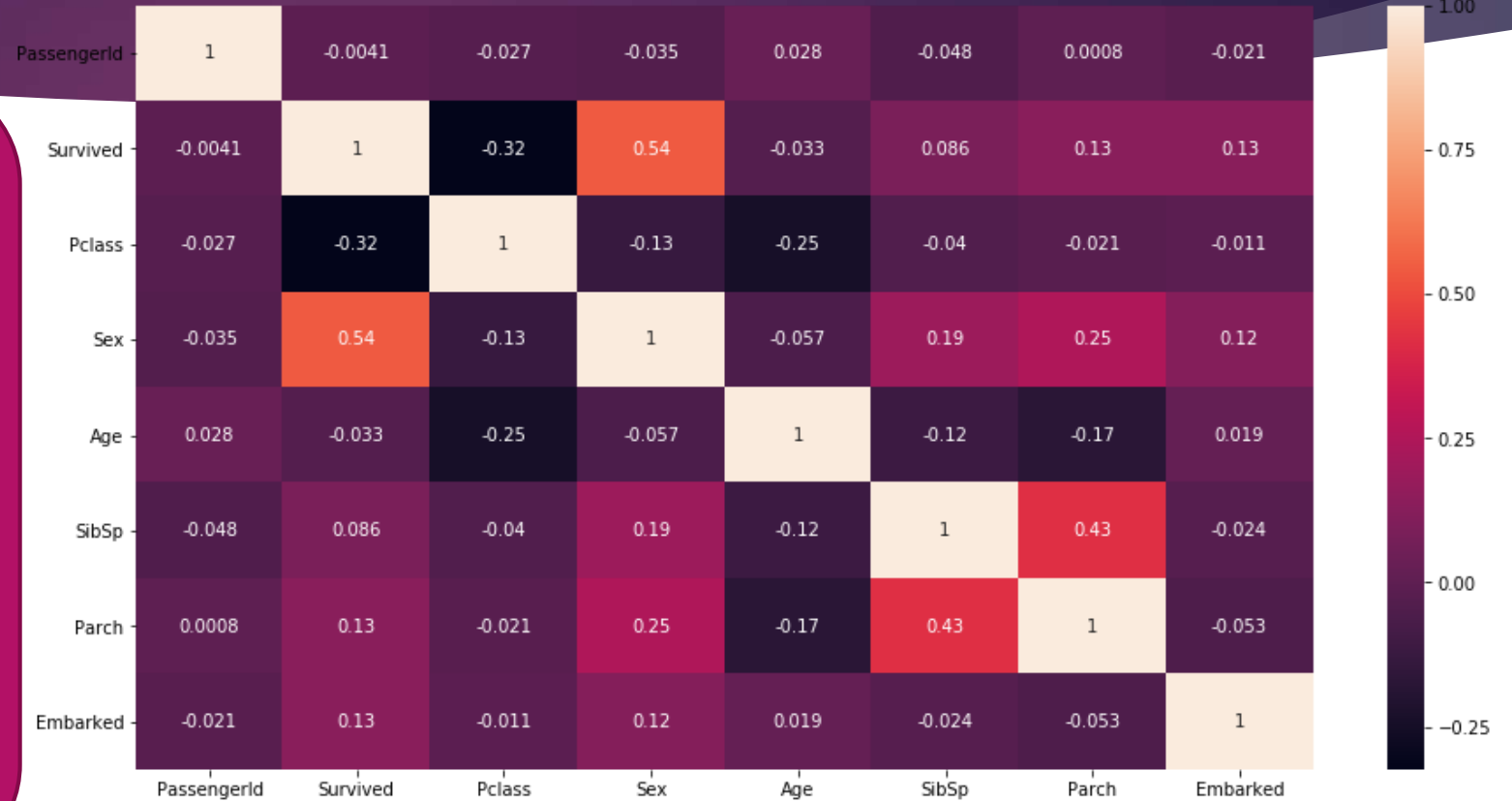


- The above histograms shows the count in each category.
- Males are more than females
- Passengers are more in between age group 25-35
- Passengers count is more in Pclass 3
- Passengers count is more from Southampton

# Heat Map

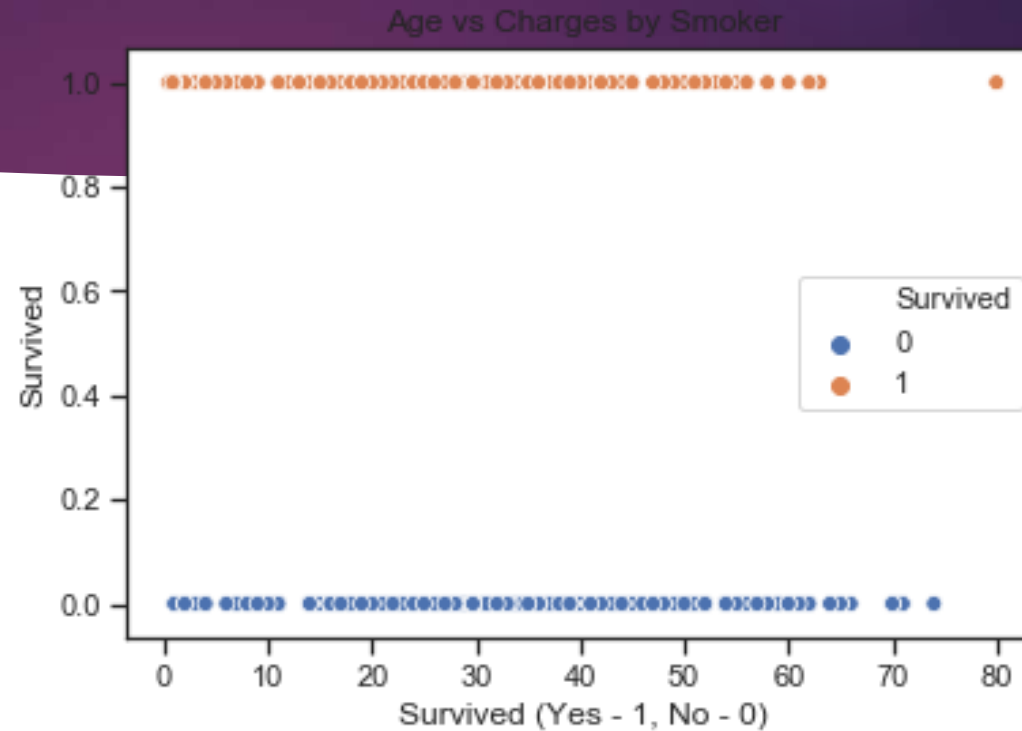
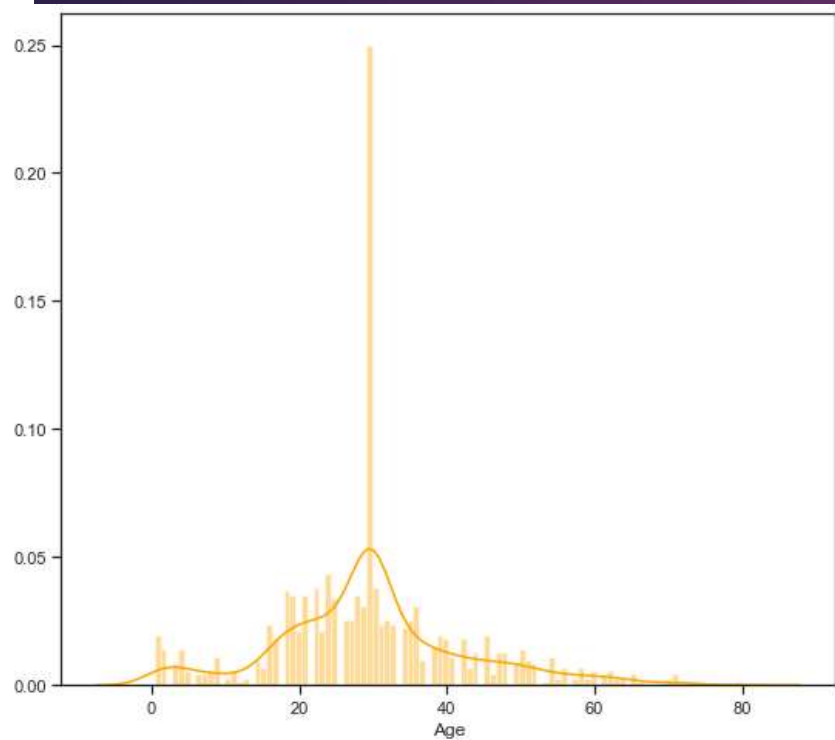
The above heatmap shows correlation between different variables

There is a high correlation between Sex and Survived and Pclass and Survived, Age and Pclass.





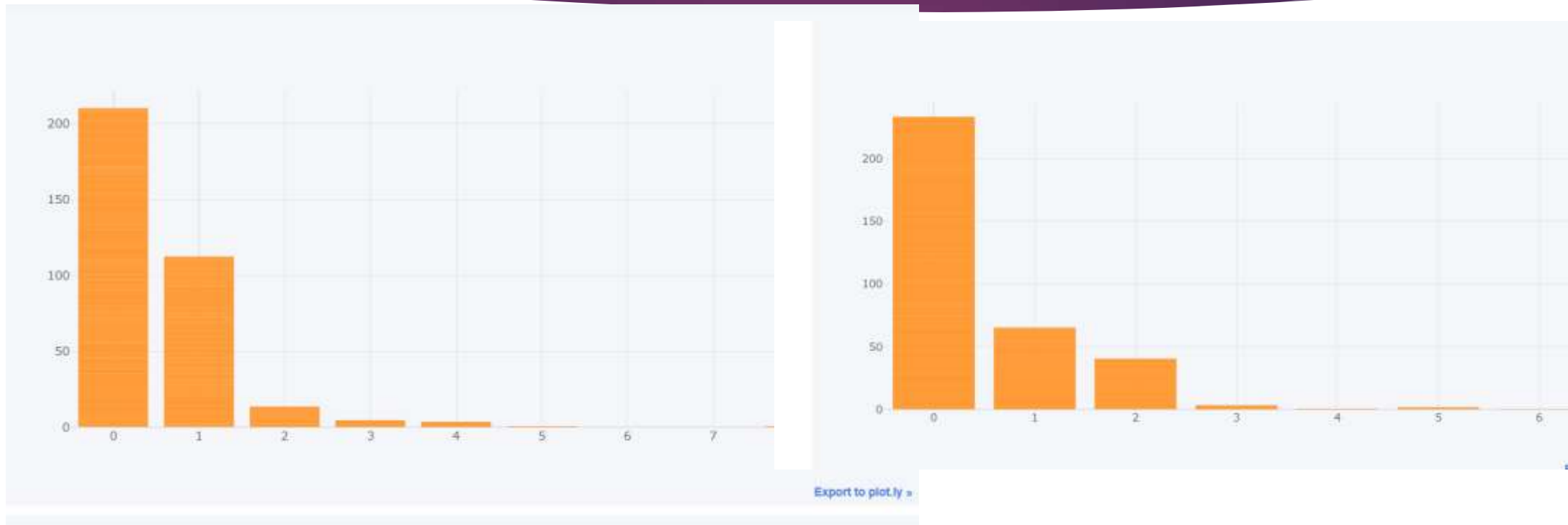
# Scatterplot & Dist plot



The above dist plot shows that there is a spike of passengers in between age group 25-35

The above scatter plot shows Survived and not survived in different age groups

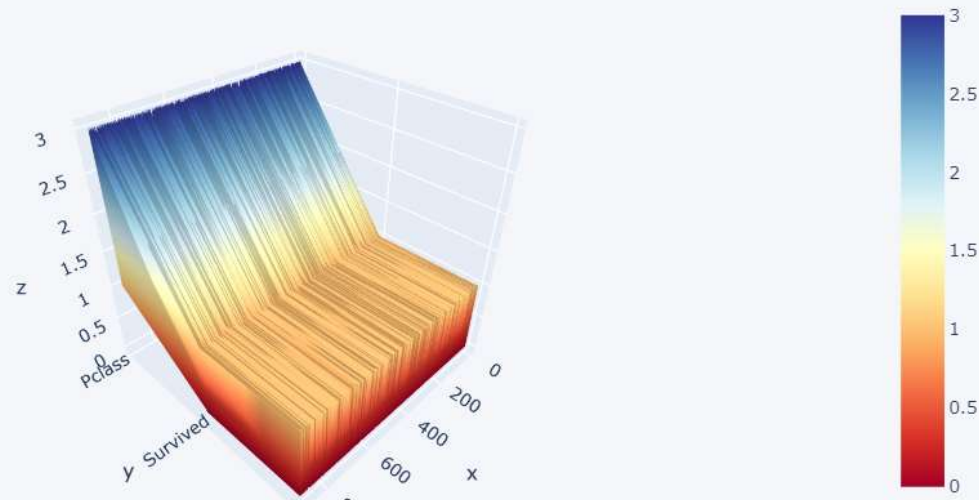
# Interactive barplot



The first fig shows that passengers travelling alone survival ratio is high compared to passengers travelling with spouse and siblings

The second fig shows that passengers travelling alone survival ratio is high compared to travelling with parents and siblings

# Interactive 3D plot



[Export to plot.ly »](#)

Interactive 3D plot  
shows Survival rate by  
Pclass and sex

# Machine Learning Models



R2 Score for Multiple Linear Regression Model is 0.42



R2 Score for Random Forest Linear Regression Model is 0.38



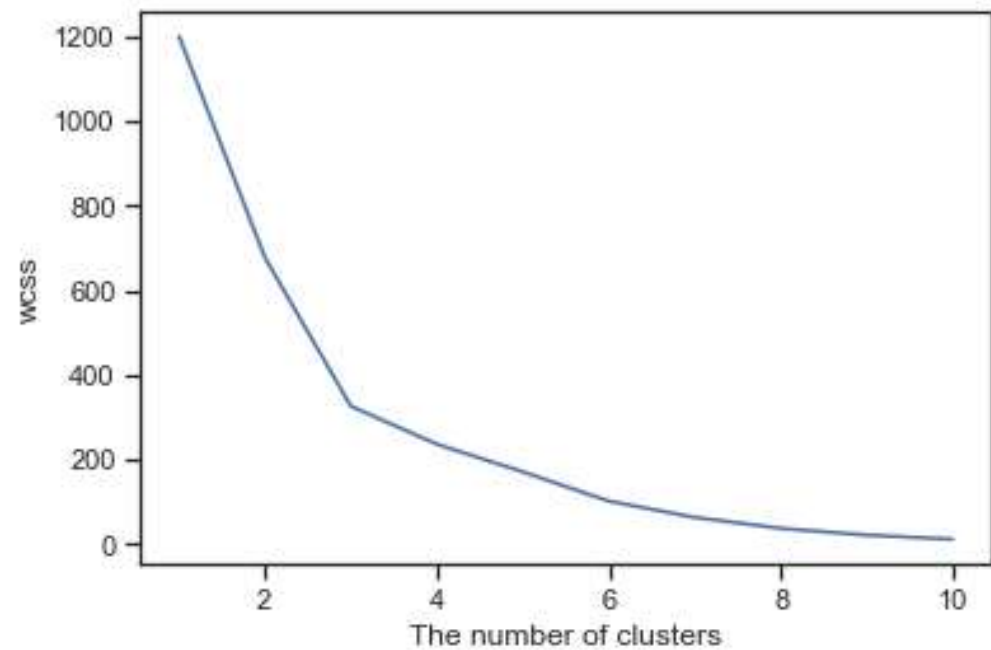
R2 Score for Decision Tree Regression Model is 0.10



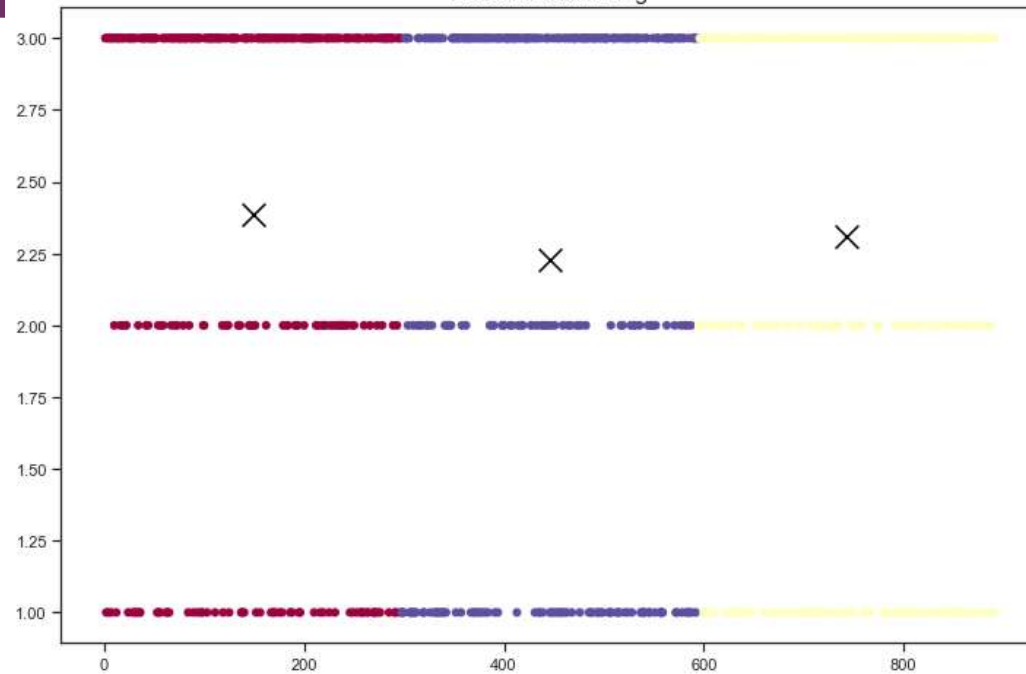
Hence Multiple Regression Model is the best model

# Kmeans Clustering

The Elbow Method



Kmeans Clustering



The elbow method is a technique to choose the most number of clusters. \*The graph above shows that there are 3 clusters in the data set.



# Conclusion :

Survival ratio is high in females than males

More number of passengers are in the age group 25-25

More number of passengers are from point of embarkation southampton

More number of passengers have survived in Pclass1

More number of passengers travelled in Pclass 3

More number of male passengers are than the female passengers

Survival rate is more in Southhampston embarkment than cherbourg and queenstown

Age group who have survived more are in the range of 25-30

Youngest traveller is below 1 yr and oldest traveller is 80 yr old

Out of total 891 passengers 342 have survived and 549 have died

Avg age of passengers is 30. Avg male age is 24 and Avg fem age is 23

Females survival rate is 74% and males survival rate is 19%

R2 score for multiple linear regression is 0.42 which is best score to test new data for predicting survival rate.

