

# R Notebook

Alex Baker, James Kitch, Jackson Smith, Matthew Sheridan

## Introduction

Sports games are designed around the idea of finding the team that plays the “best” over a given period of time. Professional baseball teams have nine innings, soccer teams ninety minutes, and American football, hockey, and basketball clubs have sixty minutes to score more points than their opponents. While officiating, weather, and “luck” inevitably play a role in deciding the outcome, the end result of a game should hopefully be evident by in-game statistics. Bill James’ arguably invented the field of sports analytics in 1977 with his first edition of *Baseball Abstract*, attempting to analyze past player performance and predict future team success based on quantitative measurements rather than qualitative scouting reports. This quantitative view of player and team statistics, now known as “sabermetrics”, exploded in baseball where it is very easy to generate a large number of individual statistics for every player.

The National Football League (NFL) was slower to adopt sabermetrics, but the rise of remote sensing software in recent years has made it much easier for players and individuals to acquire quantitative game statistics. Football presents an interesting area of analysis since it is naturally discretized into distinct plays which can be measured quantitatively. How many net yards were gained? Was it a run or a rush play? Was it a missed field goal attempt, and if so from how many yards? Motivated to explore these questions, in this project we sought to explore the relationship between NFL in-game statistics and the final score differential.

R is a natural environment in which to analyze this as it contains the `nflfastR` package, which contains accumulated play-by-play data from 1999 through 2021, with additional predictors beginning in 2006. We hope to create a parsimonious, streamlined model capable of predicting the score differential of an NFL game. While we will have some variables with relatively clear-cut relationships with score differential, our analysis will also focus on how more “strategic” variables relate to success. More turnovers in a game and more total yards should correlate strongly with score differential, but we are also interested in variables that have a more unclear relationship with the final outcome. Are missed extra point attempts indicative of the team as a whole having a bad day? What about third down conversion rate? And quarterback hits? Is it better if a quarterback is able to spread his passes to different receivers, or when one receiver dominates the game? These are the kinds of interactions and trends we hope to expose in our model. By using metrics that are not generally used to predict game outcomes, we hope to find insight into predicting wins that are not conventionally expected.

In our project, we hope to answer the following questions:

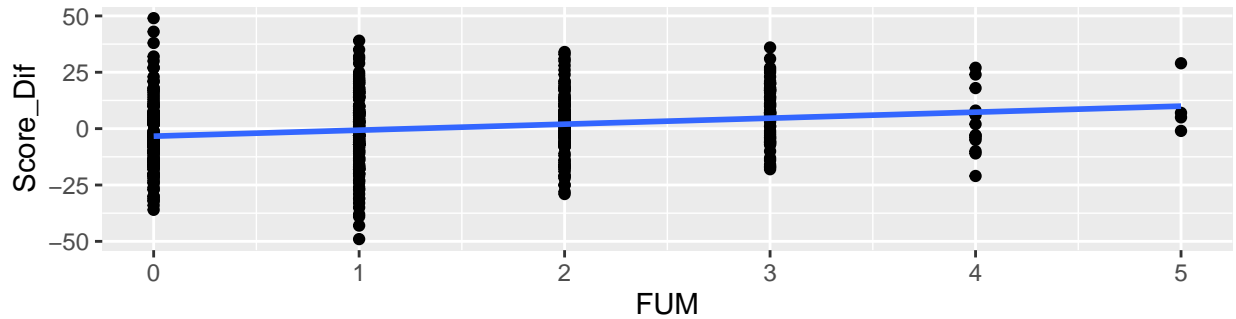
How do in-game statistics relate to the final score differential of a game? Which predictors are most significantly associated with score differential? What about predictors that aren’t classically associated with offense?

Do some predictors have a different impact as the game progresses? For instance, the first half might be more indicative of team strategy and the second half of the current score in the game.

## Data Cleaning, EDA

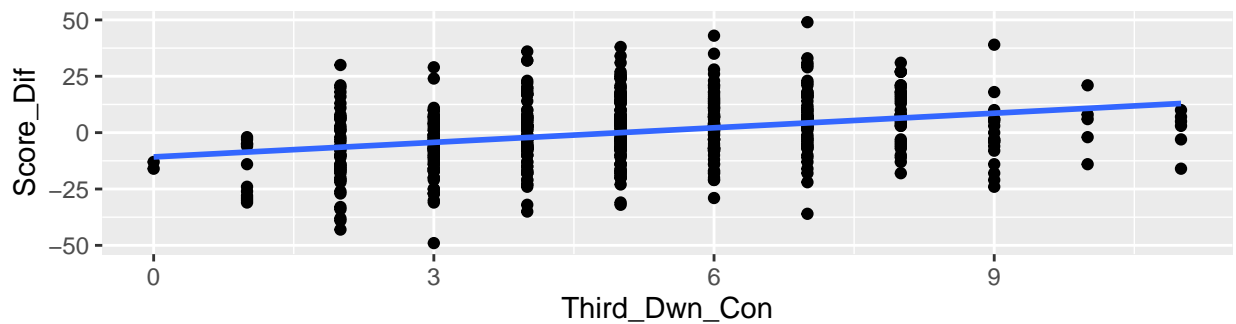
Number of Fumbles Caused to Score Difference

2019 NFL Games



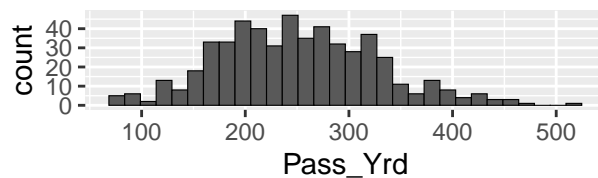
Third Down Conversions vs Score Difference

2019 NFL Games



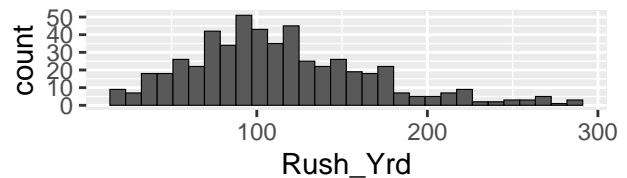
Total Passing Yards

2019 NFL Games



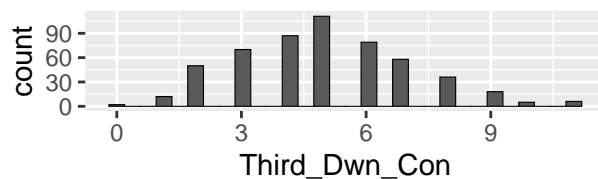
Total Rushing Yards

2019 NFL Games



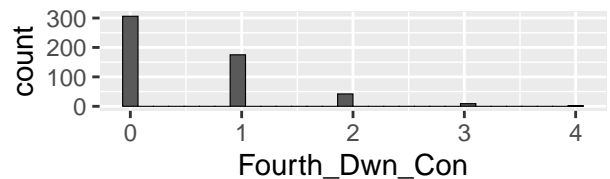
Third Down Conversions

2019 NFL Games



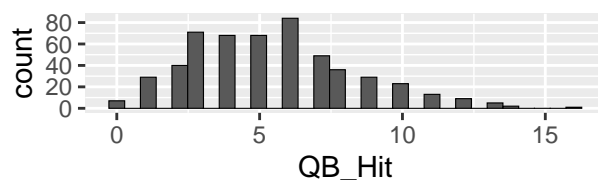
Fourth Down Conversions

2019 NFL Games



QB Hits

2019 NFL Games



The data originally obtained from nflfastR is very granular, with one row for every play in a season. We acquired the play-by-play data for 2017, 2018, 2019, and 2021, skipping 2020 because that was the “COVID season” where game dynamics may have been different without fans in the stadium. The play-by-play level was too granular for our purposes, as it would have led to a sparse design matrix due to the high variability of plays and the low frequency of certain game events such as quarterback hits, touchdowns, and turnovers. In order to explore the relationship between total in-game statistics and score differential, we aggregated the data by team and by game. For example, we would sum the passing yards over all of one team’s plays in a game to get their total passing yards in the game. Aggregating the plays to the game-level allowed us to create a richer data matrix that would provide more interpretable, useful results. We chose a wide range of predictors from the nflfastR dataset, encompassing predictors which we thought would directly relate to score differential such as touchdowns and interceptions as well as variables whose relationship was less intuitively clear such as unique receivers, quarterback hits, and missed extra point attempts. After aggregation and preliminary variable selection our dataset looks as follows:

##	Game_id	Team_Name	Home	TFL	QB_Hit	Pass_Yrd	Rush_Yrd	FG_attempt
## 1	2017_01_ARI_DET	DET	1	1	4	292	82	1
## 2	2017_01_ARI_DET	ARI	0	4	7	268	45	2
## 3	2017_01_ATL_CHI	CHI	1	2	10	213	125	1
## 4	2017_01_ATL_CHI	ATL	0	1	6	321	64	3
## 5	2017_01_BAL_CIN	CIN	1	5	6	170	77	0
## 6	2017_01_BAL_CIN	BAL	0	2	3	121	157	2
##	XP_missed	FG_missed	FUM	INT	Fourth_Dwn_Con	Third_Dwn_Con	Pen	RedZone_Plays
## 1	0	0	2	3	0	8	4	15
## 2	0	1	1	1	1	6	5	20
## 3	0	0	0	0	1	5	2	11
## 4	0	0	2	0	0	5	4	10
## 5	0	0	0	1	0	4	4	10
## 6	0	0	2	4	0	6	7	14
##	Unique_Receivers	Shotgun_Plays	Score_Dif	year				
## 1	7	50	12	2017				
## 2	8	43	-12	2017				
## 3	8	40	-6	2017				
## 4	8	20	6	2017				
## 5	9	42	-20	2017				
## 6	9	27	20	2017				

The included predictors are:

**Game\_id:** ID to distinguish games.

**Team\_Name:** What team are statistics being recorded for

**Home:** dummy variable, 1 if the given team was at home for the game and 0 if they were away.

**TFL:** the number of offensive plays from a team that are Tacked For a Loss (i.e. tackles behind the line of scrimmage) in a game.

**QB\_hit:** the number of hits on a given team’s quarterback in a game.

**Pass\_Yrd:** Total game passing yards from a given team.

**Rush\_Yrd:** Total game rushing yards from a given team.

**FG\_Attempt:** Total field goal attempts in a game by a given team.

**XP\_missed:** Total missed extra-point attempts in a game

**FG\_missed:** Total missed field-goal attempts in a game

**FUM:** Total team fumbles.

**INT:** Total team interceptions.

**Fourth\_Dwn\_Con:** Total fourth-down conversions.

**Third\_Dwn\_Con:** Total third-down conversions

**Pen:** Total penalties committed by a given team

**RedZone\_Plays:** Total number of plays run within the 20 yard line for a given team.

**Unique\_Receivers:** Number of different receivers that caught a ball for a team in a game.

**Shotgun\_Plays:** Number of plays run out of the shotgun formation for a team.

**Half1\_RunPass\_Ratio:** Runs/Passes in the first half for a team.

**Third\_and\_Longs:** Third downs converted from over 5 yards away for a team.

**Half1\_AirYrds:** Total air yards (ball travels in air to meet receiver) in the first half.

**Score\_Dif:** Score Differential (positive means win negative means lose).

Below, we plotted the distributions for a number of different potential predictors, as well as their relationship with the response variable of interest. Most predictors are slightly right-skewed, as they have a positive support and have a slight “bell” shape for the most common values but are stretched out by the few exceptional games where a team throws for 500 yards or rushes for 250 yards. However, it is mild right-skewness as it is not to the extent that log-transforming them would make the distributions more symmetric. For example, if we log-transform `Pass_Yrd` we actually find that the resulting distribution is left-skewed, suggesting that the log transformation was actually too strong! While we found that a square root transformation would make the `Pass_Yrd` and `Rush_Yrd` distributions more symmetric, it would also make interpretations much less intuitive. Since the primary goal of this project is interpretation-focused and the relationship between offensive yards and score differential is still relatively linear (and the distributions of predictors still relatively normal despite very slight skew), we elected to not do any transformations here.

Furthermore, we can do a preliminary investigation of the relationships between some of the predictors and the outcome variable of score differential in 2019. Notably, we can see that the number of fumbles caused by a team is positively associated with score differential and that it is a relatively linear trend. The number of third down conversions is also positively associated with score differential, but this seems more complicated than a “simple” linear trend. Teams with a very large number ( $> 8$ ) of third down conversions in a game appear to have a more negative score differential than teams with 5-7 third down conversions. This suggests that in future models investigating the effects of quadratic or higher order polynomial terms could be reasonable.

## Initial Modeling:

Our baseline model includes all available predictors to get a sense for the relationships between variables and to final score difference. It is very unlikely to be our best-performing final model, but will be a good starting point for our analysis.

As discussed previously in our EDA, the assumptions of linear regression do generally hold for our data. The baseline model includes predicting score differential from the main effects of all available predictors. The distribution of each of the individual predictors is found to be relatively unskewed (without need for transformations). There also is a shown linear relationship between our predictor variables and score differential. There is no clear fanning (constant variance holds). Lastly, there is likely some dependence among the data (i.e., a team that has many quarterback hits is more likely to have more interceptions), but these should not greatly affect our models.

## Analysis

### Stepwise Models

We also ran stepwise variable selection backwards and in both directions, per the basis of our modeling project. We have seen firsthand how often forward selection leaves out significant predictors, and since one of our primary aims is to find less-well-known predictors, we believed the forward selection method to be too risky since it messes with the fundamental goal of our project. Forward selection also handicaps the AIC analysis, giving us further reason to avoid it for better inferential, predictive, and reasonable models. Our preliminary stepwise variable selection with simple predictors from a baseline linear model yielded a linear model with the following formula:

```
## Score_Dif ~ Home + QB_Hit + Pass_Yrd + Rush_Yrd + FG_attempt +
##      XP_missed + FG_missed + FUM + INT + Fourth_Dwn_Con + Third_Dwn_Con +
##      RedZone_Plays + Unique_Receivers + Shotgun_Plays
```

We can note that this immediately removes TLF and Pen; additionally, this model has a Multiple R-squared = 0.5439, an F-statistic of 135.2 and a p-value < 0.00000000000000022. In the coefficients table in the appendix, we can also see that XP\_missed has a p-value > 0.05, indicating overfitting, which does not bode well for the more complex stepwise model.

With the both directional stepwise variable selection model starting from the baseline linear model all the way through the model with all interaction effects, we find the following formula:

```
## Score_Dif ~ Home + QB_Hit + Pass_Yrd + Rush_Yrd + FG_attempt +
##      XP_missed + FG_missed + FUM + INT + Fourth_Dwn_Con + Third_Dwn_Con +
##      RedZone_Plays + Unique_Receivers + Shotgun_Plays + FG_attempt:RedZone_Plays +
##      Pass_Yrd:RedZone_Plays + QB_Hit:RedZone_Plays + Fourth_Dwn_Con:Shotgun_Plays +
##      Pass_Yrd:Fourth_Dwn_Con + Fourth_Dwn_Con:RedZone_Plays +
##      Home:Fourth_Dwn_Con + Pass_Yrd:Unique_Receivers + QB_Hit:Pass_Yrd +
##      XP_missed:RedZone_Plays + QB_Hit:FG_attempt + QB_Hit:FUM +
##      INT:Shotgun_Plays + FG_attempt:Third_Dwn_Con
```

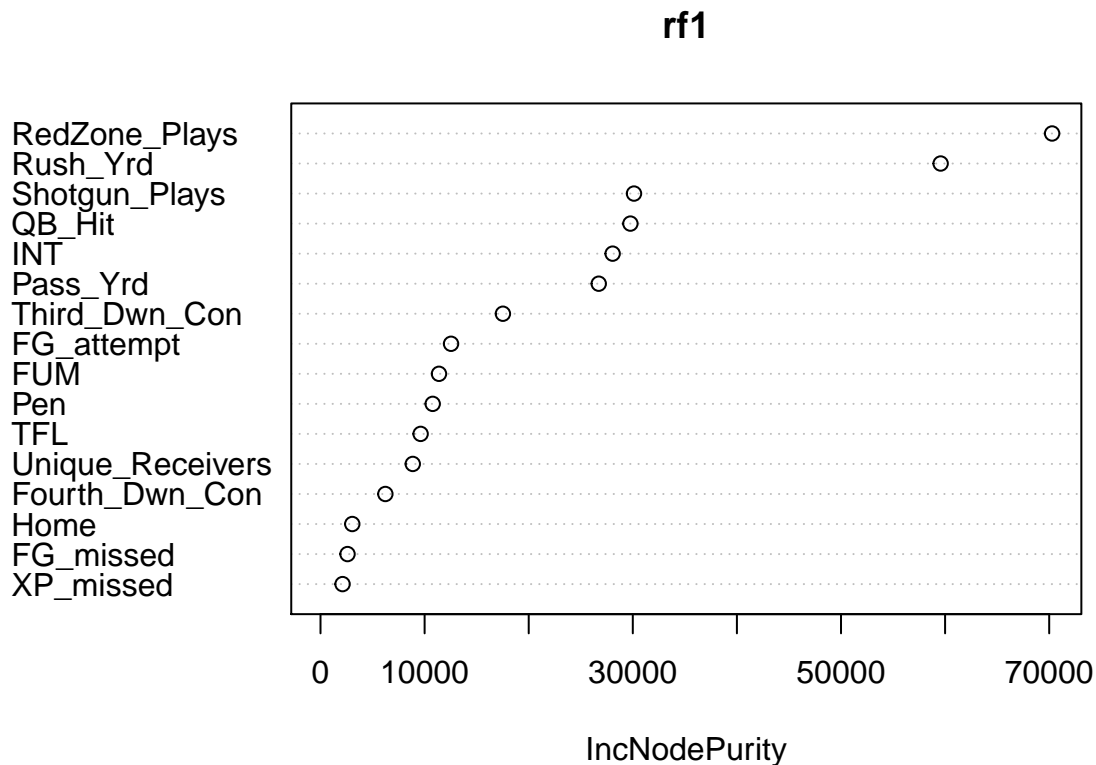
However, although the model appropriately disposes of countless insignificant predictors, it still includes several that are insignificant by their p-value, namely: Home, Pass\_Yrd, Fourth\_Dwn\_Con, Pass\_Yrd:Unique\_Receivers, QB\_Hit:Pass\_Yrd, XP\_missed:RedZone\_Plays, QB\_Hit:FG\_attempt, QB\_Hit:FUM, INT:Shotgun\_Plays, and FG\_attempt:Third\_Dwn\_Con, most of which are interaction terms that correspond to specific in-game situations or patterns. For example, while a home field advantage helps, the reasons for the advantage are often beyond numerical quantities such as being used to the weather or loud fans. Passing yards are obviously correlated with the number of unique receivers or QB hits since QB's cannot throw as much when hit more often/with less time in the pocket, and a smaller number of targetable receivers limits the number of potential yards since the team is more susceptible to defensive strategies. And of course the number of third down conversions inversely relates to the number of field goal attempts. That relates to simple football rules and strategies.

### Random Forest

A random forest was fit using the combined game data from years 2017-2019. One advantage to using a random forest in this case is that it represents a good way to measure the predictive power of the predictors chosen, due to the non-parametric nature of the random forest and lack of assumptions needed to fit the model. After fitting a random forest, we can plot the Variable Importance of the predictors, to see how much predictive power is lost when the predictor is removed from the tree, and the greater the disparity the more important said predictor is. The one disadvantage to this method is that while we are able to tell that a certain predictor is very important, we are unable to interpret whether the relationship between the predictor and the response is positive or negative. For example, we see that QB\_hits (the number of hits a team's quarterback takes) is considered important. We may assume, as logic follows, that minimizing the amount of times your quarterback is hit is most beneficial for a large score differential. However, we cannot determine that this is the case from only this plot.

A random forest was fit over a range of different hyper-parameters. Different values of mtry, ranging from 1 (essentially a random forest) to 12, reflect the number of variables randomly sampled at each split within the random forest, where higher values add more complexity. Secondly is maxnodes, which represents the maximum number of terminal nodes present in the resulting forest. The number of trees within our forest was fixed at 200 to maintain a reasonable run time. After cross-validating the RMSE between predicted and true values among training (2017-2019 data) and testing (2021 data) sets, the lowest test set RMSE was selected, and it was found that the optimal set of hyper parameters from our tests is mtry = 8 and maxnodes = 1000. Looking at the calculated differences in RMSEs on training and testing data, no models appear to be egregiously over-fit. Our selected hyper-parameters do not produce one of the largest magnitude RMSE differences, so it appears that our tuned model is acceptable.

Using the aforementioned hyper-parameters, the model was fit again and a variable importance plot was generated. This plot shows which variables have the most predictive power in the random forest. We can see that the highest predictors in terms of variable importance are Red Zone Plays and rushing yards, which appear to be in a tier of their own. This is followed by another tier that includes shotgun plays, QB hits, interceptions and passing yards. The rest of the predictors are of much lower importance based on this plot. This suggests that these mentioned predictors have the most important relationships with score differential for any given game.



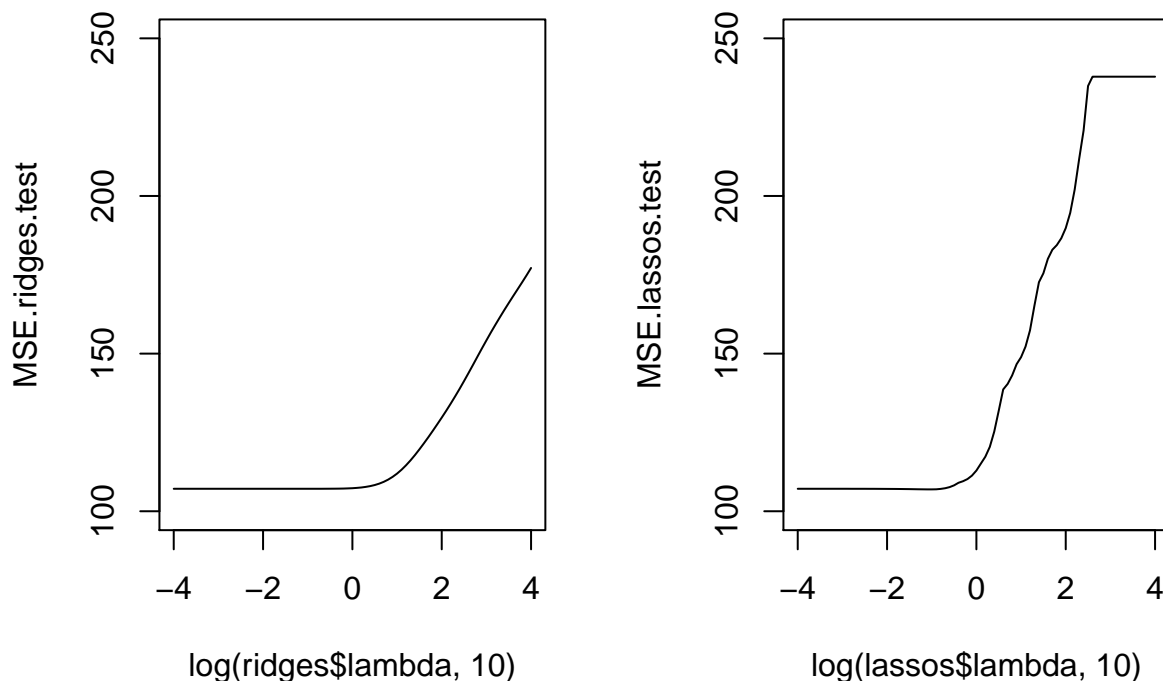
In the broader context of football, it makes sense that red zone plays would be considered highly important for predicting score differential because, considering that the red zone is the area within 20 yards of the opponent's end zone, this is where points are scored, and the more plays that are run in this area the more points should be scored. It also makes logical sense that rushing yards should be almost as important, both because teams that are leading by large amounts tend to run much more in order to possess the ball, and conversely teams that have great success running tend to have more control of the game and time of possession. In terms of the "second tier" predictors, QB\_hits makes sense as a solid predictor. Presumably, the fewer hits your quarterback is taking, then the better chance you have of putting up more points than your opponent. The number of interceptions a team gets gives more opportunities to score while having to drive fewer yards down the field than if the ball was obtained otherwise (such as a punt, or a kickoff). The number of passing yards and shotgun plays are slightly less interpretable in terms of their importance. Passing yards could be a positive predictor because a team that amasses lots of passing yardage likely scores many points. However, a team that is behind by a large amount may gain many passing yards towards the end of the game in a meaningless effort to catch up to the winning team (this is referred to as "garbage time" statistics, when the outcome of the game is relatively determined, but play still occurs and players amass stats). Similarly, the number of stats taken out of the shotgun could be the mark of a team being successfully aggressive in their passing game, or that a team is struggling behind and trying to play catch-up. It would be best to investigate the nature of the relationship of these predictors with our chosen response in some sort of other linear model in order to examine the sign and size of the associated coefficient.

### Ridge Regression and LASSO

The next modeling methods we chose to use were Ridge and Lasso regressions. Both of these models shrink our  $\beta_j$  coefficients towards 0 in an effort to expand on the ordinary least squares regression's unbiased estimates by reducing the Type I error rate through a penalizing  $\lambda$ . It is critical to note that the  $\beta_0$  coefficient is unaffected by either penalized loss function, which means the intercept could theoretically remain consistent if the model weren't to change.

First, we fit a Ridge regression model, which uses a squared  $\lambda$  penalty factor, then the Lasso regression with its absolute value  $\lambda$  followed. The Lasso typically has an advantage over the Ridge in that it punishes the  $\beta_j$  coefficients more harshly, pushing the insignificant and unimportant predictors all the way to 0 if they are truly irrelevant.

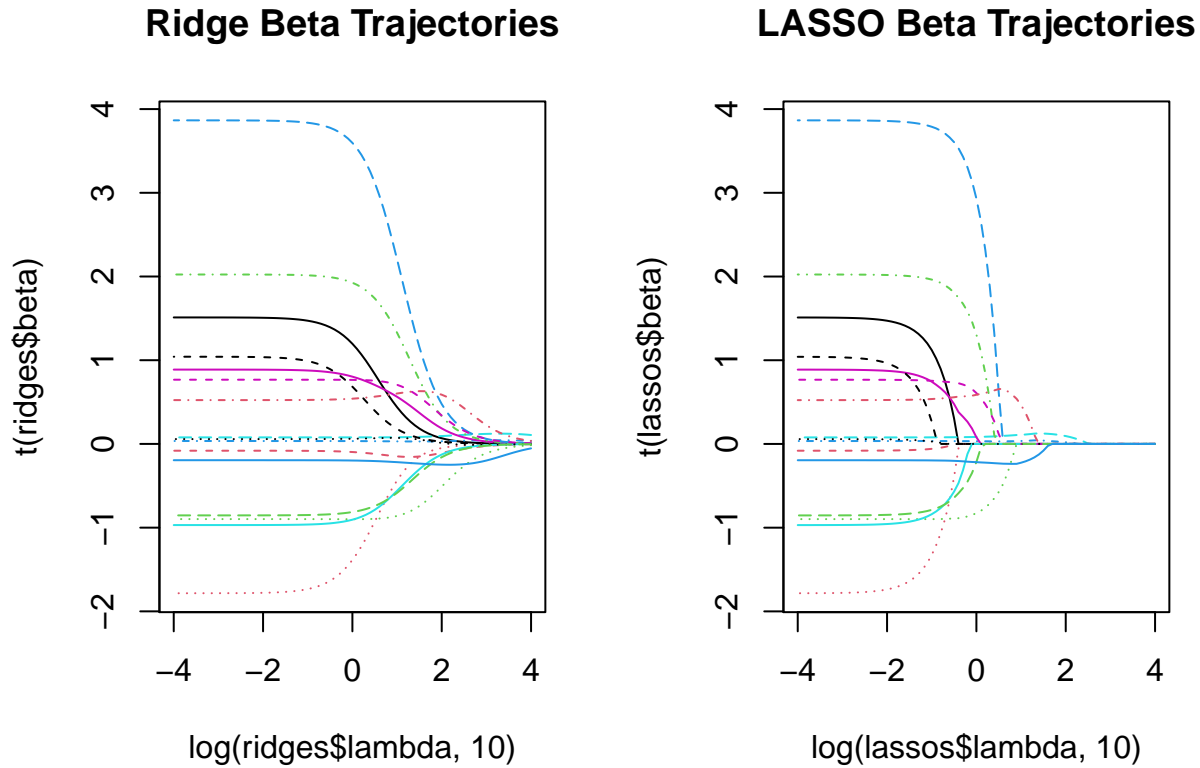
```
par(mfrow=c(1,2))
plot(MSE.ridges.test~log(ridges$lambda,10),type="l", ylim = c(100, 250))
plot(MSE.lassos.test~log(lassos$lambda,10),type="l", ylim = c(100, 250))
```



The MSE in the test plot demonstrates that this shrinkage, in both the Ridge and Lasso regressions, greatly improves the model estimates. Accordingly, the Ridge's best  $\lambda$  turns out to be  $\log_{10}(\lambda) \approx 0.5$  and  $\lambda = 3.162278$ , which matches the graph's minimum MSE well, from a visual inspection. The corresponding MSE for the best  $\lambda$ 's Ridge is 109.7959. In a similar plot for the Lasso regression, the optimal  $\lambda$  came out as  $\log_{10}(\lambda) = -1$  so  $\lambda = 0.1$ . The Lasso also performed better than the Ridge regression, as expected, with an MSE of 107.4165, which is verifiably less than the Ridge's MSE ( $107.4165 < 109.7959$ ). That means that the Lasso regression model is  $\sim 2.2\%$  better than the Ridge regression model in terms of MSE, and we logically know that the  $\lambda$ 's cannot be compared due to the difference in calculation methods. This makes sense since our data does appear to have several unimportant predictors mixed in with a few significant ones, as is common in the real world beyond pure statistical problems. It is also crucial to note that the Lasso MSE tops out at 240 while the Ridge MSE does not appear to stop growing beyond our chosen  $\lambda$  bounds. For the Lasso, we have already found our optimal  $\lambda = 0.01$ , but the MSE boundary shows that MSE does not increase when  $\lambda > \sim 2.3$ .

Looking at the Ridge model's  $\beta_j$  trajectory plots, we can see wide variability for many of the predictors, with some shrinking to 0 incredibly quickly while others more gradually or steeply further along the  $\log(\lambda)$ 's).

In the Lasso's version of the test  $\beta_j$  trajectory plots, the trajectories are much more enticing. all of the predictors are significant, until they aren't, dropping off sharply after specific lambda values, and almost all of them have shrunk to 0 by the time  $\lambda = 0.1$  or shortly thereafter.



We can be thankful for this result, albeit expected, since the Lasso regression's ability to shrink coefficients to 0 often makes it more conservative and more interpretable, leading to an easy and efficient modeling choice.

## LMER Modeling

Next, we fit a Linear Mixed-Effects model with random intercepts by team to see if that had any impact on our predictions. Given that observations can be grouped by team and there are 32 teams, a mixed-effects model was a natural next step in our modeling approach. We fit two different types of mixed-effects models: one with all available predictors, and one with the predictors chosen by sequential variable selection. We found that the inclusion of a random intercept did not improve prediction substantially, with the random intercept term only accounting for  $\approx 9\%$  of the total variance in the model. Most coefficient estimates were very similar to the standard OLS run previously. Some of this may be due to the fact that there are a relatively large number of observations per team. Although there are 32 teams, each team has 48 - 60 "measurements" or games in the dataset and thus the benefit of linear mixed-effects modeling in shrinking outlier intercepts towards the mean is less pronounced.

Nonetheless, a random intercept model with only the predictors returned by the stepwise variable selection process as outlined previously decreased the AIC and maintained the impact of the random intercept. Notably, the coefficients in the linear mixed-effects model did not change substantially from the ordinary least squares model with the same predictors. While most of the coefficients were slightly closer to zero in the linear mixed-effects model, the  $t$  statistics (i.e. level of significance of the coefficients) generally stayed the same as the addition of the random intercept helped decrease the variability in the model.

```
df <- data.frame(lm_coefficients=summary(step.both)$coefficients[,1],
                 lmer_coefficients=summary(lmer_step)$coefficients[,1])
df$difference <- df$lm_coefficients - df$lmer_coefficients
df <- sapply(df, signif, 3)
rownames(df) <- rownames(summary(step.both)$coefficients)
knitr::kable(df)
```

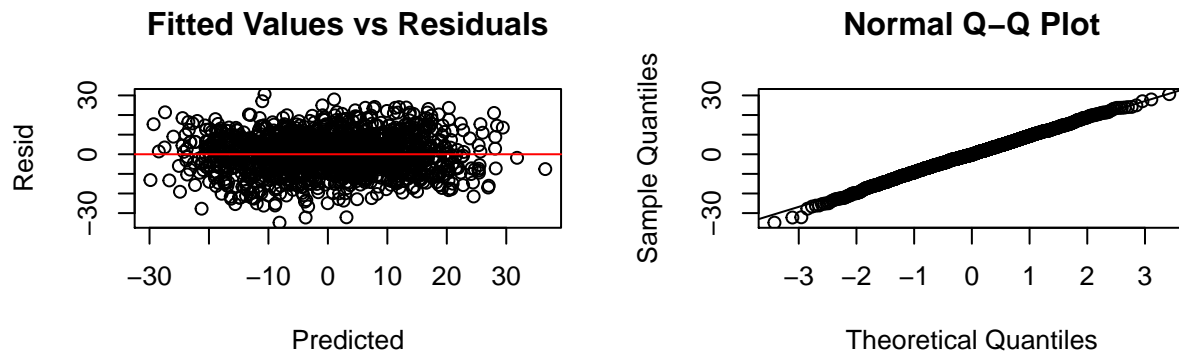
	lm_coefficients	lmer_coefficients	difference
(Intercept)	-15.2000	-12.5000	-2.790000



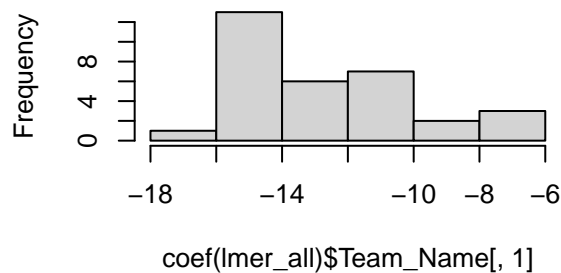
	lm_coefficients	lmer_coefficients	difference
Home	1.5200	1.5900	-0.069000
QB_Hit	-0.8980	-0.8070	-0.091400
Pass_Yrd	0.0350	0.0359	-0.000954
Rush_Yrd	0.0786	0.0744	0.004130
FG_attempt	0.8820	0.8790	0.003230
XP_missed	1.0400	1.0500	-0.007260
FG_missed	-1.7900	-1.6600	-0.131000
FUM	2.0300	2.0100	0.020400
INT	3.8800	3.7000	0.181000
Fourth_Dwn_Con	-0.9770	-0.8510	-0.127000
Third_Dwn_Con	0.7680	0.7740	-0.006470
RedZone_Plays	0.5240	0.4920	0.031300
Unique_Receivers	-0.8510	-0.7810	-0.069700
Shotgun_Plays	-0.1930	-0.2710	0.078100

We also fit a linear mixed-effects model with a random intercept by team and random slope for the relationship of total rushing yards with final score differential. In theory, if some teams were primarily a “passing” team, total rushing yards might have less of an outcome on final score differential than if their primary game strategy relied on the run. However, we find that this random slope model actually worsens the model AIC suggesting that this is not a particularly important trend.

The fitted values vs residuals plot and Normal Q-Q plots are both reasonable for this random-intercept mixed-effects model, demonstrating that constant variance and normality of residuals are acceptable assumptions for this model. However, there is some concern with the distribution of fitted random intercept estimates—the histogram is neither normal-shaped or symmetric.



**Histogram of Random Intercept Estim**



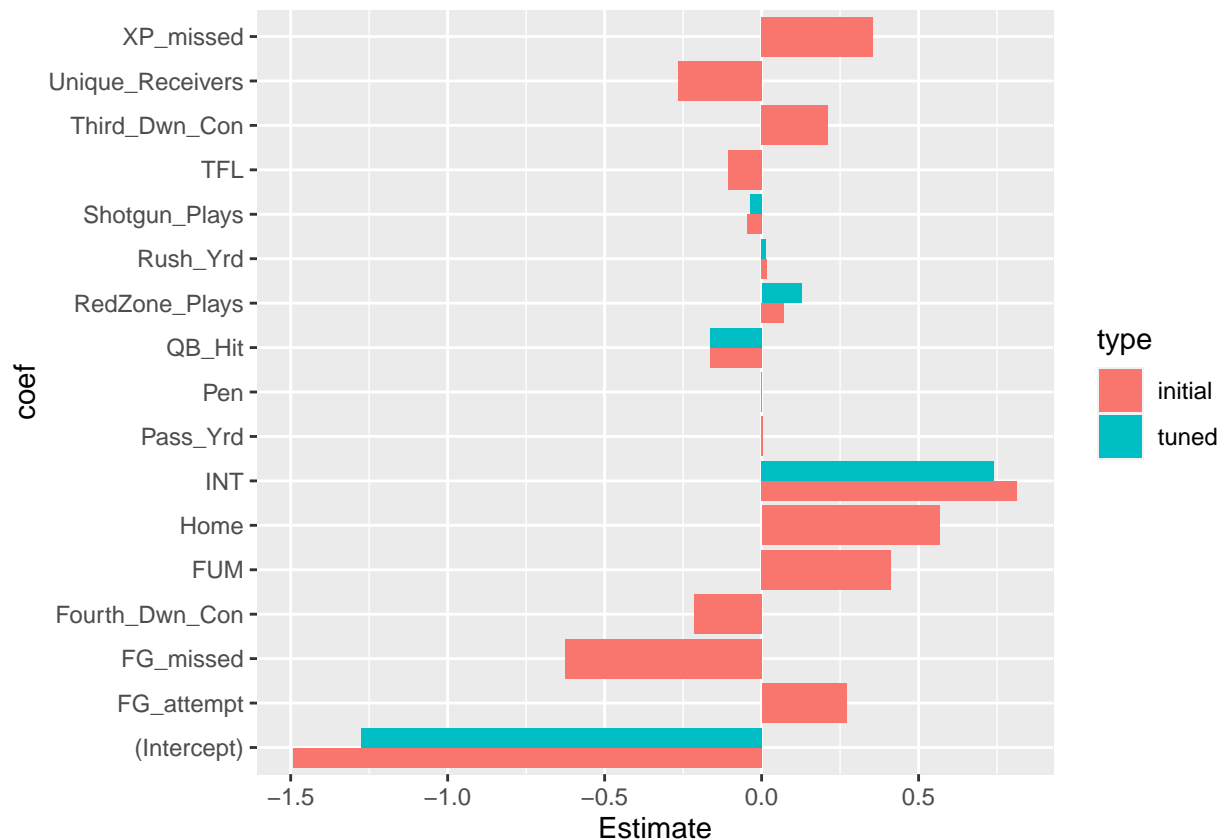
	RMSE.Train	RMSE.Test
lmer	9.360256	10.60430
lmer_stepwise	9.371261	10.57038
lmer_stepwise_rslope	9.321312	10.58357
lm.all	9.851668	10.35031
lm.step	9.853232	10.35228

	RMSE.Train	RMSE.Test
lm.step.interaction	9.563922	10.37966
Ridge	11.595004	12.23013
Lasso	10.441002	10.99376
Random Forest	10.492822	10.50129

	AIC	BIC
lmer	11875.70	11977.90
lmer_stepwise	11869.88	11961.32
lmer_stepwise_rslope	11871.62	11973.82
lm.all	11911.88	12008.70
lm.step	11908.39	11994.45
lm.step.interaction	11840.90	12002.27

## Logistic applications:

Although our main focus was predicting score differential, we thought that predicting the probability of a win would be an interesting extension to our project, and would highlight some of the oversights of our models as well as some interesting scenarios that just can't be handled meaningfully (looking at you, Bill Belichick). In this section, we fit a glm model with the parameter of being in the binomial family for the response, which in this case is wins. Our assumptions are minimal, as seen in the lectures, and all we need is independent observations which we generally have. We recorded each teams' wins for each game, to be used as a response. Then, we created two glm's, one with the top 5 predictors from the random forest variable importance (Rush Yards, RedZone Plays, Qb Hits, Shotgun plays, and interceptions) as well as a model with all predictors except for the game id, year, team, and score differential. Using these models, we could create models showing what things contributed to the probability of a win. Doing this, we could compare the two models' coefficients and see what changed from model to model:



Our plot shows most notably that the intercept lost magnitude when predictors were removed, and as a result changed the

magnitude of the kept predictors. Most importantly of which was rushing yards, whose coefficient decreased a lot from the initial model (which may not look like a ton on this plot). The decrease from the initial model to the tuned model was 0.016 to 0.013, and is significant because rushing yards are usually on the scale of about 120 to 180 per team per game.

Interestingly, both models had a worse train accuracy than test accuracy, which was in the years 2017-2019, as seen in the table below. Showing that our model wasn't too overfit and did have some predictive power.

```
##   Train_Acc_Init Train_Acc_Best Test_Acc_Init Test_Acc_Best
## 1      0.7940075      0.7509363      0.806701      0.7860825
```

Using these models, we were able to explore some interesting extensions, such as the Patriots week 4 bout versus the Buccaneers. In this game, the Patriots exceeded expectations and held the defending super bowl champions to 6 points in the first half, and 13 in the second, good for a win in many weeks. Importantly, they played them close all game and the Bucs won only on a last second field goal. What is strange, however, is that our model showed the Patriots as having a 3.78% probability of a win in the tuned model, and a 1.21% chance in the initial model. The Bucs were given a 70.2% and 86.5% probability in the tuned and untuned model, respectively. So, why did our model completely miss on a game that, honestly, should have been won by the Patriots? Taking a closer look, we can see that the Pats had -1 rush yards (no, not an error). Our model heavily favors rush yards and gives passing yards almost no weights, especially in the tuned model where passing yards meant nothing. The Patriots often perform unusual game plans, such as only throwing 3 times in a 14-10 win versus the Buffalo Bills on December 6, 2021. These game plans clearly throw off our model, showing that it really isn't robust to model a win probability logistically off of in game data (for the Patriots, at least). There's something about the "eye test" (a theoretical, informal "test" that involves years of watching and analyzing games, making snap judgements about outcomes and decisions while watching the game). Someone watching the week 4 Patriots-Bucs game surely wouldn't be giving the Patriots a ~4% chance of winning after seeing them absolutely dominate in the first half, yet our model says that not having any rush yards is highly correlated with losing the game.

Importantly, the model did actually get that prediction correct, but odd scenarios like this can mess with the model's accuracy (off the top of my head, the falcon's game versus the saints where they had 34 rushing yards and won is an example of a false prediction). Overall, rating one metric too highly can mess with the model, just as in any model, so we must be careful about determining direct causation.

In the long run, we'd be able to claim that rushing for more yards likely means you are winning and don't have to pass, which points to there being a win, but a team could just rush for every play and rack up the yards and still lose, which is why we cannot say that rushing yards increase your chances of winning, but rather we associate higher amounts of rushing yards with a higher chance of winning.

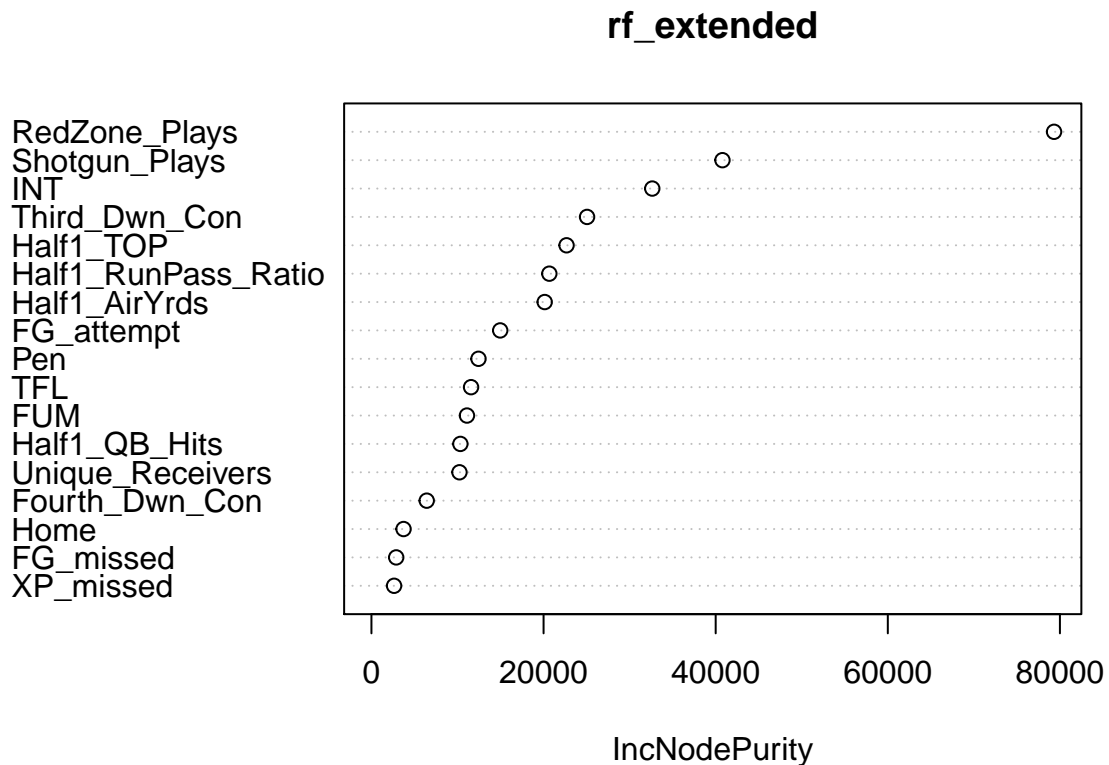
```
##   NE_Chance_Tuned NE_Chance_Untuned TB_Chance_Tuned TB_Chance_Untuned
##      0.0378281      0.01212202      0.7022904      0.8652279
```

## Extension: 1st Half Predictors

Many strategies during NFL games are determined by the context of the game as it unfolds. For instance, a team that is trailing greatly in the 4th quarter will be more inclined to throw risky passes downfield, resulting in meaningless higher potential passing yards, interceptions, quarterback hits, unique receivers, and other related predictors. Similarly, a team with a sizable lead late in the game is incentivized to be more conservative with their play calling; these teams are more likely to run the ball and attempt to let the clock expire. A team's play makeup in the first half, however, may be more resulting from premeditated strategy than the context of the game itself. We examine multiple first-half predictors, including run-pass play ratio, air yards, and time of possession. All of these new predictors were found to have symmetrical distributions, except for the run-pass ratio, which was right skewed and fixed with a log transformation.

These predictors are incorporated into previously generated models to weigh their comparative importance with older predictors, while removing predictors similar to our first half ones (rushing yards, passing yards, QB hits).

A random forest was fitted using these new predictors. It was run before log-transforming the Run/Pass ratio data, as there are no underlying distributional assumptions. The variable importance plot produced with this new data in the presence of old predictors shares the same general structure as before. The new 1st half predictors can be seen grouped together in what appears to be their own "tier" of importance, below the same top 6 important predictors as before. This may be a result of multicollinearity between the newly introduced predictors. Half1\_QB\_Hits appears to be a relatively unimportant predictor.



We also conducted backwards/both selection based on AIC to see if the new 1st half variables are present. We see that backward a

## Conclusion

National Football League data provides a ripe opportunity for statistical analysis, especially when able to access play-by-play level data using the `nflfastR` package. While most variables have a slight right skew to them and only take on positive values, this did not present tremendous challenges for our models. Going forward, we're interested to continue investigating the impact of variables that provide a clear relationship towards "success" in a game—such as total passing yards or total rushing yards—as well as more "strategic" variables such as the number of third and fourth down conversions. Given the violations of symmetric distributions in many of the predictors and the potential for non-linear relationships to arise, examining the performance and interpretations of Random Forest models will also be a focus in our analysis.

## Appendix

### Packages and Cleaning:

```
library(tidyverse)
library(ggrepel)
library(ggimage)
library(nflfastR)
library(dplyr)
library(gridExtra)
library(randomForest)
library(lme4)
library(xfun)
library(lubridate)
options(scipen = 9999)
```

```

data17 <- load_pbp(2017)
data18 <- load_pbp(2018)
data19 <- load_pbp(2019)
data20 <- load_pbp(2020)
data21 <- load_pbp(2021)
#head(data1$game_id,100)

#unique(data1$game_id)
clean = function(data1){
  names = c("Game_id", "Team_Name", "Home", "TFL", "QB_Hit", "Pass_Yrd", "Rush_Yrd", "FG_attempt", "XP_missed", "FG_missed", "FUM", "INT")

  df <- data.frame()

  for(i in unique(data1$game_id)) {
    temp <- subset(data1, game_id == i)
    for(j in c(temp$home_team[1], temp$away_team[1])) {
      to_add <- c(i, #Game id
                  j, #team name
                  1*(j==temp$home_team[1]), #team is home
                  sum(temp$stackled_for_loss[temp$defteam == j] == 1, na.rm=TRUE), #tackles for loss
                  sum(temp$qb_hit[temp$posteam == j] == 1, na.rm=TRUE), #how many times team's QB was hit
                  sum(temp$passing_yards[temp$posteam == j], na.rm=TRUE), #team passing yards
                  sum(temp$rushing_yards[temp$posteam == j], na.rm=TRUE), #team rushing yards
                  sum(temp$field_goal_attempt[temp$posteam == j], na.rm=TRUE), #team field goal attempts
                  sum(temp$extra_point_result[temp$posteam == j]=="failed" |
                      temp$extra_point_result[temp$posteam == j]=="blocked", na.rm=TRUE), #exp missed
                  sum(temp$field_goal_result[temp$posteam == j]=="missed" |
                      temp$field_goal_result[temp$posteam == j]=="blocked", na.rm=TRUE), #fg missed
                  sum(temp$fumble[temp$defteam == j], na.rm=TRUE), #number of fumbles
                  sum(temp$interception[temp$defteam == j], na.rm=TRUE), #number of interceptions
                  sum(temp$fourth_down_converted[temp$posteam == j], na.rm=TRUE), #4th down conversions
                  sum(temp$third_down_converted[temp$posteam == j], na.rm=TRUE), #3rd down conversions
                  sum(temp$penalty[temp$posteam == j & temp$penalty_team == j], na.rm=TRUE), #number of penalties
                  sum(temp$yardline_100[temp$posteam == j] < 20, na.rm=TRUE), #number of Red Zone Plays
                  length(unique(temp$receiver_player_id[temp$posteam == j])) - 1, #number of unique receivers (remove NA)
                  sum(temp$shotgun[temp$posteam == j], na.rm=TRUE), #number of plays out of the Shotgun Formation
                  temp$result[1]*ifelse(j==temp$away_team[1], -1, 1) #score differential
            )
      df <- rbind(df, to_add)
    }
  }

  colnames(df) = names

  df = df %>%
  mutate_at(c("Home", "TFL", "QB_Hit", "Pass_Yrd", "Rush_Yrd", "FG_attempt", "XP_missed", "FG_missed", "FUM", "INT", "Fourth_Dwn_Con"),
            return(df)
}

RMSE <- function(y.obs, y.pred){
  return(sqrt(mean((y.obs-y.pred)^2)))
}

#DO EDA HERE
data_2017 = clean(data17)
data_2018 = clean(data18)
data_2019 = clean(data19)
data_2021 = clean(data21)
data_2021$year <- 2021
data_2019$year <- 2019
data_2018$year <- 2018
data_2017$year <- 2017

data_2018_2019 <- rbind(data_2018, data_2019)
data_2017_19 <- rbind(data_2017, data_2018, data_2019)

```

```

g.score <- ggplot(data_2019, aes(x=Score_Dif)) +
  geom_histogram(binwidth =1, color="black", size=0.2) +
  ggtitle("Score Differential", subtitle="2019 NFL Games")

g.pass <- ggplot(data_2019, aes(x=Pass_Yrd)) +
  geom_histogram(bins=30, color="black", size=0.2) +
  ggtitle("Total Passing Yards", subtitle="2019 NFL Games")

# ggplot(data_2018_2019, aes(x=Pass_Yrd, group=as.factor(year), fill=as.factor(year))) +
#   geom_histogram(bins=30, color="black", size=0.2, position="identity", alpha=0.6) +
#   ggtitle("Total Passing Yards", subtitle="2018 and 2019 NFL Games")

g.rush <- ggplot(data_2019, aes(x=Rush_Yrd)) +
  geom_histogram(bins=30, color="black", size=0.2) +
  ggtitle("Total Rushing Yards", subtitle="2019 NFL Games")

g.ps_rsh <- ggplot(data_2019, aes(x=Pass_Yrd, y=Rush_Yrd)) +
  geom_point() +
  ggtitle("Relationship of Passing Yards to Rushing Yards", subtitle="2019 NFL Games")

g.fum_sc <- ggplot(data_2019, aes(x=FUM, y=Score_Dif)) +
  geom_point() +
  stat_smooth(method="lm", se=F) +
  ggtitle("Number of Fumbles Caused to Score Difference", subtitle="2019 NFL Games")

g.3_sc <- ggplot(data_2019, aes(x=Third_Dwn_Con, y=Score_Dif)) +
  geom_point() +
  stat_smooth(method="lm", se=F) +
  ggtitle("Third Down Conversions vs Score Difference", subtitle="2019 NFL Games")

g.pass_sc <- ggplot(data_2019, aes(x=Pass_Yrd, y=Score_Dif)) +
  geom_point() +
  ggtitle("Relationship of Passing Yards to Score Difference", subtitle="2019 NFL Games")

g.3 <- ggplot(data_2019, aes(x=Third_Dwn_Con)) +
  geom_histogram(bins=30, color="black", size=0.2) +
  ggtitle("Third Down Conversions", subtitle="2019 NFL Games")

g.4 <- ggplot(data_2019, aes(x=Fourth_Dwn_Con)) +
  geom_histogram(bins=30, color="black", size=0.2) +
  ggtitle("Fourth Down Conversions", subtitle="2019 NFL Games")

g.qb <- ggplot(data_2019, aes(x=QB_Hit)) +
  geom_histogram(bins=30, color="black", size=0.2) +
  ggtitle("QB Hits", subtitle="2019 NFL Games")

#gridExtra::grid.arrange(g.pass, g.rush, nrow=1)
gridExtra::grid.arrange(g.fum_sc, g.3_sc)
gridExtra::grid.arrange(g.pass, g.rush, g.3, g.4, g.qb, nrow=3)

```

Stepwise:

```
lm1 = lm(Score_Dif~. -Team_Name -year, data=data_2017_19[,seq(2, ncol(data_2017_19))])
```

```

## Score_Dif ~ Home + QB_Hit + Pass_Yrd + Rush_Yrd + FG_attempt +
##   XP_missed + FG_missed + FUM + INT + Fourth_Dwn_Con + Third_Dwn_Con +
##   RedZone_Plays + Unique_Receivers + Shotgun_Plays

```

*#Stepwise Both*

```

step.both = step(lm1, scope = c(lower = formula(Score_Dif ~ 1), upper = formula(Score_Dif ~ (. - Game_id - Team_Name)^2)), direction="both")

```

Random Forest:

```

maxnodes <- c(100, 500, 1000)
mtry <- c(1,5,8,10,12)

```

```

combs <- expand.grid(mtry, maxnodes)
colnames(combs) <- c("mtry", "maxnodes")
combs$RMSE_train <- rep(NA, nrow(combs))
combs$RMSE_test <- rep(NA, nrow(combs))

set.seed(139)
for(i in 1:nrow(combs)) {
  temp_forest <- randomForest(Score_Dif~.-Game_id-Team_Name-year,data=data_2017_19, mtry = combs$mtry[i], maxnodes = combs$maxnodes)
  combs$RMSE_train[i] <- RMSE(data_2017_19$Score_Dif, temp_forest$predicted)
  combs$RMSE_test[i] <- RMSE(data_2021$Score_Dif, predict(temp_forest, new = data_2021))
}
combs$RMSE_difference <- combs$RMSE_test - combs$RMSE_train

rf1 <- randomForest(Score_Dif~.-Game_id-Team_Name-year,
  data=data_2017_19,
  mtry = combs[which.min(combs$RMSE_test),]$mtry,
  maxnodes = combs[which.min(combs$RMSE_test),]$maxnodes,
  ntree=200)
rf.varimp <- varImpPlot(rf1)

```

### Ridge Regression and LASSO:

```

library(glmnet)
library(Rcpp)
library(Matrix)
#Baseline Linear Models
X = model.matrix(lm1)[,-1] #drop the intercept

#Ridge Model
ridges = glmnet(X, data_2017_19$Score_Dif, alpha = 0, lambda = 10^seq(-4,4,0.1), standardize = F)
#coef(ridges)
#ridges$lambda

n.test = nrow(data_2021)
lm.test = lm(formula(lm1), data = data_2021)
Xtest = model.matrix(lm.test)[,-1] #drop the intercept again
yhats.ridges = predict(ridges, Xtest)
residuals.ridge = (data_2021$Score_Dif - yhats.ridges)^2
MSE.ridges.test = apply(residuals.ridge, 2, sum)/n.test

plot(MSE.ridges.test~log(ridges$lambda,10),type="l")
ridges$lambda[which.min(MSE.ridges.test)]
min(MSE.ridges.test)

matplot(log(ridges$lambda, 10), t(ridges$beta),type="l")

#Lasso Model
lassos = glmnet(X, data_2017_19$Score_Dif, alpha = 1, lambda = 10^seq(-4,4,0.1), standardize = F)
#coef(lassos)
#lassos$lambda

yhats.lassos = predict(lassos, Xtest)
residuals.lasso = (data_2021$Score_Dif - yhats.lassos)^2
MSE.lassos.test = apply(residuals.lasso, 2, sum)/n.test

plot(MSE.lassos.test~log(lassos$lambda,10),type="l")
lassos$lambda[which.min(MSE.lassos.test)]
min(MSE.lassos.test)
matplot(log(lassos$lambda, 10), t(lassos$beta),type="l")

#Stepwise Models
#Stepwise Backward
step.back = step(lm1, direction = "backward", k = 2, trace = F)
formula(step.back)

#Stepwise Both

```

```

step.both = step(step.back,
  scope = c(lower = Score_Dif ~ . - Game_id - Team_Name -year, upper = formula(Score_Dif ~ (. - Game_id - Team_Name -year)^2)),
  direction = "both", trace = F)
step.both.int = step(step.back,
  scope = c(lower = Score_Dif ~1,
    upper = formula(Score_Dif ~ (. - Game_id - Team_Name -year)^2)),
  direction = "both", trace = F)

```

## LMER Modeling:

```

lmer_all <- lmer(Score_Dif ~ Home +TFL + QB_Hit + Pass_Yrd + Rush_Yrd + FG_attempt +
  XP_missed + FG_missed + FUM + INT + Fourth_Dwn_Con + Third_Dwn_Con + Pen +
  RedZone_Plays + Unique_Receivers + Shotgun_Plays + (1|Team_Name),data=data_2017_19)

```

```

lmer_step <- lmer(Score_Dif ~ Home + QB_Hit + Pass_Yrd + Rush_Yrd + FG_attempt +
  XP_missed + FG_missed + FUM + INT + Fourth_Dwn_Con + Third_Dwn_Con +
  RedZone_Plays + Unique_Receivers + Shotgun_Plays + (1|Team_Name),
  data=data_2017_19)

```

```

lmer_step_slope <- lmer(Score_Dif ~ Home + QB_Hit + Pass_Yrd + Rush_Yrd + FG_attempt +
  XP_missed + FG_missed + FUM + INT + Fourth_Dwn_Con + Third_Dwn_Con +
  RedZone_Plays + Unique_Receivers + Shotgun_Plays + (1+Rush_Yrd|Team_Name),
  data=data_2017_19)

```

```

df <- data.frame(lm_coefficients=summary(step.both)$coefficients[,1],
  lmer_coefficients=summary(lmer_step)$coefficients[,1])
df$difference <- df$lm_coefficients - df$lmer_coefficients
df <- sapply(df, signif, 3)
rownames(df) <- rownames(summary(step.both)$coefficients)
knitr::kable(df)

```

```

par(mfrow=c(2, 2))
plot(resid(lmer_all)~predict(lmer_all), main="Fitted Values vs Residuals", xlab="Predicted", ylab="Resid")
abline(h=0, col = "red")
#normality of residuals
qqnorm(resid(lmer_all))
qqline(resid(lmer_all))
#check normality of random effects - non-normal
hist(coef(lmer_all)$Team_Name[,1], main="Histogram of Random Intercept Estimates")

```

```

list.models <- list(lmer_all, lmer_step, lmer_step_slope, lm1, step.both, step.both.int)
list.all.models <- list(lmer_all, lmer_step, lmer_step_slope, lm1, step.both, step.both.int)
model.results <- data.frame(unlist(lapply(list.models, AIC)),
  unlist(lapply(list.models, BIC)),
  unlist(lapply(list.models, predict, newdata=data_2017_19) %>% lapply(RMSE, data_2017_19$Score_Dif)),
  unlist(lapply(list.models, predict, newdata=data_2021) %>% lapply(RMSE, data_2021$Score_Dif))) %>% d

names(model.results) <- c("AIC", "BIC", "RMSE.Train", "RMSE.Test")
rownames(model.results) <- c("lmer", "lmer_stepwise", "lmer_stepwise_rslope", "lm.all", "lm.step", "lm.step.interaction")
model.rmse <- rbind(c(RMSE(predict(lassos, X), data_2017_19$Score_Dif),
  RMSE(predict(lassos, Xtest), data_2021$Score_Dif)),
  c(RMSE(predict(ridges, X), data_2017_19$Score_Dif),
  RMSE(predict(ridges, Xtest), data_2021$Score_Dif)),
  c(RMSE(rf1$predicted, data_2017_19$Score_Dif),
  RMSE(predict(rf1, data_2021), data_2021$Score_Dif))) %>% data.frame()
names(model.rmse) <- c("RMSE.Train", "RMSE.Test")

rmse.models <- rbind(model.results[,c(3, 4)], model.rmse)
rownames(rmse.models) <- c("lmer", "lmer_stepwise", "lmer_stepwise_rslope", "lm.all", "lm.step", "lm.step.interaction", "Ridge", "Ridge")
knitr::kable(rmse.models)

knitr::kable(model.results[,1:2])

```



## Logistic applications:

```
data_2017_19_test = data.frame(data_2017_19)

data_2017_19_test$wins = ifelse(data_2017_19$Score_Dif > 0, 1, 0)

win.log = glm(wins~.-Score_Dif-year, data=data_2017_19_test[,seq(3, ncol(data_2017_19_test))], family="binomial")

best.log = glm(wins~RedZone_Plays + Rush_Yrd+Shotgun_Plays +QB_Hit+INT, data = data_2017_19_test, family="binomial")

init_coef_log = data.frame(summary(win.log)$coefficients)
best_coef_log = data.frame(summary(best.log)$coefficients)

init_coef_log$coef = row.names(init_coef_log)
best_coef_log$coef = row.names(best_coef_log)
init_coef_log$type = "initial"
best_coef_log$type = "tuned"

ggplot(rbind(init_coef_log, best_coef_log), aes(fill=type, y=Estimate, x=coef)) +
  geom_bar(position="dodge", stat="identity") + coord_flip()
```

```
new_2021 = data.frame(data_2021)
new_2021$year = 2021

new_2021$wins = ifelse(data_2021$Score_Dif > 0, 1, 0)

predicted_wins_initial = 1*(predict(win.log, type = "response")>0.5)
predicted_wins_best = 1*(predict(best.log, type = "response")>0.5)

Train_Acc_Init = sum(predicted_wins_initial == data_2017_19_test$wins)/nrow(data_2017_19_test)
Train_Acc_Best = sum(predicted_wins_best == data_2017_19_test$wins)/nrow(data_2017_19_test)

pred_test_init = 1*(predict(win.log, newdata=new_2021, type = "response")>0.5)
pred_test_best = 1*(predict(best.log, newdata=new_2021, type="response")>0.5)

Test_Acc_Init = sum(pred_test_init == new_2021$wins)/nrow(new_2021)
Test_Acc_Best = sum(pred_test_best == new_2021$wins)/nrow(new_2021)

data.frame(Train_Acc_Init, Train_Acc_Best, Test_Acc_Init, Test_Acc_Best)
```

```
ne.game.init = predict(win.log, newdata = new_2021[new_2021$Game_id=="2021_04_TB_NE" & new_2021$Home==1,], type="response")
tb.game.init = predict(win.log, newdata = new_2021[new_2021$Game_id=="2021_04_TB_NE" & new_2021$Home==0,], type="response")
ne.game.best = predict(best.log, newdata = new_2021[new_2021$Game_id=="2021_04_TB_NE" & new_2021$Home==1,], type="response")
tb.game.best = predict(best.log, newdata = new_2021[new_2021$Game_id=="2021_04_TB_NE" & new_2021$Home==0,], type="response")
data.frame(NE_Chance_Tuned = ne.game.best, NE_Chance_Untuned = ne.game.init, TB_Chance_Tuned = tb.game.best, TB_Chance_Untuned =
```

## Extension: 1st Half Predictors:

```
##### FIRST HALF PREDICTORS EXTENSION #####

temp.df <- c()
for(i in c("17","18","19")) {
  temp.raw <- get(paste0("data", i))
  for(j in unique(temp.raw$game_id)) {
    temp <- subset(temp.raw, game_id == j)
    for(k in c(temp$home_team[1], temp$away_team[1])) {
      Half1_RunPass_Ratio <-
        sum(temp$rush_attempt[temp$posteam == k & temp$game_half == "Half1"], na.rm=TRUE) /
        sum(temp$pass_attempt[temp$posteam == k & temp$game_half == "Half1"], na.rm=TRUE)
```

```

Half1_AirYrds <- sum(temp$air_yards[temp$posteam == k & temp$game_half == "Half1"], na.rm=TRUE)
Half1_QB_Hits <- sum(temp$qb_hit[temp$posteam == k & temp$game_half == "Half1"] == 1, na.rm=TRUE)
## Calculate Time of Possession
pos <- period_to_seconds(ms(temp$time[temp$posteam == k & temp$game_half == "Half1"]))
Half1_TOP <- abs(sum(c(0, lag(pos)[2:length(pos)]) - pos , na.rm=TRUE))
## build final DF with data to add
temp.df <- as.data.frame(rbind(cbind(Half1_RunPass_Ratio,
                                     Half1_AirYrds,
                                     Half1_QB_Hits,
                                     Half1_TOP),
                               temp.df))
  }
}
}

Extended_Data <- cbind(data_2017_19[,!(colnames(data_2017_19) %in% c("QB_Hit", "Rush_Yrd", "Pass_Yrd"))], temp.df)

rf_extended <- randomForest(Score_Dif~.-Game_id-Team_Name-year,
                             data=Extended_Data,
                             mtry = combs[which.min(combs$RMSE_test),]$mtry,
                             maxnodes = combs[which.min(combs$RMSE_test),]$maxnodes,
                             ntree=200)
varImpPlot(rf_extended)

step.backward.extended = step(lm(Score_Dif~.- Game_id - Team_Name -year, data=Extended_Data), scope = c(lower = formula(Score_Di

step.both.extended = step(lm(Score_Dif~.- Game_id - Team_Name -year, data=Extended_Data), scope = c(lower = formula(Score_Dif ~1

```