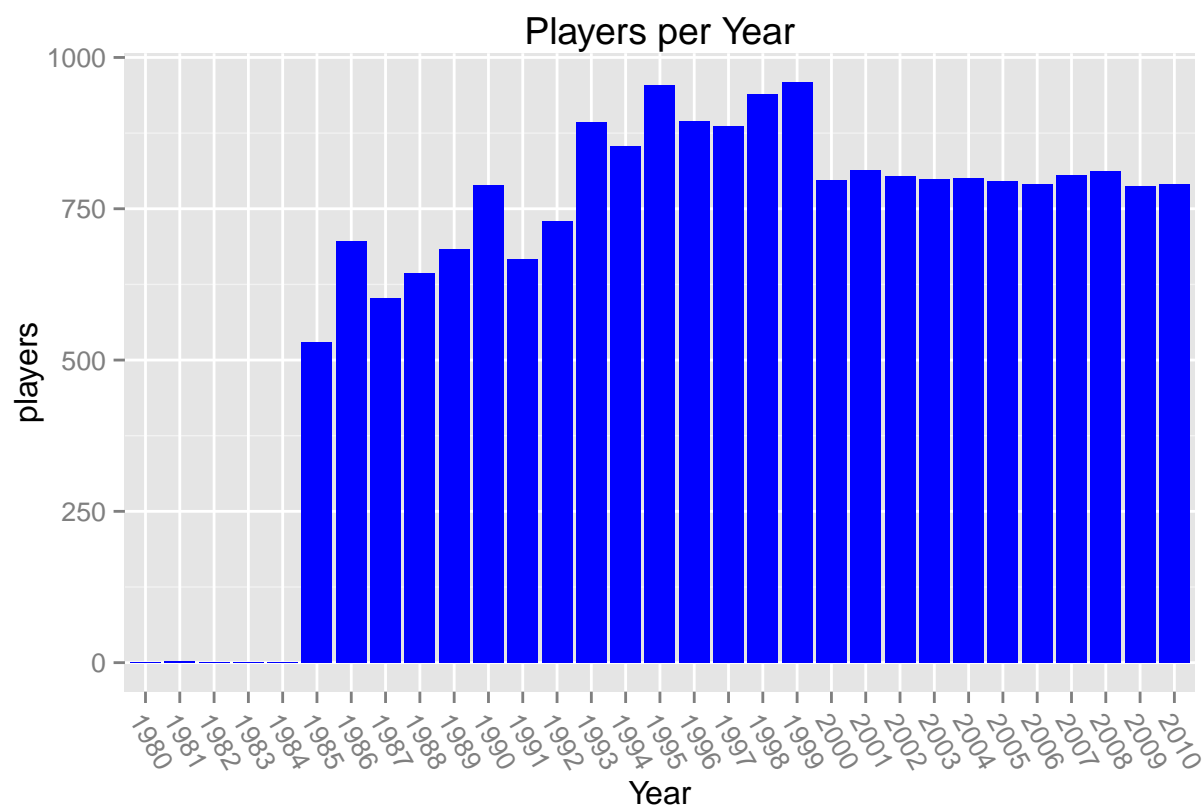


# Baseball Data Analysis

## Exploratory analysis of the data

The data set contains baseball player data from 1980 to 2010. Records are maintained on each player, for each year, team, league, stint, and fielding position they played. Overall, 4055 unique players are recorded in the data set, covering 20524 man years.

```
ggplot(by_year, aes(x = as.factor(yearID), y = players)) +  
  geom_bar(stat = "identity", fill = "Blue") +  
  xlab("Year") +  
  ggtitle("Players per Year") +  
  theme(axis.text.x = element_text(angle = -60, hjust = -0))
```



## Batting Data Exploration

Batting data was recorded including:

- b\_G -
- b\_G\_batting
- b\_AB - At Bats
- b\_R - Runs
- b\_H - Hits
- b\_2B - Doubles

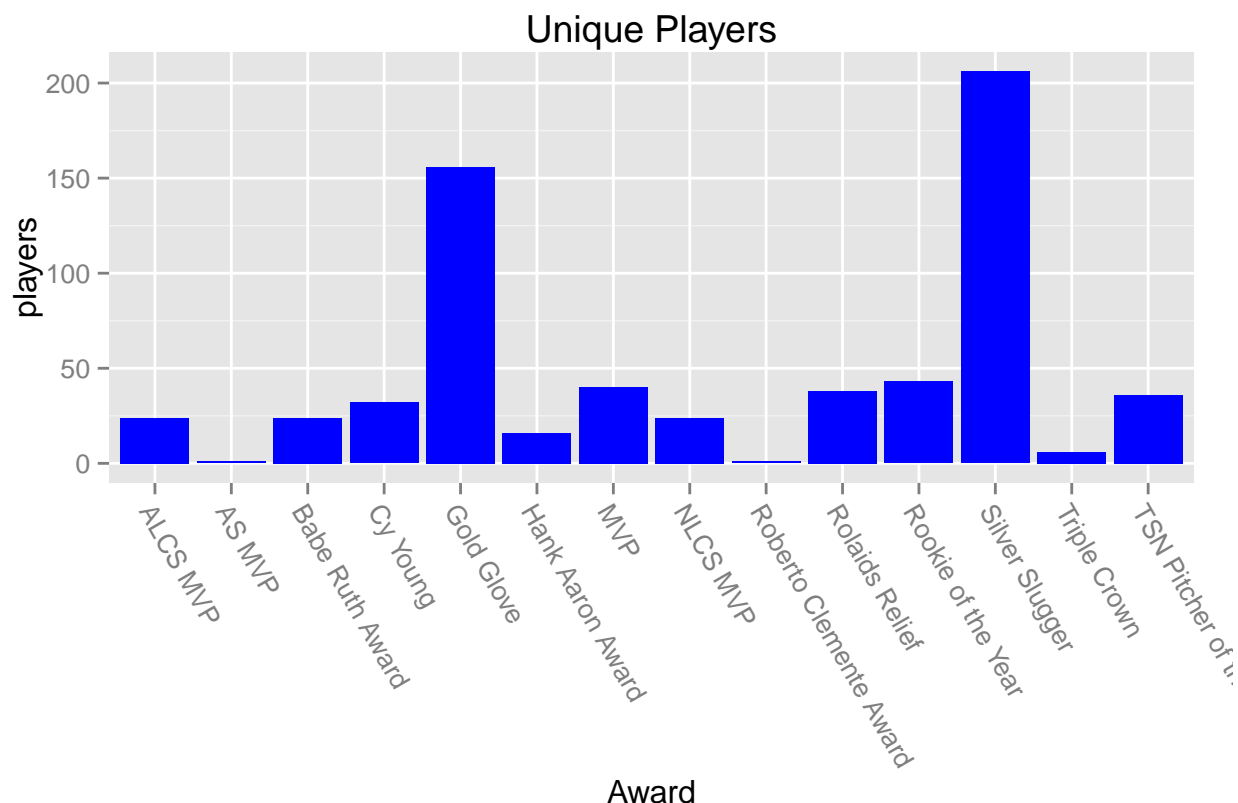
- b\_3B - Triples
- b\_HR - Home Runs
- b\_RBI - Runs Batted In
- b\_SB - Stolen Bases
- b\_CS - Caught Stealing
- b\_BB - Bases on Balls (Walks)
- b\_SO - Strikeouts
- b\_IBB - Intentional Bases on Balls (Walks)
- b\_HBP - Hit By Pitch
- b\_SH - Sacrifice Hits (Bunts)
- b\_SF - Sacrifice Flies
- b\_GIDP
- b\_G\_old

In addition to the baseball statistics provided, several additional ratios were computed to normalise the data:

- b\_hits\_per\_AB - Hits per At Bat
- b\_runs\_per\_AB - Runs per At Bat
- b\_runs\_per\_H - Runs per Hit
- b\_home\_runs\_per\_AB - Home Runs per At Bat
- b\_balls\_per\_AB - Balls per At Bat
- b\_RBI\_per\_H - Runs Batted In per Hit
- b\_HBP\_per\_AB - Hit By Pitch per At Bat

Finally, award data was considered:

```
ggplot(filter(award_winners, !is.na(awardID)), aes(x = as.factor(awardID), y = players)) +
  geom_bar(stat = "identity", fill = "Blue") +
  xlab("Award") +
  ggtitle("Unique Players") +
  theme(axis.text.x = element_text(angle = -60, hjust = -0))
```



As we can see, the silver slugger award has the most winners contained in the data set. As a result of this, it was chosen as the dependent variable.

The dataset also contained team, league, & fielding data that was not considered.

## Model selection

To model the likelihood that a player would received the Silver Slugger Award a Logistic Regression was performed. Since the Silver Slugger is a batting award recognising players of high offensive value, the batting data set was used. According to wikipedia, the baseball coaches vote for players on other teams to win based on several batting ratios that we were able to derive from our raw count data.

The full model contained:

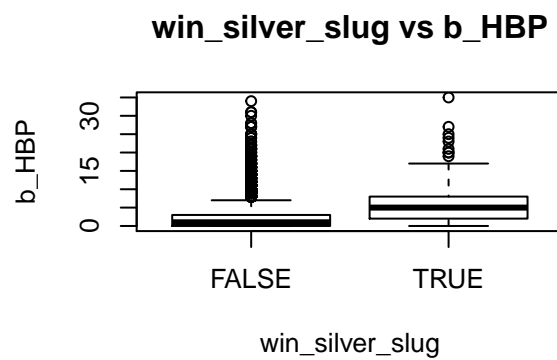
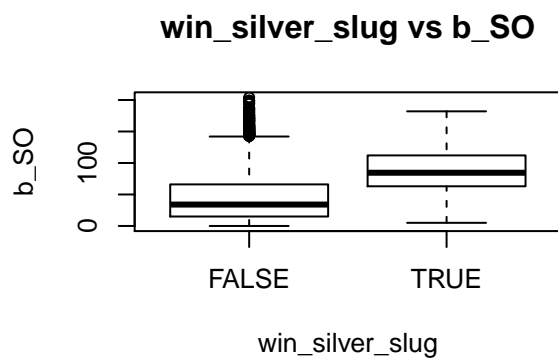
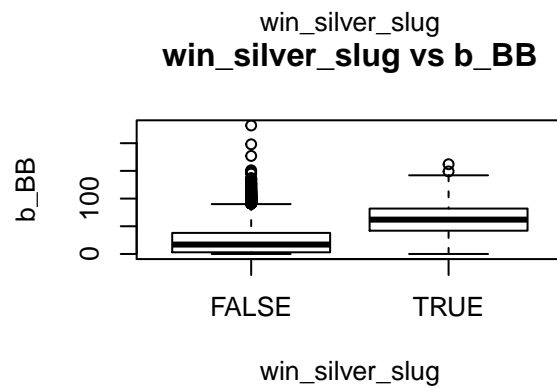
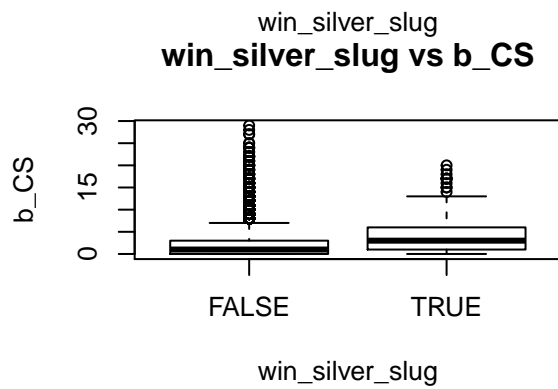
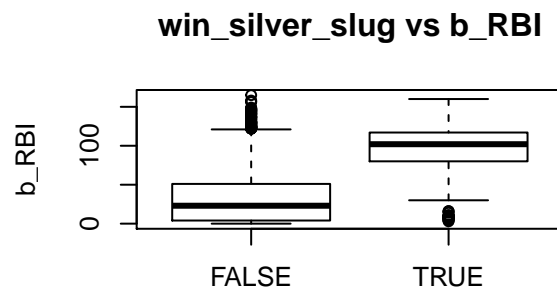
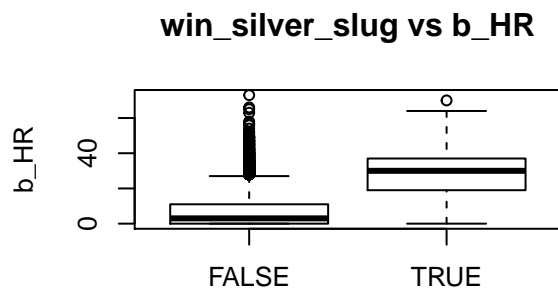
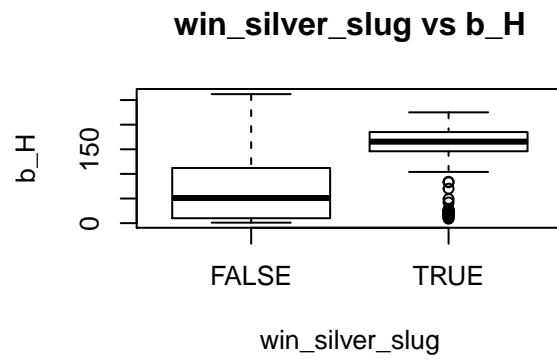
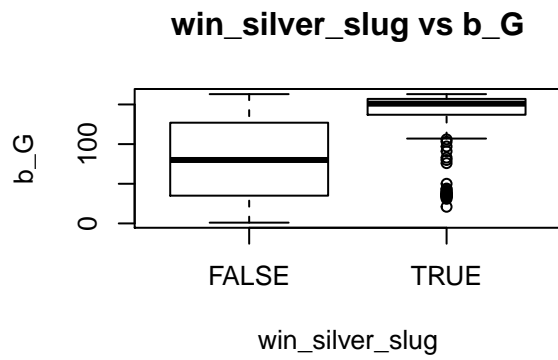
$$\begin{aligned} \text{win\_silver\_slug} = & \text{salary} + b\_G + b\_G\_batting + b\_AB + b\_R + b\_H + b\_2B + b\_3B \\ & + b\_HR + b\_RBI + b\_SB + b\_CS + b\_BB + b\_SO + b\_IBB + b\_HBP + b\_SH + \\ & b\_SF + b\_GIDP + b\_G\_old + b\_hits\_per\_AB + b\_runs\_per\_AB + b\_runs\_per\_H + \\ & b\_home\_runs\_per\_H + b\_balls\_per\_AB + b\_RBI\_per\_H + b\_HBP\_per\_AB \end{aligned}$$

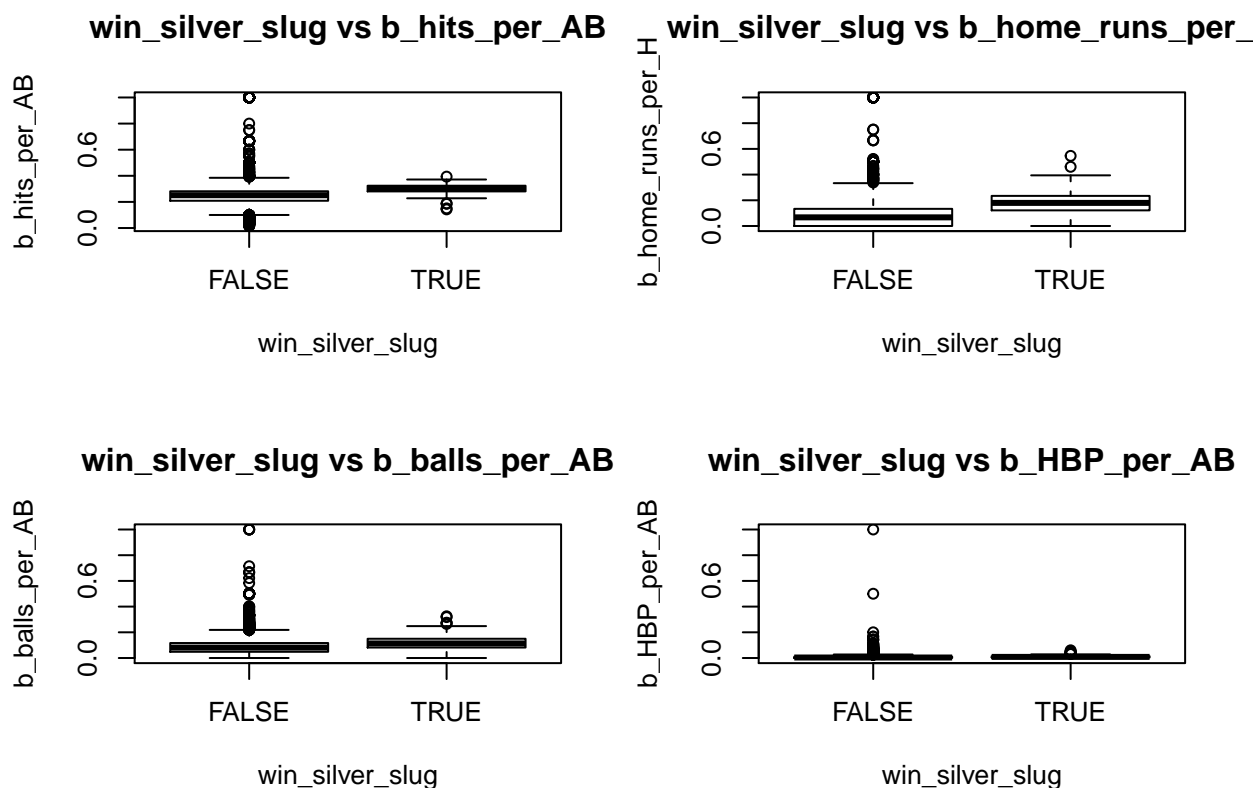
## Model Selection

After backward selection the model that was chosen was:  $\log(\text{odds}(\text{win\_silver\_slug})) = -5.10697989512792 + -0.0437813071772097 * b\_G + 0.0298363646334512 * b\_H + 0.0353237814406278 * b\_HR + 0.0235952761309906 * b\_RBI + 0.03690307794777 * b\_CS + 0.0283786990755013 * b\_BB + -0.00719867486183699 * b\_SO + 0.115025302963441 * b\_HBP + 1.90522965303078 * b\_hits\_per\_AB + 2.46316604034702 * b\_home\_runs\_per\_H + -10.2429175488287 * b\_balls\_per\_AB + -46.0399398528793 * b\_HBP\_per\_AB$

## Check Fit

```
batting_data[which(colnames(batting_data) %in% c("win_silver_slug", names(model_backwards_selection$coe  
plot_all_box(which(colnames(.) == "win_silver_slug"), .)
```





```
vif(model_backwards_selection)
```

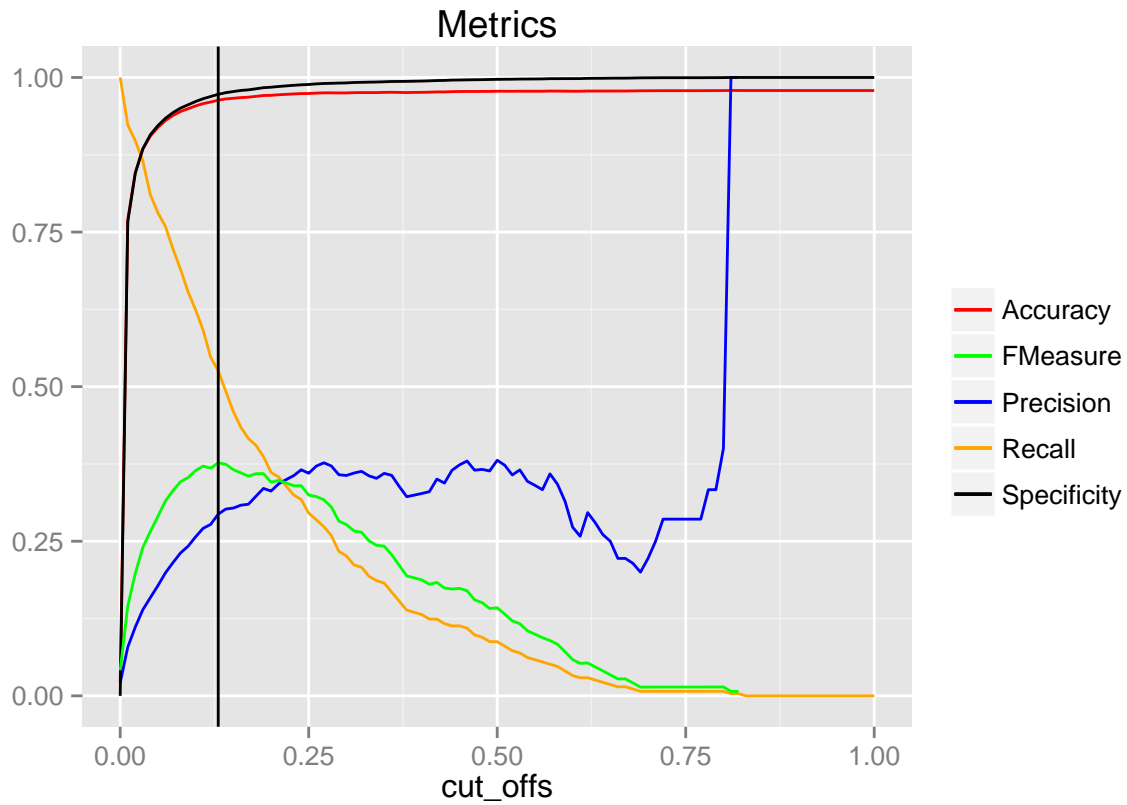
```
##          b_G          b_H          b_HR          b_RBI
##    11.457211    11.286429     7.332786    10.183930
##          b_CS          b_BB          b_SO          b_HBP
##     1.370876    17.147657     2.581509    19.221669
## b_hits_per_AB b_home_runs_per_H b_balls_per_AB b_HBP_per_AB
##     1.164119     1.841251    12.686127    17.834542
```

```
lm.beta(model_backwards_selection)
```

```
##
## Call:
## glm(formula = win_silver_slug ~ b_G + b_H + b_HR + b_RBI + b_CS +
##    b_BB + b_SO + b_HBP + b_hits_per_AB + b_home_runs_per_H +
##    b_balls_per_AB + b_HBP_per_AB, family = binomial(), data = batting_data)
##
## Standardized Coefficients::
##      (Intercept)          b_G          b_H          b_HR
##      0.0000000    -14.7591199    12.2880563     2.3867899
##          b_RBI          b_CS          b_BB          b_SO
##      5.3694099     0.8367252     5.0736515    -1.8274445
##          b_HBP b_hits_per_AB b_home_runs_per_H b_balls_per_AB
##      2.5745607     1.2437720     1.6719370    -4.4004964
## b_HBP_per_AB
##    -4.6642504
```

## Finding a Optimal Cut Off

```
plot_of_cut_offs
```



After searching for a good cut off value, 0.13 was chosen.

```
best_cut_off
```

```
## cut_offs true_positive true_negative false_positive false_negative
## 1 0.13 144 12345 346 130
## recall accuracy precision specificity f_measure prior lift
## 1 0.5255474 0.9632858 0.2938776 0.9727366 0.3769634 0.02113382 13.90556
```