

Task4

rm

2025-04-10

```
#install and load imp library
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
library(ggplot2)
install.packages("GGally")

## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.4'
## (as 'lib' is unspecified)
library(GGally)

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2
library(readr)
install.packages("corrplot")

## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.4'
## (as 'lib' is unspecified)
library(corrplot)

## corrplot 0.95 loaded
```

Load dataset

```
titanic <- read_csv("train.csv")

## Rows: 891 Columns: 12
## -- Column specification -----
## Delimiter: ","
## chr (5): Name, Sex, Ticket, Cabin, Embarked
## dbl (7): PassengerId, Survived, Pclass, Age, SibSp, Parch, Fare
```

```
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

Structure and info

```
str(titanic)
```

```
## spc_tbl_ [891 x 12] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
## $ PassengerId: num [1:891] 1 2 3 4 5 6 7 8 9 10 ...
## $ Survived   : num [1:891] 0 1 1 1 0 0 0 0 1 1 ...
## $ Pclass     : num [1:891] 3 1 3 1 3 3 1 3 3 2 ...
## $ Name       : chr [1:891] "Braund, Mr. Owen Harris" "Cumings, Mrs. John Bradley (Florence Briggs T
## $ Sex        : chr [1:891] "male" "female" "female" "female" ...
## $ Age        : num [1:891] 22 38 26 35 35 NA 54 2 27 14 ...
## $ SibSp      : num [1:891] 1 1 0 1 0 0 0 3 0 1 ...
## $ Parch      : num [1:891] 0 0 0 0 0 0 0 1 2 0 ...
## $ Ticket     : chr [1:891] "A/5 21171" "PC 17599" "STON/O2. 3101282" "113803" ...
## $ Fare       : num [1:891] 7.25 71.28 7.92 53.1 8.05 ...
## $ Cabin      : chr [1:891] NA "C85" NA "C123" ...
## $ Embarked   : chr [1:891] "S" "C" "S" "S" ...
## - attr(*, "spec")=
## .. cols(
## ..   PassengerId = col_double(),
## ..   Survived = col_double(),
## ..   Pclass = col_double(),
## ..   Name = col_character(),
## ..   Sex = col_character(),
## ..   Age = col_double(),
## ..   SibSp = col_double(),
## ..   Parch = col_double(),
## ..   Ticket = col_character(),
## ..   Fare = col_double(),
## ..   Cabin = col_character(),
## ..   Embarked = col_character()
## .. )
## - attr(*, "problems")=<externalptr>
```

```
glimpse(titanic)
```

```
## Rows: 891
## Columns: 12
## $ PassengerId <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,~
## $ Survived    <dbl> 0, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1~
## $ Pclass      <dbl> 3, 1, 3, 1, 3, 3, 1, 3, 3, 2, 3, 1, 3, 3, 3, 2, 3, 2, 3, 3~
## $ Name        <chr> "Braund, Mr. Owen Harris", "Cumings, Mrs. John Bradley (Fl~
## $ Sex         <chr> "male", "female", "female", "female", "male", "male", "mal~
## $ Age         <dbl> 22, 38, 26, 35, 35, NA, 54, 2, 27, 14, 4, 58, 20, 39, 14, ~
## $ SibSp       <dbl> 1, 1, 0, 1, 0, 0, 0, 3, 0, 1, 1, 0, 0, 1, 0, 0, 4, 0, 1, 0~
## $ Parch       <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 1, 2, 0, 1, 0, 0, 5, 0, 0, 1, 0, 0, 0~
## $ Ticket      <chr> "A/5 21171", "PC 17599", "STON/O2. 3101282", "113803", "37~
## $ Fare        <dbl> 7.2500, 71.2833, 7.9250, 53.1000, 8.0500, 8.4583, 51.8625,~
## $ Cabin       <chr> NA, "C85", NA, "C123", NA, NA, "E46", NA, NA, NA, "G6", "C~
## $ Embarked    <chr> "S", "C", "S", "S", "S", "Q", "S", "S", "S", "C", "S", "S"~
```

Summary statistics

```
summary(titanic)
```

```
## 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
##
```

Value counts

```
table(titanic$Survived)
```

```
##
## 0 1
## 549 342
```

```
table(titanic$Sex)
```

```
##
## female male
## 314 577
```

```
table(titanic$Pclass)
```

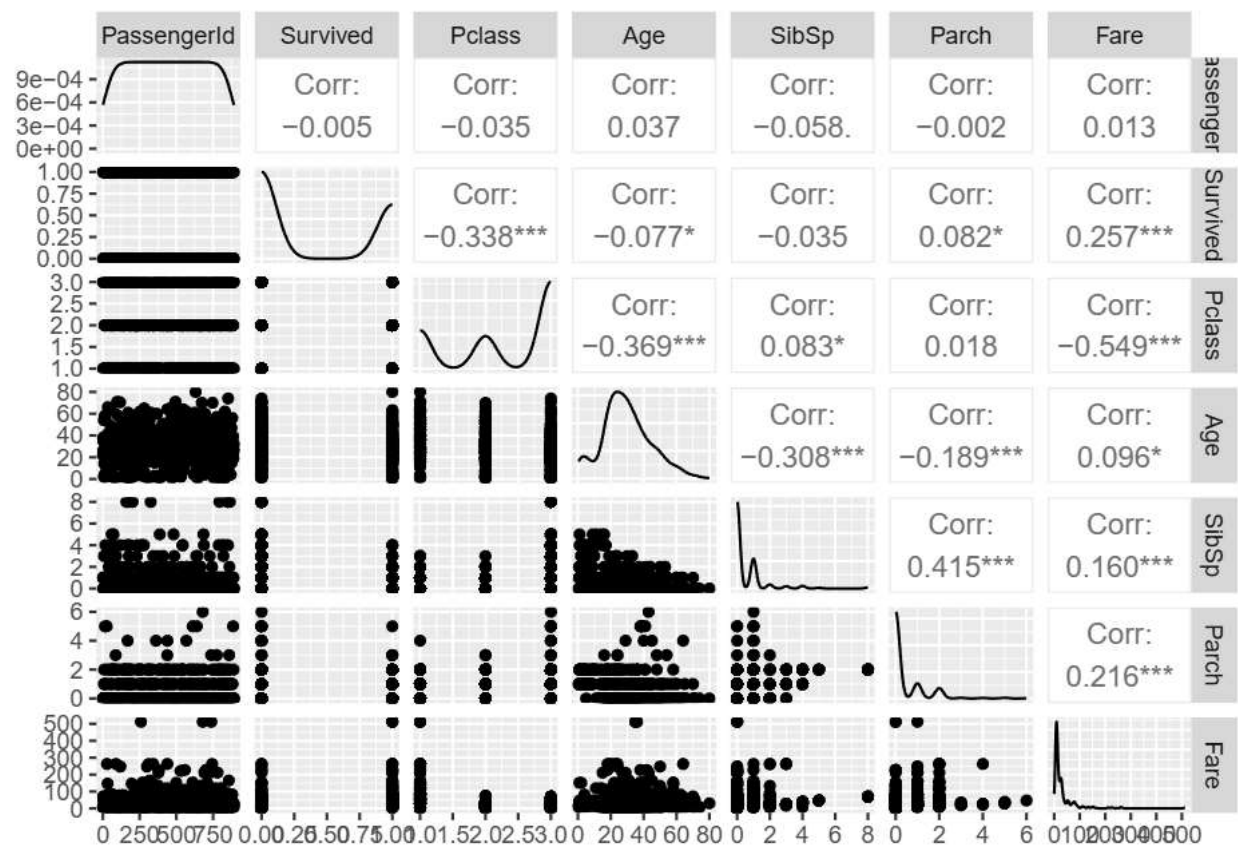
```
##
## 1 2 3
## 216 184 491
```

Use GGally to create pair plot

```
numeric_cols <- titanic %>%
  select_if(is.numeric)
```

```
ggpairs(numeric_cols)
```

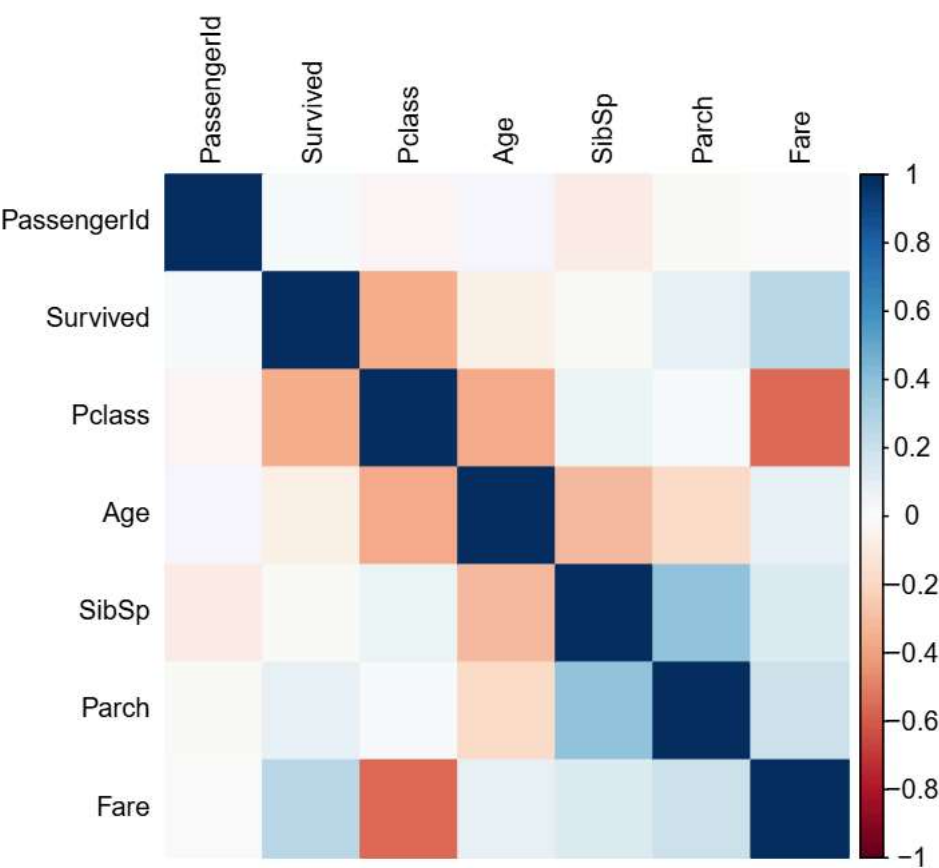
```
## Warning in ggally_statistic(data = data, mapping = mapping, na.rm = na.rm, :  
## Removed 177 rows containing missing values  
## Warning in ggally_statistic(data = data, mapping = mapping, na.rm = na.rm, :  
## Removed 177 rows containing missing values  
## Warning in ggally_statistic(data = data, mapping = mapping, na.rm = na.rm, :  
## Removed 177 rows containing missing values  
  
## Warning: Removed 177 rows containing missing values or values outside the scale range  
## (`geom_point()`).  
## Removed 177 rows containing missing values or values outside the scale range  
## (`geom_point()`).  
## Removed 177 rows containing missing values or values outside the scale range  
## (`geom_point()`).  
  
## Warning: Removed 177 rows containing non-finite outside the scale range  
## (`stat_density()`).  
  
## Warning in ggally_statistic(data = data, mapping = mapping, na.rm = na.rm, :  
## Removed 177 rows containing missing values  
## Warning in ggally_statistic(data = data, mapping = mapping, na.rm = na.rm, :  
## Removed 177 rows containing missing values  
## Warning in ggally_statistic(data = data, mapping = mapping, na.rm = na.rm, :  
## Removed 177 rows containing missing values  
  
## Warning: Removed 177 rows containing missing values or values outside the scale range  
## (`geom_point()`).  
## Removed 177 rows containing missing values or values outside the scale range  
## (`geom_point()`).  
## Removed 177 rows containing missing values or values outside the scale range  
## (`geom_point()`).
```

Clean NA for correlation

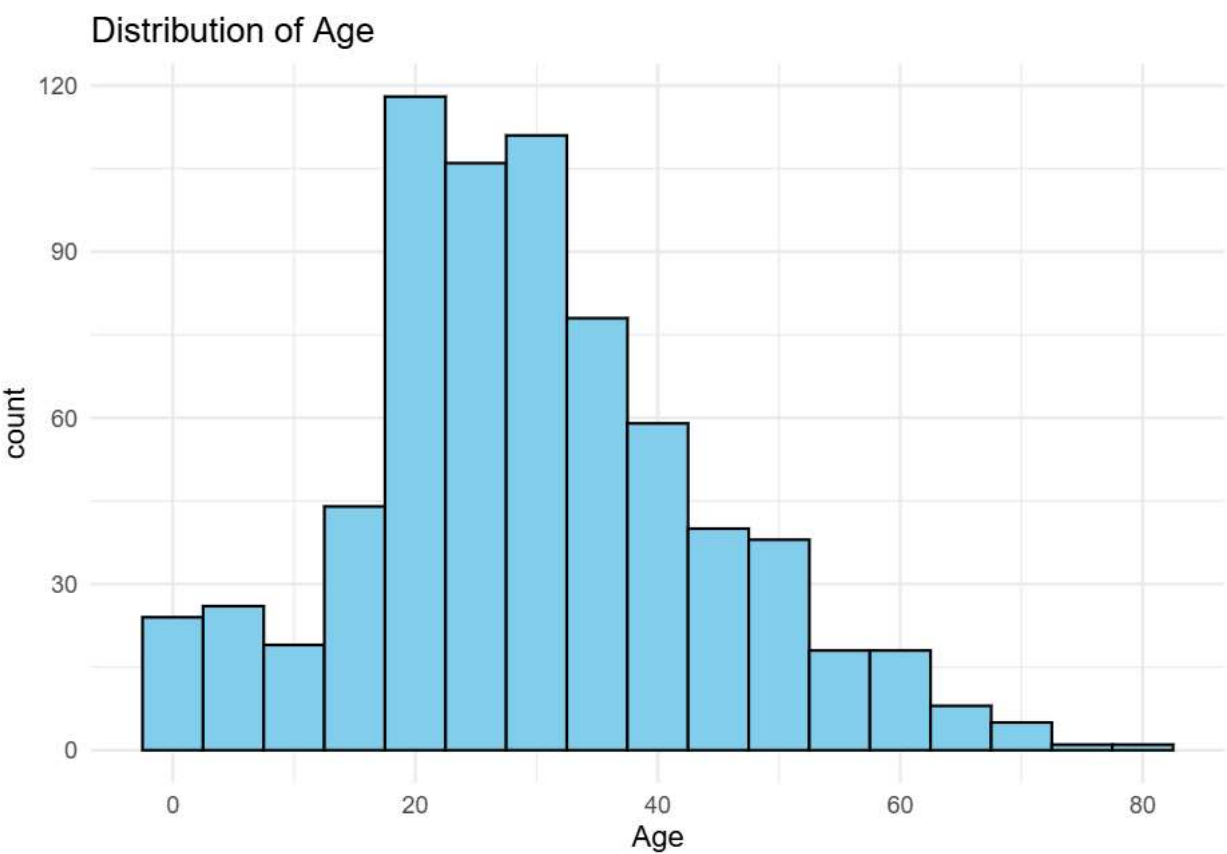
```
cor_data <- na.omit(numeric_cols)
cor_matrix <- cor(cor_data)

#coorelation heatmap
corrplot(cor_matrix, method = "color", tl.col = "black", tl.cex = 0.8)
```



```
#Histograms
ggplot(titanic, aes(x = Age)) +
  geom_histogram(binwidth = 5, fill = "skyblue", color = "black") +
  theme_minimal() +
  labs(title = "Distribution of Age")
```

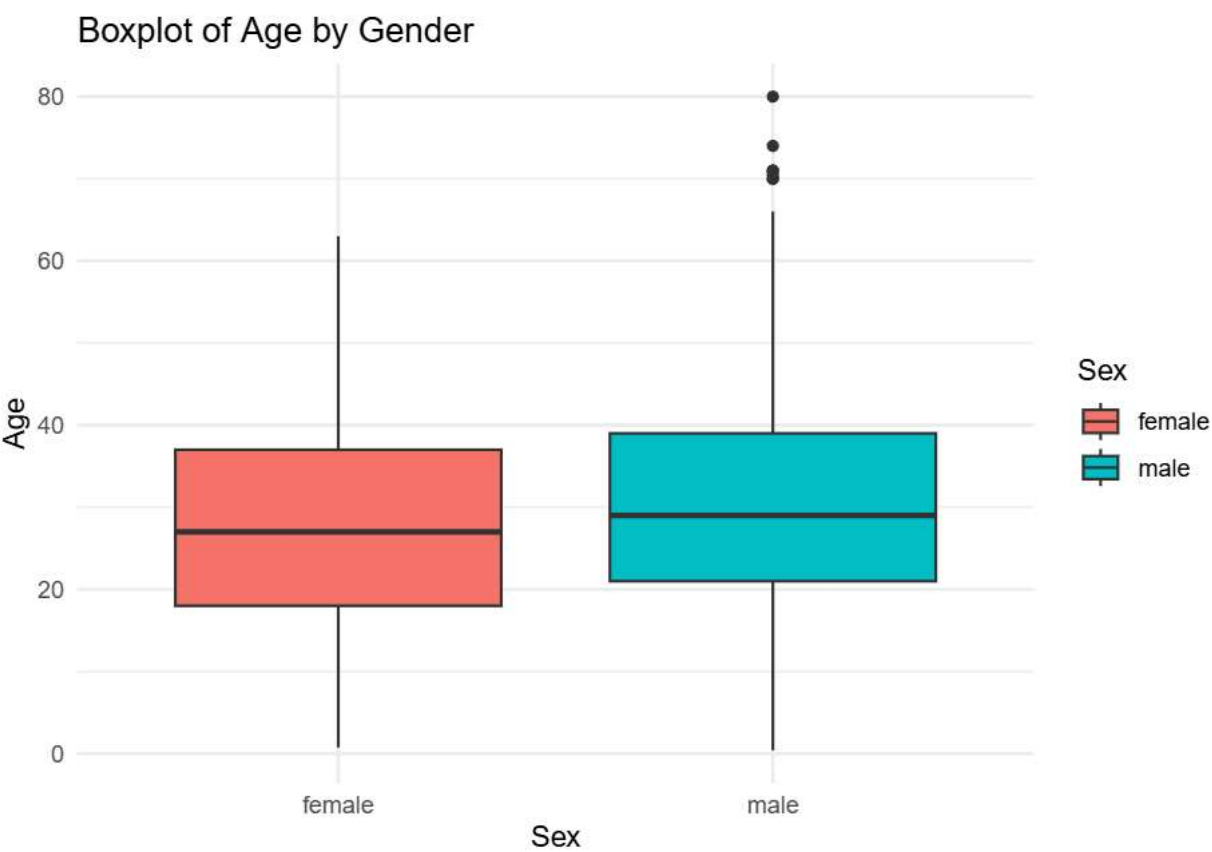
Warning: Removed 177 rows containing non-finite outside the scale range
(`stat_bin()`).



Boxplots

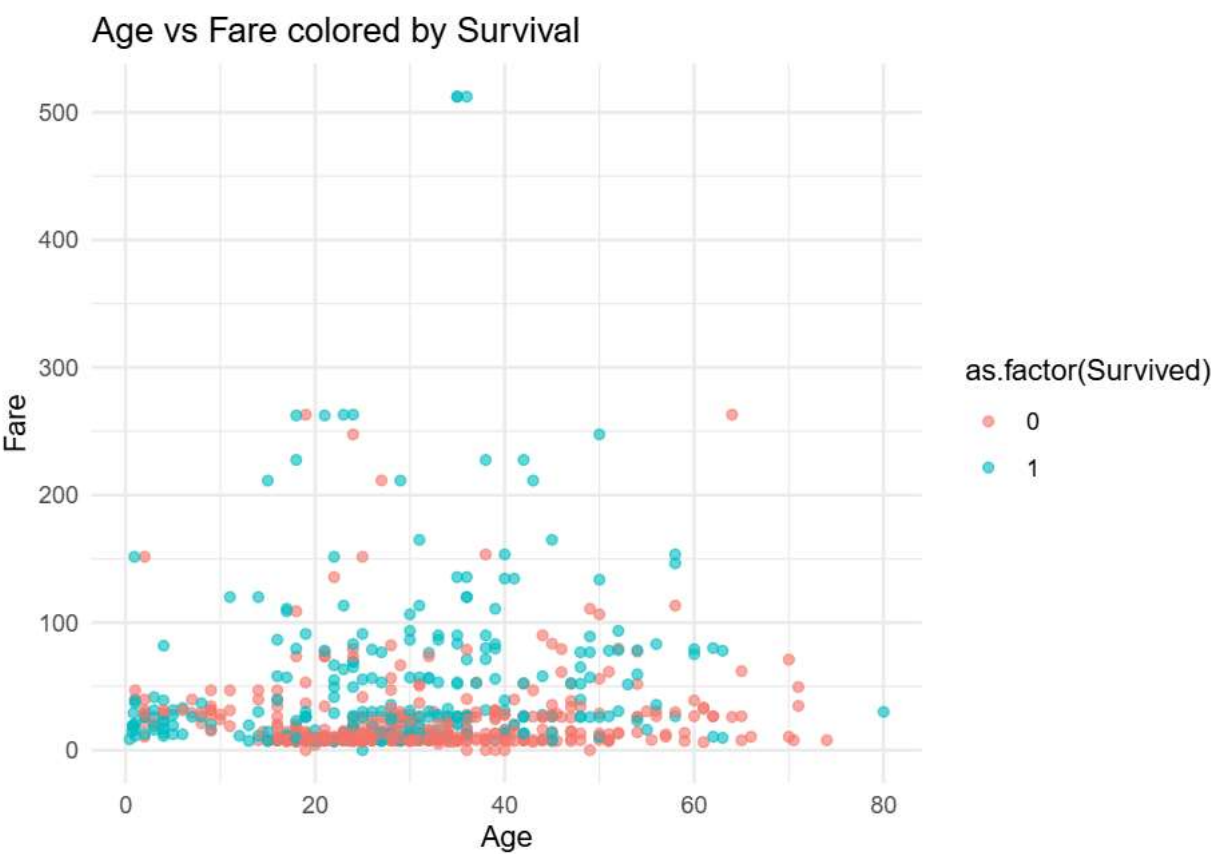
```
ggplot(titanic, aes(x = Sex, y = Age, fill = Sex)) +  
  geom_boxplot() +  
  theme_minimal() +  
  labs(title = "Boxplot of Age by Gender")
```

Warning: Removed 177 rows containing non-finite outside the scale range
(`stat_boxplot()`).

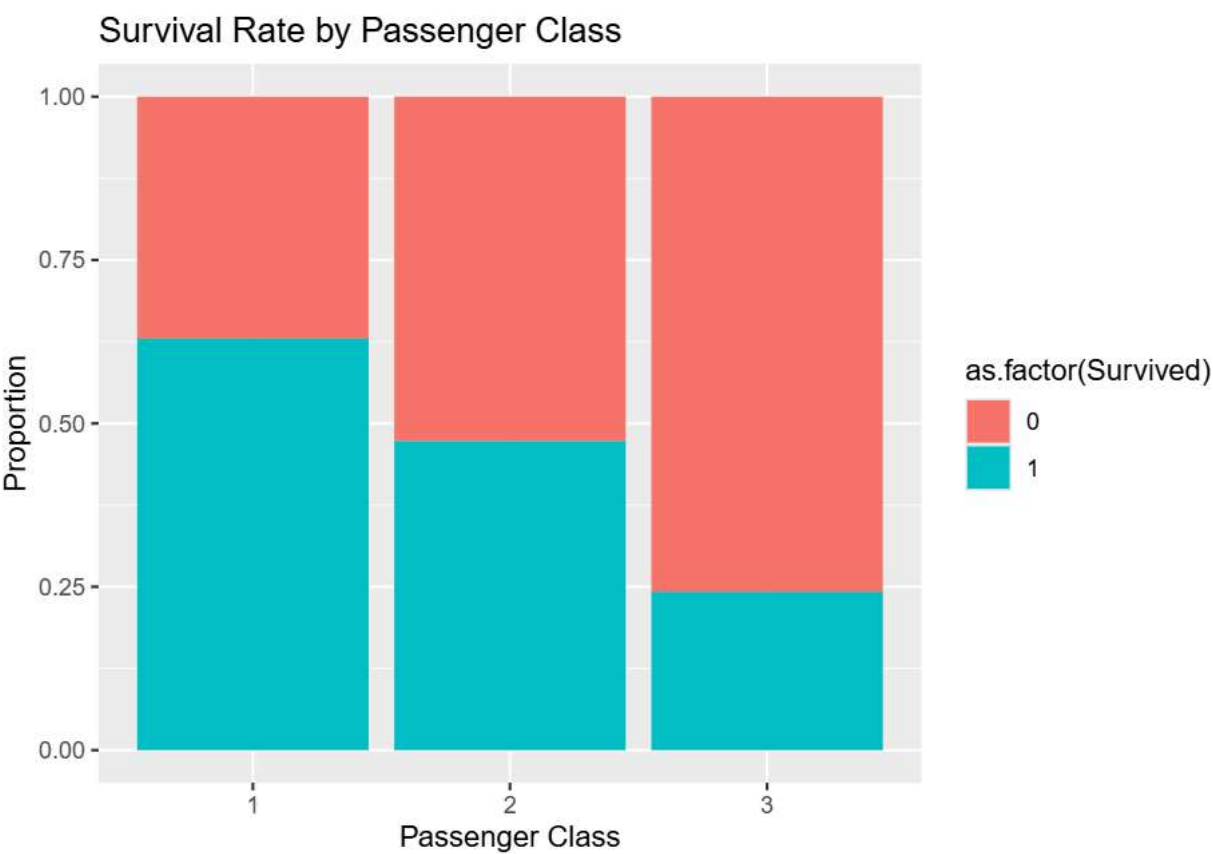


```
#Scatterplot
ggplot(titanic, aes(x = Age, y = Fare, color = as.factor(Survived))) +
  geom_point(alpha = 0.6) +
  theme_minimal() +
  labs(title = "Age vs Fare colored by Survival")

## Warning: Removed 177 rows containing missing values or values outside the scale range
## (`geom_point()`).
```

```
#Bar plot of Survival vs Class
ggplot(titanic, aes(x = as.factor(Pclass), fill = as.factor(Survived))) +
  geom_bar(position = "fill") +
  labs(title = "Survival Rate by Passenger Class", x = "Passenger Class", y = "Proportion")
```



#Majority of passengers were in 3rd class. #Females had higher survival rate than males. #Younger passengers had slightly better survival. #Higher fare and 1st class were associated with survival. #Age has some missing values that may need imputation.