

## MSDS 6372 Group Project I

### **INTRODUCTION:**

For our group project, we decided to look at the professional football team New England Patriots. Our goal is twofold:

- 1) To determine which, if any, variables explain the number of regular season wins for the Patriots since Tom Brady has been starting quarterback (last 18 seasons).
- 2) To determine if there is a statistical difference in Tom Brady's performance over his wins and losses across his playoff career.

### **DATA DESCRIPTION:**

We used two datasets, both obtained from [www.pro-football-reference.com](http://www.pro-football-reference.com). One dataset includes combined offensive and defensive stats, 70 variables total, for each of the Patriots' last 18 seasons (2001 through 2018). This data was compiled from multiple seasons on the [www.pro-football-reference.com](http://www.pro-football-reference.com) website, but the compiled dataset can be found at the "season stats" link below. The second dataset consists of metrics describing Tom Brady's playoff performances during his tenure as starting quarterback for the New England Patriots. The data focuses on specifically passing metrics during both playoff and super bowl performances. The raw data can be found at the "raw data" link below and the cleaned data can be found at the project GitHub at the link below.

#### **Season Stats (all regular season statistics):**

<https://github.com/newtgunslinger/6372.404.AS.Project1/blob/master/PatriotsYearlyStats.csv>

#### **Tom Brady Playoff Stats:**

**Raw Data:** [https://www.pro-football-reference.com/players/B/BradTo00.htm#all\\_passing\\_playoffs](https://www.pro-football-reference.com/players/B/BradTo00.htm#all_passing_playoffs)

**Cleaned Data:** <https://github.com/newtgunslinger/6372.404.AS.Project1/blob/master/BradyStats.csv>

Additional variables were created to support our analysis and some of the [www.pro-football-reference.com](http://www.pro-football-reference.com) require a little more explanation (see below).

#### **Season Stats (all regular, per season statistics):**

RegSeasonWins: Total wins per season.

BradyPasserRating: Tom Brady's average passer rating.

PointsFor: Total points scored by the New England Patriots.

PointsAgainst: Total points allowed by the New England Patriots.

PointsDifferential: Total points scored minus total points allowed.

MarginOfVictory: Average point differential (points scored minus points allowed) per game.

StrengthOfSchedule: A metric describing average quality of opponent, measured by simple rating system.

SimpleRatingSystem: A metric describing team quality relative to average.

OffSimpleRatingSys: A metric describing offensive quality relative to average.

DefSimpleRatingSys: A metric describing defensive quality relative to average.

Yards: Total offensive yards made.

Plays: Total offensive plays.  
YardsPerPlay: Total offensive yards per play.  
Turnovers: Total turnovers.  
FumblesLost: Total fumbles.  
FirstDowns: Total number of first downs.  
PassCompletions: Total number of pass completions.  
PassAttempts: Total number of pass attempts.  
CompletionPercentage: A ratio of PassCompletions to PassAttempts.  
PassYards: Total number of passing yards.  
PassTouchdowns: Total number of passing touchdowns.  
PassInterceptions: Total number of interceptions.  
NetYardsPerPass: Ratio of passing yards minus sack yards to passing attempts plus times sacked.  
PassFirstDowns: Total number of passing first downs.  
RushAttempts: Total number of rushing attempts.  
RushYards: Total number of rushing yards.  
RushTouchdowns: Total number of rushing touchdowns.  
RushYardsPerAttempt: Average rushing yards per attempts.  
RushFirstDowns: Total number of rushing first downs.  
Penalties: Total number of penalties.  
PenaltyYards: Total number of penalty yards.  
PenaltyFirstDowns: Total number of New England Patriots penalties resulting in a first down.  
NumberDrives: Total number of offensive drives.  
DriveScorePercent: Percentage of an offensive drive resulting in a score.  
DriveTurnoverPercent: Percentage of an offensive drive resulting in an offensive turnover.  
AvgStartingPosition: Average yardage marker starting position for offense.  
AvgDriveTime: Average amount of time run off the clock per drive.  
AvgDrivePlays: Average number of plays per offensive drive.  
AvgDriveYards: Average number of yards consumed per drive.  
AvgDrivePoints: Average number of points scored per offensive drive.

**Tom Brady Playoff Stats (only three variables used):**

AgeBin: A categorical variable describing the point in Tom Brady's career - early, mid, and late (1, 2, and 3 respectively).  
Rate: Tom Brady's passer rating for the individual playoff and Super Bowl games.  
WonLost: Game won or lost by the New England Patriots.

**EXPLORATORY ANALYSIS:**

We used both SAS and R in our exploratory analysis and ultimately used SAS output for report quality diagrams. In both SAS and R, scatter plot matrix diagrams were produced to identify any variables that were colinear, identify non-linear trends in the data ideal for transformation, and highlight which variables had highest correlation with our response variable, RegSeasonWins. It was determined that opponent statistics were duplicative when considering strength of schedule. Additionally, it was found that many of the offensive variables were also duplicative, pass completions and passing attempts versus completion percentage, for instance. Using the scatter plot matrix, we were able to cherry-pick a handful of variables for our analysis.

## **OBJECTIVE 1:**

### **Problem:**

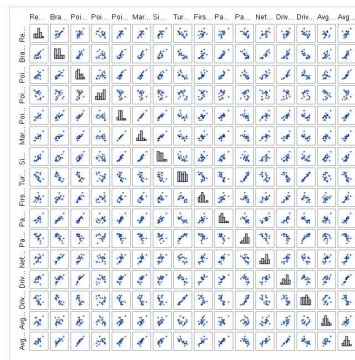
Which variables correlate to regular season wins for Tom Brady's Patriots?

### **Overall Approach:**

Using the aggregated regular season statistics, we will determine which variables correlate with the New England Patriots' regular season record and subsequently build a model using those parameters to interpret the relationship between those parameters and the New England Patriots' regular season record.

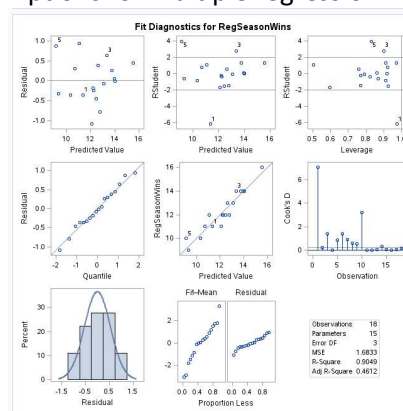
### **Determining Predictors:**

First, we loaded the dataset into SAS.<sup>[1]</sup> Then, we performed exploratory data analysis on the data by creating multiple scatterplot matrices to identify which, if any, showed a correlation with regular season wins.<sup>[2]</sup> We identified 16 variables that showed a correlation with regular season wins and it was determined that none of the variables required a transformation. Our next step was to create a scatterplot matrix with these 16 variables.<sup>[3]</sup>



### **Checking Assumptions:**

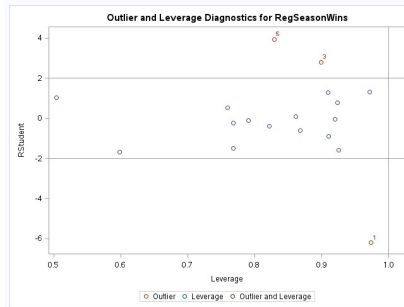
We then checked to see if the assumptions for multiple regression were met:<sup>[4]</sup>



1) The residuals are normally distributed (predictors and response variables don't have to be).

2) There is constant variance.

3) The observations (different seasons) are independent from and of one another.



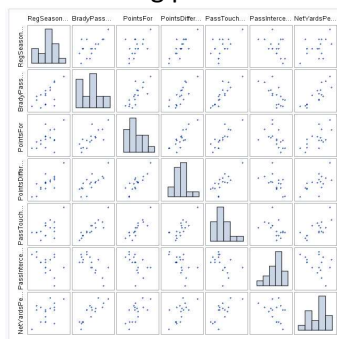
4) Checking through the residual diagnostics, there are three outliers and no leverage points. We will determine if any or all the three outliers need to be removed.

Correlation																
Variable	BradyPasserRating	PointsFor	PointsAgainst	PointsDifferential	MarginOfVictory	SimpleRatingSystem	Turnovers	FirstDowns	PassTouchdowns	PassInterceptions	NetYardsPerPass	DriveScorePercent	DriveTurnoverPercent	AvgDriveTime	AvgDrivePoints	RegressionWins
BradyPasserRating	1.0000	0.7628	-0.0060	0.7617	0.7608	0.7621	-0.7044	0.4764	0.8830	-0.7038	0.9801	0.7591	-0.7108	0.3882	0.8288	0.8313
PointsFor	0.7628	1.0000	0.0041	0.8484	0.8481	0.7807	-0.8485	0.7947	0.8730	-0.8391	0.8833	0.8844	-0.8426	0.3812	0.8432	0.8381
PointsAgainst	-0.0060	0.0041	1.0000	-0.2801	-0.2804	-0.2888	-0.2483	0.3022	0.1831	0.0746	0.0807	0.3851	-0.2241	-0.2844	0.2185	-0.4838
PointsDifferential	0.7607	0.8484	-0.2801	1.0000	0.8328	-0.8481	0.8788	0.7782	-0.8883	0.8844	0.7328	-0.8742	0.8239	0.2142	0.8610	0.8310
MarginOfVictory	0.7608	0.8481	-0.2804	0.8328	1.0000	0.8240	-0.8387	0.8784	0.7788	-0.8821	0.8838	0.7327	-0.8768	0.8210	0.8138	0.8312
SimpleRatingSystem	0.7621	0.7807	-0.2888	0.8240	0.8240	1.0000	-0.4188	0.4473	0.7850	-0.4131	0.8789	0.8320	-0.3828	0.8388	0.7582	0.7328
Turnovers	-0.7044	-0.8485	-0.2483	-0.8481	-0.8387	-0.4188	1.0000	-0.8185	-0.8879	0.8883	-0.8717	-0.7482	0.8778	-0.1389	-0.7318	-0.4218
FirstDowns	0.4764	0.7947	0.3022	0.8788	0.8784	0.4473	-0.8185	1.0000	0.8344	-0.8344	0.8128	0.8318	-0.4503	0.2473	0.7718	0.2144
PassTouchdowns	0.8830	0.8730	0.1831	0.7782	0.7788	0.7850	-0.8879	0.8344	1.0000	-0.4882	0.7480	0.7448	-0.8485	0.3713	0.8488	0.8734
PassInterceptions	-0.7038	-0.8391	0.0746	-0.8883	-0.8821	-0.4131	0.8883	-0.8344	-0.4882	1.0000	-0.4871	0.8841	0.8888	-0.1786	-0.8778	-0.8384
NetYardsPerPass	0.9801	0.8833	0.0807	0.8844	0.8838	0.8788	-0.8717	0.8128	0.7480	-0.4871	1.0000	0.8821	-0.8718	0.4282	0.7721	0.4798
DriveScorePercent	0.7591	0.8844	0.3851	0.7328	0.7327	0.8320	-0.7482	0.8318	0.7448	-0.8841	0.8821	1.0000	-0.8844	0.8842	0.8838	0.4413
DriveTurnoverPercent	-0.7108	-0.8426	-0.2241	-0.8742	-0.8768	-0.3828	0.8778	-0.4503	-0.8485	0.8888	-0.8718	-0.8844	1.0000	-0.2178	-0.8388	-0.4401
AvgDriveTime	0.3882	0.3812	-0.2844	0.8209	0.8210	0.8388	-0.1389	0.2473	0.3713	-0.1786	0.4282	0.8382	-0.8718	1.0000	0.8881	0.2843
AvgDrivePoints	0.8288	0.8432	0.2185	0.8138	0.8138	0.7582	-0.7318	0.7718	0.8488	-0.8778	0.7721	0.8833	-0.8388	0.8881	1.0000	0.4887
RegressionWins	0.8313	0.8381	-0.4838	0.8310	0.8312	0.7328	-0.4218	0.2144	0.8734	-0.8324	0.4798	0.4413	-0.4401	0.2843	0.4887	1.0000

Parameter Estimates					
Variable	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	8	80.42854	24.88588	1.76	0.1181
BradyPasserRating	8	-0.28342	0.28985	-1.05	0.3052
PointsFor	8	2.40425	1.01048	1.33	0.2162
PointsAgainst	8	-2.44858	1.02742	-1.34	0.2122
PointsDifferential	8	0	0	0	0
MarginOfVictory	8	-38.18877	28.10228	-1.31	0.2035
SimpleRatingSystem	1	-0.01437	0.28970	-0.54	0.5928
Turnovers	1	-0.24879	0.71458	-0.32	0.7512
FirstDowns	1	-0.00588	0.02883	-1.78	0.1188
PassTouchdowns	8	0.24841	0.37287	0.68	0.5058
PassInterceptions	1	0.01807	0.88848	0.13	0.8978
NetYardsPerPass	1	2.84188	2.37728	1.20	0.2518
DriveScorePercent	1	0.00480	0.33003	0.01	0.9893
DriveTurnoverPercent	8	0.22480	1.08232	0.17	0.8746
AvgDriveTime	1	-108.78774	117.31184	-1.35	0.2083
AvgDrivePoints	8	4.84759	10.12838	0.48	0.6348

5) Looking at the correlation, we removed PointsAgainst, Turnovers, FirstDowns, DriveTurnoverPercent, and AvgDriveTime for having low correlations with RegSeasonWins. Looking through the VIFs, we also removed MarginOfVictory, SimpleRatingSystem, and AvgDrivePoints for having correlation with other variables (multicollinearity assumption).

We created a scatterplot matrix with the remaining predictors:[5]



## Model Selection:

We obtained our model through a series of feature selection tools: LARS, LASSO, stepwise, forward, and leave-one-out cross validation which was a k-fold cross validation taken to its extreme.<sup>[6]</sup>

## Final Model:

### Leave-One-Out Cross Validation:

The GLMSELECT Procedure  
Selected Model

The selected model is the model at the last step (Step 2).

Effects: Intercept PointsFor PointsDifferential

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	2	37.74542	18.87271	18.42
Error	15	15.36550	1.02436	
Corrected Total	17	53.11111		

Root MSE	1.01212
Dependent Mean	12.22222
R-Square	0.7107
Adj R-Sq	0.6721
AIC	23.15178
AICC	28.22070
PRESS	22.67804
SBC	5.62289

Parameter Estimates				
Parameter	DF	Estimate	Standard Error	t Value
Intercept	1	13.36450	2.34635	5.70
PointsFor	1	-0.01331	0.007036	-1.89
PointsDifferential	1	0.03262	0.006981	4.67

## Model Selection Interpretation:

The LARS model gave us a R-square of 0.6281, the LASSO model gave us a R-square of 0.6281, the stepwise gave us an adjusted R-square of 0.6721, the partial with 2 parameters gave us an adjusted R-square of 0.5940, the partial with 1 parameter gave us a R-square of 0.6417, and the leave-one-out validation gave us an adjusted R-square of 0.6721.

The LARS model gave us an AIC of 25.67176, the LASSO model gave us an AIC of 25.67176, the stepwise gave us an AIC of 23.15178, and the leave-one-out validation gave us an AIC of 23.15178.

Out of all the models, we selected the leave-one-out validation as the best model because our dataset only contains 18 observations. The k-fold cross validation is best when there is a low number of observations because it makes use of all the observations through an iterative process where it removes one observation for testing validation. The R-Square and AIC results we obtained support this.

## Parameter Interpretation and Confidence Intervals:

We ran a test to determine the confidence intervals for the variables.<sup>[7]</sup>

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t	95% Confidence Limits
Intercept	1	13.36450	2.34635	5.70	<.0001	8.36337 18.36583
PointsFor	1	-0.01331	0.00704	-1.89	0.0780	-0.02831 0.00169
PointsDifferential	1	0.03262	0.00698	4.67	0.0003	0.01774 0.04750

$$\text{RegSeasonWins} = 13.3645 - 0.0133(\text{PointsFor}) + 0.0326(\text{PointsDifferential})$$

95% Confidence Limits for Intercept: (8.3634, 18.3656)

95% Confidence Limits for PointsFor: (-0.02831, 0.00169)

95% Confidence Limits for PointsDifferential: (0.01774, 0.04750)

If PointsFor and PointsDifferential are both zero, the Patriots will win between 8.3634 and 18.36560 with a 95% confidence.

Holding PointsDifferential constant, for every 1 increase in PointsFor, the Patriots will win between -0.02831 and 0.00169 more games with a 95% confidence.

Holding PointsFor constant, for every 1 increase in PointsDifferential, the Patriots will win between 0.01774 and 0.04750 more games with a 95% confidence.

### Interpretation:

If the New England Patriots finish the regular season with a 0-point differential, they will win 13.3645 games, keeping PointsFor constant. Logically, you would expect that if a team scores the same points as their opponents did you would see an 8-8 record. However, the tests we ran indicate that even when the Patriots score the same amount as their opponent, in the long run they will have a record above .500 (13-3). We've attributed this to parts of football that were not captured in our dataset. These attributes include coaching strategy, defensive and offensive schemes, and player management; all of which do not show up in the raw, end of season statistics we studied.

### Model with Outliers Removed Interpretation:

We created a leave-one-out validation model removing the 3 outliers to see if it would result in a better model.<sup>[8]</sup>

#### Leave-One-Out Cross Validation w/ Outliers Removed:

The GLMSELECT Procedure Selected Model				
The selected model is the model at the last step (Step 1).				
Effects: Intercept PointsDifferential				
Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	1	28.05032	28.05032	29.32
Error	13	11.54968	0.88844	
Corrected Total	14	37.60000		
Root MSE				
Dependent Mean				0.94257
R-Square				0.6692
Adj R-Sq				0.6692
AIC				17.07912
AICC				19.28094
PRESS				14.45028
SBC				1.49522
Parameter Estimates				
Parameter	DF	Estimate	Standard Error	t Value
Intercept	1	8.925914	0.680511	13.01
PointsDifferential	1	0.021998	0.004007	5.41

Removing the 3 outliers resulted in a model with a R-square of 0.6692 which is in fact lower than the R-square for the model with the 3 outliers left in (0.6721). Therefore, we concluded that it would be better to include the 3 outliers in our model. Because of the low number of observations, it is best to have as many observations as possible.

### Final Conclusion from the Analyses of Objective 1:

Against our expectation, it was determined that the regular season record was not directly affected by the individual performance of Tom Brady as Tom Brady's individual statistics were not used to build the strongest model. Bill Belichick is the real goat, not Tom Brady.

## **OBJECTIVE 2:**

### **Problem:**

Is Tom Brady's performance directly responsible for the New England Patriots' success? After exploring what regular season stats contributed to the Patriots' success over the last 18 seasons, we sought to discover if there was any difference in Tom Brady's passer rating between the New England Patriots' wins and losses over three stages of Tom Brady's career – early, mid, and late.

### **Overall Approach:**

Since we are using a two-way ANOVA, we decided to bin every playoff and Super Bowl game by its end results (win or loss) and his career stage, identified by the quarterback's age between 3 bins (24-30, 30-36, and 36-42). In these 3 different stages, we are trying to determine if Tom Brady's passer rating was higher in games that the Patriots won versus the games that the Patriots lost.

### **Determining Predictors:**

We decided to use Brady's passer rating as opposed to individual passing statistics such as touchdowns, yards, and interceptions because the passer rating statistic encompasses those metrics.

### **Two-Way ANOVA Interpretation:**

In our model  $\text{Rate} = \text{AgeBin} + \text{WonLost}$ , the p-value for AgeBin is 0.77, which indicates that there is no difference in means between levels of Tom Brady's age. The p-value for WonLost is 0.0578, although that is on the cusp of our 95% confidence level, we have decided to consider this variable significant and run a one-way ANOVA test with the lone variable as WonLost.<sup>[9]</sup>

### **One-Way ANOVA Interpretation:**

In our model  $\text{Rate} = \text{WonLost}$ , the p-value for the WonLost variable drops a little further to 0.0561, further indicating that there is a significant difference between the mean of Tom Brady's passer rating between the Patriots' wins and losses in playoff and Super Bowl matches.<sup>[10]</sup>

### **Conclusion:**

As we can see from the interaction plot<sup>[9]</sup>, there was a slight increase in Brady's passer rating if we compare his late and mid-career passer rating to his early career passer rating. The comparison between his wins versus losses passer rating is more evident. Additionally, the parallel quality of the lines in the interaction plot indicate that there is no interaction between AgeBin and WonLost. The Q-Q plot indicates a normal distribution and our residual plot implies normality and equal variance among observations.

Section 10 of the appendix has some graphics regarding the one-way ANOVA performed after age is taken out of the dataset. The box and whisker plot show the evidence that Tom Brady's performance does have a significant effect on their results in playoff games. This kind of strays from the evidence we saw in part one of our analysis which focused on the Patriots and Tom Brady's regular season performance. We were unable to find evidence that Tom Brady's individual performance had an effect in the regular season, but this seems to shift when the Patriots get to the post season. In conclusion, the Patriots have needed Tom Brady to play at his best in playoff games for them to come out on top in February.

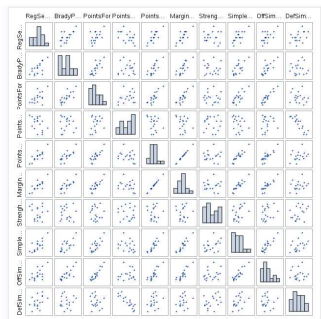
APPENDIX:

```
[1]
PROC IMPORT OUT= WORK.pats
    DATAFILE= "/home/daveknockwin0/PatriotsYearlyStats.csv"
    DBMS=CSV REPLACE;
    GETNAMES=YES;
    DATAROW=2;
RUN;

/* Print dataset */
proc print data=pats;
run;
```

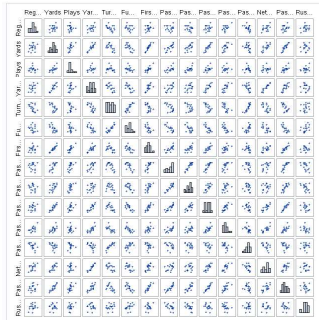
Obs	Year	RegSeasonWins	BradyPasserRating	PointsFor	PointsAgainst	PointsDifferential	MarginOfVictory	StrenghOfSchedule	SimpleRatingSystem	OffSimpleRatingSys	DefSimpleRatingSys	Yards	YardsPerPlay	Turnovers	FumblesLost	FumblesRecovered	PassCompletions	PassAttempts	
1	2002	9	88.7	301	246	55	2.2	1.8	4.3	1.2	0.1	4002	1001	4.9	26	12	292	376	402
2	2003	14	88.9	248	230	118	6.9	0.1	6.9	2.1	4.9	1038	1042	4.9	24	11	204	220	527
3	2004	14	82.9	427	280	147	11.1	1.8	12.9	8.4	6.9	1722	1038	9.9	27	13	344	380	485
4	2005	10	82.9	378	330	48	2.9	9.9	3.1	3.1	20.9	882	1031	9.9	29	8	318	382	394
5	2006	12	87.9	308	237	148	9.2	1	10.2	4.3	9.9	1089	1088	9.1	27	19	330	329	527
6	2007	16	107.2	508	274	234	19.7	9.9	20.7	19.9	4.3	882	1088	9.2	12	8	380	453	589
7	2008	11	83.9	412	339	101	9.2	-2.4	3.9	2.3	1.9	1047	1088	9.3	21	10	388	339	534
8	2009	10	88.9	407	288	119	9.9	2.3	11.2	9.7	4.9	1017	1078	9.9	20	9	373	386	482
9	2010	14	101	919	319	209	12.9	2.9	10.9	12.9	2.9	882	999	9.9	10	8	378	321	527
10	2011	13	108.9	513	342	171	10.7	-1.4	9.3	9.4	-0.1	1048	1082	9.3	17	8	399	402	512
11	2012	12	86.7	307	327	289	14.1	-1.4	12.9	12.2	2.9	1048	1199	9.7	19	7	444	462	541
12	2013	12	87.9	494	338	156	6.9	-0.7	9.9	4.9	1.4	1192	1199	9.4	20	9	379	380	629
13	2014	12	97.9	499	319	189	9.7	1.9	10.9	7.9	3.9	1047	1073	9.9	12	4	391	382	629
14	2015	12	102.2	499	319	180	9.4	-2.4	7	9.3	1.7	1091	1080	9.7	14	7	340	424	629
15	2016	10	112.2	441	280	161	11.9	-2.7	9.3	4.3	0	1192	1088	9.9	11	4	351	388	580
16	2017	12	102.9	493	280	182	10.1	-1.2	9.9	9.3	2.9	1027	1073	9.9	12	4	399	389	527
17	2018	11	97.7	438	329	111	9.9	-1.9	9.2	3.1	0.1	1099	1073	9.9	19	7	399	379	574

```
[2]
/* Scatterplot matrices to determine predictors */
proc sgscatter data=pats;
matrix RegSeasonWins BradyPasserRating PointsFor PointsAgainst PointsDifferential MarginOfVictory
StrenghOfSchedule
SimpleRatingSystem OffSimpleRatingSys DefSimpleRatingSys
/ diagonal=(histogram);
run;
```

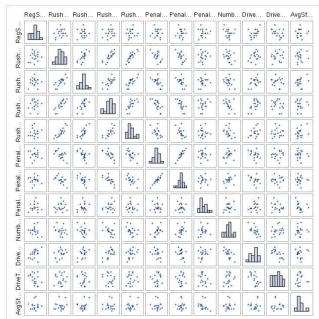




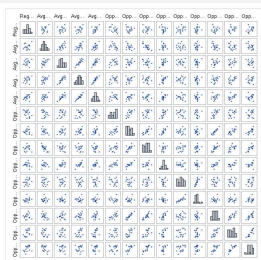
```
proc sgscatter data=pats;
matrix RegSeasonWins Yards Plays YardsPerPlay Turnovers FumblesLost FirstDowns PassCompletions
PassAttempts PassYards
PassTouchdowns PassInterceptions NetYardsPerPass PassFirstDowns RushAttempts
/ diagonal=(histogram);
run;
```



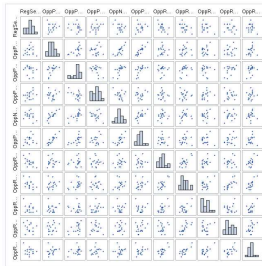
```
proc sgscatter data=pats;
matrix RegSeasonWins RushYards RushTouchdowns RushYardsPerAttempt RushFirstDowns Penalties
PenaltyYards PenaltyFirstDowns
NumberDrives DriveScorePercent DriveTurnoverPercent AvgStartingPosition
/ diagonal=(histogram);
run;
```



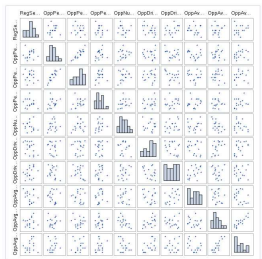
```
proc sgscatter data=pats;
matrix RegSeasonWins AvgDriveTime AvgDrivePlays AvgDriveYards AvgDrivePoints OppPointsFor
OppYards OppPlays OppYardsPerPlay
OppTurnovers OppFumblesLost OppFirstDowns OppPassCompletions OppPassAttempts
/ diagonal=(histogram);
run;
```



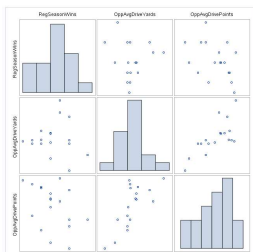
```
proc sgscatter data=pats;
matrix RegSeasonWins OppPassYards OppPassTouchdowns OppPassInterceptions OppNetYardsPerPass
OppPassFirstDowns OppRushAttempts
OppRushYards OppRushTouchdowns OppRushYardsPerAttempt OppRushFirstDowns
/ diagonal=(histogram);
run;
```



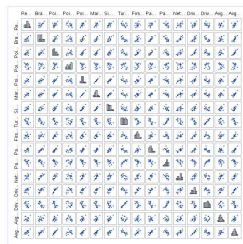
```
proc sgscatter data=pats;
matrix RegSeasonWins OppPenalties OppPenaltyYards OppPenaltyFirstDowns OppNumberDrives
OppDriveScorePercent OppDriveTurnoverPerent
OppAvgStartingPosition OppAvgDriveTime OppAvgDrivePlays
/ diagonal=(histogram);
run;
```



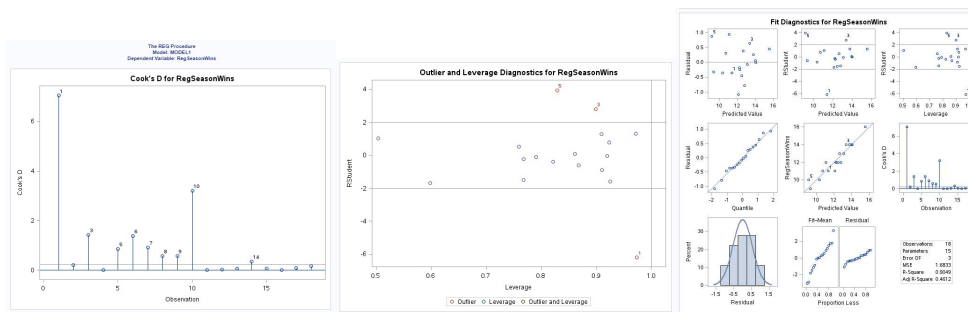
```
proc sgscatter data=pats;
matrix RegSeasonWins OppAvgDriveYards OppAvgDrivePoints/ diagonal=(histogram);
run;
```



```
/* Scatterplot matrix of the predictors */
proc sgscatter data=pats;
matrix RegSeasonWins BradyPasserRating PointsFor PointsAgainst PointsDifferential MarginOfVictory
SimpleRatingSystem Turnovers FirstDowns
PassTouchdowns PassInterceptions NetYardsPerPass DriveScorePercent DriveTurnoverPercent
AvgDriveTime AvgDrivePoints / diagonal=(histogram);
run;
```



```
/* Checking assumptions including outliers and leverage points */
proc reg data=pats plots(labels) = (rstudentleverage cooks);
model RegSeasonWins = BradyPasserRating PointsFor PointsAgainst PointsDifferential MarginOfVictory
SimpleRatingSystem Turnovers FirstDowns
PassTouchdowns PassInterceptions NetYardsPerPass DriveScorePercent DriveTurnoverPercent
AvgDriveTime AvgDrivePoints;
run; quit;
```

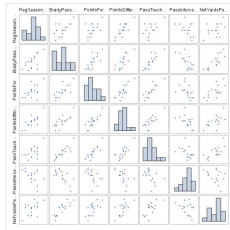


```
proc reg data=pats corr plots(label)=(rstudentleverage cooksdi);
model RegSeasonWins = BradyPasserRating PointsFor PointsAgainst PointsDifferential MarginOfVictory
SimpleRatingSystem Turnovers FirstDowns
PassTouchdowns PassInterceptions NetYardsPerPass DriveScorePercent DriveTurnoverPercent
AvgDriveTime AvgDrivePoints / VIF;
run; quit;
```

[illegible][illegible]

[5]

```
/* Scatterplot matrix of the predictors minus the low correlation and multicollinear ones */
proc sgscatter data=pats;
matrix RegSeasonWins BradyPasserRating PointsFor PointsDifferential
PassTouchdowns PassInterceptions NetYardsPerPass / diagonal=(histogram);
run;
```



[6]

```
/* LARS model */
proc GLMSELECT data=pats;
model RegSeasonWins = BradyPasserRating PointsFor PointsDifferential PassTouchdowns
PassInterceptions NetYardsPerPass / selection = LARS;
run; quit;
```

The GLMSELECT Procedure

Selected Model

The selected model is the model at the last step (Step 1).

Effects:	Intercept PointsDifferential
----------	------------------------------

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	1	33.35926	33.35926	27.02
Error	16	19.75185	1.23449	
Corrected Total	17	53.11111		

Root MSE	1.11108
Dependent Mean	12.22222
R-Square	0.6281
Adj R-Sq	0.6049
AIC	25.67178
AICC	27.38605
SBC	7.45250

Parameter Estimates		
Parameter	DF	Estimate
Intercept	1	9.529105
PointsDifferential	1	0.018327

```

/* LASSO model */
proc GLMSELECT data=pats;
model RegSeasonWins = BradyPasserRating PointsFor PointsDifferential PassTouchdowns
PassInterceptions NetYardsPerPass / selection = LASSO;
run; quit;

```

The GLMSELECT Procedure  
Selected Model

The selected model is the model at the last step (Step 1).

Effects: Intercept PointsDifferential

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	1	33.35626	33.35626	27.02
Error	16	19.75185	1.23449	
Corrected Total	17	53.11111		

Root MSE	1.11103
Dependent Mean	12.22222
R-Square	0.6281
Adj R-Sq	0.6049
AIC	28.67178
AICC	27.38805
SBC	7.45250

Parameter Estimates		
Parameter	DF	Estimate
Intercept	1	9.529105
PointsDifferential	1	0.016327

```

/* Stepwise model */
proc GLMSELECT data=pats;
model RegSeasonWins = BradyPasserRating PointsFor PointsDifferential PassTouchdowns
PassInterceptions NetYardsPerPass / selection = stepwise;
run; quit;

```

The GLMSELECT Procedure  
Selected Model

The selected model is the model at the last step (Step 2).

Effects: Intercept PointsFor PointsDifferential

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	2	37.74542	18.87271	18.42
Error	15	15.36589	1.02438	
Corrected Total	17	53.11111		

Root MSE	1.01212
Dependent Mean	12.22222
R-Square	0.7107
Adj R-Sq	0.6721
AIC	23.15178
AICC	20.22870
SBC	9.62239

Parameter Estimates				
Parameter	DF	Estimate	Standard Error	t Value
Intercept	1	13.354501	2.348353	5.70
PointsFor	1	-0.013311	0.007036	-1.89
PointsDifferential	1	0.032623	0.005981	4.67

```

/* Partial model with PointsFor and PointsDifferential */
proc reg data=pats;
model RegSeasonWins = PointsFor PointsDifferential /partial;
run;

```

The REG Procedure  
Model: MODEL1  
Dependent Variable: RegSeasonWins

Number of Observations Read 18  
Number of Observations Used 18

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	2	37.74542	18.87271	18.42
Error	15	15.36589	1.02438	
Corrected Total	17	53.11111		

Root MSE	1.01212	R-Square	0.7107
Dependent Mean	12.22222	Adj R-Sq	0.6721
Coeff Var	8.28095		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	13.35450	2.34835	5.70	<.0001
PointsFor	1	-0.01331	0.00704	-1.89	0.0780
PointsDifferential	1	0.03262	0.00598	4.67	0.0003

```
/* Partial model with PointsDifferential */
proc reg data=pats;
model RegSeasonWins = PointsDifferential /partial;
run;
```

The REG Procedure  
Model MODEL1  
Dependent Variable: RegSeasonWins

Number of Observations Read	18
Number of Observations Used	18

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	34.07901	34.07901	28.05	<.0001
Error	16	19.03211	1.18951		
Corrected Total	17	53.11111			

Root MSE	1.09055	R-Square	0.6417
Dependent Mean	12.22222	Adj R-Sq	0.6193
Coeff Var	8.92345		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	1	9.07117	0.64238	14.12	<.0001
PointsDifferential	1	0.02144	0.00401	5.35	<.0001

```
/* Leave-one-out validation model */
proc GLMSELECT data=pats;
model RegSeasonWins = BradyPasserRating PointsFor PointsDifferential PassTouchdowns
PassInterceptions NetYardsPerPass / selection=forward(STOP=Press);
run;
```

The GLMSELECT Procedure  
Selected Model

The selected model is the model at the last step (Step 2).

Effects: Intercept PointsFor PointsDifferential

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	2	37.74542	18.87271	18.42
Error	16	15.36569	0.96035	
Corrected Total	17	53.11111		

Root MSE	0.98000
Dependent Mean	12.22222
R-Square	0.7107
Adj R-Sq	0.6721
AIC	23.15178
AICC	26.22870
PRESS	22.87504
SBC	9.82359

Parameter Estimates				
Parameter	DF	Estimate	Standard Error	t Value
Intercept	1	13.38450	2.34035	5.70
PointsFor	1	-0.01331	0.00704	-1.89
PointsDifferential	1	0.03262	0.00698	4.67

```
[7]
/* Confidence intervals for the final model */
proc reg data=pats;
model RegSeasonWins = PointsFor PointsDifferential / clb;
run;
```

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t	95% Confidence Limits
Intercept	1	13.38450	2.34035	5.70	<.0001	8.38337 18.38563
PointsFor	1	-0.01331	0.00704	-1.89	0.0780	-0.02831 0.00169
PointsDifferential	1	0.03262	0.00698	4.67	0.0003	0.01774 0.04750

[8]

```
/* New dataset without outliers */
data pats2;
set pats;
if _n_=1 then delete;
if _n_=2 then delete;
if _n_=3 then delete;
run;

/* Model for K-fold cross validation (leave-one-out) on the new dataset*/
proc GLMSELECT data=pats2;
model RegSeasonWins = BradyPasserRating PointsFor PointsDifferential PassTouchdowns
PassInterceptions NetYardsPerPass / selection=forward(STOP=Press);
run;
```

The GLMSELECT Procedure  
Selected Model

The selected model is based on the best step (Step 1).

Effects Intercept PointsDifferential

Analysis of Variance				
Source	DF	Sum of Squares	Mean Square	F Value
Model	1	28.3832	28.3832	28.32
Error	35	13.5668	0.3876	
Corrected Total	36	41.9500		

Parameter Estimates				
Parameter	DF	Estimate	Standard Error	F Value
Intercept	1	0.00000	0.00000	1.00
PointsDifferential	1	0.00000	0.00000	0.00

[9]

```
proc anova data=work.import;
class AgeBin WonLost;
model Rate = AgeBin WonLost;
run;
```

```
proc glm data=work.import plots=all;
class AgeBin WonLost;
model Rate = AgeBin WonLost / clm;
run;
```



[10]

```
proc anova data = work.import;  
class WonLost;  
model Rate = WonLost;  
run;
```

```
proc glm data = work.import;  
class WonLost;  
model Rate = WonLost;  
run;
```

