**Final Report Jackson Bandow** 3/23/2021 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) library(plotly) ## Attaching package: 'plotly' ## The following object is masked from 'package:ggplot2': last plot ## The following object is masked from 'package:stats': filter ## The following object is masked from 'package:graphics': layout library(data.table) ## Attaching package: 'data.table' ## The following objects are masked from 'package:dplyr': between, first, last library(Lahman) library(stringr) library(car) ## Loading required package: carData ## Attaching package: 'carData' ## The following object is masked from 'package:Lahman': Salaries ## Attaching package: 'car' ## The following object is masked from 'package:dplyr': recode library(pROC) ## Type 'citation("pROC")' for a citation. ## Attaching package: 'pROC' ## The following objects are masked from 'package:stats': cov, smooth, var

## ##

## ##

## ##

## ##

##

##

## ##

##

## ##

get\_prob <- function(input, coefs){</pre> inp <- t(as.matrix(cbind(rep(1, dim(input)[1]),input)))</pre> cf <- t(as.matrix(coefs))</pre> prob\_new <- (exp(cf %\*% inp)) / (1 + exp(cf %\*% inp))</pre> names(prob\_new) <- NULL</pre> return(as.vector(prob\_new)) mc rate <- function(probs, actual) {</pre> mean(round(probs, 0) != actual) low out <- function(col){</pre> stats <- summary(col)</pre> iqr <- stats[5] - stats[2]</pre> res <- stats[2] - 1.5\*iqr names(res) <- NULL</pre> return(res) standardize <- function(col) {</pre> return((col - mean(col)) / sd(col)) baseball <- fread("baseball.csv")</pre> baseball\$careerLength <- as.numeric(as.Date(baseball\$finalGame) - as.Date(baseball\$debut)) / 365.25</pre> baseball\$debutDecade <- round(as.numeric(format(baseball\$debut, "%Y")), -1)</pre> baseball\$retireRank <- floor((as.numeric(format(baseball\$finalGame, "%Y")) - 4) / 10) - 191 baseball\$FP <- (baseball\$PO + baseball\$A) / (baseball\$PO + baseball\$A + baseball\$E)</pre> HOF batters <- baseball %>% filter(inducted & position != "P" & votedBy != "Negro League" & str sub(playerID, -2, -1) != "99" & finalGame < as.Date("2004-01-01") & finalGame > as.Date("1924-01-01")) batters <- baseball %>% filter(position != "P" & votedBy != "Negro League" & str sub(playerID, -2, -1) != "99" & finalGame < as.Date("2004-01-01") & finalGame > as.Date("1924-01-01")) ggplot(batters, aes(retireRank, AB, color = inducted)) + geom jitter() + theme bw() + scale color manual(values = c("yellow", "green")) + labs(title = "Number of At-Bats by Retirement Decade", x = "Retirement Decade", y = "ABs" , color = "Inducted to HOF?") Number of At-Bats by Retirement Decade 10000 -Inducted to HOF? FALSE TRUE 8 Retirement Decade paste0("Composition of HOF batters in data: ", round(mean(batters\$inducted)\*100, 1), "%.") ## [1] "Composition of HOF batters in data: 2.2%." # Limit data by ABs by retirement decade new\_batters <- data.frame()</pre> batters <- batters %>% filter(rWAR >= 10) for (i in unique(batters\$retireRank)) { temp <- batters %>% filter(retireRank == i) new batters <- rbind(new batters, temp[order(temp\$AB)[(dim(temp)[1]-59):dim(temp)[1]],])</pre> ggplot(new\_batters, aes(finalGame, AB, color = inducted)) + geom\_jitter() + theme\_bw() + scale\_color\_manual(value s = c("yellow", "green")) + labs(title = "Number of At-Bats by Retirement Decade", x = "Retirement Decade", y = " ABs", color = "Inducted to HOF?") Number of At-Bats by Retirement Decade 12000 -Inducted to HOF? FALSE TRUE 1940 1960 1980 2000 Retirement Decade paste0("Composition of HOF batters in data: ", round(mean(new\_batters\$inducted)\*100, 1), "%.") ## [1] "Composition of HOF batters in data: 24%." # Make k-fold indeces rand\_ind <- sample(1:dim(new\_batters)[1], dim(new\_batters)[1], replace = FALSE)</pre> sets <- NULL fold\_size <- dim(new\_batters)[1] / 5</pre> for (i in 1:5) { sets <- append(sets, list(rand\_ind[((i-1)\*fold\_size + 1):(i\*fold\_size)]))</pre> # Model 1: WAR mcr1 <- NULL auc1 <- NULL for (i in 1:length(sets)) { mod <- glm(data = new\_batters[!sets[[i]],], formula = inducted ~ rWAR, family = "binomial")</pre> mcr1 <- c(mcr1, mc\_rate(get\_prob(new\_batters[sets[[i]],c("rWAR")], mod\$coefficients), new\_batters[sets[[i]],]\$i</pre> nducted)) auc1 <- c(auc1, roc(new\_batters[sets[[i]],]\$inducted, get\_prob(new\_batters[sets[[i]],c("rWAR")], mod\$coefficien</pre> ts))\$auc) ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> print(paste0("Model 1 -- Misclassification Rate: ", round(mean(mcr1)\*100, 1), "%; Area Under the Curve: ", round( mean(auc1), 3))) ## [1] "Model 1 -- Misclassification Rate: 13.8%; Area Under the Curve: 0.915" # Model 2: Standard Statistics mcr2 <- NULL auc2 <- NULL for (i in 1:length(sets)) { mod <- glm(data = new batters[!sets[[i]],], formula = inducted ~ G.B + AB + BB + H + HR + R + RBI + SB + FP + O BP + SLG + BA, family = "binomial") mcr2 <- c(mcr2, mc\_rate(get\_prob(new\_batters[sets[[i]],c("G.B", "AB", "BB", "H", "HR", "R", "RBI", "SB", "FP", "OBP", "SLG", "BA")], mod\$coefficients), new batters[sets[[i]],]\$inducted)) auc2 <- c(auc2, roc(new\_batters[sets[[i]],]\$inducted, get\_prob(new\_batters[sets[[i]],c("G.B", "AB", "BB", "H",</pre> "HR", "R", "RBI", "SB", "FP", "OBP", "SLG", "BA")], mod\$coefficients))\$auc) ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE

mcr3 <- NULL auc3 <- NULL for (i in 1:length(sets)) { mod <- glm(data = snew\_batters[!sets[[i]],], formula = inducted ~ rWAR, family = "binomial")</pre> mcr3 <- c(mcr3, mc\_rate(get\_prob(snew\_batters[sets[[i]],c("rWAR")], mod\$coefficients), snew\_batters[sets[[i]],]</pre> \$inducted)) auc3 <- c(auc3, roc(snew batters[sets[[i]],]\$inducted, get\_prob(snew\_batters[sets[[i]],c("rWAR")], mod\$coeffici</pre> ents))\$auc) ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> print(paste0("Standardized Model 1 -- Misclassification Rate: ", round(mean(mcr3)\*100, 1), "%; Area Under the Cur ve: ", round(mean(auc3), 3))) ## [1] "Standardized Model 1 -- Misclassification Rate: 10.6%; Area Under the Curve: 0.937" # Standardized Model 2: Standard Statistics mcr4 <- NULL auc4 <- NULL for (i in 1:length(sets)) { mod <- glm(data = snew\_batters[!sets[[i]],], formula = inducted ~ G.B + AB + BB + H + HR + R + RBI + SB + FP +</pre> OBP + SLG + BA, family = "binomial") mcr4 <- c(mcr4, mc rate(get prob(snew batters[sets[[i]],c("G.B", "AB", "BB", "H", "HR", "R", "RBI", "SB", "FP", "OBP", "SLG", "BA")], mod\$coefficients), snew batters[sets[[i]],]\$inducted)) auc4 <- c(auc4, roc(snew\_batters[sets[[i]],]\$inducted, get\_prob(snew\_batters[sets[[i]],c("G.B", "AB", "BB", "H"</pre> , "HR", "R", "RBI", "SB", "FP", "OBP", "SLG", "BA")], mod\$coefficients))\$auc) ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> print(paste0("Standardized Model 2 -- Misclassification Rate: ", round(mean(mcr4)\*100, 1), "%; Area Under the Cur ve: ", round(mean(auc4), 3))) ## [1] "Standardized Model 2 -- Misclassification Rate: 14.2%; Area Under the Curve: 0.892"

mod <- glm(data = snew\_batters[!sets[[i]],], formula = inducted ~ retireRank\*rWAR, family = "binomial")</pre>

mcr5 <- c(mcr5, mc\_rate(get\_prob(temp, mod\$coefficients), snew\_batters[sets[[i]],]\$inducted))</pre> auc5 <- c(auc5, roc(snew\_batters[sets[[i]],]\$inducted, get\_prob(temp, mod\$coefficients))\$auc)</pre>

temp <- cbind(snew batters[sets[[i]],c("retireRank", "rWAR")], snew batters[sets[[i]],]\$retireRank\*snew batters</pre>

print(paste0("Standardized Model 1 with Time Interactions -- Misclassification Rate: ", round(mean(mcr5)\*100, 1),

## [1] "Standardized Model 1 with Time Interactions -- Misclassification Rate: 11.2%; Area Under the Curve: 0.944

mod <- glm(data = snew\_batters[!sets[[i]],], formula = inducted ~ retireRank\*G.B + retireRank\*AB + retireRank\*B</pre> B + retireRank\*H + retireRank\*HR + retireRank\*R + retireRank\*RBI + retireRank\*SB + retireRank\*FP + retireRank\*OBP

temp <- cbind(snew\_batters[sets[[i]],c("retireRank", "G.B", "AB", "BB", "H", "HR", "R", "RBI", "SB", "FP", "OBP ", "SLG", "BA")], snew\_batters[sets[[i]],c("G.B", "AB", "BB", "H", "HR", "R", "RBI", "SB", "FP", "OBP", "SLG", "B

mcr6 <- c(mcr6, mc\_rate(get\_prob(temp, mod\$coefficients), snew\_batters[sets[[i]],]\$inducted))</pre> auc6 <- c(auc6, roc(snew\_batters[sets[[i]],]\$inducted, get\_prob(temp, mod\$coefficients))\$auc)</pre>

print(paste0("Model 2 -- Misclassification Rate: ", round(mean(mcr2)\*100, 1), "%; Area Under the Curve: ", round(

snew\_batters <- new\_batters %>% select(rWAR, AB, H, HR, BB, G.B, BA, R, RBI, SB, OBP, SLG, FP) %>% log() %>% cbin

snew\_batters[snew\_batters\$retireRank == i,c("rWAR", "AB", "H", "HR", "BB", "G.B", "BA", "R", "RBI", "SB", "OBP" , "SLG", "FP")] <- data.frame(apply(snew\_batters[snew\_batters\$retireRank == i,c("rWAR", "AB", "H", "HR", "BB", "G

## [1] "Model 2 -- Misclassification Rate: 12.5%; Area Under the Curve: 0.899"

## Setting direction: controls < cases</pre>

## Setting direction: controls < cases</pre>

## Setting direction: controls < cases</pre>

# Standardize data by retirement decade

# Standardized Model 1: WAR

d(new batters[,c("inducted", "steriods", "retireRank")])

.B", "BA", "R", "RBI", "SB", "OBP", "SLG", "FP")], 2, standardize))

for (i in sort(unique(snew batters\$retireRank))) {

# Standardized Model 1 with Time Interactions: WAR

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE

"%; Area Under the Curve: ", round(mean(auc5), 3)))

+ retireRank\*SLG + retireRank\*BA, family = "binomial")

A")]\*snew\_batters[sets[[i]],]\$retireRank)

# Standardized Model 2 with Time Interactions: Standard Statistics

## Setting direction: controls < cases</pre>

mcr6 <- NULL auc6 <- NULL

## ##

##

margin

combine

for (i in 1:length(sets)) {

tters[sets[[i]],]\$inducted))

, round(mean(auc1), 3)))

mcr2 <- NULL auc2 <- NULL

# Model 1: WAR mcr1 <- NULL auc1 <- NULL

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

for (i in 1:length(sets)) {

mcr5 <- NULL auc5 <- NULL

[sets[[i]],]\$rWAR)

for (i in 1:length(sets)) {

mean(auc2), 3)))

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls < cases</pre> print(paste0("Standardized Model 2 with Time Interactions -- Misclassification Rate: ", round(mean(mcr6)\*100, 1), "%; Area Under the Curve: ", round(mean(auc6), 3))) ## [1] "Standardized Model 2 with Time Interactions -- Misclassification Rate: 15.8%; Area Under the Curve: 0.881 library(randomForest) ## randomForest 4.6-14 ## Type rfNews() to see new features/changes/bug fixes. ## Attaching package: 'randomForest' ## The following object is masked from 'package:ggplot2':

auc1 <- c(auc1, roc(new\_batters[sets[[i]],]\$inducted, predict(mod, new\_batters[sets[[i]],c("rWAR", "retireRank"</pre> )], type = "prob")[,1])\$auc) ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls > cases ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls > cases ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls > cases ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls > cases ## Setting levels: control = FALSE, case = TRUE ## Setting direction: controls > cases

print(paste0("Random Forest 1 -- Misclassification Rate: ", round(mean(mcr1)\*100, 1), "%; Area Under the Curve: "

print(paste0("Random Forest 2 -- Misclassification Rate: ", round(mean(mcr2)\*100, 1), "%; Area Under the Curve: "

## [1] "Random Forest 2 -- Misclassification Rate: 9.8%; Area Under the Curve: 0.924"

## [1] "Random Forest 1 -- Misclassification Rate: 10.8%; Area Under the Curve: 0.924"

mcr1 <- c(mcr1, mc\_rate(predict(mod, new\_batters[sets[[i]],c("rWAR", "retireRank")], type = "prob")[,2], new\_ba</pre>

mod <- randomForest(factor(inducted) ~ rWAR + retireRank, data = new\_batters[!sets[[i]],])</pre>

## Setting direction: controls > cases

, round(mean(auc2), 3)))

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE

## Setting levels: control = FALSE, case = TRUE