

Moneyball - CUNY Data Science 621

Eric Hirsch

2/20/2021

Description of the Dataset

XXXXXX

An issue with the data is hidden groupings. Records may not be independent of each other, as team data in one year will be related to team data in the next year. We know that if some records were adjusted to match a longer season, there may be an “eras of baseball” effect as teams from earlier years behave differently from later ones. Finally, within the record, columns may not be independent. In particular, teams with high offensive stats (like hitting) may have lower defensive stats (like pitching), as the teams on limited budgets make strategic choices between the two. We will attempt to address some of these issues in this analysis.

1. Data Exploration

All of the columns in the dataset are numeric. We begin by examining their means, medians and distributions.

##	INDEX	TARGET_WINS	BATTING_H	BATTING_2B
##	Min. : 1.0	Min. : 0.00	Min. : 891	Min. : 69.0
##	1st Qu.: 630.8	1st Qu.: 71.00	1st Qu.:1383	1st Qu.:208.0
##	Median :1270.5	Median : 82.00	Median :1454	Median :238.0
##	Mean :1268.5	Mean : 80.79	Mean :1469	Mean :241.2
##	3rd Qu.:1915.5	3rd Qu.: 92.00	3rd Qu.:1537	3rd Qu.:273.0
##	Max. :2535.0	Max. :146.00	Max. :2554	Max. :458.0
##				
##	BATTING_3B	BATTING_HR	BATTING_BB	BATTING_SO
##	Min. : 0.00	Min. : 0.00	Min. : 0.0	Min. : 0.0
##	1st Qu.: 34.00	1st Qu.: 42.00	1st Qu.:451.0	1st Qu.: 548.0
##	Median : 47.00	Median :102.00	Median :512.0	Median : 750.0
##	Mean : 55.25	Mean : 99.61	Mean :501.6	Mean : 735.6
##	3rd Qu.: 72.00	3rd Qu.:147.00	3rd Qu.:580.0	3rd Qu.: 930.0
##	Max. :223.00	Max. :264.00	Max. :878.0	Max. :1399.0
##				NA's :102
##	BASERUN_SB	BASERUN_CS	BATTING_HBP	PITCHING_H
##	Min. : 0.0	Min. : 0.0	Min. :29.00	Min. : 1137
##	1st Qu.: 66.0	1st Qu.: 38.0	1st Qu.:50.50	1st Qu.: 1419
##	Median :101.0	Median : 49.0	Median :58.00	Median : 1518
##	Mean :124.8	Mean : 52.8	Mean :59.36	Mean : 1779
##	3rd Qu.:156.0	3rd Qu.: 62.0	3rd Qu.:67.00	3rd Qu.: 1682
##	Max. :697.0	Max. :201.0	Max. :95.00	Max. :30132
##	NA's :131	NA's :772	NA's :2085	
##	PITCHING_HR	PITCHING_BB	PITCHING_SO	FIELDING_E
##	Min. : 0.0	Min. : 0.0	Min. : 0.0	Min. : 65.0

```
## 1st Qu.: 50.0    1st Qu.: 476.0    1st Qu.: 615.0    1st Qu.: 127.0
## Median :107.0    Median : 536.5    Median : 813.5    Median : 159.0
## Mean   :105.7    Mean   : 553.0    Mean   : 817.7    Mean   : 246.5
## 3rd Qu.:150.0    3rd Qu.: 611.0    3rd Qu.: 968.0    3rd Qu.: 249.2
## Max.    :343.0    Max.    :3645.0    Max.    :19278.0   Max.    :1898.0
##
##           NA's    :102
## FIELDING_DP
## Min.      : 52.0
## 1st Qu.   :131.0
## Median    :149.0
## Mean      :146.4
## 3rd Qu.   :164.0
## Max.      :228.0
## NA's      :286
```

We note that a number of columns have NAs. Batting_SO and Pitching_SO have the same number of NA's and may be related.

We more closely examine the distribution of columns in the dataset (fig. 1):

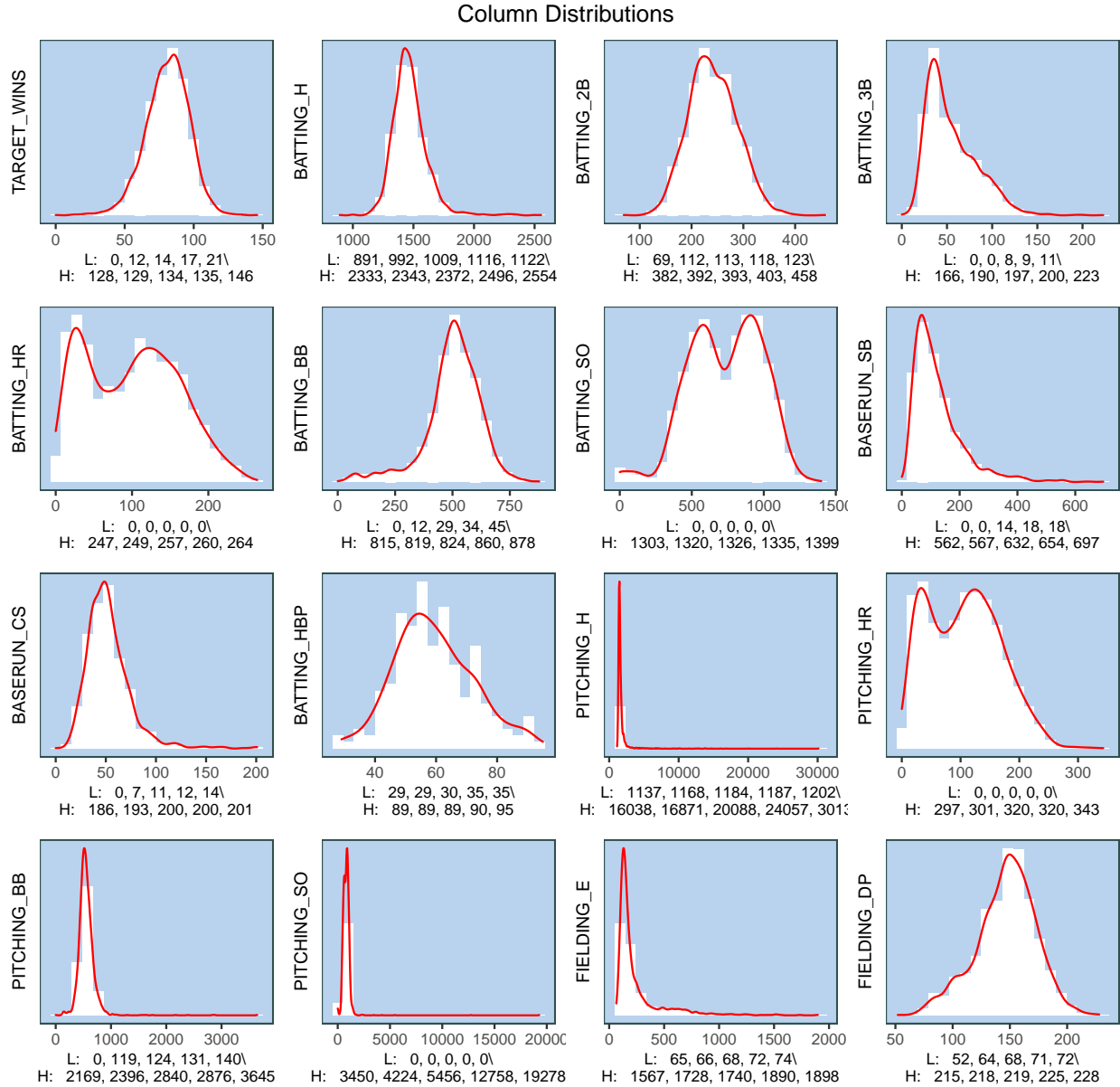


Fig. 1

Our dependent variable (Target Wins) appears to be normally distributed. However, a number of columns are severely skewed (Errors, Strikeouts, Pitching_H, etc.) A few columns (Batting SO, Pitching_HR and Batting_HR) have a bimodal distribution. This might point to some hidden groupings in the dataset.

Boxplots help us identify outliers (fig. 2):

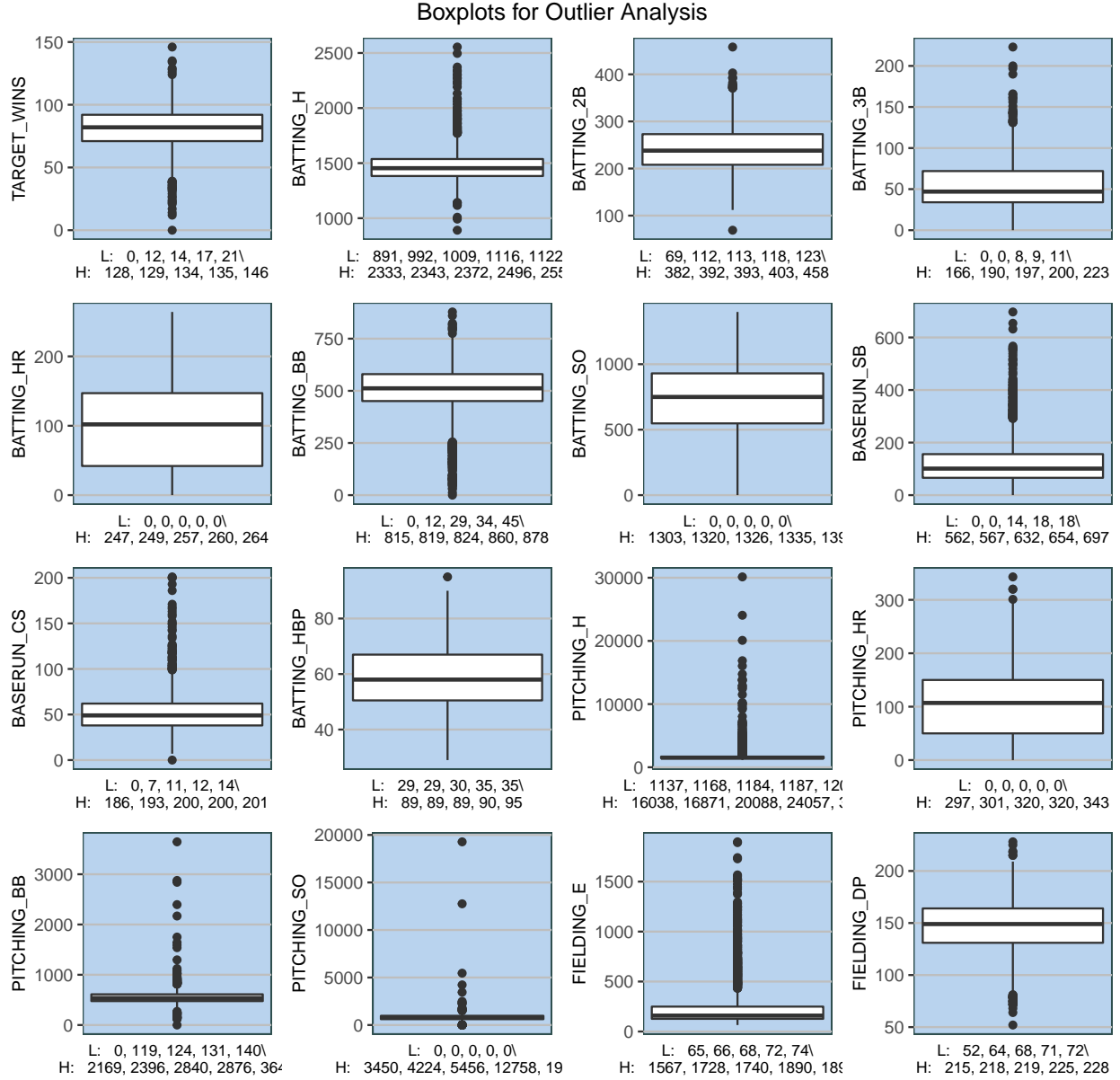


Fig. 2

There a number of outliers, both high and low. For example, there are many zeros, which may be implausible. In addition, many of the ranges appear extreme, such as giving up between 3,500 hits and 19,000 hits, or getting from 12 to over 800 walks.

We investigate correlations in the dataset, both between the dependent variable and the other variables (fig. 3), and between the dependent variables and each other (fig. 4).

Scatterplots Against TARGET_WINS

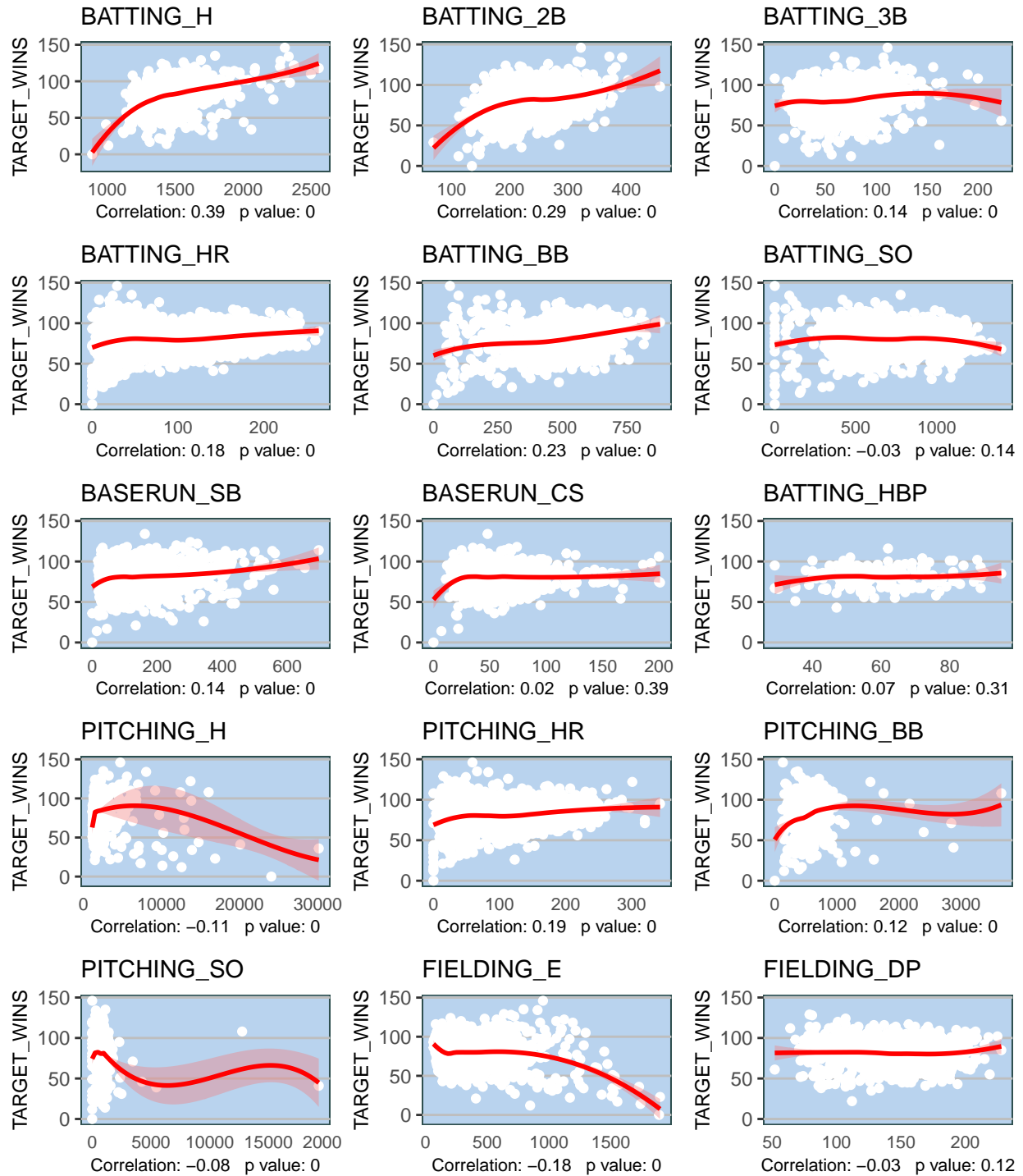


Fig.3

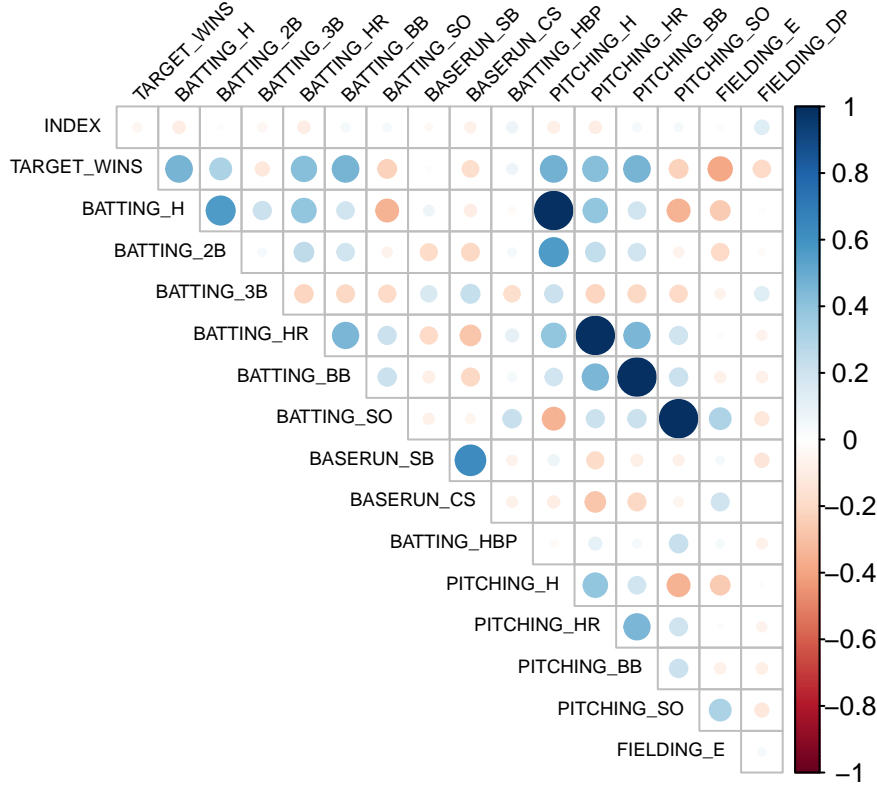
Here we see a number of puzzles, mainly among the pitching correlations. Hits should show a much stronger negative correlation, and in fact appear positive for a portion. Making double plays is surprisingly neutral, as are strikeouts. Pitching_HR is also positive when we would expect negative.

We do need to acknowledge here the possibility of strategy groupings (defense and offense) which may contribute to these anomalies. In other words, a team with poor pitching may have strong hitting, which

then wins games.

We can look for evidence of this possibility by examining multicollinearity:

Correlations, Fig. 4



Indeed, the pitching categories are strongly correlated with their hitting counterparts. All four of the pitching categories follow this pattern.

2. Data Preparation

We begin by devising a strategy for the NAs. We can eliminate the Batting_HBP and Baserun_CS columns because they have too many NA's. We also create flags for the other columns with significant NA's.

We are particularly interested in the SO columns because they do not appear random, and investigation establishes that they have complete overlap with each other and significantly overlap Baserun_SB as well. While not MCAR (missing completely at random), if they are nonetheless MAR (missing at random), we can simply eliminate these rows, as there are not so many (5% of the total).

One way to investigate the randomness of this missing cohort is to look for interactions between the cohort and other dataset columns. In fact, there are a number of columns with strong, even extreme interactions (see fig. 5).

Selected Interactions with Missing Batting_SO

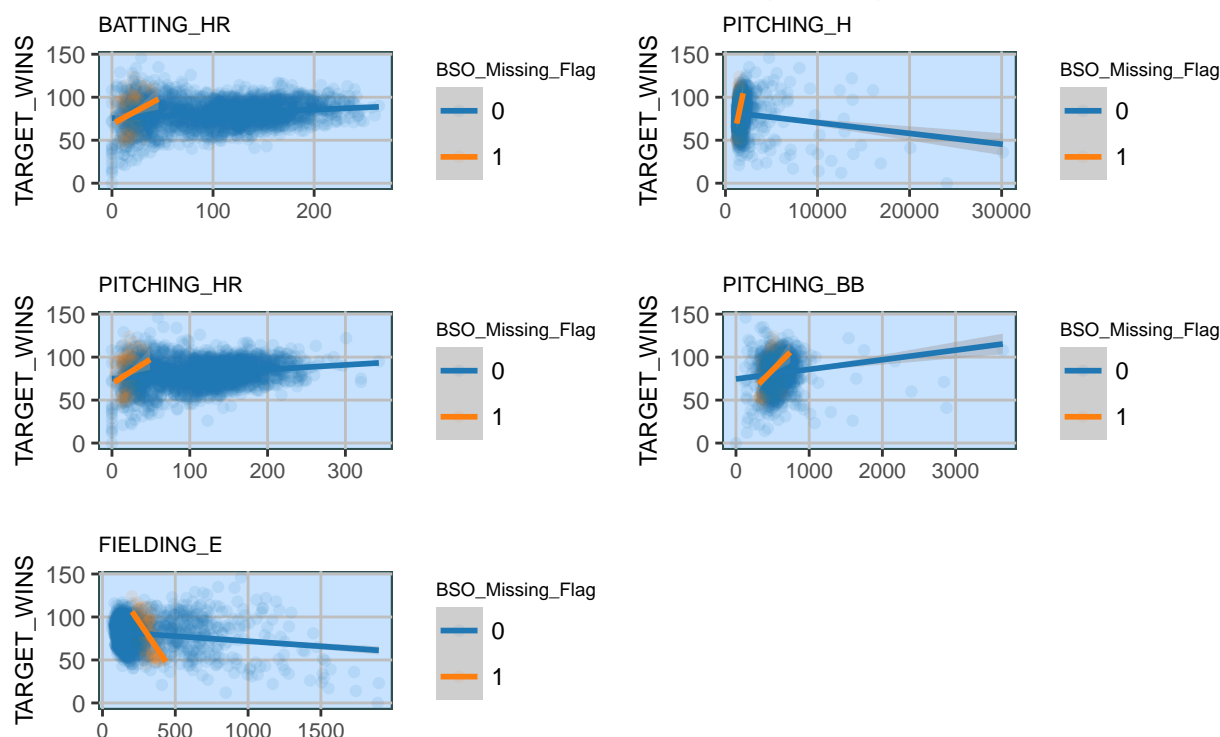


Fig. 5

It is possible this cohort represents a different baseball era when such statistics were not collected. In any case, we cannot eliminate these rows without losing critical data, so we employ the following strategy: 1) retain the rows and impute a value, 2) retain a “missing” flag to keep track of the cohort, and 2) add interaction terms where appropriate.

Before we address imputation, we want to work with the implausible zeros in the dataset. In particular, we note that the 0s in Pitching_SO and Batting_SO are a complete overlap, and that the jump between 0 and the next lowest values is not smooth, and so we will treat them as NA's. We do the same with HR, since there is also a jump up after zero which suggests it is being used as an indicator of missing value.

Just so we have some reasonable criteria for imputation strategy, we compare the r-squared of three regressions - with NA's imputed as means, with NA's imputed as medians, and with NA rows eliminated altogether.

```
## [1] "type:" "mean"
## [1] "r2mean:" "0.4031"
## [1] "r2median:" "0.403"
## [1] "r2omit" "0.4019"
```

The mean and median have the same r-squared, while the elimination of the rows has a smaller r-squared. We therefore choose to impute the mean.

Not surprisingly, the evaluation dataset shows the same results:

```
## [1] "type:" "mean"
## [1] "r2mean:" "0.4031"
## [1] "r2median:" "0.403"
## [1] "r2omit" "0.4019"
```

Although outliers and possible bad data appear in a number of places, without domain knowledge I am reluctant to eliminate outliers or influential points without good reason. We don't know if extreme numbers are necessarily implausible. Therefore the outliers will remain.

3. Data Modeling

1. We create a flag for hits under 1500

As previously noted, Pitching_H is surprisingly weak in its relationship to wins, and in fact appears positive for a large portion of its distribution. We examine more closely the relationship between pitching hits and wins, paying particular attention to the portion of the relationship where hits are below 3,000 (fig. 6).

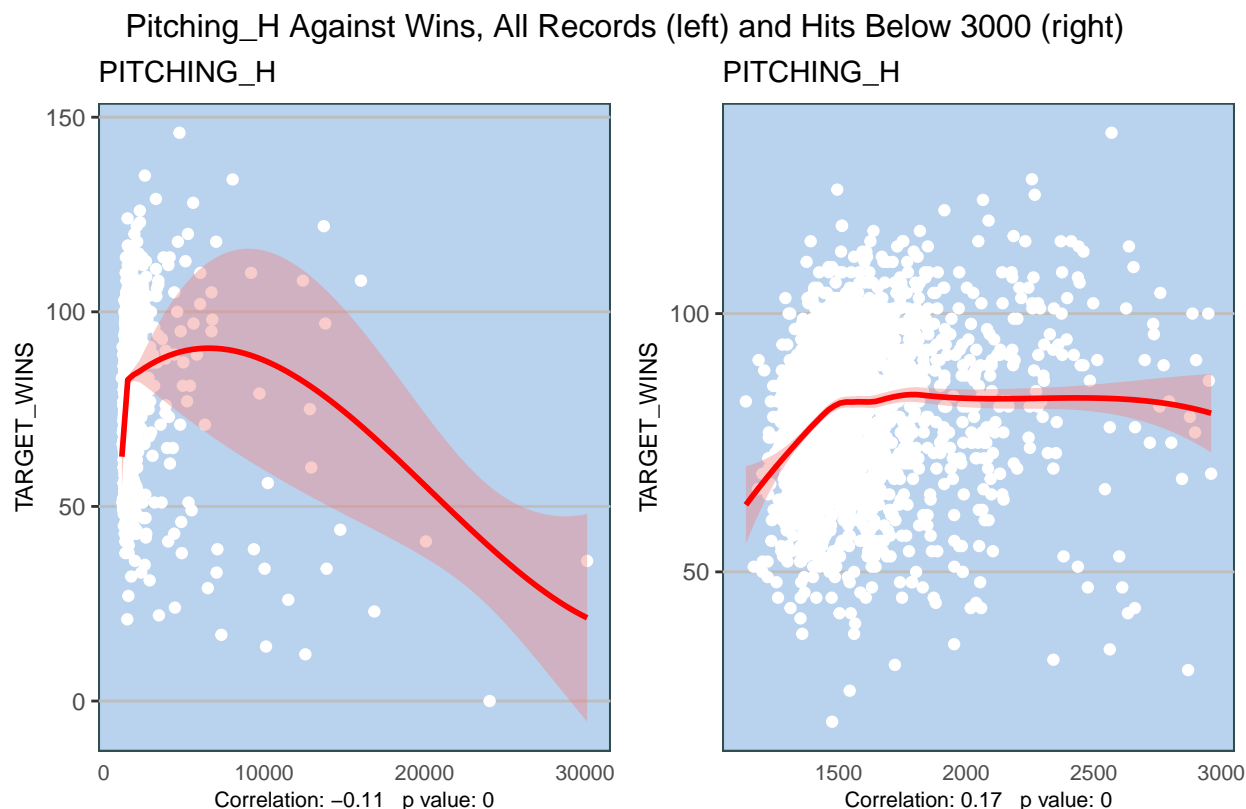


Fig.6

We can see here the positive correlation between pitching_h and wins. While we can't explain the phenomenon, we can account for it statistically by adding a binary flag for records with hits under 1500.

2. We create an interaction between Fielding_DP and hits.

The Fielding_DP correlation with Target Wins is surprising, since making double plays should help a team win. On the other hand, a team that makes double plays is also a team that gives up hits.

We therefore create an interaction term for Fielding_DP and Pitching_H.

3. We drop PITCHING_HR because it is an implausibly close match with HITTING_HR.

Like many pitching columns, Pitching_HR is unexpectedly positively correlated with wins. However, what makes this column truly implausible is how close a match it is with BATTING_HR. The scatterplot below (Fig. 7) shows that the vast majority of the figures for pitching HR are exactly the same or within 2 or 3 of Batting HR. We therefore drop it since this makes no sense.

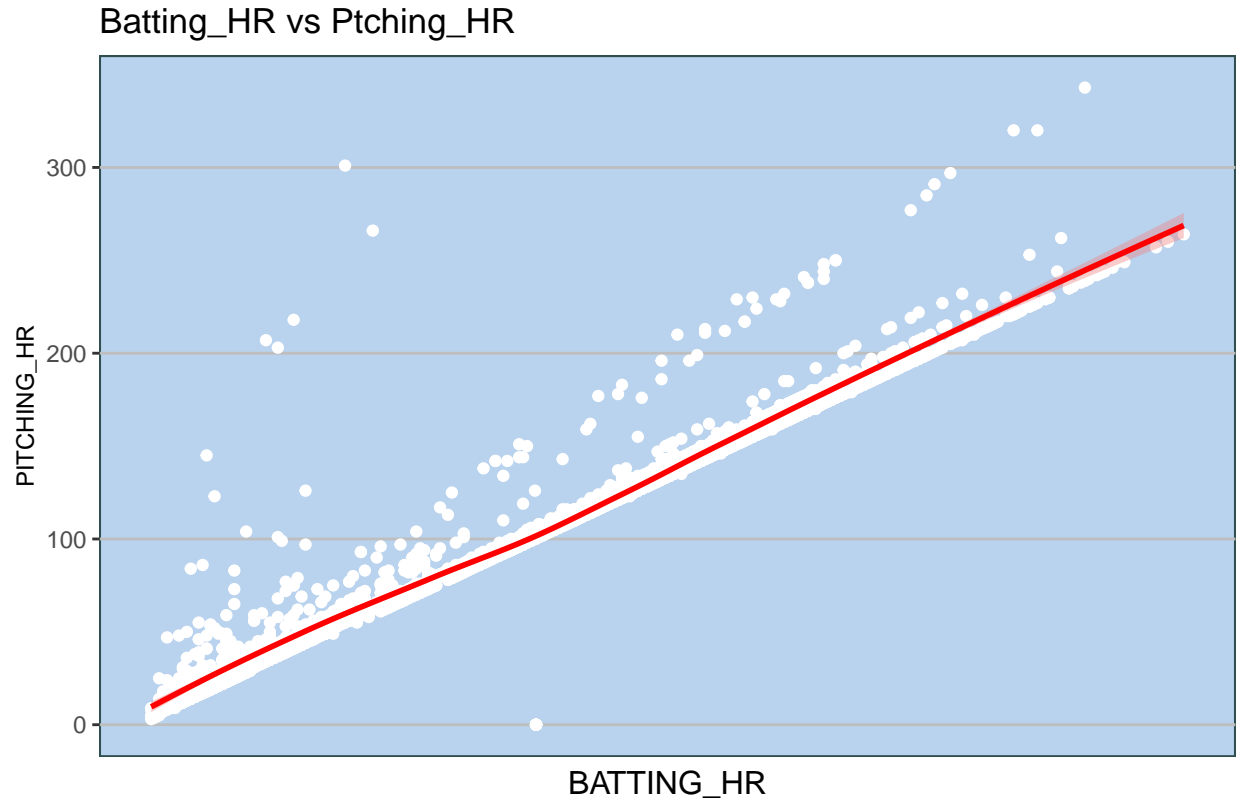


Fig. 7

4. We create a flag to account for the bimodal distribution of Batting HR.

Batting HR has a bimodal distribution (see Fig. 8). We don't explain this, but speculate that it may be related to different eras of baseball. Therefore, we create a flag to separate records with less than 80 HR from those with more.

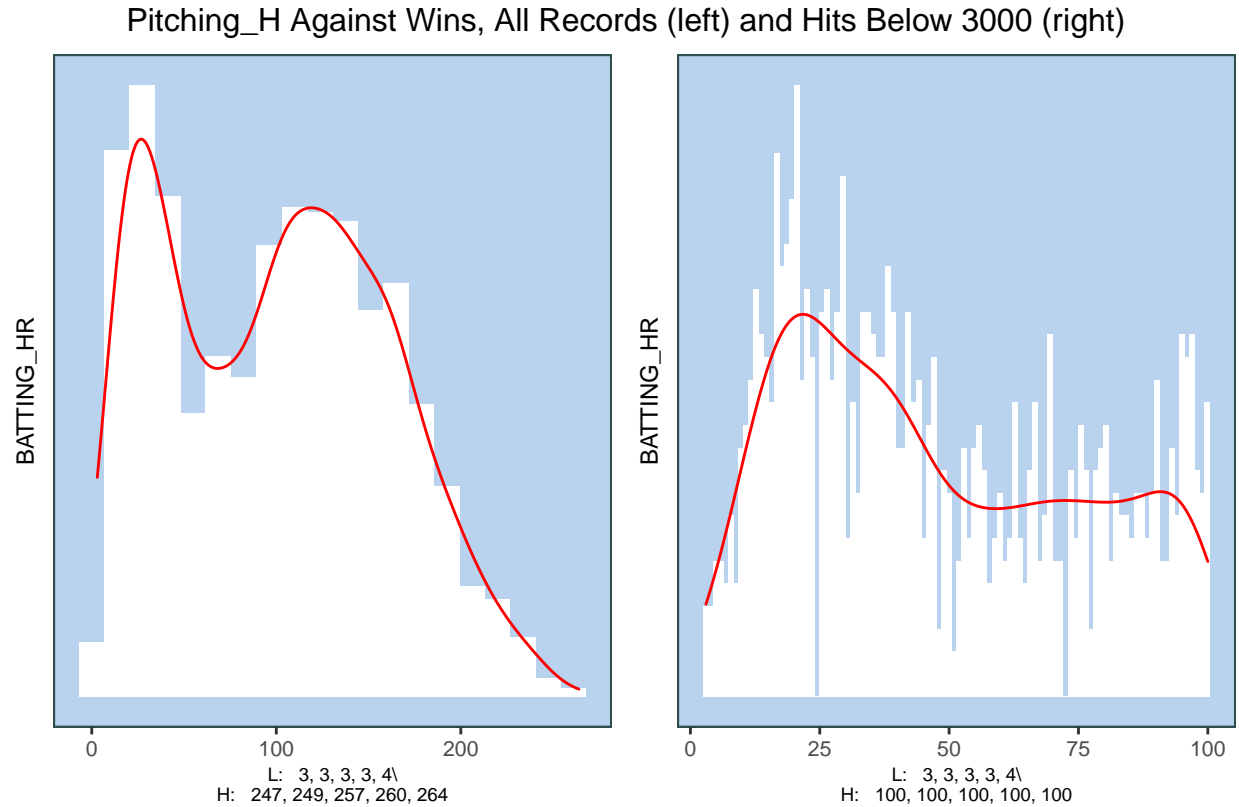


Fig.8

```
## [1] 0.02231354
```

```
## [1] 0.03503598
```

5. We transform the error variable.

```
## [1] 0.03072081
```

```
## [1] 0.04825783
```

Create Missing “cohort” interactions

Build Models

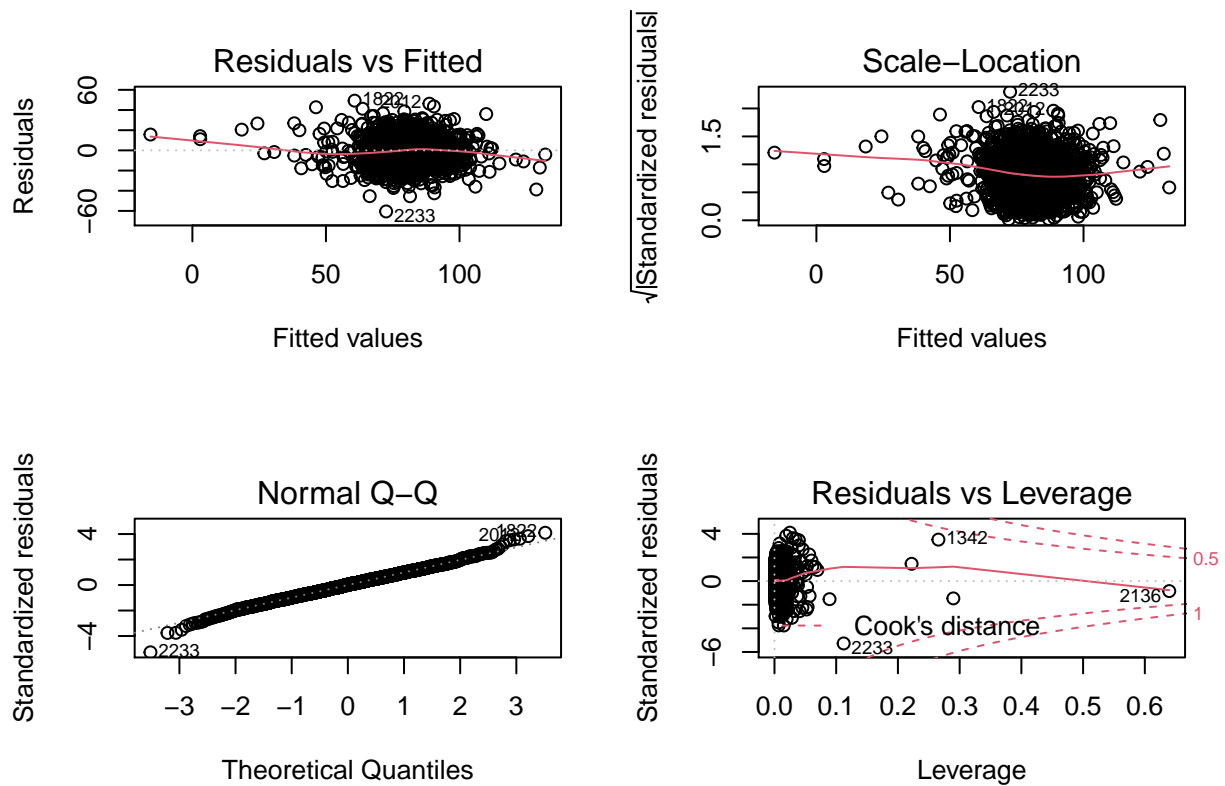
Create regression 1 - no transformations except missing flags

```
##
## Call:
## lm(formula = TARGET_WINS ~ BATTING_H + BATTING_2B + BATTING_3B +
##   BATTING_HR + BATTING_BB + BATTING_SO + BASERUN_SB + PITCHING_H +
##   PITCHING_SO + FIELDING_E + FIELDING_DP + BSO_Missing_Flag +
##   BRBSB_Missing_Flag + FDP_Missing_Flag, data = df)
##
```

```

## Residuals:
##      Min       1Q   Median       3Q      Max
## -60.531  -8.063   0.330   8.075  49.266
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    13.7948052   5.0143117   2.751  0.00599 **
## BATTING_H       0.0521109   0.0033520  15.546 < 2e-16 ***
## BATTING_2B     -0.0401259   0.0086621  -4.632 3.82e-06 ***
## BATTING_3B      0.0537762   0.0158617   3.390 0.00071 ***
## BATTING_HR      0.0595856   0.0089648   6.647 3.75e-11 ***
## BATTING_BB      0.0260490   0.0032618   7.986 2.20e-15 ***
## BATTING_SO     -0.0066440   0.0022278  -2.982 0.00289 **
## BASERUN_SB      0.0477764   0.0046194  10.343 < 2e-16 ***
## PITCHING_H      0.0018926   0.0003398   5.569 2.86e-08 ***
## PITCHING_SO    -0.0013966   0.0006654  -2.099 0.03593 *
## FIELDING_E     -0.0560670   0.0033748 -16.613 < 2e-16 ***
## FIELDING_DP    -0.0969459   0.0134629  -7.201 8.10e-13 ***
## BSO_Missing_Flag  8.3474206   1.4721894   5.670 1.61e-08 ***
## BRSB_Missing_Flag 34.1064444   1.8484454  18.451 < 2e-16 ***
## FDP_Missing_Flag  4.2303099   1.4669785   2.884 0.00397 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.17 on 2261 degrees of freedom
## Multiple R-squared:  0.4068, Adjusted R-squared:  0.4031
## F-statistic: 110.7 on 14 and 2261 DF, p-value: < 2.2e-16
##
## [1] "VIF Analysis"
##      BATTING_H      BATTING_2B      BATTING_3B      BATTING_HR
##      3.608349      2.524443      3.016545      4.444255
##      BATTING_BB      BATTING_SO      BASERUN_SB      PITCHING_H
##      2.459248      4.131567      2.380761      3.511045
##      PITCHING_SO      FIELDING_E      FIELDING_DP      BSO_Missing_Flag
##      1.946738      9.076248      1.674220      1.425731
##      BRSB_Missing_Flag FDP_Missing_Flag
##      2.848146      3.633432

```



```
## NULL
```

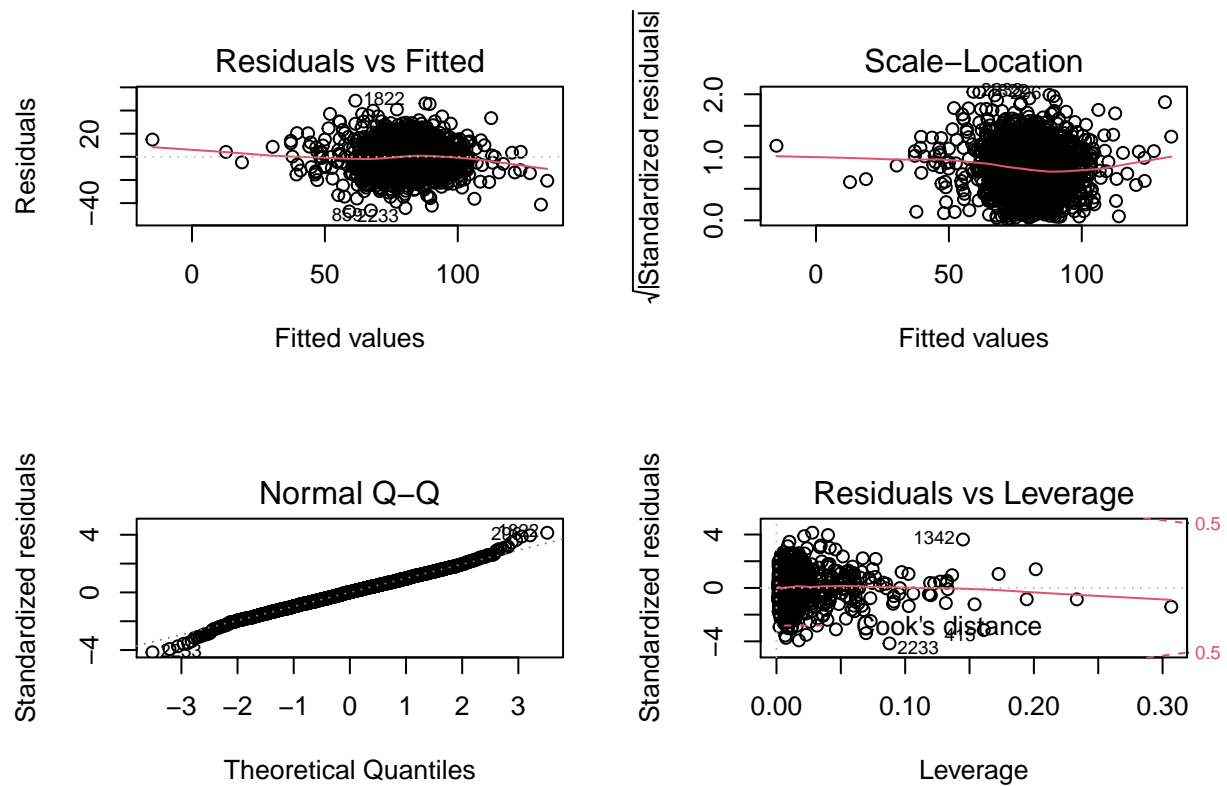
Model 2, all the transformations:

```
##
## Call:
## lm(formula = TARGET_WINS ~ BATTING_H + BATTING_2B + BATTING_3B +
##   BATTING_HR + BATTING_BB + BATTING_SO + BASERUN_SB + PITCHING_H +
##   FIELDING_E + FIELDING_DP + BSO_Missing_Flag + BRSB_Missing_Flag +
##   FDP_Missing_Flag + Pitch_h_Under1500 + Prod_DP_H + E_sq +
##   Inter_bb_Cohort + Inter_E_Cohort + Inter_bhr_Cohort + Inter_bbb_Cohort +
##   Inter_bs_Cohort, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -47.202  -7.806   0.193   7.821  48.504
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.472e+01  6.702e+00   3.688 0.000231 ***
## BATTING_H       5.622e-02  3.302e-03  17.023 < 2e-16 ***
## BATTING_2B     -4.125e-02  8.586e-03  -4.805 1.65e-06 ***
## BATTING_3B      6.743e-02  1.610e-02   4.188 2.93e-05 ***
## BATTING_HR      5.825e-02  8.978e-03   6.488 1.06e-10 ***
## BATTING_BB      2.593e-02  3.247e-03   7.984 2.23e-15 ***
```

```

## BATTING_SO      -1.223e-02  2.218e-03  -5.512  3.95e-08  ***
## BASERUN_SB      5.238e-02  4.795e-03  10.923  < 2e-16  ***
## PITCHING_H      -4.629e-03  2.995e-03  -1.546  0.122287
## FIELDING_E      -8.282e-02  7.453e-03  -11.112  < 2e-16  ***
## FIELDING_DP     -1.646e-01  3.571e-02  -4.610  4.25e-06  ***
## BSO_Missing_Flag  5.042e+01  1.190e+01  4.237  2.36e-05  ***
## BRBSB_Missing_Flag 3.794e+01  2.023e+00  18.752  < 2e-16  ***
## FDP_Missing_Flag  5.282e+00  1.713e+00  3.084  0.002064  **
## Pitch_h_Under1500 2.214e+00  6.829e-01  3.242  0.001206  **
## Prod_DP_H        3.671e-05  2.040e-05  1.799  0.072094  .
## E_sq             2.143e-05  4.284e-06  5.002  6.11e-07  ***
## Inter_bb_Cohort  1.336e-01  8.560e-02  1.560  0.118847
## Inter_E_Cohort   -1.938e-01  2.809e-02  -6.899  6.77e-12  ***
## Inter_bhr_Cohort  3.652e-01  1.546e-01  2.362  0.018285  *
## Inter_bbb_Cohort -1.397e-01  9.536e-02  -1.465  0.143105
## Inter_bs_Cohort  3.896e-02  2.653e-02  1.469  0.142097
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.89 on 2254 degrees of freedom
## Multiple R-squared:  0.4357, Adjusted R-squared:  0.4305
## F-statistic: 82.88 on 21 and 2254 DF,  p-value: < 2.2e-16
##
## [1] "VIF Analysis"
##      BATTING_H      BATTING_2B      BATTING_3B      BATTING_HR
##      3.670470      2.599487      3.258269      4.671215
##      BATTING_BB      BATTING_SO      BASERUN_SB      PITCHING_H
##      2.554246      4.292413      2.688323      285.743188
##      FIELDING_E      FIELDING_DP      BSO_Missing_Flag      BRBSB_Missing_Flag
##      46.392975      12.345591      97.619815      3.575494
##      FDP_Missing_Flag      Pitch_h_Under1500      Prod_DP_H      E_sq
##      5.189563      1.863892      282.770320      24.339858
##      Inter_bb_Cohort      Inter_E_Cohort      Inter_bhr_Cohort      Inter_bbb_Cohort
##      1095.501873      50.949700      7.645153      1173.699962
##      Inter_bs_Cohort
##      21.438192

```



NULL

second model explains better, but does not necessarily perform lot better.

Third model, categories of power - batting power and pitching weakness categories

The two are correlated

[1] 0.3967352

These boxplots show the stronger relationship with batting power

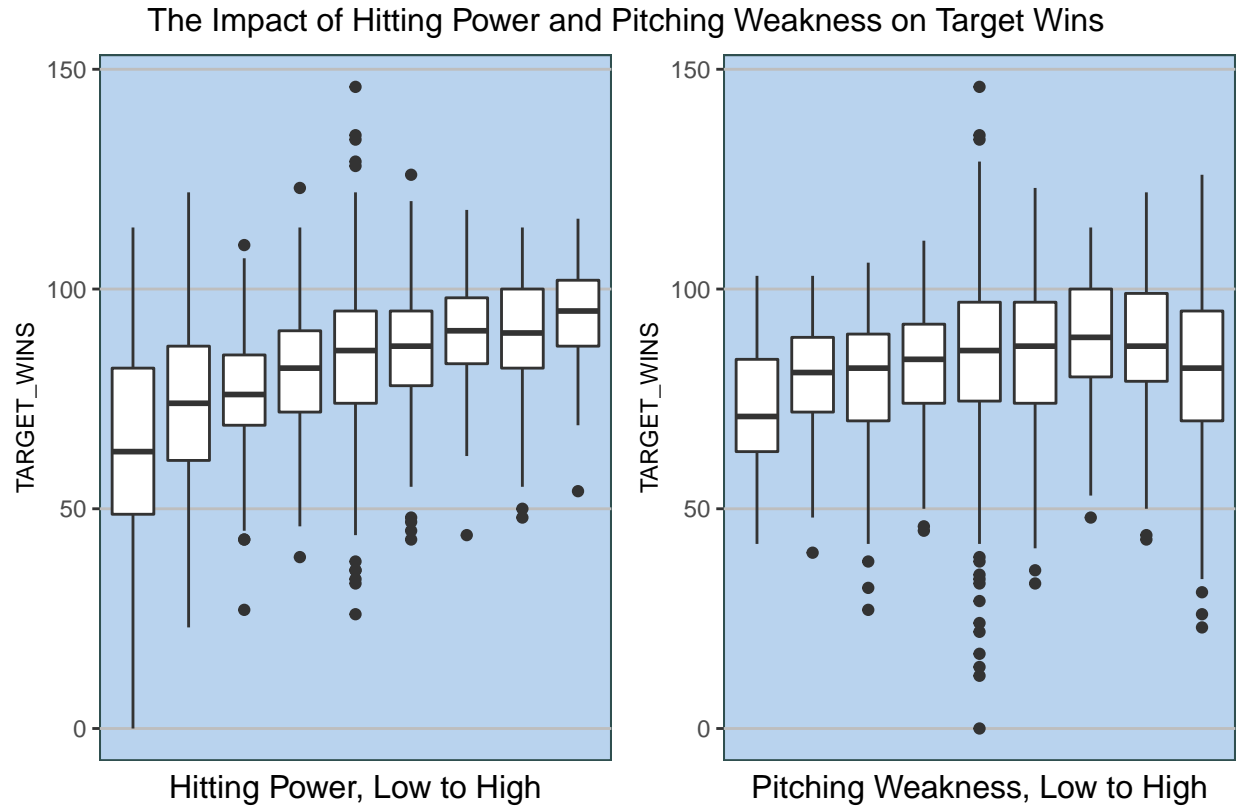


Fig. 8

we run the regressions

```
##
## Call:
## lm(formula = TARGET_WINS ~ Total_Power, data = dfCat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -74.784 -10.648   1.216  10.333  63.294
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  82.7062    0.4841  170.840  <2e-16 ***
## Total_Power   1.9804    0.2154   9.196  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 16.78 on 1200 degrees of freedom
## Multiple R-squared:  0.06583,    Adjusted R-squared:  0.06505
## F-statistic: 84.57 on 1 and 1200 DF,  p-value: < 2.2e-16

##
## Call:
## lm(formula = TARGET_WINS ~ Hitting_Power + Pitching_Weakness,
##     data = dfCat)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -68.817  -9.239   0.898  10.008  63.261
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    64.2645     1.6750  38.367  <2e-16 ***
## Hitting_Power     3.4805     0.2429  14.328  <2e-16 ***
## Pitching_Weakness -0.4014     0.2467  -1.627   0.104
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 15.94 on 1199 degrees of freedom
## Multiple R-squared:  0.1579, Adjusted R-squared:  0.1565
## F-statistic: 112.4 on 2 and 1199 DF,  p-value: < 2.2e-16

##
## Call:
## lm(formula = TARGET_WINS ~ category_PH + category_PBB + category_BH +
##     category_BBB + category_BHR, data = dfCat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -70.171  -8.976   1.030  10.344  61.136
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    62.5490     2.0517  30.487  < 2e-16 ***
## category_PH     -0.1106     0.4845  -0.228  0.81953
## category_PBB    -0.3362     0.6472  -0.520  0.60350
## category_BH      4.0065     0.3883  10.319  < 2e-16 ***
## category_BBB     2.5059     0.7097   3.531  0.00043 ***
## category_BHR     0.6661     0.4149   1.605  0.10871
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 15.89 on 1196 degrees of freedom
## Multiple R-squared:  0.1656, Adjusted R-squared:  0.1621
## F-statistic: 47.47 on 5 and 1196 DF,  p-value: < 2.2e-16
```

Analysis shows good batting and weak pitching are correlated. Poor r squared but significant batting.

Select models

Now we make predictions