# Titanic survival rate classification report

## Phase(1):Discovery

Motivation for developing this model is to offer improvement and risk assessment for ship insurance. Insurance companies provide coverage for ships and vessels, including passenger ships like the Titanic. Accurate risk assessment is crucial for determining appropriate insurance premiums and coverage for such ships. By analyzing the Titanic survival dataset, insurance companies can gain insights into the factors that influenced survival rates, which can help improve risk assessment and underwriting processes for passenger ships.

By utilizing Titanic Survival Dataset for Risk Modeling through:

1. Feature Analysis: The Titanic survival dataset contains information about passengers' characteristics, such as age, gender, passenger class, and family size. Insurance companies can analyze this data to identify which features were correlated with higher or lower survival rates. For example, they may find that fare class or gender class played a significant role in survival.

2. Risk Modeling: Based on the analysis, insurance companies can develop risk models that take into account the relevant features identified in the dataset. These models can be used to estimate the risk of future passenger ships based on similar characteristics.

# Phase(2):Data Preparation

- After exploring the available dataset it was identified that it consists of the following classes survival '0=>No// 1=>Yes', pclass 'Ticket class', Sex, Age, sibsp 'No of siblings/spouses who were abroad', parch 'No of parents/ children were abroad', ticket 'ticket number', fare 'passenger fare', cabin 'cabin number', embarked 'port of embarkation'.

- The following transformations were done on dataset which are:

  Considered as a feature engineering step, It is suggested that the suitable column/ target variables for developing such a model will be focusing on ( pclass, sex, age, fare, survived). These features are considered to be important features that will contribute in the model development to offer the company desired results.
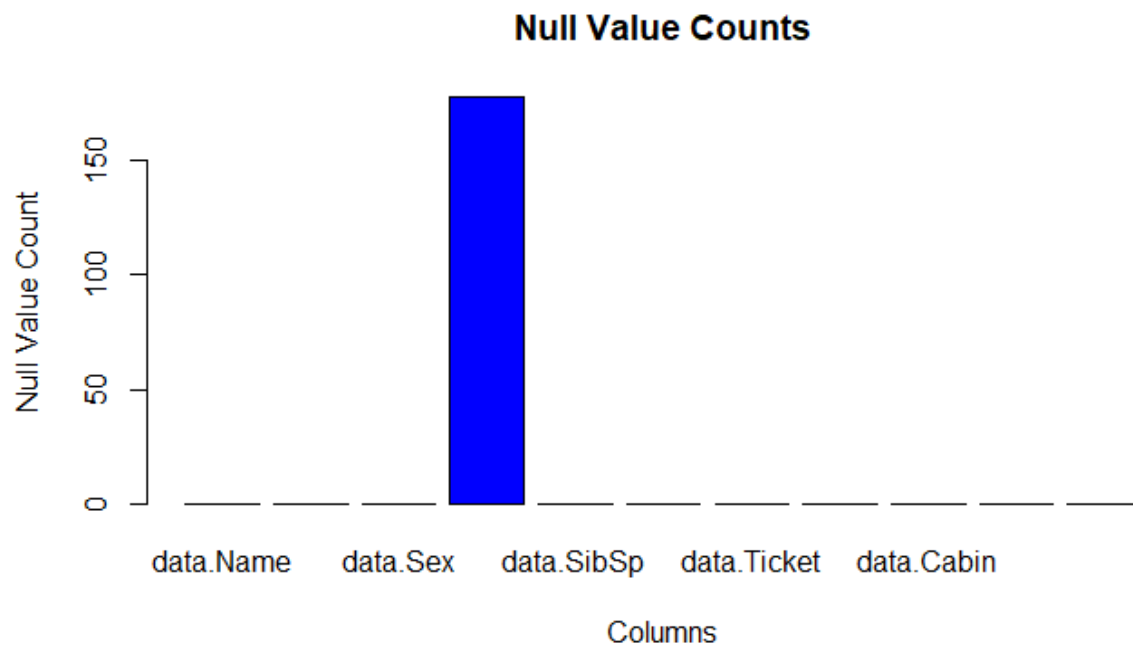
  Therefore, all columns were dropped from the data frame and previously mentioned columns were kept.

- Since it's hard to deal with categorical data as shown in sex column, therefore male and female classes were changed into binary classes of 0=> Male and 1=>Female.

- On viewing the whole data it found nulls, especially in the Age column, we had two options. Either to omit these records or fill the empty records with the mean. The later solution was much better in order not to lose any of our data and to provide a good model.

## Null Value Counts
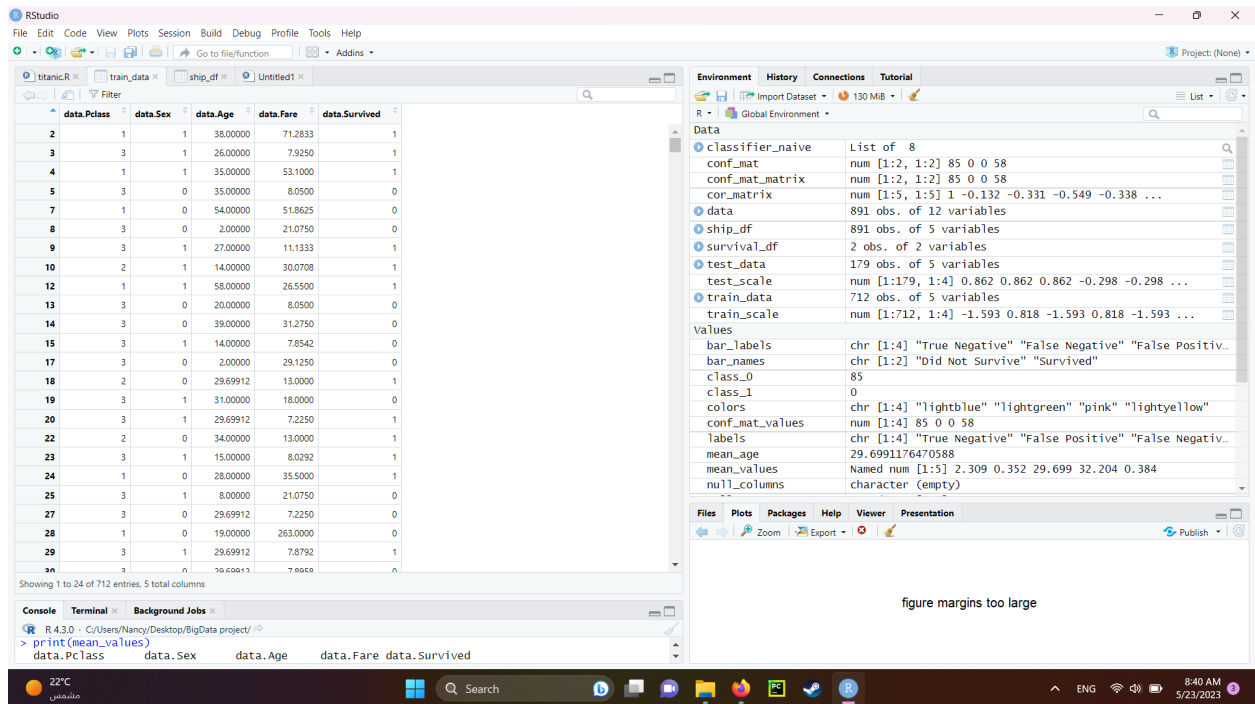


The following bar chart is a representation for nulls in the Age column. After applying the previously explained step.

## Null Value Counts

Null Value Count

data.Pclass    data.Sex    data.Age    data.Fare    data.Survived

Columns

## Mean Values

Mean

data.Pclass    data.Sex    data.Age    data.Fare    data.Survived
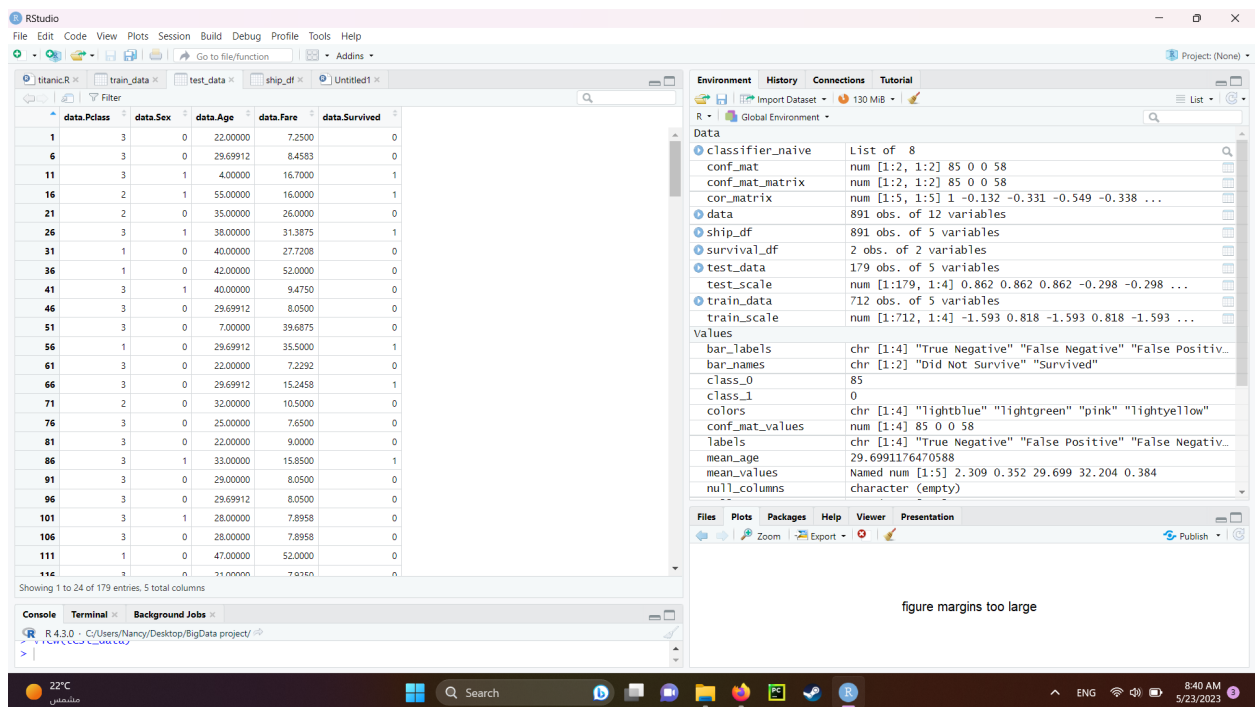
Columns

**std Values**



# Phase(3):Model Planning

Suggested techniques to be applied for such a problem is, since we need the model to predict based on inserted data if it's most likely to survive or not from the titanic based on inserted data. We have applied the Naive Bayes model. By splitting our dataset into test and training sets.

Our training set data.



Our testing set data.

# Phase(4):Model Building

- The suitable models to use for such problem type will be the following:

  1) Naive Bayes model which provides probabilistic results
  2) Logistic regression algorithm
  3) Decision
- Training and testing data have been splitted of ratio 80% train and 20% test.
- All columns are set as X variable and Survival column is set as Y variable (Target variable)..

# Phase(5):Communicate results

```
Confusion Matrix and Statistics

      Predicted
Actual  0  1
     0 85  0
     1  0 58

             Accuracy : 1
               95% CI : (0.9745, 1)
  No Information Rate : 0.5944
  P-Value [Acc > NIR] : < 2.2e-16

                Kappa : 1

 Mcnemar's Test P-Value : NA

          Sensitivity : 1.0000
          Specificity : 1.0000
       Pos Pred Value : 1.0000
       Neg Pred Value : 1.0000
           Prevalence : 0.5944
       Detection Rate : 0.5944
 Detection Prevalence : 0.5944
    Balanced Accuracy : 1.0000

     'Positive' Class : 0
```
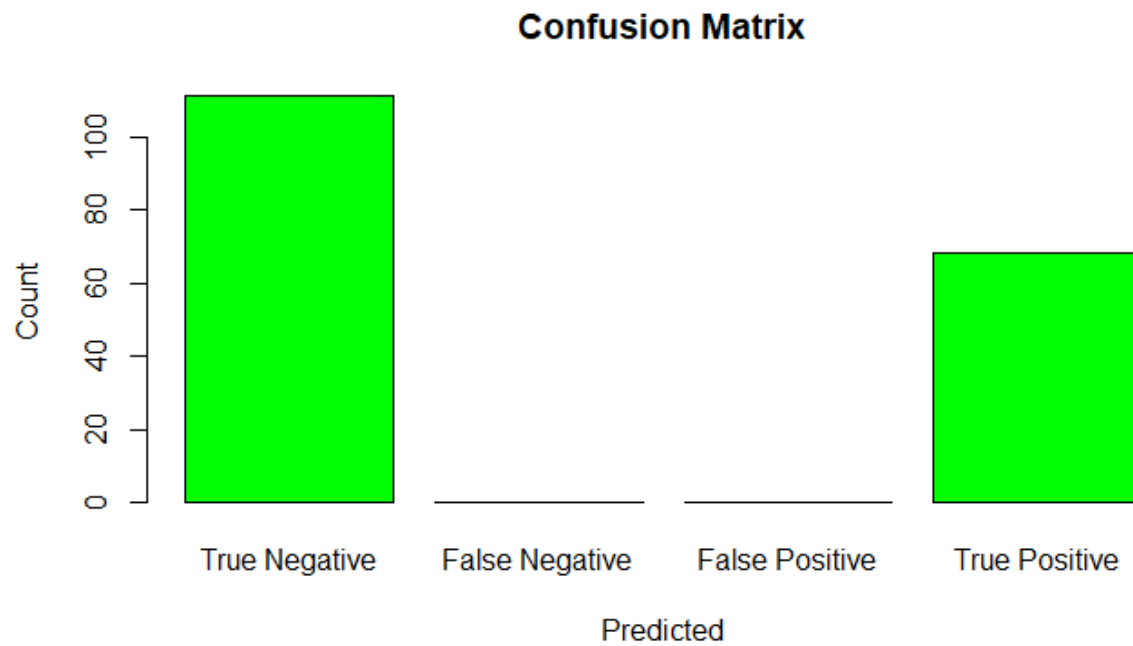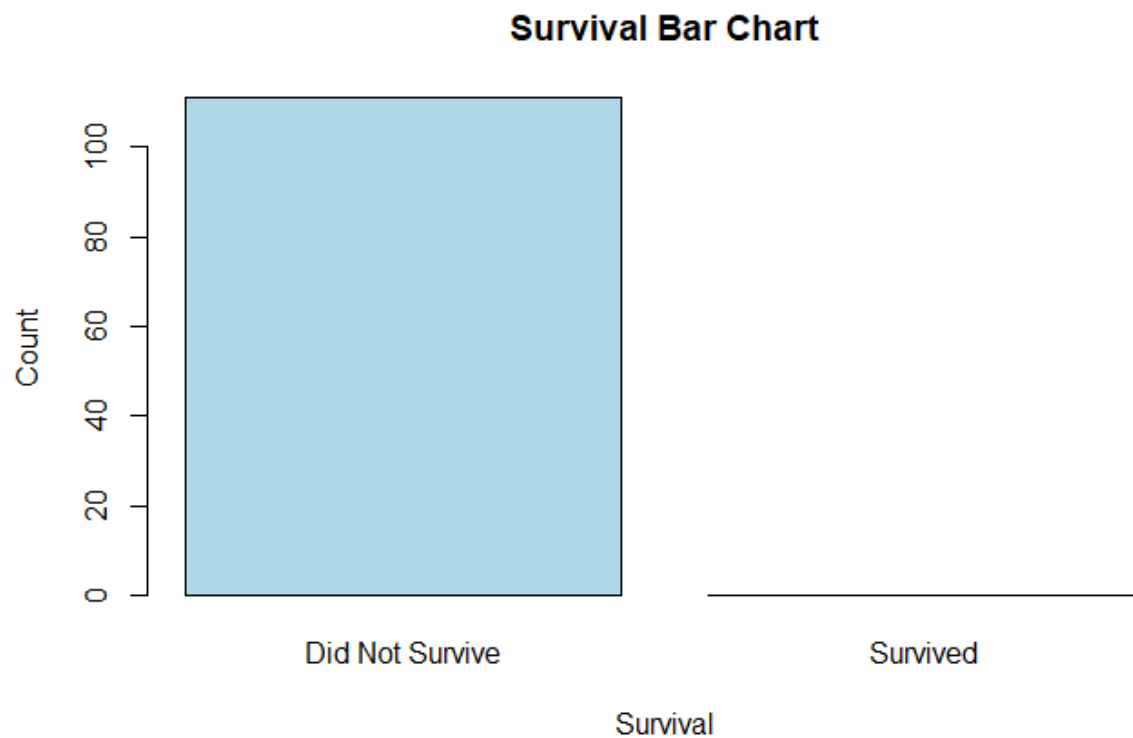
## Confusion Matrix



From the following bar chart it's clear that our model has predicted all entries true as per True negative(Meaning that it has predicted this much who won't survive) and true positive (Meaning that it predicted this much that are most likely to survive).

## Survival Bar Chart

# Phase(6):Operationalize