

# PSTAT131HW03

Yifei Zhang

2022-04-17

## Contents

Question 1 . . . . .	2
Question 2 . . . . .	3
Question 3 . . . . .	4
Question 4 . . . . .	5
Question 5 . . . . .	6
Question 6 . . . . .	6
Question 7 . . . . .	7
Question 8 . . . . .	7
Question 9 . . . . .	7
Question 10 . . . . .	10

```
library(ggplot2)
library(tidyverse)
library(tidymodels)
library(corrplot)
library(ggthemes)
library(discrim)
library(poissonreg)
library(corr)
library(klaR)
titanic <- read.csv("titanic.csv")
```

```
titanic $ survived <- as_factor(titanic $ survived)

titanic $ pclass <- as_factor(titanic $ pclass)

titanic %>%
  head()
```

```
##   passenger_id survived pclass
## 1           1      No      3
## 2           2     Yes      1
## 3           3     Yes      3
## 4           4     Yes      1
```

```
## 5          5      No      3
## 6          6      No      3
##
##              name      sex age sib_sp parch
## 1              Braund, Mr. Owen Harris   male  22      1      0
## 2 Cumings, Mrs. John Bradley (Florence Briggs Thayer) female  38      1      0
## 3              Heikkinen, Miss. Laina female  26      0      0
## 4 Futrelle, Mrs. Jacques Heath (Lily May Peel) female  35      1      0
## 5              Allen, Mr. William Henry   male  35      0      0
## 6              Moran, Mr. James         male  NA      0      0
##      ticket      fare cabin embarked
## 1      A/5 21171  7.2500 <NA>      S
## 2      PC 17599 71.2833  C85      C
## 3 STON/O2. 3101282  7.9250 <NA>      S
## 4      113803 53.1000  C123      S
## 5      373450  8.0500 <NA>      S
## 6      330877  8.4583 <NA>      Q
```

```
titanic %>%
  arrange(survived) %>%
  head() # I dont know how to make "Yes" come first
```

```
##      passenger_id survived pclass              name      sex age sib_sp
## 1              1      No      3      Braund, Mr. Owen Harris male  22      1
## 2              5      No      3      Allen, Mr. William Henry male  35      0
## 3              6      No      3              Moran, Mr. James male  NA      0
## 4              7      No      1      McCarthy, Mr. Timothy J male  54      0
## 5              8      No      3 Palsson, Master. Gosta Leonard male   2      3
## 6             13      No      3 Saunderclock, Mr. William Henry male  20      0
##      parch      ticket      fare cabin embarked
## 1      0 A/5 21171  7.2500 <NA>      S
## 2      0  373450  8.0500 <NA>      S
## 3      0  330877  8.4583 <NA>      Q
## 4      0  17463 51.8625  E46      S
## 5      1  349909 21.0750 <NA>      S
## 6      0 A/5. 2151  8.0500 <NA>      S
```

## Question 1

Split the data, stratifying on the outcome variable, **survived**. You should choose the proportions to split the data into. Verify that the training and testing data sets have the appropriate number of observations. Take a look at the training data and note any potential issues, such as missing data.

The training and testing data sets seem to have the appropriate number of observations. First of all we have a lot of missing data in our dataset especially for the cabin section which I personally think is very important for our analysis, secondly in the ticket section the records are in combination of characters and numbers and that can be hard to analyze.

Why is it a good idea to use stratified sampling for this data?

It is a good idea to use stratified sampling for this data because we have a variety of attributes, make all the subgroups present in our sample and our result will be more accurate. We will be able to see fit better for the general population on their likelihood to survive on the sinking titanic.

```
set.seed(1010)

titanic_split <- initial_split(titanic, prop = 0.50,
                               strata = survived)

titanic_train <- training(titanic_split)

titanic_test <- testing(titanic_split)

count(titanic_train, "passenger_id") # see how many rows
```

```
##  "passenger_id"  n
## 1  passenger_id 445
```

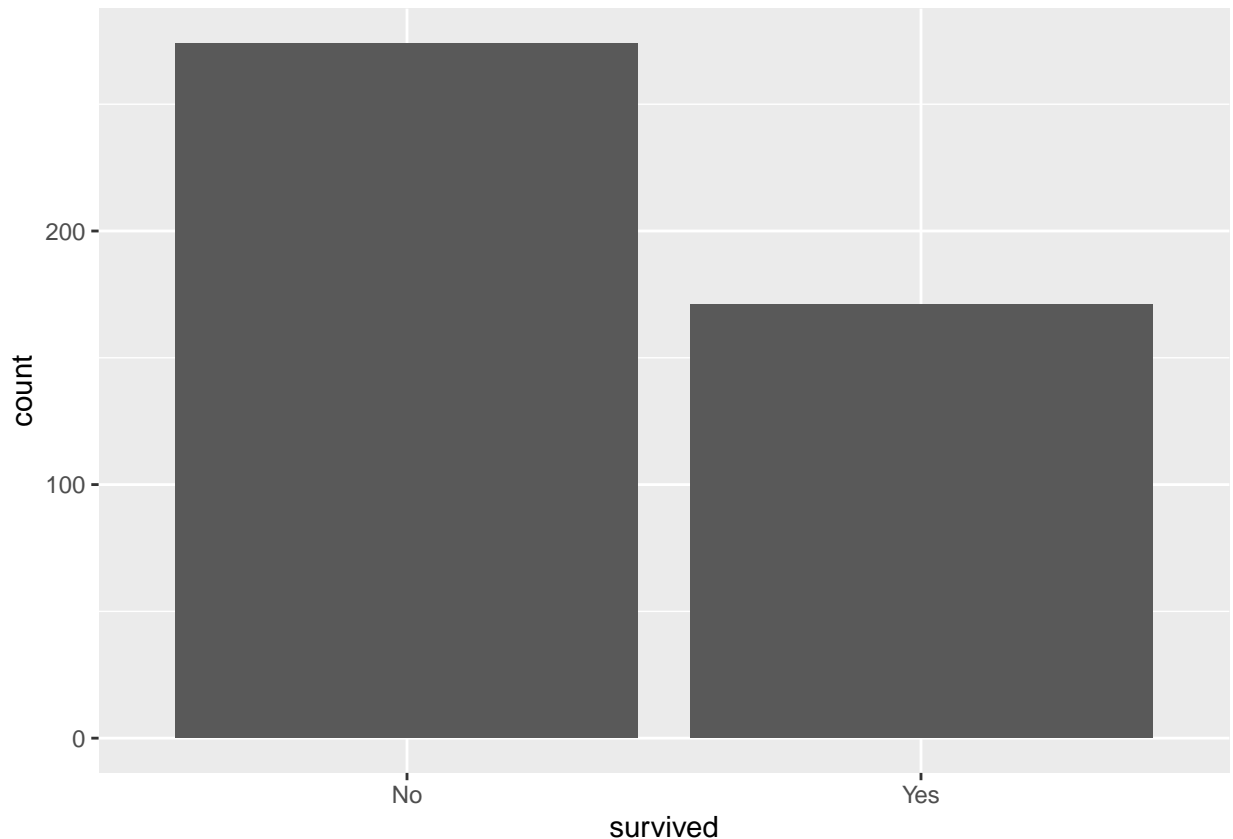
```
count(titanic_test, "passenger_id") # see how many rows
```

```
##  "passenger_id"  n
## 1  passenger_id 446
```

## Question 2

Using the **training** data set, explore/describe the distribution of the outcome variable **survived**.

```
titanic_train %>%
  ggplot(aes(x = survived)) +
  geom_bar()
```



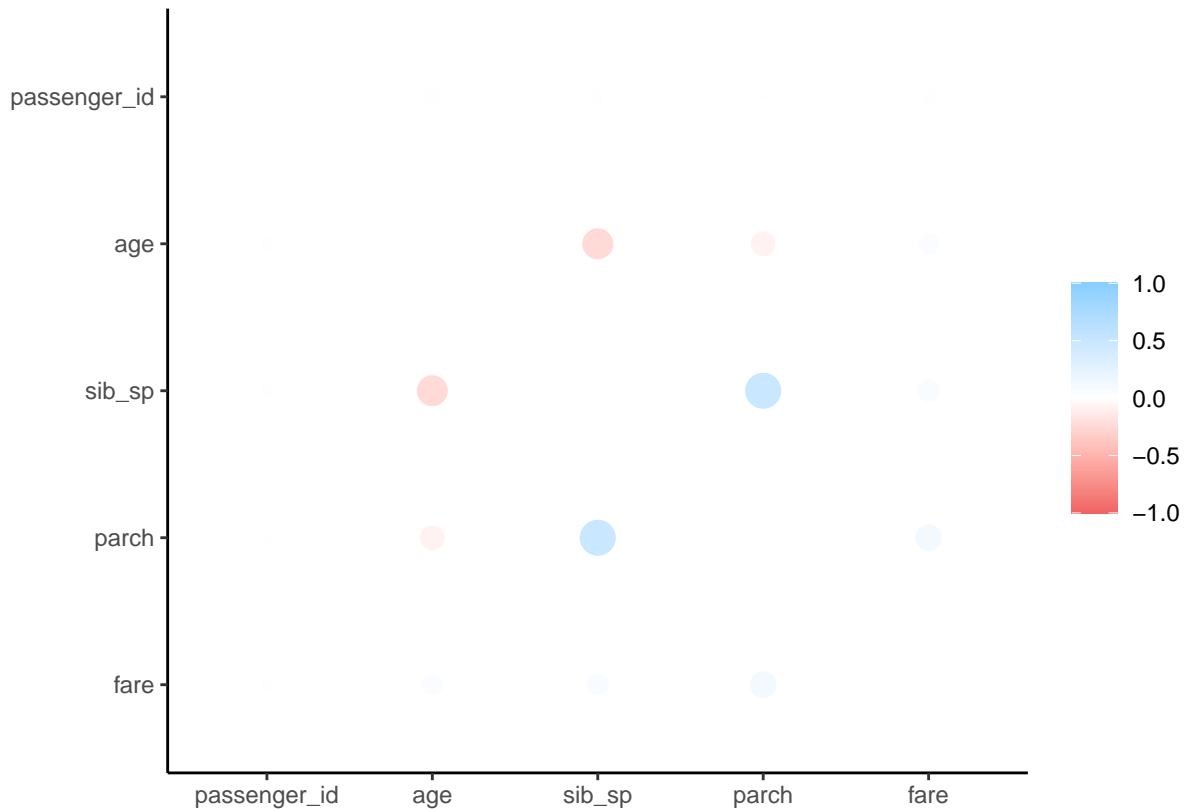
### Question 3

Using the **training** data set, create a correlation matrix of all continuous variables. Create a visualization of the matrix, and describe any patterns you see. Are any predictors correlated with each other? Which ones, and in which direction?

The numeric variables are passenger\_id, age, sib\_sp, parch, fare, although pclass seems like one of them, pclass is a factor and not numeric. A pattern I see is that a lot of them are not correlated at all. The number of siblings/ spouses aboard the Titanic has a relatively strong positive correlation with the number of parents/ children aboard the Titanic. The number of parents/ children aboard the Titanic and fare has a slightly positive correlation. The number of siblings/ spouses aboard the Titanic and fare also has a slightly positive correlation. The number of siblings/ spouses aboard the Titanic has a relatively strong negative correlation with age. There is a slightly negative correlation between the number of parents/ children aboard the Titanic and age.

```
cor_titanic_train <- titanic_train %>%
  dplyr::select(-c(survived, pclass, name, sex, ticket, cabin, embarked)) %>%
  correlate(use = "pairwise.complete.obs", method = "pearson")

rplot(cor_titanic_train)
```



#### Question 4

Using the **training** data, create a recipe predicting the outcome variable **survived**. Include the following predictors: ticket class, sex, age, number of siblings or spouses aboard, number of parents or children aboard, and passenger fare.

Recall that there were missing values for **age**. To deal with this, add an imputation step using `step_impute_linear()`. Next, use `step_dummy()` to **dummy** encode categorical predictors. Finally, include interactions between:

- Sex and passenger fare, and
- Age and passenger fare.

You'll need to investigate the `tidymodels` documentation to find the appropriate step functions to use.

```
titanic_train_recipe <- recipe(survived ~ pclass + sex +
                                age + sib_sp + parch + fare,
                                data = titanic_train) %>%
  step_impute_linear(age) %>%
  step_dummy(all_nominal_predictors()) %>%
  step_interact(terms = ~ sex_male : fare) %>%
  step_interact(terms = ~ age : fare)

titanic_train_recipe
```

```
## Data Recipe
##
## Inputs:
##
##   role #variables
##   outcome      1
##   predictor      6
##
## Operations:
##
## Linear regression imputation for age
## Dummy variables from all_nominal_predictors()
## Interactions with sex_male:fare
## Interactions with age:fare
```

### Question 5

Specify a **logistic regression** model for classification using the "glm" engine. Then create a workflow. Add your model and the appropriate recipe. Finally, use `fit()` to apply your workflow to the **training** data.

*Hint: Make sure to store the results of `fit()`. You'll need them later on.*

```
log_reg <- logistic_reg() %>%
  set_engine("glm") %>%
  set_mode("classification")

log_wkflow <- workflow() %>%
  add_model(log_reg) %>%
  add_recipe(titanic_train_recipe)

log_fit <- fit(log_wkflow, titanic_train)

log_fit %>%
  tidy()
```

```
## # A tibble: 10 x 5
##   term                estimate std.error statistic    p.value
##   <chr>              <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)        3.96      0.828      4.79 0.00000170
## 2 age               -0.0652    0.0162     -4.02 0.0000590
## 3 sib_sp             -0.502     0.152     -3.29 0.000988
## 4 parch             -0.271     0.211     -1.29 0.198
## 5 fare                0.0214    0.0168      1.27 0.204
## 6 pclass_X2          -0.801     0.454     -1.77 0.0773
## 7 pclass_X3          -1.91     0.469     -4.08 0.0000455
## 8 sex_male           -2.14     0.403     -5.30 0.000000114
## 9 sex_male_x_fare    -0.0275    0.0152     -1.81 0.0704
## 10 age_x_fare         0.000308  0.000225      1.37 0.171
```

### Question 6

**Repeat Question 5**, but this time specify a linear discriminant analysis model for classification using the "MASS" engine.

```
lda_mod <- discrim_linear() %>%
  set_mode("classification") %>%
  set_engine("MASS")

lda_wkflow <- workflow() %>%
  add_model(lda_mod) %>%
  add_recipe(titanic_train_recipe)

lda_fit <- fit(lda_wkflow, titanic_train)
```

### Question 7

**Repeat Question 5**, but this time specify a quadratic discriminant analysis model for classification using the "MASS" engine.

```
qda_mod <- discrim_quad() %>%
  set_mode("classification") %>%
  set_engine("MASS")

qda_wkflow <- workflow() %>%
  add_model(qda_mod) %>%
  add_recipe(titanic_train_recipe)

qda_fit <- fit(qda_wkflow, titanic_train)
```

### Question 8

**Repeat Question 5**, but this time specify a naive Bayes model for classification using the "klaR" engine. Set the `usekernel` argument to `FALSE`.

```
nb_mod <- naive_Bayes() %>%
  set_mode("classification") %>%
  set_engine("klaR") %>%
  set_args(usekernel = FALSE)

nb_wkflow <- workflow() %>%
  add_model(nb_mod) %>%
  add_recipe(titanic_train_recipe)

nb_fit <- fit(nb_wkflow, titanic_train)
```

### Question 9

Now you've fit four different models to your training data.

Use `predict()` and `bind_cols()` to generate predictions using each of these 4 models and your **training** data. Then use the *accuracy* metric to assess the performance of each of the four models.

Which model achieved the highest accuracy on the training data?

According to our results so far, it looks like the logistic regression model has the highest accuracy on the training data with an accuracy that's roughly 81.2%.

```
predict(log_fit, new_data = titanic_train, type = "prob")
```

```
## # A tibble: 445 x 2
##   .pred_No .pred_Yes
##   <dbl>    <dbl>
## 1  0.911    0.0890
## 2  0.892    0.108
## 3  0.800    0.200
## 4  0.934    0.0658
## 5  0.436    0.564
## 6  0.880    0.120
## 7  0.647    0.353
## 8  0.958    0.0421
## 9  0.742    0.258
## 10 0.411    0.589
## # ... with 435 more rows
```

```
predict(lda_fit, new_data = titanic_train, type = "prob")
```

```
## # A tibble: 445 x 2
##   .pred_No .pred_Yes
##   <dbl>    <dbl>
## 1  0.950    0.0504
## 2  0.923    0.0765
## 3  0.885    0.115
## 4  0.949    0.0507
## 5  0.326    0.674
## 6  0.932    0.0677
## 7  0.732    0.268
## 8  0.969    0.0310
## 9  0.793    0.207
## 10 0.309    0.691
## # ... with 435 more rows
```

```
predict(qda_fit, new_data = titanic_train, type = "prob")
```

```
## # A tibble: 445 x 2
##   .pred_No .pred_Yes
##   <dbl>    <dbl>
## 1  0.996 0.00442
## 2  0.999 0.000710
## 3  0.989 0.0108
## 4  1.00 0.0000172
## 5  0.993 0.00695
## 6  0.994 0.00615
## 7  0.931 0.0695
## 8  0.998 0.00199
## 9  0.962 0.0378
## 10 0.875 0.125
## # ... with 435 more rows
```



```
predict(nb_fit, new_data = titanic_train, type = "prob")
```

```
## # A tibble: 445 x 2
##   .pred_No .pred_Yes
##   <dbl>    <dbl>
## 1  0.976 0.0240
## 2  0.999 0.00114
## 3  0.970 0.0301
## 4  1.00 0.0000375
## 5  0.954 0.0464
## 6  0.974 0.0263
## 7  0.927 0.0729
## 8  0.900 0.100
## 9  0.789 0.211
## 10 0.713 0.287
## # ... with 435 more rows
```

```
log_reg_acc <- augment(log_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
lda_acc <- augment(lda_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
qda_acc <- augment(qda_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
nb_acc <- augment(nb_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
```

```
log_reg_acc <- augment(log_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
lda_acc <- augment(lda_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
qda_acc <- augment(qda_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)
nb_acc <- augment(nb_fit, new_data = titanic_train) %>%
  accuracy(truth = survived, estimate = .pred_class)

log_reg_acc
```

```
## # A tibble: 1 x 3
##   .metric .estimator .estimate
##   <chr>    <chr>      <dbl>
## 1 accuracy binary      0.816
```

```
lda_acc
```

```
## # A tibble: 1 x 3
##   .metric .estimator .estimate
##   <chr>    <chr>      <dbl>
## 1 accuracy binary      0.802
```

```
qda_acc
```

```
## # A tibble: 1 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>       <dbl>
## 1 accuracy binary      0.804
```

```
nb_acc
```

```
## # A tibble: 1 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>       <dbl>
## 1 accuracy binary      0.793
```

### Question 10

Fit the model with the highest training accuracy to the **testing** data. Report the accuracy of the model on the **testing** data.

The model with the highest training accuracy got a 84.35754% accuracy on the testing data.

Again using the **testing** data, create a confusion matrix and visualize it. Plot an ROC curve and calculate the area under it (AUC).

How did the model perform? Compare its training and testing accuracies. If the values differ, why do you think this is so?

The model did pretty well, the model worked better on the training data than the training data which are 0.8157303 to 0.8004484. They are pretty close. I think the reason why the accuracy is higher for the training data than the test data maybe because of over fitting. And we should get higher accuracy on the training data anyway.

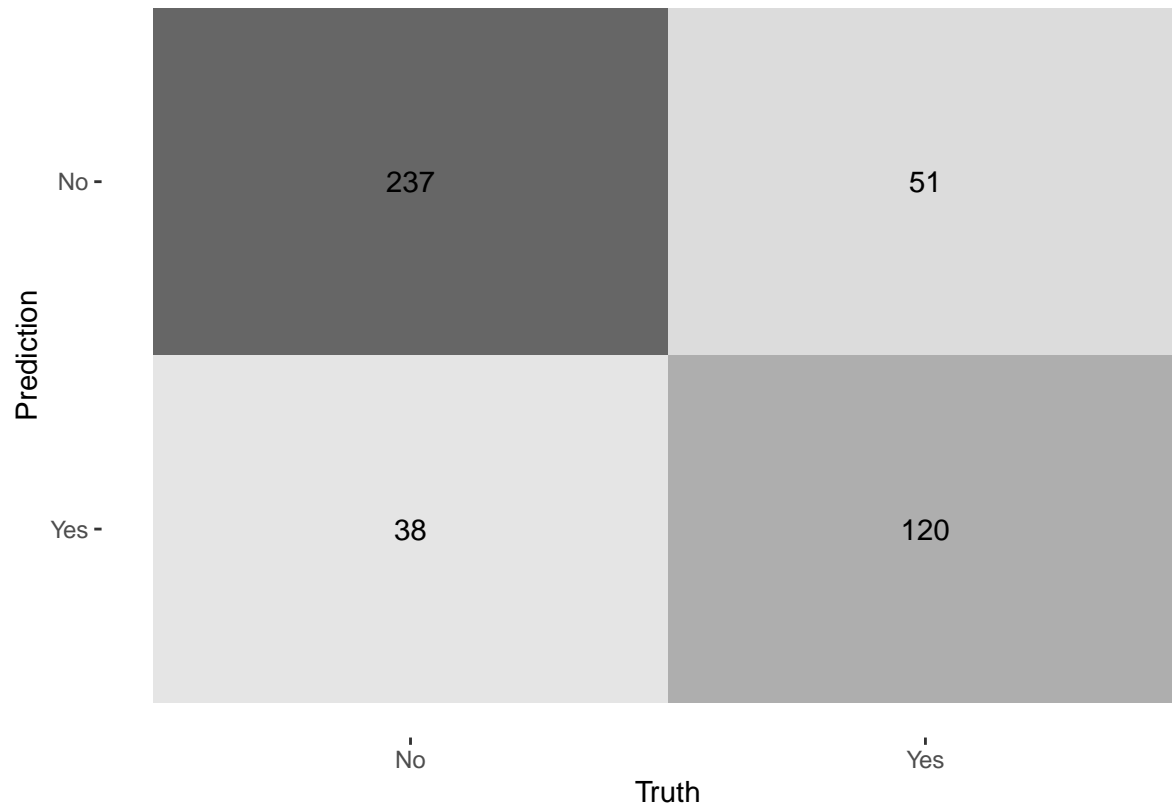
```
predict(log_fit, new_data = titanic_test, type = "prob")
```

```
## # A tibble: 446 x 2
##   .pred_No .pred_Yes
##   <dbl>   <dbl>
## 1 0.882 0.118
## 2 0.356 0.644
## 3 0.880 0.120
## 4 0.759 0.241
## 5 0.479 0.521
## 6 0.198 0.802
## 7 0.227 0.773
## 8 0.987 0.0134
## 9 0.207 0.793
## 10 0.452 0.548
## # ... with 436 more rows
```

```
log_reg_acc <- augment(log_fit, new_data = titanic_test) %>%
  accuracy(truth = survived, estimate = .pred_class)
log_reg_acc
```

```
## # A tibble: 1 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>       <dbl>
## 1 accuracy binary      0.800
```

```
augment(log_fit, new_data = titanic_test) %>%
  conf_mat(truth = survived, estimate = .pred_class) %>%
  autoplot(type = "heatmap")
```



```
augment(log_fit, new_data = titanic_test) %>%
  roc_curve(survived, .pred_No) %>%
  autoplot()
```

