# Home Work Assignment - 01

*Critical Thinking Group 5*

## Contents

# Overview

The data set contains approximately 2200 records. Each record represents a professional baseball team from the years 1871 to 2006 inclusive. Each record has the performance of the team for the given year, with all of the statistics adjusted to match the performance of a 162 game season. We will be exploring, analyzing, and modeling the data set to predict a number of wins for a team using Ordinary Least Square (OLS).

To attain our objective, we will be following the below best practice steps and guidelines:

1 -Data Exploration
2 -Data Preparation
3 -Build Models
4 -Select Models

# 1 Data Exploration Analysis

In section we will explore and gain some insights into the dataset by pursuing the below high level steps and inquiries:

-Variable identification
-Variable Relationships
-Data summary analysis
-Outliers and Missing Values Identification

## 1.1 Variable identification

First let's display and examine the data dictionary or the data columns as shown in table 1.

Table 1: Variable Definition

| VARIABLE_NAME | DEFINITION | THEORETICAL_EFFECT |
|---|---|---|
| INDEX | Identification Variable (do not use) | None |
| TARGET_WINS | Number of wins | Target |
| TEAM_BATTING_H | Base Hits by batters (1B,2B,3B,HR) | Positive Impact on Wins |
| TEAM_BATTING_2B | Doubles by batters (2B) | Positive Impact on Wins |
| TEAM_BATTING_3B | Triples by batters (3B) | Positive Impact on Wins |
| TEAM_BATTING_HR | Homeruns by batters (4B) | Positive Impact on Wins |
| TEAM_BATTING_BB | Walks by batters | Positive Impact on Wins |
| TEAM_BATTING_HBP | Batters hit by pitch (get a free base) | Positive Impact on Wins |
| TEAM_BATTING_SO | Strikeouts by batters | Negative Impact on Wins |
| TEAM_BASERUN_SB | Stolen bases | Positive Impact on Wins |
| TEAM_BASERUN_CS | Caught stealing | Negative Impact on Wins |
| TEAM_FIELDING_E | Errors | Negative Impact on Wins |
| TEAM_FIELDING_DP | Double Plays | Positive Impact on Wins |
| TEAM_PITCHING_BB | Walks allowed | Negative Impact on Wins |
| TEAM_PITCHING_H | Hits allowed | Negative Impact on Wins |
| TEAM_PITCHING_HR | Homeruns allowed | Negative Impact on Wins |

| VARIABLE_NAME | DEFINITION | THEORETICAL_EFFECT |
|---|---|---|
| TEAM_PITCHING_SO | Strikeouts by pitchers | Positive Impact on Wins |

We notice that all variables are numeric. The variable names seem to follow certain naming pattern to highlight certain arithmetic relationships. In other words, we can compute the number of '1B' hits by taking the difference between overall hits and '2B', '3B', 'HR'. Although such naming and construct is not recommended in normalized database design ( as it violates third normal form), it is very frequent practice in the data analytics.

Our predictor input is made of 15 variables. And our dependent variable is one variable called TARGET_WINS.

Please note that we will not be using INDEX variable as it serves as just an identifier for each row. And has no relationships to other variables.

## 1.3 Data Summary Analysis

In this section, we will create summary data to better understand the initial relationship variables have with our dependent variable using correlation, central tendency, and dispersion As shown in table 2.

Table 2: Data Summary

|                    | mean       | sd         | median | trimmed    |
|--------------------|-----------:|-----------:|-------:|-----------:|
| TARGET_WINS        | 80.79086   | 15.75215   | 82.0   | 81.31229   |
| TEAM_BATTING_H     | 1469.26977 | 144.59120  | 1454.0 | 1459.04116 |
| TEAM_BATTING_2B    | 241.24692  | 46.80141   | 238.0  | 240.39627  |
| TEAM_BATTING_3B    | 55.25000   | 27.93856   | 47.0   | 52.17563   |
| TEAM_BATTING_HR    | 99.61204   | 60.54687   | 102.0  | 97.38529   |
| TEAM_BATTING_BB    | 501.55888  | 122.67086  | 512.0  | 512.18331  |
| TEAM_BATTING_SO    | 735.60534  | 248.52642  | 750.0  | 742.31322  |
| TEAM_BASERUN_SB    | 124.76177  | 87.79117   | 101.0  | 110.81188  |
| TEAM_BASERUN_CS    | 52.80386   | 22.95634   | 49.0   | 50.35963   |
| TEAM_BATTING_HBP   | 59.35602   | 12.96712   | 58.0   | 58.86275   |
| TEAM_PITCHING_H    | 1779.21046 | 1406.84293 | 1518.0 | 1555.89517 |
| TEAM_PITCHING_HR   | 105.69859  | 61.29875   | 107.0  | 103.15697  |
| TEAM_PITCHING_BB   | 553.00791  | 166.35736  | 536.5  | 542.62459  |
| TEAM_PITCHING_SO   | 817.73045  | 553.08503  | 813.5  | 796.93391  |
| TEAM_FIELDING_E    | 246.48067  | 227.77097  | 159.0  | 193.43798  |
| TEAM_FIELDING_DP   | 146.38794  | 26.22639   | 149.0  | 147.57789  |

Table 3: Missing Data and Data Correlation

|                    | Missing | Correlation |
|--------------------|--------:|------------:|
| TARGET_WINS        | 0       | 1.0000000   |
| TEAM_BATTING_H     | 0       | 0.3887675   |
| TEAM_BATTING_2B    | 0       | 0.2891036   |
| TEAM_BATTING_3B    | 0       | 0.1426084   |
| TEAM_BATTING_HR    | 0       | 0.1761532   |
| TEAM_BATTING_BB    | 0       | 0.2325599   |
| TEAM_BATTING_SO    | 102     | -0.0317507  |
| TEAM_BASERUN_SB    | 131     | 0.1351389   |
| TEAM_BASERUN_CS    | 772     | 0.0224041   |
| TEAM_BATTING_HBP   | 2085    | 0.0735042   |
| TEAM_PITCHING_H    | 0       | -0.1099371  |
| TEAM_PITCHING_HR   | 0       | 0.1890137   |
| TEAM_PITCHING_BB   | 0       | 0.1241745   |
| TEAM_PITCHING_SO   | 102     | -0.0784361  |
| TEAM_FIELDING_E    | 0       | -0.1764848  |
| TEAM_FIELDING_DP   | 286     | -0.0348506  |

Based on table 2 and Table 3, we can make the below observations:

1.Some of the variables like TEAM_PITCHING_H, TEAM_PITCHING_SO and TEAM_FIELDING_E seem to have outliers which is evident from the mean, median and trimmed mean values.

2.TEAM_BATTING_HBP and TEAM_BASERUN_CS seems to be missing a lot of values which casts doubt on its usefulness as a predictor. Maybe a flag for presense or absense of TEAM_BATTING_HBP and TEAM_BASERUN_CS might be a better predictor. Also given the fact that there is low correlation, we decided to exclude these 2 variables from any missing value or outlier treatment.
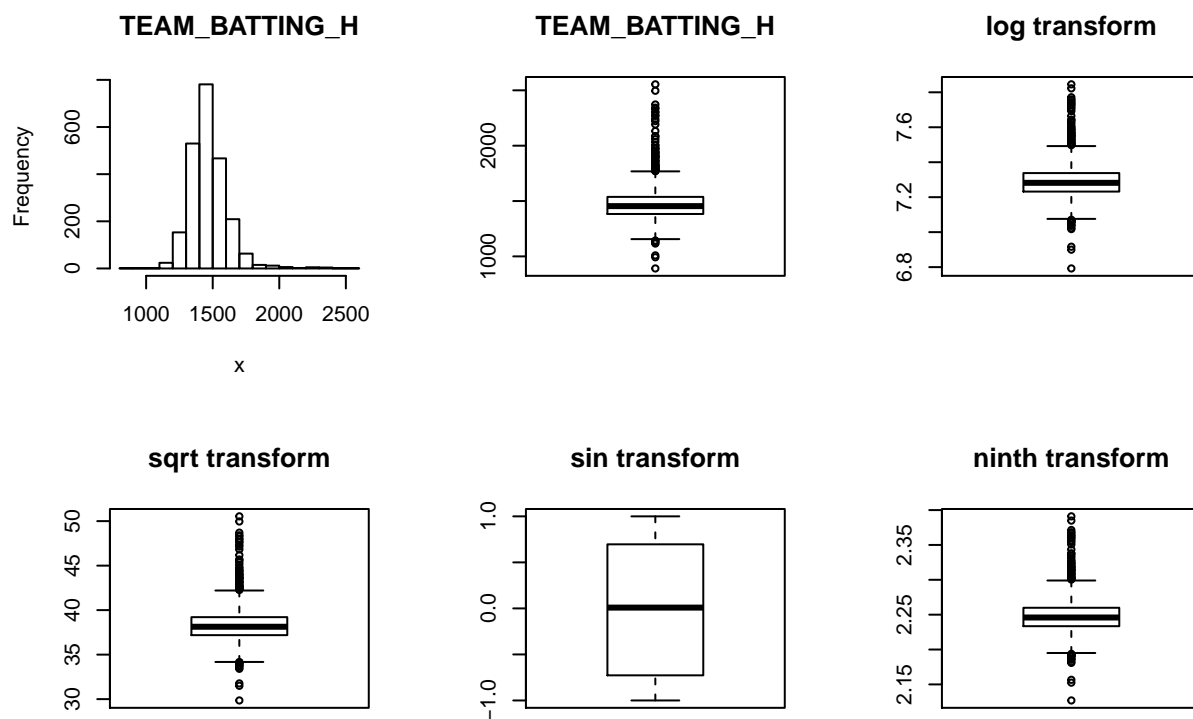
3.Most of the variables seem to indicate a positive / negative correlation in line with the theoretical effect. However, the following stand out as they show a correlation opposite to the theoretical impact: TEAM_BASERUN_CS, TEAM_PITCHING_HR, TEAM_PITCHING_BB, TEAM_PITCHING_SO and TEAM_FIELDING_DP. Lets evaluate these variables further once we fix any missing values or outliers.

4. We will impute the missing values in TEAM_BATTING_SO, FIELDING_DP, BASERUN_SB and TEAM_PITCHING_SO since it has lesser missing values even though there is low correlation. So we will create new variables that will have the respective missing values handled.
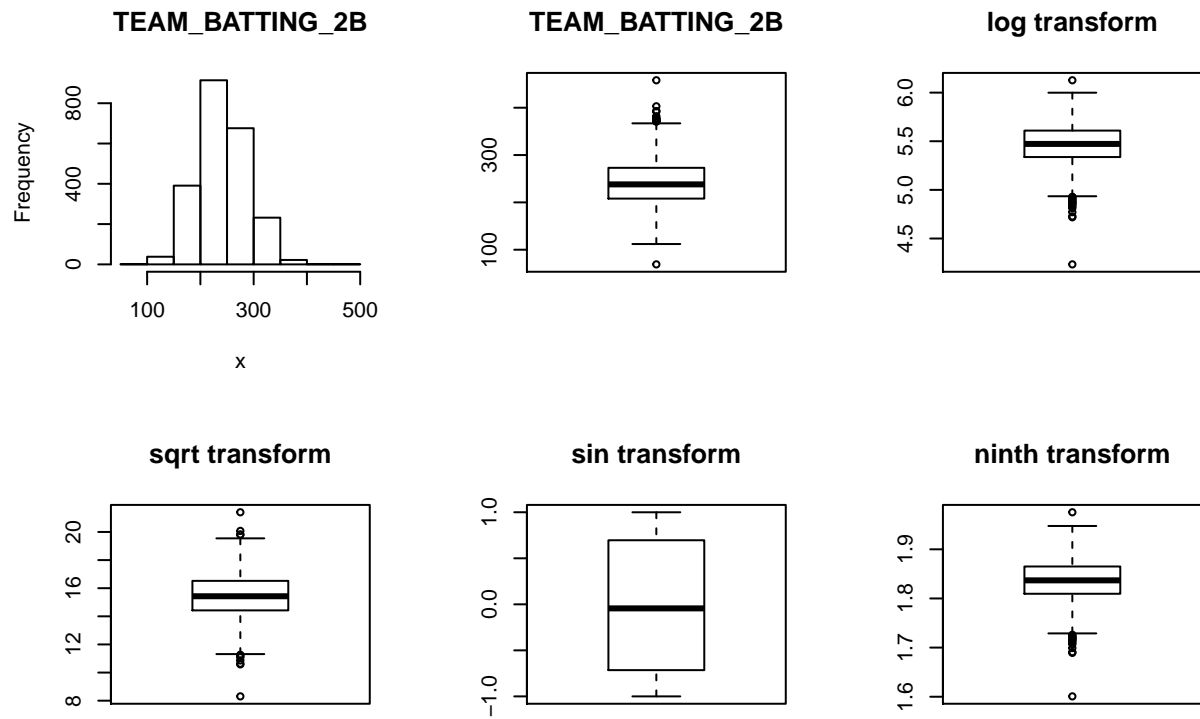
## 1.4 Outliers and Missing Values Identification

In this section we look at boxplots to determine the outliers in variables and decide on whether to act on the outliers.
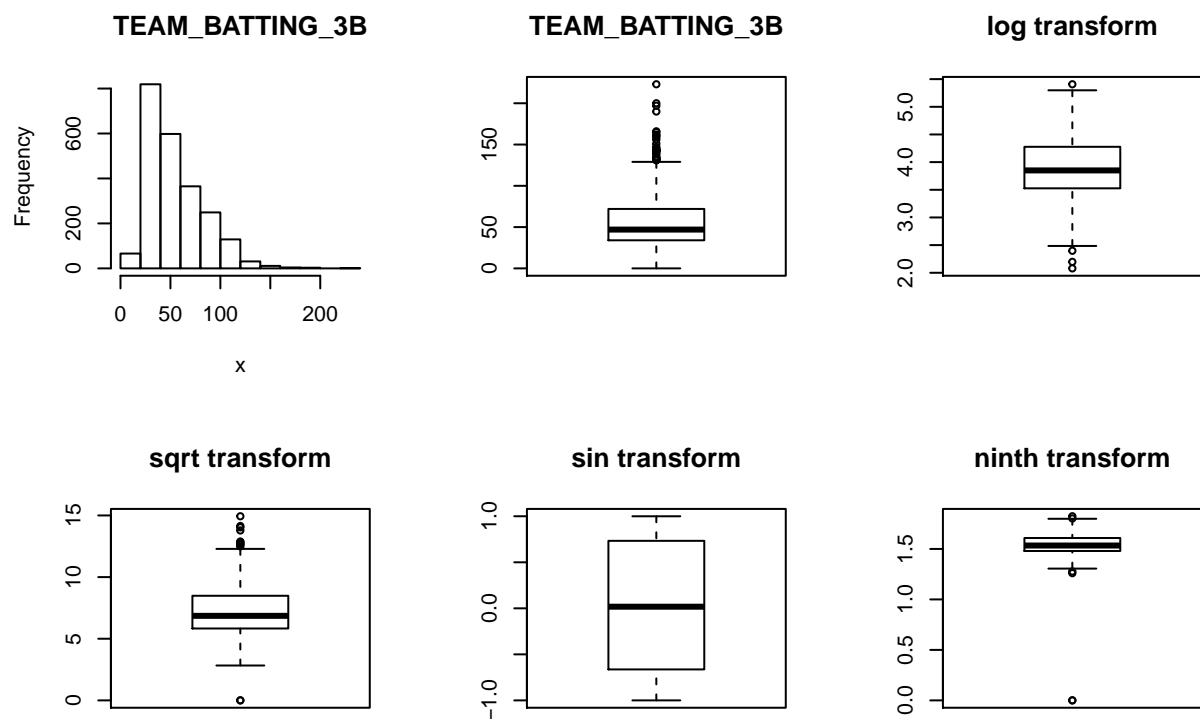
Lets do some univariate analysis. We will look at the Histogram and Boxplot for each variable to detect outliers if any and treat it accordingly.
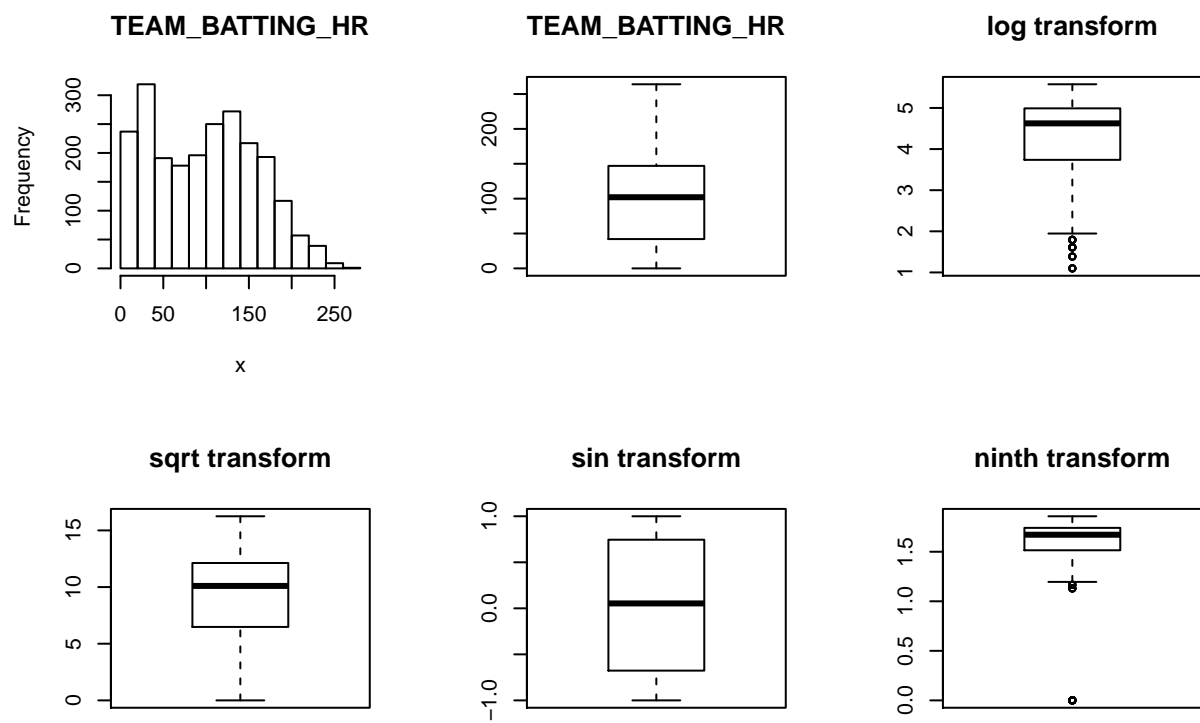


For TEAM_BATTING_H, we can see that there are quite a few outliers, both at the upper and lower end. Accordingly, we decide to create a new variable that will have the outlier fixed.
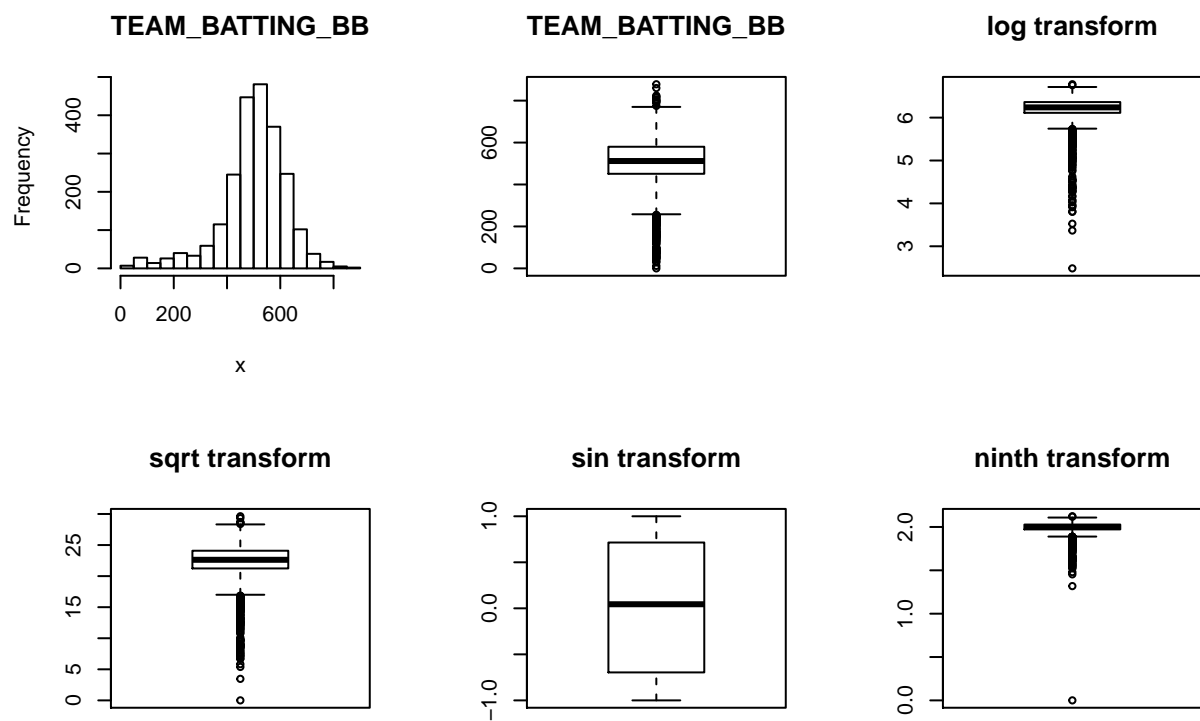
5

For TEAM_BATTING_2B, we can see that there are quite a few outliers, both at the upper and a single outlier at the lower end. For this variable we decide to create a new variable that will have the outliers fixed.
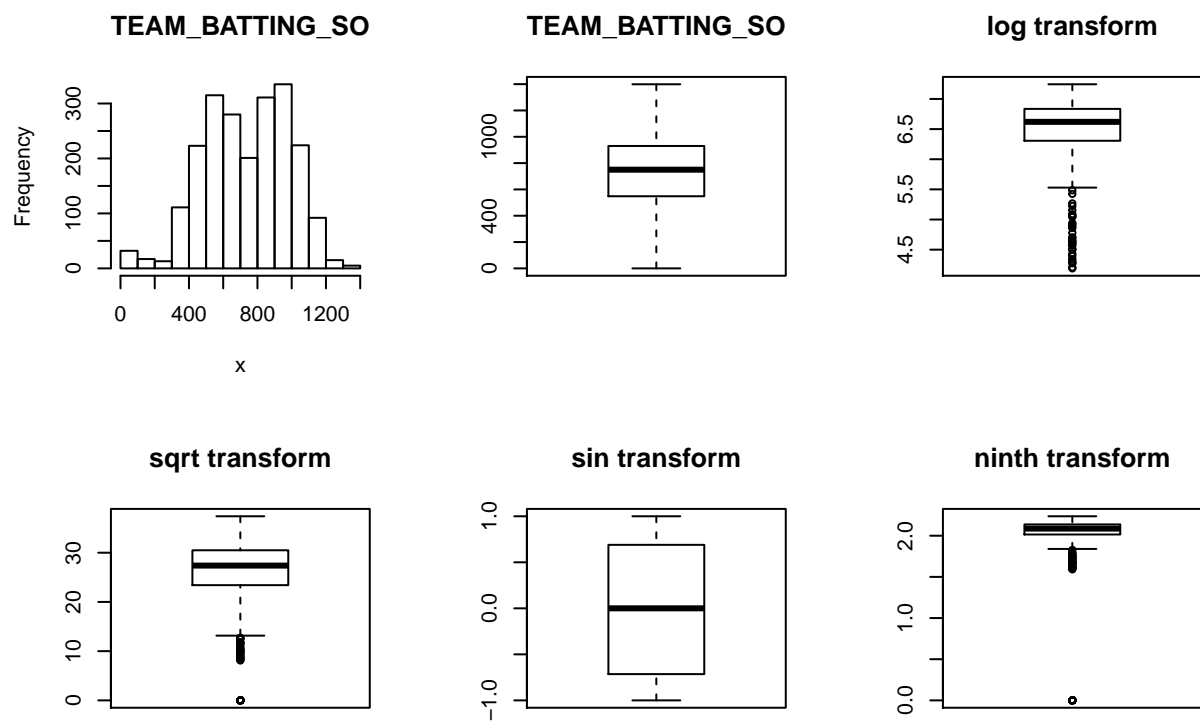
**TEAM_BATTING_3B**

**TEAM_BATTING_3B**

**log transform**

**sqrt transform**

**sin transform**

**ninth transform**

For TEAM_BATTING_3B, we can see that there are quite a few outliers at the upper end. For this variable we decide to create a new variable that will have the outliers fixed.
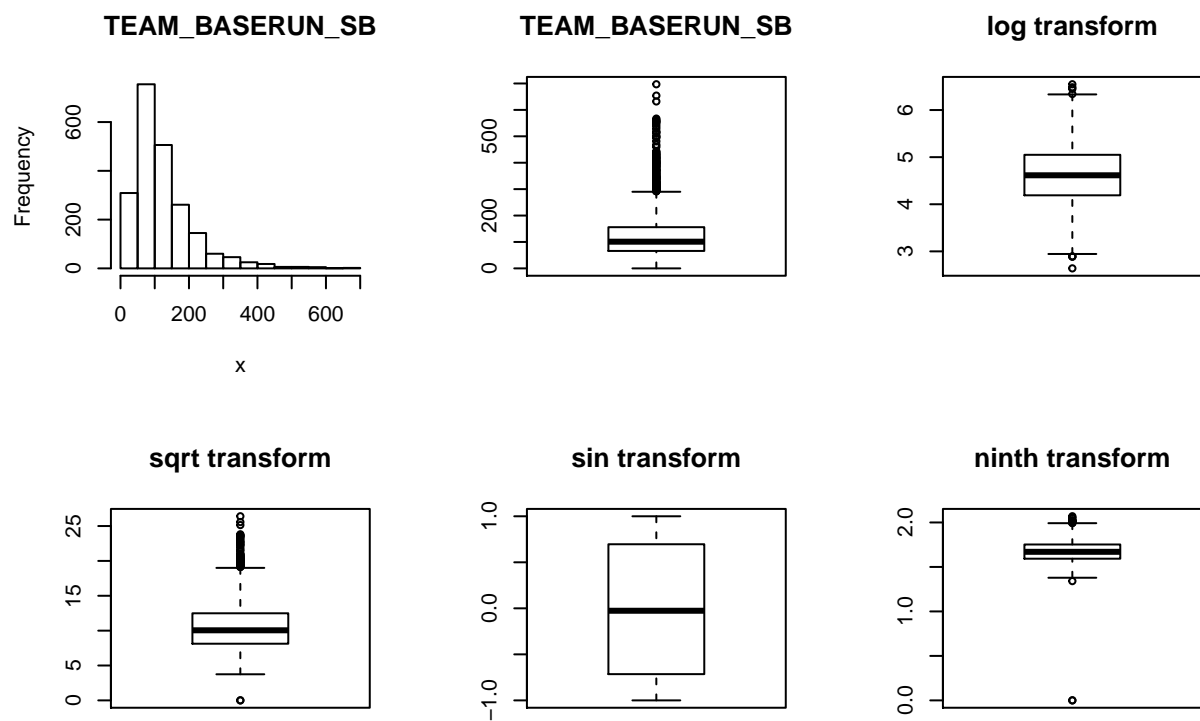
**TEAM_BATTING_HR**   **TEAM_BATTING_HR**   **log transform**



**sqrt transform**   **sin transform**   **ninth transform**



For TEAM_BATTING_HR, we can see that there are no outliers.

**TEAM_BATTING_BB**

**TEAM_BATTING_BB**

**log transform**
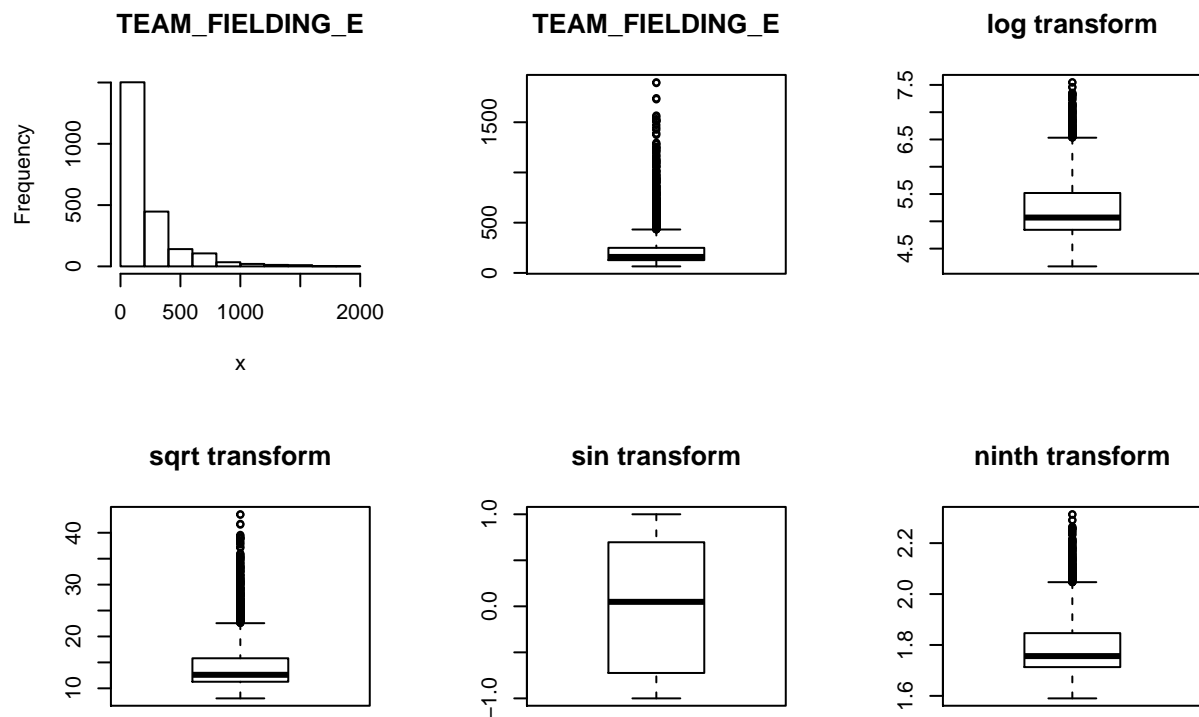
**sqrt transform**

**sin transform**

**ninth transform**

For TEAM_BATTING_BB, we can see that there are quite a few outliers, both at the upper and lower end. For this variable we decide to create a new variable that will have the outlier fixed.

**TEAM_BATTING_SO**  **TEAM_BATTING_SO**  **log transform**

**sqrt transform**  **sin transform**  **ninth transform**

For TEAM_BATTING_SO, we can see that there are no outliers. No further action needed for this variable.

**TEAM_BASERUN_SB**        **TEAM_BASERUN_SB**        **log transform**

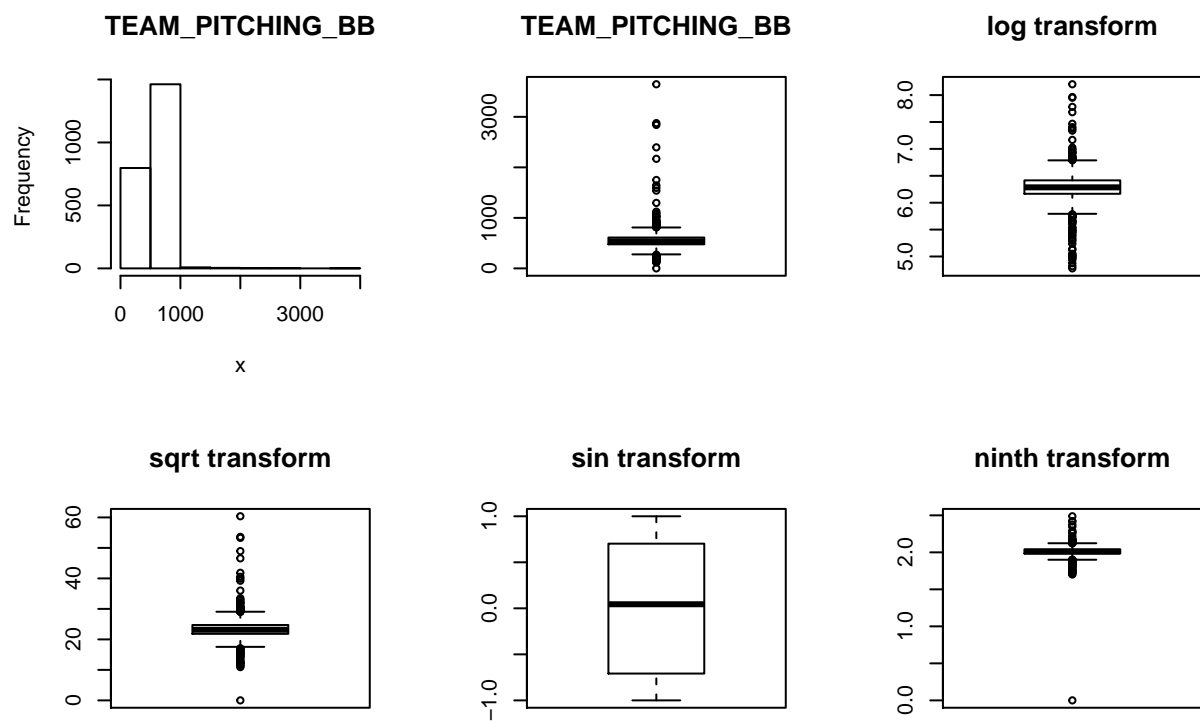**sqrt transform**        **sin transform**        **ninth transform**

For TEAM_BASERUN_SB, we can see that there are quite a few outliers at the upper end. For this variable we decide to create a new variable that will have the outlier fixed.
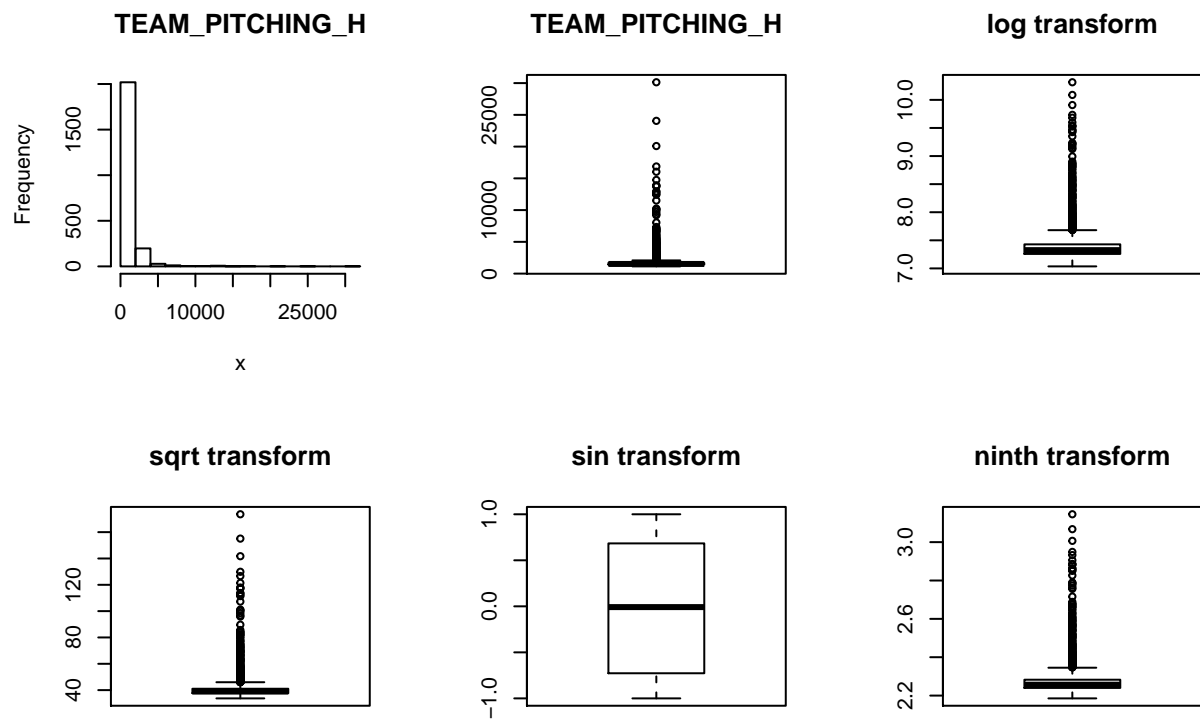
**TEAM_FIELDING_E**      **TEAM_FIELDING_E**      **log transform**

**sqrt transform**      **sin transform**      **ninth transform**

For TEAM_FIELDING_E, we can see that there are quite a few outliers at the upper end. For this variable we decide to create a new variable that will have the outlier fixed.
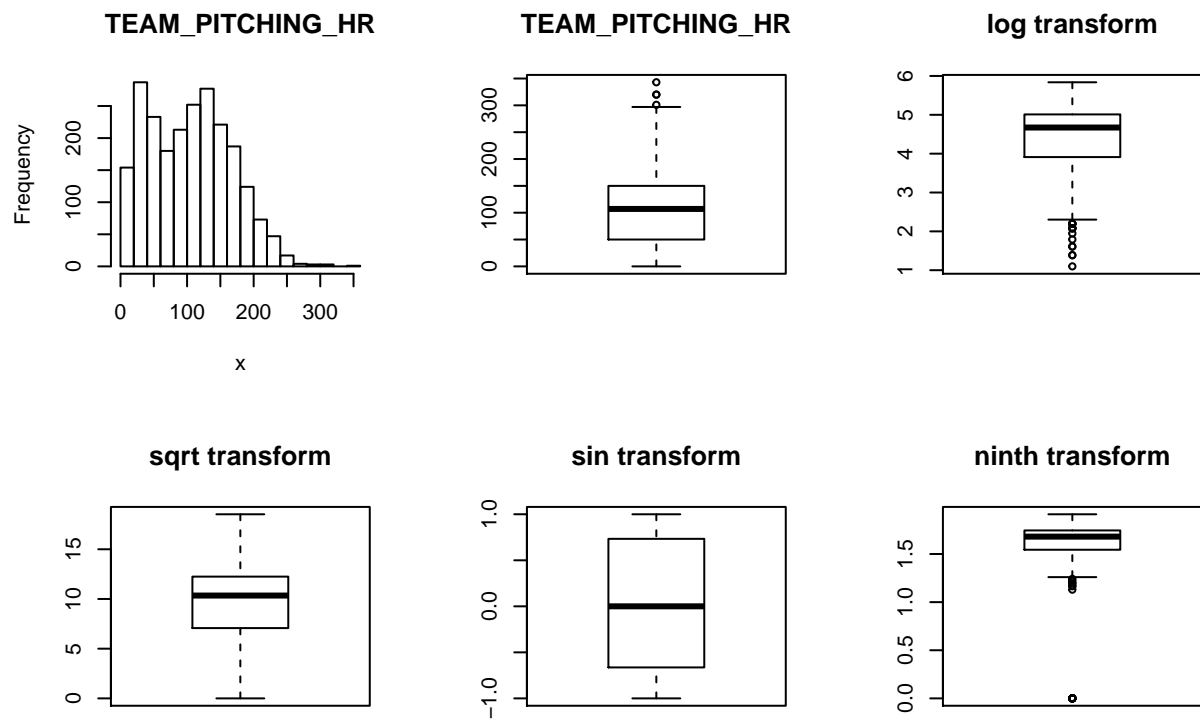
```
show_charts(moneyball2$TEAM_FIELDING_DP)
```

For TEAM_FIELDING_DP, we can see that there are quite a few outliers, both at the upper and lower end. For this variable we decide to create a new variable that will have the outlier fixed.

**TEAM_PITCHING_BB**

**TEAM_PITCHING_BB**

**log transform**

**sqrt transform**

**sin transform**

**ninth transform**

For TEAM_PITCHING_BB, we can see that there are quite a few outliers, both at the upper and lower end. For this variable we decide to create a new variable that will have the outlier fixed.
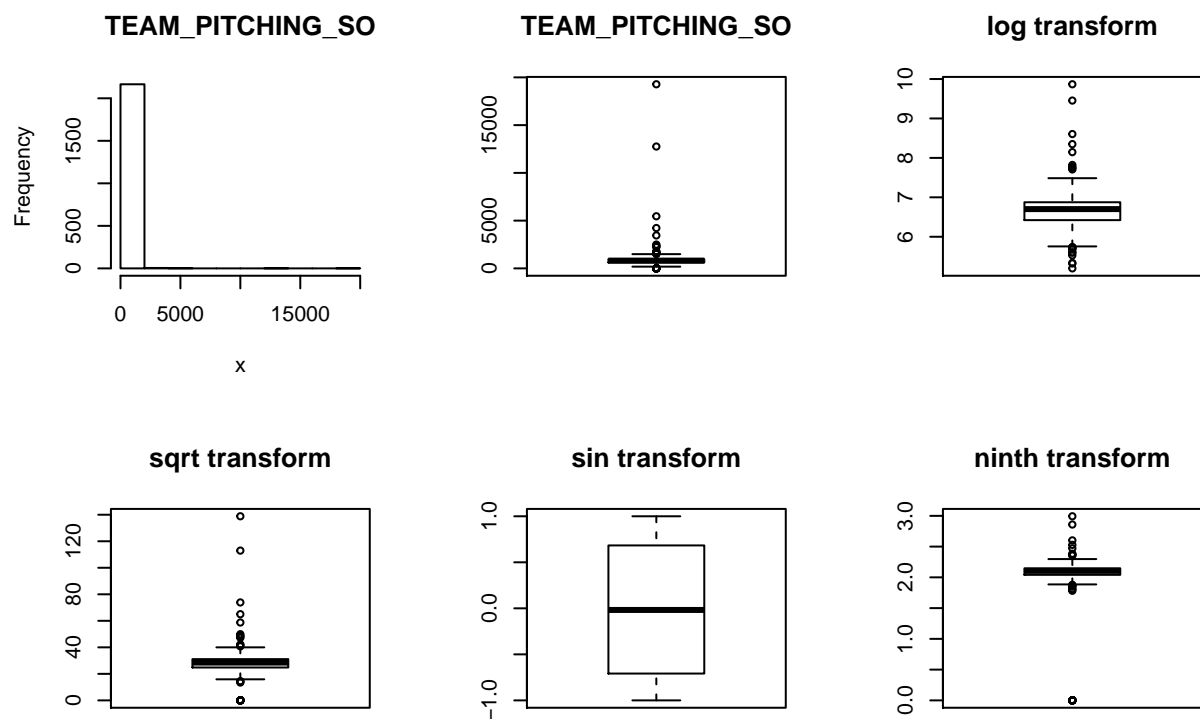
**TEAM_PITCHING_H**      **TEAM_PITCHING_H**      **log transform**

**sqrt transform**      **sin transform**      **ninth transform**

For TEAM_PITCHING_H, we can see that there are quite a few outliers at the upper end. For this variable we decide to create a new variable that will have the outlier fixed.

**TEAM_PITCHING_HR**   **TEAM_PITCHING_HR**   **log transform**

**sqrt transform**   **sin transform**   **ninth transform**

For TEAM_PITCHING_HR, we can see that there only 3 outliers at the upper end. For this variable we decide to create a new variable that will have the outlier fixed.

**TEAM_PITCHING_SO**     **TEAM_PITCHING_SO**     **log transform**

**sqrt transform**     **sin transform**     **ninth transform**

For TEAM_PITCHING_SO, we can see that there are quite a few outliers at the upper and a single outlier on the lower end. For this variable we decide to create a new variable that will have the outlier fixed.

** In most of the cases above, we see that a SIN transformation seems to work well to take care of the outliers. We will go ahead and create these new variables respectively.**

# 2. Data Preparation

Now that we have completed the preliminary analysis, we will be cleaning and consolidating data into one dataset for use in analysis and modeling. We will be puring the belwo steps as guidlines:
- Outliers treatment
- Missing values treatment
- Data transformation

## 2.1 Outliers treatment
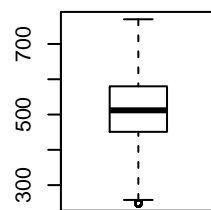
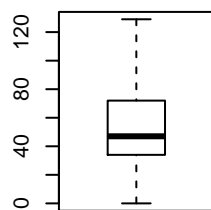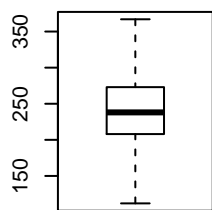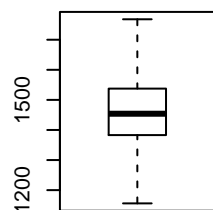For outliers, we will create 2 sets of variables.

The first set uses the capping method. In this method, we will replace all outliers that lie outside the 1.5 times of IQR limits. We will cap it by replacing those observations less than the lower limit with the value of 5th %ile and those that lie above the upper limit with the value of 95th %ile.

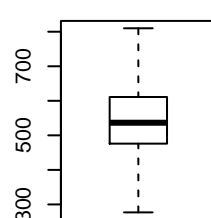Accordingly we create the following new variables while retaining the original variables.
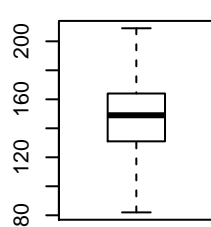
TEAM_BATTING_H_NEW
TEAM_BATTING_2B_NEW
TEAM_BATTING_3B_NEW
TEAM_BATTING_BB_NEW
TEAM_BASERUN_SB_NEW
TEAM_FIELDING_E_NEW
TEAM_FIELDING_DP_NEW
TEAM_PITCHING_BB_NEW
TEAM_PITCHING_H_NEW
TEAM_PITCHING_HR_NEW
TEAM_PITCHING_SO_NEW

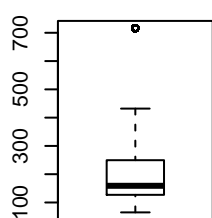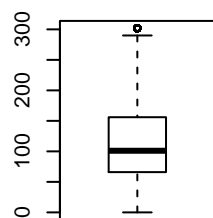Lets see how the new variables look in boxplots.

TEAM_BATTING_H_NE    TEAM_BATTING_2B_NE    TEAM_BATTING_3B_NE    TEAM_BATTING_BB_NI

TEAM_BASERUN_SB_N    TEAM_FIELDING_E_NE    TEAM_FIELDING_DP_N    TEAM_PITCHING_BB_N

**TEAM_PITCHING_H_NE TEAM_PITCHING_HR_N TEAM_PITCHING_SO_N**

In the second set, we will use the sin transformation and create the following variables:

TEAM_BATTING_H_SIN
TEAM_BATTING_2B_SIN
TEAM_BATTING_3B_SIN
TEAM_BATTING_BB_SIN
TEAM_BASERUN_SB_SIN
TEAM_FIELDING_E_SIN
TEAM_FIELDING_DP_SIN
TEAM_PITCHING_BB_SIN
TEAM_PITCHING_H_SIN
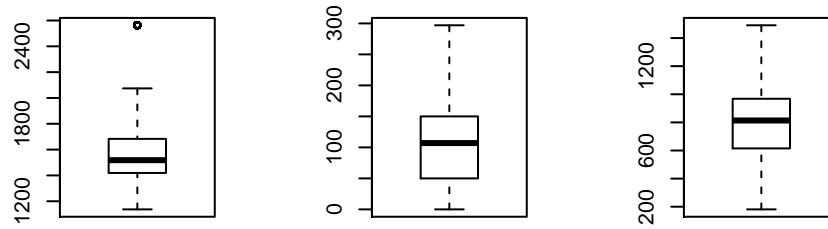TEAM_PITCHING_HR_SIN
TEAM_PITCHING_SO_SIN

TEAM_BATTING_H_S|    TEAM_BATTING_2B_S    TEAM_BATTING_3B_S    TEAM_BATTING_BB_S

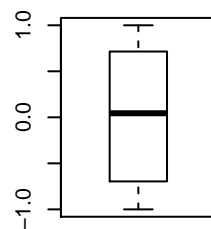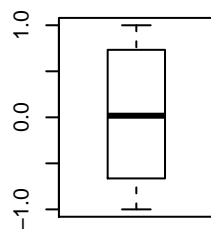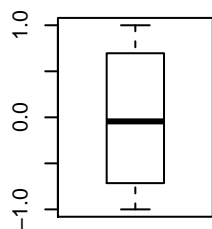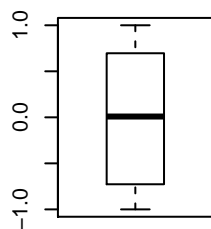TEAM_BASERUN_SB_S    TEAM_FIELDING_E_S|    TEAM_FIELDING_DP_S    TEAM_PITCHING_BB_S

**TEAM_PITCHING_H_S  TEAM_PITCHING_HR_S  TEAM_PITCHING_SO_S**



## 2.2 Missing values treatment

Next we impute missing values. Since we have handled outliers, we can go ahead and use the mean as impute values. As with outliers, we will go ahead and create new variables for the following:

TEAM_BATTING_SO_NEW

We will re-use the already created new variables for fixing the missing values for the below:

TEAM_PITCHING_SO_NEW TEAM_BASERUN_SB_NEW TEAM_FIELDING_DP_NEW

## 2.3 Additional Variables ???

Lets now create some additional variables that might help us in out analysis.

## 2.4 Missing Flags ??? 0 and 1 flags interchanged?

First we create flag variables to indicate whether TEAM_BATTING_HBP and TEAM_BASERUN_CS and missing. If the value is missing, we code it with 1 and if the value is present we code it with 0.
We will name our missing flag variables as follow:
TEAM_BATTING_HBP_Missing
TEAM_BASERUN_CS_Missing

## 2.5 Ratios

Next we create some additional variables, that we think may be useful with the prediction. Here we create the following ratios:

Hits_R = TEAM_BATTING_H/TEAM_PITCHING_H
Walks_R = TEAM_BATTING_BB/TEAM_PITCHING_BB
HomeRuns_R = TEAM_BATTING_HR/TEAM_PITCHING_HR
Strikeout_R = TEAM_BATTING_SO/TEAM_PITCHING_SO

## 2.6 Calculated Variables

Finally, we will also create calculated variables as below:

1. TEAM_BATTING_EB (Extra Base Hits) = 2B + 3B + HR
2. TEAM_BATTING_1B (Singles by batters) = TEAM_BATTING_H - TEAM_BATTING_EB

## 2.7 Correlation for new variables

Lets see how the new variables stack up against wins.

| | |
|---|---|
| TEAM_BATTING_HBP_Missing | 0.0026106 |
| TEAM_BASERUN_CS_Missing | 0.0048642 |
| Hits_R | 0.0958000 |
| Walks_R | 0.0836602 |
| HomeRuns_R | 0.0134410 |
| Strikeout_R | 0.0631939 |
| TEAM_BATTING_EB | 0.3449581 |
| TEAM_BATTING_1B | 0.2174301 |

All new variables seem to have a positive correlation with wins. However, some of them do not seem to have a strong correlation. Lets see how they perform while modeling.

# 3 Build Models

In this phase, we will build four models. The models independent variables will be based initially on the original data set variables, derived dataset variables, transformed dataset variables, and all variables in the dataset. In addition, for each model, we will perform a stepwise selection and stop at a point where we retain only those variables that have lower AIC (Akaike An Information Criterion). Recall (AIC) is a measure of the relative quality of statistical models for a given set of data. Given a collection of models for the data, AIC estimates the quality of each model, relative to each of the other models. Lower AIC leads to better quality model.

Below is a summary table showing models and their respective variables.

| VARIABLE_NAME | Comments | Theoretical.Effect | Model1 | Model2 | Model3 | Model4 |
|---|---|---|---|---|---|---|
| TEAM_BATTING_H | Given | Positive | Y | | | Y |
| TEAM_BATTING_2B | Given | Positive | Y | | | Y |
| TEAM_BATTING_3B | Given | Positive | Y | | | Y |
| TEAM_BATTING_HR | Given | Positive | Y | | | Y |
| TEAM_BATTING_BB | Given | Positive | Y | | | Y |
| TEAM_BATTING_HBP | Given | Positive | Y | | | |
| TEAM_BATTING_SO | Given | Negative | Y | | | Y |
| TEAM_BASERUN_SB | Given | Positive | Y | | | Y |
| TEAM_BASERUN_CS | Given | Negative | Y | | | |
| TEAM_FIELDING_E | Given | Negative | Y | | | Y |
| TEAM_FIELDING_DP | Given | Positive | Y | | | Y |
| TEAM_PITCHING_BB | Given | Negative | Y | | | Y |
| TEAM_PITCHING_H | Given | Negative | Y | | | Y |
| TEAM_PITCHING_HR | Given | Negative | Y | | | Y |
| TEAM_PITCHING_SO | Given | Positive | Y | | | Y |
| TEAM_BATTING_H_NEW | Derived | Positive | | Y | | Y |
| TEAM_BATTING_2B_NEW | Derived | Positive | | Y | | Y |
| TEAM_BATTING_3B_NEW | Derived | Positive | | Y | | Y |
| TEAM_BATTING_BB_NEW | Derived | Positive | | Y | | Y |
| TEAM_BASERUN_SB_NEW | Derived | Positive | | Y | | Y |
| TEAM_FIELDING_E_NEW | Derived | Negative | | Y | | Y |
| TEAM_FIELDING_DP_NEW | Derived | Positive | | Y | | Y |
| TEAM_PITCHING_BB_NEW | Derived | Negative | | Y | | Y |
| TEAM_PITCHING_H_NEW | Derived | Negative | | Y | | Y |
| TEAM_PITCHING_HR_NEW | Derived | Negative | | Y | | Y |
| TEAM_PITCHING_SO_NEW | Derived | Positive | | Y | | Y |
| TEAM_BATTING_H_SIN | Derived | Positive | | | Y | Y |
| TEAM_BATTING_2B_SIN | Derived | Positive | | | Y | Y |
| TEAM_BATTING_3B_SIN | Derived | Positive | | | Y | Y |
| TEAM_BATTING_BB_SIN | Derived | Positive | | | Y | Y |
| TEAM_BASERUN_SB_SIN | Derived | Positive | | | Y | Y |
| TEAM_FIELDING_E_SIN | Derived | Negative | | | Y | Y |
| TEAM_FIELDING_DP_SIN | Derived | Positive | | | Y | Y |
| TEAM_PITCHING_BB_SIN | Derived | Negative | | | Y | Y |
| TEAM_PITCHING_H_SIN | Derived | Negative | | | Y | Y |
| TEAM_PITCHING_HR_SIN | Derived | Negative | | | Y | Y |
| TEAM_PITCHING_SO_SIN | Derived | Positive | | | Y | Y |
| TEAM_BATTING_HBP_Missing | Derived | | | | Y | Y |
| TEAM_BASERUN_CS_Missing | Derived | | | | Y | Y |
| Hits_R | Derived | | | | Y | Y |
| Walks_R | Derived | | | | Y | Y |
| HomeRuns_R | Derived | | | | Y | Y |
| Strikeout_R | Derived | | | | Y | Y |
| TEAM_BATTING_EB | Derived | | | | Y | Y |
| TEAM_BATTING_1B | Derived | | | | Y | Y |

## 3.1 Model One

In this model, we will be using the original variables. We will create model and we will highlight the variables that being recommended using the AIC value.
First we will produce the summary model as per below:

```
##                  Estimate Std. Error t value  Pr(>|t|)
## (Intercept)      60.288263  19.678417  3.0637  0.002532
## TEAM_BATTING_H    1.913476   2.761394  0.6929  0.489267
## TEAM_BATTING_2B   0.026388   0.030290  0.8712  0.384844
## TEAM_BATTING_3B  -0.101176   0.077507 -1.3054  0.193477
## TEAM_BATTING_HR  -4.843707  10.508511 -0.4609  0.645420
## TEAM_BATTING_BB  -4.459691   3.636241 -1.2265  0.221675
## TEAM_BATTING_HBP  0.082473   0.049600  1.6628  0.098152
## TEAM_BATTING_SO   0.341963   2.598759  0.1316  0.895462
## TEAM_BASERUN_SB   0.033044   0.028673  1.1524  0.250708
## TEAM_BASERUN_CS  -0.011044   0.071431 -0.1546  0.877303
## TEAM_FIELDING_E  -0.172042   0.041404 -4.1552 5.076e-05
## TEAM_FIELDING_DP -0.108192   0.036541 -2.9609  0.003494
## TEAM_PITCHING_BB  4.510891   3.633720  1.2414  0.216120
## TEAM_PITCHING_H  -1.890957   2.760946 -0.6849  0.494317
## TEAM_PITCHING_HR  4.930432  10.506645  0.4693  0.639462
## TEAM_PITCHING_SO -0.373645   2.597052 -0.1439  0.885767
##
## n = 191, p = 16, Residual SE = 8.46704, R-Squared = 0.55
```

We notice that model 1 has the following summary characteristics:
-The Residual standard error is 8.467
-Degrees of freedom: 175
-Deleted observations due missing data: 2085.
-Multiple R-squared: 0.5501
-Adjusted R-squared: 0.5116
-F-statistic: 14.27 on 15 and 175 DF
-p-value: $< 2.2e-16$
Next. let's step thru this model (model 1) and retain only those variables that have the most impact. below the relevant varuibale for model 1:

|                  | Coefficients |
|------------------|--------------|
| (Intercept)      | 60.4020741   |
| TEAM_PITCHING_H  | 0.0257684    |
| TEAM_PITCHING_BB | 0.0565953    |
| TEAM_FIELDING_E  | -0.1728308   |
| TEAM_FIELDING_DP | -0.1183189   |
| TEAM_BATTING_SO  | -0.0313641   |
| TEAM_PITCHING_HR | 0.0895863    |
| TEAM_BATTING_HBP | 0.0867904    |

## 3.2 Model Two

In this model (model2), we will be using the adjusted values based on our outlier treatment process. We will create model and we will highlight the variables that being recommended using the AIC value. First we will

produce the summary model as per below:

```
##                        Estimate Std. Error t value  Pr(>|t|)
## (Intercept)          58.8339822 19.3451178  3.0413 0.0027097
## TEAM_BATTING_H_NEW   -0.1019444  0.2050405 -0.4972 0.6196643
## TEAM_BATTING_2B_NEW   0.0256621  0.0307208  0.8353 0.4046437
## TEAM_BATTING_3B_NEW  -0.1255273  0.0756929 -1.6584 0.0989930
## TEAM_BATTING_BB_NEW   0.0367410  0.0849861  0.4323 0.6660312
## TEAM_BASERUN_SB_NEW   0.0313690  0.0227078  1.3814 0.1688727
## TEAM_FIELDING_E_NEW  -0.1771414  0.0404787 -4.3762 2.049e-05
## TEAM_FIELDING_DP_NEW -0.1037738  0.0365712 -2.8376 0.0050705
## TEAM_PITCHING_BB_NEW  0.0176296  0.0831672  0.2120 0.8323653
## TEAM_PITCHING_H_NEW   0.1260298  0.2053907  0.6136 0.5402520
## TEAM_PITCHING_HR_NEW  0.0905404  0.0256367  3.5317 0.0005252
## TEAM_PITCHING_SO_NEW -0.0296146  0.0073102 -4.0511 7.586e-05
##
## n = 191, p = 12, Residual SE = 8.46938, R-Squared = 0.54
```

Lets now step thru this model and retain only those variables that have the most impact.

## 3.3 Model Three

In this model (model3), we will be using the derived values based on our variable transformation process. We will create model and we will highlight the variables that being recommended using the AIC value. First we will produce the summary model as per below:

```
##
## Coefficients: (2 not defined because of singularities)
##                        Estimate  Std. Error t value  Pr(>|t|)
## (Intercept)          4.5570e+02  4.0350e+02  1.1294 0.260305
## TEAM_BATTING_H_SIN  -6.1473e-01  2.0827e+00 -0.2952 0.768229
## TEAM_BATTING_2B_SIN  8.2353e-02  1.0852e+00  0.0759 0.939598
## TEAM_BATTING_3B_SIN  5.8838e-01  1.1441e+00  0.5143 0.607716
## TEAM_BATTING_BB_SIN -2.3186e+00  2.2911e+00 -1.0120 0.312953
## TEAM_BASERUN_SB_SIN -2.1823e+00  1.1036e+00 -1.9775 0.049570
## TEAM_FIELDING_E_SIN  5.0539e-01  1.0992e+00  0.4598 0.646255
## TEAM_FIELDING_DP_SIN 2.3551e+00  1.1145e+00  2.1131 0.036023
## TEAM_PITCHING_BB_SIN 4.7164e-01  2.2460e+00  0.2100 0.833924
## TEAM_PITCHING_H_SIN  7.7264e-01  2.0684e+00  0.3735 0.709202
## TEAM_PITCHING_HR_SIN -1.6959e+00 1.1009e+00 -1.5405 0.125257
## TEAM_PITCHING_SO_SIN 7.7773e-01  1.1130e+00  0.6988 0.485643
## Hits_R               4.7994e+02  9.6173e+03  0.0499 0.960256
## Walks_R             -1.0068e+04  5.2166e+03 -1.9300 0.055242
## HomeRuns_R           3.9477e+03  2.0061e+03  1.9679 0.050679
## Strikeout_R          5.1720e+03  8.5650e+03  0.6039 0.546730
## TEAM_BATTING_EB      1.0204e-01  1.7555e-02  5.8124 2.902e-08
## TEAM_BATTING_1B      4.2866e-02  1.2977e-02  3.3033 0.001161
```

```
##
## n = 191, p = 18, Residual SE = 10.13987, R-Squared = 0.36
```

Lets now step thru this model and retain only those variables that have the most impact.

## 3.4 Model Four

In this model (model4), we will be using all variables original, adjusted, and derived values. We will create model and we will highlight the variables that being recommended using the AIC value. First we will produce the summary model as per below:

```
##
## Coefficients: (14 not defined because of singularities)
##                        Estimate  Std. Error t value  Pr(>|t|)
## (Intercept)            1.5803e+04  1.6334e+04  0.9675   0.33478
## TEAM_BATTING_H         1.4193e+01  1.1598e+01  1.2238   0.22286
## TEAM_BATTING_2B        6.9092e-02  1.2259e-01  0.5636   0.57382
## TEAM_BATTING_3B       -6.7019e-02  8.0344e-02 -0.8341   0.40546
## TEAM_BATTING_HR       -3.1382e+01  2.2076e+01 -1.4216   0.15712
## TEAM_BATTING_BB        1.4865e+01  8.5333e+00  1.7420   0.08345
## TEAM_BATTING_SO       -7.5314e+00  3.9562e+00 -1.9037   0.05877
## TEAM_BASERUN_SB        2.8975e-02  3.0059e-02  0.9639   0.33655
## TEAM_BASERUN_CS       -2.5752e-02  7.2907e-02 -0.3532   0.72440
## TEAM_BATTING_HBP       8.8491e-02  5.0536e-02  1.7510   0.08188
## TEAM_PITCHING_H       -1.4180e+01  1.1597e+01 -1.2227   0.22327
## TEAM_PITCHING_HR       3.1466e+01  2.2073e+01  1.4255   0.15597
## TEAM_PITCHING_BB      -1.4843e+01  8.5379e+00 -1.7385   0.08408
## TEAM_PITCHING_SO       7.4970e+00  3.9540e+00  1.8961   0.05978
## TEAM_FIELDING_E       -1.8995e-01  4.3293e-02 -4.3875 2.087e-05
## TEAM_FIELDING_DP      -9.8295e-02  3.8185e-02 -2.5742   0.01097
## TEAM_BATTING_2B_NEW   -4.6960e-02  1.2464e-01 -0.3768   0.70686
## TEAM_BATTING_BB_NEW    3.1093e-02  8.6932e-02  0.3577   0.72107
## TEAM_BATTING_H_SIN    -8.1802e-01  1.8861e+00 -0.4337   0.66508
## TEAM_BATTING_2B_SIN   -6.8111e-01  9.2777e-01 -0.7341   0.46395
## TEAM_BATTING_3B_SIN   -4.1022e-01  9.8554e-01 -0.4162   0.67780
## TEAM_BATTING_BB_SIN   -1.0084e+00  1.9831e+00 -0.5085   0.61182
## TEAM_BASERUN_SB_SIN   -2.3013e+00  9.3403e-01 -2.4638   0.01482
## TEAM_FIELDING_E_SIN   -4.9238e-01  9.2782e-01 -0.5307   0.59639
## TEAM_FIELDING_DP_SIN   1.7662e+00  9.5433e-01  1.8507   0.06608
## TEAM_PITCHING_BB_SIN  -7.4780e-02  1.9432e+00 -0.0385   0.96935
## TEAM_PITCHING_H_SIN    1.0784e+00  1.8692e+00  0.5770   0.56479
## TEAM_PITCHING_HR_SIN  -9.5148e-01  9.3622e-01 -1.0163   0.31104
## TEAM_PITCHING_SO_SIN  -9.0822e-01  9.6051e-01 -0.9456   0.34581
## Hits_R                -2.3615e+04  2.0532e+04 -1.1502   0.25181
## Walks_R               -1.8042e+04  9.1272e+03 -1.9767   0.04981
## HomeRuns_R             1.1879e+04  6.1102e+03  1.9441   0.05367
## Strikeout_R            1.4052e+04  8.9098e+03  1.5772   0.11675
##
## n = 191, p = 33, Residual SE = 8.27202, R-Squared = 0.61
```

Lets now step thru this model and retain only those variables that have the most impact.

Discuss the coefficients in the models, do they make sense? For example, if a team hits a lot of Home Runs, it would be reasonably expected that such a team would win more games. However, if the coefficient is negative (suggesting that the team would lose more games), then that needs to be discussed. Are you keeping the model even though it is counter intuitive? Why? The boss needs to know.
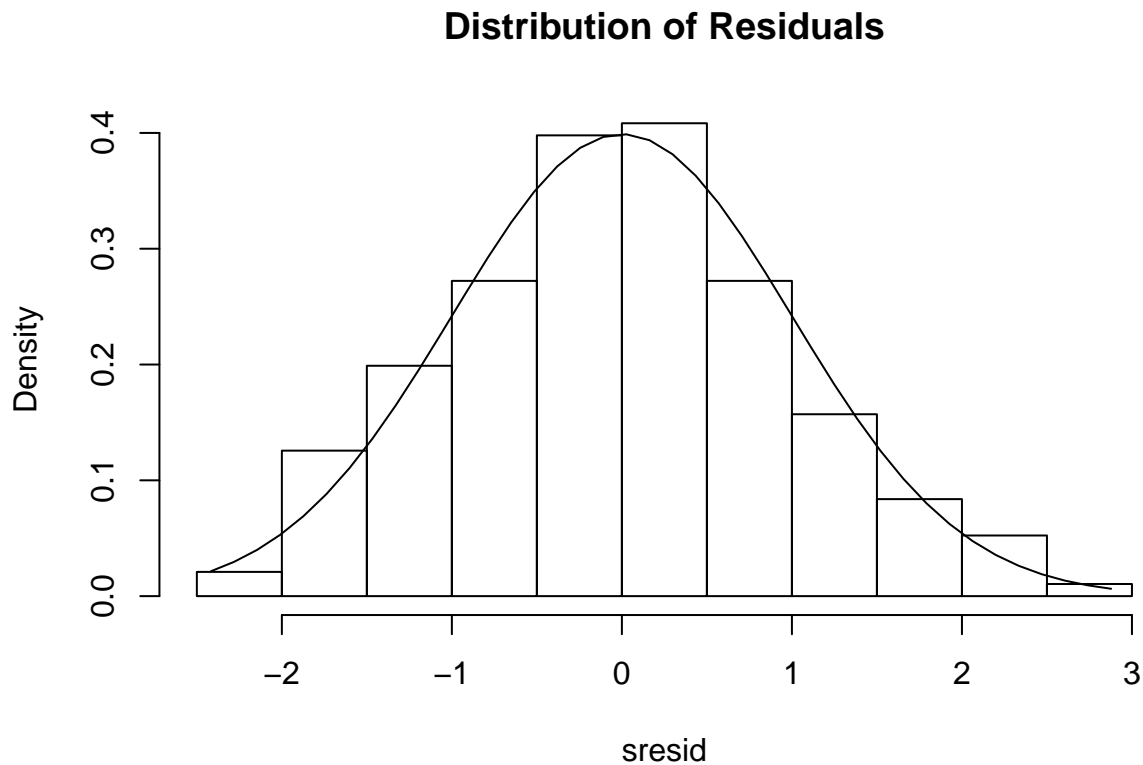
## Select Models

Decide on the criteria for selecting the best multiple linear regression model. Will you select a model with slightly worse performance if it makes more sense or is more parsimonious? Discuss why you selected your model. For the multiple linear regression model, will you use a metric such as Adjusted R2, RMSE, etc.? Be sure to explain how you can make inferences from the model, discuss multi-collinearity issues (if any), and discuss other relevant model output. Using the training data set, evaluate the multiple linear regression model based on (a) mean squared error, (b) R2, (c) F-statistic, and (d) residual plots. Make predictions using the evaluation data set.

**Model One with original data**

```
##
## Call:
## lm(formula = TARGET_WINS ~ TEAM_BATTING_H + TEAM_BATTING_2B +
##     TEAM_BATTING_3B + TEAM_BATTING_HR + TEAM_BATTING_BB + TEAM_BATTING_HBP +
##     TEAM_BATTING_SO + TEAM_BASERUN_SB + TEAM_BASERUN_CS + TEAM_FIELDING_E +
##     TEAM_FIELDING_DP + TEAM_PITCHING_BB + TEAM_PITCHING_H + TEAM_PITCHING_HR +
##     TEAM_PITCHING_SO, data = moneyball2)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -19.8708  -5.6564  -0.0599   5.2545  22.9274
##
## Coefficients:
##                   Estimate Std. Error t value Pr(>|t|)
## (Intercept)       60.28826   19.67842   3.064  0.00253 **
## TEAM_BATTING_H     1.91348    2.76139   0.693  0.48927
## TEAM_BATTING_2B    0.02639    0.03029   0.871  0.38484
## TEAM_BATTING_3B   -0.10118    0.07751  -1.305  0.19348
## TEAM_BATTING_HR   -4.84371   10.50851  -0.461  0.64542
## TEAM_BATTING_BB   -4.45969    3.63624  -1.226  0.22167
## TEAM_BATTING_HBP   0.08247    0.04960   1.663  0.09815 .
## TEAM_BATTING_SO    0.34196    2.59876   0.132  0.89546
## TEAM_BASERUN_SB    0.03304    0.02867   1.152  0.25071
## TEAM_BASERUN_CS   -0.01104    0.07143  -0.155  0.87730
## TEAM_FIELDING_E   -0.17204    0.04140  -4.155 5.08e-05 ***
## TEAM_FIELDING_DP  -0.10819    0.03654  -2.961  0.00349 **
## TEAM_PITCHING_BB   4.51089    3.63372   1.241  0.21612
## TEAM_PITCHING_H   -1.89096    2.76095  -0.685  0.49432
## TEAM_PITCHING_HR   4.93043   10.50664   0.469  0.63946
## TEAM_PITCHING_SO  -0.37364    2.59705  -0.144  0.88577
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.467 on 175 degrees of freedom
##   (2085 observations deleted due to missingness)
## Multiple R-squared:  0.5501, Adjusted R-squared:  0.5116
## F-statistic: 14.27 on 15 and 175 DF,  p-value: < 2.2e-16
```

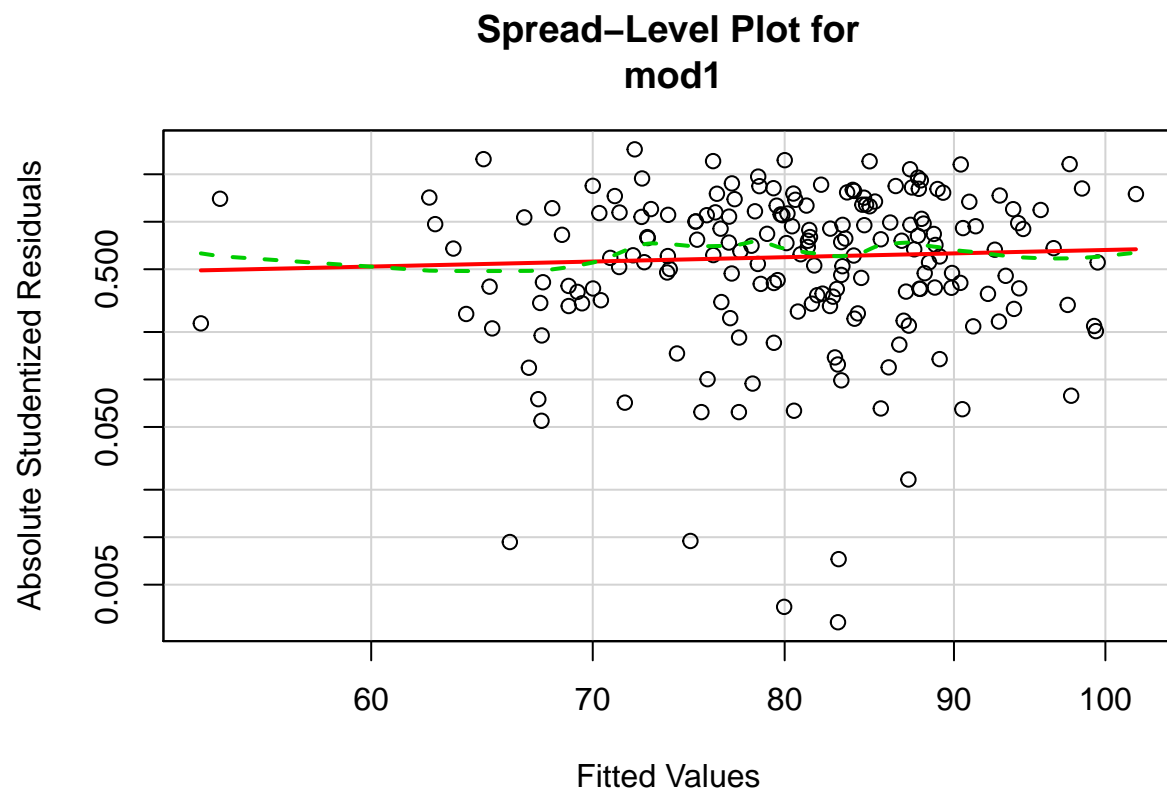**Normality check of Residuals**

## Distribution of Residuals



The residuals are normally distributed, this indicates That the mean of the difference between our predictions

and the actual values is close to 0 which is good for our analysis.

Also, it's unlikely that no relationship exists between **TEAM_FIELDING_E** and **TAR-GET_WINS.**

homoscedasticity check or non-constant error variance test

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.03994848    Df = 1      p = 0.8415813
```

## Spread–Level Plot for
## mod1



```
##
## Suggested power transformation:  0.5267026
```

The test confirms the non-constant error variance test. It also has a p-value higher than a significance level of **0.05**.

Therefore we can accept the null hypothesis that the variance of the residuals is constant and infer that heteroscedasticity is not present.

**Collinearity Check**

```
##    TEAM_BATTING_H  TEAM_BATTING_2B  TEAM_BATTING_3B  TEAM_BATTING_HR
##      1.171824e+05     1.685623e+00     1.302198e+00     3.074804e+05
##   TEAM_BATTING_BB TEAM_BATTING_HBP  TEAM_BATTING_SO  TEAM_BASERUN_SB
##      1.962853e+05     1.096334e+00     1.941752e+05     1.950069e+00
##   TEAM_BASERUN_CS  TEAM_FIELDING_E TEAM_FIELDING_DP TEAM_PITCHING_BB
##      1.914415e+00     1.256819e+00     1.097611e+00     1.964039e+05
##   TEAM_PITCHING_H TEAM_PITCHING_HR TEAM_PITCHING_SO
##      1.160417e+05     3.069624e+05     1.946316e+05
```

```
##                 Collinearity
## TEAM_BATTING_H          TRUE
## TEAM_BATTING_2B        FALSE
```

```
## TEAM_BATTING_3B        FALSE
## TEAM_BATTING_HR          TRUE
## TEAM_BATTING_BB          TRUE
## TEAM_BATTING_HBP       FALSE
## TEAM_BATTING_SO          TRUE
## TEAM_BASERUN_SB        FALSE
## TEAM_BASERUN_CS        FALSE
## TEAM_FIELDING_E        FALSE
## TEAM_FIELDING_DP       FALSE
## TEAM_PITCHING_BB         TRUE
## TEAM_PITCHING_H          TRUE
## TEAM_PITCHING_HR         TRUE
## TEAM_PITCHING_SO         TRUE
```

# Test for Autocorrelated Errors

durbinWatsonTest(mod1)

```
##  lag Autocorrelation D-W Statistic p-value
##   1        0.2128921      1.567453        0
##  Alternative hypothesis: rho != 0
```

**goodness of fit of your model**

**using R-squared and adjusted R-squared, our model is about 55% predicts the TARGET_WINS**