

Moneyball

We are about to recreate the famous moneyball model.

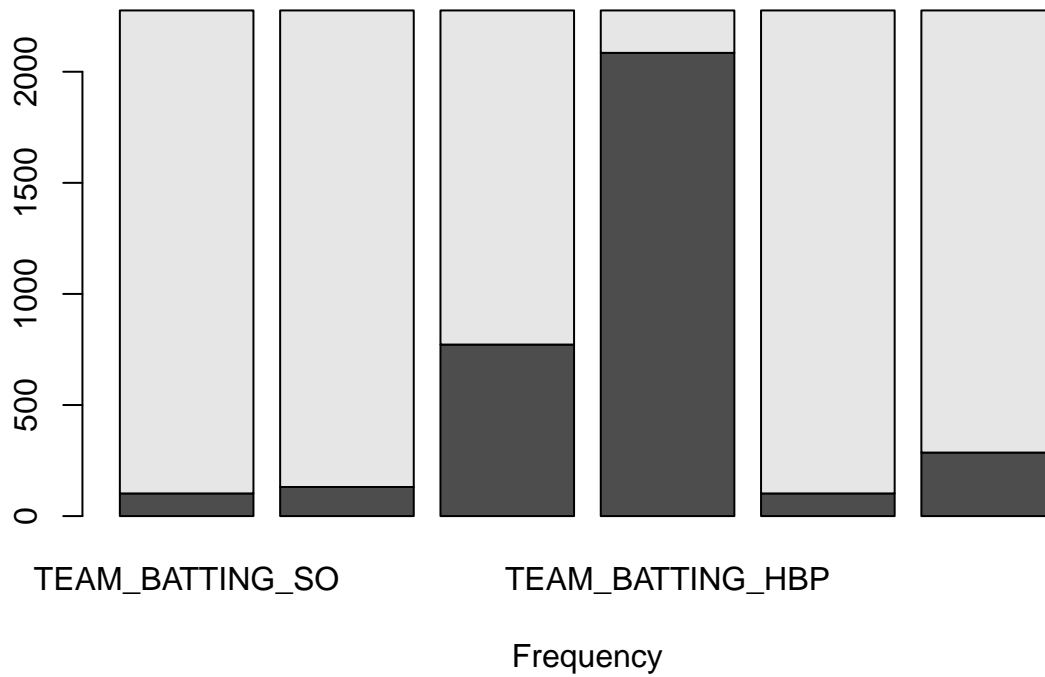
VARIABLE NAME	DEFINITION	THEORETICAL EFFECT
INDEX	Identification Variable (do not use)	None
TARGET_WINS	Number of wins	
TEAM_BATTING_H	Base Hits by batters (1B,2B,3B,HR)	Positive Impact on Wins
TEAM_BATTING_2B	Doubles by batters (2B)	Positive Impact on Wins
TEAM_BATTING_3B	Triples by batters (3B)	Positive 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	Strike outs 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
TEAM_PITCHING_SO	Strike outs by pitchers	Positive Impact on Wins

DATA EXPLORATION and PREPARATION

Missing information

```
barplot(tr.binmat[,imp+1], main = "Missing information, training set",  
        xlab = "Frequency")
```

Missing information, training set

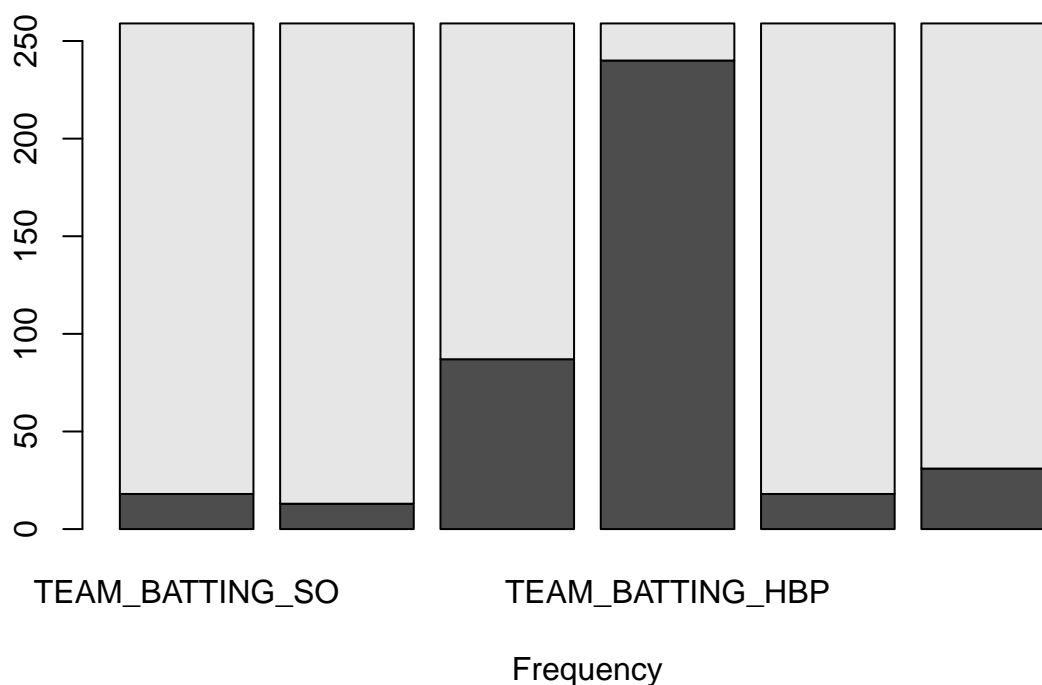


```
colSums(is.na(tr))
```

```
##          INDEX      TARGET_WINS  TEAM_BATTING_H  TEAM_BATTING_2B
##           0           0           0           0
## TEAM_BATTING_3B TEAM_BATTING_HR  TEAM_BATTING_BB  TEAM_BATTING_SO
##           0           0           0           102
## TEAM_BASERUN_SB TEAM_BASERUN_CS TEAM_BATTING_HBP  TEAM_PITCHING_H
##          131          772          2085           0
## TEAM_PITCHING_HR TEAM_PITCHING_BB TEAM_PITCHING_SO  TEAM_FIELDING_E
##           0           0           102           0
## TEAM_FIELDING_DP
##          286
```

```
barplot(ev.binmat[,imp], main = "Missing information, evaluation set",
        xlab = "Frequency")
```

Missing information, evaluation set



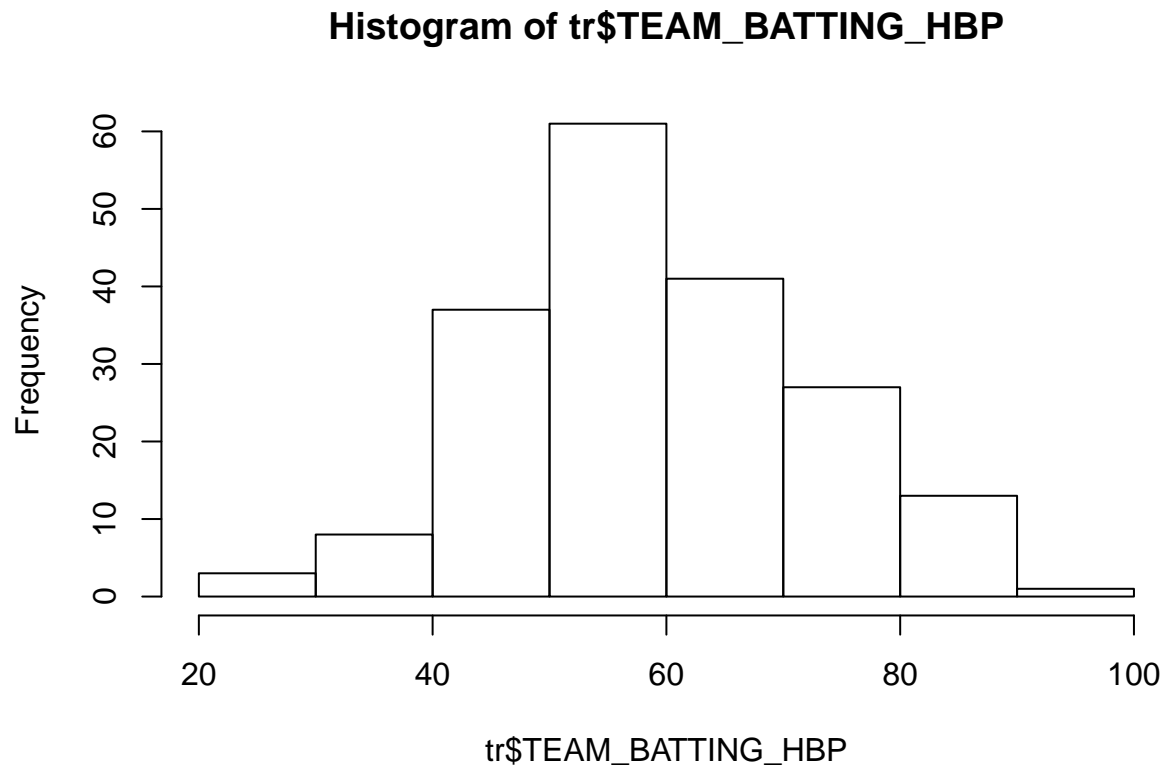
```
colSums(is.na(ev))
```

```
##          INDEX  TEAM_BATTING_H  TEAM_BATTING_2B  TEAM_BATTING_3B
##           0          0          0          0
## TEAM_BATTING_HR  TEAM_BATTING_BB  TEAM_BATTING_SO  TEAM_BASERUN_SB
##           0          0          18          13
## TEAM_BASERUN_CS  TEAM_BATTING_HBP  TEAM_PITCHING_H  TEAM_PITCHING_HR
##           87          240          0          0
## TEAM_PITCHING_BB  TEAM_PITCHING_SO  TEAM_FIELDING_E  TEAM_FIELDING_DP
##           0          18          0          31
```

TEAM_BATTING_HBP is missing almost every value, so we may have to destroy this variable. For now we will add a flag and impute the values. . . After the other things are filled. This variable represents batters hit by a pitch, and none of the values are zero, so these are truly missing.

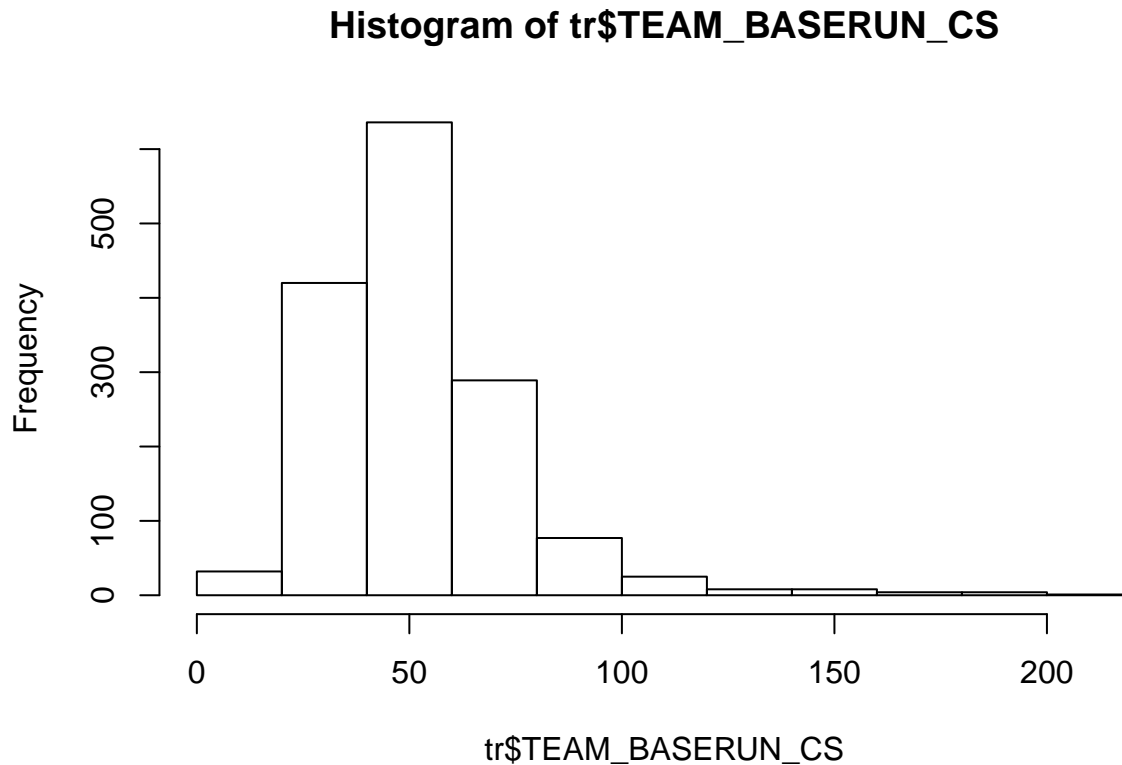
The main question for this dataset, is whether or not imputing this hige missing column will help or hurt the models.

```
hist(tr$TEAM_BATTING_HBP)
```



Apart from that, TEAM_BASERUN_CS is missing about a third of its values, and it seems like a good idea to look at the distribution of known values and add a flag for it, too. This is players caught stealing bases. This information also makes sense. This appears to be a skewed normal distribution ranging from zero off near 200. A good candidate for imputation.

```
hist(tr$TEAM_BASERUN_CS)
```



The other variables have relatively low amounts of missing data. These two exhibit very cooperative distributions, and we will impute around them.

The distributions in the training and evaluation sets are nearly identical, so we can make a generic function for this.

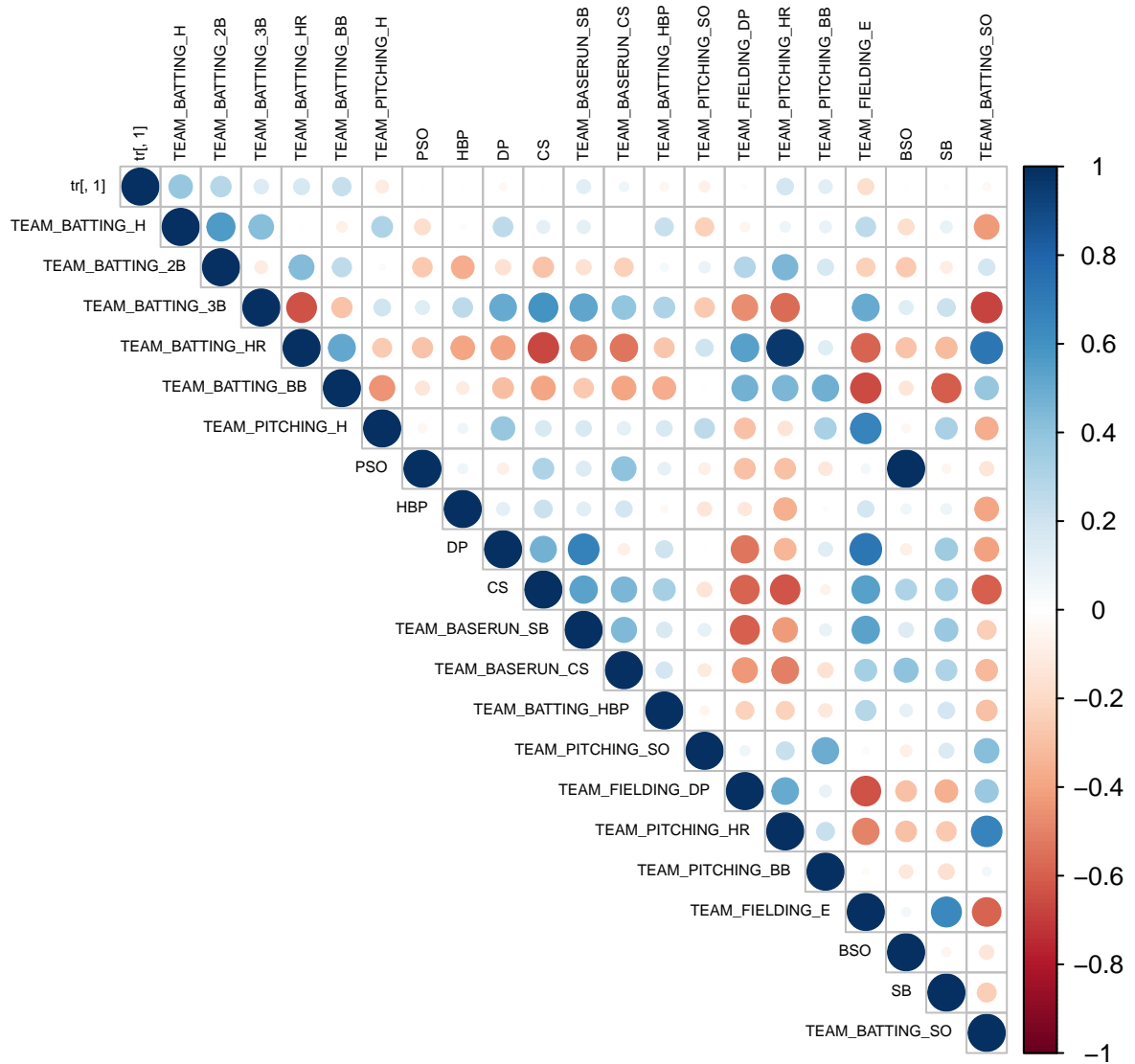
Munging

Our data munging process mostly just deals with adding flags for missing information and imputing these points based on normal distributions. The information presented is 100% numeric, all integers, so we don't have to transform it much.

We will start by imputing the data in the training set alone because it contains the target variable. Then we impute the data from the evaluation set using the filled training set.

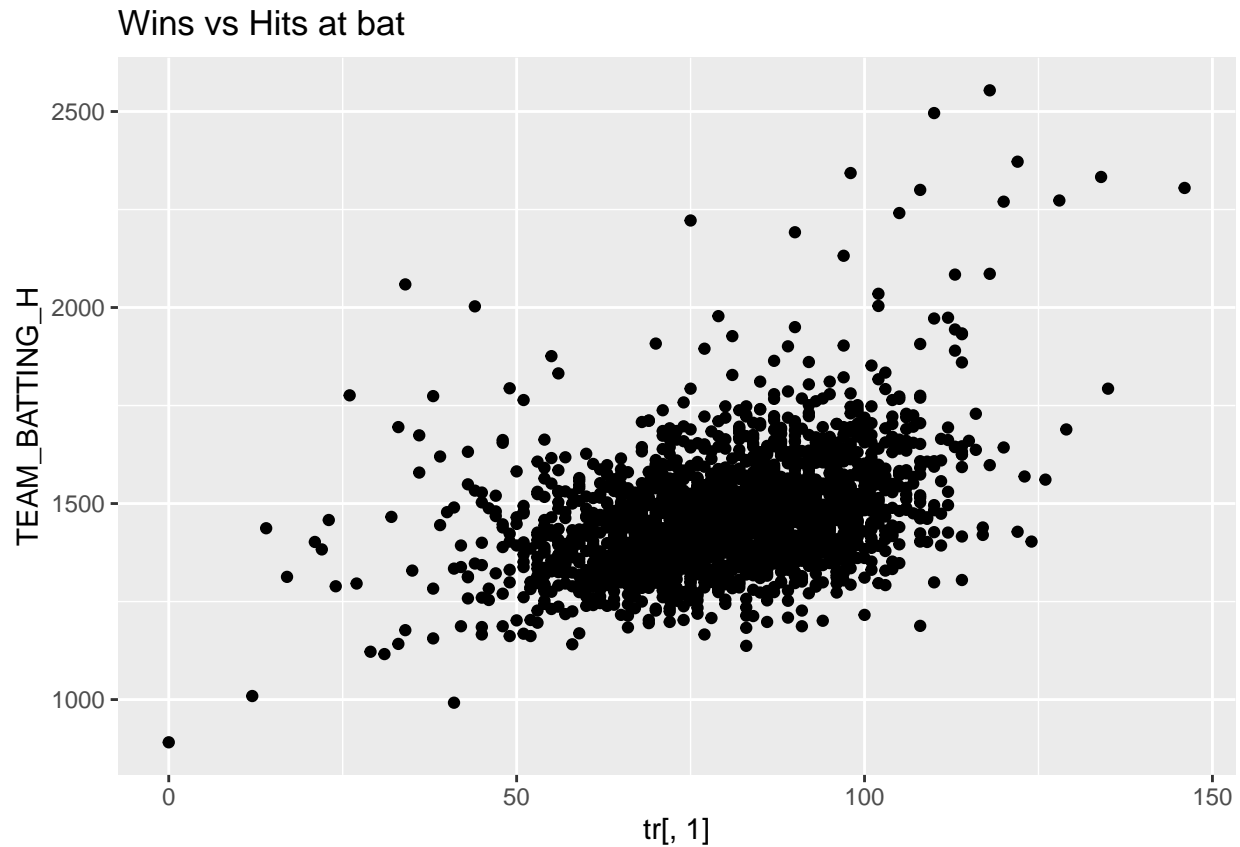
Correlations

```
corrplot(cor(tr.imp),  
         type = 'upper',  
         tl.col = 'black',  
         tl.cex = 0.5)
```



This is a little bit scary. Not a single one of the dependent variables is strongly correlated to the target. There are some interesting pairs, `TEAM_PITCHING_HR` and `TEAM_BATTING_HR` and the flags for `TEAM_PITCHING_SO` and `TEAM_BATTING_SO`, which could be an artifact of the imputation. The correlations look pretty complex and full of information.

```
ggplot(tr.imp, aes(x = tr[, 1], y = TEAM_BATTING_H)) +
  geom_point() + ggtitle(label = "Wins vs Hits at bat")
```



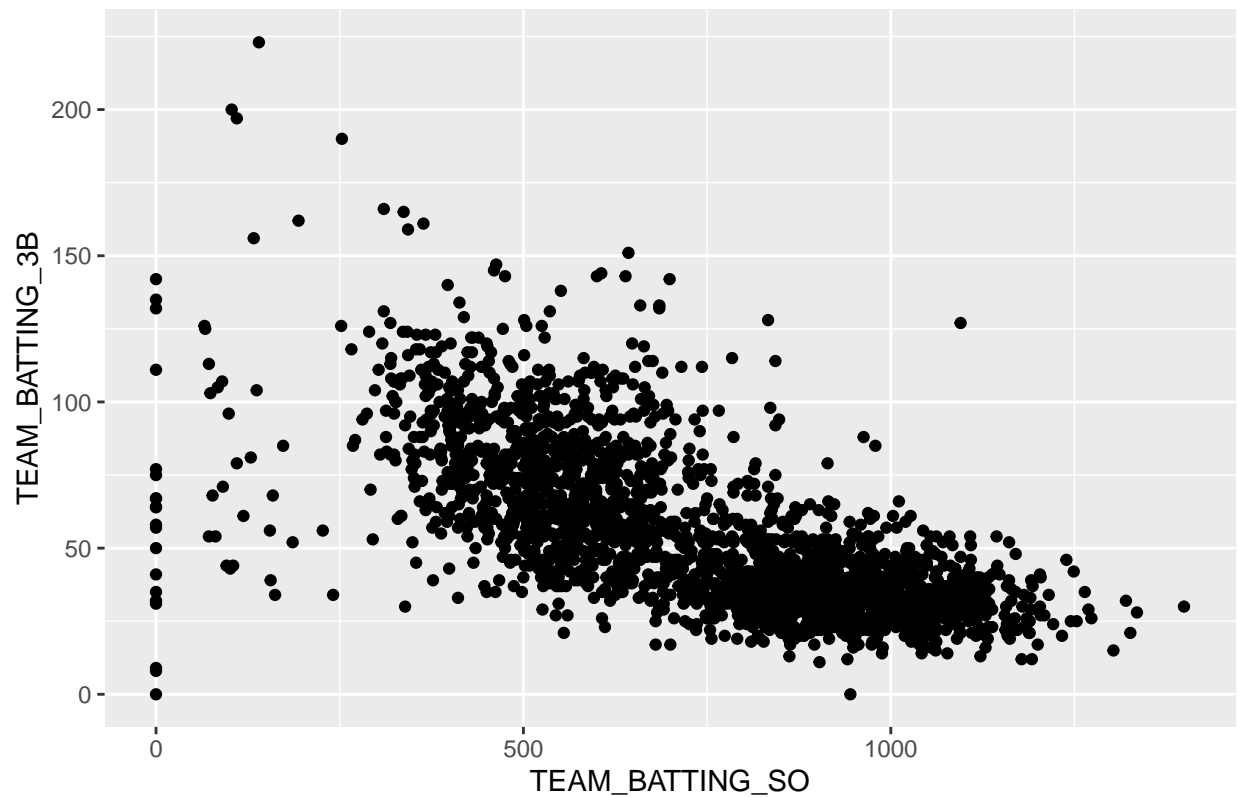
Batting hits is the only datapoint that correlates noticeably with wins. There is a lot more to the story.

Some key distributions

Of the most correlated and anticorrelated, here is a selection:

```
ggplot(tr.imp, aes(x = TEAM_BATTING_SO, y = TEAM_BATTING_3B)) +  
  geom_point() + ggtitle("Batting strikeouts vs batting triples")
```

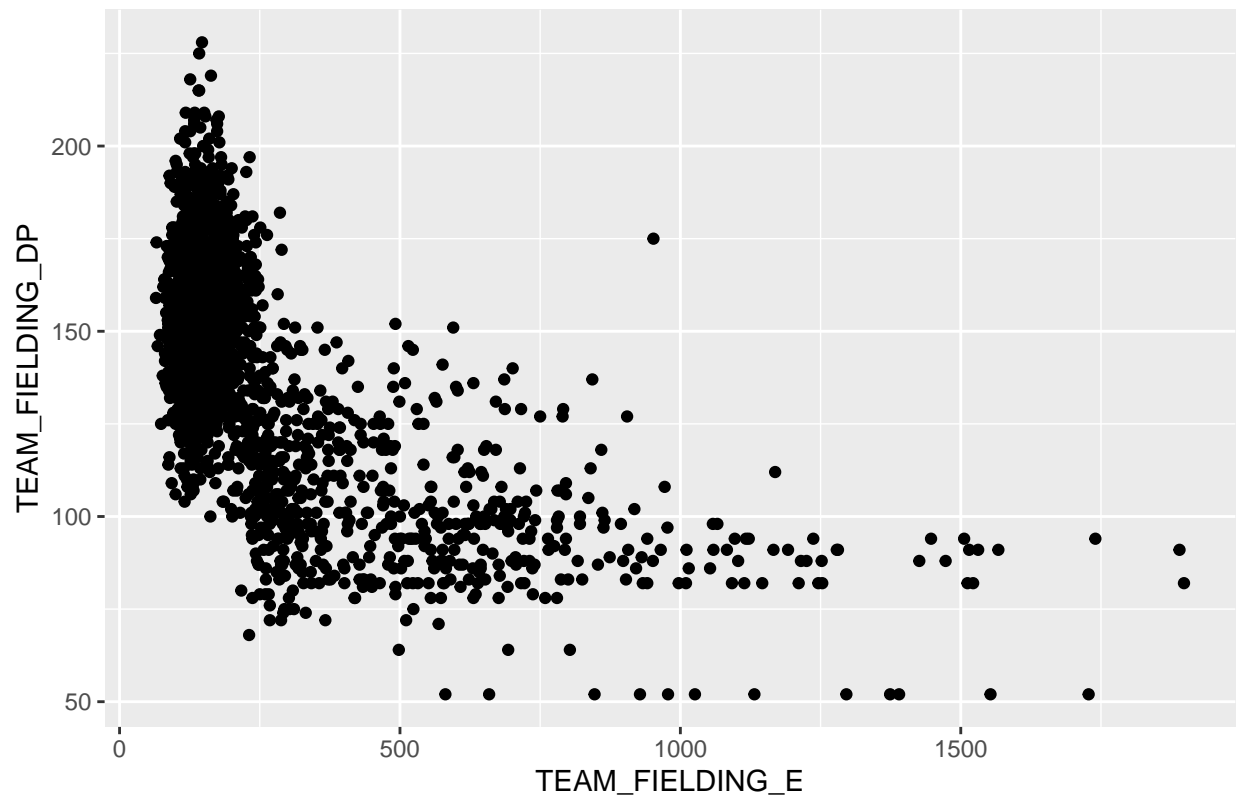
Batting strikeouts vs batting triples



It makes a lot of sense that strikeouts and triples would be anticorrelated. This is good information. There is changing variance in this relationship, so it may be a good feature.

```
ggplot(tr.imp, aes(x = TEAM_FIELDING_E, y = TEAM_FIELDING_DP)) +  
  geom_point() + ggtitle("Fielding errors vs. fielding double plays")
```

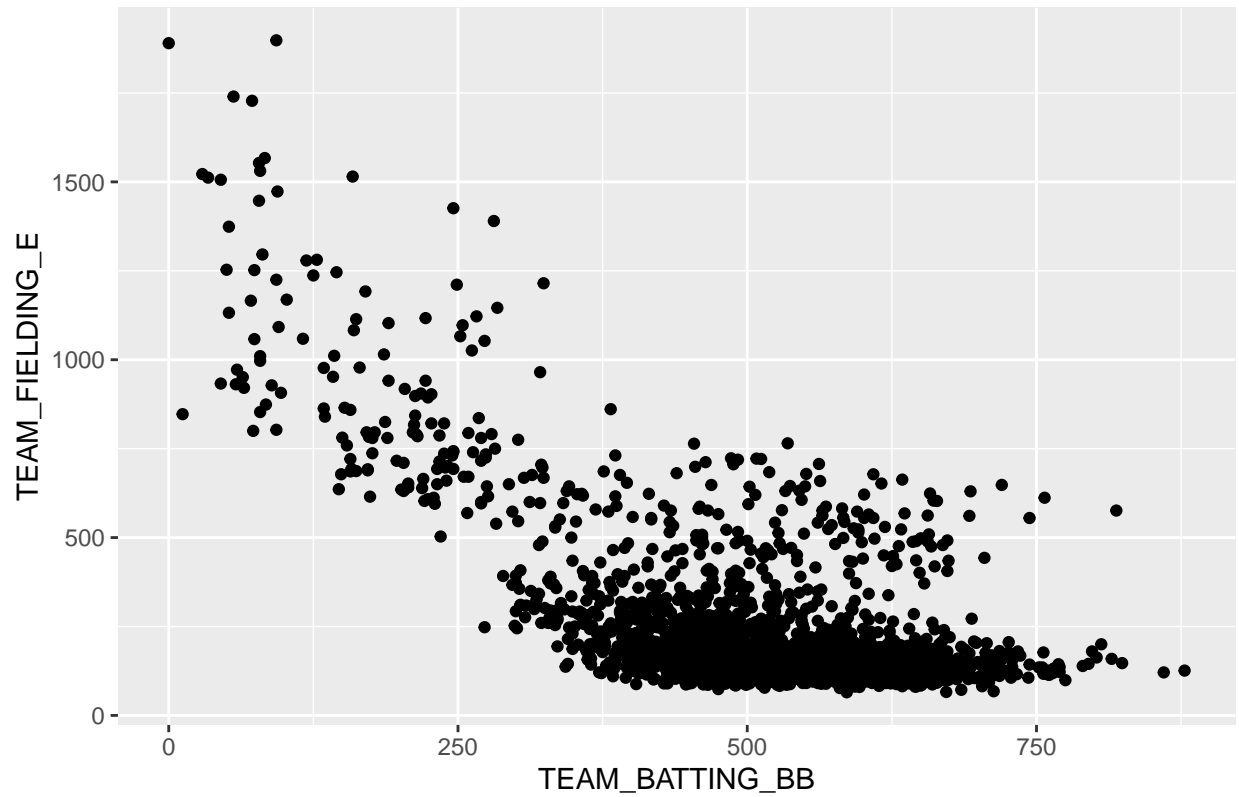

Fielding errors vs. fielding double plays



This may be a very strong feature because the relationship between errors and double plays tells a lot about a team's ability to respond to wild hits.

```
ggplot(tr.imp, aes(x = TEAM_BATTING_BB, y = TEAM_FIELDING_E)) +  
  geom_point() + ggtitle("Walks by batters vs. Fielding error")
```

Walks by batters vs. Fielding error



It makes a lot of sense that teams that walk a lot at bat are also less likely to make fielding errors. This shows a defensive mentality in some teams, and an offensive one in others.

BUILD MODELS

We need to split the data for training and test sets internally, since we don't know the target values in the evaluation set.

The all-in model

```
summary(m1 <- glm(TARGET_WINS ~ . -BSO, tr.tr, family = 'gaussian'))

##
## Call:
## glm(formula = TARGET_WINS ~ . - BSO, family = "gaussian", data = tr.tr)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -44.918  -7.946   0.151   7.855  47.677
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    34.9337395   6.2381941   5.600 2.44e-08 ***
## TEAM_BATTING_H     0.0441621   0.0036461  12.112 < 2e-16 ***
## TEAM_BATTING_2B   -0.0371464   0.0095895  -3.874 0.000111 ***
## TEAM_BATTING_3B    0.0543852   0.0166114   3.274 0.001079 **
## TEAM_BATTING_HR    0.0925289   0.0279140   3.315 0.000934 ***
## TEAM_BATTING_BB    0.0281209   0.0063392   4.436 9.67e-06 ***
## TEAM_PITCHING_H    0.0022894   0.0004144   5.525 3.73e-08 ***
## PSO                5.6998139   1.4861004   3.835 0.000129 ***
## HBP                1.7304979   1.2014381   1.440 0.149925
## DP                 7.8871356   2.8296544   2.787 0.005365 **
## CS                -0.5071417   0.9115006  -0.556 0.578013
## TEAM_BASERUN_SB    0.0271490   0.0069015   3.934 8.65e-05 ***
## TEAM_BASERUN_CS    0.0585306   0.0140037   4.180 3.05e-05 ***
## TEAM_BATTING_HBP  -0.0875536   0.0272910  -3.208 0.001357 **
## TEAM_PITCHING_SO   0.0002461   0.0015404   0.160 0.873090
## TEAM_FIELDING_DP  -0.1287003   0.0137193  -9.381 < 2e-16 ***
## TEAM_PITCHING_HR   0.0036741   0.0245541   0.150 0.881071
## TEAM_PITCHING_BB  -0.0041643   0.0045926  -0.907 0.364660
## TEAM_FIELDING_E   -0.0626614   0.0038189 -16.408 < 2e-16 ***
## SB                25.4011184   1.8919487  13.426 < 2e-16 ***
## TEAM_BATTING_SO   -0.0160581   0.0029869  -5.376 8.51e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 138.3495)
##
##      Null deviance: 483484  on 1999  degrees of freedom
## Residual deviance: 273794  on 1979  degrees of freedom
## AIC: 15558
##
## Number of Fisher Scoring iterations: 2
```

Trimmed feature model

```
summary(m2 <- glm(TARGET_WINS ~ . -HBP -DP -CS -TEAM_PITCHING_SO -TEAM_PITCHING_HR -TEAM_PITCHING_BB -B

##
## Call:
## glm(formula = TARGET_WINS ~ . - HBP - DP - CS - TEAM_PITCHING_SO -
##     TEAM_PITCHING_HR - TEAM_PITCHING_BB - BSO, family = "gaussian",
##     data = tr.tr)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -41.356  -7.916   0.193   7.943  47.484
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    36.4451998    5.7154746     6.377 2.25e-10 ***
## TEAM_BATTING_H     0.0462031    0.0035621    12.971 < 2e-16 ***
## TEAM_BATTING_2B   -0.0460110    0.0090645    -5.076 4.22e-07 ***
## TEAM_BATTING_3B    0.0601254    0.0162778     3.694 0.000227 ***
## TEAM_BATTING_HR    0.0956870    0.0096149     9.952 < 2e-16 ***
## TEAM_BATTING_BB    0.0239204    0.0035201     6.795 1.42e-11 ***
## TEAM_PITCHING_H    0.0019118    0.0003147     6.074 1.49e-09 ***
## PSO              5.3254663    1.4664578     3.632 0.000289 ***
## TEAM_BASERUN_SB    0.0407073    0.0042310     9.621 < 2e-16 ***
## TEAM_BASERUN_CS    0.0272926    0.0076021     3.590 0.000339 ***
## TEAM_BATTING_HBP  -0.0917886    0.0265629    -3.456 0.000561 ***
## TEAM_FIELDING_DP  -0.1303499    0.0133908    -9.734 < 2e-16 ***
## TEAM_FIELDING_E   -0.0561953    0.0030411   -18.478 < 2e-16 ***
## SB               24.2112119    1.7592288    13.762 < 2e-16 ***
## TEAM_BATTING_SO   -0.0164127    0.0023347    -7.030 2.83e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 138.6702)
##
##      Null deviance: 483484  on 1999  degrees of freedom
## Residual deviance: 275260  on 1985  degrees of freedom
## AIC: 15557
##
## Number of Fisher Scoring iterations: 2
```

Trimmed model with generated features

We add the three new relationships from the data exploration section as features to the trimmed model.

```
summary(m3 <- glm(TARGET_WINS ~ . -HBP -DP -CS -TEAM_PITCHING_SO -TEAM_PITCHING_HR -TEAM_PITCHING_BB -B

##
## Call:
## glm(formula = TARGET_WINS ~ . - HBP - DP - CS - TEAM_PITCHING_SO -
```

```

##      TEAM_PITCHING_HR - TEAM_PITCHING_BB - BSO + TEAM_BATTING_SO *
##      TEAM_BATTING_3B + TEAM_FIELDING_E * TEAM_FIELDING_DP + TEAM_BATTING_BB *
##      TEAM_FIELDING_E, family = "gaussian", data = tr.tr)
##
## Deviance Residuals:
##      Min        1Q      Median        3Q        Max
## -38.489   -7.757    0.156    7.777   50.643
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.562e+01  5.823e+00   4.399 1.14e-05 ***
## TEAM_BATTING_H    4.618e-02  3.649e-03  12.656 < 2e-16 ***
## TEAM_BATTING_2B   -4.741e-02  8.950e-03  -5.297 1.31e-07 ***
## TEAM_BATTING_3B    8.229e-02  3.124e-02   2.634 0.00851 **
## TEAM_BATTING_HR    8.036e-02  9.710e-03   8.277 2.30e-16 ***
## TEAM_BATTING_BB    4.100e-02  4.600e-03   8.913 < 2e-16 ***
## TEAM_PITCHING_H    9.248e-04  3.355e-04   2.756 0.00590 **
## PSO              5.073e+00  1.446e+00   3.509 0.00046 ***
## TEAM_BASERUN_SB    5.158e-02  4.634e-03  11.130 < 2e-16 ***
## TEAM_BASERUN_CS    1.872e-02  7.655e-03   2.446 0.01454 *
## TEAM_BATTING_HBP   -6.501e-02  2.640e-02  -2.462 0.01390 *
## TEAM_FIELDING_DP   -8.449e-02  1.948e-02  -4.337 1.52e-05 ***
## TEAM_FIELDING_E   -1.739e-02  6.780e-03  -2.564 0.01041 *
## SB                2.207e+01  1.831e+00  12.053 < 2e-16 ***
## TEAM_BATTING_SO   -1.727e-02  2.997e-03  -5.761 9.68e-09 ***
## TEAM_BATTING_3B:TEAM_BATTING_SO  2.722e-05  4.576e-05   0.595 0.55201
## TEAM_FIELDING_DP:TEAM_FIELDING_E -2.549e-04  6.739e-05  -3.782 0.00016 ***
## TEAM_BATTING_BB:TEAM_FIELDING_E -5.938e-05  1.106e-05  -5.370 8.78e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 134.6641)
##
##      Null deviance: 483484  on 1999  degrees of freedom
## Residual deviance: 266904  on 1982  degrees of freedom
## AIC: 15501
##
## Number of Fisher Scoring iterations: 2

```

SELECT MODELS

To select our models, we should make predictions on the internal evaluation set, where we know the target values.

Mean squared error

The mean of the square of the difference between real and predicted values.

It looks like the trimmed model with added features based on correlations has a slightly lower mean-squared error than the all-in and trimmed models.

Coefficient of determination

F-statistic

Residuals

EXPORT PREDICTIONS

CODE APPENDIX