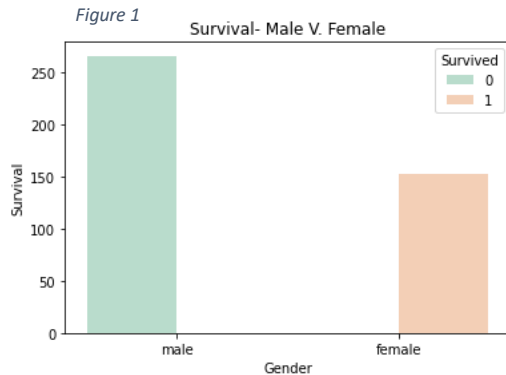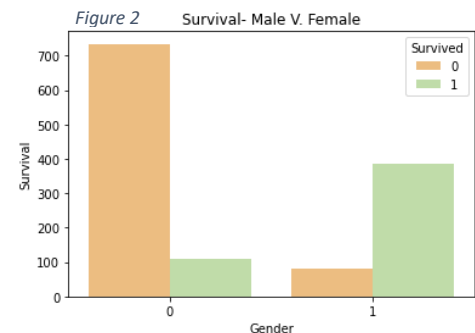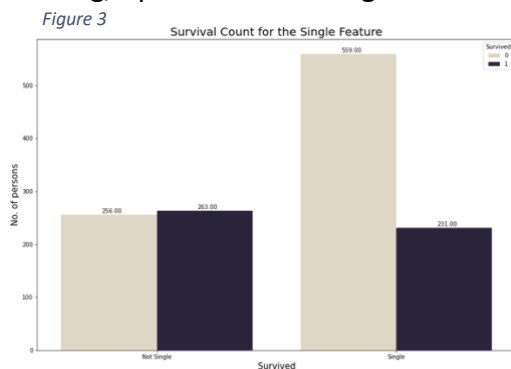## RISK ASSESSMENT ON THE TITANIC

To begin with the analysis, the summary statistics of train and test data was prepared. As initial step I checked the relationship between various features and the survival rate using bar graphs and correlation matrix. Upon checking the distribution of survival and sex of the test data, I



Figure 1

could find that all females survived while all the males died(fig1). This indicated that there must be a bias in the splitting of data into train and test sets. And sex was highly correlated feature with survival. Therefore, I combined the test and train data. As we can see in fig2

(0-Male, 1-Female) the survival is now distributed unbiased. I filled in the missing values



Figure 2

in age and fare with mean values and missing values in embarked with mode. The variables that were just for the individual and that had no predictive power were removed, which are Name, Passenger Id, Cabin, Ticket type. I combined parent child feature and sibling/ spouse feature to generate a new



Figure 3



Figure 4

variable to see if travelling alone or not affected the chances of survival (Fig3). The features that I used after cleaning and EDA are in fig4. I ran a decision tree model to check how many survived on the basis of the chosen features. After cost complexity pruning, I chose the alpha that would make optimal amount of splits to predict the model with utmost accuracy which is 0.0025(Fig5 &6). I also ran a probit model to predict the survival, however decision tree(85.9%) was more accurate than probit model(84.7%). Fig 7 & 8 is the confusion matrix presenting the actual and predicted value by decision tree and probit respectively
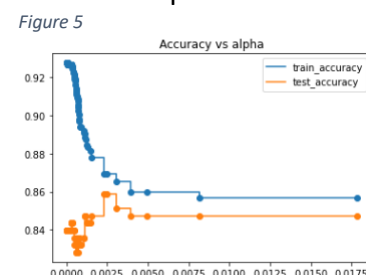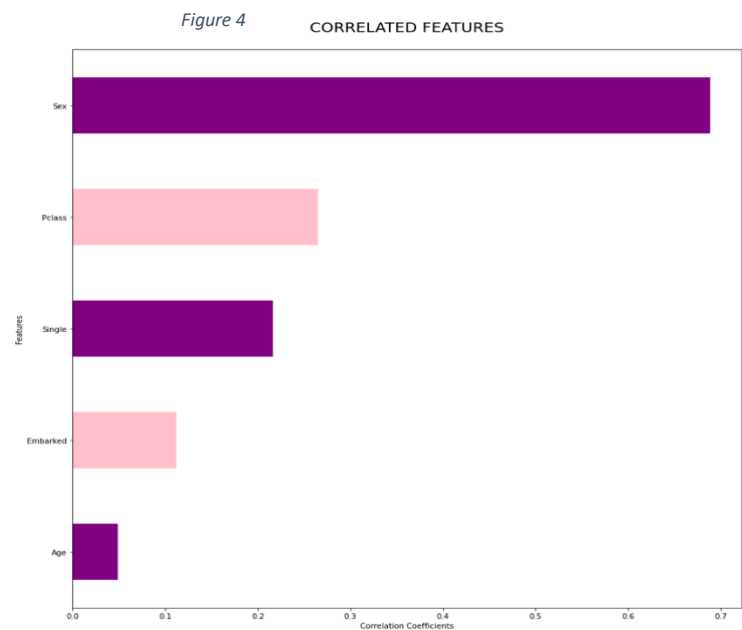


Figure 5

*Figure 6*