

Stats-Thesis2.Rmd

Jackson Bandow

3/9/2021

```
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)
```

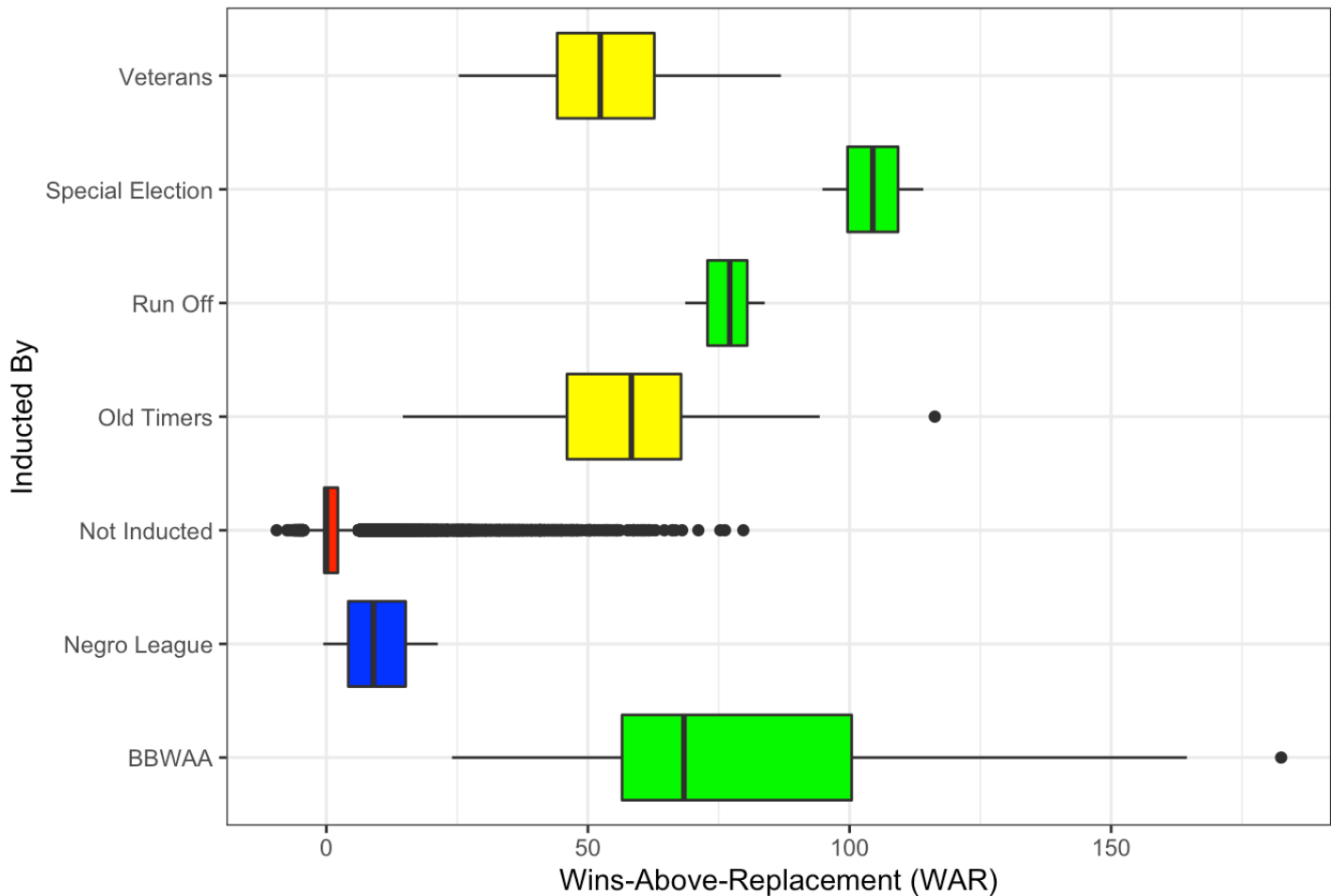
```
baseball <- fread("baseball.csv")
```

```
baseball[baseball$votedBy == "", "votedBy"] <- "Not Inducted"
baseball$votedBy2 <- baseball$votedBy
baseball[baseball$votedBy2 == "Old Timers", "votedBy2"] <- "Veterans"
baseball[baseball$votedBy2 == "Run Off", "votedBy2"] <- "BBWAA"
baseball[baseball$votedBy2 == "Special Election", "votedBy2"] <- "BBWAA"
baseball$preHOFera <- baseball$finalGame < as.Date("1922-01-01")
#Create a training and testing set
bb_train <- baseball %>% filter(finalGame < as.Date("1991-01-01"))
bb_test <- baseball %>% filter(finalGame > as.Date("1991-01-01"))
```

I've used the era from 1871 through 1990 as my training dataset because it is approximately 80% of the data available. Also, I avoid most players that are still in consideration for the Hall of Fame, so I am unlikely to falsely assume that they won't be inducted because they have not yet been.

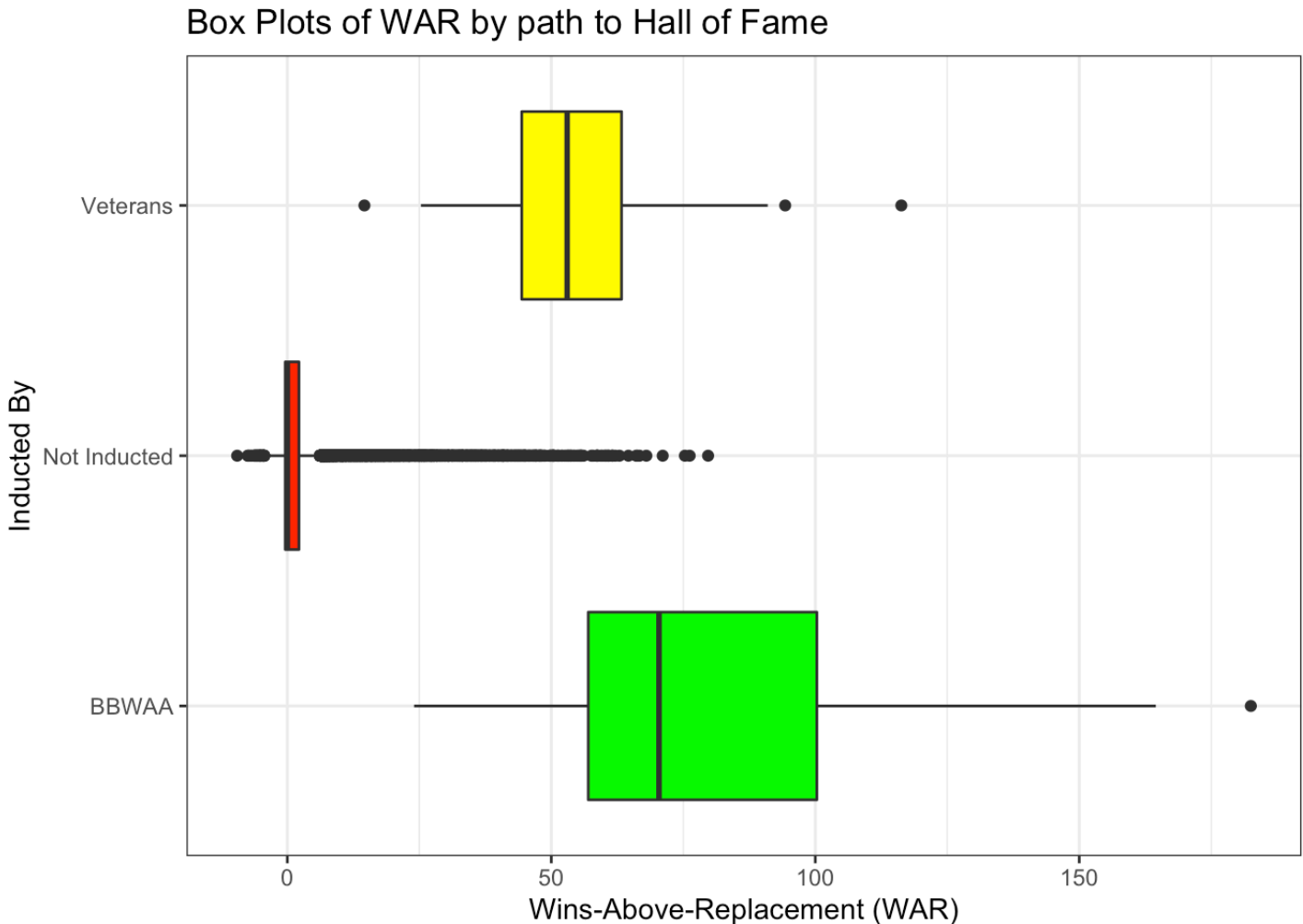
```
ggplot(bb_train, aes(rWAR, votedBy)) + geom_boxplot(fill = c("green", "blue", "red",
"yellow", "green", "green", "yellow")) + theme_bw() + labs(title = "Box Plots of WAR by
path to Hall of Fame", y = "Inducted By", x = "Wins-Above-Replacement (WAR)")
```

Box Plots of WAR by path to Hall of Fame



Although there are a lot of committees for induction to the Hall of Fame, they can be condensed. The BBWAA is the organization of sportswriters that vote in players that have retired 5-15 years prior. On two occasions (Lou Gehrig and Roberto Clemente), there have been a special election, which can be included under the BBWAA umbrella. Additionally, there is occasionally a run off where only the top voting result is inducted. This process can also be included under BBWAA. Generally after a player is passed on by the BBWAA (this rule has changed historically), they are still eligible for election under the Veterans committee. The first election for the Hall of Fame was in 1936, so a subcommittee of the Veterans committee, the Old Timers committee, was formed to fill in the gaps of older players. These categories can also be combined. Finally, although the Negro League is also a subcommittee of Veterans, I will consider them separately. Because most of them did not play in the MLB, most of them do not have a WAR statistic. Those that joined the MLB during the integration era have a WAR, but joined the league part-way through their career, so their WAR only reflects a fraction of their career. For those reasons, I will remove the Negro League inductees for now.

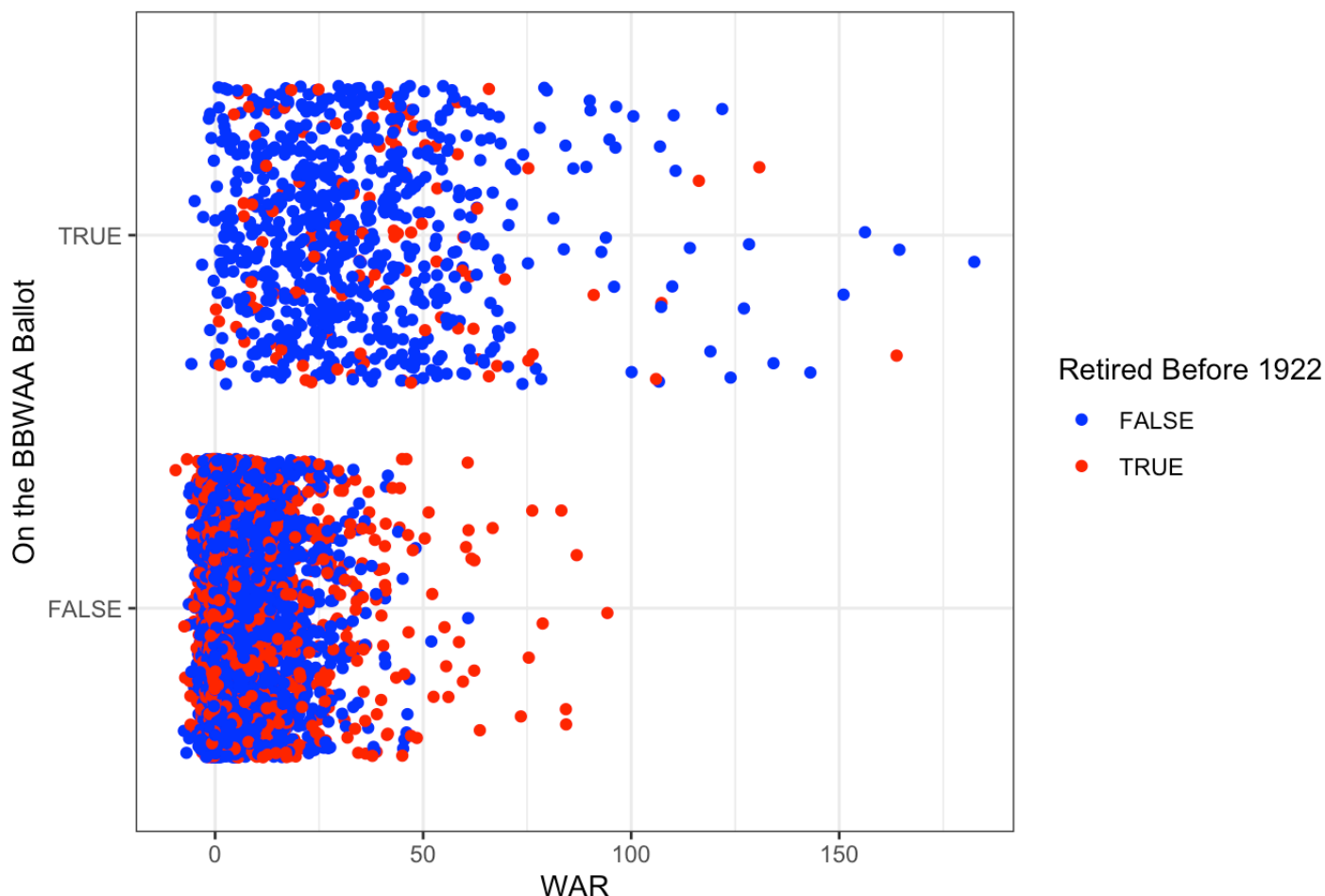
```
# Excludes Negro League players
bb_train <- bb_train %>% filter(!is.na(rWAR) & votedBy != "Negro League")
bb_test <- bb_test %>% filter(!is.na(rWAR) & votedBy != "Negro League")
ggplot(bb_train, aes(rWAR, votedBy2)) + geom_boxplot(fill = c("green", "red", "yellow")) + theme_bw() + labs(title = "Box Plots of WAR by path to Hall of Fame", y = "Inducted By", x = "Wins-Above-Replacement (WAR)")
```



WAR seems to generally be a good indicator of Hall of Fame status, but there is clearly more to the story. There are definitely some outliers on either side that need some explaining. Perhaps they deserve to be in the Hall of Fame or perhaps there is more to a player's Hall of Fame candidacy than their value, described by WAR. The first step to induction is to be nominated for the BBWAA Hall of Fame ballot (this is determined by a subcommittee of the association nowadays).

```
ggplot(bb_train, aes(rWAR, onBallot, color = preHOFera)) + geom_jitter() + theme_bw()
+ labs(title = "WAR by Appearance on BBWAA Ballot", x = "WAR", y = "On the BBWAA Ball
ot", color = "Retired Before 1922") + scale_color_manual(values = c("blue", "red"))
```

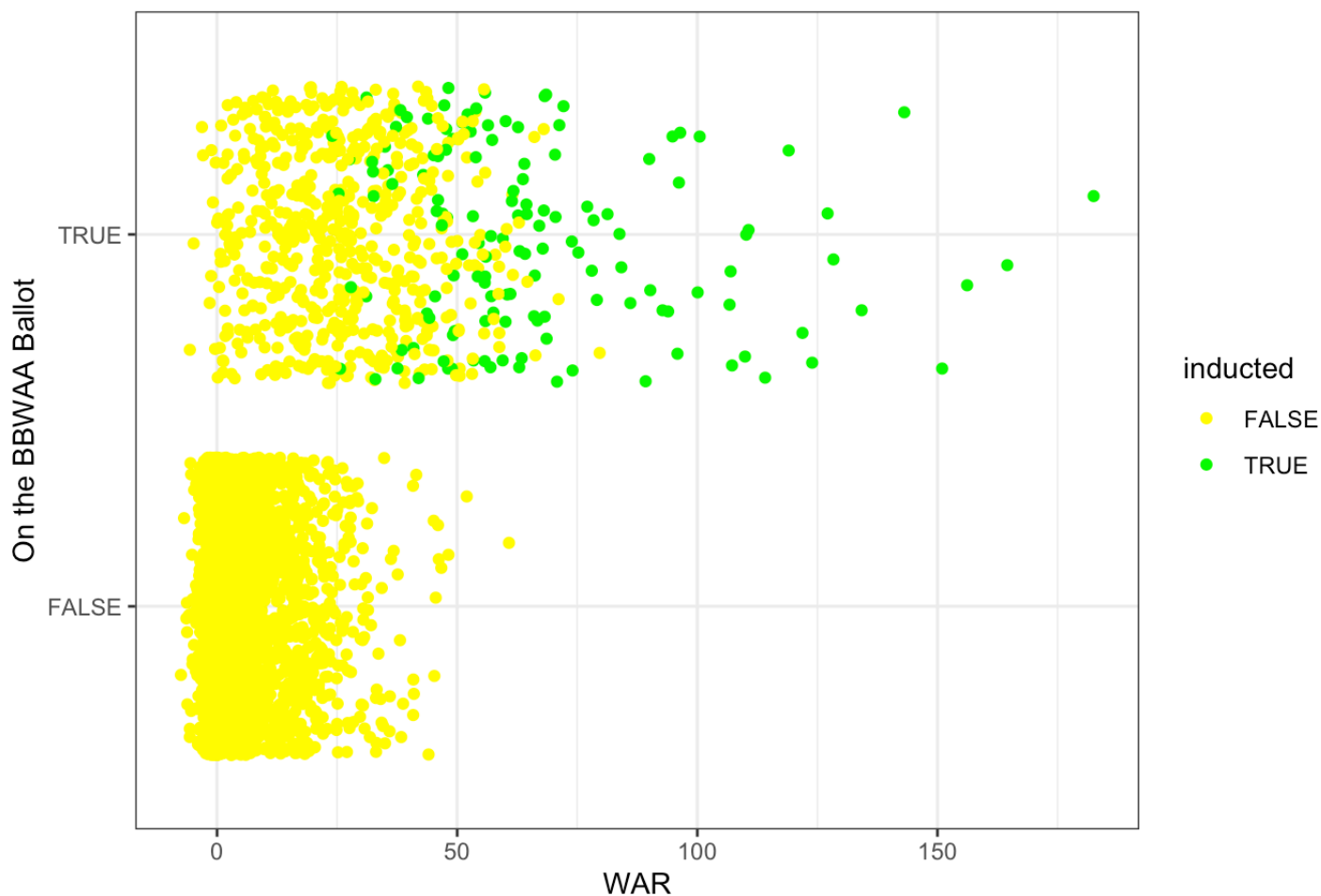
WAR by Appearance on BBWAA Ballot



The blue points represent players that retired before the 1922 season. These players would be ineligible for the first HOF ballot in 1936 under today's standards. Some players from this era were considered, as the rules were not solidified until much later. However, the Old Timers committee was created to address the issue of underrepresentation among the older era, so that would explain why the majority of the higher outliers for WAR that were not included on any BBWAA ballot are from this older era.

```
bb_train.HOFera <- bb_train %>% filter(!preHOFera)
ggplot(bb_train.HOFera, aes(rWAR, onBallot, color = inducted)) + geom_jitter() + theme_bw() + labs(title = "WAR by Appearance on BBWAA Ballot", x = "WAR", y = "On the BBWAA Ballot") + scale_color_manual(values = c("yellow", "green"))
```

WAR by Appearance on BBWAA Ballot



With the pre Hall of Fame players eliminated, the plot is slightly improved. There is a pretty clear difference between those eventually inducted in the Hall of Fame and those that were not included on the ballot. However, there are still a significant amount of mediocre to just bad players, considered by WAR that were included on the ballot. Why was Tommy Thevenow included with a career 5.7 WAR? Let's see if a logistic regression can pick up on interesting trends.

```
mod1 <- glm(onBallot ~ rWAR, data = bb_train.HOFera, family = "binomial")
summary(mod1)
```

```
##
## Call:
## glm(formula = onBallot ~ rWAR, family = "binomial", data = bb_train.HOFera)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5704  -0.1892  -0.1690  -0.1605   3.2378
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -4.241419   0.092316  -45.95  <2e-16 ***
## rWAR         0.174565   0.005216   33.47  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 4808.9  on 7978  degrees of freedom
## Residual deviance: 2170.2  on 7977  degrees of freedom
## AIC: 2174.2
##
## Number of Fisher Scoring iterations: 7
```

```
b0 <- mod1$coefficients[1]
b1 <- mod1$coefficients[2]

get_prob_mod1 <- function(input){
  prob_new <- (exp(b0 + b1 * input)) / (1 + exp(b0 + b1 * input))
  names(prob_new) <- NULL
  return(prob_new)
}

plot_prob_mod1 <- function(input) {
  get_prob_mod1(input) + 1
}

(misclassification_rate <- sum(round(get_prob_mod1(bb_train.HOFera$rWAR), 0) != bb_train.HOFera$onBallot) / dim(bb_train.HOFera)[1])
```

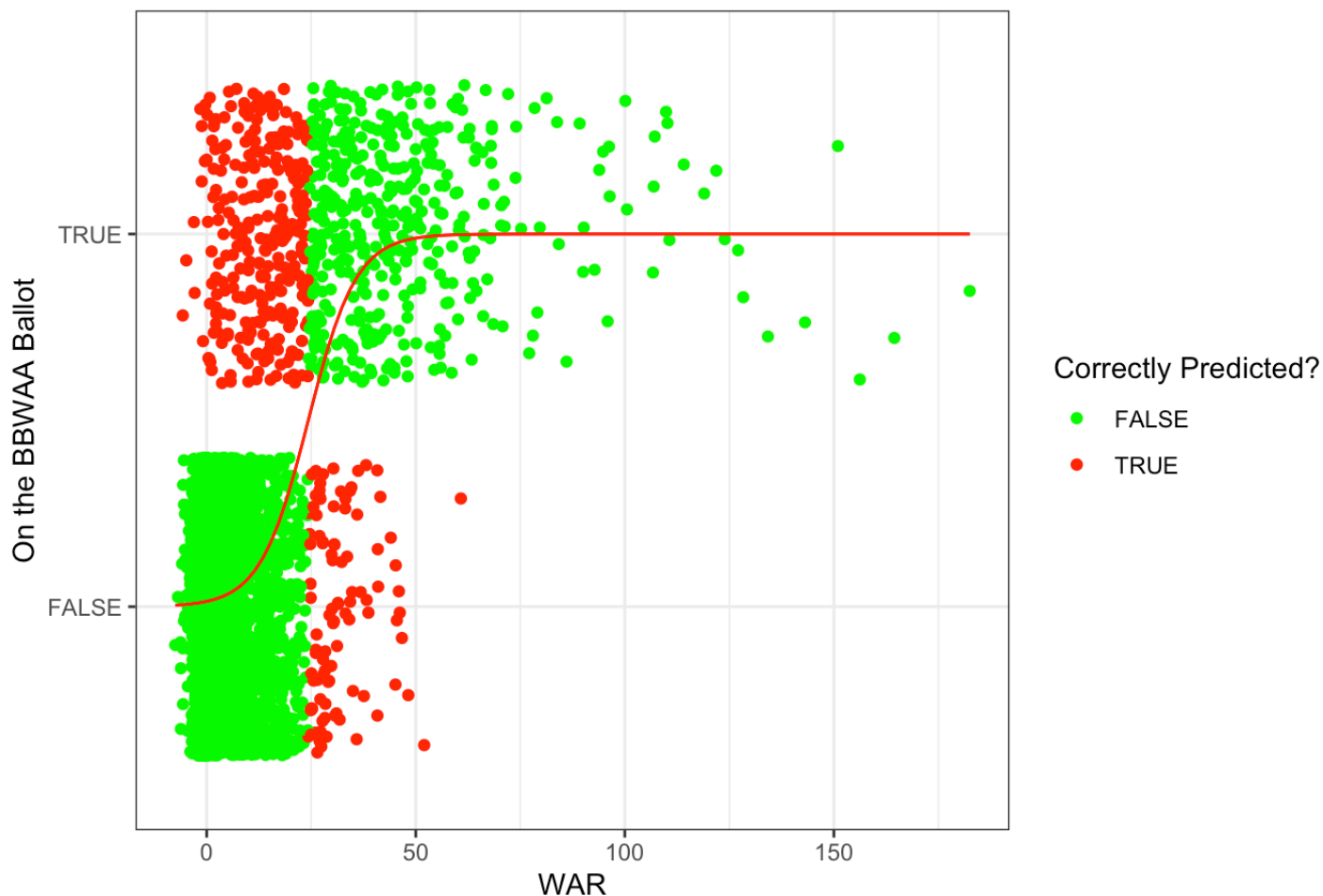
```
## [1] 0.048001
```

```
bb_train.HOFera$pred1_cor <- round(get_prob_mod1(bb_train.HOFera$rWAR), 0) != bb_train.HOFera$onBallot
```

```
ggplot(bb_train.HOFera, aes(rWAR, onBallot, color = bb_train.HOFera$pred1_cor)) + geom_jitter() + theme_bw() + labs(title = "WAR by Appearance on BBWAA Ballot", x = "WAR", y = "On the BBWAA Ballot", color = "Correctly Predicted?") + stat_function(fun = plot_prob_mod1, color = "red") + scale_color_manual(values = c("green", "red"))
```

```
## Warning: Multiple drawing groups in `geom_function()`. Did you use the correct  
## `group`, `colour`, or `fill` aesthetics?
```

WAR by Appearance on BBWAA Ballot



The logistic regression is a decent predictor of ballot inclusion, but it could be improved. Let's add in features and observe changes in predictive ability.


```

# Need to figure out how to deal with missing data. All WAR data is present but need
to find a way to fill in data for those that do not have any. Discuss more on Friday

# mod2 <- glm(onBallot ~ rWAR.B + AB + H + HR + BA + R + RBI + SB + OBP + SLG + OPS +
rWAR.P + W + L + ERA + G + GS + SV + IP + SO + WHIP, data = bb_train.HOFera, family =
"binomial")
# summary(mod2)
# betas.2 <- mod2$coefficients
#
# get_prob_mod2 <- function(input){
#   prob_new <- (exp(betas.2[1] + betas.2[2:length(betas.2)] * input)) / (1 + exp(bet
as.2[1] + betas.2[2:length(betas.2)] * input))
#   names(prob_new) <- NULL
#   return(prob_new)
# }
#
# plot_prob_mod2 <- function(input) {
#   get_prob_mod1(input) + 1
# }
#
# (misclassification_rate <- sum(round(get_prob_mod2(bb_train.HOFera[,c("rWAR.B", "AB
", "H", "HR", "BA", "R", "RBI", "SB", "OBP", "SLG", "OPS", "rWAR.P", "W", "L", "ERA",
"G", "GS", "SV", "IP", "SO", "WHIP")] , 0) != bb_train.HOFera$onBallot) / dim(bb_trai
n.HOFera)[1])
#
# bb_train.HOFera$pred1_cor <- round(get_prob_mod1(bb_train.HOFera$rWAR), 0) != bb_tr
ain.HOFera$onBallot
#
# ggplot(bb_train.HOFera, aes(rWAR, onBallot, color = bb_train.HOFera$pred1_cor)) + g
eom_jitter() + theme_bw() + labs(title = "WAR by Appearance on BBWAA Ballot", x = "WA
R", y = "On the BBWAA Ballot", color = "Correctly Predicted?") + stat_function(fun =
plot_prob_mod1, color = "red") + scale_color_manual(values = c("green", "red"))

```