# HW 4

Samuel Pell

03/18/2024

This homework is designed to give you practice fitting a logistic regression and working with statistical/philosophical measures of fairness. We will work with the `titanic` dataset which we have previously seen in class in connection to decision trees.

Below I will preprocess the data precisely as we did in class. You can simply refer to `data_train` as your training data and `data_test` as your testing data.

```r
#this is all of the preprocessing done for the decision trees lecture.

path <- 'https://raw.githubusercontent.com/guru99-edu/R-
Programming/master/titanic_data.csv'
titanic <-read.csv(path)
head(titanic)

##   x pclass survived                                           name    sex
## 1 1      1        1                   Allen, Miss. Elisabeth Walton female
## 2 2      1        1                  Allison, Master. Hudson Trevor   male
## 3 3      1        0                   Allison, Miss. Helen Loraine female
## 4 4      1        0         Allison, Mr. Hudson Joshua Creighton   male
## 5 5      1        0 Allison, Mrs. Hudson J C (Bessie Waldo Daniels) female
## 6 6      1        1                            Anderson, Mr. Harry   male
##       age sibsp parch ticket     fare   cabin embarked
## 1      29     0     0  24160 211.3375      B5        S
## 2 0.9167     1     2 113781   151.55 C22 C26        S
## 3      2     1     2 113781   151.55 C22 C26        S
## 4     30     1     2 113781   151.55 C22 C26        S
## 5     25     1     2 113781   151.55 C22 C26        S
## 6     48     0     0  19952    26.55     E12        S
##                      home.dest
## 1                   St Louis, MO
## 2 Montreal, PQ / Chesterville, ON
## 3 Montreal, PQ / Chesterville, ON
## 4 Montreal, PQ / Chesterville, ON
## 5 Montreal, PQ / Chesterville, ON
## 6                   New York, NY

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

#replace ? with NA
replace_question_mark <- function(x) {
  if (is.character(x)) {
    x <- na_if(x, "?")
  }
  return(x)
}

titanic <- titanic %>%
  mutate_all(replace_question_mark)

set.seed(678)
shuffle_index <- sample(1:nrow(titanic))
head(shuffle_index)

## [1]   57  774  796 1044  681  920

titanic <- titanic[shuffle_index, ]
head(titanic)

##           x pclass survived                                        name
## 57       57      1        1                      Carter, Mr. William Ernest
## 774     774      3        0                              Dimic, Mr. Jovan
## 796     796      3        0                        Emir, Mr. Farred Chehab
## 1044   1044      3        1                    Murphy, Miss. Margaret Jane
## 681     681      3        0                              Boulos, Mr. Hanna
## 920     920      3        0 Katavelas, Mr. Vassilios ('Catavelas Vassilios')
##           sex   age sibsp parch ticket    fare   cabin embarked      home.dest
## 57       male    36     1     2 113760     120 B96 B98        S Bryn Mawr, PA
## 774      male    42     0     0 315088  8.6625    <NA>        S          <NA>
## 796      male  <NA>     0     0   2631   7.225    <NA>        C          <NA>
## 1044   female  <NA>     1     0 367230    15.5    <NA>        Q          <NA>
## 681      male  <NA>     0     0   2664   7.225    <NA>        C         Syria
## 920      male  18.5     0     0   2682  7.2292    <NA>        C          <NA>

library(dplyr)
# Drop variables
clean_titanic <- titanic %>%
select(-c(home.dest, cabin, name, x, ticket)) %>%
#Convert to factor level
    mutate(pclass = factor(pclass, levels = c(1, 2, 3), labels = c('Upper',
'Middle', 'Lower')),
    survived = factor(survived, levels = c(0, 1), labels = c('No', 'Yes')))
```

```
  %>%
  na.omit()
  #previously were characters
  clean_titanic$age <- as.numeric(clean_titanic$age)
  clean_titanic$fare <- as.numeric(clean_titanic$fare)
  glimpse(clean_titanic)

  ## Rows: 1,043
  ## Columns: 8
  ## $ pclass   <fct> Upper, Lower, Lower, Middle, Lower, Middle, Lower, Lower,
  Upp…
  ## $ survived <fct> Yes, No, No, No, No, No, No, No, Yes, No, Yes, No, No,
  Yes, N…
  ## $ sex      <chr> "male", "male", "male", "male", "female", "female",
  "male", "…
  ## $ age      <dbl> 36.0, 42.0, 18.5, 44.0, 19.0, 26.0, 23.0, 28.5, 64.0,
  36.5, 4…
  ## $ sibsp    <int> 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0,
  0, 0…
  ## $ parch    <int> 2, 0, 0, 0, 0, 1, 0, 0, 2, 2, 1, 0, 0, 0, 0, 1, 0, 0, 0,
  0, 0…
  ## $ fare     <dbl> 120.0000, 8.6625, 7.2292, 13.0000, 16.1000, 26.0000,
  7.8542, …
  ## $ embarked <chr> "S", "S", "C", "S", "S", "S", "S", "S", "C", "S", "S",
  "S", "…

  create_train_test <- function(data, size = 0.8, train = TRUE) {
      n_row = nrow(data)
      total_row = size * n_row
      train_sample <- 1: total_row
      if (train == TRUE) {
          return (data[train_sample, ])
      } else {
          return (data[-train_sample, ])
      }
  }
  data_train <- create_train_test(clean_titanic, 0.8, train = TRUE)
  data_test <- create_train_test(clean_titanic, 0.8, train = FALSE)
```

Create a table reporting the proportion of people in the training set surviving the Titanic.
Do the same for the testing set. Comment on whether the current training-testing partition
looks suitable.

```
#Student Input
round(prop.table(table(data_train$survived)),3)*100
```

```
## 
##    No  Yes
## 60.2 39.8
```

```
round(prop.table(table(data_test$survived)),3)*100
```

```
## 
##    No  Yes
## 55.5 44.5
```

*student input*

Based on the proportions given above, I would suggest that perhaps a different training/testing split should be used to evaluate the data. While the survival rates are not extremely different, there is a marked difference in the two proportions, which could be indicative of some skew or failure to properly sample.

Use the `glm` command to build a logistic regression on the training partition. `survived` should be your response variable and `pclass`, `sex`, `age`, `sibsp`, and `parch` should be your response variables.

```
#student input
glm1 <- glm(survived ~ pclass+sex+age+sibsp+parch, family = binomial(link =
"logit"), data = data_train)
```

We would now like to test whether this classifier is *fair* across the sex subgroups. It was reported that women and children were prioritized on the life-boats and as a result survived the incident at a much higher rate. Let us see if our model is able to capture this fact.

Subset your test data into a male group and a female group. Then, use the `predict` function on the male testing group to come up with predicted probabilities of surviving the Titanic for each male in the testing set. Do the same for the female testing group.

```
#student input
male_test <- data_test[data_test$sex == "male",]
female_test <- data_test[data_test$sex != "male",]
male_probs <- predict(glm1, male_test)
female_probs <- predict(glm1, female_test)
```

Now recall that for this logistic *regression* to be a true classifier, we need to pair it with a decision boundary. Use an `if-else` statement to translate any predicted probability in the male group greater than 0.5 into `Yes` (as in Yes this individual is predicted to have survived). Likewise an predicted probability less than 0.5 should be translated into a `No`.

Do this for the female testing group as well, and then create a confusion matrix for each of the male and female test set predictions. You can use the `confusionMatrix` command as seen in class to expidite this process as well as provide you necessary metrics for the following questions.

```r
library(caret)

## Warning: package 'caret' was built under R version 4.3.2

## Loading required package: ggplot2

## Warning: package 'ggplot2' was built under R version 4.3.2

## Loading required package: lattice

#student input
fitted.results.m <- ifelse(male_probs > 0.5, "Yes", "No")
fitted.results.f <- ifelse(female_probs > 0.5, "Yes", "No")

matrix.m <- confusionMatrix(as.factor(fitted.results.m), male_test$survived,
positive = "Yes")
matrix.f <- confusionMatrix(as.factor(fitted.results.f),
female_test$survived, positive = "Yes")

matrix.m

## Confusion Matrix and Statistics
##
##           Reference
## Prediction No Yes
##        No  97  30
##        Yes  0   2
##
##                Accuracy : 0.7674
##                  95% CI : (0.6849, 0.8373)
##     No Information Rate : 0.7519
##     P-Value [Acc > NIR] : 0.3859
##
##                   Kappa : 0.0911
##
##  Mcnemar's Test P-Value : 1.192e-07
##
##             Sensitivity : 0.0625
##             Specificity : 1.0000
```

```
##            Pos Pred Value : 1.0000
##            Neg Pred Value : 0.7638
##                Prevalence : 0.2481
##            Detection Rate : 0.0155
##      Detection Prevalence : 0.0155
##         Balanced Accuracy : 0.5312
##
##          'Positive' Class : Yes
##
```

matrix.f

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction No Yes
##        No  10   3
##        Yes  9  58
##
##                Accuracy : 0.85
##                  95% CI : (0.7526, 0.92)
##     No Information Rate : 0.7625
##     P-Value [Acc > NIR] : 0.03896
##
##                   Kappa : 0.5353
##
##  Mcnemar's Test P-Value : 0.14891
##
##             Sensitivity : 0.9508
##             Specificity : 0.5263
##          Pos Pred Value : 0.8657
##          Neg Pred Value : 0.7692
##              Prevalence : 0.7625
##          Detection Rate : 0.7250
##    Detection Prevalence : 0.8375
##       Balanced Accuracy : 0.7386
##
##          'Positive' Class : Yes
##
```

We can see that indeed, at least within the testing groups, women did seem to survive at a higher proportion than men (24.8% to 76.3% in the testing set). Print a summary of your trained model and interpret one of the fitted coefficients in light of the above disparity.

```
#student input
summary(glm1)
```

```
##
## Call:
## glm(formula = survived ~ pclass + sex + age + sibsp + parch,
##     family = binomial(link = "logit"), data = data_train)
##
## Coefficients:
##                 Estimate Std. Error z value Pr(>|z|)
## (Intercept)     3.903165   0.409280   9.537  < 2e-16 ***
## pclassMiddle   -1.291506   0.257421  -5.017 5.25e-07 ***
## pclassLower    -2.404084   0.262022  -9.175  < 2e-16 ***
## sexmale        -2.684206   0.200130 -13.412  < 2e-16 ***
## age            -0.036776   0.007494  -4.907 9.24e-07 ***
## sibsp          -0.395584   0.118587  -3.336  0.00085 ***
## parch           0.032494   0.111916   0.290  0.77155
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1121.27  on 833  degrees of freedom
## Residual deviance:  757.87  on 827  degrees of freedom
## AIC: 771.87
##
## Number of Fisher Scoring iterations: 5
```

*Student Input*

In light of the disparity of male and female survival rates, it is only more fitting to see sex accounted for within my model. The coefficient, -2.684206, suggests that for every observation which is male (male = 1), the log odds of their survival decreased by 2.684206.

Now let's see if our model is *fair* across this explanatory variable. Calculate five measures (as defined in class) in this question: the Overall accuracy rate ratio between females and males, the disparate impact between females and males, the statistical parity between females and males, and the predictive equality as well as equal opportunity between females and males (collectively these last two comprise equalized odds). Set a reasonable $\epsilon$ each time and then comment on which (if any) of these five criteria are met.

```
#Student Input
epsilon = 0.8
#this is based upon legal precedent, rather than any more substantial
reasoning for selecting a value of 0.8. It is likely that, although the
conclusions may differ with a different value, it would be difficult to argue
that any given epsilon would not suffice or be reasonable, especially within
typical statistical boundaries (such as 0.5 to 0.95, which are often used as
decision boundaries).
#Accuracy Rate Ratio
```

```r
misClasificError.m <- mean(fitted.results.m != male_test$survived)
misClasificError.f <- mean(fitted.results.f != female_test$survived)
accuracyRateRatio <- misClasificError.f / misClasificError.m
accuracyRateRatio
```

```
## [1] 0.645
```

```r
accuracyRateRatio > epsilon
```

```
## [1] FALSE
```

```r
#Disparate Impact
dispImpact <- mean(fitted.results.m == "Yes") / mean(fitted.results.f ==
"Yes")
dispImpact
```

```
## [1] 0.01851209
```

```r
dispImpact > 1 - epsilon
```

```
## [1] FALSE
```

```r
#Statistical Parity
statParity <- abs(mean(fitted.results.m == "Yes") - mean(fitted.results.f ==
"Yes"))
statParity
```

```
## [1] 0.8219961
```

```r
statParity < epsilon
```

```
## [1] FALSE
```

```r
#Predictive Equality
predEquality <- min(abs(mean((fitted.results.m == "Yes" & male_test$survived
== "No")) - mean((fitted.results.f == "Yes" & female_test$survived ==
"No"))), abs(mean((fitted.results.m == "Yes" & male_test$survived == "Yes"))
- mean((fitted.results.f == "Yes" &  female_test$survived == "Yes"))))
predEquality
```

```
## [1] 0.1125
```

```r
predEquality < epsilon
```

```
## [1] TRUE
```

```r
#Equal Opportunity
equalOpp <- abs(mean(fitted.results.f == "Yes" & female_test$survived ==
"Yes") - mean(fitted.results.m == "Yes" & male_test$survived == "Yes"))
equalOpp
```

```
## [1] 0.7094961
```

```r
equalOpp < epsilon
```

```
## [1] TRUE
```

*Student Input.*

It is always important for us to interpret our results in light of the original data and the context of the analysis. In this case, it is relevant that we are analyzing a historical event post-facto and any disparities across demographics identified are unlikely to be replicated. So even though our model fails numerous of the statistical fairness criteria, I would argue we need not worry that our model could be misused to perpetuate discrimination in the future. After all, this model is likely not being used to prescribe a preferred method of treatment in the future.

Even so, provide a *philosophical* notion of justice or fairness that may have motivated the Titanic survivors to act as they did. Spell out what this philosophical notion or principle entails?

*Student Input*

The notion, at least from a virtue ethics persepective, was to save the most virtuous passengers. This could appeal to virtues such as tenacity or caring, as women were often seen in a much more traditional light at the time, and may have been seen as the more deserving of the lifeboats/available resources. Since virtue ethics encourages people to pursue the most virtuous efforts and to preserve virtue as the medium between extremes, it makes sense that the more virtuous passengers would receive more efforts and affects with which to save themselves.