

# Data Analysis

*Pipe It Up!: Nagaprasad Rudrapatna, Karen Deng, Jackson Muraika, Anna Zolotor*

*2020-04-15*

```
library(tidyverse)
library(broom)
library(stringr)
library(knitr)
library(tidyverse)
library(gridExtra)
library(leaps)

nba_social_power <- read_csv("../data/nba.csv")

nba_2016_2017_100 <- nba_social_power

nba_social_power_mod <- nba_social_power %>%
  filter(TWITTER_HANDLE != "0") %>%
  select(PLAYER_NAME,
         TEAM_ABBREVIATION,
         AGE,
         W_PCT,
         OFF_RATING,
         DEF_RATING,
         NET_RATING,
         AST_RATIO,
         REB_PCT,
         USG_PCT,
         PIE,
         SALARY_MILLIONS,
         ACTIVE_TWITTER_LAST_YEAR,
         TWITTER_FOLLOWER_COUNT_MILLIONS,
         PTS)
```

## Research Question and Objective

The research question explored in this analysis is: Is there a relationship between measures of athletic success (win percentage, offensive and defensive ratings, etc.) and the internet "popularity" of NBA athletes, measured in the number of Twitter followers?

The objective of this analysis is to predict Twitter follower counts of NBA players using measures of athletic success.

## The Data

The dataset we use in this analysis includes on-court performance data for NBA players in the 2016-2017 season, along with their salaries and Twitter engagement. Because we are examining the relationship between player stats and the number of twitter followers, we filtered for players who had an active twitter account, by filtering for values where TWITTER\_HANDLE is not "NA" (0). After filtering, we have 95 observations.

The response variable is `TWITTER_FOLLOWER_COUNT_MILLIONS`, which measures players' Twitter follower counts at the time the data was collected.

## Exploratory Data Analysis

### Univariate

First, we will do univariate EDA on the dataset. Player name will be used to refer to observations in our dataset, but since each player name is distinct we do not need to do EDA on the `PLAYER_NAME` variable.

Here, we'll take a look at how many players there are from each team in the dataset:

```
nba_social_power_mod %>%  
  count(Team_Abbreviation) %>%  
  arrange(n)
```

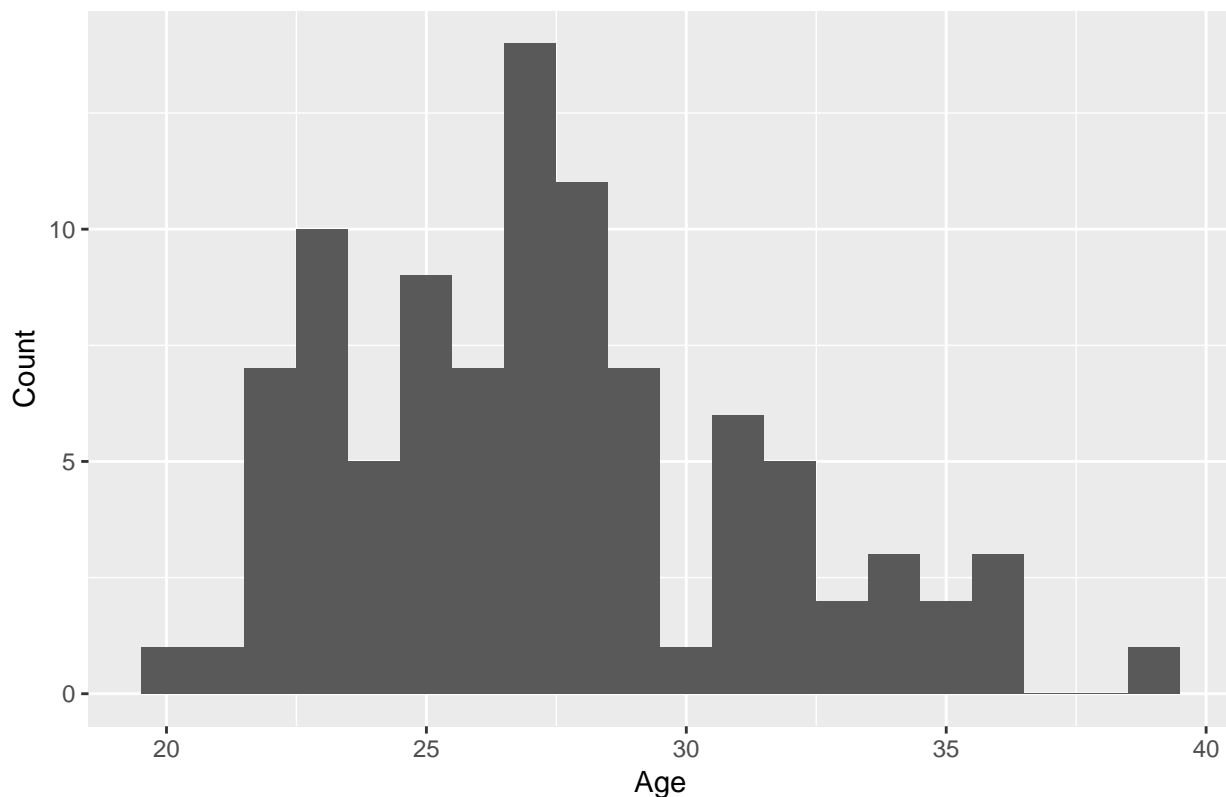
```
## # A tibble: 30 x 2  
##   Team_Abbreviation      n  
##   <chr>              <int>  
## 1 SAC                  1  
## 2 CHI                  2  
## 3 IND                  2  
## 4 LAL                  2  
## 5 MIA                  2  
## 6 MIN                  2  
## 7 ORL                  2  
## 8 WAS                  2  
## 9 ATL                  3  
## 10 BKN                 3  
## # ... with 20 more rows
```

As we can see from the output, there is only one team that is represented just once in the dataset: SAC, the Sacramento Kings. The greatest number of times teams are represented in the dataset is 5. GSW (Golden State Warriors), LAC (Los Angeles Clippers), and SAS (San Antonio Spurs) are all represented 5 times.

Now, we'll explore the distribution of the `AGE` variable in the dataset:

```
ggplot(data = nba_social_power_mod, aes(x = AGE)) +  
  geom_histogram(binwidth = 1) +  
  labs(x = "Age", y = "Count", title = "Distribution of Age")
```

### Distribution of Age



```
nba_social_power_mod %>%
  summarise(mean = mean(AGE), min= min(AGE), Q1 = quantile(AGE, .25), median = median(AGE),
            Q3 = quantile(AGE, .75), max = max(AGE))
```

```
## # A tibble: 1 x 6
##   mean   min    Q1 median    Q3   max
##   <dbl> <dbl> <dbl>   <dbl> <dbl> <dbl>
## 1  27.4    20  24.5     27    29    39
```

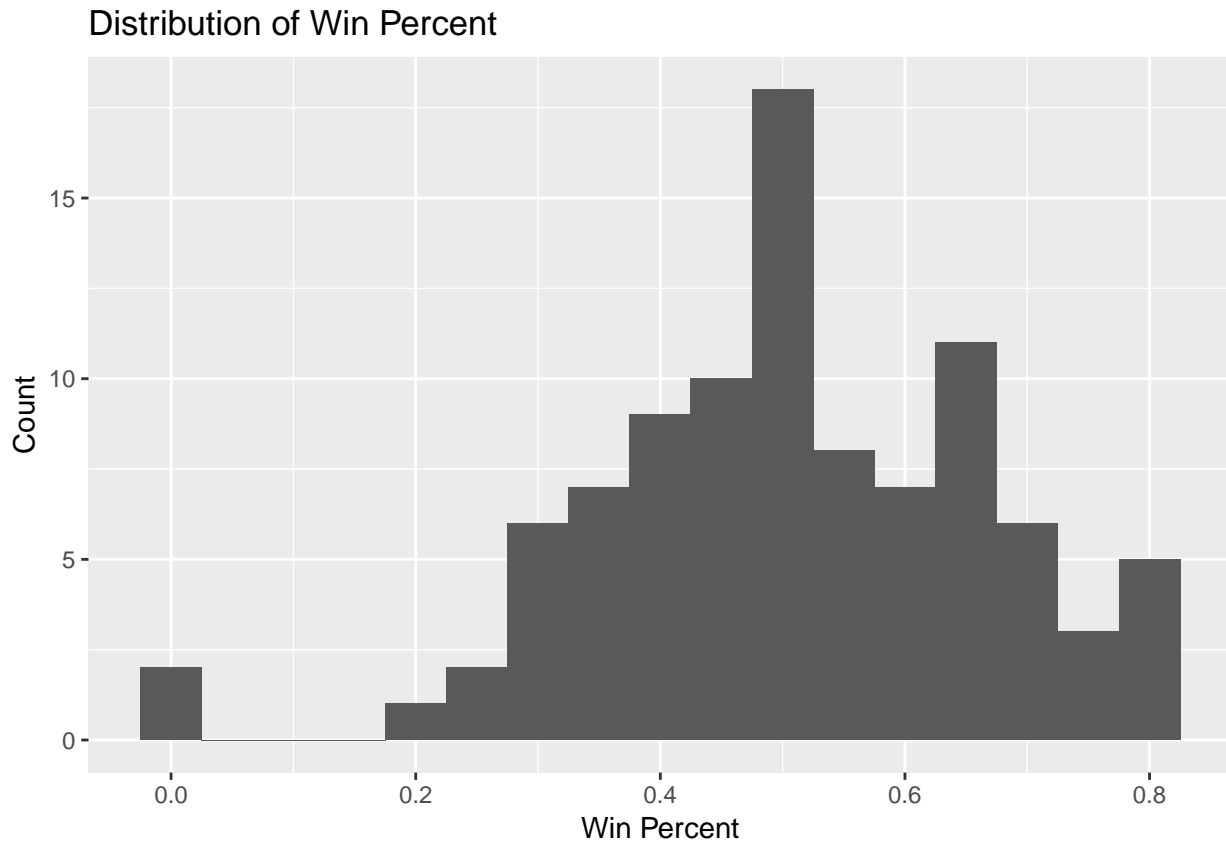
```
nba_social_power_mod %>%
  arrange(desc(AGE)) %>%
  head(1)
```

```
## # A tibble: 1 x 15
##   PLAYER_NAME TEAM_ABBREVIATI~ AGE W_PCT OFF_RATING DEF_RATING NET_RATING
##   <chr>         <chr>         <dbl> <dbl>     <dbl>     <dbl>     <dbl>
## 1 Dirk Nowit~ DAL             39 0.426     105.     106.     -1.7
## # ... with 8 more variables: AST_RATIO <dbl>, REB_PCT <dbl>,
## #   USG_PCT <dbl>, PIE <dbl>, SALARY_MILLIONS <dbl>,
## #   ACTIVE_TWITTER_LAST_YEAR <dbl>, TWITTER_FOLLOWER_COUNT_MILLIONS <dbl>,
## #   PTS <dbl>
```

As we can see from the histogram, age is somewhat normally distributed in the dataset, with a mode around 27 and a surprisingly low number of 30-year olds. The mean age, 27.39, and median age, 27, are very close together, indicating little skew. The lowest age is 20 and the highest is 39. The oldest player by far, at 39, is Dirk Nowitzki.

Now, we'll examine the distribution of win percent, W\_PCT:

```
ggplot(data = nba_social_power_mod, aes(x= W_PCT)) +
  geom_histogram(binwidth = .05) +
  labs(x = "Win Percent", y = "Count", title = "Distribution of Win Percent")
```



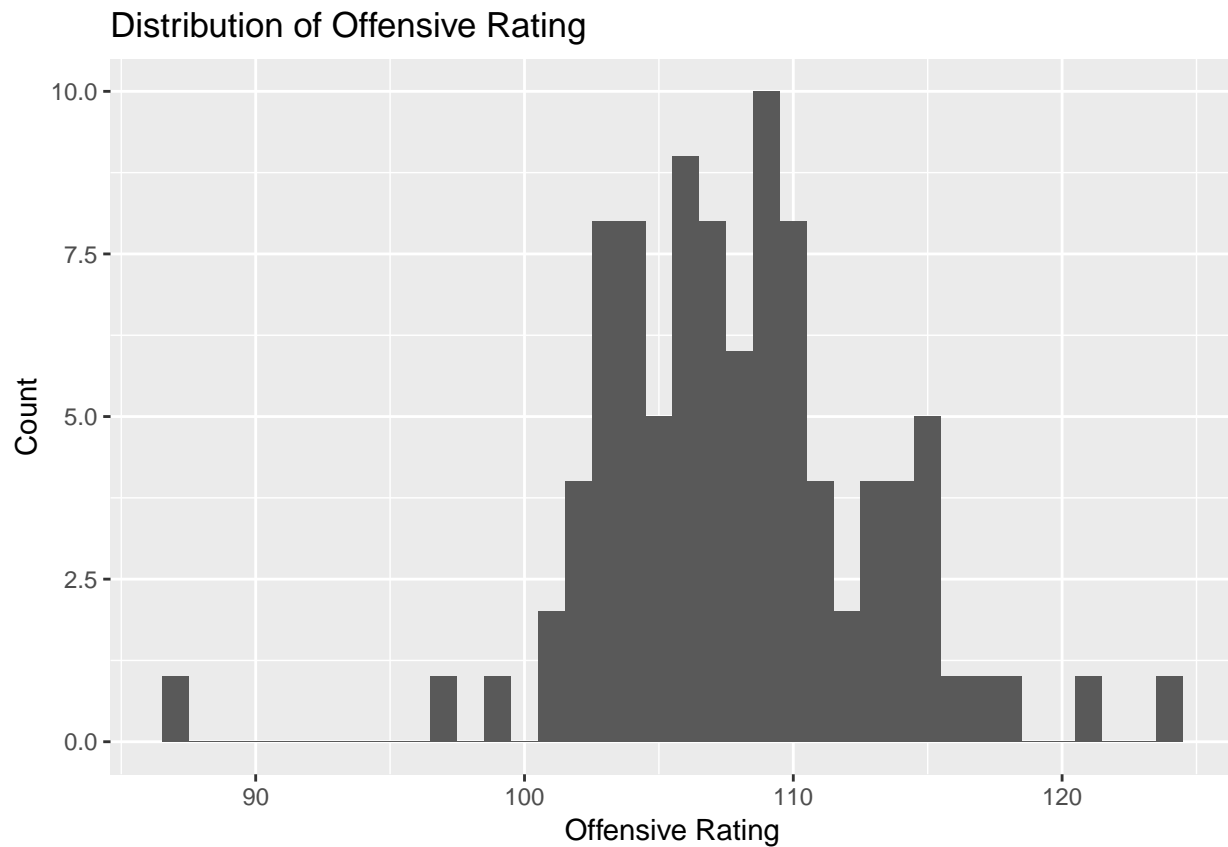
```
nba_social_power_mod %>%
  summarise(mean = mean(W_PCT), min= min(W_PCT), Q1 = quantile(W_PCT, .25),
            median = median(W_PCT), Q3 = quantile(W_PCT, .75), max = max(W_PCT))
```

```
## # A tibble: 1 x 6
##   mean   min    Q1 median    Q3   max
##   <dbl> <dbl> <dbl>   <dbl> <dbl> <dbl>
## 1 0.511     0 0.418  0.507  0.63 0.824
```

As we can see from the histogram, win percent is also somewhat normally distributed, with a mode around 50 percent. The minimum win percent in the dataset is 0, while the maximum is 82.4. The median of 50.7 is very similar to the mean of 51%. The fact that the mean and median win percents in the dataset fall so close to 50% indicate good randomness in the dataset, b/c the mean and median win percents for all nba players are 50%.

Next, we'll look at the distributions for offensive rating, `OFF_RATING` and defensive rating `DEF_RATING`, as well as the distribution for `NET_RATING`, which is the average of the offensive and defensive rating:

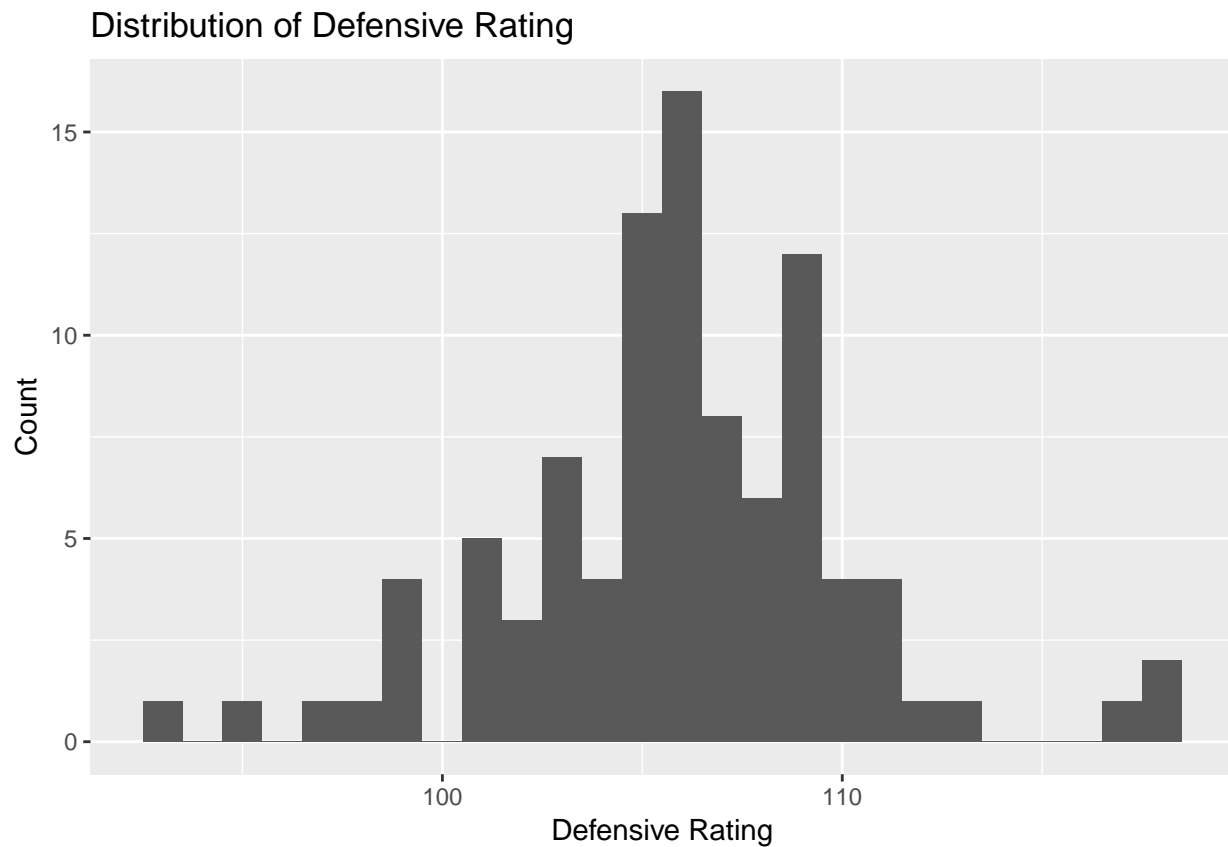
```
ggplot(data = nba_social_power_mod, aes(x= OFF_RATING)) +
  geom_histogram(binwidth = 1) +
  labs(x = "Offensive Rating", y = "Count",
       title = "Distribution of Offensive Rating")
```



```
nba_social_power_mod %>%
  summarise(min= min(OFF_RATING), median = median(OFF_RATING), max = max(OFF_RATING))

## # A tibble: 1 x 3
##   min median  max
##   <dbl> <dbl> <dbl>
## 1  86.8  108.  124.

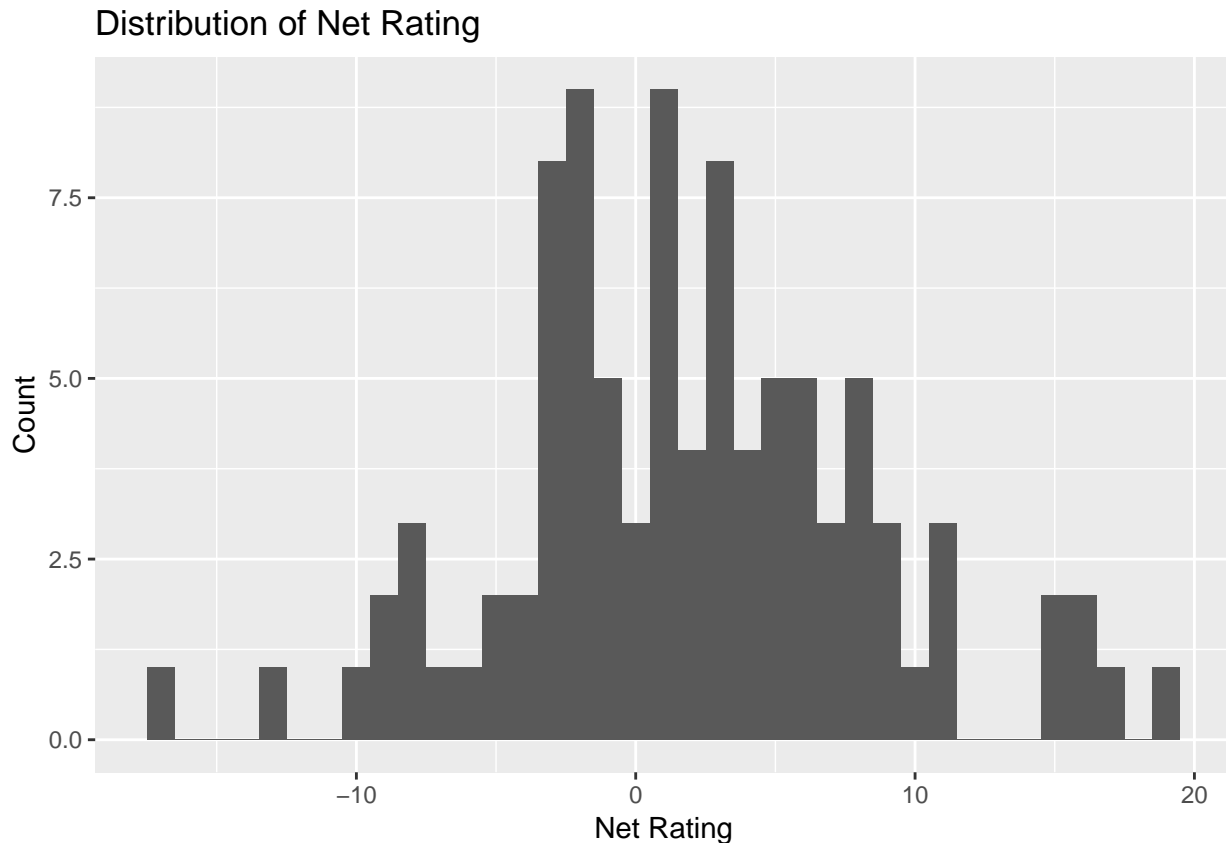
ggplot(data = nba_social_power_mod, aes(x= DEF_RATING)) +
  geom_histogram(binwidth = 1) +
  labs(x = "Defensive Rating", y = "Count", title = "Distribution of Defensive Rating")
```



```
nba_social_power_mod %>%
  summarise(min= min(DEF_RATING), median = median(DEF_RATING), max = max(DEF_RATING))
```

```
## # A tibble: 1 x 3
##   min median  max
##   <dbl>  <dbl> <dbl>
## 1    93   106  118.
```

```
ggplot(data = nba_social_power_mod, aes(x= NET_RATING)) +
  geom_histogram(binwidth = 1) +
  labs(x = "Net Rating", y = "Count", title = "Distribution of Net Rating")
```



```
nba_social_power_mod %>%
  summarise(min= min(NET_RATING), median = median(NET_RATING), max = max(NET_RATING))
```

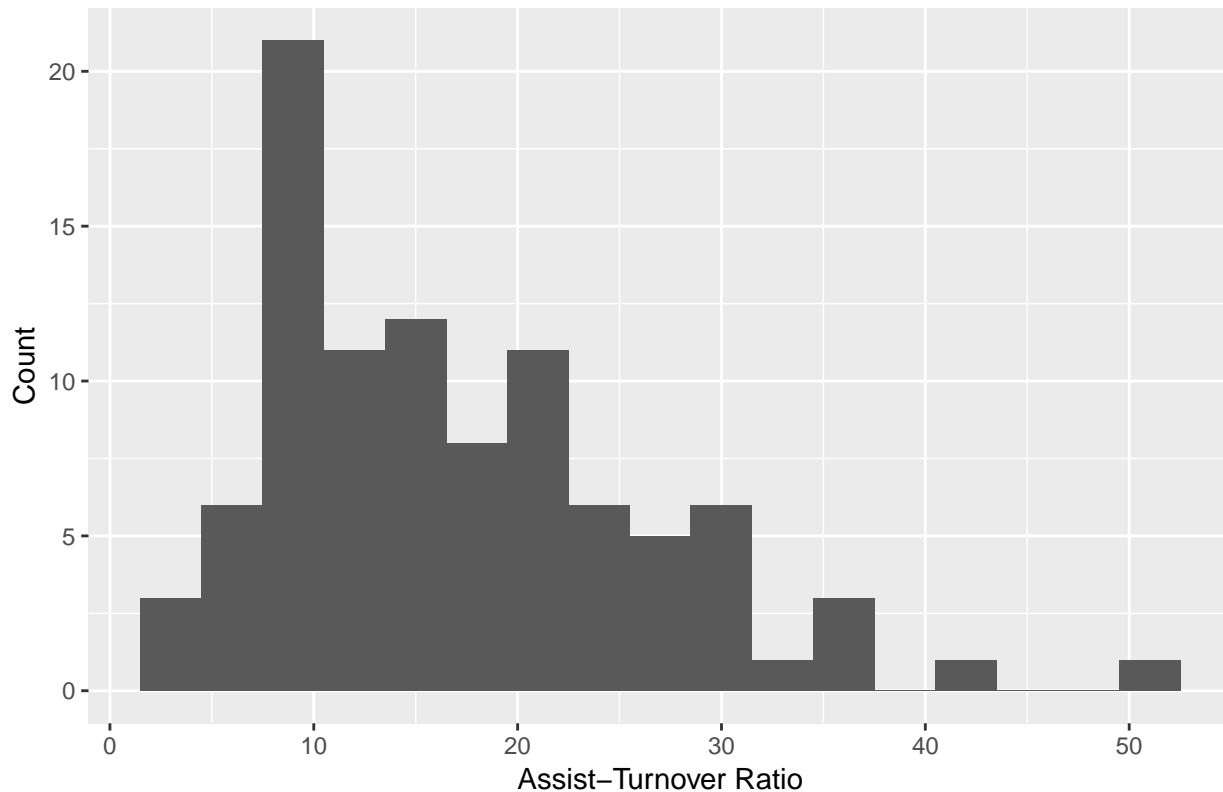
```
## # A tibble: 1 x 3
##   min median  max
##   <dbl> <dbl> <dbl>
## 1 -17.2   1.5  18.7
```

Defensive rating, offensive rating, and net rating do not stray far from normally distributed. Offensive rating varies from 86.8 to 124.2, with a median of 107.6. Defensive rating varies from 93 to 118.3, with a median of 106. Thus, the dataset contains a larger range in terms of offensive rating, and the median is also slightly higher for defensive rated players. The distribution of net rating has multiple nearly equal modes around -2 to -3 and around 1 and 3. The median net rating is 1.5, and the net ratings in the dataset vary from -17.2 to 18.7.

Next, we'll look at the distribution of the assist-to-turnovers ratio, `AST_RATIO`:

```
ggplot(data = nba_social_power_mod, aes(x= AST_RATIO)) +
  geom_histogram(binwidth = 3) +
  labs(x = "Assist-Turnover Ratio", y = "Count", title = "Distribution of Assist-Turnover Ratio")
```

### Distribution of Assist–Turnover Ratio



```
nba_social_power_mod %>%
  summarise(mean = mean(AST_RATIO), min= min(AST_RATIO), Q1 = quantile(AST_RATIO, .25), median = median(AST_RATIO), max = max(AST_RATIO))
```

```
## # A tibble: 1 x 6
##   mean   min    Q1 median    Q3   max
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  17.1    4  9.75    15  22.2  51.5
```

```
nba_social_power_mod %>%
  arrange(desc(AST_RATIO)) %>%
  head(2)
```

```
## # A tibble: 2 x 15
##   PLAYER_NAME TEAM_ABBREVIATION AGE W_PCT OFF_RATING DEF_RATING NET_RATING
##   <chr>        <chr>          <dbl> <dbl>    <dbl>    <dbl>    <dbl>
## 1 Jarnell St~ DEN             23 0      115.     118.    -3.1
## 2 Ricky Rubio MIN             26 0.373  109.     110.    -1
## # ... with 8 more variables: AST_RATIO <dbl>, REB_PCT <dbl>,
## #   USG_PCT <dbl>, PIE <dbl>, SALARY_MILLIONS <dbl>,
## #   ACTIVE_TWITTER_LAST_YEAR <dbl>, TWITTER_FOLLOWER_COUNT_MILLIONS <dbl>,
## #   PTS <dbl>
```

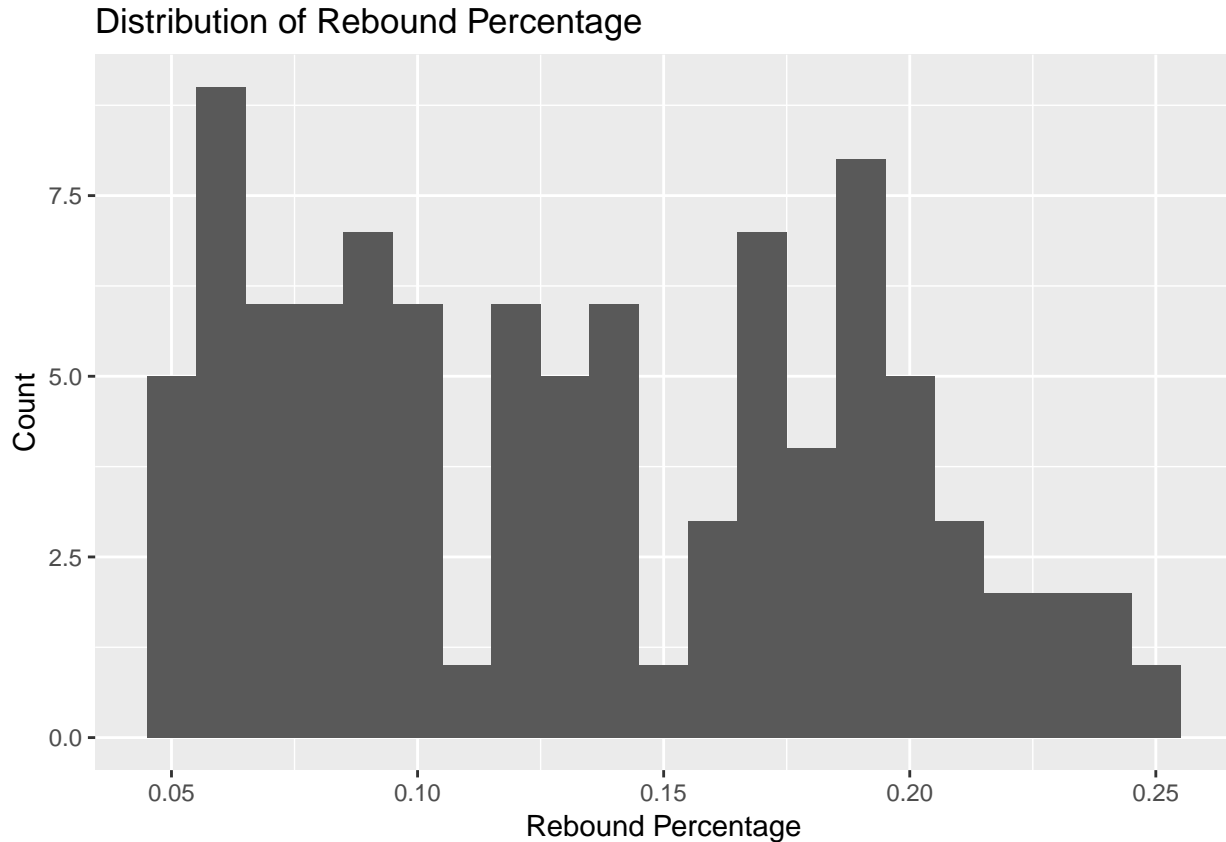
As we can see from the histogram, the assist-turnover ratio is very right skewed. The mode is at around 10, even though the median is at 15, and the mean is 17.12526, all of which are summary statistics that emphasize the right skew. This means that while most players in the dataset had a very high assist-turnover ratio (meaning they had many more assists than turnovers), there is a wider variation among players with a high ratio and the players with lower ratios are concentrated around a few numbers. The dataset minimum ratio of 4 means that there were no players with more turnovers than assists. Notably, this is the first



variable we've examined so far with a significantly non-normal distribution. The two players with very high assist-turnover ratios, 51.5 and 41.3, are Jarnell Stokes and Ricky Rubio, respectively.

Next, we'll examine the variation of REB\_PCT, the percent of rebounds a player makes:

```
ggplot(data = nba_social_power_mod, aes(x= REB_PCT)) +
  geom_histogram(binwidth = .01) +
  labs(x = "Rebound Percentage", y = "Count", title = "Distribution of Rebound Percentage")
```



```
nba_social_power_mod %>%
  summarise(mean = mean(REB_PCT), min= min(REB_PCT), Q1 = quantile(REB_PCT, .25), median = median(REB_PCT), Q3 = quantile(REB_PCT, .75), max = max(REB_PCT))

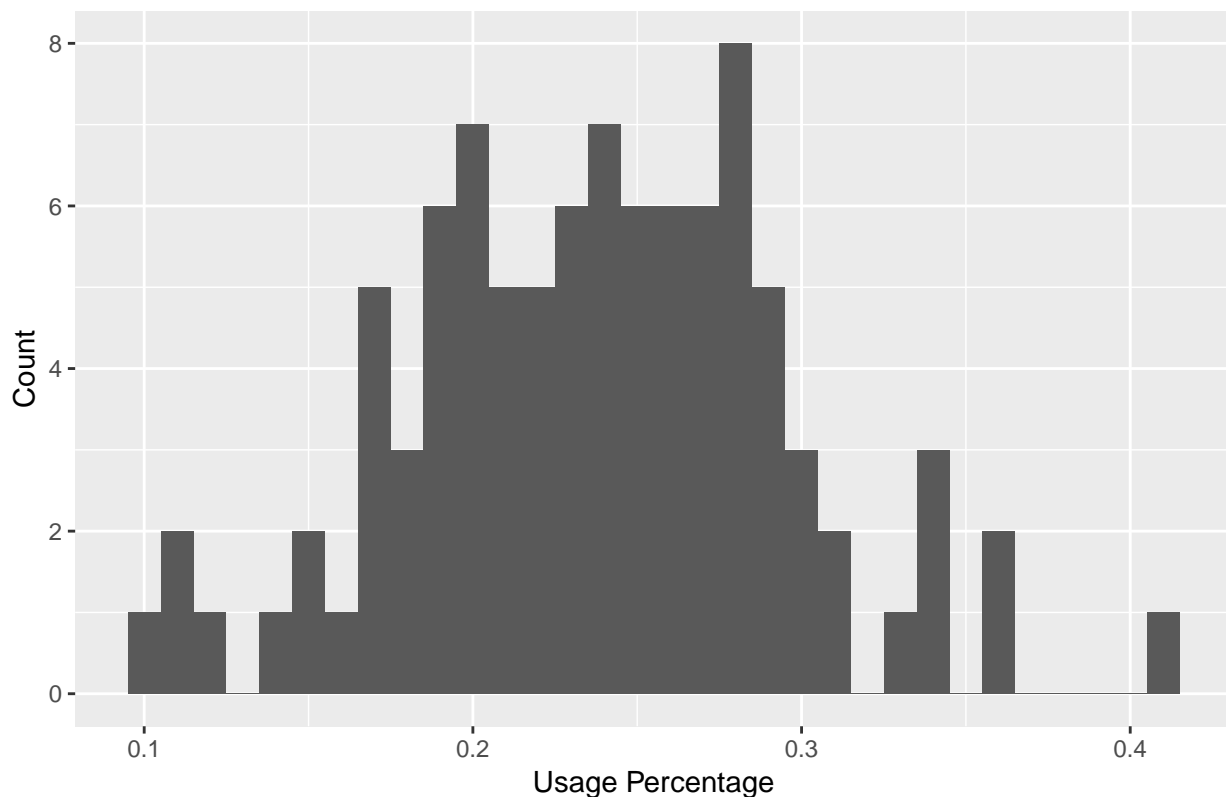
## # A tibble: 1 x 6
##   mean   min    Q1 median    Q3   max
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 0.133 0.045 0.0825 0.127 0.180 0.252
```

The distribution of rebound percentage has a minimum of 0.045 and a maximum of 0.252. The distribution is not very skewed one way or another, as supported by the similar mean of .133 and median of .127. However, the distribution is not normal in that it does not resemble a bell curve; with exceptions, the data is somewhat evenly distributed from the minimum to near the maximum (although there is some trail-off towards the right side of the distribution). This non-normal spread is likely partially an indication of the fact that the dataset contains both offensive and defensive players, because whether a player is on offense or defense has a significant effect on their rebound percentage.

Next, we'll look at USG\_PCT, usage percentage, which is an estimate of how often a player makes team plays:

```
ggplot(data = nba_social_power_mod, aes(x= USG_PCT)) +
  geom_histogram(binwidth = .01) +
  labs(x = "Usage Percentage", y = "Count", title = "Distribution of Usage Percentage")
```

## Distribution of Usage Percentage



```
nba_social_power_mod %>%
  summarise(mean = mean(USG_PCT), min= min(USG_PCT), Q1 = quantile(USG_PCT, .25), median = median(USG_PCT), Q3 = quantile(USG_PCT, .75), max = max(USG_PCT))

## # A tibble: 1 x 6
##   mean   min   Q1 median   Q3   max
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 0.238 0.101   0.2  0.242 0.276 0.408

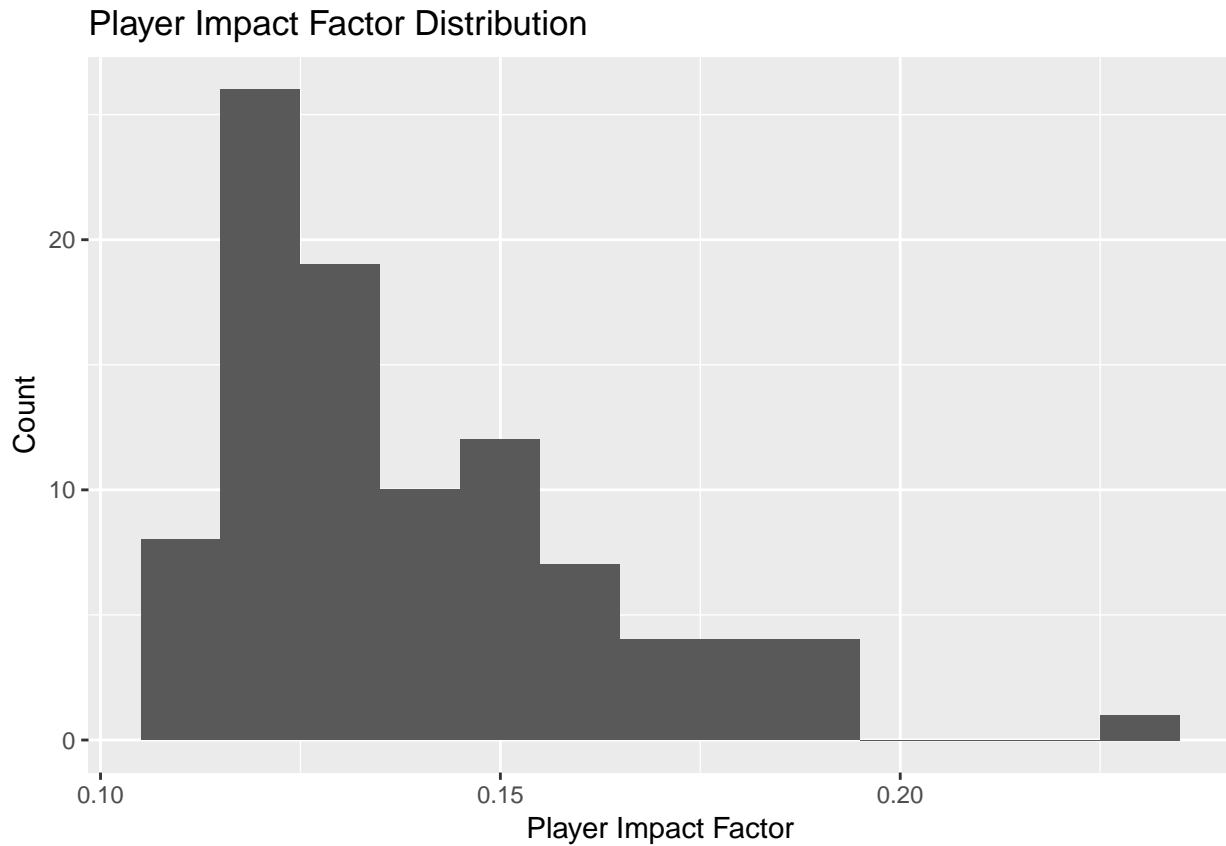
nba_social_power_mod %>%
  arrange(desc(USG_PCT)) %>%
  head(1)

## # A tibble: 1 x 15
##   PLAYER_NAME TEAM_ABBREVIATION AGE W_PCT OFF_RATING DEF_RATING NET_RATING
##   <chr>      <chr>          <dbl> <dbl>    <dbl>    <dbl>    <dbl>
## 1 Russell Westbrook OKC          28 0.568    108.     105.     3.3
## # ... with 8 more variables: AST_RATIO <dbl>, REB_PCT <dbl>,
## #   USG_PCT <dbl>, PIE <dbl>, SALARY_MILLIONS <dbl>,
## #   ACTIVE_TWITTER_LAST_YEAR <dbl>, TWITTER_FOLLOWER_COUNT_MILLIONS <dbl>,
## #   PTS <dbl>
```

The distribution of usage percentage, with a minimum of .101 and a maximum of .408, is fairly normally distributed. The mean, .238, and median, .242, are similar. The fairly wide spread may indicate that the dataset contains a decent sampling of players- some 'star player' types and others that are not the centerpieces of their teams. The maximum of .408, while perhaps not quite an outlier, is separated from most of the other points; this usage percentage belongs to Russell Westbrook.

Next, we'll examine PIE, player impact factor, a statistic roughly measuring a player's impact on the games that they play that's used by nba.com:

```
ggplot(data = nba_social_power_mod, aes(x= PIE)) +
  geom_histogram(binwidth = .01) +
  labs(x = "Player Impact Factor", y = "Count", title = "Player Impact Factor Distribution")
```



```
nba_social_power_mod %>%
  summarise(mean = mean(PIE), min= min(PIE), Q1 = quantile(PIE, .25), median = median(PIE), Q3 = quantile(PIE, .75))
```

```
## # A tibble: 1 x 6
##   mean   min    Q1 median    Q3    max
##   <dbl> <dbl> <dbl>  <dbl> <dbl> <dbl>
## 1 0.139 0.112 0.122  0.131 0.152  0.23
```

```
nba_social_power_mod %>%
  arrange(desc(PIE)) %>%
  select(PLAYER_NAME, PIE) %>%
  head(10)
```

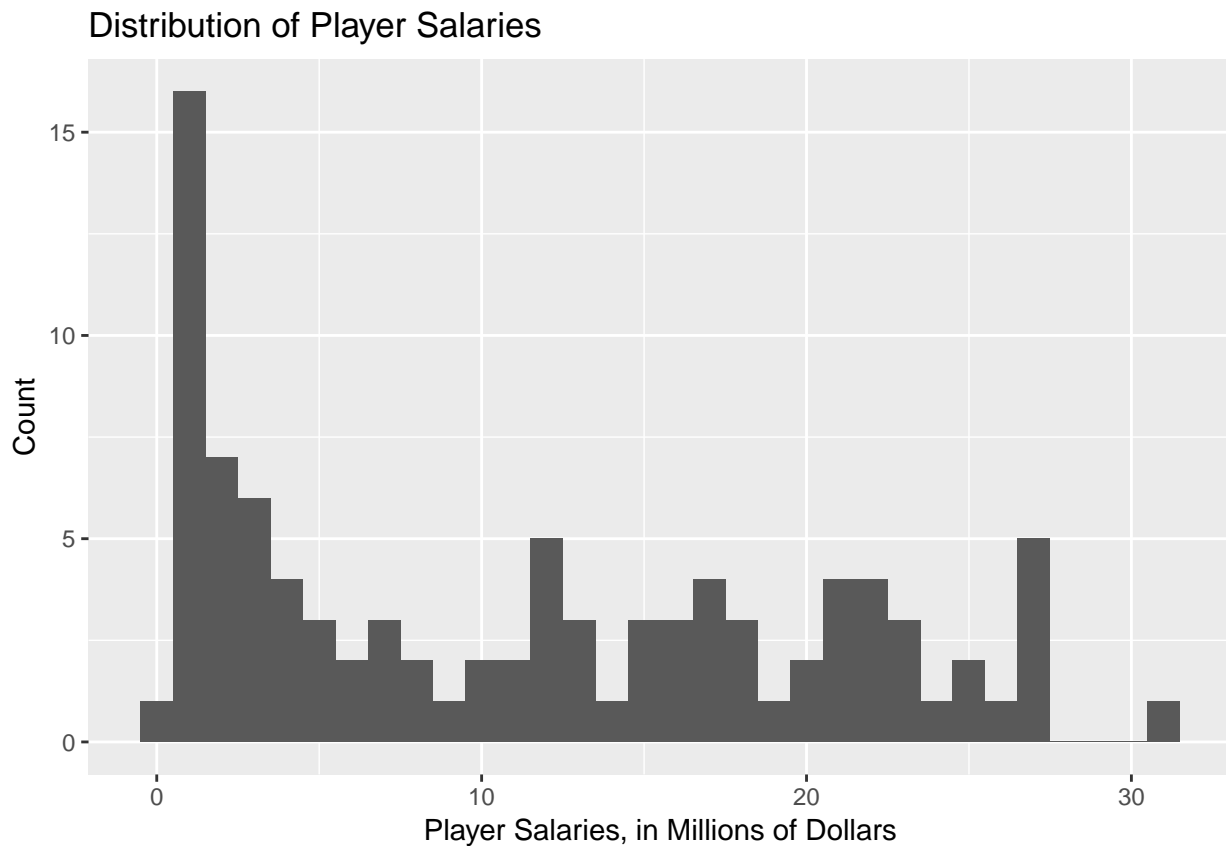
```
## # A tibble: 10 x 2
##   PLAYER_NAME      PIE
##   <chr>          <dbl>
## 1 Russell Westbrook 0.23
## 2 Demetrius Jackson 0.194
## 3 Anthony Davis    0.192
## 4 James Harden     0.19
## 5 Kevin Durant     0.186
## 6 LeBron James     0.183
## 7 Chris Paul       0.182
## 8 DeMarcus Cousins 0.178
```

```
## 9 Giannis Antetokounmpo 0.176
## 10 Kawhi Leonard 0.174
```

As we can see from the histogram, the player impact factor, with a minimum of .112 and a maximum of .23, is quite right-skewed. The median player impact factor is .131, and the mean is .139, evidence of the right skew. The mode is around the median. The maximum, .23, is a significant outlier, and is that of Russell Westbrook, the same player who had by far the highest usage percentage; clearly, his data will need to be examined more closely later to see if it ultimately affects our model.

Next, we'll look into players' salaries, in millions of dollars:

```
ggplot(data = nba_social_power_mod, aes(x= SALARY_MILLIONS)) +
  geom_histogram(binwidth = 1) +
  labs(x = "Player Salaries, in Millions of Dollars", y = "Count", title = "Distribution of Player Salaries")
```



```
nba_social_power_mod %>%
  summarise(mean = mean(SALARY_MILLIONS), min= min(SALARY_MILLIONS), Q1 = quantile(SALARY_MILLIONS, .25)
```

```
## # A tibble: 1 x 6
##   mean   min   Q1 median   Q3   max
##   <dbl> <dbl> <dbl>   <dbl> <dbl> <dbl>
## 1  11.3  0.31  2.47   11.3  18.5  31.0
```

```
nba_social_power_mod %>%
  arrange(desc(SALARY_MILLIONS)) %>%
  select(PLAYER_NAME, SALARY_MILLIONS) %>%
  head(10)
```

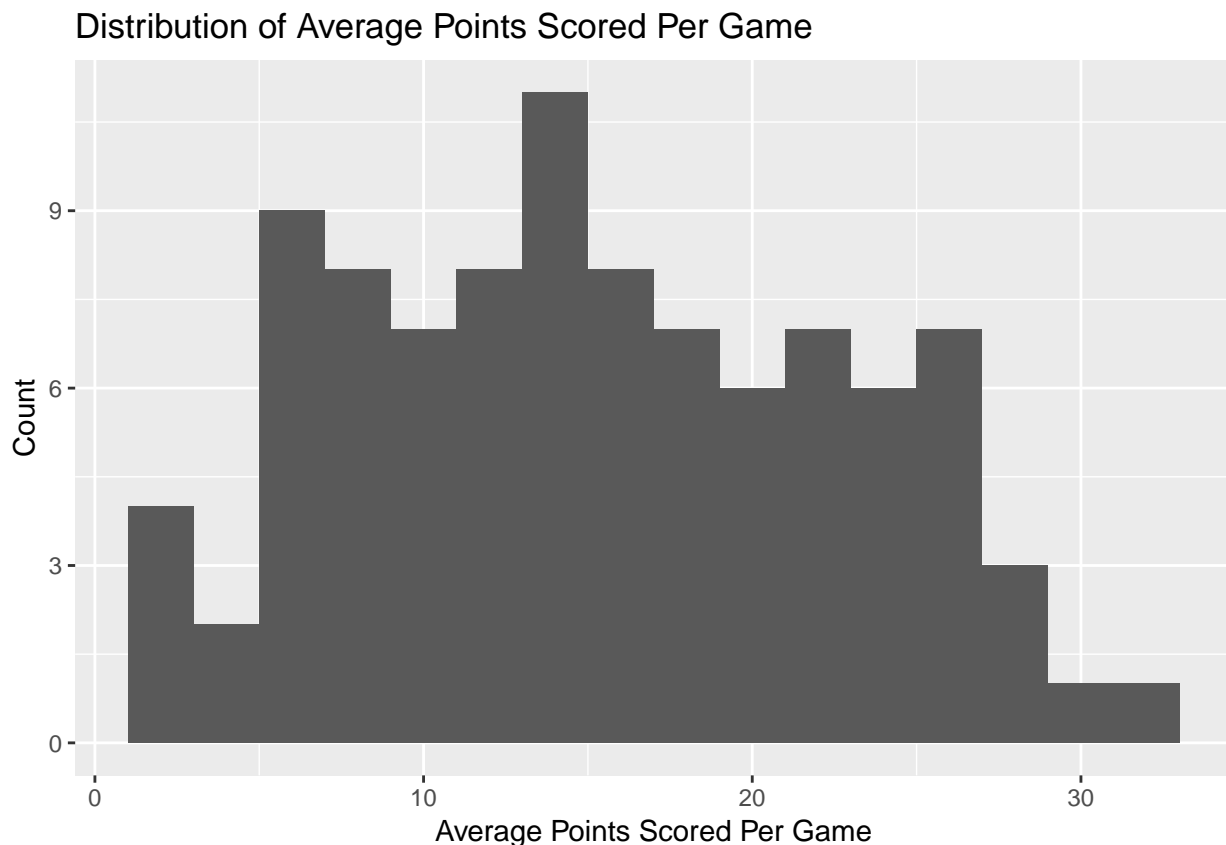
```
## # A tibble: 10 x 2
##   PLAYER_NAME      SALARY_MILLIONS
```

```
##      <chr>                <dbl>
## 1 LeBron James           31.0
## 2 Russell Westbrook      26.5
## 3 Kevin Durant           26.5
## 4 Mike Conley             26.5
## 5 DeMar DeRozan          26.5
## 6 Al Horford              26.5
## 7 James Harden           26.5
## 8 Dirk Nowitzki           25
## 9 Carmelo Anthony         24.6
## 10 Damian Lillard         24.3
```

As we can see from the histogram, the distribution of salaries is somewhat right-skewed, with most of the players making less than 20 million dollars a year. The mean salary is 11.3 million dollars a year. The player who earns the most, at 30.96 million dollars per year, is LeBron James. On the other hand, Russell Westbrook, Kevin Durant, Mike Conley, DeMar DeRozan, and Al Horford each earn 26.54 million dollars per year.

Next, we'll look at average points scored per game:

```
ggplot(data = nba_social_power_mod, aes(x= PTS)) +
  geom_histogram(binwidth = 2) +
  labs(x = "Average Points Scored Per Game", y = "Count", title = "Distribution of Average Points Scored Per Game")
```



```
nba_social_power_mod %>%
  summarise(mean = mean(PTS), min= min(PTS), Q1 = quantile(PTS, .25), median = median(PTS), Q3 = quantile(PTS, .75), max = max(PTS))

## # A tibble: 1 x 6
##   mean  min  Q1 median  Q3  max
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
```

```
##    <dbl> <dbl> <dbl>    <dbl> <dbl> <dbl>
## 1   15.3   1.5   9.5     14.6  21.4  31.6
```

As we can see from the histogram, the distribution of average points scored is slightly normal, with some obvious departures from normality. The median number of points scored per game is 14.6, and the mean is 15.28232. The maximum is 31.6, but this does not seem to be an obvious outlier.

Next, we'll examine the variable `ACTIVE_TWITTER_LAST_YEAR`, which tells us whether or not each player posted on Twitter the year before the data was collected:

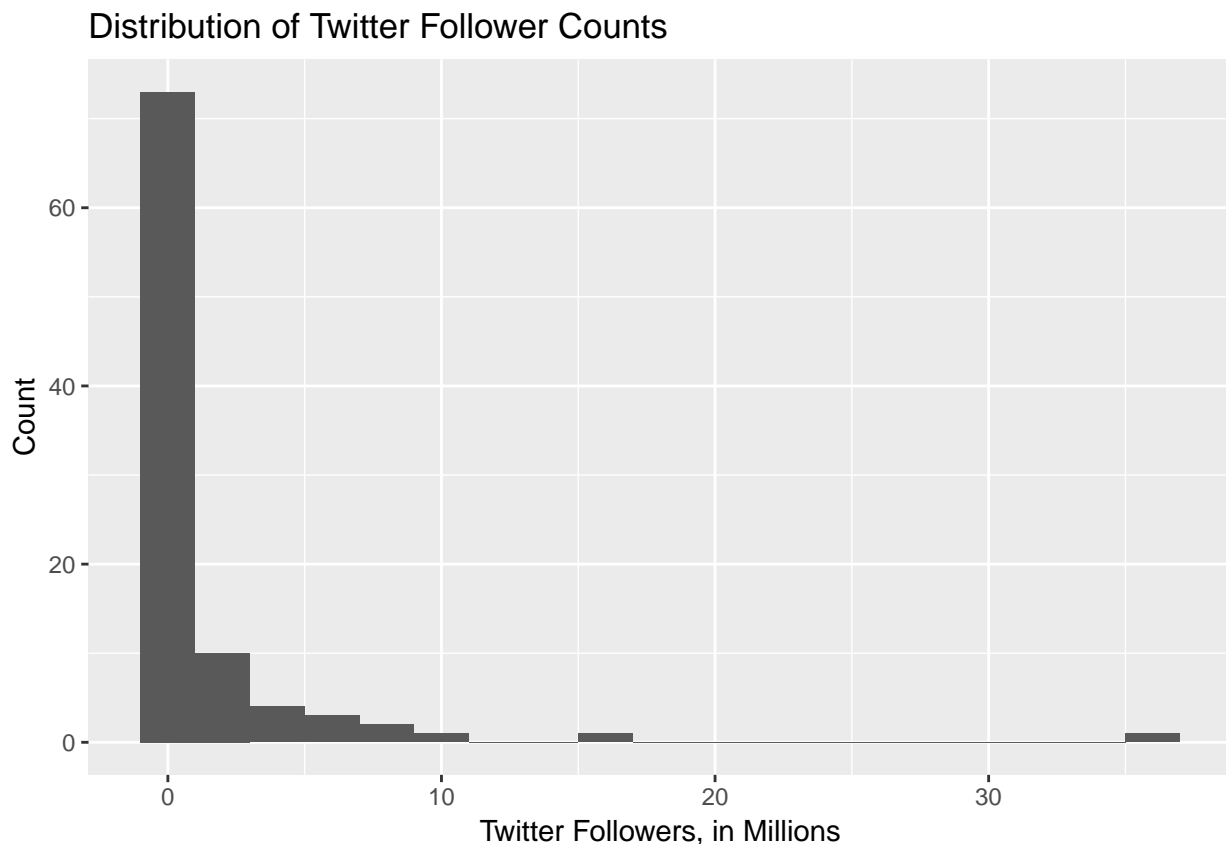
```
nba_social_power_mod %>%
  count(ACTIVE_TWITTER_LAST_YEAR)
```

```
## # A tibble: 2 x 2
##   ACTIVE_TWITTER_LAST_YEAR     n
##               <dbl> <int>
## 1                     0     2
## 2                     1    93
```

Out of the 95 players in our modified dataset, 2 were not active on Twitter the year before the data was collected and 93 were.

Finally, we'll look into the response variable, `TWITTER_FOLLOWER_COUNT_MILLIONS`:

```
ggplot(data = nba_social_power_mod, aes(x= TWITTER_FOLLOWER_COUNT_MILLIONS)) +
  geom_histogram(binwidth = 2) +
  labs(x = "Twitter Followers, in Millions", y = "Count", title = "Distribution of Twitter Follower Counts")
```



```
nba_social_power_mod %>%
  summarise(mean = mean(TWITTER_FOLLOWER_COUNT_MILLIONS), min= min(TWITTER_FOLLOWER_COUNT_MILLIONS), Q1
```

```
## # A tibble: 1 x 6
##   mean   min     Q1 median   Q3    max
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  1.60 0.002 0.0595 0.246 0.912   37

nba_social_power_mod %>%
  arrange(desc(TWITTER_FOLLOWER_COUNT_MILLIONS)) %>%
  head(2)
```

```
## # A tibble: 2 x 15
##   PLAYER_NAME TEAM_ABBREVIATI~ AGE W_PCT OFF_RATING DEF_RATING NET_RATING
##   <chr>         <chr>         <dbl> <dbl>      <dbl>      <dbl>      <dbl>
## 1 LeBron Jam~ CLE             32 0.689      115.      107.       7.7
## 2 Kevin Dura~ GSW             28 0.823      117.      101.       16
## # ... with 8 more variables: AST_RATIO <dbl>, REB_PCT <dbl>,
## #   USG_PCT <dbl>, PIE <dbl>, SALARY_MILLIONS <dbl>,
## #   ACTIVE_TWITTER_LAST_YEAR <dbl>, TWITTER_FOLLOWER_COUNT_MILLIONS <dbl>,
## #   PTS <dbl>
```

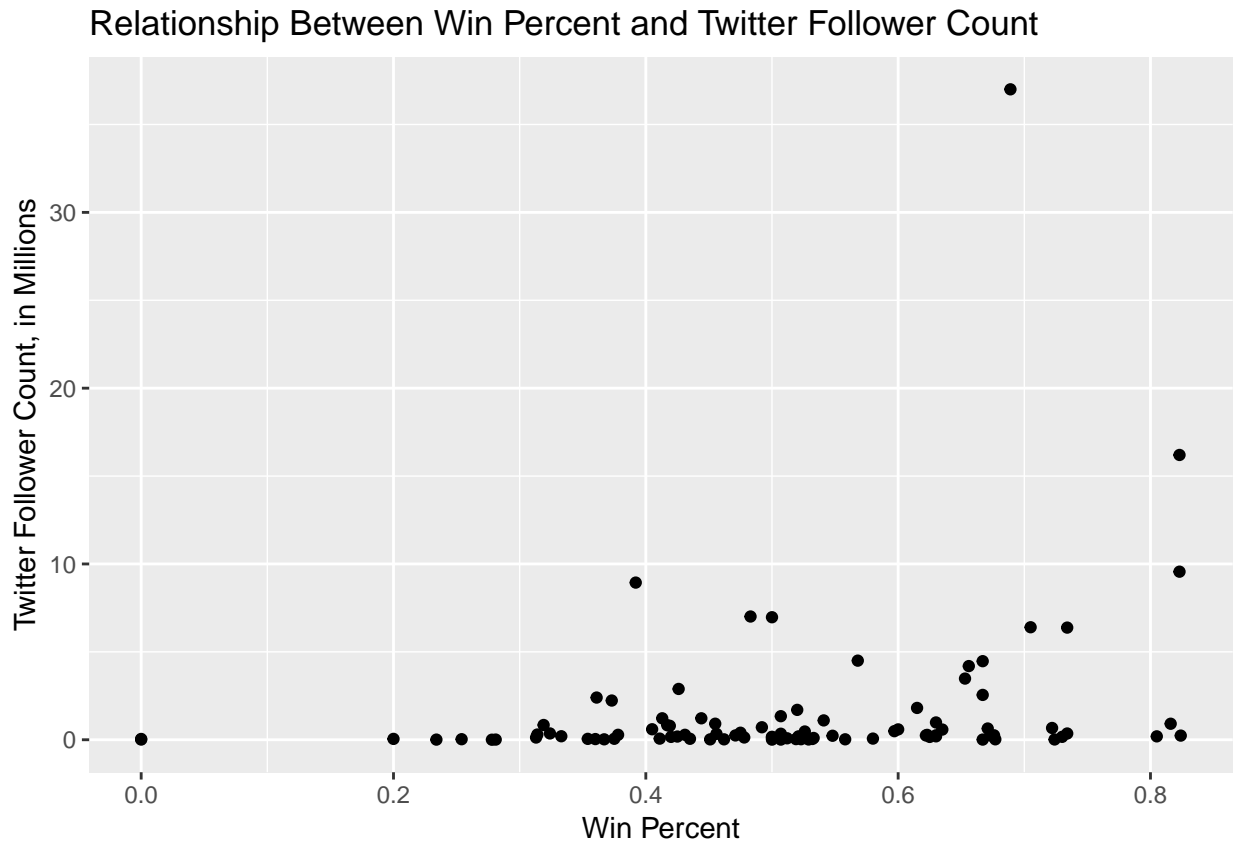
From the histogram, we can see that the distribution of Twitter follower counts is extremely right-skewed. The number of twitter followers ranges from .002 million to 37 million, with a mean of 1.6 million and a median of .246 million. There are two obvious outliers: Kevin Durant, with 16.2 million followers, and LeBron James, with 37 million followers.

## Bivariate

Next, we will do bivariate EDA, looking into the relationships of some of the predictor variables with the response variables. We won't do bivariate EDA on player name, Twitter handle, age, team abbreviation, or whether the players were active on Twitter last year, instead focusing on terms we believe may play a more nuanced / important role in predicting Twitter followers.

First, we'll look for a relationship between win percent and twitter followers in millions:

```
ggplot(nba_social_power_mod, aes(x = W_PCT, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +
  geom_point() +
  labs(title = "Relationship Between Win Percent and Twitter Follower Count",
       x = "Win Percent", y = "Twitter Follower Count, in Millions")
```

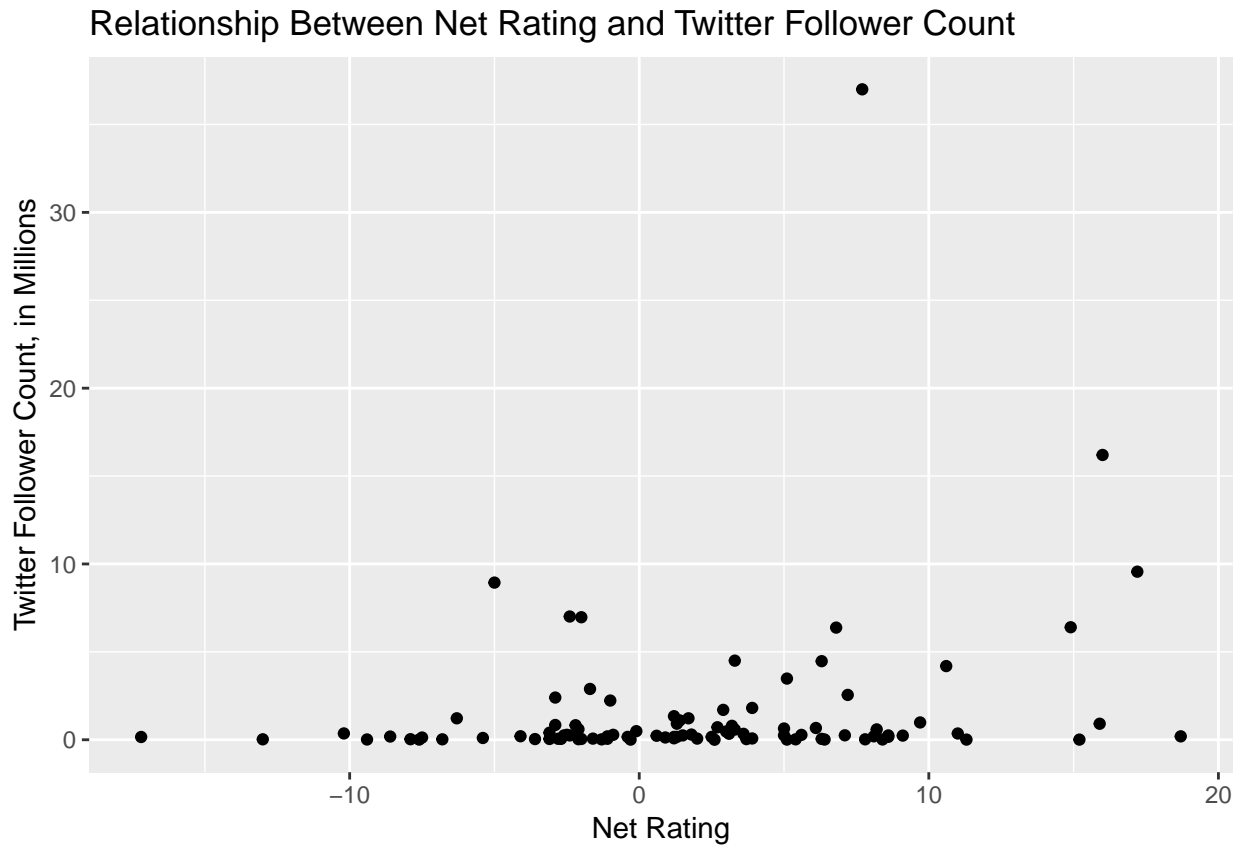


From the above plot, we can see that win percent and twitter follower count may have a very weak positive correlations. The players with significantly higher-than-average Twitter follower counts tend to have higher win percentages; however, this relationship is very weak.

Next, we'll look at the relationship between net rating, which combines offensive and defensive rating, and Twitter follower count:

```
ggplot(nba_social_power_mod, aes(x = NET_RATING, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +  
  geom_point() +  
  labs(title = "Relationship Between Net Rating and Twitter Follower Count",  
        x = "Net Rating", y = "Twitter Follower Count, in Millions")
```



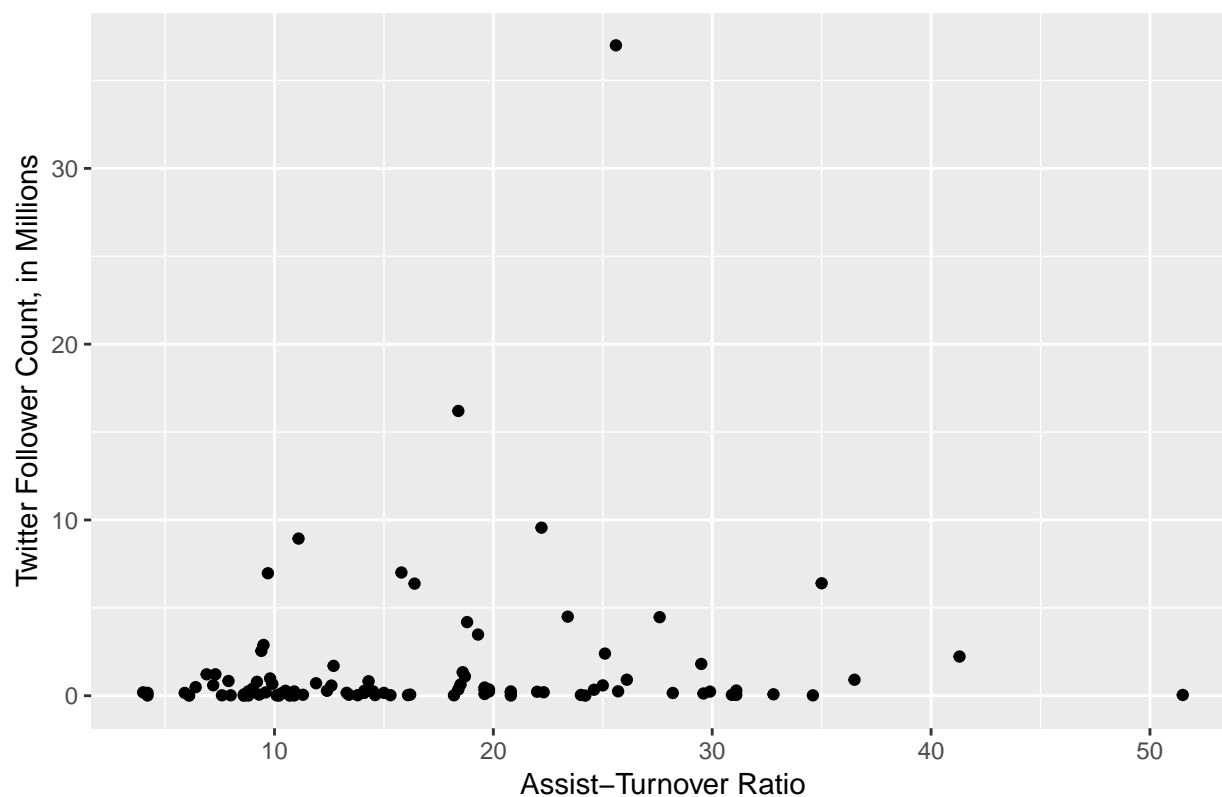


Similarly to win percent, Twitter follower count and net rating seem to have a very weak positive relationship. The players with the highest net ratings more often have higher-than-average Twitter follower counts, and more players with positive net ratings have high Twitter follower counts than those with negative net ratings.

Next, we'll look at the relationship between assist-to-turnovers ratio and Twitter follower count:

```
ggplot(nba_social_power_mod, aes(x = AST_RATIO, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +  
  geom_point() +  
  labs(title = "Relationship Between Assist-Turnover Ratio and Twitter Follower Count",  
        x = "Assist-Turnover Ratio", y = "Twitter Follower Count, in Millions")
```

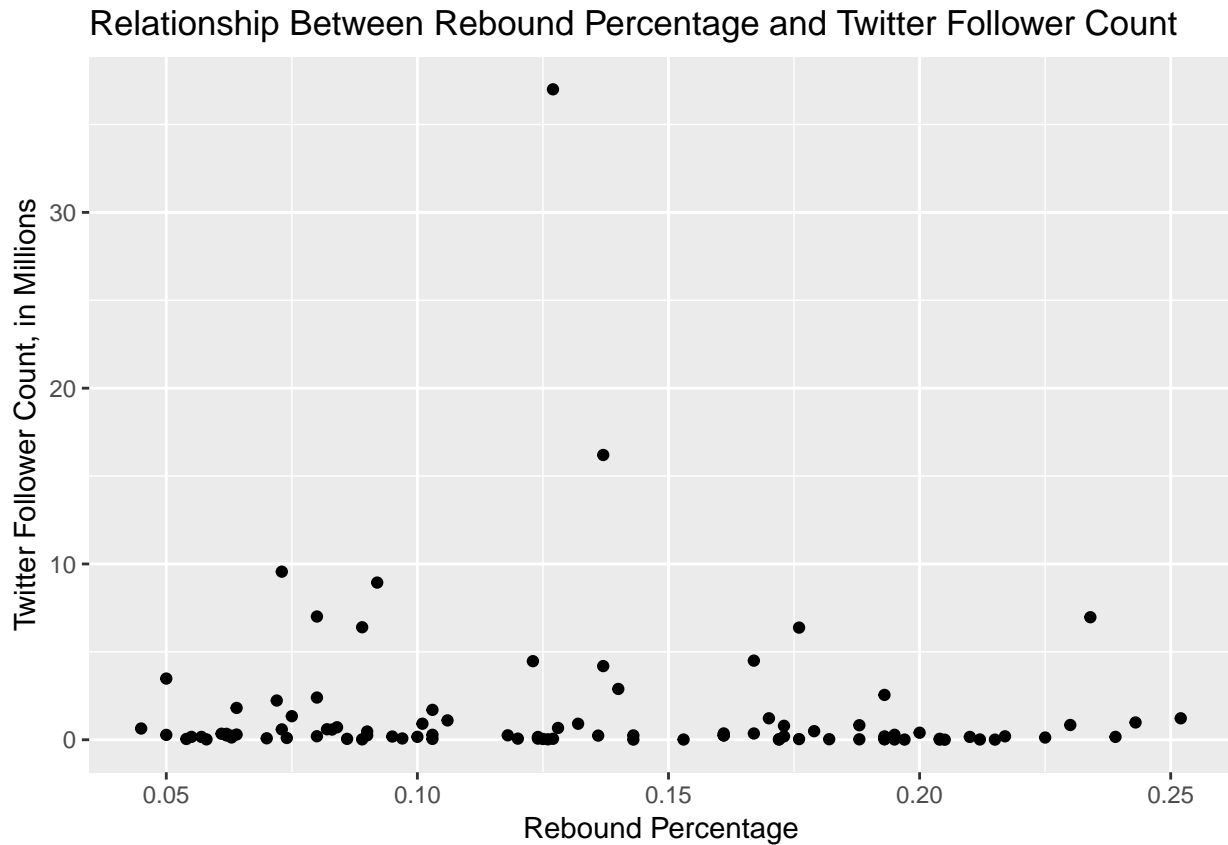
Relationship Between Assist–Turnover Ratio and Twitter Follower Count



There is no evident relationship between assist–turnover ratio and Twitter follower count.

Next, we'll examine whether there is a relationship between rebound percentage and Twitter follower count:

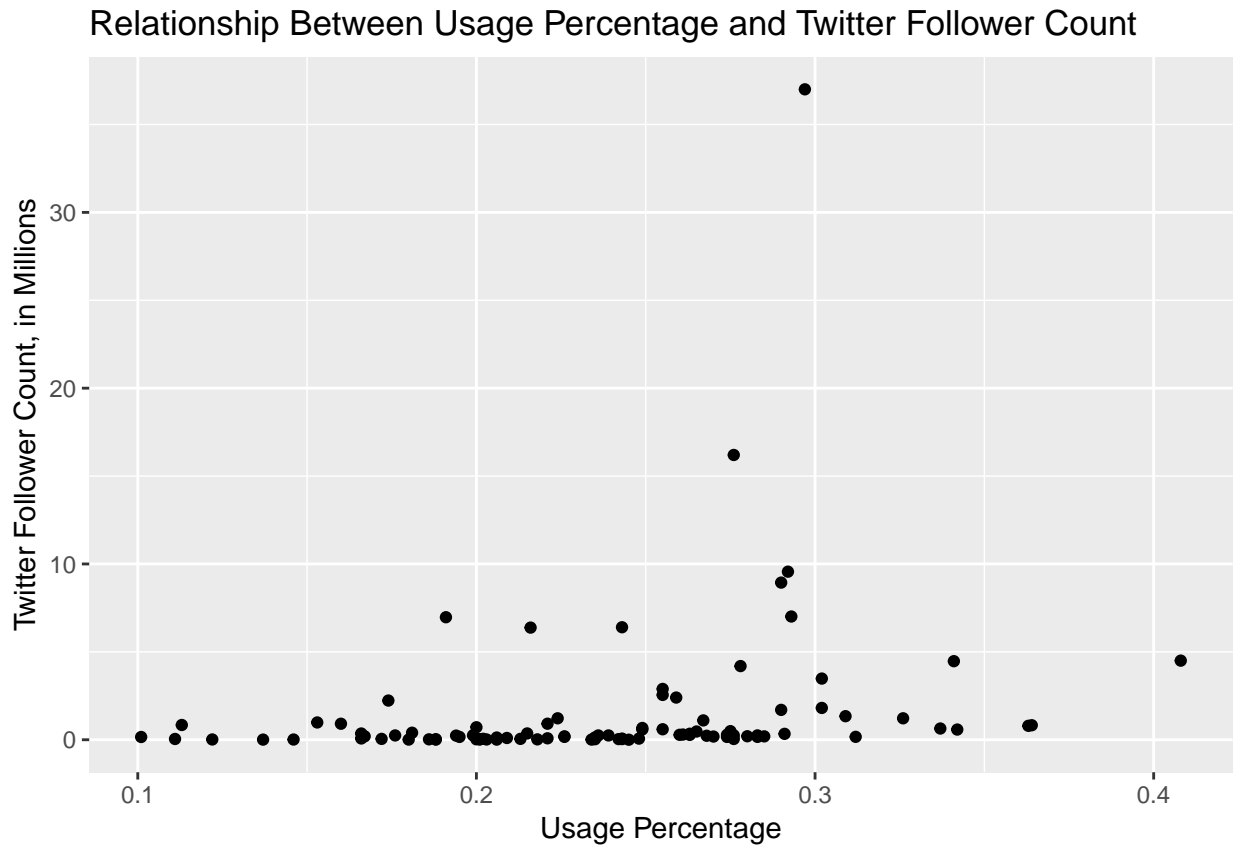
```
ggplot(nba_social_power_mod, aes(x = REB_PCT, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +  
  geom_point() +  
  labs(title = "Relationship Between Rebound Percentage and Twitter Follower Count",  
        x = "Rebound Percentage", y = "Twitter Follower Count, in Millions")
```



There is no evident relationship between rebound percentage and Twitter follower count.

Next, we'll examine whether there is a relationship between usage percentage and Twitter follower count:

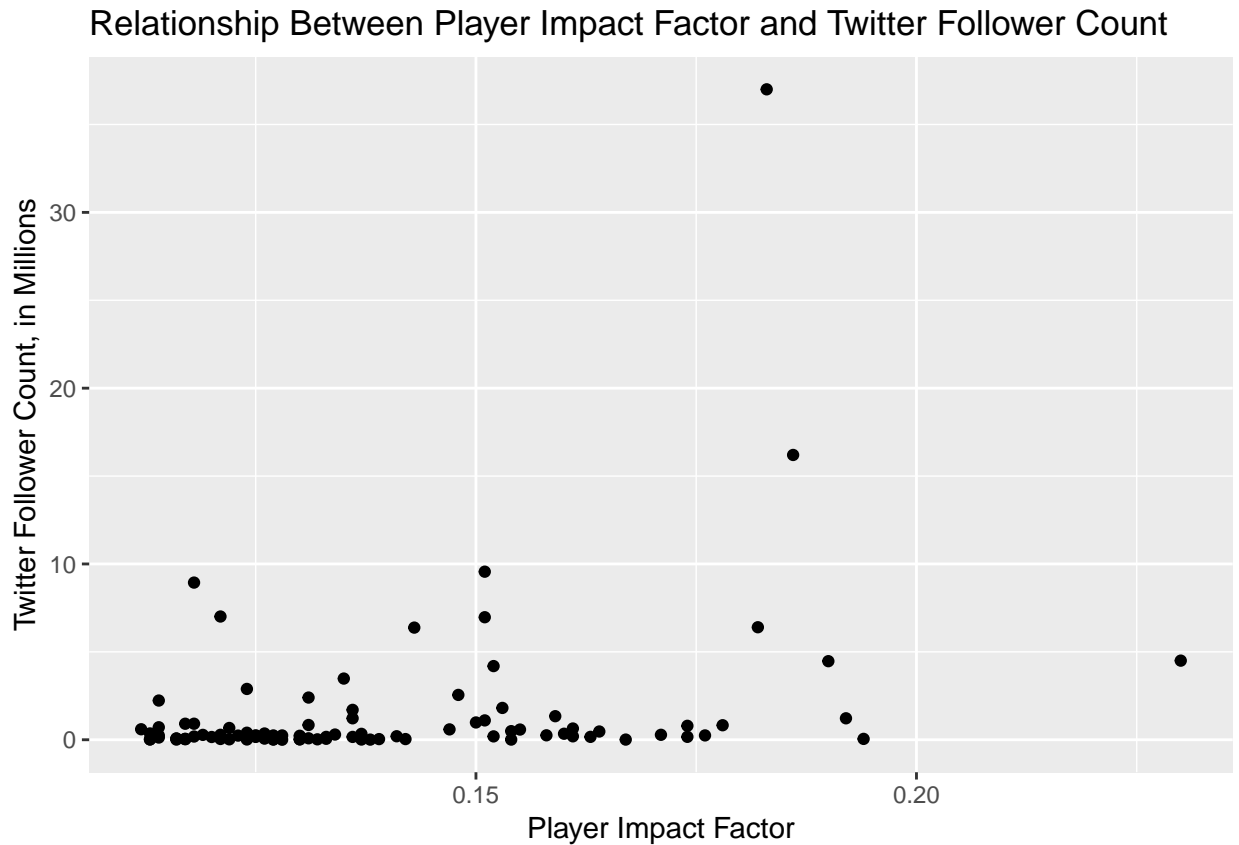
```
ggplot(nba_social_power_mod, aes(x = USG_PCT, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +  
  geom_point() +  
  labs(title = "Relationship Between Usage Percentage and Twitter Follower Count",  
        x = "Usage Percentage", y = "Twitter Follower Count, in Millions")
```



There appears to be a very weak positive correlation between usage percentage and Twitter follower count; players with a high usage percentage tend to have more Twitter followers, on average, than those with a lower usage percentage.

Next, we'll examine whether there is a relationship between player impact factor and Twitter follower count:

```
ggplot(nba_social_power_mod, aes(x = PIE, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +
  geom_point() +
  labs(title = "Relationship Between Player Impact Factor and Twitter Follower Count",
       x = "Player Impact Factor", y = "Twitter Follower Count, in Millions")
```

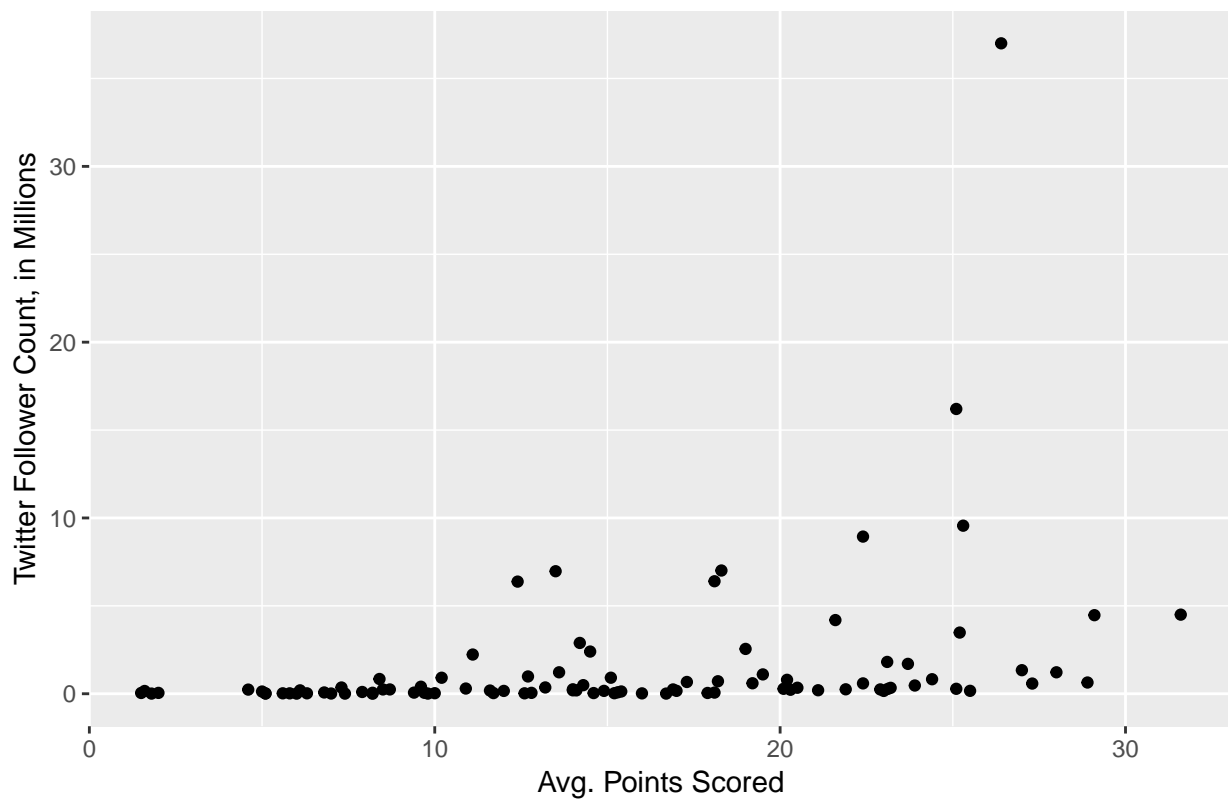


There appears to be a somewhat positive correlation between player impact factor and Twitter follower count; players with a high player impact factor tend to have more Twitter followers, on average, than those with a lower player impact factor.

Now, we'll look for a relationship between average points scored per game and Twitter follower count:

```
ggplot(nba_social_power_mod, aes(x = PTS, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +  
  geom_point() +  
  labs(title = "Relationship Between Avg. Points Scored and Twitter Follower Count",  
        x = "Avg. Points Scored", y = "Twitter Follower Count, in Millions")
```

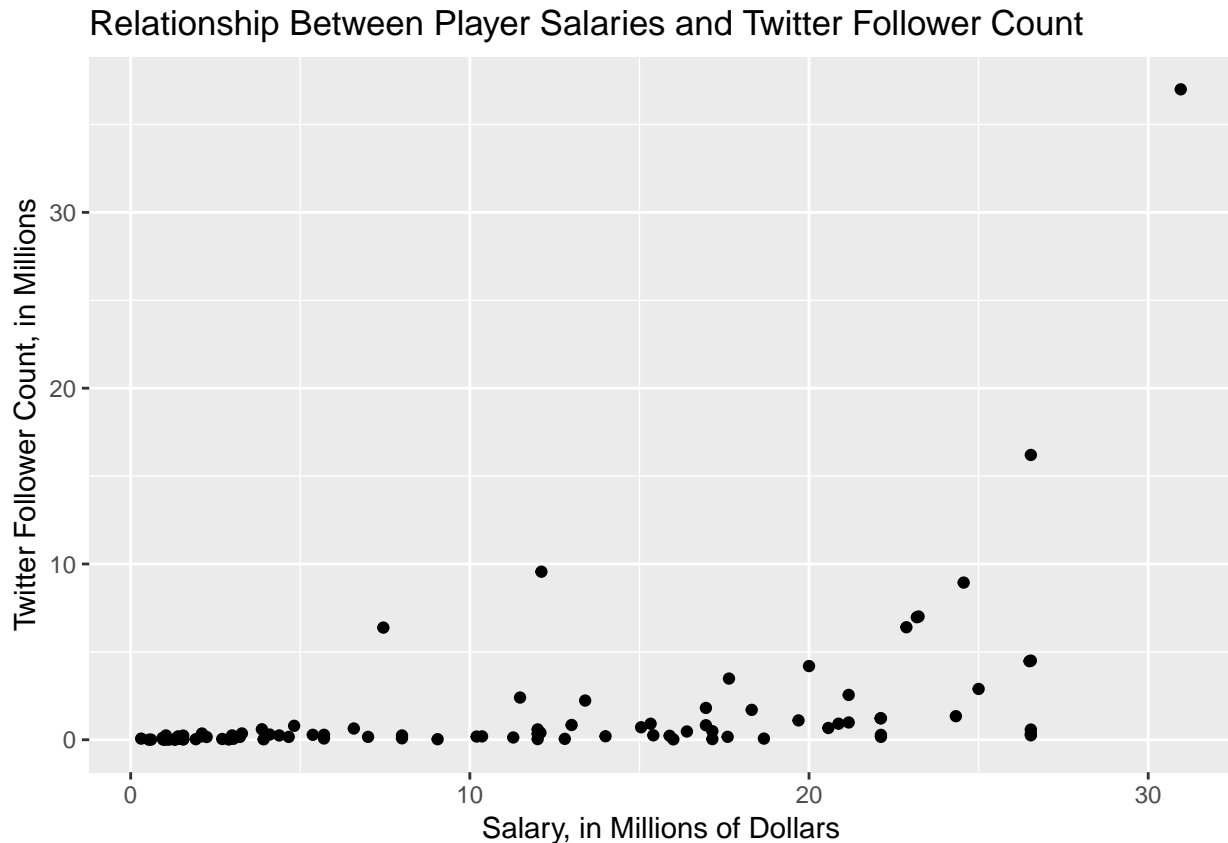
### Relationship Between Avg. Points Scored and Twitter Follower Count



There appears to be a weak positive correlation between average points scored and Twitter follower count; players with higher avg. points scored tend to have more Twitter followers than those with lower avg. points scored.

Finally, we'll look for a relationship between player salaries and points scored:

```
ggplot(nba_social_power_mod, aes(x = SALARY_MILLIONS, y = TWITTER_FOLLOWER_COUNT_MILLIONS)) +  
  geom_point() +  
  labs(title = "Relationship Between Player Salaries and Twitter Follower Count",  
        x = "Salary, in Millions of Dollars", y = "Twitter Follower Count, in Millions")
```



There appears to be a positive correlation between salary and Twitter follower count; the players with higher salaries tend to have higher Twitter follower counts.

In sum, there appear to be weak to moderate positive correlations between Twitter follower count and win percentage, net rating, usage percentage, player impact factor, avg. points per game, and player salaries; there is no obvious relationship between Twitter follower count and assist-to-turnover ratio or rebound percentage.

patchwork: multinomial logistic lecture code

## Multivariate Data Analysis

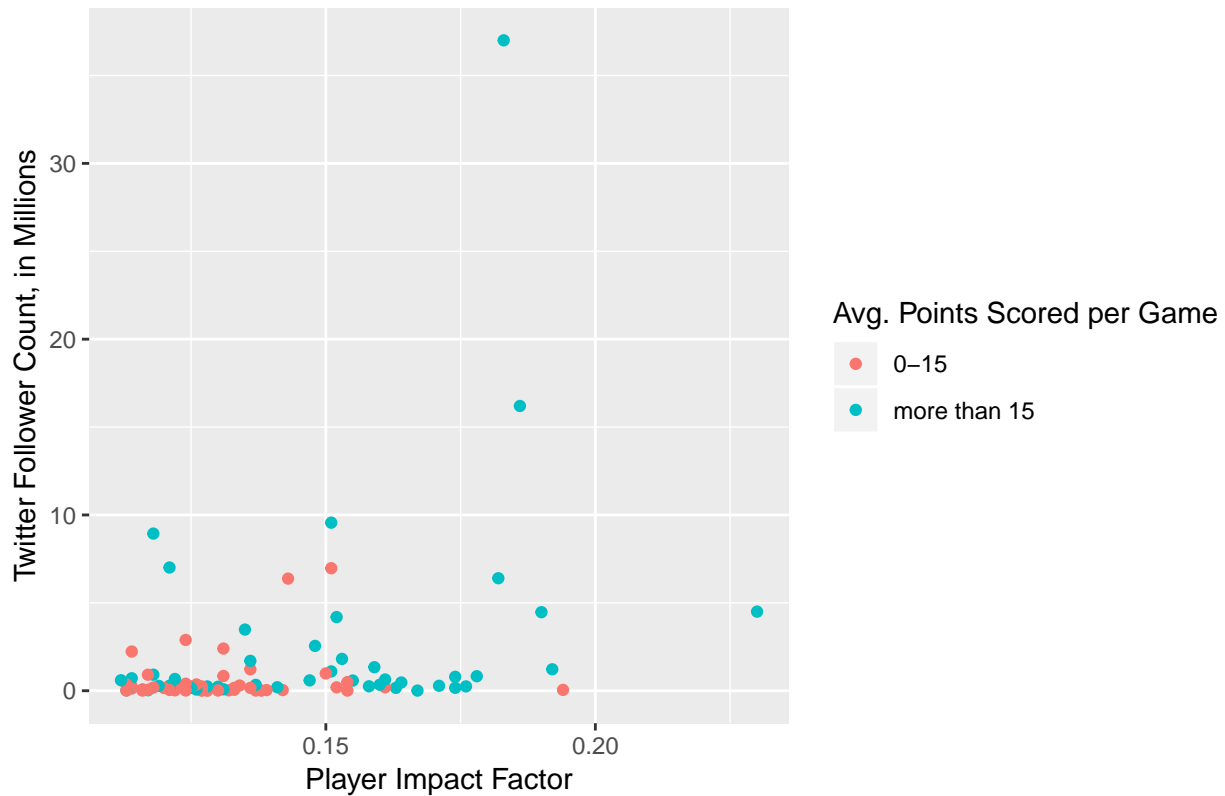
Now, we'll do some multivariate analysis. In this section, we are looking for predictor variables that may affect the way other predictor variables relate to the response variable.

First, we'll look to see if points scored affects the way player impact factor (PIE) relates to Twitter follower count:

```
nba_social_power_mod1 <- nba_social_power_mod %>%
  mutate(PTS_CAT = case_when(
    PTS <= 15 ~ "0-15" ,
    PTS > 15 ~ "more than 15"
  ))

ggplot(data = nba_social_power_mod1, aes(x = PIE, y = TWITTER_FOLLOWER_COUNT_MILLIONS, color = PTS_CAT)) +
  geom_point() +
  labs(x = "Player Impact Factor", y = "Twitter Follower Count, in Millions",
       color = "Avg. Points Scored per Game",
       title = "Relationship Between Pts Scored, Twitter Followers, and Player Impact Factor")
```

## Relationship Between Pts Scored, Twitter Followers, and Player Impact Fact



As we can see from the above color-coded boxplot, many of the players with the most points scored have higher player impact factors, and player impact factor values have a weak positive correlation with the Twitter follower count. This could be an opportunity for an interaction term.

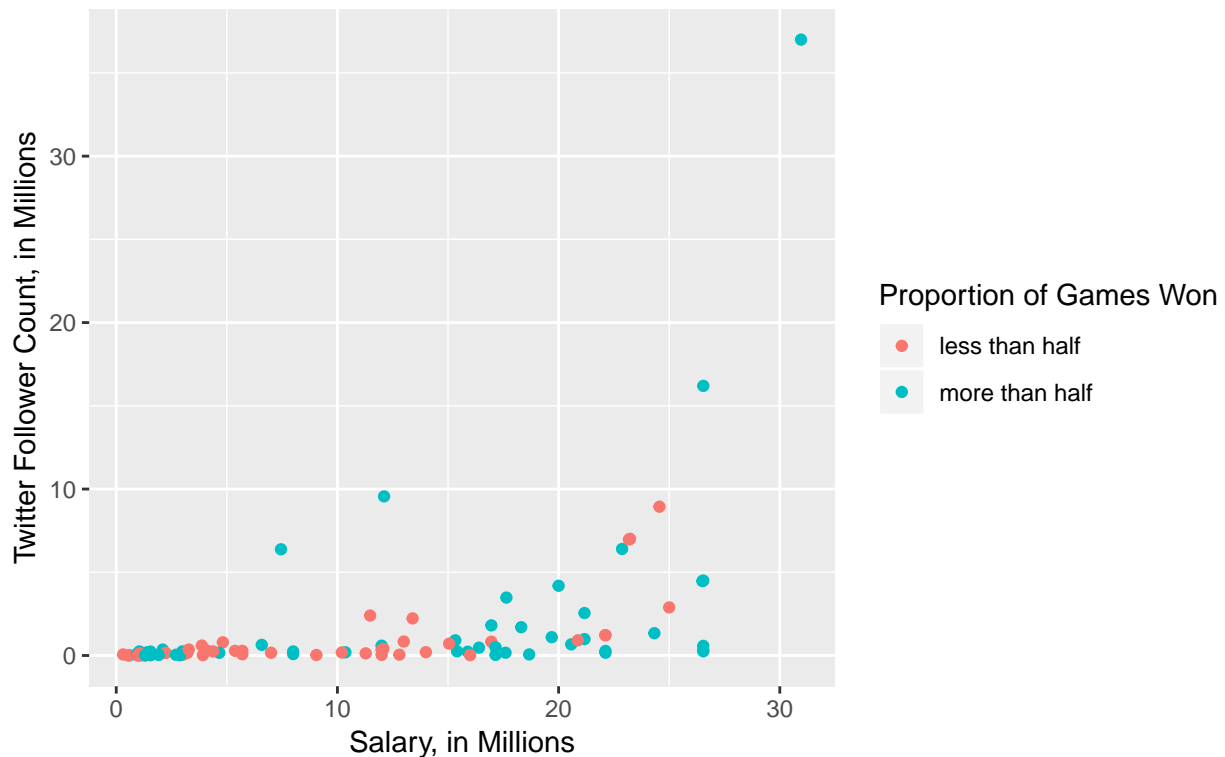
Next, we'll try to determine whether win percentage affects the way salary relates to the Twitter follower count:

```
nba_social_power_mod1 <- nba_social_power_mod %>%
  mutate(W_PCT_CAT = case_when(
    W_PCT <= .5 ~ "less than half" ,
    W_PCT > .5 ~ "more than half"
  ))
```

```
ggplot(data = nba_social_power_mod1, aes(x = SALARY_MILLIONS, y = TWITTER_FOLLOWER_COUNT_MILLIONS, color =
  W_PCT_CAT)) +
  geom_point() +
  labs(x = "Salary, in Millions", y = "Twitter Follower Count, in Millions",
    color = "Proportion of Games Won",
    title = "Relationship Between Win Percentage,
    Twitter Followers, and Salary")
```



## Relationship Between Win Percentage, Twitter Followers, and Salary



As we can see from the scatterplot, players with higher win percentages tend to be paid more, and high salary has a weak positive correlation with the Twitter follower count. This could also be an opportunity for an interaction term.

### Adjustments (Mean-Centering)

Before we proceed with the analysis, we will mean-center the quantitative variables:

```
nba_social_power_mod <- nba_social_power_mod %>%
  mutate(ageCent = AGE - mean(AGE),
         ast_ratioCent = AST_RATIO - mean(AST_RATIO),
         off_ratingCent = OFF_RATING - mean(OFF_RATING),
         def_ratingCent = DEF_RATING - mean(DEF_RATING),
         net_ratingCent = NET_RATING - mean(NET_RATING),
         pieCent = PIE - mean(PIE),
         reb_pctCent = REB_PCT - mean(REB_PCT),
         usg_pctCent = USG_PCT - mean(USG_PCT),
         salary_millionsCent = SALARY_MILLIONS - mean(SALARY_MILLIONS),
         w_pctCent = W_PCT - mean(W_PCT),
         ptsCent = PTS - mean(PTS))
```

### Modeling Approach

As our response, `TWITTER_FOLLOWER_COUNT_MILLIONS`, is a continuous numerical variable, we will use a multiple linear regression model (obviously, we plan to include two or more predictors, so simple linear regression would not apply!).

In regard to model selection, we will begin by fitting a multiple linear regression model with main effects, as well as interaction effects. We are considering `salary_millionsCent * w_pctCent` and `PIECent * ptsCent` as potential interaction terms for our model since multivariate EDA highlighted strong positive relationships between win percentage and player salaries and points scored and PIE (player impact factor).

Next, we will perform two iterations of backward selection on the first model: (i) using AIC as the selection criterion and (ii) using adjusted R-squared as the selection criterion. We decided against trying BIC as the selection criterion because we would prefer more terms in the final model as our objective is to predict the Twitter follower counts of NBA players using measures of athletic success (and predictions are generally more accurate with more relevant predictor variables).

After completing the two iterations of backward selection, we will compare the resulting models and see whether certain terms are removed in both (which would suggest those terms are not particularly useful).

After obtaining a final model, we will examine the impact of including versus excluding prominent athletes (e.g. LeBron James) in our model since our objective is to design a model with the best predictive accuracy. We choose to analyze this at the end of the model selection phase because prominent athletes are also included in the population we want to understand when fitting the multiple linear regression model.

## Model 1

We will simply fit a model with eleven main effect terms—`ageCent`, `ast_ratioCent`, `off_ratingCent`, `def_ratingCent`, `PIECent`, `reb_pctCent`, `usg_pctCent`, `salary_millionsCent`, `w_pctCent`, `ptsCent`, and `ACTIVE_TWITTER_LAST_YEAR`—and two interaction terms (`salary_millionsCent * w_pctCent` and `PIECent * ptsCent`). We decided to use `off_ratingCent` and `def_ratingCent` as opposed to `net_ratingCent` because the two individual components potentially provide more important information in predicting the number of Twitter followers.

```
m1 <- lm(TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent +
  ast_ratioCent +
  off_ratingCent +
  def_ratingCent +
  PIECent +
  reb_pctCent +
  usg_pctCent +
  salary_millionsCent +
  w_pctCent +
  ptsCent +
  ACTIVE_TWITTER_LAST_YEAR +
  salary_millionsCent * w_pctCent +
  PIECent * ptsCent,
  data = nba_social_power_mod)
tidy(m1, conf.int = TRUE) %>%
  kable(format = "markdown", digits = 3)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-2.386	2.744	-0.869	0.387	-7.845	3.074
ageCent	0.262	0.123	2.131	0.036	0.017	0.507
ast_ratioCent	0.017	0.064	0.269	0.788	-0.109	0.144
off_ratingCent	0.026	0.106	0.247	0.805	-0.185	0.237
def_ratingCent	0.086	0.118	0.732	0.467	-0.149	0.321
PIECent	13.104	27.617	0.474	0.636	-41.846	68.053
reb_pctCent	7.689	12.120	0.634	0.528	-16.426	31.803
usg_pctCent	0.052	14.017	0.004	0.997	-27.838	27.942
salary_millionsCent	0.086	0.070	1.240	0.219	-0.052	0.225

term	estimate	std.error	statistic	p.value	conf.low	conf.high
w_pctCent	7.095	3.834	1.850	0.068	-0.534	14.725
ptsCent	0.063	0.132	0.481	0.632	-0.199	0.325
ACTIVE_TWITTER_LAST_YEAR	3.523	2.758	1.278	0.205	-1.964	9.011
salary_millionsCent:w_pctCent	0.999	0.344	2.901	0.005	0.314	1.684
PIECent:ptsCent	1.293	2.293	0.564	0.574	-3.269	5.855

Based on the output, the equation of the linear model is:  $\text{TWITTER\_FOLLOWER\_COUNT\_MILLIONS-hat} = -2.386 + 0.262 * \text{ageCent} + 0.017 * \text{ast\_ratioCent} + 0.026 * \text{off\_ratingCent} + 0.086 * \text{def\_ratingCent} + 13.104 * \text{PIECent} + 7.689 * \text{reb\_pctCent} + 0.052 * \text{usg\_pctCent} + 0.086 * \text{salary\_millionsCent} + 7.095 * \text{w\_pctCent} + 0.063 * \text{ptsCent} + 3.523 * \text{ACTIVE\_TWITTER\_LAST\_YEAR} + 0.999 * \text{salary\_millionsCent} * \text{w\_pct\_Cent} + 1.293 * \text{PIECent} * \text{ptsCent}.$

### Backward Selection (Iteration 1)

We will now perform the first iteration of backward selection using AIC as the selection criterion:

```
m1_aic <- step(m1, direction = "backward")
```

```
## Start: AIC=265.14
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + ast_ratioCent + off_ratingCent +
##   def_ratingCent + PIECent + reb_pctCent + usg_pctCent + salary_millionsCent +
##   w_pctCent + ptsCent + ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent *
##   w_pctCent + PIECent * ptsCent
##
##               Df Sum of Sq    RSS    AIC
## - usg_pctCent    1      0.000 1153.0 263.14
## - off_ratingCent  1      0.869 1153.8 263.21
## - ast_ratioCent   1      1.033 1154.0 263.23
## - PIECent:ptsCent  1      4.525 1157.5 263.51
## - reb_pctCent     1      5.729 1158.7 263.61
## - def_ratingCent  1      7.619 1160.6 263.77
## - ACTIVE_TWITTER_LAST_YEAR  1     23.232 1176.2 265.04
## <none>                                1153.0 265.14
## - ageCent         1     64.667 1217.6 268.32
## - salary_millionsCent:w_pctCent  1    119.780 1272.7 272.53
##
## Step: AIC=263.14
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + ast_ratioCent + off_ratingCent +
##   def_ratingCent + PIECent + reb_pctCent + salary_millionsCent +
##   w_pctCent + ptsCent + ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent:w_pctCent +
##   PIECent:ptsCent
##
##               Df Sum of Sq    RSS    AIC
## - off_ratingCent  1      0.879 1153.8 261.21
## - ast_ratioCent   1      1.119 1154.1 261.23
## - PIECent:ptsCent  1      4.673 1157.6 261.52
## - reb_pctCent     1      6.349 1159.3 261.66
## - def_ratingCent  1      7.812 1160.8 261.78
## - ACTIVE_TWITTER_LAST_YEAR  1     23.398 1176.4 263.05
## <none>                                1153.0 263.14
## - ageCent         1     65.008 1218.0 266.35
## - salary_millionsCent:w_pctCent  1    120.208 1273.2 270.56
```

```

##
## Step: AIC=261.21
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + ast_ratioCent + def_ratingCent +
##     PIECent + reb_pctCent + salary_millionsCent + w_pctCent +
##     ptsCent + ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent:w_pctCent +
##     PIECent:ptsCent
##
##
## Df Sum of Sq RSS AIC
## - ast_ratioCent 1 1.551 1155.4 259.34
## - PIECent:ptsCent 1 3.809 1157.7 259.53
## - reb_pctCent 1 6.572 1160.4 259.75
## - def_ratingCent 1 10.226 1164.1 260.05
## - ACTIVE_TWITTER_LAST_YEAR 1 23.162 1177.0 261.10
## <none> 1153.8 261.21
## - ageCent 1 64.179 1218.0 264.36
## - salary_millionsCent:w_pctCent 1 144.592 1298.4 270.43
##
## Step: AIC=259.34
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + def_ratingCent +
##     PIECent + reb_pctCent + salary_millionsCent + w_pctCent +
##     ptsCent + ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent:w_pctCent +
##     PIECent:ptsCent
##
##
## Df Sum of Sq RSS AIC
## - PIECent:ptsCent 1 5.074 1160.5 257.76
## - reb_pctCent 1 5.330 1160.7 257.78
## - def_ratingCent 1 10.283 1165.7 258.18
## - ACTIVE_TWITTER_LAST_YEAR 1 24.144 1179.5 259.31
## <none> 1155.4 259.34
## - ageCent 1 64.470 1219.9 262.50
## - salary_millionsCent:w_pctCent 1 156.501 1311.9 269.41
##
## Step: AIC=257.76
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + def_ratingCent +
##     PIECent + reb_pctCent + salary_millionsCent + w_pctCent +
##     ptsCent + ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent:w_pctCent
##
##
## Df Sum of Sq RSS AIC
## - ptsCent 1 6.132 1166.6 256.26
## - reb_pctCent 1 8.021 1168.5 256.41
## - def_ratingCent 1 10.161 1170.6 256.58
## - PIECent 1 12.392 1172.9 256.77
## - ACTIVE_TWITTER_LAST_YEAR 1 23.951 1184.4 257.70
## <none> 1160.5 257.76
## - ageCent 1 70.150 1230.6 261.33
## - salary_millionsCent:w_pctCent 1 164.153 1324.6 268.32
##
## Step: AIC=256.26
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + def_ratingCent +
##     PIECent + reb_pctCent + salary_millionsCent + w_pctCent +
##     ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent:w_pctCent
##
##
## Df Sum of Sq RSS AIC
## - reb_pctCent 1 3.143 1169.7 254.51

```

```
## - def_ratingCent      1      11.861 1178.5 255.22
## <none>                  1166.6 256.26
## - ACTIVE_TWITTER_LAST_YEAR 1      25.646 1192.2 256.32
## - PIECent              1      31.477 1198.1 256.79
## - ageCent              1      64.560 1231.2 259.37
## - salary_millionsCent:w_pctCent 1 158.085 1324.7 266.33
##
## Step: AIC=254.51
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + def_ratingCent +
##   PIECent + salary_millionsCent + w_pctCent + ACTIVE_TWITTER_LAST_YEAR +
##   salary_millionsCent:w_pctCent
##
##              Df Sum of Sq  RSS    AIC
## - def_ratingCent      1      11.079 1180.8 253.41
## <none>                  1169.7 254.51
## - ACTIVE_TWITTER_LAST_YEAR 1      25.350 1195.1 254.55
## - PIECent              1      36.705 1206.5 255.45
## - ageCent              1      62.548 1232.3 257.46
## - salary_millionsCent:w_pctCent 1 155.368 1325.1 264.36
##
## Step: AIC=253.41
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + PIECent + salary_millionsCent +
##   w_pctCent + ACTIVE_TWITTER_LAST_YEAR + salary_millionsCent:w_pctCent
##
##              Df Sum of Sq  RSS    AIC
## - ACTIVE_TWITTER_LAST_YEAR 1      25.051 1205.9 253.40
## <none>                  1180.8 253.41
## - PIECent              1      32.005 1212.8 253.95
## - ageCent              1      54.622 1235.4 255.70
## - salary_millionsCent:w_pctCent 1 162.024 1342.8 263.62
##
## Step: AIC=253.4
## TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent + PIECent + salary_millionsCent +
##   w_pctCent + salary_millionsCent:w_pctCent
##
##              Df Sum of Sq  RSS    AIC
## <none>                  1205.9 253.40
## - PIECent              1      29.646 1235.5 253.71
## - ageCent              1      60.024 1265.9 256.02
## - salary_millionsCent:w_pctCent 1 158.457 1364.3 263.13
tidy(m1_aic, conf.int = TRUE) %>%
  kable(format = "markdown", digits = 3)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	1.175	0.397	2.957	0.004	0.386	1.964
ageCent	0.225	0.107	2.105	0.038	0.013	0.437
PIECent	27.461	18.565	1.479	0.143	-9.427	64.349
salary_millionsCent	0.111	0.051	2.158	0.034	0.009	0.213
w_pctCent	6.294	2.761	2.280	0.025	0.808	11.780
salary_millionsCent:w_pctCent	1.019	0.298	3.420	0.001	0.427	1.611

Based on the output displayed above from the first iteration of backward selection (using AIC as the selec-

tion criterion), seven main effect terms (`ast_ratioCent`, `off_ratingCent`, `def_ratingCent`, `reb_pctCent`, `usg_pctCent`, `ptsCent`, and `ACTIVE_TWITTER_LAST_YEAR`) and one interaction term (`PIECent * ptsCent`) were removed.

Hence, the equation of the selected linear model is:  $\text{TWITTER\_FOLLOWER\_COUNT\_MILLIONS-hat} = 1.175 + 0.225 * \text{ageCent} + 27.461 * \text{PIECent} + 0.111 * \text{salary\_millionsCent} + 6.294 * \text{w\_pctCent} + 1.019 * \text{salary\_millionsCent} * \text{w\_pctCent}$ .

However, `PIECent`, which has the highest slope coefficient in the linear model, has a high standard error and a high p-value (0.143). Furthermore, the confidence interval for the slope coefficient is extremely wide (-9.427 to 64.349), not to mention the fact that it includes zero. So, we can reasonably infer `PIECent` may be a particularly troublesome predictor in the model (especially since it appears to not significantly affect the response, `TWITTER_FOLLOWER_COUNT_MILLIONS`).

We will proceed with the second iteration of backward selection and will revisit the issue of `PIECent` after viewing the selected linear regression model based on adjusted R-squared.

## Backward Selection (Iteration 2)

Next, we will perform the second iteration of backward selection using adjusted R-squared as the selection criterion:

```
m1_adjR <- regsubsets(TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent +
  ast_ratioCent +
  off_ratingCent +
  def_ratingCent +
  PIECent +
  reb_pctCent +
  usg_pctCent +
  salary_millionsCent +
  w_pctCent +
  ptsCent +
  ACTIVE_TWITTER_LAST_YEAR +
  salary_millionsCent * w_pctCent +
  PIECent * ptsCent,
  data = nba_social_power_mod,
  method = "backward")

select_summary <- summary(m1_adjR)
coef(m1_adjR, which.max(select_summary$adjr2))
```

```
##              (Intercept)              ageCent
##              -2.3794214              0.2146898
##              PIECent              salary_millionsCent
##              28.5598883              0.1148099
##              w_pctCent              ACTIVE_TWITTER_LAST_YEAR
##              6.6375685              3.6257226
## salary_millionsCent:w_pctCent
##              1.0306857
```

Based on the output displayed above from the second iteration of backward selection (using adjusted R-squared as the selection criterion), six main effect terms (`ast_ratioCent`, `off_ratingCent`, `def_ratingCent`, `reb_pctCent`, `usg_pctCent`, and `ptsCent`) and one interaction term (`PIECent * ptsCent`) were removed.

Hence, the equation of the selected linear model is:  $\text{TWITTER\_FOLLOWER\_COUNT\_MILLIONS-hat} = -2.379 + 0.215 * \text{ageCent} + 28.560 * \text{PIECent} + 0.115 * \text{salary\_millionsCent} + 6.638 * \text{w\_pctCent} + 3.626 *$

```
ACTIVE_TWITTER_LAST_YEAR + 1.031 * salary_millionsCent * w_pct_Cent.
```

### Model Comparison: AIC vs. Adjusted R-squared

We notice that the model selected using adjusted R-squared as the selection criterion included an additional term (`ACTIVE_TWITTER_LAST_YEAR`) which was omitted in the model selected based on the first iteration of backward selection (using AIC as the selection criterion).

If we closely investigate the results of the `step()` function, we can see that `ACTIVE_TWITTER_LAST_YEAR` was the last predictor removed during the first iteration of backward selection. In fact, the difference in AIC between the final two steps is only  $253.41 - 253.4 = 0.01$ . Again, recalling our objective in this analysis—to predict Twitter follower counts of NBA players—we choose to leave `ACTIVE_TWITTER_LAST_YEAR` in the final model because the information provided by this categorical variable—whether the player maintained an active Twitter account during the 2015-2016 season—is clearly relevant to the response, `TWITTER_FOLLOWER_COUNT_MILLIONS` (which measures players’ total Twitter follower counts). Also, when designing a model for predictive purposes, including more predictors (which are relevant) is generally accepted. However, we must consider how `TWITTER_FOLLOWER_COUNT_MILLIONS` will affect our model interpretations since it has a high p-value and a wide confidence interval including zero.

Now, as promised, we return to the issue of `PIECent`. Unfortunately, both selected models included `PIECent`; thus, we must decide whether to keep the variable in the model, or ignore the results from the two iterations of backward selection and remove it.

To answer this question, we will compare the AIC and adjusted R-squared values for the model selected based on the second iteration of backward selection (`m1_adjR`, which includes `ACTIVE_TWITTER_LAST_YEAR`) and the same model except without `PIECent`.

```
m1_modV1 <- lm(TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent +
  salary_millionsCent +
  PIECent +
  w_pctCent +
  ACTIVE_TWITTER_LAST_YEAR +
  salary_millionsCent * w_pctCent,
  data = nba_social_power_mod)
```

```
m1_modV2 <- lm(TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent +
  salary_millionsCent +
  w_pctCent +
  ACTIVE_TWITTER_LAST_YEAR +
  salary_millionsCent * w_pctCent,
  data = nba_social_power_mod)
```

```
glance(m1_modV1)$AIC
```

```
## [1] 525.0067
```

```
glance(m1_modV1)$adj.r.squared
```

```
## [1] 0.3207863
```

```
glance(m1_modV2)$AIC
```

```
## [1] 525.5473
```

```
glance(m1_modV2)$adj.r.squared
```

```
## [1] 0.3102152
```

Based on the output of the `glance()` function, the AIC of the model with `PIECent` is 525.0067. Conversely, the AIC of the model without `PIECent` is 525.5473.

Moreover, the adjusted R-squared value for the model with `PIECent` is roughly 0.321, whereas the adjusted R-squared value for the model without `PIECent` is about 0.310.

Therefore, the model with `PIECent` maximizes adjusted R-squared and minimizes AIC. Hence, despite our qualms with the large standard error and high p-value, we will keep `PIECent` in the model. Of course, we will need to consider this information carefully when interpreting our model and evaluating its predictive capabilities.

## Final Model

Thus, our final model is:

```
tidy(m1_modV1, conf.int = TRUE) %>%
  kable(format = "markdown", digits = 3)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-2.379	2.631	-0.904	0.368	-7.608	2.850
ageCent	0.215	0.106	2.018	0.047	0.003	0.426
salary__millionsCent	0.115	0.051	2.244	0.027	0.013	0.216
PIECent	28.560	18.493	1.544	0.126	-8.190	65.310
w_pctCent	6.638	2.759	2.406	0.018	1.155	12.120
ACTIVE_TWITTER_LAST_YEAR	3.626	2.654	1.366	0.175	-1.648	8.899
salary__millionsCent:w_pctCent	1.031	0.297	3.475	0.001	0.441	1.620

So, the equation of the final model is:  $\text{TWITTER\_FOLLOWER\_COUNT\_MILLIONS-hat} = -2.379 + 0.215 * \text{ageCent} + 28.560 * \text{PIECent} + 0.115 * \text{salary\_millionsCent} + 6.638 * \text{w\_pctCent} + 3.626 * \text{ACTIVE\_TWITTER\_LAST\_YEAR} + 1.031 * \text{salary\_millionsCent} * \text{w\_pct\_Cent}$ .

## Impact of Prominent Players

Lastly, before proceeding to discuss assumptions, we would like to examine the impact of including versus excluding prominent athletes (e.g. LeBron James) in our model:

```
nba_social_power_mod %>%
  arrange(desc(SALARY_MILLIONS)) %>%
  select(PLAYER_NAME, SALARY_MILLIONS) %>%
  head(10)
```

```
## # A tibble: 10 x 2
##   PLAYER_NAME      SALARY_MILLIONS
##   <chr>           <dbl>
## 1 LeBron James    31.0
## 2 Russell Westbrook 26.5
## 3 Kevin Durant    26.5
## 4 Mike Conley      26.5
## 5 DeMar DeRozan    26.5
## 6 Al Horford       26.5
## 7 James Harden     26.5
## 8 Dirk Nowitzki     25
## 9 Carmelo Anthony  24.6
```



```
## 10 Damian Lillard                24.3
```

```
nba_social_power_mod %>%  
  arrange(desc(TWITTER_FOLLOWER_COUNT_MILLIONS)) %>%  
  select(PLAYER_NAME, TWITTER_FOLLOWER_COUNT_MILLIONS) %>%  
  head(10)
```

```
## # A tibble: 10 x 2  
##   PLAYER_NAME      TWITTER_FOLLOWER_COUNT_MILLIONS  
##   <chr>                <dbl>  
## 1 LeBron James          37  
## 2 Kevin Durant         16.2  
## 3 Stephen Curry         9.56  
## 4 Carmelo Anthony       8.94  
## 5 Dwyane Wade         7.01  
## 6 Dwight Howard         6.97  
## 7 Chris Paul           6.4  
## 8 Pau Gasol             6.38  
## 9 Russell Westbrook     4.5  
## 10 James Harden         4.47
```

Based on the tables above, it is clear LeBron James is an outlier, both in regard to his annual salary and Twitter follower count. So, we will remove LeBron from the dataset and see how the model changes:

```
nba_social_power_mod2 <- nba_social_power_mod %>%  
  filter(PLAYER_NAME != "LeBron James")  
  
glimpse(nba_social_power_mod2)
```

```
## Observations: 94  
## Variables: 26  
## $ PLAYER_NAME      <chr> "Russell Westbrook", "Demetriu...  
## $ TEAM_ABBREVIATION <chr> "OKC", "BOS", "NOP", "HOU", "G...  
## $ AGE              <dbl> 28, 22, 24, 27, 28, 32, 26, 22...  
## $ W_PCT            <dbl> 0.568, 0.200, 0.413, 0.667, 0....  
## $ OFF_RATING       <dbl> 107.9, 124.2, 104.2, 113.6, 11...  
## $ DEF_RATING       <dbl> 104.6, 117.8, 102.5, 107.3, 10...  
## $ NET_RATING       <dbl> 3.3, 6.3, 1.7, 6.3, 16.0, 14.9...  
## $ AST_RATIO        <dbl> 23.4, 31.1, 7.3, 27.6, 18.4, 3...  
## $ REB_PCT          <dbl> 0.167, 0.103, 0.170, 0.123, 0....  
## $ USG_PCT          <dbl> 0.408, 0.172, 0.326, 0.341, 0....  
## $ PIE              <dbl> 0.230, 0.194, 0.192, 0.190, 0....  
## $ SALARY_MILLIONS  <dbl> 26.54, 1.45, 22.12, 26.50, 26....  
## $ ACTIVE_TWITTER_LAST_YEAR <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, ...  
## $ TWITTER_FOLLOWER_COUNT_MILLIONS <dbl> 4.500, 0.049, 1.220, 4.470, 16...  
## $ PTS              <dbl> 31.6, 2.0, 28.0, 29.1, 25.1, 1...  
## $ ageCent          <dbl> 0.6105263, -5.3894737, -3.3894...  
## $ ast_ratioCent    <dbl> 6.274737, 13.974737, -9.825263...  
## $ off_ratingCent   <dbl> -0.009473684, 16.290526316, -3...  
## $ def_ratingCent   <dbl> -1.39368421, 11.80631579, -3.4...  
## $ net_ratingCent   <dbl> 1.3810526, 4.3810526, -0.21894...  
## $ PIECent          <dbl> 0.09073684, 0.05473684, 0.0527...  
## $ reb_pctCent      <dbl> 0.033778947, -0.030221053, 0.0...  
## $ usg_pctCent      <dbl> 0.170, -0.066, 0.088, 0.103, 0...  
## $ salary_millionsCent <dbl> 15.2351368, -9.8548632, 10.815...  
## $ w_pctCent        <dbl> 0.056589474, -0.311410526, -0....
```

```
## $ ptsCent <dbl> 16.3176842, -13.2823158, 12.71...
```

Based on the `glimpse()` output, we can see LeBron has been removed from the dataset (94 observations remaining).

Now, to assess the impact of his absence on the final model:

```
m1_modV3 <- lm(TWITTER_FOLLOWER_COUNT_MILLIONS ~ ageCent +
  salary_millionsCent +
  PIECent +
  w_pctCent +
  ACTIVE_TWITTER_LAST_YEAR +
  salary_millionsCent * w_pctCent,
  data = nba_social_power_mod2)

tidy(m1_modV3, conf.int = TRUE) %>%
  kable(format = "markdown", digits = 3)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-1.358	1.485	-0.915	0.363	-4.310	1.594
ageCent	0.128	0.060	2.129	0.036	0.009	0.248
salary_millionsCent	0.082	0.029	2.844	0.006	0.025	0.140
PIECent	16.105	10.463	1.539	0.127	-4.691	36.902
w_pctCent	4.081	1.566	2.606	0.011	0.968	7.194
ACTIVE_TWITTER_LAST_YEAR	2.490	1.498	1.662	0.100	-0.488	5.467
salary_millionsCent:w_pctCent	0.472	0.172	2.742	0.007	0.130	0.814

The equation of the final model without LeBron James is:  $\text{TWITTER\_FOLLOWER\_COUNT\_MILLIONS-hat} = -1.358 + 0.128 * \text{ageCent} + 16.105 * \text{PIECent} + 0.082 * \text{salary\_millionsCent} + 4.081 * \text{w\_pctCent} + 2.490 * \text{ACTIVE\_TWITTER\_LAST\_YEAR} + 0.472 * \text{salary\_millionsCent} * \text{w\_pct\_Cent}$ .

Comparing the two equations, we notice a huge discrepancy in the slope coefficient of `PIECent`. This makes sense since LeBron has the sixth-highest PIE in the league (0.183), so eliminating him from the dataset dramatically affects the mean PIE (as well as the spread, or standard deviation).

Since our objective is to accurately predict the Twitter follower counts of NBA players, it is probably best to exclude LeBron to avoid overestimating for less prominent athletes. Hence, we will continue with our analysis using the multiple linear regression model without LeBron James.

## Discussion of Assumptions

\*use model without LeBron (`m1_modV3`)