

DATA 621—Business Analytics and Data Mining

Fall 2020—Group 2—Homework #1

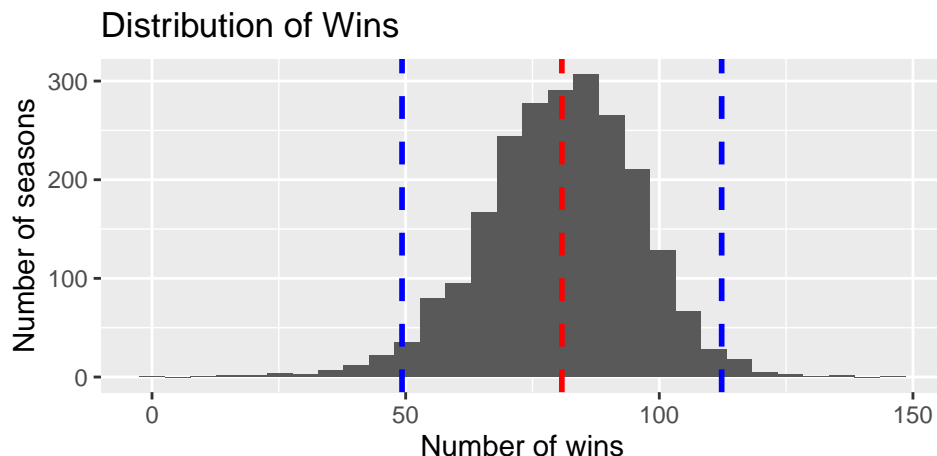
Avraham Adler, Samantha Deokinanan, Amber Ferger, John Kellogg, Bryan Persaud, Jeff Shamp

9/27/2020

DATA EXPLORATION

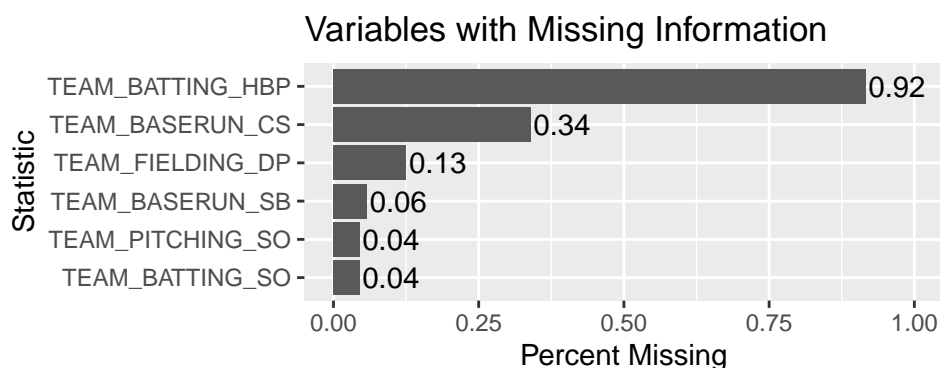
How often does the team win?

We are given a data set of 2,276 records containing 15 seasonal statistics and the total number of wins a team had in a given year. On average, about 50% of games played are won (81 games out of 162), with the best individual season having 146 wins and the worst season having 0 wins. The data is normally distributed and most years have between 49 and 112 wins (blue lines below). The nature of the distribution means there aren't too many extreme seasons where wins are significantly higher or lower than usual. This serves as a good gut-check for our final predictions; if the predicted wins are too high or too low, we know something in our model is probably off.



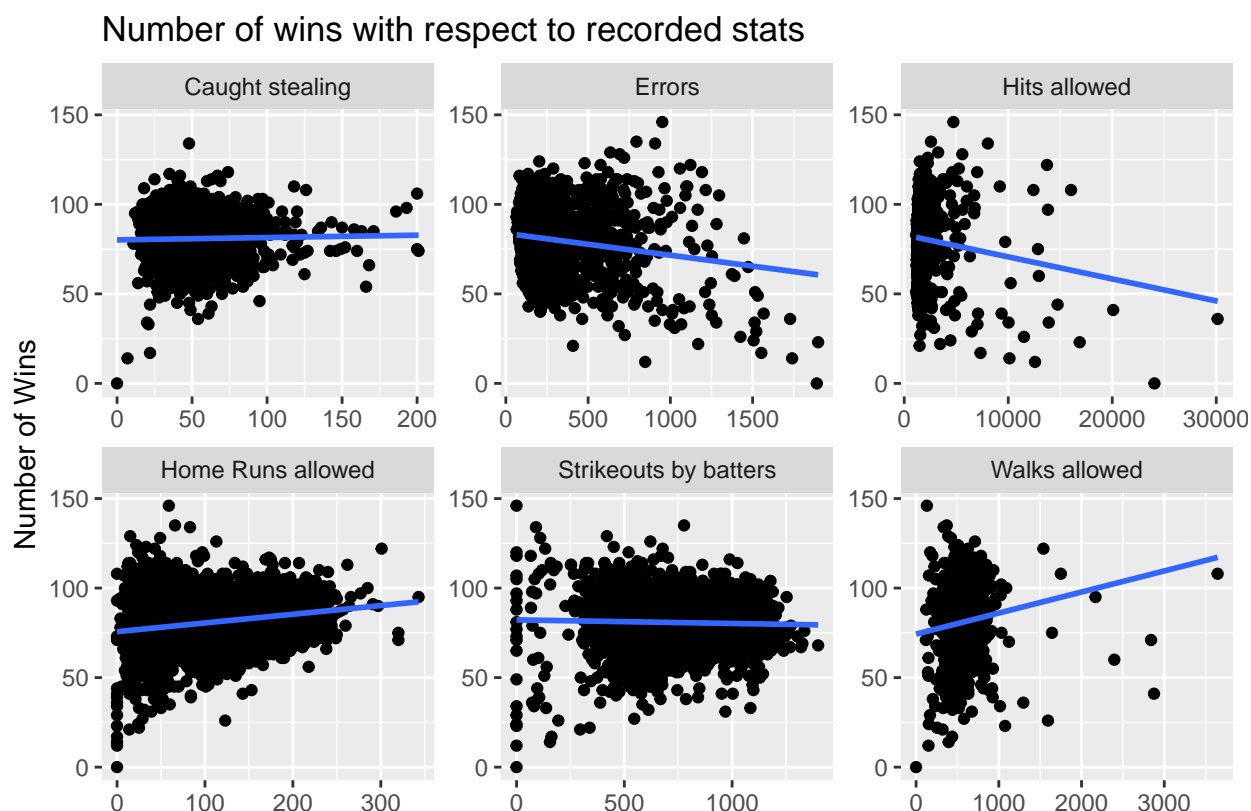
What's missing?

A first look at the data shows that only about 8% of the records have a full set of information. The good news is that most of the missing values come from statistics that don't happen too often: hit-by-pitch (`TEAM_BATTING_HBP`, 92% missing!), caught stealing (`TEAM_BASERUN_CS`, 34% missing), and double plays (`TEAM_FIELDING_DP`, 13% missing). Since we have so little hit-by-pitch data, we expect that it doesn't contribute much to overall wins and will eliminate it from a few of the models we propose. The other two stats have less than half of the data missing, so we'll need to think of a clever way to fill in these values. The remaining missing information is from a combination of stolen bases and strikeouts by pitchers and batters (`TEAM_BASERUN_SB`, `TEAM_PITCHING_SO`, `TEAM_BATTING_SO`). **It seems completely unreasonable** to have zero strike outs in a season, so this is something we'll most certainly have to impute.



Do the individual stats affect winning?

Stats with an expected negative impact: Intuitively, we expect that Caught stealing, Errors, Hits allowed, Home Runs allowed, Strikeouts by batters, and Walks allowed would all have a **negative** impact on the total wins. In other words, as these values increase, we expect that the team is less likely to win.

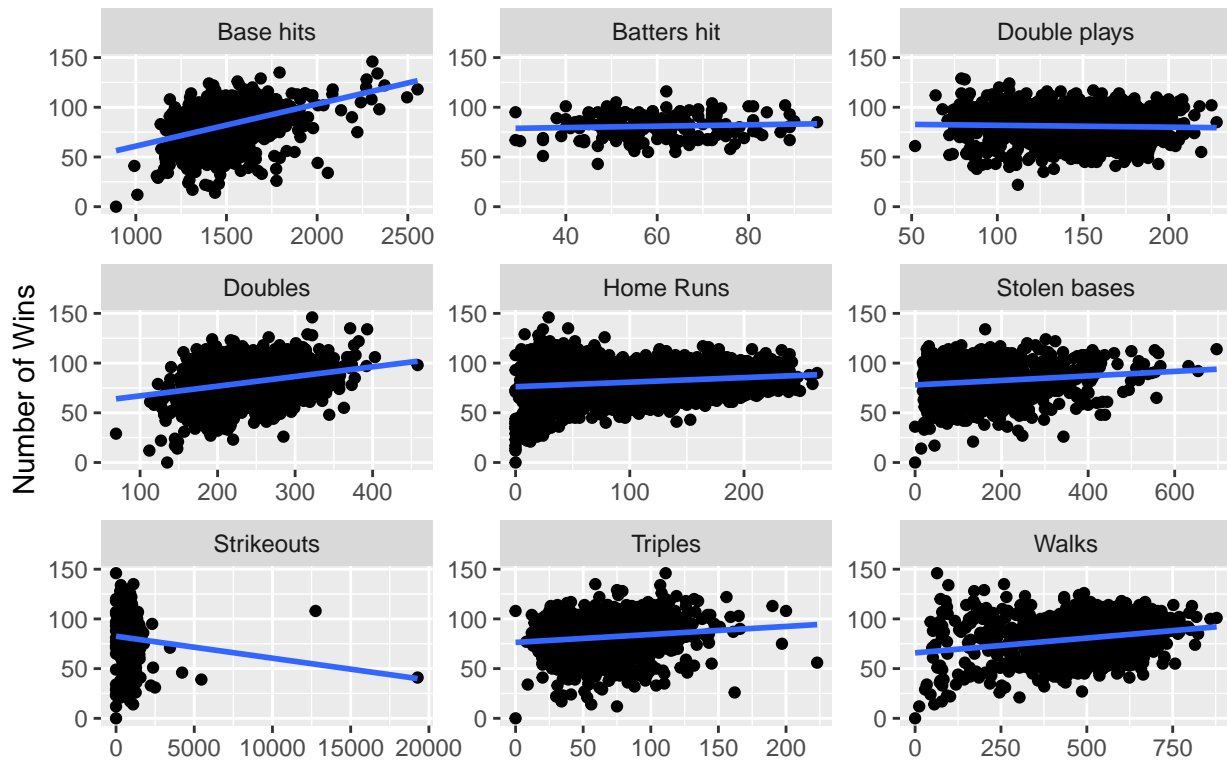


When we take a closer look at the data, these negative relationships aren't obvious. In fact, only **Errors** and **Hits allowed** seem to have a negative impact on wins. **Caught stealing** and **Strikeouts by batters** appear to be random; this means that whether the stat for a particular season is high or low doesn't affect the overall number of wins.

Even more interestingly, **Home Runs allowed** and **Walks allowed** have the *opposite* effect; as these stats increase, so do the number of wins!

Stats with an expected positive impact: We can look at the same information for the stats that we expect to have a **positive** effect on wins: Base hits, Doubles, Triples, Home Runs, Walks, Batters getting hit by pitches, Stolen bases, Double Plays, and Strikeouts by pitchers.

Number of Wins with Respect to Recorded Stats

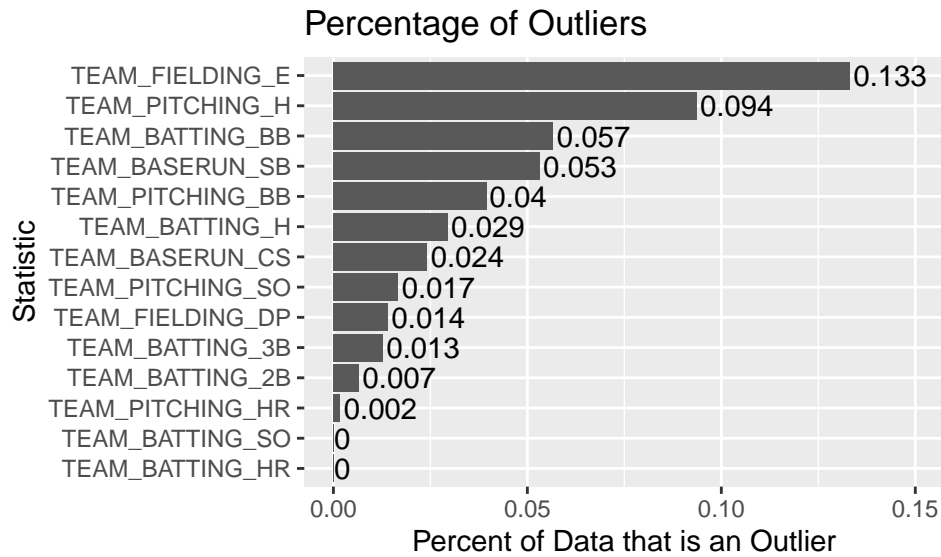


Many of these stats *do* seem to have an effect on the number of wins, most notably, **Base hits** and **Walks**. We see weaker positive relationships for **Home runs**, **Doubles**, **Triples**, and **Stolen bases**. This makes sense when we think about it; these things tend to happen less often in games than pure base hits and walks, so they don't have as much of an effect on winning. Finally, **Double plays** and **Batters hit** don't appear to have any correlation with the number of wins. Once again, this intuitively makes sense because they are less likely to happen in a game.

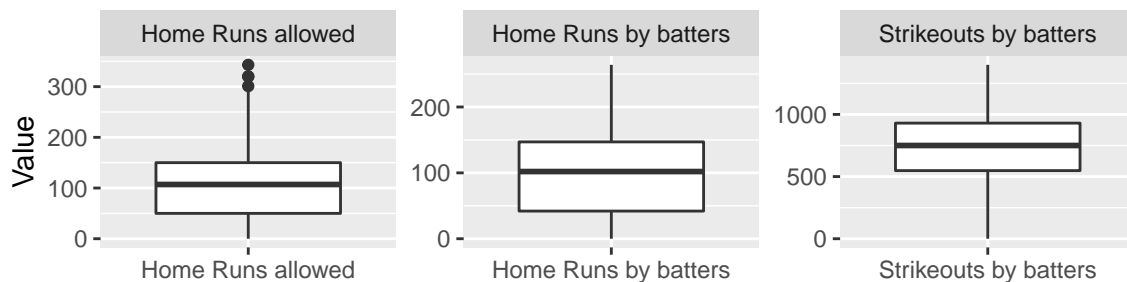
One thing to note is the number of strikeouts compared to the number of wins. We can see that there are a few outliers (abnormally high numbers of strikeouts in a season). This should be taken with caution, as they don't represent a typical season's stats.

Are some stats more skewed than others?

Before using any of the statistics in a model, we need to take a closer look at the variation in the data. We call out-of-the-ordinary values (exceptionally high or low values) **outliers**. We need to take these into account in our modeling because we want to make sure our predictions aren't skewed because of them.



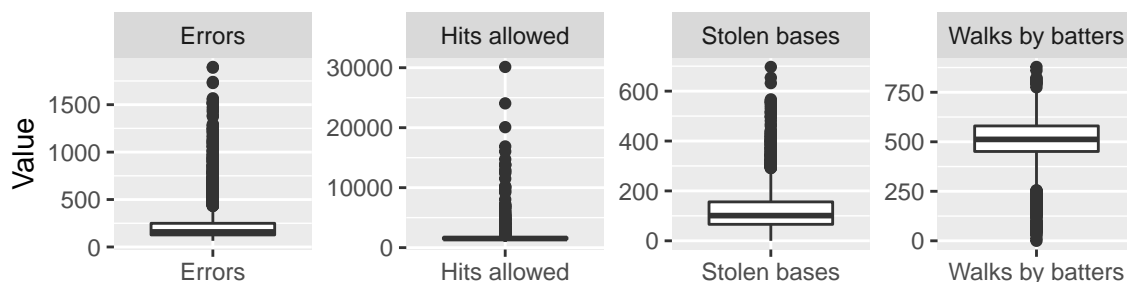
As we can see, some of the provided statistics are well-balanced in the sense that there are very *few* (or no) extreme values. **Home Runs allowed** (TEAM_PITCHING_HR), **Strikeouts by batters** (TEAM_BATTING_SO), and **Home Runs by batters** (TEAM_BATTING_HR) are examples of this.



Some things to note about each of these statistics:

- **Home Runs allowed** (average = ~100/year) and **Home Runs by batters** (average = ~106/year) have a very similar mid-range distribution (50% of the data lies between ~50 and 150). The slight difference in average stats means that teams tend to have a higher number of Home Runs than the opposition team.
- The only thing that stands out about **Strikeouts by batters** (average = ~736/year) is how nearly perfectly normal it is. 50% of the data is between about 500 and 1000 and there are absolutely no outliers in the dataset! This means that there were no surprisingly high or low seasons.

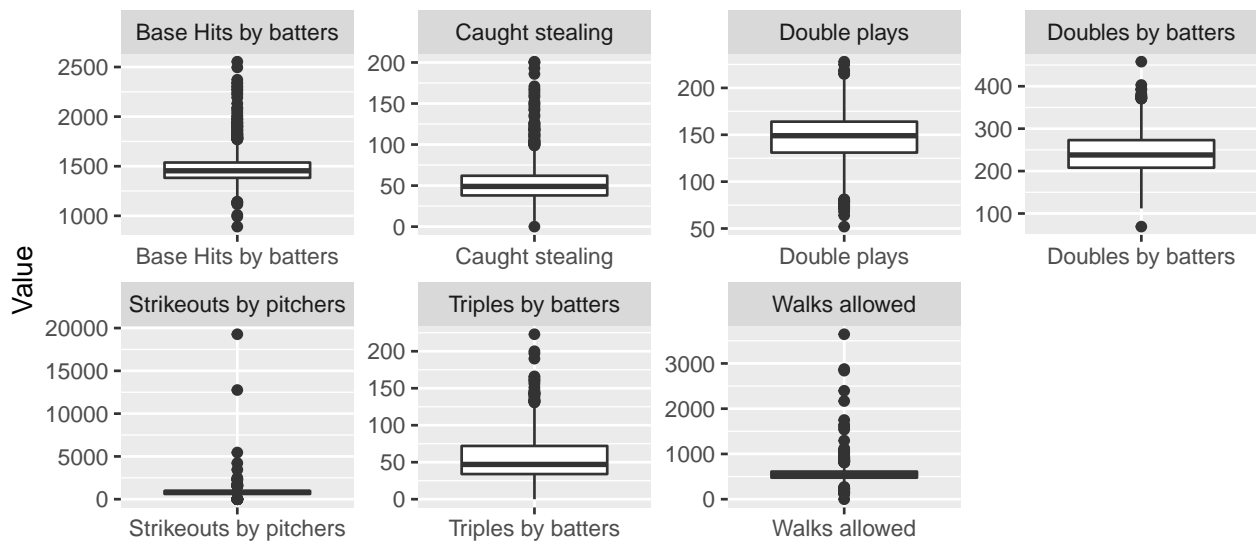
Conversely, some of the stats have a very *high* number of outliers, indicating that there are some seasons with some abnormally high or low values. **Errors** (TEAM_FIELDING_E), **Hits allowed** (TEAM_PITCHING_H), **Walks by batters** (TEAM_BATTING_BB), and **Stolen bases** (TEAM_BASERUN_SB) are examples of this.



Some things to note about each of these statistics:

- All of the outliers for **Errors**, **Hits allowed**, and **Stolen bases** are above the upper tail of the data set. This is further illustrated by the mean and median values for these stats; in all instances, the mean per year (Errors = ~246/year, Hits allowed = ~1779/year, Stolen bases = ~125/year) are higher than the median per year (Errors = ~159/year, Hits allowed = ~1518/year, Stolen bases = ~101/year). This means that some seasons with exceptionally high values skew the dataset.
- There are a few *very* extreme outliers for **Hits allowed**. The maximum value is 30,132, which is over 16 times the average number of hits allowed per season!
- There are outliers both above *and* below the tails of the data for the **Walks by batters** stat. This means that we have exceptionally low (min = 0!) and exceptionally high (max = 878) seasons.

The remaining stats, **Walks allowed** (TEAM_PITCHING_BB), **Base Hits by batters** (TEAM_BATTING_H), **Caught stealing** (TEAM_BASERUN_CS), **Strikeouts by pitchers** (TEAM_PITCHING_SO), **Double plays** (TEAM_FIELDING_DP), **Triples by batters** (TEAM_BATTING_3B), and **Doubles by batters** (TEAM_BATTING_2B) have between 29 and 99 outliers.

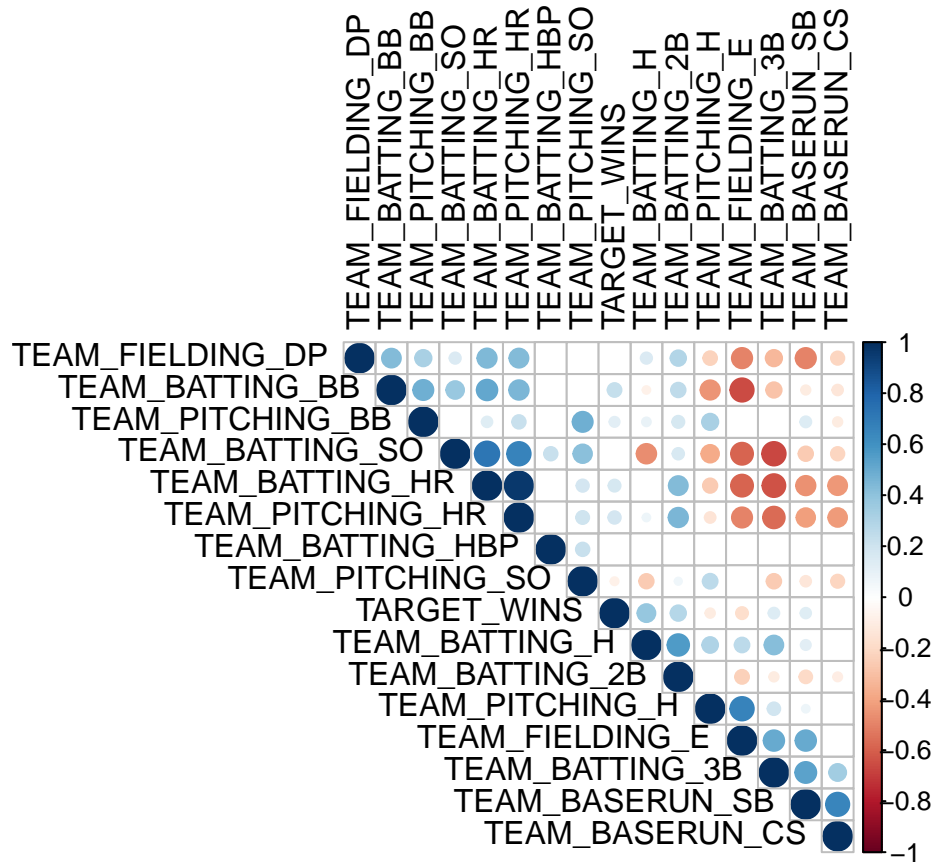


Some things to note:

- All variables are very narrowly distributed, meaning that most of the data falls within a small range.
- **Strikeouts by pitchers** and **Walks allowed** have a few very extreme outliers; these represent seasons that have abnormally high values for the statistics.
- The average number of pure **Base Hits** (1469/season) is greater than the average number of Doubles (~241/season) and Triples (~55/season). This isn't at all surprising, but serves as a good gut check on the validity of the data.

Are stats correlated?

We would expect that a few things in the dataset might be correlated: perhaps number of errors and hits/homeruns allowed or the number of base hits by batters and doubles/triples/homeruns. We can visualize the correlations between the statistics to determine if there is a significant relationship between them: blue dots represent positively correlated variables (as one increases, so does the other) and red dots represent negatively correlated variables (as one increases, the other decreases).



Some noteworthy relationships (coincidental or not):

- **Errors** are highly, negatively correlated with walks by batters, strikeouts by batters, and homeruns (both by batters and allowed).
- **Triples by batters** are highly, negatively correlated with strikeouts by batters and homeruns (both by batters and allowed).
- **Strikeouts by batters** are highly, positively correlated with homeruns (both by batters and allowed).
- **Homeruns by batters** and **Homeruns allowed** are both positively correlated with walks by batters.
- As expected, **basehits** are positively correlated with doubles and triples.

We can keep these correlations in mind when developing our models: if we have correlated statistics, there could be in-built redundancy in the features, and we may be able to create a simpler, more accurate model by eliminating some.

DATA PREPARATION

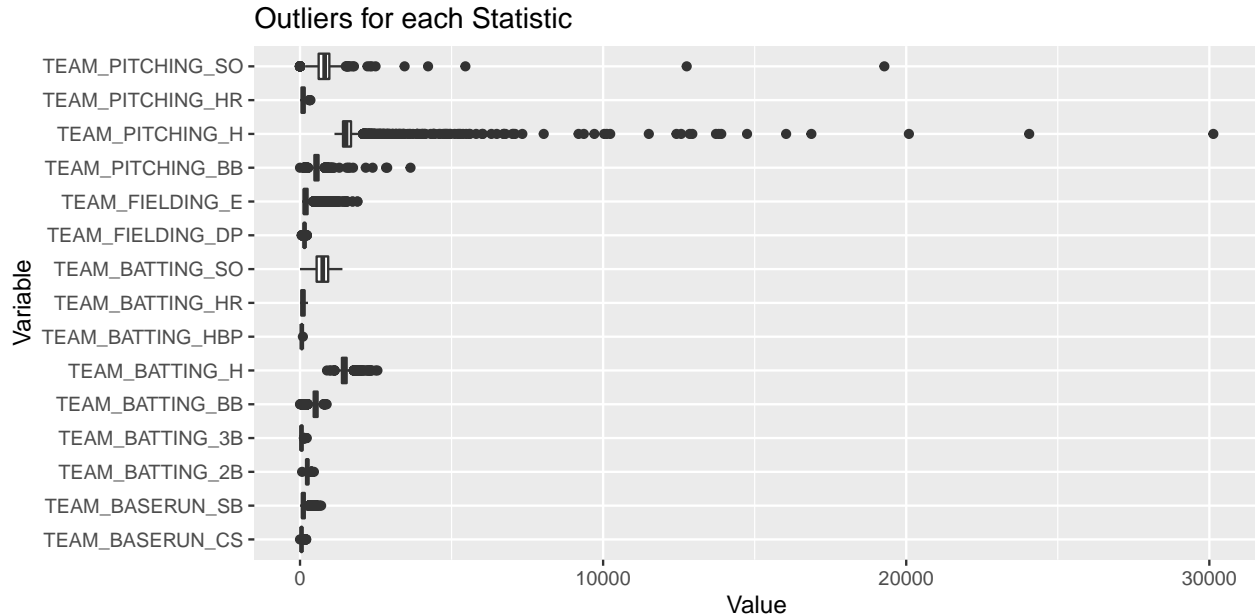
Now we have a good idea of the data we are looking at we can take the next steps to prepare it for building a solid model.

Outlier Removal

As we saw in the data analysis, there are some outlier concerns for some of the variables, so we will need to account for this in our modeling. When performing the processes of outlier removal, a cautious approach is always best. Each outlier is evaluated to ensure:

- It is clearly from incorrectly or mis-centered data.
- Its removal does not affect later assumptions.
- It creates a significant association/relation that is eliminated with its removal.

We can take another look at the outliers for each variable:

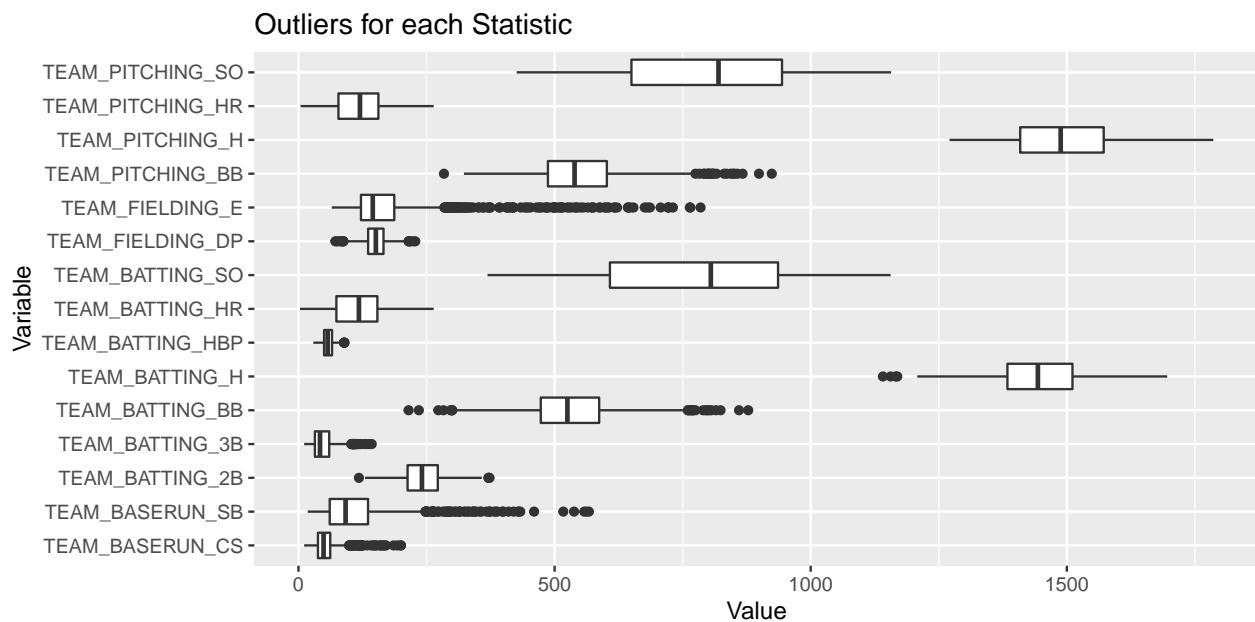


Off the bat, there are some clear issues with two of the variables: **TEAM_PITCHING_SO** and **TEAM_PITCHING_H**.

- **TEAM_PITCHING_SO**: The largest outlier is close to 20,000, which averages ~123 strikeouts per game. The MLB team pitching strikeout record is 1,450 and set in 2014 (CLE). There are a minimum of 27 at-bats per game and most teams average ~35 at-bats, thus making 123 strikeouts impossible.
- **TEAM_PITCHING_H**: The largest outlier for this variable is over 30,000 which is highly unlikely.

Since both appear with a heavy right skew, we will use median and IQR to remove the outliers.

We can take a look at the distributions of these variables after outlier removal:



This has eliminated what appears to be the most extreme outliers. There are still large outlier sets for Errors and stolen bases, but none that seem to dwarf the other variables. According to the box plot, things appear

to be on a scale that seems logical given the data and source of information that we have.

“MICE” imputation method

From our early exploration of the data, the vast majority of the data is complete with only a few variables with missing values. We will use the “mice” imputation method, specifically predictive mean matching method.

MICE (Multivariate imputation by chained equations) is a principled method of dealing with missing data. It creates multiple imputations, as opposed to single imputations and accounts for the statistical uncertainty in the imputations. It creates predictive values for the mean instead of imputing the IQR values.

Feature Engineering: Single base hits

With outliers removed and missing values imputed via MICE, we can create a few new variables. The first will be single base hits, and it will be derived from some of the variables we do have. We know that Team Batting Hits (`TEAM_BATTING_H`) is a combination of *all* hits for the season, so we can create Single Base Hits as follows:

$$Singles = Total\ Hits - (Doubles + Triples + Homeruns)$$

Feature Engineering: Slugging Percentage

The second variable we will create is **Slugging Percentage**. Slugging is an offensive statistic that is a good predictor of winning and it tends to have better variance and correlation than most other variables. It is composed of singles, doubles, triples, home runs and at-bats.

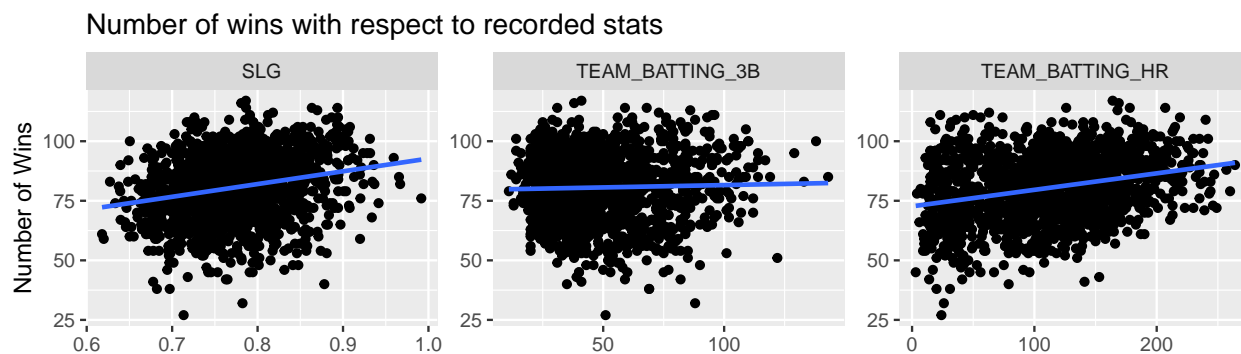
Slugging is calculated in the following way:

$$SLG = \frac{(1B) + (2 \times 2B) + (3 \times 3B) + (4 \times HR)}{AB}$$

Because we don’t have a statistic for at-bats, we will approximate it as follows:

$$\sim SLG = \frac{(1B) + (2 \times 2B) + (3 \times 3B) + (4 \times HR)}{SO_{batting} + H + BB}$$

As a gut-check, we can take a look at Slugging in comparison to Triples (`TEAM_BATTING_3B`) and Homeruns (`TEAM_BATTING_HR`). All of these variables should have a positive correlation with the total number of wins.



BUILD MODELS

Base Model

We will start with a simple linear model to serve as a baseline. This includes all variables in the dataset.

Model 1

Model 1 includes all variables except Singles, Triples, Base Hits, Strikeouts by pitchers, Walks allowed, and Batters Hit by Pitch. The final variables were chosen as a result of backwards dplyr::selection based on null hypothesis testing for non-zero slope.

Further data transformation:

We will use Cook's distance to remove outliers that are influencing the fit of the model above. We will use a cutoff of $\frac{4}{N}$.

Coefficient Discussion: First, we are keeping the coefficient for TEAM_BATTING_HR since it's p-value is marginally above the general threshold and knowledge of the game suggests it is important. There are some counter intuitive results, which is expected given that baseball is a messy, imprecise game that has evolved over time. To that end, we should expect some seemingly strange results from algorithmic regression. We used a combination of domain knowledge and data analysis to justify retaining features.

Fielding double plays and batting doubles both appear to have negative impacts on wins even though they *should* be positive impacts. Turning double plays, while a good for the defensive team, may suggest a larger, negative issue. Namely, weak pitching that leads to runners on base. Similarly, batting doubles, leaves runners open to double plays. Allowed hits by pitching, caught stealing, and batting strike-outs are all counter intuitive results as well, but they seem to be small contributors and these are events that happen regularly in every game.

Slugging as an approximation is the major predictor in this regression, by far. The other predictors other than the intercept are orders of magnitude less in predictive value. It should be noted that slugging alone is not a good predictor for wins overall.

Model 2

This is a model allowing pairwise interactions within the data types of baserunning, batting, pitching, and fielding, fit using a forward and backwards stepwise regression.

Coefficient Discussion:

Table 1: Step Model with Pairwise Class Interactions Output

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-93.708	43.050	-2.177	0.030
TEAM_BASERUN_CS	-0.032	0.018	-1.785	0.075
TEAM_BASERUN_SB	0.031	0.012	2.650	0.008
TEAM_BATTING_2B	0.091	0.144	0.627	0.531
TEAM_BATTING_3B	0.337	0.077	4.370	0.000
TEAM_BATTING_BB	0.043	0.117	0.367	0.714
TEAM_BATTING_H	0.152	0.052	2.903	0.004
TEAM_BATTING_HR	1.951	1.319	1.479	0.139
TEAM_BATTING_SO	-0.109	0.052	-2.085	0.037
TEAM_FIELDING_DP	0.000	0.034	-0.003	0.997
TEAM_FIELDING_E	-0.025	0.025	-0.994	0.321
TEAM_PITCHING_BB	0.130	0.096	1.359	0.174
TEAM_PITCHING_H	-0.048	0.039	-1.208	0.227
TEAM_PITCHING_HR	-2.054	1.302	-1.577	0.115
TEAM_PITCHING_SO	0.177	0.051	3.444	0.001

	Estimate	Std. Error	t value	Pr(> t)
TEAM_BASERUN_CS:TEAM_BASERUN_SB	0.000	0.000	4.447	0.000
TEAM_BATTING_2B:TEAM_BATTING_BB	0.000	0.000	2.426	0.015
TEAM_BATTING_2B:TEAM_BATTING_H	0.000	0.000	-1.421	0.156
TEAM_BATTING_2B:TEAM_BATTING_HR	0.000	0.000	-1.646	0.100
TEAM_BATTING_2B:TEAM_BATTING_SO	0.000	0.000	-1.658	0.097
TEAM_BATTING_3B:TEAM_BATTING_SO	0.000	0.000	-2.264	0.024
TEAM_BATTING_BB:TEAM_BATTING_H	0.000	0.000	-2.261	0.024
TEAM_BATTING_BB:TEAM_BATTING_HR	0.001	0.001	2.024	0.043
TEAM_BATTING_H:TEAM_BATTING_HR	-0.001	0.000	-1.527	0.127
TEAM_FIELDING_DP:TEAM_FIELDING_E	-0.001	0.000	-3.336	0.001
TEAM_PITCHING_BB:TEAM_PITCHING_HR	-0.001	0.001	-1.688	0.092
TEAM_PITCHING_BB:TEAM_PITCHING_SO	0.000	0.000	-2.852	0.004
TEAM_PITCHING_H:TEAM_PITCHING_HR	0.001	0.000	1.953	0.051
TEAM_PITCHING_HR:TEAM_PITCHING_SO	0.000	0.000	-2.825	0.005

Most of the coefficients are reasonable within the context of the game of baseball. These two statements are axiomatic:

- The only way to win is to have the higher score at the end of the game.
- The only way to score is for a baserunner to cross home plate.

With those in mind, we can make the following observations about the linear non-interactive coefficients.

- Reasonable
 - Caught stealing removes baserunners; negative coefficient makes sense.
 - Stolen bases get a runner closer to home plate; positive coefficient makes sense.
 - Triples get a runner very close to home plate; positive coefficient makes sense.
 - Walks are free baserunners; positive coefficient makes sense.
 - Hits increase the number of baserunners; positive coefficient makes sense.
 - Hitting home runs directly increase the score; positive coefficient makes sense.
 - Striking out reduces the number of baserunners; negative coefficient makes sense.
 - Errors allow the other team free baserunners; negative coefficient makes sense.
 - Giving up home runs gives the opponent scores; negative coefficient makes sense.
 - Getting strikeouts reduces the opponents baserunners; positive coefficient makes sense.
- Curious
 - Hitting doubles get runners on base; why negative?
 - Turning double plays reduce baserunners; why negative?
 - Giving up walks allows the other team free baserunners; negative coefficient makes sense but why insignificant?
 - Giving up hits allows the other team baserunners; why positive although insignificant?

An absolutely fascinating observation. In the first draft of this exercise, `TEAM_BATTING_HR` was removed due to its high correlation with `TEAM_PITCHING_HR`. Note the coefficients, batting is almost +12 and pitching is -9. However, this is such a strong indicator, that in the first run, pitching was given a *positive* coefficient of around +3. In hindsight this is because *it was being used as an indicator for **batting**!!* The correlation allowed the use of pitching HRs as an indicator for the hidden batting HRs! Once both were restored to the model, the logical coefficients surfaced. Another reason why models should not be trusted out of the box, but all model results should be reviewed for sanity and sense!!

The curiosities above can *possibly* be resolved by looking at the interaction terms. The interaction between giving up hits and strikeouts, `TEAM_PITCHING_H:TEAM_PITCHING_SO`, is highly negative and significant. It probably captures most of the giving up hits information making the singleton less relevant.

A possible explanation for the negative coefficient for getting double plays, is that double plays require at

least two people on base. That means that the opponent has a lot of base runners, which is very highly correlated with scoring.

The behavior of doubles remains confusing. Unless its hiding something like a team’s propensity to strand runners on base. It would be interesting to see a breakdown between the American and National leagues, as the latter tends to be somewhat better at “small ball” and moving the runners along.

Model 3

The final model is a Stepwise Regression with Repeated k-fold Cross-Validation, and higher order polynomials variables were introduced into the full model.

A stepwise variable selection model is conducted to determine what are the variables that can help predict the number of wins for the team. The method allows variables to be added one at a time to the model, as long as the F-statistic is below the specified α , in this case $\alpha = 0.05$. However, variables already in the model do not necessarily stay in. The steps evaluate all of the variables already included in the model and remove any variable that has an insignificant F-statistic. Only after this test ends, is the best model found, that is when none of the variables can be excluded and every variable included in the model is significant.

Here, the dependent variable is the continuous variable, **TARGET_WINS**, and the independent variables are the full model to identify the most contributing predictors. In addition, a robust method for estimating the accuracy of a model, the k-fold cross-validation method, was performed evaluate the model performance on different subset of the training data and then calculate the average prediction error rate.

Coefficient Discussion:

Table 2: K-fold Step Model with Higher Order Polynomials Output

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	129.936	13.477	9.641	0.000
TEAM_BATTING_2B	0.102	0.044	2.316	0.021
TEAM_BATTING_3B	0.227	0.041	5.560	0.000
TEAM_BATTING_HR	0.113	0.009	12.883	0.000
TEAM_BATTING_BB	-0.180	0.014	-12.594	0.000
TEAM_BATTING_SO	0.036	0.007	4.930	0.000
TEAM_BASERUN_SB	0.088	0.012	7.444	0.000
TEAM_BASERUN_CS	-0.067	0.035	-1.942	0.052
TEAM_FIELDING_E	-0.080	0.006	-13.861	0.000
TEAM_FIELDING_DP	-0.373	0.073	-5.087	0.000
TEAM_BATTING_1B	-0.056	0.018	-3.105	0.002
TEAM_BATTING_2B_1	0.000	0.000	-1.740	0.082
TEAM_BATTING_3B_1	-0.001	0.000	-2.613	0.009
TEAM_BATTING_BB_1	0.000	0.000	13.069	0.000
TEAM_BATTING_SO_1	0.000	0.000	-6.693	0.000
TEAM_BASERUN_SB_1	0.000	0.000	-2.981	0.003
TEAM_BASERUN_CS_1	0.000	0.000	2.232	0.026
TEAM_PITCHING_BB_1	0.000	0.000	4.819	0.000
TEAM_FIELDING_E_1	0.000	0.000	4.668	0.000
TEAM_FIELDING_DP_1	0.001	0.000	3.523	0.000
TEAM_BATTING_1B_1	0.000	0.000	6.063	0.000

Studying the coefficients of the model suggest that winning is in favor if the team batting hits more doubles, triples and home runs. Moreover, increase in the number of stolen bases, and a decrease in caught steals, double plays, error, and walks allowed would all lead to a win for the batting team. It is noteworthy that the model suggests that a decrease in single hits by batter and an increase in strikeouts by batters which seems

counter intuitive. But these variables were kept because when a batter steps to the plate, the player is more likely to strike out than to get a hit. Trying to hit the ball out of the park will come with strikeouts but it will also increase the chances of hitting home runs (even 1B, 2B, 3B), and that is pretty good exchange that most teams are willing carry out.

SELECTING A MODEL

In order to select the best model to make predictions, we looked at some measurements that tells us how well each model fits the training. These include the 1) R^2 , which represents the proportion of the variance explained by a model; 2) $adjR^2$, which is a modified version of R^2 ; 3) *Root Mean Squared Error* (RMSE), which is the square root of the mean squared error; and 4) *Akaike Information Criterion* (AIC), which is an estimator of out-of-sample prediction error.

When comparing the three models, it was interesting to select which would be our best model given that their performance statistics and error measurements were not significantly different from each other. It is apparent that Model #1: Backwards accounts for nearly 44% of the variation in the dependent variable with the independent variables, which is acceptable as a good model. But it's explanatory power based on an $adjR^2 = 0.414$ suggests it is no different from Model #3: K-fold. In addition, Model #1: Backwards has a highest RSME of all the models, which ranks as one of the major criteria we are using to decide on a model. Now, Model #2: Pairwise resulted in the $R^2 = 0.427$, $adjR^2 = 0.417$, and the lowest $RSME = 10.193$. Because the RMSE and adjusted R^2 statistics already include a minor adjustment for the number of coefficients estimated, to evaluate the model complexity, we compared the AIC. While a smaller AIC is deemed less complex, Model #1: Backwards would have been the way to go, however, it was deemed that the AIC differences among the models are quite insignificant. Thus, it is unanimous that Model #2: Pairwise is our final model.

Table 3: Performance Statistics & Error Measure of Models

Models	R.squared	adj.R.squared	AIC	RSME
Model #1: Backwards	0.436	0.429	11516.86	9.509
Model #2: Pairwise	0.427	0.417	12389.22	10.193
Model #3: K-fold	0.434	0.429	17231.34	11.756

PREDICTIONS

Using the test data and the selected final model, Model #2: Pairwise, a comparison in the prediction statistic was conducted.

Table 4: Prediction Comparison with Model #1: Backwards

dataset	n	mean	sd	median	trimmed	min	max	skew	kurtosis	se
Training Data	1648	80.57	13.47	82.00	80.98	27.00	117.00	-0.30	-0.06	0.33
Test Prediction	259	82.31	28.34	80.61	80.52	27.76	444.43	8.58	102.96	1.76

Table 5: Prediction Comparison with Model #2: Pairwise

dataset	n	mean	sd	median	trimmed	min	max	skew	kurtosis	se
Training Data	1648	80.57	13.47	82.00	80.98	27.00	117.00	-0.30	-0.06	0.33
Test Prediction	259	-93.71	0.67	-93.64	-93.65	-99.13	-92.71	-3.92	25.95	0.04

Table 6: Prediction Comparison with Model #3: K-fold

dataset	n	mean	sd	median	trimmed	min	max	skew	kurtosis	se
Training Data	1648	80.57	13.47	82.00	80.98	27.00	117.00	-0.30	-0.06	0.33
Test Prediction	259	79.87	11.90	80.46	80.23	25.52	121.74	-0.67	4.04	0.74

CONCLUSIONS

CODE APPENDIX

The code chunks below represent the R code called in order during the analysis. They are reproduced in the appendix for review and comment.