

Titanic Analysis And ML

Clinton Mwachia

2022-10-01

Contents

The Data	1
Train Data	1
Test Data	2
Data transformation	2
Missing Data	2
EDA: Most of the survivors were female	4
EDA: Passengers with first class ticket have higher survival rates	6
EDA: AGE, FARE BY SURVIVAL	7
EDA: Check for outliers	8

The Data

Train Data

```
## PassengerId      Survived  Pclass      Name
## Min.   : 1.0      Min.   :0.0000  Min.   :1.000  Length:891
## 1st Qu.:223.5      1st Qu.:0.0000  1st Qu.:2.000  Class :character
## Median :446.0      Median :0.0000  Median :3.000  Mode  :character
## Mean   :446.0      Mean   :0.3838  Mean   :2.309
## 3rd Qu.:668.5      3rd Qu.:1.0000  3rd Qu.:3.000
## Max.   :891.0      Max.   :1.0000  Max.   :3.000
##
##      Sex      Age      SibSp      Parch
## Length:891  Min.   : 0.42  Min.   :0.000  Min.   :0.0000
## Class :character  1st Qu.:20.12  1st Qu.:0.000  1st Qu.:0.0000
## Mode  :character  Median :28.00  Median :0.000  Median :0.0000
##                      Mean   :29.70  Mean   :0.523  Mean   :0.3816
##                      3rd Qu.:38.00  3rd Qu.:1.000  3rd Qu.:0.0000
##                      Max.   :80.00  Max.   :8.000  Max.   :6.0000
##                      NA's   :177
##      Ticket      Fare      Cabin      Embarked
## Length:891      Min.   : 0.00  Length:891  Length:891
## Class :character  1st Qu.: 7.91  Class :character  Class :character
```

```
## Mode :character Median : 14.45 Mode :character Mode :character
## Mean : 32.20
## 3rd Qu.: 31.00
## Max. :512.33
##
```

Test Data

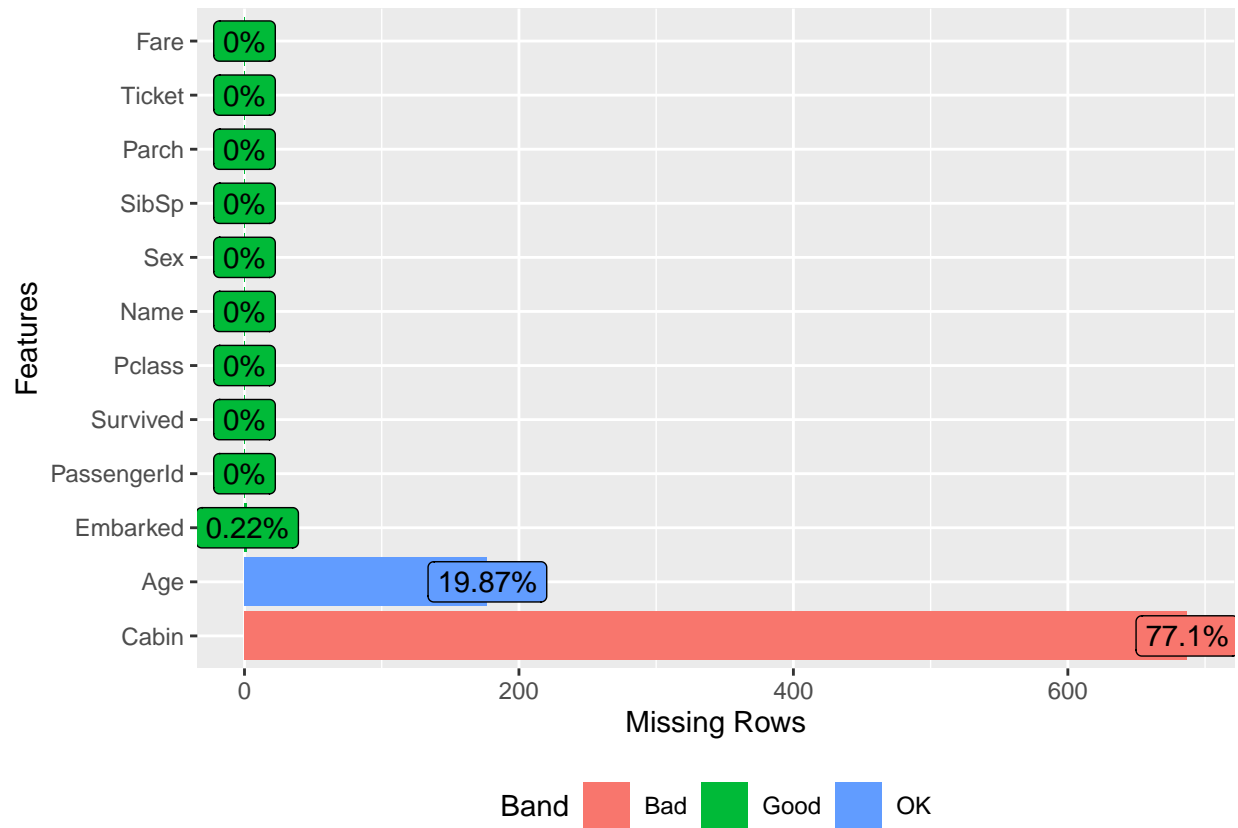
```
## PassengerId      Pclass      Name      Sex
## Min. : 892.0 Min. :1.000 Length:418 Length:418
## 1st Qu.: 996.2 1st Qu.:1.000 Class :character Class :character
## Median :1100.5 Median :3.000 Mode :character Mode :character
## Mean :1100.5 Mean :2.266
## 3rd Qu.:1204.8 3rd Qu.:3.000
## Max. :1309.0 Max. :3.000
##
## Age      SibSp      Parch      Ticket
## Min. : 0.17 Min. :0.0000 Min. :0.0000 Length:418
## 1st Qu.:21.00 1st Qu.:0.0000 1st Qu.:0.0000 Class :character
## Median :27.00 Median :0.0000 Median :0.0000 Mode :character
## Mean :30.27 Mean :0.4474 Mean :0.3923
## 3rd Qu.:39.00 3rd Qu.:1.0000 3rd Qu.:0.0000
## Max. :76.00 Max. :8.0000 Max. :9.0000
## NA's :86
## Fare      Cabin      Embarked
## Min. : 0.000 Length:418 Length:418
## 1st Qu.: 7.896 Class :character Class :character
## Median : 14.454 Mode :character Mode :character
## Mean : 35.627
## 3rd Qu.: 31.500
## Max. :512.329
## NA's :1
```

Most of the variables are not in the appropriate data types, lets transform them. survived, pclass, sex etc and all character should be factor

Data transformation

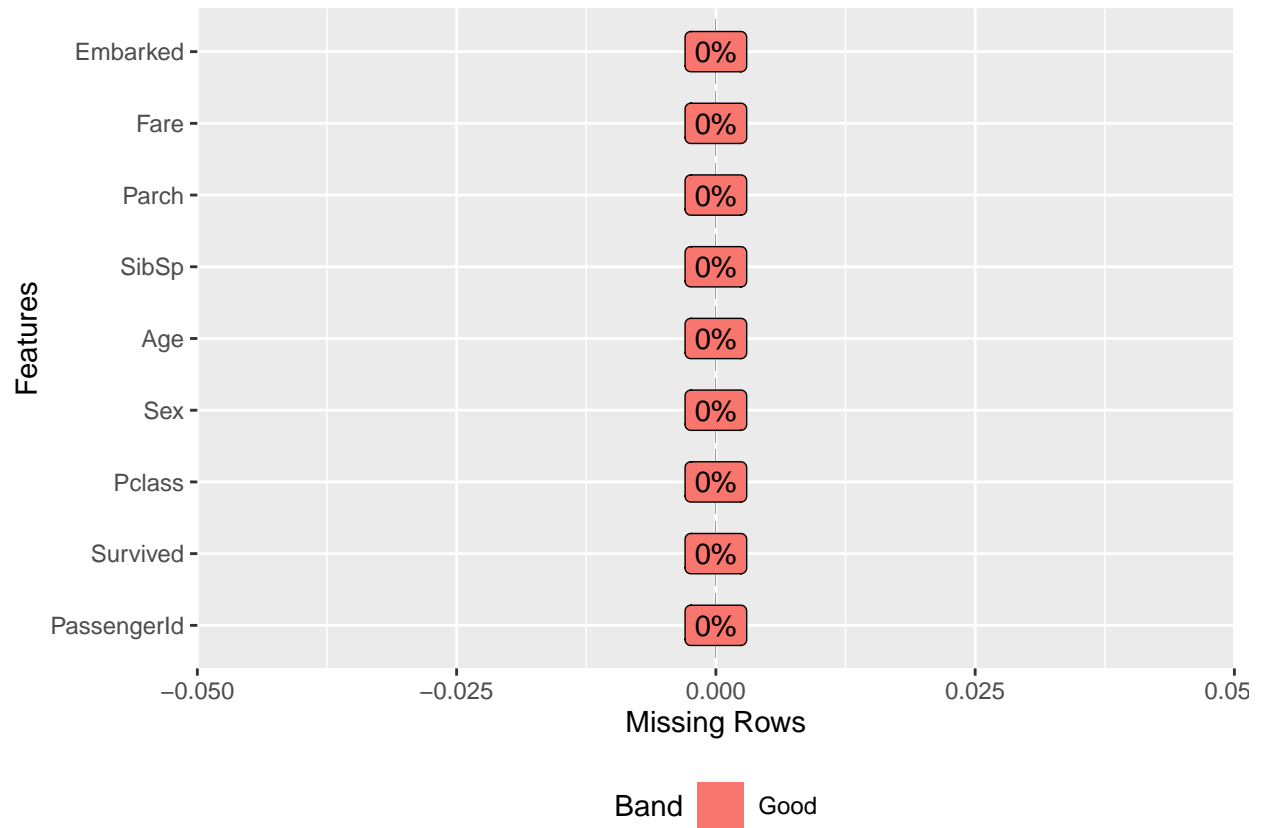
Missing Data

```
## PassengerId  Survived  Pclass  Name      Sex      Age
##          0          0          0          0          0      177
## SibSp      Parch      Ticket  Fare      Cabin  Embarked
##          0          0          0          0      687          2
```



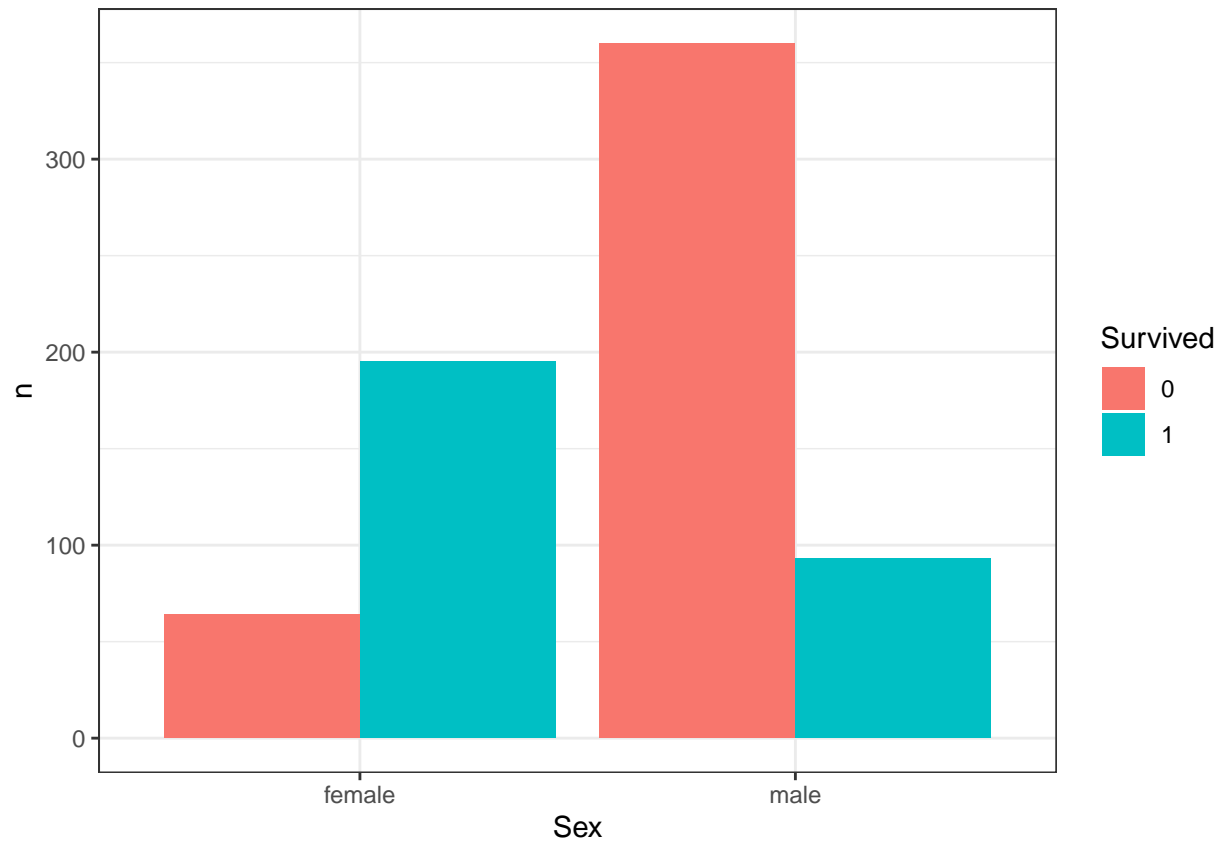
Lets drop column transform, too many missing data points. We will also drop name and ticket columns. Lets also drop missing data points,(You can try imputing them and then compare results to see if it improves the model).

```
## PassengerId  Survived  Pclass
##           0         0         0
##           0         0         0
##           0         0         0
##           0         0         0
```



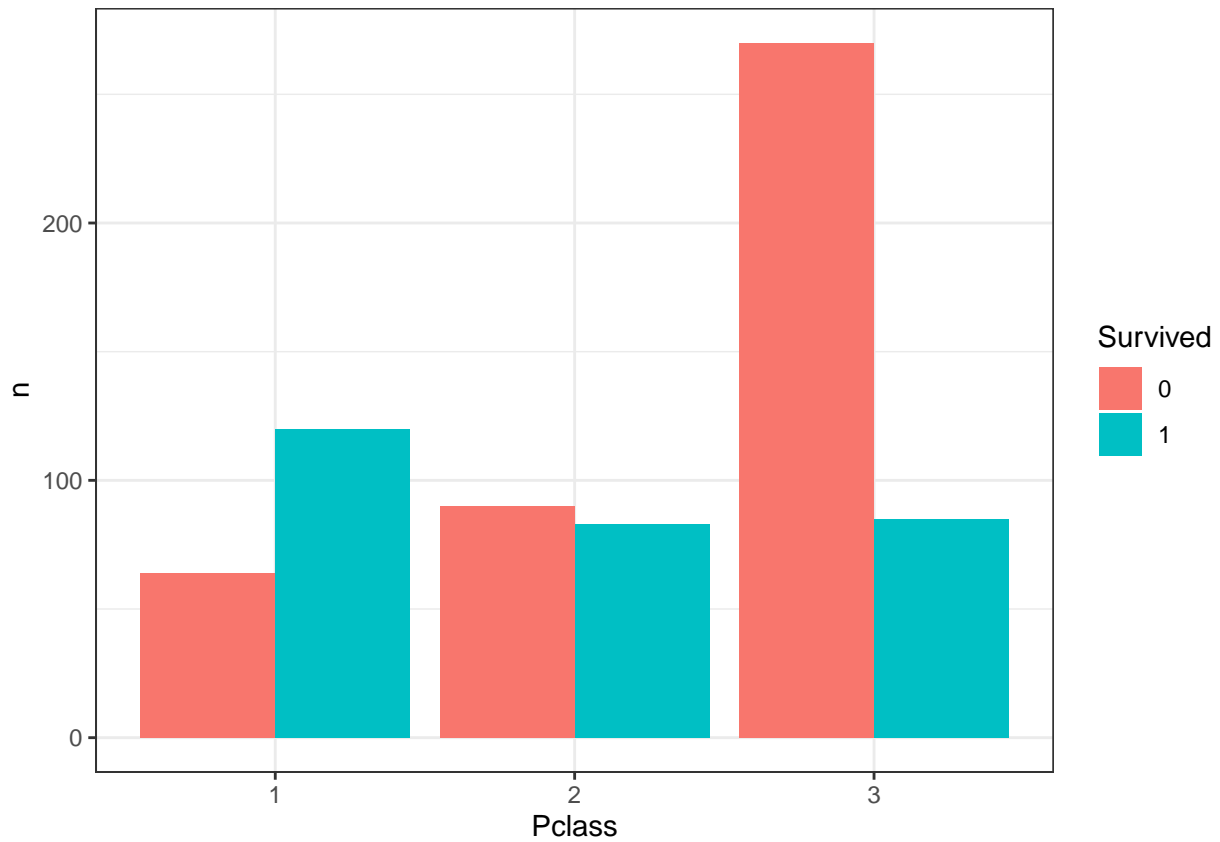
EDA: Most of the survivors were female

Lets find out if this is true



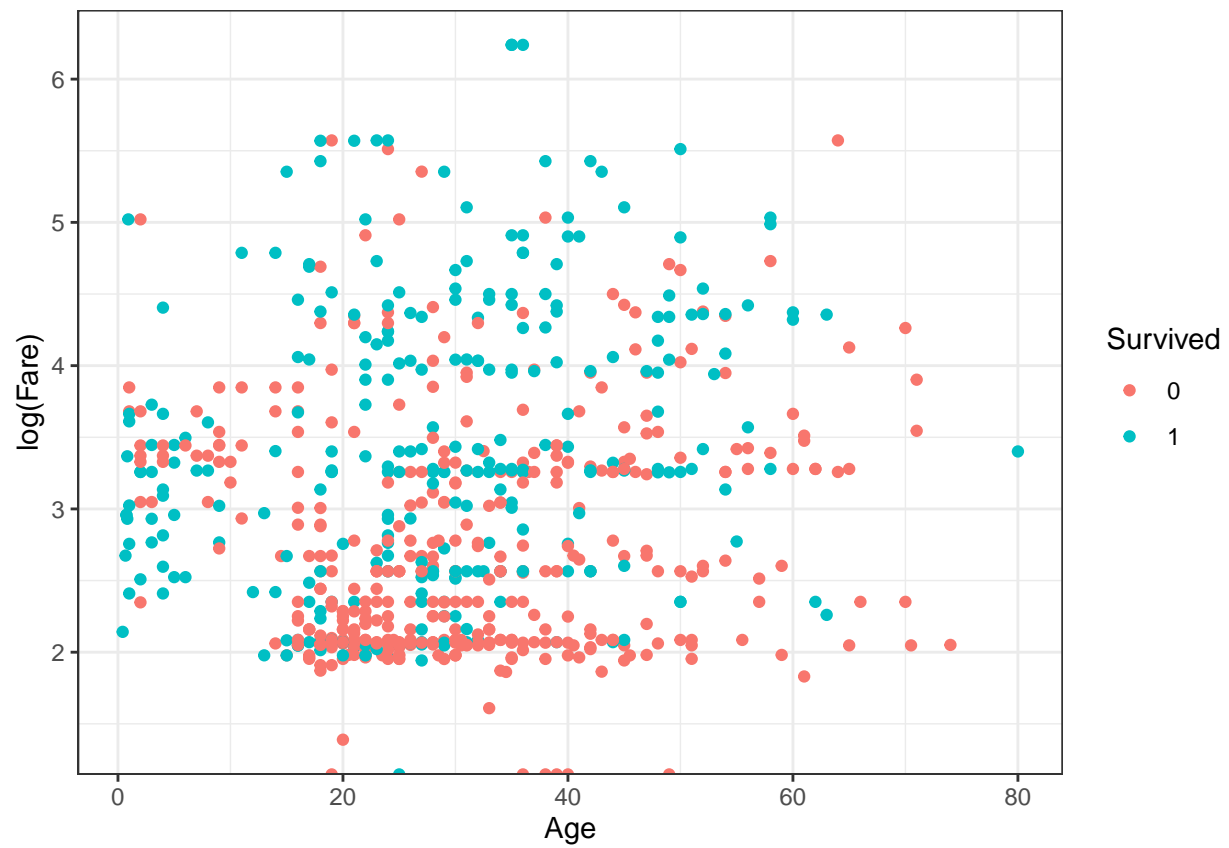
- Most of the survivors were female.
- Most of the male passengers died.
- Female passengers have higher chances of survival.

EDA: Passengers with first class ticket have higher survival rates



- Most survivors were from the first class.
- Most of the deaths were from the 3rd class.
- 1st class passengers had the highest survival rates.
- 3rd class passengers had the highest death rates.

EDA: AGE, FARE BY SURVIVAL



EDA: Check for outliers

