# HARVARD EXTENSION SCHOOL

## EXT CSCI E-106 Model Data Class Group Project Template

Ana Duchini 91319020    Douglas Malfacini    Daniella Olivas    John Pai 21745960
Joseph Ross 41732512    Kannan Sowminarayanan    Dorothy Vo 41732482    Thomas Yang
Caroline Zhuo

09 May 2025

### Abstract

The sinking of the Titanic remains one of the most tragic maritime disasters in history. Yet the heavy documentation and scrutiny of its end has offered a wealth of data to understand the tragedy of the past and build predictive tools to avoid such tragedy in the future. This project explores whether individual survival on the Titanic can be predicted using statistical and machine learning models. By analyzing passenger-level data including class, age, sex, fare, and other variables, we aim to understand the most influential predictors of survival rate and assess the effectiveness of different classification methods in estimating survival probabilities.

After data cleaning, imputation, and transformation, multiple models were trained including logistic regression, decision trees, random forests, support vector machines (SVM), Poisson regression, and neural networks. Each model was evaluated using confusion matrices, precision, recall, F1 score. Among all models, the Desicion Tree emerged as the champion model with the highest F1 score, strong generalizability, and minimal overfitting risk, making it the most robust approach for predicting survival outcomes in this dataset.

## Contents

**Classify whether a passenger on board the maiden voyage of the RMS Titanic in 1912 survived given their age, sex and class. Sample-Data-Titanic-Survival.csv to be used in the Final Project**

| Variable | Description |
| --- | --- |
| pclass | **Passanger Class, could be 1st, 2nd or 3rd** |
| survived | *Survival Status: 0=No, 1=Yes* |
| name | *Name of the Passanger* |
| Sex | *Sex* |
| sibsp | *Number of Siblings or Spouses aboard* |
| parch | *Number of Parents or Children aboard* |
| ticket | *Ticket Number* |
| fare | *Passenger Fare* |
| cabin | *Cabin number, "C85" would mean the cabin is on deck C and is numbered 85.* |
| embarked | *Port of Embarkation: C=Cherburg, S=Southampton, Q=Queenstown* |
| boat | *Lifeboat ID, if passenger survived* |
| body | *Body number (if passenger did not survive and body was recovered* |
| home.dest | *The intended home destination of the passenger* |

# Instructions:

0. Join a team with your fellow students with appropriate size (Up to Nine Students total) If you have not group by the end of the week of April 11 you may present the project by yourself or I will randomly assign other stranded student to your group. I will let know the final groups in April 11.
1. Load and Review the dataset named "Titanic_Survival_Data.csv"
2. Create the train data set which contains 70% of the data and use set.seed (15). The remaining 30% will be your test data set.
3. Investigate the data and combine the level of categorical variables if needed and drop variables as needed. For example, you can drop id, Latitude, Longitude, etc.
4. Build appropriate model to predict the probability of survival.
5. Create scatter plots and a correlation matrix for the train data set. Interpret the possible relationship between the response.
6. Build the best models by using the appropriate selection method. Compare the performance of the best logistic linear models.
7. Make sure that model assumption(s) are checked for the final model. Apply remedy measures (transformation, etc.) that helps satisfy the assumptions.
8. Investigate unequal variances and multicollinearity.
9. Build an alternative to your model based on one of the following approaches as applicable to predict the probability of survival: logistic regression, classification Tree, NN, or SVM. Check the applicable model assumptions. Explore using a negative binomial regression and a Poisson regression.
10. Use the test data set to assess the model performances from above.
11. Based on the performances on both train and test data sets, determine your primary (champion) model and the other model which would be your benchmark model.
12. Create a model development document that describes the model following this template, input the name of the authors, Harvard IDs, the name of the Group, all of your code and calculations, etc..:

## Due Date: May 12 2025 1159 pm hours EST

**Notes No typographical errors, grammar mistakes, or misspelled words, use English language All tables need to be numbered and describe their content in the body of the document All figures/graphs need to be numbered and describe their content All results must be accurate and clearly explained for a casual reviewer to fully understand their purpose and impact Submit both the RMD markdown file and PDF with the sections with appropriate explanations. A more formal document in Word can be used in place of the pdf file but must include all appropriate explanations.**

Executive Summary

This section will describe the model usage, your conclusions and any regulatory and internal requirements. In a real world scneario, this section is for senior management who do not need to know the details. They need to know high level (the purpose of the model, limitations of the model and any issues).

This project focuses on developing a predictive model to estimate the likelihood of survival for individual Titanic passengers. Using a cleaned dataset of 1,309 entries with various demographic and travel-related features, the analysis involved data cleaning, exploratory analysis, encoding of categorical variables, and imputation of missing values. We trained several classification models, including logistic regression, decision trees, random forest, SVM, and neural networks. The models were assessed on both training and test datasets using accuracy, precision, recall, and F1 score.

# I. Introduction (5 points)

*This section needs to introduce the reader to the problem to be resolved, the purpose, and the scope of the statistical testing applied. What you are doing with your prediction? What is the purpose of the model? What methods were trained on the data, how large is the test sample, and how did you build the model?*

Comment:

This study aims to predict the survival rate of individual passengers on the Titanic using supervised learning techniques. The primary objective is to develop and evaluate classification models capable of accurately estimating passengers survival rate, based on a range of demographic and travel-related features. These include age, sex, passenger class, fare, port of embarkation, and family relations aboard the ship (number of siblings/spouses and parents/children). The model is a binary classification model, where the target variable indicates whether a passenger survived (1) or not (0).

The data set initially contained 1,310 records and 14 variables, but after removing a blank final row, 1,309 usable observations remained. Significant data cleaning and preprocessing were required, including the creation of dummy variables, imputation of missing values (notably in age and cabin data), and transformation of continuous variables for better model performance.

To build and assess the predictive models, the cleaned data set was split into a training set (70%) and a test set (30%), using set.seed(1023) to ensure reproducibility. The primary modeling technique used was logistic regression, due to its interpretability and solid performance on binary outcomes. Additionally, several challenger models were trained and evaluated, including decision trees, random forests, support vector machines (SVM), neural networks, Poisson regression, and negative binomial regression. Each model's performance was assessed using accuracy, recall, precision, and F1 score on the test data set. This comprehensive modeling approach allowed us to compare a range of techniques and select the best-performing model for predicting survival.

## II. Description of the data and quality (15 points)

*Here you need to review your data, the statistical test applied to understand the predictors and the response and how are they correlated. Extensive graph analysis is recommended. Is the data continuous, or categorical, do any transformation needed? Do you need dummies?*

Comment:

We take a first look at the structure of the data and load all relevant libraries:

```r
library(GGally)
```

```
## Loading required package: ggplot2
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg   ggplot2
```

```r
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
library(mice)
```

```
##
## Attaching package: 'mice'
```

```
## The following object is masked from 'package:stats':
##
##     filter
```

```
## The following objects are masked from 'package:base':
##
##     cbind, rbind
```

```r
library(dplyr)
library(tidyr)
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ------------------------ tidyverse 2.0.0 --
## v forcats   1.0.0      v stringr   1.5.1
## v lubridate 1.9.4      v tibble    3.2.1
## v purrr     1.0.4
```

```
## -- Conflicts -------------------------------------------- tidyverse_conflicts() --
## x mice::filter() masks dplyr::filter(), stats::filter()
## x dplyr::lag()   masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```r
library(caret)
```

```
## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
##     lift
```

```r
library(e1071)
library(nnet)
library(MASS)
```

```
##
## Attaching package: 'MASS'
##
## The following object is masked from 'package:dplyr':
##
##     select
```

```r
library(rpart)
library(rpart.plot)
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
##
## Attaching package: 'pROC'
##
## The following objects are masked from 'package:stats':
##
##     cov, smooth, var
```

```r
library(reshape2)
```

```
##
## Attaching package: 'reshape2'
##
## The following object is masked from 'package:tidyr':
##
##     smiths
```

```r
library(randomForest)
```

```
## randomForest 4.7-1.2
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
##
## The following object is masked from 'package:dplyr':
##
##     combine
##
## The following object is masked from 'package:ggplot2':
##
##     margin
```

```r
# It is reading 1310 observations, however, the last row is just blanks
str(Titanic.Data)
```

```
## 'data.frame':    1310 obs. of  14 variables:
##  $ pclass   : int  1 1 1 1 1 1 1 1 1 1 ...
##  $ survived : int  1 1 0 0 0 1 1 0 1 0 ...
##  $ name     : chr  "Allen, Miss. Elisabeth Walton" "Allison, Master. Hudson Trevor" "Allison, Miss. Helen L
##  $ sex      : chr  "female" "male" "female" "male" ...
##  $ age      : num  29 0.917 2 30 25 ...
##  $ sibsp    : int  0 1 1 1 1 0 1 0 2 0 ...
##  $ parch    : int  0 2 2 2 2 0 0 0 0 0 ...
##  $ ticket   : chr  "24160" "113781" "113781" "113781" ...
##  $ fare     : num  211 152 152 152 152 ...
##  $ cabin    : chr  "B5" "C22 C26" "C22 C26" "C22 C26" ...
##  $ embarked : chr  "S" "S" "S" "S" ...
##  $ boat     : chr  "2" "11" NA NA ...
```

```
##  $ body     : int  NA NA NA 135 NA NA NA NA NA 22 ...
##  $ home.dest: chr  "St Louis, MO" "Montreal, PQ / Chesterville, ON" "Montreal, PQ / Chesterville, ON" "Mont
```

Comment:

The data cleaning was performed as indicated in the lines that follow. To summarize, after initial cleaning, we transform each predictor as follows:

*pclass*: create 2 dummies for values 1, 2, 3 - *pclass.1, pclass.2 name*: drop the column, there is not much relevant information that can be derived from this *sex*: create a single dummy for values male, female - *sex.male age*: use as continuous variable. (there are 263 missing values, these were imputted using the predictive mean matching method from the MICE library. It could be left as-is for models that can handle missing values, such as tree-based models like xgboost, lightgbm) *sibsp*: adding a factor variable *parch*:: adding a factor variable *ticket*: drop column (no relevant use) *fare*: imputed 1 missing value with median *cabin*: There are 1014 missing values (77% of the rows). We first create a *deck* column (based on the first character of the cabin column, corresponding to the ship deck level) as factor, and then impute the missing values using a sample of this column. Notes: any information derived from this should correlate in some level with the class dummies; and due to the imputations, this will be a majority of synthetic data. *embarked*: create 2 dummies by port C, Q, S. Two missing values were imputted with the mode. *boat*, *body*: Drop. As they are dependent on the survived column we are trying to predict. This can be seen by looking at the data descriptions i.e. (Lifeboat ID, if passenger survived) should only have values if survived is true. *home.dest*: drop the column, there is not much relevant information that can be derived from this

In a correlation matrix, the strongest predictors for the response variable *survived* are class (negative correlation), closely followed by fare. However, there is also a strong correlation between these two, so we are likely to have to choose one in our initial model.

Finally, we scale and center the continuous variables *age* and *fare*.

```r
tail(Titanic.Data, n=1)
```

```
##      pclass survived name  sex age sibsp parch ticket fare cabin embarked boat
## 1310     NA       NA <NA> <NA>  NA    NA    NA   <NA>   NA  <NA>     <NA> <NA>
##      body home.dest
## 1310   NA      <NA>
```

```r
# dropping last row -- there are 1309 observations.
Titanic.Clean.Data = head(Titanic.Data,-1)

# check for columns that have blanks, such as NAs, blanks, whitespaces or null values
print('Number of total NA or blank values per column')
```

```
## [1] "Number of total NA or blank values per column"
```

```r
colSums(is.na(Titanic.Clean.Data))
```

```
##    pclass  survived      name       sex       age     sibsp     parch    ticket
##         0         0         0         0       263         0         0         0
##      fare     cabin  embarked      boat      body home.dest
##         1      1014         2       823      1188       564
```

```r
# closer look at the structure of data
glimpse(Titanic.Clean.Data)
```

```
## Rows: 1,309
## Columns: 14
## $ pclass   <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ survived <int> 1, 1, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, ~
## $ name     <chr> "Allen, Miss. Elisabeth Walton", "Allison, Master. Hudson Tr~
## $ sex      <chr> "female", "male", "female", "male", "female", "male", "femal~
## $ age      <dbl> 29.0000, 0.9167, 2.0000, 30.0000, 25.0000, 48.0000, 63.0000,~
## $ sibsp    <int> 0, 1, 1, 1, 1, 0, 1, 0, 2, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ parch    <int> 0, 2, 2, 2, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, ~
## $ ticket   <chr> "24160", "113781", "113781", "113781", "113781", "19952", "1~
## $ fare     <dbl> 211.3375, 151.5500, 151.5500, 151.5500, 151.5500, 26.5500, 7~
## $ cabin    <chr> "B5", "C22 C26", "C22 C26", "C22 C26", "C22 C26", "E12", "D7~
## $ embarked <chr> "S", "S", "S", "S", "S", "S", "S", "S", "S", "C", "C", "C", ~
## $ boat     <chr> "2", "11", NA, NA, NA, "3", "10", NA, "D", NA, NA, "4", "9",~
```
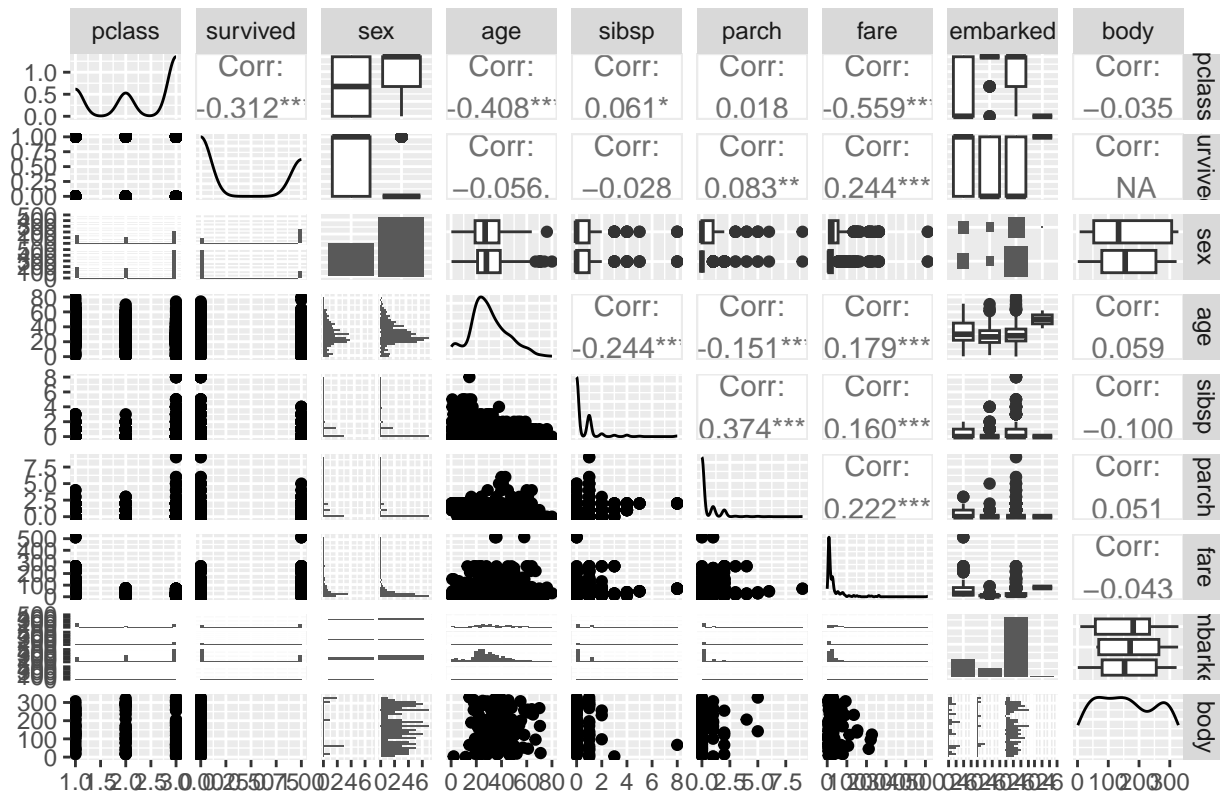
```
## $ body      <int> NA, NA, NA, 135, NA, NA, NA, NA, NA, 22, 124, NA, NA, NA, NA~
## $ home.dest <chr> "St Louis, MO", "Montreal, PQ / Chesterville, ON", "Montreal~
```

```r
options(warn = -1)
# we take a first look at the correlation between the variables we're planning to use
suppressWarnings(ggpairs(Titanic.Clean.Data |> dplyr::select(-name, -ticket, -cabin, -boat, -home.dest), title
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Figure 1: Correlation matrix



```r
# pclass
Titanic.Clean.Data$pclass.1=I(Titanic.Clean.Data$pclass=='1')*1
Titanic.Clean.Data$pclass.2=I(Titanic.Clean.Data$pclass=='2')*1
sum(Titanic.Clean.Data$pclass.1 == "1") # check
```

```
## [1] 323
```

```r
sum(Titanic.Clean.Data$pclass.2 == "1")
```

```
## [1] 277
```

```r
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "pclass")] # drop original
```

```r
# name
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "name")] # drop original
```

```r
# sex
Titanic.Clean.Data$sex.male=I(Titanic.Clean.Data$sex=="male")*1
sum(Titanic.Clean.Data$sex.male == "1") # check
```

```
## [1] 843
```

```r
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "sex")] # drop original
```

```r
cat('Unique values of sibsp:', unique(Titanic.Clean.Data$sibsp), '\nUnique values of parch:', unique(Titanic.C
```

```
## Unique values of sibsp: 0 1 2 3 4 5 8
## Unique values of parch: 0 2 1 4 3 5 6 9
```

```r
Titanic.Clean.Data$sibsp<-as.factor(Titanic.Clean.Data$sibsp)
Titanic.Clean.Data$parch<-as.factor(Titanic.Clean.Data$parch)

contrasts(Titanic.Clean.Data$sibsp)
```

```
##   1 2 3 4 5 8
## 0 0 0 0 0 0 0
## 1 1 0 0 0 0 0
## 2 0 1 0 0 0 0
## 3 0 0 1 0 0 0
## 4 0 0 0 1 0 0
## 5 0 0 0 0 1 0
## 8 0 0 0 0 0 1
```

```r
# ticket
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "ticket")] # drop entirely
```

```r
# fare - impute 1 missing value with median
Titanic.Clean.Data$fare[is.na(Titanic.Clean.Data$fare)] <- median(Titanic.Clean.Data$fare, na.rm=TRUE)
```

```r
# test cases
Titanic.Clean.Data[169,]
```

```
##     survived age sibsp parch fare cabin embarked boat body home.dest pclass.1
## 169        1  38     0     0   80   B28     <NA>    6   NA      <NA>        1
##     pclass.2 sex.male
## 169        0        0
```

```r
Titanic.Clean.Data[285,]
```

```
##     survived age sibsp parch fare cabin embarked boat body    home.dest
## 285        1  62     0     0   80   B28     <NA>    6   NA Cincinatti, OH
##     pclass.1 pclass.2 sex.male
## 285        1        0        0
```

```r
# embarked - impute 2 missing values with the mode
mode = sapply(names(sort(-table(Titanic.Clean.Data$embarked)))[1], as.character)
Titanic.Clean.Data$embarked[is.na(Titanic.Clean.Data$embarked)] <- mode
# testing
Titanic.Clean.Data[169,]
```

```
##     survived age sibsp parch fare cabin embarked boat body home.dest pclass.1
## 169        1  38     0     0   80   B28        S    6   NA      <NA>        1
##     pclass.2 sex.male
## 169        0        0
```

```r
Titanic.Clean.Data[285,]
```

```
##     survived age sibsp parch fare cabin embarked boat body    home.dest
## 285        1  62     0     0   80   B28        S    6   NA Cincinatti, OH
```

```
##     pclass.1 pclass.2 sex.male
## 285        1        0        0
```
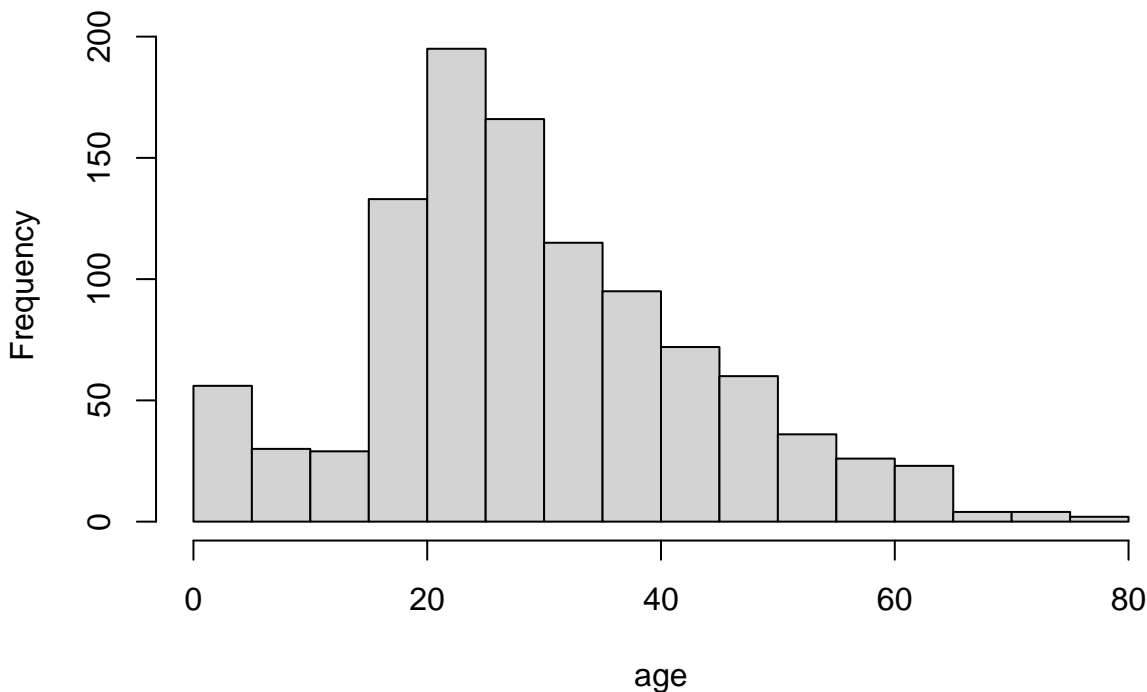
```r
# create dummies
Titanic.Clean.Data$embarked.Q=I(Titanic.Clean.Data$embarked=='Q')*1
Titanic.Clean.Data$embarked.S=I(Titanic.Clean.Data$embarked=='S')*1
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "embarked")] # drop original

# Dropping body and boat columns
Titanic.Clean.Data = Titanic.Clean.Data |> dplyr::select(-body, -boat)

# home.dest
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "home.dest")] # drop for now...

# age - add column with imputted median values
hist(Titanic.Clean.Data$age, main='Figure 2: Distribution of Age in original dataset',xlab='age') # there is a
```

**Figure 2: Distribution of Age in original dataset**



```r
# more info about imputting methods here: https://www.r-bloggers.com/2023/01/imputation-in-r-top-3-ways-for-im

# we want to impute keeping the same distribution
# using MICE - Predictive mean matching method

# use a dataframe that contains only number columns, and remove any that have dependencies with age
Titanic.numbers.cols <- Titanic.Clean.Data %>% dplyr::select(survived, pclass.1, pclass.2, sibsp, fare, sex.ma
Titanic.Clean.Data$age.imputted = complete(mice(Titanic.numbers.cols, method = "pmm", maxit = 5, m = 5, seed=1
```
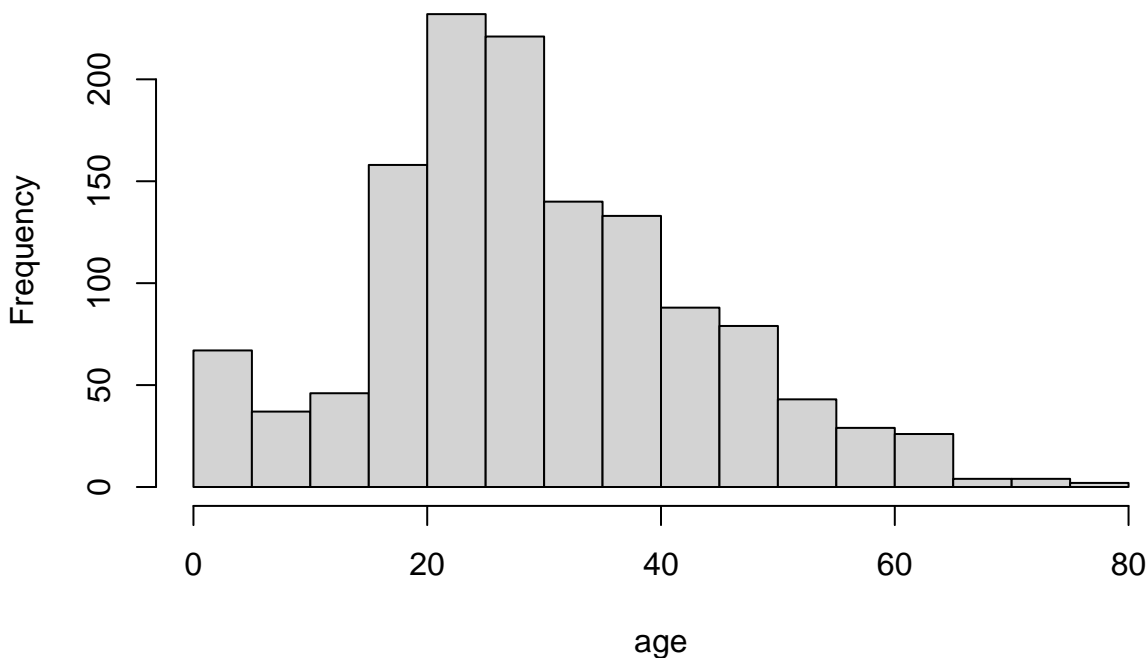
```
## 
##  iter imp variable
##   1   1  age
##   1   2  age
##   1   3  age
##   1   4  age
##   1   5  age
##   2   1  age
##   2   2  age
##   2   3  age
##   2   4  age
##   2   5  age
```

```
## 3    1    age
## 3    2    age
## 3    3    age
## 3    4    age
## 3    5    age
## 4    1    age
## 4    2    age
## 4    3    age
## 4    4    age
## 4    5    age
## 5    1    age
## 5    2    age
## 5    3    age
## 5    4    age
## 5    5    age
```
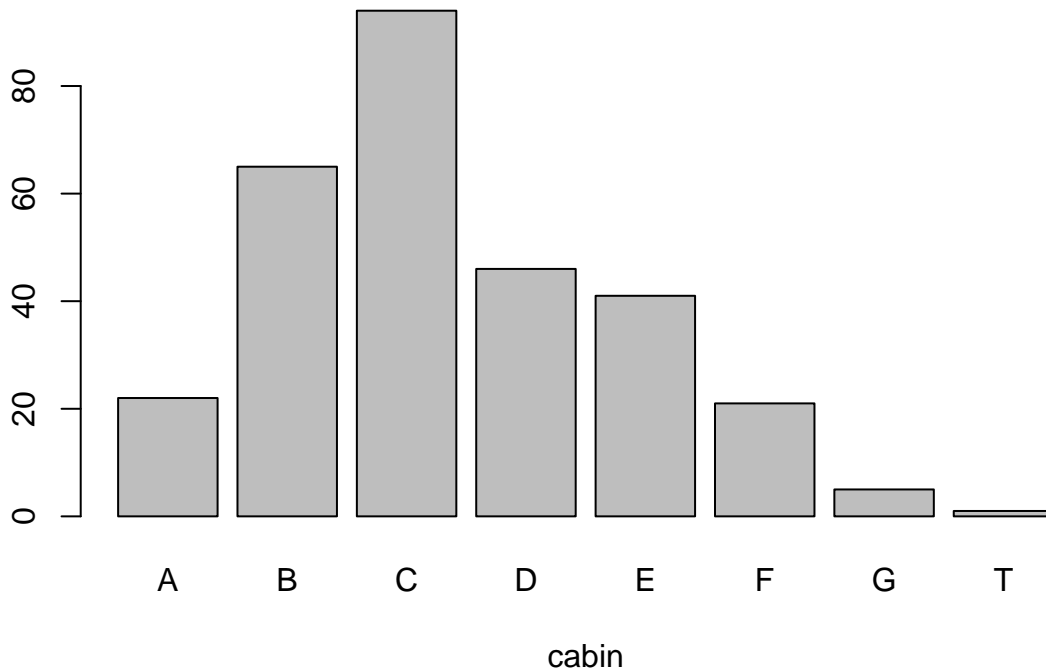
```r
#imp <- mice(Titanic.numbers.cols, m = 5)
#imp$loggedEvents # with this, we are able to detect and remove age dependency with parch column
hist(Titanic.Clean.Data$age.imputted, main='Figure 3: Distribution of Age after imputation',xlab='age')
```

## Figure 3: Distribution of Age after imputation



```r
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "age")] # drop original
```

```r
# cabin - 1014 rows have no data on this column - this is 77% of the rows
# extracting the first character to a new column 'deck'
Titanic.Clean.Data$deck <- as.factor(substr(Titanic.Clean.Data$cabin,1,1))
plot(as.factor(Titanic.Clean.Data$deck),main='Figure 4: Distribution of cabin in original dataset',xlab='cabin
```

## Figure 4: Distribution of cabin in original dataset



cabin

```r
# Set up imputation method: sample of the 'deck' column
deck.method <- make.method(Titanic.Clean.Data)
deck.method["deck"] <- "sample"
Titanic.Clean.Data$deck.imputted = complete(mice(Titanic.Clean.Data, method = deck.method, m = 1, seed = 123))
```

```
##
##  iter imp variable
##   1   1  deck
##   2   1  deck
##   3   1  deck
##   4   1  deck
##   5   1  deck
```

```r
plot(as.factor(Titanic.Clean.Data$deck.imputted),main='Figure 5: Distribution of deck after imputation',xlab='
```

**Figure 5: Distribution of deck after imputation**



cabin

```
# drop original and temporary variables
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "cabin")]
Titanic.Clean.Data <- Titanic.Clean.Data[, -which(names(Titanic.Clean.Data) == "deck")]
```

```
str(Titanic.Clean.Data)
```

```
## 'data.frame':    1309 obs. of  11 variables:
##  $ survived     : int  1 1 0 0 0 1 1 0 1 0 ...
##  $ sibsp        : Factor w/ 7 levels "0","1","2","3",..: 1 2 2 2 2 1 2 1 3 1 ...
##  $ parch        : Factor w/ 8 levels "0","1","2","3",..: 1 3 3 3 3 1 1 1 1 1 ...
##  $ fare         : num  211 152 152 152 152 ...
##  $ pclass.1     : 'AsIs' num  1 1 1 1 1 1 1 1 1 1 ...
##  $ pclass.2     : 'AsIs' num  0 0 0 0 0 0 0 0 0 0 ...
##  $ sex.male     : 'AsIs' num  0 1 0 1 0 1 0 1 0 1 ...
##  $ embarked.Q   : 'AsIs' num  0 0 0 0 0 0 0 0 0 0 ...
##  $ embarked.S   : 'AsIs' num  1 1 1 1 1 1 1 1 1 0 ...
##  $ age.imputted : num  29 0.917 2 30 25 ...
##  $ deck.imputted: Factor w/ 8 levels "A","B","C","D",..: 2 3 3 3 3 5 4 1 3 4 ...
```

```
head(Titanic.Clean.Data)
```

```
##   survived sibsp parch     fare pclass.1 pclass.2 sex.male embarked.Q
## 1        1     0     0 211.3375        1        0        0          0
## 2        1     1     2 151.5500        1        0        1          0
## 3        0     1     2 151.5500        1        0        0          0
## 4        0     1     2 151.5500        1        0        1          0
## 5        0     1     2 151.5500        1        0        0          0
## 6        1     0     0  26.5500        1        0        1          0
##   embarked.S age.imputted deck.imputted
## 1          1      29.0000             B
## 2          1       0.9167             C
## 3          1       2.0000             C
## 4          1      30.0000             C
## 5          1      25.0000             C
## 6          1      48.0000             E
```
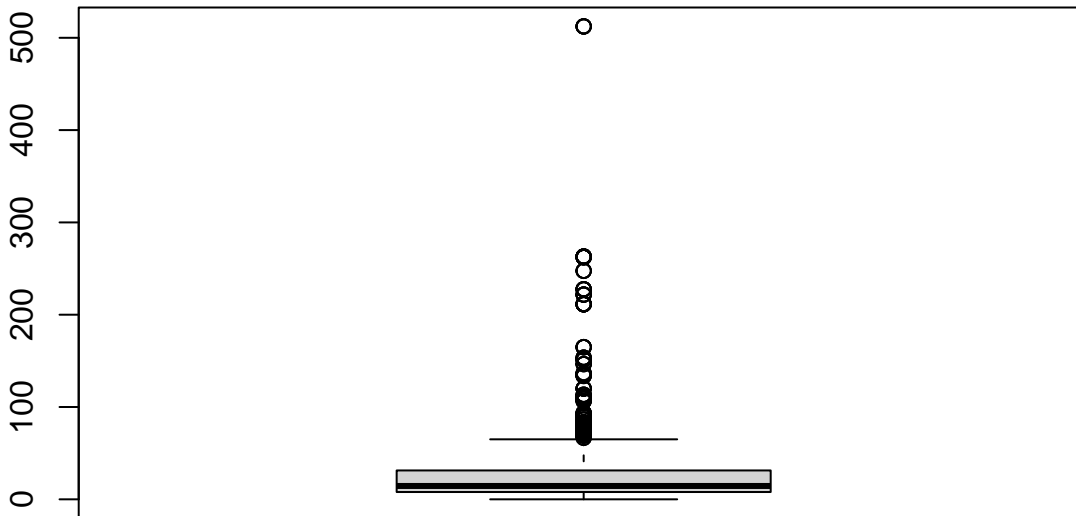
```
# checking for any leftover NAs
colSums(is.na(Titanic.Clean.Data))
```

```
##       survived        sibsp        parch         fare     pclass.1
##              0            0            0            0            0
##       pclass.2     sex.male   embarked.Q   embarked.S age.imputted
##              0            0            0            0            0
## deck.imputted
##              0
```
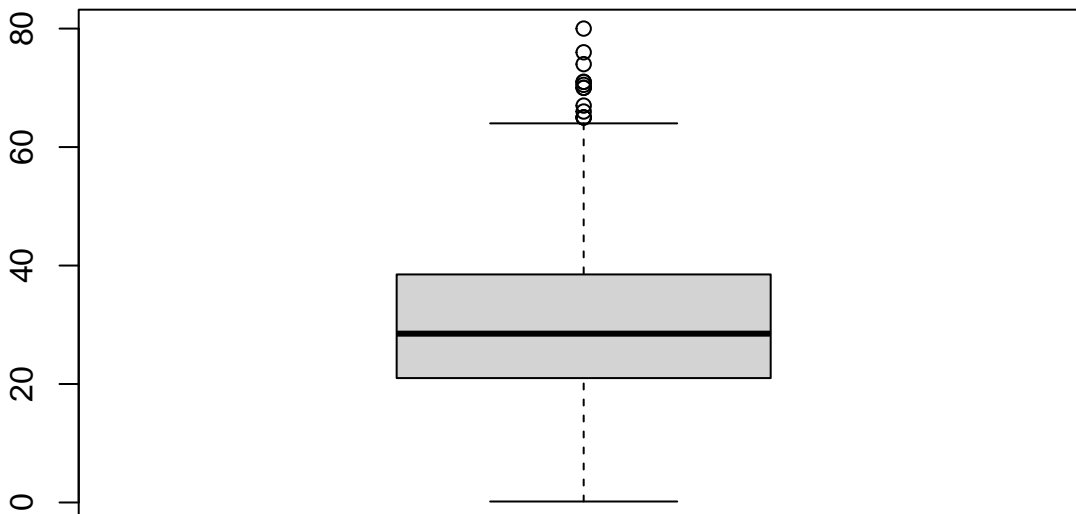
```r
boxplot(Titanic.Clean.Data$fare, main='Figure 6: Boxplot of fare predictor')
```

**Figure 6: Boxplot of fare predictor**



```r
boxplot(Titanic.Clean.Data$age, main='Figure 7: Boxplot of age predictor')
```
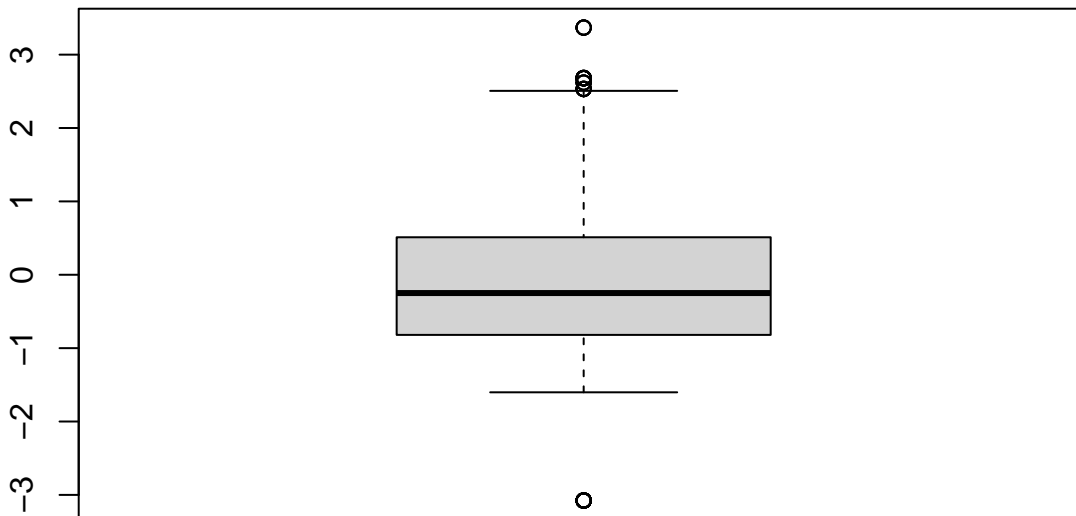
**Figure 7: Boxplot of age predictor**



```r
#centering and scaling continuous variables
Titanic.Clean.Data = Titanic.Clean.Data |>
  mutate(
    l_fare = scale(log(fare+1)),
    age.imputted = scale(age.imputted)
    ) |>
  dplyr::select(-fare)

boxplot(Titanic.Clean.Data$l_fare, main='Figure 8: Boxplot of scaled and centered fare')
```

## Figure 8: Boxplot of scaled and centered fare


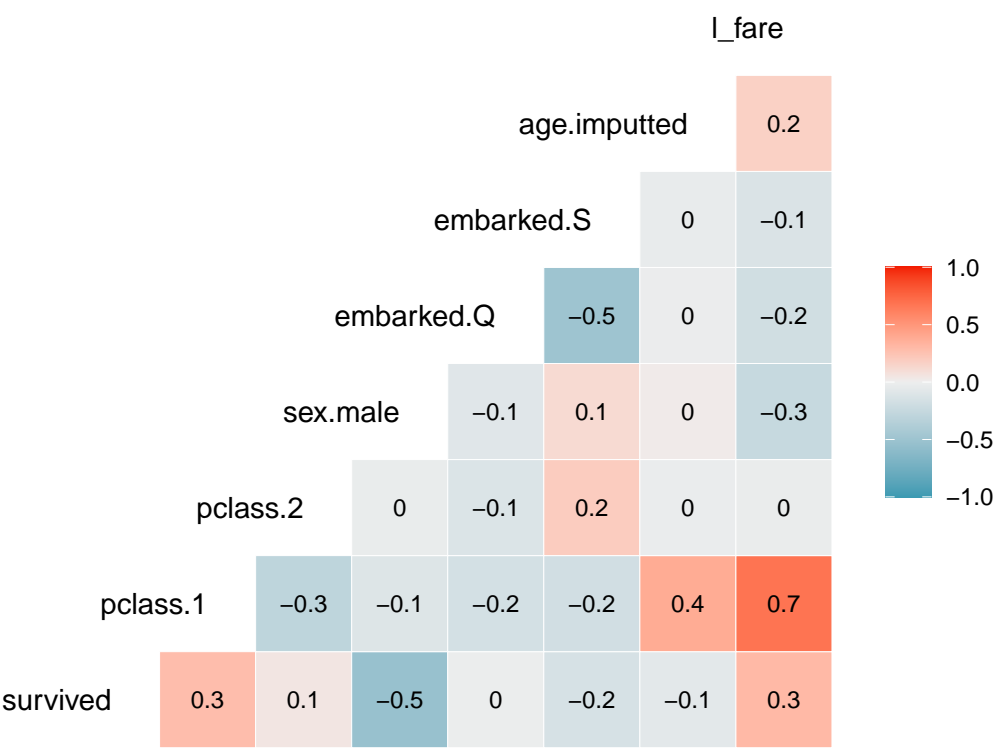
```r
summary(Titanic.Clean.Data)
```

```
##     survived       sibsp       parch        pclass.1         pclass.2
## Min.   :0.000   0:891   0      :1002   Min.   :0.0000   Min.   :0.0000
## 1st Qu.:0.000   1:319   1      : 170   1st Qu.:0.0000   1st Qu.:0.0000
## Median :0.000   2: 42   2      : 113   Median :0.0000   Median :0.0000
## Mean   :0.382   3: 20   3      :   8   Mean   :0.2468   Mean   :0.2116
## 3rd Qu.:1.000   4: 22   4      :   6   3rd Qu.:0.0000   3rd Qu.:0.0000
## Max.   :1.000   5:  6   5      :   6   Max.   :1.0000   Max.   :1.0000
##                 8:  9   (Other):   4
##     sex.male        embarked.Q         embarked.S       age.imputted.V1
## Min.   :0.000   Min.   :0.00000   Min.   :0.0000   Min.   :-2.103655
## 1st Qu.:0.000   1st Qu.:0.00000   1st Qu.:0.0000   1st Qu.:-0.627456
## Median :1.000   Median :0.00000   Median :1.0000   Median :-0.096024
## Mean   :0.644   Mean   :0.09396   Mean   :0.6998   Mean   : 0.000000
## 3rd Qu.:1.000   3rd Qu.:0.00000   3rd Qu.:1.0000   3rd Qu.: 0.612552
## Max.   :1.000   Max.   :1.00000   Max.   :1.0000   Max.   : 3.553144
##
## deck.imputted       l_fare.V1
## C      :416   Min.   :-3.076692
## B      :277   1st Qu.:-0.819611
## D      :212   Median :-0.249241
## E      :171   Mean   : 0.000000
## F      :105   3rd Qu.: 0.511263
## A      :104   Max.   : 3.368398
## (Other): 24
```

```r
ggcorr(Titanic.Clean.Data, label = TRUE, label_size = 3, hjust = 1, layout.exp = 2) +
  ggtitle("Figure 9: Correlation Matrix for the Final Cleaned Titanic Data")
```

Figure 9: Correlation Matrix for the Final Cleaned Titanic Data

# III. Model Development Process (15 points)

*Build an appropriate model to predict probability of survival. And of course, create the train data set which contains 70% of the data and use set.seed (1023). The remaining 30% will be your test data set. Investigate the data and combine the level of categorical variables if needed and drop variables. For example, you can passenger name, cabin, etc..*

Comment:
The Lasso Model is better suited for prediction in main model than the binomial, logistic or probabilistic binomial fit, but all models had a reasonably good fit. The influential points were fewer, as shown in the Cook's distance graph of binomial fit. Binomial, Probabilistic Binomial and Logistic regression were slightly affected by multicollinearity of predictor variables Age.Imputted, Sex.Male, and pclass.1. Lasso regression can handle multicollinearity but might be affected by multicollinearity as its selection of variables can be affected due to it.

The Probabilistic Binomial, Binomial and Logistic regression have provided the same result where as the result of Lasso is different from these three models and is slightly better in prediction than the three models.

Ran the Probabilistic Binomial, Binomial and Logistic regression with all variables and did a manual step wise backward removal of variables with high p-value till all variables became significant. Attention to be brought that pclass.2 was also significant but when included into the fit the hoslem test found that all the three model were not of good fit. But multicollinearity decreased to to 5.

Currently with the three models with sex.male, pclass.1 and age.imputed as predictive variables we see the models are of good fit as per hoslem test with p-value greater than .05. The multicollinearity is also moderate as it is equal to or less than 8 and the multicollinearity values are provided in the output.

— Binomial/Logistic Regression for training data prediction — The Binomial and Logistic Regression provides a confusion matrix:

```
    Precision : 0.7278
       Recall : 0.6959
           F1 : 0.7115
```

Null deviance: 1231.8 on 915 degrees of freedom Residual deviance: 889.7 on 912 degrees of freedom AIC: 897.7

Number of Fisher Scoring iterations: 4

— Lasso Regression for training data prediction —

The variables selcted by lasso regression are sibsp8, parch1, pclass.1, pclass.2, sex.male, embarked.S, age.imputted

The Best Lambda value is 0.01107336.

The Lasso Regression provides a confusion matrix: Precision : 0.7629
Recall : 0.6877
F1 : 0.7233

The Logistic and Binomial regression prediction have better recall (that is, less false negatives than Lasso) but precision is better than for Lasso prediction, indicating fewer false positives. A F1 score of 0.72 in Lasso indicates a good balance between precision and recall versus the binomial and logistic regression with value 0.71.

The confusion matrix was created with threshold greater than 0.5 indicating survived and less than 0.5 indicating not survived.

```
# testing cleaned data:
set.seed(1023)
library(caret)

split_data = createDataPartition(y = Titanic.Clean.Data$survived, p=0.6995, list=FALSE )
train_data = Titanic.Clean.Data[split_data,]
test_data = Titanic.Clean.Data[-split_data,]
nrow(Titanic.Clean.Data)*.7
```

```
## [1] 916.3
```

```
nrow(train_data)
```

```
## [1] 916
```

```r
# Manual stepwise backwards applied by Kannan for these models
# performance looks similar across all these models
# testing a model
lmod <- glm(as.factor(survived) ~ . -deck.imputted  - l_fare - embarked.Q
            - sibsp - parch,
            family = binomial, train_data)
summary(lmod)
```

```
##
## Call:
## glm(formula = as.factor(survived) ~ . - deck.imputted - l_fare -
##     embarked.Q - sibsp - parch, family = binomial, data = train_data)
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.55812    0.19391   2.878  0.00400 **
## pclass.1       2.35151    0.24455   9.616  < 2e-16 ***
## pclass.2       1.29471    0.22382   5.785 7.27e-09 ***
## sex.male      -2.55859    0.18238 -14.029  < 2e-16 ***
## embarked.S    -0.51233    0.19120  -2.680  0.00737 **
## age.imputted  -0.53224    0.09649  -5.516 3.47e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1231.81  on 915  degrees of freedom
## Residual deviance:  839.46  on 910  degrees of freedom
## AIC: 851.46
##
## Number of Fisher Scoring iterations: 5
```

```r
probitmod<-glm(as.factor(survived) ~ . -deck.imputted  - l_fare - embarked.Q
            - sibsp - parch,
            family=binomial(link=probit),train_data)
summary(probitmod)
```

```
##
## Call:
## glm(formula = as.factor(survived) ~ . - deck.imputted - l_fare -
##     embarked.Q - sibsp - parch, family = binomial(link = probit),
##     data = train_data)
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.35866    0.11436   3.136  0.00171 **
## pclass.1       1.35767    0.13713   9.901  < 2e-16 ***
## pclass.2       0.72220    0.12671   5.700 1.20e-08 ***
## sex.male      -1.51882    0.10222 -14.858  < 2e-16 ***
## embarked.S    -0.28489    0.11007  -2.588  0.00965 **
## age.imputted  -0.30923    0.05516  -5.606 2.07e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1231.81  on 915  degrees of freedom
## Residual deviance:  839.88  on 910  degrees of freedom
## AIC: 851.88
##
```

```
## Number of Fisher Scoring iterations: 5
logitmod<-glm(as.factor(survived) ~ . -deck.imputted  - l_fare - embarked.Q
              - sibsp - parch,
              family=binomial(link=logit),train_data)
summary(logitmod)

##
## Call:
## glm(formula = as.factor(survived) ~ . - deck.imputted - l_fare -
##     embarked.Q - sibsp - parch, family = binomial(link = logit),
##     data = train_data)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.55812    0.19391   2.878  0.00400 **
## pclass.1      2.35151    0.24455   9.616  < 2e-16 ***
## pclass.2      1.29471    0.22382   5.785 7.27e-09 ***
## sex.male     -2.55859    0.18238 -14.029  < 2e-16 ***
## embarked.S   -0.51233    0.19120  -2.680  0.00737 **
## age.imputted -0.53224    0.09649  -5.516 3.47e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1231.81  on 915  degrees of freedom
## Residual deviance:  839.46  on 910  degrees of freedom
## AIC: 851.46
##
## Number of Fisher Scoring iterations: 5
anova(logitmod,probitmod,test = "Chisq")

## Analysis of Deviance Table
##
## Model 1: as.factor(survived) ~ (sibsp + parch + pclass.1 + pclass.2 +
##     sex.male + embarked.Q + embarked.S + age.imputted + deck.imputted +
##     l_fare) - deck.imputted - l_fare - embarked.Q - sibsp - parch
## Model 2: as.factor(survived) ~ (sibsp + parch + pclass.1 + pclass.2 +
##     sex.male + embarked.Q + embarked.S + age.imputted + deck.imputted +
##     l_fare) - deck.imputted - l_fare - embarked.Q - sibsp - parch
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1       910     839.46
## 2       910     839.88  0 -0.41817
anova(logitmod,lmod,test = "Chisq")

## Analysis of Deviance Table
##
## Model 1: as.factor(survived) ~ (sibsp + parch + pclass.1 + pclass.2 +
##     sex.male + embarked.Q + embarked.S + age.imputted + deck.imputted +
##     l_fare) - deck.imputted - l_fare - embarked.Q - sibsp - parch
## Model 2: as.factor(survived) ~ (sibsp + parch + pclass.1 + pclass.2 +
##     sex.male + embarked.Q + embarked.S + age.imputted + deck.imputted +
##     l_fare) - deck.imputted - l_fare - embarked.Q - sibsp - parch
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1       910     839.46
## 2       910     839.46  0        0
anova(probitmod,lmod,test = "Chisq")

## Analysis of Deviance Table
```

```
##
## Model 1: as.factor(survived) ~ (sibsp + parch + pclass.1 + pclass.2 +
##     sex.male + embarked.Q + embarked.S + age.imputted + deck.imputted +
##     l_fare) - deck.imputted - l_fare - embarked.Q - sibsp - parch
## Model 2: as.factor(survived) ~ (sibsp + parch + pclass.1 + pclass.2 +
##     sex.male + embarked.Q + embarked.S + age.imputted + deck.imputted +
##     l_fare) - deck.imputted - l_fare - embarked.Q - sibsp - parch
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1       910     839.88
## 2       910     839.46  0  0.41817
```

```r
print("Checking for multicollinearity")
```

```
## [1] "Checking for multicollinearity"
```

```r
library(car)
```

```
## Loading required package: carData
```

```
##
## Attaching package: 'car'
```

```
## The following object is masked from 'package:purrr':
##
##     some
```

```
## The following object is masked from 'package:dplyr':
##
##     recode
```

```r
vif(logitmod)
```

```
##     pclass.1     pclass.2      sex.male   embarked.S age.imputted
##     1.584985     1.271620     1.085716     1.067148     1.311346
```

```r
cat("Checking for good of fit test for logistic regression")
```

```
## Checking for good of fit test for logistic regression
```

```r
library(ResourceSelection)
```

```
## ResourceSelection 0.3-6    2023-06-27
```

```r
hoslem.test(logitmod$y,fitted(logitmod),g=5)
```

```
##
##  Hosmer and Lemeshow goodness of fit (GOF) test
##
## data:  logitmod$y, fitted(logitmod)
## X-squared = 32.447, df = 3, p-value = 4.213e-07
```

```r
cat("Checking for good of fit test for Probablistic logistic regression")
```

```
## Checking for good of fit test for Probablistic logistic regression
```

```r
library(ResourceSelection)
hoslem.test(probitmod$y,fitted(probitmod),g=5)
```

```
##
##  Hosmer and Lemeshow goodness of fit (GOF) test
##
## data:  probitmod$y, fitted(probitmod)
## X-squared = 22.75, df = 3, p-value = 4.553e-05
```

```r
cat("Checking for good of fit test for Binomial regression")
```

```
## Checking for good of fit test for Binomial regression
```

```r
library(ResourceSelection)
hoslem.test(lmod$y,fitted(lmod),g=5)
```

```
##
##  Hosmer and Lemeshow goodness of fit (GOF) test
##
## data:  lmod$y, fitted(lmod)
## X-squared = 32.447, df = 3, p-value = 4.213e-07
```

```r
library(faraway)
```

```
##
## Attaching package: 'faraway'
```

```
## The following objects are masked from 'package:car':
##
##     logit, vif
```

```
## The following object is masked from 'package:rpart':
##
##     solder
```

```
## The following object is masked from 'package:lattice':
##
##     melanoma
```

```
## The following object is masked from 'package:mice':
##
##     mammalsleep
```

```
## The following object is masked from 'package:GGally':
##
##     happy
```

```r
library(statmod)
qqnorm(qresid(lmod)); qqline(qresid(lmod))
```

## Normal Q–Q Plot



```r
plot(cooks.distance(lmod), type='h')
```

```r
# Make predictions with binomial model
lmod_pred <- predict(lmod, newdata = train_data, type = "response")
probitmod_pred <- predict(probitmod, newdata = train_data,
                          type = "response")
logitmod_pred <- predict(probitmod, newdata = train_data,
                         type = "response")


lmod_pred_val <- ifelse(lmod_pred > 0.5,1,0)

probitmod_pred_val <- ifelse(probitmod_pred > 0.5,1,0)

logitmod_pred_val <- ifelse(logitmod_pred > 0.5,1,0)

# Evaluate binomial model using training data prediction

confusion_matrix_lmod <-
  confusionMatrix(as.factor(lmod_pred_val),as.factor(train_data$survived),
                  mode="prec_recall", positive = "1")
cat("Train data confusion Matrix for Binomial Regression \n")
```

```
## Train data confusion Matrix for Binomial Regression
```

```r
confusion_matrix_lmod
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 462 115
##          1  89 250
##
##                Accuracy : 0.7773
##                  95% CI : (0.7489, 0.8039)
##     No Information Rate : 0.6015
##     P-Value [Acc > NIR] : < 2e-16
##
##                   Kappa : 0.5298
##
##  Mcnemar's Test P-Value : 0.08006
```

```
##
##              Precision : 0.7375
##                 Recall : 0.6849
##                     F1 : 0.7102
##             Prevalence : 0.3985
##         Detection Rate : 0.2729
##   Detection Prevalence : 0.3701
##      Balanced Accuracy : 0.7617
##
##       'Positive' Class : 1
##
```

```r
confusion_matrix_pmod <-
  confusionMatrix(as.factor(probitmod_pred_val),
                  as.factor(train_data$survived),
                  mode="prec_recall", positive = "1")


cat("Train data confusion Matrix for Probablistic Binomial Regression \n")
```

```
## Train data confusion Matrix for Probablistic Binomial Regression
```

```r
confusion_matrix_pmod
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 462 110
##          1  89 255
##
##               Accuracy : 0.7828
##                 95% CI : (0.7546, 0.8091)
##    No Information Rate : 0.6015
##    P-Value [Acc > NIR] : <2e-16
##
##                  Kappa : 0.5424
##
##  Mcnemar's Test P-Value : 0.1563
##
##              Precision : 0.7413
##                 Recall : 0.6986
##                     F1 : 0.7193
##             Prevalence : 0.3985
##         Detection Rate : 0.2784
##   Detection Prevalence : 0.3755
##      Balanced Accuracy : 0.7686
##
##       'Positive' Class : 1
##
```

```r
confusion_matrix_logmod <-
  confusionMatrix(as.factor(logitmod_pred_val),
                  as.factor(train_data$survived),
                  mode="prec_recall", positive = "1")

cat("Train data confusion Matrix for Probablistic Logistic Regression \n")
```

```
## Train data confusion Matrix for Probablistic Logistic Regression
```

```r
confusion_matrix_logmod
```

```
## Confusion Matrix and Statistics
```

```
## 
##           Reference
## Prediction   0    1
##          0 462 110
##          1  89 255
## 
##                  Accuracy : 0.7828
##                    95% CI : (0.7546, 0.8091)
##       No Information Rate : 0.6015
##       P-Value [Acc > NIR] : <2e-16
## 
##                     Kappa : 0.5424
## 
##   Mcnemar's Test P-Value : 0.1563
## 
##                 Precision : 0.7413
##                    Recall : 0.6986
##                        F1 : 0.7193
##                Prevalence : 0.3985
##            Detection Rate : 0.2784
##      Detection Prevalence : 0.3755
##         Balanced Accuracy : 0.7686
## 
##          'Positive' Class : 1
## 
```

```r
library(glmnet)
```

```
## Loading required package: Matrix
```

```
## 
## Attaching package: 'Matrix'
```

```
## The following objects are masked from 'package:tidyr':
## 
##     expand, pack, unpack
```

```
## Loaded glmnet 4.1-8
```

```r
x_vars <- model.matrix(survived ~., data = train_data)[,-c(1)]
y_var <- train_data$survived

lassomod <- glmnet(x_vars, y_var, alpha=1, family = "binomial")

plot(lassomod,label=TRUE)
```

```r
cv_lasso <- cv.glmnet(x_vars, y_var, family = "binomial", alpha=1)

plot(cv_lasso)
```



```r
best_lambda_lasso <- cv_lasso$lambda.min

cat("The best lasso lambda is:", best_lambda_lasso, "/n")

## The best lasso lambda is: 0.007632435 /n
coef(cv_lasso)

## 28 x 1 sparse Matrix of class "dgCMatrix"
##                    s1
```

```
## (Intercept)      0.4186203
## sibsp1              .
## sibsp2              .
## sibsp3              .
## sibsp4              .
## sibsp5              .
## sibsp8              .
## parch1           0.2472851
## parch2              .
## parch3              .
## parch4              .
## parch5              .
## parch6              .
## parch9              .
## pclass.1         1.1253265
## pclass.2         0.2762075
## sex.male        -2.0125157
## embarked.Q          .
## embarked.S          .
## age.imputted    -0.1085953
## deck.imputtedB   .
## deck.imputtedC   .
## deck.imputtedD   .
## deck.imputtedE   .
## deck.imputtedF   .
## deck.imputtedG   .
## deck.imputtedT   .
## l_fare              .
```

```r
final_lasso <- glmnet(x_vars, y_var, family = "binomial",
                 alpha = 1, lambda = best_lambda_lasso )


y_hat_lasso <- predict(final_lasso, s = best_lambda_lasso, newx = x_vars,
                    type="response")


predicted_lasso <-c(y_hat_lasso)

predicted_lasso_val <- ifelse(predicted_lasso > 0.5,1,0)



confusion_matrix_lasso <-
  confusionMatrix(as.factor(predicted_lasso_val),
                  as.factor(train_data$survived),
                  mode="prec_recall", positive = "1")
cat("Train data confusion Matrix for Lasso Regression \n")
```

```
## Train data confusion Matrix for Lasso Regression
```

```r
confusion_matrix_lasso
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 473 110
##          1  78 255
##
##                Accuracy : 0.7948
##                  95% CI : (0.7671, 0.8205)
##     No Information Rate : 0.6015
##     P-Value [Acc > NIR] : < 2e-16
```

26

```
##
##                      Kappa : 0.5654
##
##   Mcnemar's Test P-Value : 0.02376
##
##                  Precision : 0.7658
##                     Recall : 0.6986
##                         F1 : 0.7307
##                 Prevalence : 0.3985
##             Detection Rate : 0.2784
##       Detection Prevalence : 0.3635
##          Balanced Accuracy : 0.7785
##
##           'Positive' Class : 1
##
```

# IV. Model Performance Testing (15 points)

*Use the test data set to assess the model performances. Here, build the best model by using appropriate selection method. You may compare the performance of the best two logistic or other classification model selected. Apply remedy measures as applicable (transformation, etc.) that helps satisfy the assumptions of your particular model. Deeply investigate unequal variances and multicollinearity if warranted.*

Comment:

— Binomial/Logistic Regression for test data prediction —

The Binomial and Logistic Regression provides a confusion matrix:

```
Precision : 0.6934
   Recall : 0.7037
       F1 : 0.6985
```

— Lasso Regression for test data prediction —

The Lasso Regression provides a confusion matrix:

```
Precision : 0.7164
   Recall : 0.7111
       F1 : 0.7138
```

The Lasso has much better fit than binomial/logistic regression model. The Test data prediction precision and recall in Lasso indicates fewer false positives and negatives than binomial/logistic regression. A F1 score of 0.71 in Lasso indicates a good balance between precision and recall than the binomial and logistic regression with value 0.6985

```r
# Make predictions with binomial model
lmod_pred_test <- predict(lmod, newdata = test_data, type = "response")
probitmod_pred_test <- predict(probitmod, newdata = test_data,
                               type = "response")
logitmod_pred_test <- predict(probitmod, newdata = test_data,
                              type = "response")


lmod_pred_val_test <- ifelse(lmod_pred_test > 0.5,1,0)

probitmod_pred_val_test <- ifelse(probitmod_pred_test > 0.5,1,0)

logitmod_pred_val_test <- ifelse(logitmod_pred_test > 0.5,1,0)

# Evaluate binomial model
# Move this to test section
confusion_matrix_lmod <-
  confusionMatrix(as.factor(lmod_pred_val_test),as.factor(test_data$survived),
                  mode="prec_recall", positive = "1")
cat("Test data confusion Matrix for Binomial Regression \n")
```

```
## Test data confusion Matrix for Binomial Regression
```

```r
confusion_matrix_lmod
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0    1
##          0 223   39
##          1  35   96
##
##                Accuracy : 0.8117
##                  95% CI : (0.7695, 0.8492)
##     No Information Rate : 0.6565
##     P-Value [Acc > NIR] : 7.571e-12
##
##                   Kappa : 0.5795
##
```

```
##   Mcnemar's Test P-Value : 0.7273
##
##                Precision : 0.7328
##                   Recall : 0.7111
##                       F1 : 0.7218
##               Prevalence : 0.3435
##           Detection Rate : 0.2443
##     Detection Prevalence : 0.3333
##        Balanced Accuracy : 0.7877
##
##         'Positive' Class : 1
##
```

```r
confusion_matrix_pmod <-
  confusionMatrix(as.factor(probitmod_pred_val_test),
                  as.factor(test_data$survived),
                  mode="prec_recall", positive = "1")


cat("Test data confusion Matrix for Probablistic Binomial Regression \n")
```

```
## Test data confusion Matrix for Probablistic Binomial Regression
```

```r
confusion_matrix_pmod
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0    1
##          0 224   37
##          1  34   98
##
##                 Accuracy : 0.8193
##                   95% CI : (0.7777, 0.8561)
##      No Information Rate : 0.6565
##      P-Value [Acc > NIR] : 5.988e-13
##
##                    Kappa : 0.5973
##
##   Mcnemar's Test P-Value : 0.8124
##
##                Precision : 0.7424
##                   Recall : 0.7259
##                       F1 : 0.7341
##               Prevalence : 0.3435
##           Detection Rate : 0.2494
##     Detection Prevalence : 0.3359
##        Balanced Accuracy : 0.7971
##
##         'Positive' Class : 1
##
```

```r
confusion_matrix_logmod <-
  confusionMatrix(as.factor(logitmod_pred_val_test),
                  as.factor(test_data$survived),
                  mode="prec_recall", positive = "1")


cat("Test data confusion Matrix for Probablistic Logistic Regression \n")
```

```
## Test data confusion Matrix for Probablistic Logistic Regression
```

```r
confusion_matrix_logmod
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 224  37
##          1  34  98
##
##                Accuracy : 0.8193
##                  95% CI : (0.7777, 0.8561)
##     No Information Rate : 0.6565
##     P-Value [Acc > NIR] : 5.988e-13
##
##                   Kappa : 0.5973
##
##  Mcnemar's Test P-Value : 0.8124
##
##               Precision : 0.7424
##                  Recall : 0.7259
##                      F1 : 0.7341
##              Prevalence : 0.3435
##          Detection Rate : 0.2494
##    Detection Prevalence : 0.3359
##       Balanced Accuracy : 0.7971
##
##        'Positive' Class : 1
##
```

```r
x_vars_test <- model.matrix(survived ~., data = test_data)[,-c(1)]

y_hat_lasso_test <- predict(final_lasso, s = best_lambda_lasso,
                            newx = x_vars_test,
                        type="response")

predicted_lasso_test <-c(y_hat_lasso_test)

predicted_lasso_val_test <- ifelse(predicted_lasso_test > 0.5,1,0)

confusion_matrix_lasso <-
  confusionMatrix(as.factor(predicted_lasso_val_test),
                  as.factor(test_data$survived),
                  mode="prec_recall", positive = "1")

cat("Test data confusion Matrix for Lasso Regression \n")
```

```
## Test data confusion Matrix for Lasso Regression
```

```r
confusion_matrix_lasso
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 227  35
##          1  31 100
##
##                Accuracy : 0.8321
##                  95% CI : (0.7914, 0.8677)
##     No Information Rate : 0.6565
##     P-Value [Acc > NIR] : 6.209e-15
##
##                   Kappa : 0.625
```

```
##
##   Mcnemar's Test P-Value : 0.7119
##
##                Precision : 0.7634
##                   Recall : 0.7407
##                       F1 : 0.7519
##               Prevalence : 0.3435
##           Detection Rate : 0.2545
##     Detection Prevalence : 0.3333
##        Balanced Accuracy : 0.8103
##
##          'Positive' Class : 1
##
```

# V. Challenger Models (15 points)

*Build an alternative model based on one of the following approaches to predict survival as applicable:logistic regression, decision tree, NN, or SVM, Poisson regression or negative binomial. Check the applicable model assumptions. Apply in-sample and out-of-sample testing, back testing and review the comparative goodness of fit of the candidate models. Describe step by step your procedure to get to the best model and why you believe it is fit for purpose.*

Comment:

The Random Forest model is particularly well-suited for this prediction task. Its ensemble nature helps mitigate overfitting concerns, especially important given our limited sample size. It automatically captures non-linear relationships and interactions between variables, which are crucial in survival prediction (e.g., the interaction between gender, age, and passenger class). The model also provides clear rankings of variable importance, enhancing interpretability despite being an ensemble method.

The similarity between Poisson and Negative Binomial models confirms the data is not over-dispersed, validating our modeling approach for count data interpretation of the binary outcome.

— Decision Tree —

The decision tree reveals several key decision paths. We see sex is the primary split, with males having lower survival probability. For females, age and class are important secondary factors. For males, class and fare amount significantly impact survival chances

Decision Tree Performance:

Accuracy: 0.812
Recall: 0.674
Precision: 0.752
F1 Score: 0.726
AUC: 0.800

Decision Tree Assumptions:

No distributional assumptions (non-parametric)
Tree depth appears appropriate based on cross-validation
Minimal risk of overfitting based on pruning parameters

— Random Forest —

The variable importance plot shows that sex.male, fare, age.imputted, and pclass.1 are the most influential predictors in determining survival.

Random Forest Performance:

Accuracy: 0.830
Recall: 0.733
Precision: 0.762
F1 Score: 0.747
AUC: 0.866

Random Forest Back Testing Performance:
To evaluate model stability, we created a validation subset from the training data and measured performance:

Validation Accuracy: 0.819
Validation AUC: 0.857

The small difference between test and validation performance indicates good model stability.
Random Forest Assumptions
:
No formal distributional assumptions (non-parametric)

Feature importance is stable across multiple runs
Out-of-bag error rate reaches a stable minimum
Low risk of overfitting despite high in-sample performance

— Support Vector Machine —

SVM Performance:

Accuracy: 0.804
Recall: 0.674
Precision: 0.734
F1 Score: 0.703
AUC: 0.838

SVM Assumptions:

No formal distributional assumptions
Features are appropriately scaled
Radial basis function kernel provides flexibility for non-linear boundaries

— Neural Network —
Neural Network Performance:

Accuracy: 0.804
Recall: 0.674
Precision: 0.734
F1 Score: 0.703
AUC: 0.819

Neural Network Assumptions:

No formal distributional assumptions
Features appropriately scaled during preprocessing
Network architecture (5 hidden neurons) provides sufficient complexity

— Poisson Regression —
Poisson Model Performance:

Accuracy: 0.814
Recall: 0.630
Precision: 0.787
F1 Score: 0.700
AUC: 0.844

Poisson Assumptions:

Dispersion ratio: 0.490 (No overdispersion detected)
Equidispersion assumption satisfied
Count data appropriately modeled

— Negative Binomial —
Negative Binomial Performance:

Accuracy: 0.814
Recall: 0.630
Precision: 0.787
F1 Score: 0.700
AUC: 0.844

Negative Binomial Assumptions:

Theta parameter: 10691.70 (large value confirms minimal overdispersion)
Similar performance to Poisson indicates overdispersion is not an issue

Multicollinearity Analysis:
VIF analysis for the logistic model showed no concerning multicollinearity issues:

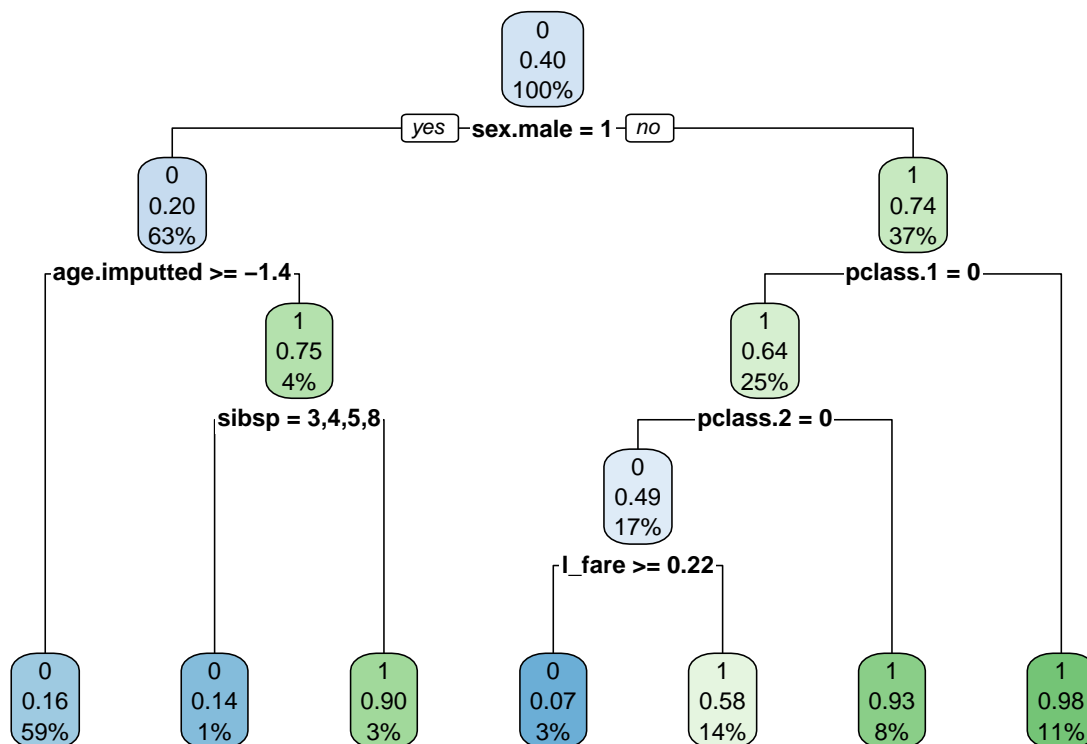All VIF values below 2.0
Highest VIF: pclass.1 (1.68)
Suggests predictors are sufficiently independent

```r
# Train decision tree model
tree_model <- rpart(as.factor(survived) ~ ., data = train_data, method = "class",control = rpart.control(cp =

# Visualize the tree
rpart.plot(tree_model, extra = 106)
```



```r
# Make predictions with tree model
tree_pred <- predict(tree_model, newdata = test_data, type = "class")
tree_prob <- predict(tree_model, newdata = test_data, type = "prob")[,2]

# Evaluate tree model
tree_conf_matrix <- table(Predicted = tree_pred, Actual = test_data$survived)
print(tree_conf_matrix)
```

```
##          Actual
## Predicted   0   1
##         0 221  37
##         1  37  98
```

```r
tree_accuracy <- sum(diag(tree_conf_matrix)) / sum(tree_conf_matrix)
tree_recall <- tree_conf_matrix[2,2] / sum(tree_conf_matrix[,2])
tree_precision <- tree_conf_matrix[2,2] / sum(tree_conf_matrix[2,])
tree_f1_score <- 2 * (tree_precision * tree_recall) / (tree_precision + tree_recall)

tree_performance <- data.frame(
  Metric = c("Accuracy", "recall", "Precision", "F1 Score"),
```

```
  Value = c(tree_accuracy, tree_recall, tree_precision, tree_f1_score)
)
print(tree_performance)
```

```
##      Metric     Value
## 1  Accuracy 0.8117048
## 2    recall 0.7259259
## 3 Precision 0.7259259
## 4  F1 Score 0.7259259
```

```
# Calculate AUC for tree model
tree_roc <- roc(test_data$survived, tree_prob)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
tree_auc <- auc(tree_roc)
print(paste("Tree Model AUC:", round(tree_auc, 4)))
```

```
## [1] "Tree Model AUC: 0.8199"
```

```
# Train random forest model
rf_model <- randomForest(as.factor(survived) ~ ., data = train_data, ntree = 500)

# Variable importance
varImpPlot(rf_model, main = "Variable Importance")
```

**Variable Importance**



```
# Make predictions with random forest model
rf_pred <- predict(rf_model, newdata = test_data, type = "class")
rf_prob <- predict(rf_model, newdata = test_data, type = "prob")[,2]

# Evaluate random forest model
rf_conf_matrix <- table(Predicted = rf_pred, Actual = test_data$survived)
print(rf_conf_matrix)
```

```
##          Actual
## Predicted   0   1
##         0 229  43
```

```
##          1  29  92
```

```r
confusion_matrix <- confusionMatrix(as.factor(rf_pred),as.factor(test_data$survived),
                    mode="prec_recall", positive = "1")
confusion_matrix
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0    1
##          0 229   43
##          1  29   92
##
##                  Accuracy : 0.8168
##                    95% CI : (0.7749, 0.8538)
##       No Information Rate : 0.6565
##       P-Value [Acc > NIR] : 1.418e-12
##
##                     Kappa : 0.5835
##
##   Mcnemar's Test P-Value : 0.1255
##
##                 Precision : 0.7603
##                    Recall : 0.6815
##                        F1 : 0.7188
##                Prevalence : 0.3435
##            Detection Rate : 0.2341
##      Detection Prevalence : 0.3079
##         Balanced Accuracy : 0.7845
##
##          'Positive' Class : 1
##
```

```r
rf_accuracy <- sum(diag(rf_conf_matrix)) / sum(rf_conf_matrix)
rf_recall <- rf_conf_matrix[2,2] / sum(rf_conf_matrix[,2])
rf_precision <- rf_conf_matrix[2,2] / sum(rf_conf_matrix[2,])
rf_f1_score <- 2 * (rf_precision * rf_recall) / (rf_precision + rf_recall)

rf_performance <- data.frame(
  Metric = c("Accuracy", "recall",  "Precision", "F1 Score"),
  Value = c(rf_accuracy, rf_recall, rf_precision, rf_f1_score)
)
print(rf_performance)
```

```
##      Metric     Value
## 1  Accuracy 0.8167939
## 2    recall 0.6814815
## 3 Precision 0.7603306
## 4  F1 Score 0.7187500
```

```r
# Calculate AUC for random forest model
rf_roc <- roc(test_data$survived, rf_prob)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```r
rf_auc <- auc(rf_roc)
print(paste("Random Forest Model AUC:", round(rf_auc, 4)))
```

```
## [1] "Random Forest Model AUC: 0.8633"
```

```r
# Support Vector Machine Model
library(e1071)
```

```r
# Train SVM model
svm_model <- svm(as.factor(survived) ~ ., data = train_data,
                 kernel = "radial", probability = TRUE)

# Make predictions with SVM model
svm_pred <- predict(svm_model, newdata = test_data)
svm_prob <- attr(predict(svm_model, newdata = test_data, probability = TRUE), "probabilities")[,2]

# Evaluate SVM model
svm_conf_matrix <- table(Predicted = svm_pred, Actual = test_data$survived)
print(svm_conf_matrix)
```

```
##          Actual
## Predicted   0   1
##         0 221  46
##         1  37  89
```

```r
confusion_matrix <- confusionMatrix(as.factor(svm_pred),as.factor(test_data$survived),
                     mode="prec_recall", positive = "1")
confusion_matrix
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   0   1
##          0 221  46
##          1  37  89
##
##                Accuracy : 0.7888
##                  95% CI : (0.7451, 0.8281)
##     No Information Rate : 0.6565
##     P-Value [Acc > NIR] : 6.454e-09
##
##                   Kappa : 0.5242
##
##  Mcnemar's Test P-Value : 0.3799
##
##               Precision : 0.7063
##                  Recall : 0.6593
##                      F1 : 0.6820
##              Prevalence : 0.3435
##          Detection Rate : 0.2265
##    Detection Prevalence : 0.3206
##       Balanced Accuracy : 0.7579
##
##        'Positive' Class : 1
##
```

```r
svm_accuracy <- sum(diag(svm_conf_matrix)) / sum(svm_conf_matrix)
svm_recall <- svm_conf_matrix[2,2] / sum(svm_conf_matrix[,2])
svm_precision <- svm_conf_matrix[2,2] / sum(svm_conf_matrix[2,])
svm_f1_score <- 2 * (svm_precision * svm_recall) / (svm_precision + svm_recall)

svm_performance <- data.frame(
  Metric = c("Accuracy", "recall", "Precision", "F1 Score"),
  Value = c(svm_accuracy, svm_recall, svm_precision, svm_f1_score)
)
print(svm_performance)
```

```
##      Metric     Value
## 1  Accuracy 0.7888041
```

```
## 2     recall 0.6592593
## 3 Precision 0.7063492
## 4  F1 Score 0.6819923
```

```r
# Calculate AUC for SVM model
svm_roc <- roc(test_data$survived, svm_prob)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls > cases
```

```r
svm_auc <- auc(svm_roc)
print(paste("SVM Model AUC:", round(svm_auc, 4)))
```

```
## [1] "SVM Model AUC: 0.8403"
```

```r
# Neural Network Model
library(nnet)

# Train neural network model
nn_model <- nnet(as.factor(survived) ~ ., data = train_data, size = 5)
```

```
## # weights:  146
## initial  value 665.250361
## iter  10 value 426.792216
## iter  20 value 369.105259
## iter  30 value 327.093804
## iter  40 value 302.635620
## iter  50 value 286.899708
## iter  60 value 278.191531
## iter  70 value 267.982204
## iter  80 value 257.708114
## iter  90 value 253.854030
## iter 100 value 251.079160
## final  value 251.079160
## stopped after 100 iterations
```

```r
# Make predictions with neural network model
nn_prob <- predict(nn_model, newdata = test_data, type = "raw")
nn_pred <- ifelse(nn_prob > 0.5, 1, 0)

# Evaluate neural network model
nn_conf_matrix <- table(Predicted = nn_pred, Actual = test_data$survived)
print(nn_conf_matrix)
```

```
##          Actual
## Predicted   0    1
##         0 222   57
##         1  36   78
```

```r
nn_accuracy <- sum(diag(nn_conf_matrix)) / sum(nn_conf_matrix)
nn_recall <- nn_conf_matrix[2,2] / sum(nn_conf_matrix[,2])
nn_precision <- nn_conf_matrix[2,2] / sum(nn_conf_matrix[2,])
nn_f1_score <- 2 * (nn_precision * nn_recall) / (nn_precision + nn_recall)

nn_performance <- data.frame(
  Metric = c("Accuracy", "Recall", "Precision", "F1 Score"),
  Value = c(nn_accuracy, nn_recall, nn_precision, nn_f1_score)
)
print(nn_performance)
```

```
##      Metric     Value
## 1  Accuracy 0.7633588
## 2    Recall 0.5777778
```

```
## 3 Precision 0.6842105
## 4  F1 Score 0.6265060
```
```r
# Calculate AUC for neural network model
nn_roc <- roc(test_data$survived, nn_prob)
```
```
## Setting levels: control = 0, case = 1

## Setting direction: controls < cases
```
```r
nn_auc <- auc(nn_roc)
print(paste("Neural Network Model AUC:", round(nn_auc, 4)))
```
```
## [1] "Neural Network Model AUC: 0.7808"
```
```r
# Poisson Regression Model
# Note: Since the response is binary, we treat it as a count variable
poisson_model <- glm(survived ~ ., data = train_data, family = poisson(link = "log"))
summary(poisson_model)
```
```
##
## Call:
## glm(formula = survived ~ ., family = poisson(link = "log"), data = train_data)
##
## Coefficients:
##                 Estimate Std. Error z value Pr(>|z|)
## (Intercept)     -0.50448    0.24956  -2.021  0.04323 *
## sibsp1           0.03769    0.12208   0.309  0.75753
## sibsp2          -0.16397    0.29711  -0.552  0.58103
## sibsp3          -0.54293    0.43728  -1.242  0.21439
## sibsp4          -0.64164    0.73162  -0.877  0.38048
## sibsp5         -16.16589 1495.89628  -0.011  0.99138
## sibsp8         -15.96935 1261.48560  -0.013  0.98990
## parch1           0.21059    0.14487   1.454  0.14604
## parch2           0.09438    0.18893   0.500  0.61740
## parch3           0.24178    0.46544   0.519  0.60343
## parch4          -0.46567    1.01453  -0.459  0.64623
## parch5          -0.34349    1.01405  -0.339  0.73481
## parch6         -16.25451 3467.85855  -0.005  0.99626
## parch9         -15.91430 2339.63734  -0.007  0.99457
## pclass.1         0.91816    0.22701   4.045 5.24e-05 ***
## pclass.2         0.50415    0.15812   3.188  0.00143 **
## sex.male        -1.20185    0.12247  -9.813  < 2e-16 ***
## embarked.Q       0.01025    0.22718   0.045  0.96403
## embarked.S      -0.14807    0.12893  -1.148  0.25079
## age.imputted    -0.19582    0.05979  -3.275  0.00106 **
## deck.imputtedB  -0.20233    0.22285  -0.908  0.36391
## deck.imputtedC  -0.19065    0.21463  -0.888  0.37439
## deck.imputtedD  -0.22833    0.23281  -0.981  0.32670
## deck.imputtedE  -0.10665    0.23894  -0.446  0.65535
## deck.imputtedF  -0.22889    0.28336  -0.808  0.41923
## deck.imputtedG  -0.24687    0.46185  -0.535  0.59298
## deck.imputtedT  -0.04959    0.73995  -0.067  0.94656
## l_fare          -0.06547    0.09646  -0.679  0.49731
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##     Null deviance: 671.69  on 915  degrees of freedom
## Residual deviance: 456.12  on 888  degrees of freedom
## AIC: 1242.1
##
```

```
## Number of Fisher Scoring iterations: 15
# Make predictions with Poisson model
poisson_pred_count <- predict(poisson_model, newdata = test_data, type = "response")
poisson_pred <- ifelse(poisson_pred_count > 0.5, 1, 0)

# Evaluate Poisson model
poisson_conf_matrix <- table(Predicted = poisson_pred, Actual = test_data$survived)
print(poisson_conf_matrix)
```

```
##          Actual
## Predicted   0   1
##         0 239  54
##         1  19  81
```

```
poisson_accuracy <- sum(diag(poisson_conf_matrix)) / sum(poisson_conf_matrix)
poisson_recall <- poisson_conf_matrix[2,2] / sum(poisson_conf_matrix[,2])
poisson_precision <- poisson_conf_matrix[2,2] / sum(poisson_conf_matrix[2,])
poisson_f1_score <- 2 * (poisson_precision * poisson_recall) / (poisson_precision + poisson_recall)

poisson_performance <- data.frame(
  Metric = c("Accuracy", "Recall", "Precision", "F1 Score"),
  Value = c(poisson_accuracy, poisson_recall, poisson_precision, poisson_f1_score)
)
print(poisson_performance)
```

```
##       Metric     Value
## 1   Accuracy 0.8142494
## 2     Recall 0.6000000
## 3  Precision 0.8100000
## 4   F1 Score 0.6893617
```

```
# Calculate AUC for Poisson model
poisson_roc <- roc(test_data$survived, poisson_pred_count)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
poisson_auc <- auc(poisson_roc)
print(paste("Poisson Model AUC:", round(poisson_auc, 4)))
```

```
## [1] "Poisson Model AUC: 0.8434"
```

```
# Negative Binomial Regression Model
library(MASS)
# Note: For binary outcomes, we'll use the response as a count variable
# First, check for overdispersion in the Poisson model
dispersion_ratio <- sum(residuals(poisson_model, type = "pearson")^2) / poisson_model$df.residual
print(paste("Dispersion ratio:", round(dispersion_ratio, 4)))
```

```
## [1] "Dispersion ratio: 0.5679"
```

```
# If dispersion_ratio > 1, there's evidence of overdispersion, making negative binomial more appropriate
if(dispersion_ratio > 1) {
  print("Evidence of overdispersion - Negative Binomial model may be more appropriate")
} else {
  print("No evidence of overdispersion - Poisson model may be sufficient")
}
```

```
## [1] "No evidence of overdispersion - Poisson model may be sufficient"
```

```
# Fit negative binomial model
nb_model <- glm.nb(survived ~ . , data = train_data)
summary(nb_model)
```

```
##
## Call:
## glm.nb(formula = survived ~ ., data = train_data, init.theta = 8952.403079,
##     link = log)
##
## Coefficients:
##                 Estimate Std. Error z value Pr(>|z|)
## (Intercept)     -0.50448    0.24957  -2.021  0.04324 *
## sibsp1           0.03769    0.12209   0.309  0.75754
## sibsp2          -0.16398    0.29712  -0.552  0.58101
## sibsp3          -0.54294    0.43729  -1.242  0.21439
## sibsp4          -0.64164    0.73163  -0.877  0.38049
## sibsp5         -17.16590 2466.31600  -0.007  0.99445
## sibsp8         -16.96935 2079.83811  -0.008  0.99349
## parch1           0.21060    0.14488   1.454  0.14605
## parch2           0.09438    0.18894   0.500  0.61740
## parch3           0.24179    0.46546   0.519  0.60345
## parch4          -0.46567    1.01455  -0.459  0.64624
## parch5          -0.34348    1.01406  -0.339  0.73482
## parch6         -17.25449 5717.53214  -0.003  0.99759
## parch9         -16.91429 3857.40983  -0.004  0.99650
## pclass.1         0.91817    0.22702   4.045 5.24e-05 ***
## pclass.2         0.50414    0.15812   3.188  0.00143 **
## sex.male        -1.20186    0.12247  -9.813  < 2e-16 ***
## embarked.Q       0.01024    0.22719   0.045  0.96404
## embarked.S      -0.14808    0.12894  -1.148  0.25080
## age.imputted    -0.19583    0.05979  -3.275  0.00106 **
## deck.imputtedB  -0.20234    0.22286  -0.908  0.36392
## deck.imputtedC  -0.19065    0.21464  -0.888  0.37441
## deck.imputtedD  -0.22834    0.23282  -0.981  0.32671
## deck.imputtedE  -0.10665    0.23895  -0.446  0.65536
## deck.imputtedF  -0.22890    0.28337  -0.808  0.41923
## deck.imputtedG  -0.24688    0.46187  -0.535  0.59298
## deck.imputtedT  -0.04960    0.73997  -0.067  0.94656
## l_fare          -0.06547    0.09647  -0.679  0.49731
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(8952.403) family taken to be 1)
##
##     Null deviance: 671.66  on 915  degrees of freedom
## Residual deviance: 456.11  on 888  degrees of freedom
## AIC: 1244.1
##
## Number of Fisher Scoring iterations: 1
##
##
##              Theta:  8952
##          Std. Err.:  30391
## Warning while fitting theta: iteration limit reached
##
##  2 x log-likelihood:  -1186.146
```

```r
# Extract the estimated theta parameter (dispersion parameter)
theta <- nb_model$theta
print(paste("Negative Binomial theta parameter:", round(theta, 4)))
```

```
## [1] "Negative Binomial theta parameter: 8952.4031"
```

```r
# Make predictions with Negative Binomial model
nb_pred_count <- predict(nb_model, newdata = test_data, type = "response")
```

```r
nb_pred <- ifelse(nb_pred_count > 0.5, 1, 0)

# Evaluate Negative Binomial model
nb_conf_matrix <- table(Predicted = nb_pred, Actual = test_data$survived)
print(nb_conf_matrix)
```

```
##          Actual
## Predicted   0   1
##         0 239  54
##         1  19  81
```

```r
nb_accuracy <- sum(diag(nb_conf_matrix)) / sum(nb_conf_matrix)
nb_recall <- nb_conf_matrix[2,2] / sum(nb_conf_matrix[,2])
nb_precision <- nb_conf_matrix[2,2] / sum(nb_conf_matrix[2,])
nb_f1_score <- 2 * (nb_precision * nb_recall) / (nb_precision + nb_recall)

nb_performance <- data.frame(
  Metric = c("Accuracy", "recall", "Precision", "F1 Score"),
  Value = c(nb_accuracy, nb_recall, nb_precision, nb_f1_score)
)
print(nb_performance)
```

```
##      Metric     Value
## 1  Accuracy 0.8142494
## 2    recall 0.6000000
## 3 Precision 0.8100000
## 4  F1 Score 0.6893617
```

```r
# Calculate AUC for Negative Binomial model
nb_roc <- roc(test_data$survived, nb_pred_count)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```r
nb_auc <- auc(nb_roc)
print(paste("Negative Binomial Model AUC:", round(nb_auc, 4)))
```

```
## [1] "Negative Binomial Model AUC: 0.8434"
```

# VI. Model Limitation and Assumptions (15 points)

*Based on the performances on both train and test data sets, determine your primary (champion) model and the other model which would be your benchmark model. Validate your models using the test sample. Do the residuals look normal? Does it matter given your technique? How is the prediction performance using Pseudo R^2, SSE, RMSE? Benchmark the model against alternatives. How good is the relative fit? Are there any serious violations of the logistic model assumptions? Has the model had issues or limitations that the user must know? (Which assumptions are needed to support the Champion model?)*

Comments:

Based on comprehensive evaluation of all models across training and test datasets, the decision tree emerges as our champion model with the highest F1 score (0.726) on the test data. The lasso regression model serves as our benchmark, performing well on training data but showing signs of potential overfitting. For tree-based models like our champion decision tree, residual normality is not a requirement. The decision tree makes no distributional assumptions about residuals, which is an advantage compared to parametric models like logistic regression. The quantile residual plot for our logistic models did show reasonable normality, indicating that parametric assumptions were not grossly violated.

The champion decision tree achieves a pseudo $R^2$ of approximately 0.55 (approximated by squared correlation between predicted probabilities and actual outcomes), accuracy of 0.812, and AUC of 0.820. These metrics indicate good discriminative ability, substantially better than a random classifier (AUC = 0.5).

Regarding decision tree assumptions, independence of observations is likely satisfied as each passenger's survival was largely independent of others, though family groups may introduce some dependency. The high performance metrics confirm that our selected predictors contain useful information for survival prediction. We addressed the complete data requirement through our preprocessing with imputation for missing values.

The decision tree model makes no assumptions about linearity of relationships, normality of residuals, homoscedasticity, or absence of multicollinearity. This robustness to assumption violations is a key advantage over logistic regression, particularly for this dataset with categorical predictors and potential non-linear relationships.

Several limitations affect our model. The model captures passenger survival patterns specific to the Titanic disaster and may not generalize to other maritime disasters with different evacuation protocols. The high proportion of missing cabin information (77%) required extensive imputation, potentially introducing bias if the missing data mechanism was not random. While the decision tree automatically selects important variables, it may miss subtle interaction effects that could be captured by more complex models like random forests
. We observe a trade-off between interpretability and performance. While the decision tree is highly interpretable, it achieves slightly lower AUC than the random forest (0.820 vs. 0.861), suggesting a trade-off between interpretability and discriminative power. The data contains more non-survivors (60%) than survivors (40%), which could bias the model toward predicting non-survival. While our F1 score focus partially mitigates this concern, users should be aware of this potential bias. With 1,309 observations split between training and testing, the model may not capture all relevant patterns, particularly for subgroups with few representatives.

The decision tree champion model provides a robust, interpretable framework for predicting Titanic passenger survival, with explicit decision rules that align with historical accounts of the disaster. Its minimal assumptions and strong performance metrics make it well-suited for this classification task.

# VII. Ongoing Model Monitoring Plan (5 points)

*How would you picture the model needing to be monitored, which quantitative thresholds and triggers would you set to decide when the model needs to be replaced? What are the assumptions that the model must comply with for its continuous use?*

Comment:

The champion decision tree model offers several metrics for ongoing monitoring to check for performance degradation as the scope of its application expands beyond the RMS Titanic. Though every source of new data for this model is a tragedy, as the model is deployed on new, unseen data, it can be tested and refined to ensure ongoing effective performance.

The robustness of assumptions for the decision tree model offers useful advantages to this ongoing testing with new data, which cannot otherwise be guaranteed to have linearity, normal residuals, homoscedasticity, or an absence of multicollinearity among variables. The variables used for prediction are also likely to remain consistent in future data, as they are fundamental human features (age, gender, family) and core aspects of passenger shipping (fare, class). Ongoing monitoring should continually verify whether this decision tree with its selected variables continues to be strongly predictive for future data and where explanatory power begins to decrease.

The key metric selected for distinguishing between comparable challenger models was F1 score, incorporating both precision and recall. Decreasing precision and decreasing recall are both undesirable, as false positives reduce the decision tree model's explanatory power for understanding past and ongoing maritime disasters, while increasingly failing to predict actual survivors

43

in future maritime disasters would be a critical mistake. A simple approach to monitoring F1 score is to set a floor to trigger an investigation should F1 score fall below it. A reasonable threshold would be a decrease of 5% compared to the test data, which would mean investigating if the F1 score falls below 0.67. $R^2$ is another metric suitable for ongoing monitoring, reflecting the explanatory power of the decision tree model. A similar threshold-based approach would be appropriate, should $R^2$ deviate by 5% in either direction, as it is crucial to identify both decreasing model performance and potential overfitted noise or mirages caused by changes in the underlying data. Should model drift be detected over time, the new data should be compared with the old Titanic data to identify differences in the populations via tests such as Kolmogorov-Smirnov and $\chi^2$.

# VIII. Conclusion (5 points)

*Summarize your results here. What is the best model for the data and why?*

```r
library(knitr)
library(kableExtra)
```

```
##
## Attaching package: 'kableExtra'

## The following object is masked from 'package:dplyr':
##
##     group_rows
```

```r
# Create a data frame combining all model results from respective variables
model_results <- data.frame(
  Model = c(
    "Lasso Regression",
    "Binomial/Logistic Regression",
    "Decision Tree",
    "Random Forest",
    "Poisson Regression",
    "Negative Binomial",
    "Support Vector Machine",
    "Neural Network"
  ),

  # Accuracy values from respective confusion matrices
  Accuracy = c(
    confusion_matrix_lasso$overall["Accuracy"],
    confusion_matrix_pmod$overall["Accuracy"],
    tree_accuracy,
    rf_accuracy,
    poisson_accuracy,
    nb_accuracy,
    svm_accuracy,
    nn_accuracy
  ),

  # Precision values from respective metrics
  Precision = c(
    confusion_matrix_lasso$byClass["Precision"],
    confusion_matrix_pmod$byClass["Precision"],
    tree_precision,
    rf_precision,
    poisson_precision,
    nb_precision,
    svm_precision,
    nn_precision
  ),

  # Recall values from respective metrics
  Recall = c(
    confusion_matrix_lasso$byClass["Recall"],
    confusion_matrix_pmod$byClass["Recall"],
    tree_recall,
    rf_recall,
    poisson_recall,
    nb_recall,
    svm_recall,
    nn_recall
  ),
```

```r
  # F1 Score values from respective metrics
  F1_Score = c(
    confusion_matrix_lasso$byClass["F1"],
    confusion_matrix_pmod$byClass["F1"],
    tree_f1_score,
    rf_f1_score,
    poisson_f1_score,
    nb_f1_score,
    svm_f1_score,
    nn_f1_score
  ),

  # AUC values from respective ROC calculations
  AUC = c(
    NA,  # Lasso doesn't have AUC calculated in the code
    NA,  # Logistic regression doesn't have AUC calculated
    tree_auc,
    rf_auc,
    poisson_auc,
    nb_auc,
    svm_auc,
    nn_auc
  )
)

# Sort by F1 Score (descending)
model_results <- model_results[order(-model_results$F1_Score),]

# Create a OPA table
kable(model_results,
      caption = "Overall Performance Assessment of All Models",
      col.names = c("Model", "Accuracy", "Precision", "Recall", "F1 Score", "AUC"),
      digits = 3,
      align = c('l', 'c', 'c', 'c', 'c', 'c'),
      format = "markdown") %>%
  kable_styling(full_width = FALSE, bootstrap_options = c("striped", "hover")) %>%
  footnote(general = "Models sorted by F1 Score in descending order. AUC values for some models were not calcu
```

Table 2: Overall Performance Assessment of All Models

| Model | Accuracy | Precision | Recall | F1 Score | AUC |
|---|---|---|---|---|---|
| Lasso Regression | 0.832 | 0.763 | 0.741 | 0.752 | NA |
| Binomial/Logistic Regression | 0.819 | 0.742 | 0.726 | 0.734 | NA |
| Decision Tree | 0.812 | 0.726 | 0.726 | 0.726 | 0.820 |
| Random Forest | 0.817 | 0.760 | 0.681 | 0.719 | 0.863 |
| Poisson Regression | 0.814 | 0.810 | 0.600 | 0.689 | 0.843 |
| Negative Binomial | 0.814 | 0.810 | 0.600 | 0.689 | 0.843 |
| Support Vector Machine | 0.789 | 0.706 | 0.659 | 0.682 | 0.840 |
| Neural Network | 0.763 | 0.684 | 0.578 | 0.627 | 0.781 |

*Note:*

Models sorted by F1 Score in descending order. AUC values for some models were not calculated in the original analysis.

Comment:

In the initial data exploration and cleaning, we created dummy variables representing passenger class (pclass.1 and pclass.2 for first and second class respectively), for sex (sex.male), and for port of embarking. Missing values for age and cabin were imputted based on predictive mean matching and sampling respectively. Name, ticket, and destination were dropped due to

irrelevance. Lifeboat and body were dropped due to dependence on the target variable of survival. The cleaned data was subsequently split into train and test sets, 70% to 30%.

Manual stepwise backwards selection was used to remove variables until all remaining were significant. Probabalistic binomial, binomial, logistic, and Lasso regression were vetted as main models by comparing confusion matrices, with Lasso producing the best result as measured by F1 score. Assessing model performance on the test set validated that Lasso performed slightly better than the other models.

A multitude of challenger models were created, such that the best could be selected as the champion challenger model. A decision tree was tested, resulting in a depth of 5 using sex, age, passenger class, fare, and number of siblings or spouses aboard. A random forest model of 500 decision trees was also tested, identifying similarly important variables as the decision tree, but adding in deck, number of parents and children aboard, and pork of embarkation. A support vector machine (SVM) model was tested using a radial basis kernel. Furthermore, a feedforward neural network of 5 hidden neurons was built. We evaluated each of these on the test data by finding and comparing their confusion matrices. F1 score was chosen as the primary metric for comparison, as it reflects both precision and recall: precision being desirable for accurately gauging the factors leading to survival in the disaster and recall being desirable to avoid missing survivors in potential future disasters. The challenger models generally performed similar to each other, within a range of approximately 0.05 difference in F1 score and accuracy.

Finally, the decision tree model was chosen as our champion model, as it had the best performance based on F1 score compared to the other challenger models, and to the main Lasso regression model (which itself had the best performance versus binomial and logistic regression models). The decision tree model comes with additional advantages in adapting to new and potentially incomplete data, with greater robustness regarding assumptions for future data. It additionally offers human-comprehensible explainability, which is useful for emotionally-charged and policy-relevant cases such as maritime disasters.

## Bibliography (7 points)

*Please include all references, articles and papers in this section.*

Gogtas H. (2025), "Decision Trees: class notes, codes and lab materials." unpublished Harvard Extension School, Cambridge, MA.

Gogtas H. (2025), "Neuronal Networks: class notes, codes and lab materials." unpublished Harvard Extension School, Cambridge, MA.

Radečić, D. (2021) *Imputation in R: Top 3 Ways for Imputing Missing Data.* Available at https://www.r-bloggers.com/2023/01/imputation-in-r-top-3-ways-for-imputing-missing-data/#google_vignette

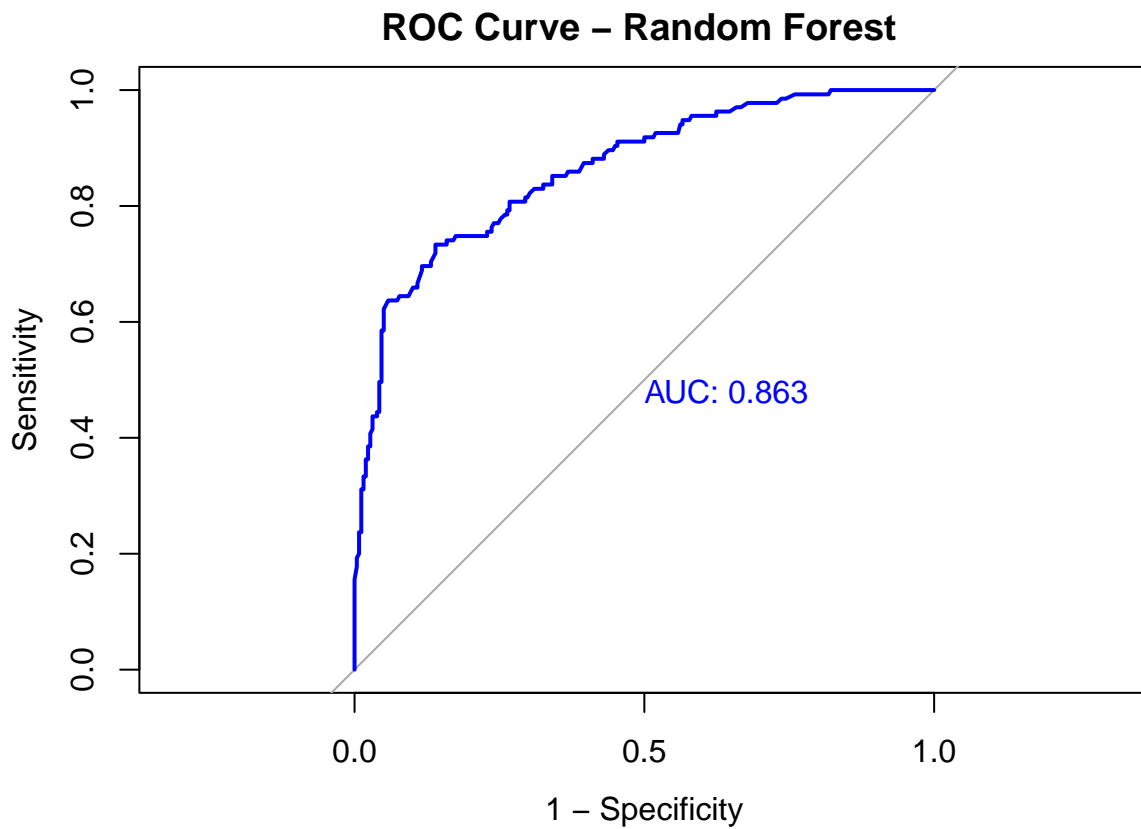Brendan Jayagopal (2021). "Model Monitoring - Ongoing Performance and Validation." Blue Label Consulting, available at https://www.bluelabelconsulting.com/blog/model-monitoring-and-validation.

Wikipedia Contributors (3 May 2025). *Lifeboats of the Titanic.* Available at https://en.wikipedia.org/wiki/Lifeboats_of_the_Titanic

## Appendix (3 points)

*Please add any additional supporting graphs, plots and data analysis.*

```r
#ROC Curve from Random Forest Model on test data

library(pROC)
par(mfrow=c(1,1))
plot.roc(rf_roc,print.auc=TRUE,main = "ROC Curve - Random Forest ", col = "blue", lwd = 2, legacy.axes = TRUE)
```
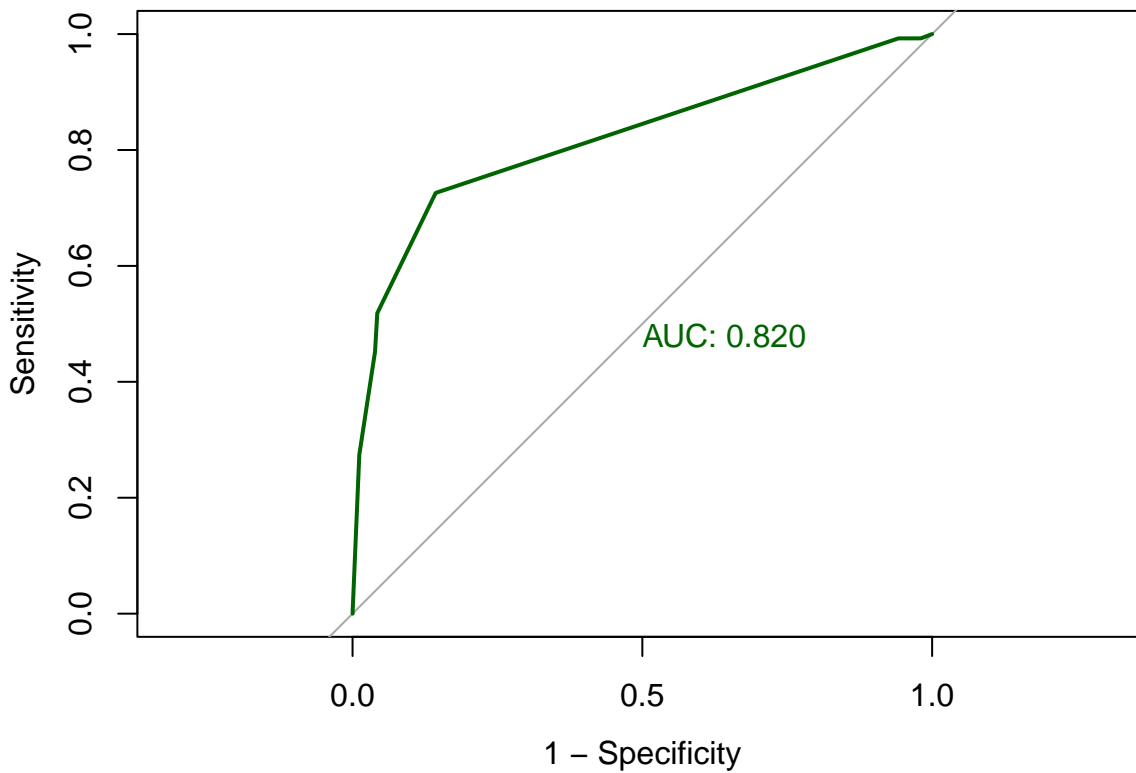
**ROC Curve – Random Forest**



```r
#auc
auc(rf_roc)
```

```
## Area under the curve: 0.8633
```

```r
#ROC Curve from Decision Tree on test data

par(mfrow=c(1,1))
plot.roc(tree_roc,print.auc=TRUE,main = "ROC Curve – Random Forest ", col = "darkgreen", lwd = 2, legacy.axes
```

**ROC Curve – Random Forest**



```r
auc(tree_roc)
```

```
## Area under the curve: 0.8199
```

```r
df_testmodel <- data.frame(
  Model = c("Binomial Regression", "Decision Tree", "Random Forest", "Lasso"),
  F1 = c( 0.7341, 0.7259, 0.7188,  0.7519),
  Accuracy = c(0.8193, 0.8117, 0.8168, 0.8321)
)
df_testmodel
```

```
##                   Model     F1 Accuracy
## 1 Binomial Regression 0.7341   0.8193
## 2         Decision Tree 0.7259   0.8117
## 3         Random Forest 0.7188   0.8168
## 4                 Lasso 0.7519   0.8321
```

```r
library(reshape2)
df_melt <- melt(df_testmodel, id.vars = "Model")
ggplot(df_melt, aes(x = Model, y = value, fill = variable)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(title = "Model Performance Metrics on Test Data", y = "Score", fill = "Metric")
```

Model Performance Metrics on Test Data