

Summary of Credit Scoring Assignment

University of Connecticut FNCE 5352

Matt McDonald

April 30, 2024

Summary

Eight submissions were collected (including one submission by Professor McDonald). Although each submission should have included a prediction for every record (37,500) in the “ConsumerCred-train.csv” file, one included predictions for fewer than that. In order to make the files scoreable, Professor McDonald assigned the average score within the file to the rows that were missing a prediction. This document summarizes the results of scoring these submissions, and includes an analysis of the distribution of the results.

Processing

The following code evaluates the submissions against the *solution* file, which is coded with a Default indicator for each of the records in the test file.

```
library(tidyverse)
library(pROC)

#Load in the solution file
solutionfile <- here::here('ConsumerCredit', 'solution.csv')
solution <- read_csv(solutionfile)

#All submissions are found in the "grading" folder
submissionsfolder <- here::here('ConsumerCredit', 'grading')

#create the "submissions" tibble
#file is the file name
submissions <- tibble(file=list.files(path=submissionsfolder, pattern='.csv', recursive = TRUE))
#csv is a list column containing a tibble with the submission from the team
submissions <- submissions %>%
  mutate(csv=map(file, ~ read_csv(paste(submissionsfolder, .x, sep='//'))))

#create a column indicating which column in the solution
#contains the score
submissions$scorecol <- c(2,2,2,2,15,2,2,2)

submissions %>% mutate(avg_score=map2_dbl(csv,
                                          scorecol,
                                          ~mean(pull(.x[,.y]),
```

```

na.rm=TRUE)))

submissions %>% mutate(map2_int(csv,
                                scorecol,
                                ~ sum(is.na(.x[,.y]))))

#one of the submissions has 7400 NA values!!!
WZY <- submissions$csv[[7]]
sum(is.na(WZY$probability))

#this function applies the "average" score for any ids that had a missing prediction
missing_score_remediation <- function(csv){
  out <- tibble(id=1:37500)
  colnames(csv) <- c('id', 'pred')
  out <- out %>%
    left_join(csv)
  avg_score <- mean(csv$pred, na.rm=TRUE)
  out[is.na(out$pred), 'pred'] <- avg_score
  out
}

submissions$csv[[7]] <- missing_score_remediation(WZY)

#this function scores the submission
getAUC <- function(csv, colnum=2){
  out <- 0

  if (nrow(csv) == nrow(solution)) {
    rocobj <- roc(
      response = solution$SeriousDlqin2yrs,
      predictor = pull(csv[,colnum])
    )
    out <- auc(rocobj)
  }
  out
}

#use my getAUC function to score the data contained in column 'csv'

submissions <- submissions %>%
  mutate(AUC = map2_dbl(csv, scorecol, getAUC))

View(submissions %>% select(-csv) %>% arrange(AUC))

ggplot(filter(submissions, AUC > 0), aes(x=AUC)) + geom_density() +geom_rug()

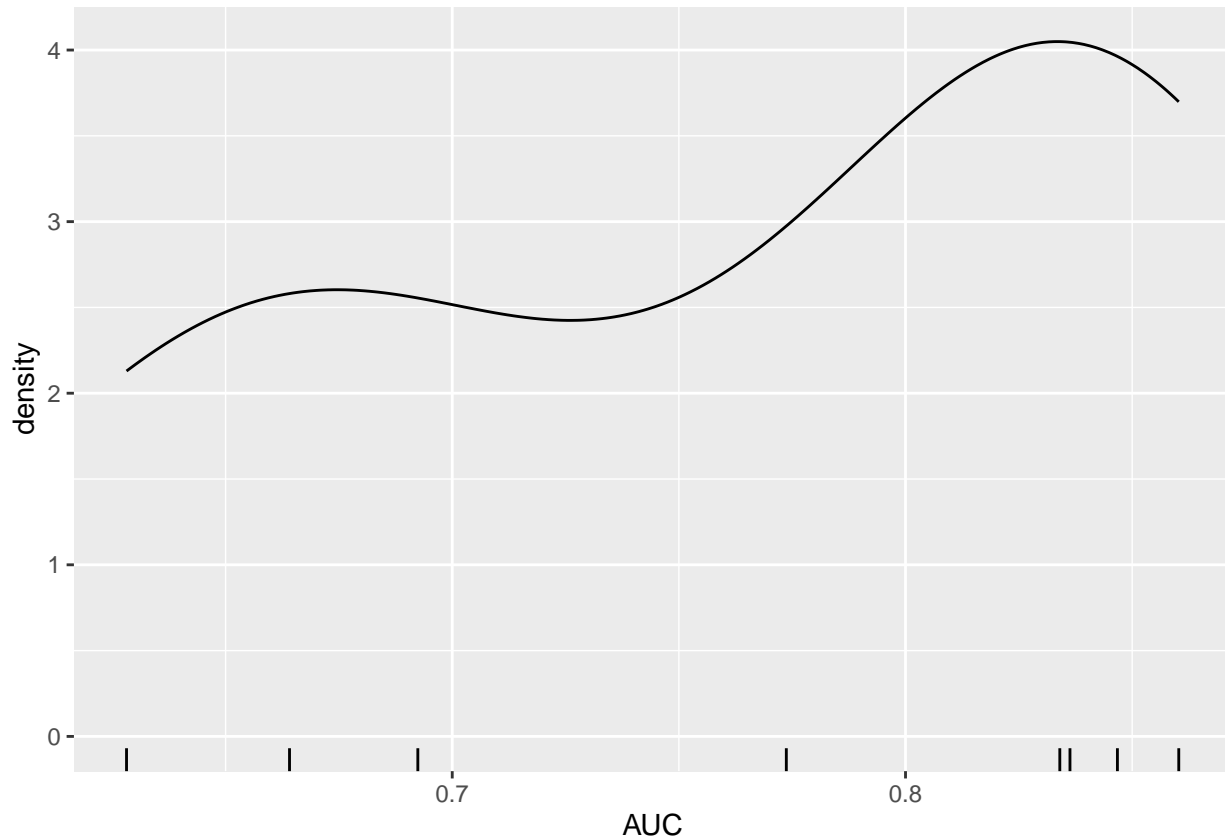
saveRDS(select(submissions, -csv), file=here::here('ConsumerCredit', 'submissions.RDS'))

```

Distribution of scores

The following graph shows the distribution of scores. All submissions have AUC values greater than 50%, which indicates that every conforming submission had some predictive power. An AUC value of 50% indicates a completely random model. During class, we created a model using only the “Age” column as a predictor and got an AUC value of 63.8%.

```
submissions <- readRDS(here::here('ConsumerCredit', 'submissions.RDS'))
ggplot(submissions, aes(x=AUC)) + geom_density() + geom_rug()
```



Notes

The data was taken from the Kaggle competition “Give Me Some Credit”, which can be found at the following link: <https://www.kaggle.com/c/GiveMeSomeCredit>

This Kaggle competition contains a lot of discussion about approaches that can be used to process the data and improve performance. As it stands, the best score obtained by the class is close to the winning score of 0.8695. However, our highest score would not have cracked the top 100 scores in this competition.

The code Professor McDonald used to generate his results can be found at https://github.com/mattmcd71/fnce5352_spring2023/tree/main/ConsumerCredit under the file names *modelingexample.R*.

This summary was written using RMarkdown, which is a useful tool in RStudio for communicating results. A helpful CheatSheet can be found at <https://rmarkdown.rstudio.com/>