Why the Packers struggled offensively in the second half of last season?

CS 170A Final Project Report

 Name:
 Jerry Liu

 ID:
 404474229

December 12, 2016

Contents

1	Intr	oduction	1
2	Pos	sible Factors	1
3	The	e NFL Play-by-Play Dataset	2
4	NF	L Savant Play-by-Play Dataset	3
5	Ana	alysis	5
	5.1	Principal Component Analysis	5
	5.2	Penalty	10
		5.2.1 Number of Penalties	10
		5.2.2 Penalties Yardage	10
	5.3	Yards Per Play	13
		5.3.1 Hypothesis Test $(t \text{ Test})$	15
		5.3.2 χ^2 Test	17
	5.4	Rushing Yards Per Play	18
		5.4.1 Hypothesis Test	21
	5.5	Sack	22

6 Conclusion

1 Introduction

As the most popular sports in the United States, NFL not only requires physicality from players on the field, but also prefers finely tuned strategies and schemes to put players on the position to win. In the past, NFL coaches only use their experience and feeling to determine the strategies and schemes against their opponents. As data analysis step in as an important factor, coaches can target opponent's weaknesses and maximize their chances to win. However, data still has not been seen as the most important measurement of performance and the most important factor to determine corresponding strategies and schemes. Analysts and scouts still refer to "eye test" to explain what a phenomenon is.

Since Green Bay Packers' embarrassing performance against Denver Broncos last year, quarterback Aaron Rodgers showed the worst performance since he became the starter for that franchise in 2008. In this report, we will use data to find out what leads to the Green Bay Packers' and Aaron Rodgers's miserable second half of the 2015 - 2016 season. Since no one has ever done any data analysis on Aaron Rodgers' and the Packers' performance, any method is unprecedented. I would like to use Principal Component Analysis to determine the factors that affect Packers' performance most. I would also like to find relevant attributes that may affect Aaron Rodgers' passing game.

2 Possible Factors

There are several possible reasons for their struggles. The most popular reason is: the receivers can't get open. People have the right to believe in this. Denver Broncos started to use press-man coverage against Packers receivers, which means the defensive backs put their hands on Packers receivers to prevent them from gaining separation. According to ESPN [1], the Packers have the League-wide slowest receiving corps. Also, Eddie Lacy, their primary running back, seemed to have overweight issues. There are also criticisms of the coaching staff, but we are unable to verify these claims with data.

3 The NFL Play-by-Play Dataset

This is a dataset available on kaggle [2]. This dataset is created by a group of Carnegie Mellon University statistical researchers, led by Maksim Horowitz [2].

The **46,129**-row and **63**-columns dataset, shown in Figure 1, contains all regular season plays from the 2015-2016 NFL season. Each play is broken down into fine details containing players involved, player positions, play results, penalties, etc. With a great amount of detail, we can use it to analyze causes for some phenomena.

	Α	В	С	D	E	F	G	н	1	J	к	L	м	N	0	Р	Q	R	S	Т	U V
1		Date	GameID	Drive	qtr	down	time	TimeUnder	TimeSecs	PlayTimeDiff Si	deofField	yrdin	yrdline100	ydstogo	ydsnet	GoalToGo	FirstDown	posteam	DefensiveTe	desc	PlayAttempt Yards.Gair
2	36	9/10/19	2015091000		1	1 NA	1	:00 1	5 3600	0 N	E	35	35	i C)	0	0 NA	PIT	NE	S.Gostkowsk	1
3	51	9/10/1	2015091000		1	1	1 1	:00 1	5 3600	0 PI	т	20	80	10	1	8	0 :	1 PIT	NE	(15:00) De.W	1
4	72	9/10/1	2015091000		1	1	1 1	:21 1	5 3561	. 39 PI	т	38	62	10	3	1	0 1	D PIT	NE	(14:21) B.Roe	1
5	101	9/10/1	5 2015091000		1	1	2 1	:04 1	5 3544	17 PI	т	47	53	1	. 3	1	0 :	1 PIT	NE	(14:04) De.W	1
6	122	9/10/1	5 2015091000		1	1	1 1	1:26 1	4 3506	i 38 N	E	49	49	10	4	5	0 :	1 PIT	NE	(13:26) B.Roo	1
7	159	9/10/1	5 2015091000		1	1	1 1	:42 1	3 3462	44 N	E	35	35	5 10	5	6	0 :	1 PIT	NE	(12:42) (Shot	1
8	180	9/10/1	5 2015091000		1	1	1 1	:05 1	3 3425	37 N	E	24	24	10	4	8	0 1	D PIT	NE	(12:05) A.Bro	1
9	199	9/10/1	5 2015091000		1	1	2 1	:20 1	2 3380	45 N	E	32	32	18	5	4	0 1	D PIT	NE	(11:20) (Shot	1
10	236	9/10/1	2015091000		1	1	2 1	1:53 1	1 3353	27 N	E	42	42	28	5	4	0 1) PIT	NE	(10:53) W.Jo	1
11	261	9/10/1	2015091000		1	1	3 1	1:28	1 3328	25 N	E	36	36	5 22	5	4	0 1	D PIT	NE	(10:28) (Shot	1
12	285	9/10/1	2015091000		1	1	4 0	:44 1	0 3284	44 N	E	26	26	5 12	5	4	0 :	1 PIT	NE	(9:44) J.Scob	1
13	305	9/10/1	2015091000		2	1	1 0	:40 1	0 3280	4 N	E	34	66	5 10	1	0	0 :	1 NE	PIT	(9:40) (Shotg	1
14	346	9/10/1	5 2015091000		2	1	1 0	:14 1	0 3254	26 N	E	32	68	3 10	1	0	0 :	1 NE	PIT	(9:14) (Shotg	1
15	371	9/10/1	5 2015091000		2	1	1 0	00:00	9 3240	14 N	E	44	56	5 10	2	3	0 :	1 NE	PIT	(9:00) (No Hu	1
16	396	9/10/1	5 2015091000		2	1	1 0	:31	9 3211	. 29 PI	т	43	43	10	2	3	0 1	D NE	PIT	(8:31) (No Hu	1
17	418	9/10/1	5 2015091000		2	1	2 0	1:27	9 3207	4 PI	т	43	43	10	2	3	0 1	D NE	PIT	(8:27) T.Brad	1
18	440	9/10/1	2015091000		2	1	3 0	1:22	9 3202	5 PI	т	43	43	10	2	3	0 1	D NE	PIT	(8:22) (Shotg	1
19	460	9/10/1	2015091000		2	1	4 0	:48	B 3168	34 PI	т	43	43	10	2	3	0 :	1 NE	PIT	(7:48) R.Aller	1
20	479	9/10/1	2015091000		3	1	1 0	:41	B 3161	. 7 PI	т	7	93	10)	6	0 1	D PIT	NE	(7:41) De.Wi	1
21	500	9/10/1	2015091000		3	1	2 0	:07	8 3127	34 PI	т	13	87	4		5	0 1	D PIT	NE	(7:07) De.Wi	1
22	521	9/10/1	2015091000		3	1	3 0	:26	7 3086	i 41 Pi	т	12	88	8 5	1	5	0 :	1 PIT	NE	(6:26) (Shotg	1
23	550	9/10/1	5 2015091000		3	1	1 0	:54	6 3054	32 PI	т	22	78	3 10	1	0	0 :	1 PIT	NE	(5:54) De.Wi	1
24	582	9/10/1	5 2015091000		3	1	1 0	:29	6 3029	25 PI	т	17	83	15	1	3	0 1	D PIT	NE	(5:29) (Shotg	1
25	606	9/10/1	5 2015091000		3	1	2 0	1:48	5 2988	41 PI	т	20	80	12		7	0 1	D PIT	NE	(4:48) B.Roet	1
26	625	9/10/1	5 2015091000		3	1	3 0	1:03	5 2943	45 PI	т	14	86	5 18	2	4	0 1	D PIT	NE	(4:03) (Shotg	1
27	650	9/10/1	2015091000		3	1	4 0	:25	4 2905	38 PI	т	31	69) 1	. 2	4	0 :	1 PIT	NE	(3:25) J.Berry	1
28	686	9/10/1	2015091000		4	1	1 0	1:14	4 2894	11 N	E	10	90	10)	8	0 1	D NE	PIT	(3:14) D.Lew	1
29	703	9/10/1	2015091000		4	1	2 0	::40	3 2860	34 N	E	18	82	2 2	1	9	0 1	D NE	PIT	(2:40) D.Lew	1
30	728	9/10/1	2015091000		4	1	3 0	:05	3 2825	35 N	E	19	81	1	. 1	0	0 :	1 NE	PIT	(2:05) T.Brad	1
31	749	9/10/1	2015091000		4	1	1 0	:14	2 2774	51 N	E	20	80	10)	0	0 :	1 NE	PIT	(1:14) D.Lew	1
32	787	9/10/1	2015091000	1	4	1	1 0	1:45	1 2745	29 N	E	10	90	20	1	9	0 1	D NE	PIT	(:45) (Shotgu	1
33	812	9/10/1	2015091000		4	1	2 0	0:12	1 2712	33 N	E	19	81	. 11	. 1	8	0 1	D NE	PIT	(:12) (Shotgu	1
34	836	9/10/1	2015091000		4	1 NA	0	0:00	0 2700	12 N	E	19	19) (1	8	0 1	0	NA	END QUARTE	1
35	852	9/10/1	2015091000		4	2	3 1	:00 1	5 2700	0 N	E	28	72	2	2	7	0	1 NE	PIT	(15:00) J.Ede	1
36	873	9/10/1	2015091000		4	2	1 1	:30 1	5 2670	30 N	E	37	63	10	2	4	0 1	D NE	PIT	(14:30) B.Bol	1
37	898	9/10/1	2015091000		4	2	2 1	:55 1	4 2635	35 N	E	34	66	5 13	3	1	0 1	D NE	PIT	(13:55) T.Bra	1
38	923	9/10/1	2015091000		4	2	3 1	:21 1	4 2601	. 34 N	E	41	59) E	3	9	0 :	1 NE	PIT	(13:21) (Shot	1
	► N	LPlaybyPla	iy2015 +	+																	

Figure 1: A screenshot of the CMU *NFL Play-by-Play Dataset*. This dataset is not numeric oriented and not very friendly to data analysis.

The data can be processed as follows:

```
[NUM, TXT, RAW] = xlsread('NFLPlaybyPlay2015.xlsx');
1
  % Turn the cell into a string matrix
2
  RAWStr = string(RAW);
3
4
  % Offensive / Defensive Team information is in column 18, 19 respectively
\mathbf{5}
6 Off = RAWStr(:, 18);
  Def = RAWStr(:, 19);
7
  PackersOffIndex = find(Off == 'GB') - 1;
8
  PackersDefIndex = find(Def == 'GB') - 1;
9
10
```

```
11 Attr = RAWStr(2, :);
12 PackersOffNum = NUM(PackersOffIndex, 2 : 66);
13 PackersDeffNum = NUM(PackersDefIndex, 2 : 66);
```

Here the NUM matrix contains all numeric data; TXT cell contains all non-numeric data and the whole dataset is represented in the cell RAW.

4 NFL Savant Play-by-Play Dataset

However, after some arduous digging of the unfriendly kaggle dataset [2], I decided to move on to the more friendly *NFL Savant* dataset [3]. This dataset is available on nfl savant [3]. This dataset consists of **46278**-row and **45**-columns of data. Most of the values are numeric, and unlike the kaggle data in Section 3, logical values are expressed in 0 / 1, which is much easier to analyze. The dataset is shown in Figure 2.

Α	В	С	D	E	F	G	н	1	J	k	L	м	N	0	Р	Q	R	S	Т	U	V
Gameld	GameDate	Quarter	Minute	Second	OffenseTear	r DefenseTear	Down	ToGo	YardLine		SeriesFirstDo	wn	NextScore	Description	TeamWin			SeasonYear	Yards	Formatio	n PlayType
*****	9/10/15		2	2	0	PIT		0	0	0	1			TWO-MINUT)		2015		0 UNDER C	ENTER
2015091300	9/13/15		3	4 5	0	GB		0	0	0	1		1	D TIMEOUT AT)		2015		0 UNDER C	ENTER
2015091300	9/13/15		4	0	0	GB		0	0	0	1			D END GAME	()		2015		0 UNDER C	ENTER
2015091301	9/13/15		2	2	0	SEA		0	0	0	1			0 TWO-MINUT	()		2015		0 UNDER C	ENTER
2015091302	9/13/15		1	0	0	CAR		0	0	0	1			D END QUARTI	()		2015		0	
2015091303	9/13/15		2	9 4	8 WAS	MIA		2	2	20	1			0 (9:48) 46-A.M	()		2015		4 UNDER C	ENT RUSH
2015091303	9/13/15		2	9 1	3 WAS	MIA		1	10	24	1			0 (9:13) 8-K.CO)		2015		22 UNDER C	ENT PASS
2015091303	9/13/15		4	2 3	3 WAS	MIA		2	11	76	0			0 (2:33) (SHOT)		2015		4 SHOTGU	N PASS
2015091304	9/13/15		1	4 5	7 BUF	IND		2	20	37	0			0 (4:57) 5-T.TA)		2015		0 UNDER C	ENT PASS
2015091304	9/13/15		2 1	2 2	4	IND		0	0	0	1			D TIMEOUT #1	(1		2015		0 UNDER C	ENT TIMEOUT
2015091304	9/13/15		3 1	4	4 BUF	IND		1	10	68	0			0 (14:04) 5-T.T)		2015		3 UNDER C	ENT PASS
2015091304	9/13/15		3 1	.3 2	7 BUF	IND		2	7	71	0			0 (13:27) 25-L.	()		2015		-1 UNDER C	ENT RUSH
2015091304	9/13/15		3	0	0	IND		0	0	0	1			D END QUARTI	()		2015		0	
2015091304	9/13/15		4 1	.3 3	3 BUF	IND		4	17	13	0			0 (13:33) 6-C.S)		2015		0 PUNT	PUNT
2015091304	9/13/15		4	2	4	IND		0	0	0	1			D TIMEOUT #2	()		2015		0 UNDER C	ENT TIMEOUT
2015091305	9/13/15		1	8 2	0 CLE	NYJ		3	1	28	1			0 (8:20) (SHOT	. ()		2015		11 SHOTGU	N PASS
2015091305	9/13/15		1	7 2	6 CLE	NYJ		2	9	40	0			0 (7:26) 34-I.C)		2015		8 UNDER C	ENT RUSH
2015091305	9/13/15		2	2	0	CLE		0	0	0	1			D TWO-MINUT	()		2015		0 UNDER C	ENTER
2015091305	9/13/15		2	2	0 NYJ	CLE		2	2	47	1			0 (2:00) (SHOT	()		2015		11 SHOTGU	N PASS
2015091305	9/13/15		2	0	0	CLE		0	0	0	1			D END QUARTI	()		2015		0	
2015091305	9/13/15		3 1	0 5	2 CLE	NYJ		3	7	37	1			0 (10:52) (SHC	()		2015		18 SHOTGU	N PASS
2015091306	9/13/15		1	0	0	кс		0	0	0	1			D END QUARTI)		2015		0	
2015091306	9/13/15		3	4 4	9 KC	HOU		1	10	35	0			0 (4:49) 11-A.S)		2015		0 UNDER C	ENT PASS
2015091306	9/13/15		3	4 4	2 KC	HOU		2	10	35	0			0 (4:42) (SHOT)		2015		2 SHOTGU	N PASS
2015091306	9/13/15		4	3 2	6 KC	HOU		2	6	24	0			0 (3:26) 11-A.S)		2015		0 UNDER C	ENT PASS
2015091306	9/13/15		4	1 2	2 KC	HOU		1	10	50	0			0 (1:22) 11-A.S)		2015		0 UNDER C	ENT OB KNEEL
2015091307	9/13/15		1	7 2	0 NO	ARI		2	1	41	0			0 (7:20) 29-K.F)		2015		0 UNDER C	ENT RUSH
2015091307	9/13/15		3	9	2 NO	ARI		1	10	43	0			0 (9:02) 29-K.F		,		2015		8 UNDER C	ENT RUSH
2015091308	9/13/15		1	0	0	DET		0	0	0	1			END OUART)		2015		0	
2015091308	9/13/15		2	5	2 SD	DET		0	0	35	1			2-LLAMBO K				2015		0 UNDER C	ENT KICK OFF
2015091308	9/13/15		4 1	4 2	8 SD	DET		2	7	50	0			0 (14:28) (SHO)		2015		-3 SHOTGU	N PASS
2015091309	9/13/15		2	2	0	BAI		0	0	0	1			TWO-MINUT)		2015		0 UNDER C	ENTER
2015091309	9/13/15		4 1	4 1	2 BAI	DEN		2	10	35	0			0 (14:12) 5-LE)		2015		0 UNDER C	ENT PASS
2015091309	9/13/15		4	3	0	BAI		0	0	0	1			TIMEOUT #1		1		2015		0 UNDER C	ENT TIMEOUT
2015091310	9/13/15		2 1	3 4	1 04K	CIN		4	6	47	0			0 (13:41) 7-M		1		2015		0 PUNT	PLINT
2015091310	9/13/15		3	7 5	0 CIN	OAK		1	10	15	0			0 (7.50) 14-4 [,		2015		8 LINDER C	ENT PASS
2015091310	9/13/15		3	0 4	7 046	CIN		2	6	32	1			0 (-47) (NO HI		1		2015		8 NO HUDE	DIE PASS
▶ ph	p-2015	+	-					-						, (110 110				LOID			

Figure 2: A screenshot of the *Play-by-Play 2015 Dataset*. This dataset is much more numeric oriented and more friendly to data analysis.

It consists of not only numeric data, but also textual data (i.e. strings). I first converted the csv file into an xlsx file, and then used xlsread to get the raw data. The code is shown on the next page:

```
1 [NUM, TXT, RAW] = xlsread('pbp-2015.xlsx');
2 RAW = string(RAW);
3 TXT = string(TXT);
4
5 % Attributes, ["GameId" "GameDate" "Quarter" "Minute" ...]
6 Attr = RAW(1, :);
7
8 % Get the col numbers that represent teams
9 OffCol = Attr == 'OffenseTeam';
10 DefCol = Attr == 'DefenseTeam';
11 OffIndex = find(RAW(:, OffCol) == 'GB') - 1;
12 DefIndex = find(RAW(:, DefCol) == 'GB') - 1;
13
14 % Filter Data
15 Off = NUM(OffIndex, :);
16 OffRaw = RAW(OffIndex + 1, :);
17 DenverGameIndex = find(OffRaw(:, DefCol) == 'DEN', 1);
18 offLength = length(Off);
19 Def = NUM(DefIndex, :);
20
21 % Get the games, in sequence
22 Games = process_game(RAW(OffIndex + 1, DefCol));
```

The source code of the function process_game is below:

```
1 % The source code of function process_game:
  function Games = process_game(GameLog)
2
       % Allocate Storage
3
       Games = GameLog(1 : 16);
4
5
       prevGame = GameLog(1);
       j = 2;
6
       for i = 2 : length(GameLog)
7
           if GameLog(i) ~= prevGame
8
9
               prevGame = GameLog(i);
               Games(j) = prevGame;
10
                j = j + 1;
11
12
           end
13
       end
14 end
```

5 Analysis

5.1 Principal Component Analysis

Since there are a lot of attributes for this dataset, we want to use *Principal Component Analysis* (PCA) to find out which attribute(s) have a significant impact on the whole data. First select relevant numeric attributes according to my football knowledge, shown in Table 1.

Attribute	Type	Meaning					
Yards	Numeric	Yards Per Play.					
Down	Numeric	The down count in the current drive.					
ToGo	Numeric	Yards needed to get a first down / score.					
YardLine	Numeric	The ball's snap position					
		(at n yard line) on the field.					
IsRush	Logical	Indicates whether the current play is a rush.					
IsPass	Logical	Indicates whether the current play is a pass.					
IsIncomplete	Logical	Indicates whether the pass is incomplete.					
IsTouchdown	Logical	Indicates whether the current play is a TD.					
IsSack	Logical	Indicates whether the passer is sacked.					
IsInterception	Logical	Indicates whether the passer is intercepted.					
IsFumble	Logical	Indicates whether the offense fumbles the ball.					
IsPenalty	Logical	Indicates whether a penalty flag is on the field.					
IsPenaltyAccepted	Logical	Indicates whether the penalty is accepted.					
IsNoPlay	Logical	Indicates whether the current play					
		is a no play (offsetting penalty).					

Table 1: Table of Projected Attributes and Corresponding Meanings

This forms a 1402×14 matrix and we want to use PCA to project it onto a 3D plane. We make the following distinction:

- If YPP(yards per play) < 5, then we categorize it as low YPP.
- If $5 \ge YPP < 15$, then we categorize it as *medium YPP*.
- If $YPP \ge 15$, then we categorize it as high YPP.

Then we use PCA (SVD algorithm) to extract the first three principal components. Figure 3 is a nice 3D plot from our projection and Figure 4 is an overview of the 3D plot.



Offense data projected on first three principal components

Figure 3: The PCA 3D Projection of our offense data

The code for projection is shown below:

```
% Select Relevant Attribute
1
\mathbf{2}
  IndexToProj = [20, 8 : 10, 23 : 26, 28, 33 : 35, 41, 43];
   % Yards, Down, ToGo, YardLine, IsRush, IsPass, IsIncomplete, IsTouchdown, IsSack, ...
3
       IsInterception, IsFumble, IsPenalty, IsPenaltyAccepted, IsNoPlay
  AttrToProj = Attr(:, IndexToProj);
4
  OffToProj = double(OffRaw(:, IndexToProj));
\mathbf{5}
6
   % Sort the data according to "Yards"
7
   [Val, I] = sort(OffToProj, 1, 'ascend');
8
   OffToProj = OffToProj(I(:, 1), :);
9
10
   % Divide Line between low, medium and high YPP
11
  LowYards = find(OffToProj(:, 1) == 5, 1) - 1;
12
  GoodYards = find(OffToProj(:, 1) == 15, 1) - 1;
13
14
```



Figure 4: Overview of the PCA 3D Projection of our offense data

```
% Normalize the data using correlation matrix
15
   OffToProjNorm = corr(OffToProj);
16
17
   % PCA
18
   [U,S,V] = svd(OffToProjNorm);
19
   n = size(OffToProj, 1);
20
21
22
   PrincipalComponent1 = U(:,1);
   PrincipalComponent2 = U(:,2);
23
   PrincipalComponent3 = U(:,3);
24
   SingularValues = diag(S(1:3,1:3));
25
26
   % Projection to 3D
27
   X = OffToProj * PrincipalComponent1;
28
   Y = OffToProj * PrincipalComponent2;
29
   Z = OffToProj * PrincipalComponent3;
30
31
   figure
32
  hold on
33
  plot3(X(1 : LowYards), Y(1 : LowYards), Z(1 : LowYards), 'g+')
34
```

```
35 plot3(X(LowYards : GoodYards), Y(LowYards : GoodYards), Z(LowYards : GoodYards), ...
       'b+' )
   plot3( X(GoodYards : n),
                              Y(GoodYards : n),
                                                   Z(GoodYards : n),
                                                                      'r+')
36
   xlabel('1st principal component (scaled)')
37
   ylabel('2nd principal component (scaled)')
38
   zlabel('3rd principal component (scaled)')
39
   title('Offense data projected on first three principal components')
40
   legend('low YPP', 'medium YPP', 'high YPP')
41
  rotate3d on
42
43 hold off
```

We can see a rough 3-cluster distribution from the graph, corresponding to low, medium and high YPP. With the first three Principal Components (PC), we want to find which attributes have significant weights in each of the three PC. The limit is quite arbitrary; for example, we can choose the attributes that exceed the mean of the absolute values of one PC vector. The values for all three vectors are shown in Table 2.

Vector Index	$1^{st} \mathbf{PC}$	2 nd PC	3 rd PC
1	-0.051190666	-0.211690262	0.603892049
2	0.143645804	-0.023413132	-0.117699317
3	0.051154443	-0.149140799	0.02702564
4	0.012284638	-0.0742752	0.114178674
5	-0.216100654	0.445428616	0.268642516
6	0.1825421	-0.634452804	0.011132822
7	0.220952173	-0.395082098	-0.302385616
8	0.010351095	-0.104089723	0.29517581
9	-0.000825734	0.17579242	-0.500473668
10	0.024150857	-0.073811184	0.048718148
11	-0.029826596	0.104427344	-0.273225556
12	0.530853587	0.183283386	0.115779429
13	0.536912381	0.21031385	0.102165359
14	0.52294893	0.184219641	0.074117185

Table 2: Table of the First Three PC Vector Values

The above manually chosen significant absolute values are typed in **boldface**, and their correspond-

ing emphasized attributes are shown in Table 3. This table also contains the emphasized values if we choose the mean of the sum of absolute values in the principal component vector as the criterion. For the first component, we can see that the first PC emphasizes values relevant to penalties, that

Method of Selection	$1^{st} \mathbf{PC}$	$2^{\mathbf{nd}} \mathbf{PC}$	3 rd PC
	IsPenalty	IsRush	Yards
Manual	IsPenaltyAccepted	IsPass	IsSack
	IsNoPlay	IsIncomplete	
	IsRush	IsRush	Yards
Mean of Absolute	IsPass	IsPass	IsRush
Values of the	IsIncomplete	IsIncomplete	IsIncomplete
PC Vectors	IsPenalty		IsTouchdown
	IsPenaltyAccepted		IsSack
	IsNoPlay		IsFumble

Table 3: Table of the Emphases of First Three PC Vector Values

the second PC emphasizes values relevant to play type and pass result (data related to passing will not be discussed in this report since they have already been analyzed by media), and that the third PC emphasizes yards per play and sack. This gives us some clue on analyzing the play-by-play data. The code of selecting values above the mean of absolute values of the PC vectors is shown below.

```
FirstEmpAttr = calc_attr(PrincipalComponent1, AttrToProj)
1
  SecondEmpAttr = calc_attr(PrincipalComponent2, AttrToProj)
2
  ThirdEmpAttr = calc_attr(PrincipalComponent3, AttrToProj)
3
4
  % calc_attr function implementation
\mathbf{5}
   function EmphasizedAttrs = calc_attr(PC, Attr)
6
       PC = abs(PC);
7
       divVal = mean(PC);
8
       EmphasizedAttrs = Attr(PC > divVal);
9
  end
10
```

5.2 Penalty

It is natural to take a look into the the emphasized category in the first PC: *penalty*. We first need to get the plays with *offensive penalties*, which is accomplished by the following code:

```
1 % Filter Data
2 PenaltyIndex = find(OffRaw(:, Attr == 'PenaltyTeam') == 'GB');
3 PenaltyRaw = OffRaw(PenaltyIndex, :);
  PenaltyYards = double(PenaltyRaw(:, Attr == 'PenaltyYards'));
4
5
  % Find out the Denver Game Index
6
  DenverGameIndex = find(PenaltyRaw(:, DefCol) == 'DEN', 1);
\overline{7}
  penLength = length(PenaltyYards);
8
9
  gameLength = length(Games);
10
  preDenLength = find(Games == 'DEN') - 1;
11
12 postDenLength = gameLength - preDenLength;
```

5.2.1 Number of Penalties

The number of penalties will definitely set the offense back a lot and it can "change the momentum". Then we need to find out whether our data suggests that there is a significant change in the number of penalties per game, which is accomplished by the code below.

```
1 PenaltyPreAvg = length(1 : DenverGameIndex - 1) / preDenLength
2 PenaltyPostAvg = length(DenverGameIndex : penLength) / postDenLength
```

It turns out that the number of offensive penalties per game is very consistent: 5 before the Denver game and 4.6 after the Denver game.

5.2.2 Penalties Yardage

Since the number of offensive penalties is very consistent, we need to take a look at the penalty yards, in that a 5-yard penalty is much easier to overcome than a serious penalties. The resulting plot is shown in Figure 5. The code is shown on the next page.



Figure 5: Average Penalty yards per play histograms.

```
1 % Penalty Yards Pre and Post Denver Game
2 PenaltyYardsPre = PenaltyYards(1 : DenverGameIndex - 1, :);
  PenaltyYardsPost = PenaltyYards(DenverGameIndex : penLength, :);
3
4
\mathbf{5}
   % Average Penalty Yards Per Game
   numberOfBins = 5;
6
  df = numberOfBins - 1;
7
   Edges = linspace(-2.5, 17.5, numberOfBins);
8
   [prePenCounts, Edges] = histcounts(PenaltyYardsPre, Edges);
9
   [postPenCounts, Edges] = histcounts(PenaltyYardsPost, Edges);
10
11
   % Get the average
12
   prePenCounts = prePenCounts / preDenLength;
13
   postPenCounts = postPenCounts / postDenLength;
14
15
16 figure
   subplot(2, 1, 1)
17
   histogram('BinEdges', Edges, 'BinCounts', prePenCounts)
18
   title('Average penalty yards in pre Denver games');
19
   subplot(2, 1, 2)
20
21 histogram('BinEdges', Edges, 'BinCounts', postPenCounts)
  title('Average penalty yards in post Denver games');
22
```

From the plot we can see an increase in the number of serious penalties: penalties of 10 and 15 yards. To see the difference better, please take a look at Figure 6. Serious penalties are much harder to overcome. For example, if on 3^{rd} down and 2, the Packers commit a 5-yard penalty, it will be 3^{rd} down and 8; 8 yards are not hard for Aaron Rodgers and company to overcome. However, if it is a 10 or even 15 yard penalty, considering how slow Packers' receivers are last season, one can easily imagine how those penalties can greatly affect the Packers' offense.



Figure 6: Difference between penalty yards between post and pre Denver game

The code for the difference is below:

```
Difference
1
  8
  diffX = linspace(0, 15, df);
2
  figure
3
  plot(diffX, postPenCounts - prePenCounts, 'bo')
4
  hold on
\mathbf{5}
  plot(diffX, zeros(df, 1), 'r--')
6
 title('Difference of average penalty yards between pre and post Denver games');
7
  axis([-5, 20, -1.4, 0.6]);
```

5.3 Yards Per Play

As the dominant value in the third PC, we have to look into the value: yards per play. First, we want to take a look at histograms of yards per play in a game. We split the data into two parts: pre-Denver games and post-Denver games. Since there are a lot more games after the Denver game, we need to use the average of yards per play.

```
1 OffYards = Off(:, Attr == 'Yards');
  gameLength = length(Games);
2
  preDenLength = find(Games == 'DEN') - 1;
3
  postDenLength = gameLength - preDenLength;
4
5 OffPre = OffYards(1 : DenverGameIndex - 1, :);
  OffPost = OffYards(DenverGameIndex : offLength, :);
6
   Edges = linspace (-12, 66, 27);
7
8
  % Use histcount function to process counts
9
   [preCounts, Edges] = histcounts(OffPre, Edges);
10
   [postCounts, Edges] = histcounts(OffPost, Edges);
11
12
  % Get the average
13
  preCountsAvg = preCounts / preDenLength;
14
  postCountsAvg = postCounts / postDenLength;
15
```

From the above code, we get two histograms of yards per play in pre Denver games and post Denver games, shown in Figure 7. From the plot we can see that there are more negative plays, and more plays with little gain.

To confirm our feeling, we plotted the difference between the post and pre Denver games in Figure 8: postCounts - preCounts. From Figure 8, we can see that there is a sharp increase in plays with small gains and significant increase in negative plays. There are also noticeable decreases for plays beyond 10 yards.

We see that the offense is not as productive after the Denver game than before the Denver game.

```
1 % Difference
```

```
2 start = mean(Edges(1 : 2));
```

```
3 \text{ ending} = \text{mean}(\text{Edges}(26 : 27));
```

4 diffX = linspace(start, ending, 26);



Figure 7: Average yards per play histograms.



Figure 8: Difference between yards per play between post and pre Denver game

```
5 plot(diffX, postCounts - preCounts, 'bo')
```

6 hold on

```
7 plot(diffX, zeros(26, 1), 'r--')
```

8 title('Difference of average yards per play between pre and post Denver games');

5.3.1 Hypothesis Test (t Test)

To see if there is enough evidence for worse yards per play in post Denver games, we want to setup a hypothesis test with significance level $\alpha = 0.05$. Let subscript 1 denote yards per play in pre Denver games and let subscript 2 denote yards per play in post Denver games.

 $H_0: \mu_{\rm pre} = \mu_{\rm post}$ $H_a: \mu_{\rm pre} > \mu_{\rm post}$

$$\begin{split} t &= \frac{\overline{x_1} - \overline{x_2} - 0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} \\ p\text{-value} &= \mathbf{P}(t > t_{df,\alpha}), \end{split} \qquad \qquad df = \min(n_1, n_2) - 1, \ \alpha = 0.05 \end{split}$$

And the result is that we <u>reject</u> the null hypothesis, which means we have enough evidence that the offense before the Denver game is better than that after the Denver game. This is not surprising as we see Aaron Rodgers and company struggle on offense. Receivers are pressed at the line of scrimmage and they cannot gain separation from average-best corners, making Rodgers scramble for himself or dumping a short check-pass a running back. The data are consistent with our eye-test, but are there any other possible cause for the debacle of their offense? Can we find, if any, from the dataset?

Since we know and see that their passing game struggled, we would like to know what else contributed to the avalanche of their offense. First we want to see whether the running game has affected the outcome. There is no significant difference between the average rushing attempts per game, as the value before the Denver game is 28.3333 and that after the Denver game is 26.4000. A 2-attempt difference is completely unnoticeable in a game. Thus, we want to see whether the efficiency of rushing has declined. The average rushing yards per attempt is 4.7471 before the Denver game and 4.2197 after the Denver game. The average rushing yards per attempt dropped for about 0.53 yards, which is 11%. The code for the hypothesis test is shown below:

```
1 % Hypothesis Test
2 [mul, s1, n1] = calc_t_attr(OffPre);
3 [mu2, s2, n2] = calc_t_attr(OffPost);
4 df = min(n1, n2) - 1;
5
6 tValProb = cdf('T', (mul - mu2) / sqrt(s1^2 / n1 + s2^2 / n2), df, 'upper');
7 Significance = 0.05;
8
9 if tValProb < Significance
10 fprintf('Reject Null Hypothesis.\n');
11 else
12 fprintf('Do not reject Null Hypothesis.\n');
13 end</pre>
```

The calc_t_attr function implementation is shown below:

```
1 function [mu, s, n] = calc_t_attr(X)
2     mu = mean(X); s = std(X); n = length(X);
3 end
```

I originally planned to only use χ^2 Test, but in some bins the number of occurrences is 0. If we simply plug these ill-formatted histogram counts into the cdf function, it will only return NaN. Also, since we are comparing the means, t test (another form of hypothesis test) seems more natural. However, we can reduce the number of bins to avoid 0s in the histogram counts. Note: reducing the number of bins can lead to a poorly represented data. For yardages where a 5-yard play differs from a 10-yard play significantly, I do think it is better to use Hypothesis Test. The code is below:

```
% Chi-Square Test
1
  numberOfBins = 5;
2
   df = numberOfBins - 1;
3
   Edges = linspace(-20, 80, numberOfBins);
\mathbf{5}
   % Use histcount function to process counts
6
   [postCounts, Edges] = histcounts(OffPost, Edges);
\overline{7}
   [seasonCounts, Edges] = histcounts(OffYards, Edges);
9
   % Get the average
10
   postCounts = postCounts / preDenLength;
11
   seasonCounts = seasonCounts / gameLength;
12
13
   ChiSquareStatistic = sum((postCounts - seasonCounts) .^ 2 ./ seasonCounts)
14
   ChiSquareProbability = cdf('Chisquare', ChiSquareStatistic, df)
15
```

The resulting χ^2 statistic is really large: ≈ 46.5266 , and the corresponding χ^2 probability is ≈ 1 , which means this situation is almost impossible. Thus, the χ^2 Test agrees with our Hypothesis Test result: the offense performed worse in post-Denver games.

5.4 Rushing Yards Per Play

Like what we did in Section 5.3, we use histograms to find out what happened to the ground game. The following code plots histograms of average rushing yards per play from both pre and post Denver games, shown in Figure 9. From the graph, we actually see that there are more small yards



Figure 9: Average rushing yards per play histograms.

and less good gains before the Denver games. The average value before the Denver game might very well be boosted by some large yardage runs at the right hand side of the histogram. The plot of the difference in Figure 10. The code is shown below.

```
1 % Filter Rushing Data
2 OffYards = Off(:, Attr == 'Yards');
3 RushIndex = double(OffRaw(:, Attr == 'IsRush')) == 1;
4 RushRaw = OffRaw(RushIndex, :);
```



Figure 10: Difference between rushing yards per play between post, pre Denver game

```
5 RushLength = length(RushRaw);
   RushYards = OffYards(RushIndex, :);
6
   DenverGameIndex = find(RushRaw(:, DefCol) == 'DEN', 1);
\overline{7}
8
   % Find the first Denver game index
9
   gameLength = length(Games);
10
   preDenLength = find(Games == 'DEN') - 1;
11
   postDenLength = gameLength - preDenLength;
12
13
   % Divide Data into two parts
14
   RushPre = RushYards(1 : DenverGameIndex - 1, :);
15
   RushPost = RushYards(DenverGameIndex : RushLength, :);
16
17
   Edges = linspace(-10, 65, 25);
18
19
```

```
20 % Use histcount function to process counts
   [preCounts, Edges] = histcounts(RushPre, Edges);
21
   [postCounts, Edges] = histcounts(RushPost, Edges);
22
23
24
25
   % Get Average
   preCountsAvg = preCounts / preDenLength;
26
   postCountsAvg = postCounts / postDenLength;
27
28
   figure
29
30
   subplot(2, 1, 1)
31 histogram('BinEdges', Edges, 'BinCounts', preCountsAvg)
32 title('Average rushing yards per play in pre Denver games');
   subplot(2, 1, 2)
33
34 histogram('BinEdges', Edges, 'BinCounts', postCountsAvg)
   title('Average rushing yards per play in post Denver games');
35
36
   % Difference
37
   start = mean(Edges(1 : 2));
38
   ending = mean(Edges(24 : 25));
39
   diffX = linspace(start, ending, 24);
40
41 figure
   plot(diffX, postCountsAvg - preCountsAvg, 'bo')
42
43 hold on
44 plot(diffX, zeros(24, 1), 'r--')
45 title('Difference of average rushing yards per play between pre, post Denver games');
```

5.4.1 Hypothesis Test (t Test)

To see if there is enough evidence for worse rushing attacks in post Denver games, we want to setup a hypothesis test with significance level $\alpha = 0.05$. Let subscript 1 denote rush yards in pre Denver games and let subscript 2 denote rush yards in post Denver games.

 $H_0: \mu_{\rm pre} = \mu_{\rm post}$ $H_a: \mu_{\rm pre} > \mu_{\rm post}$

$$t = \frac{\overline{x_1} - \overline{x_2} - 0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

p-value = $\mathbf{P}(t > t_{df,\alpha}),$ $df = \min(n_1, n_2) - 1, \ \alpha = 0.05$

Thus we use the following code to do the test:

```
% Hypothesis Test
1
2 [mu1, s1, n1] = calc_t_attr(RushPre);
   [mu2, s2, n2] = calc_t_attr(RushPost);
3
   df = min(n1, n2) - 1;
4
\mathbf{5}
   tValProb = cdf('T', (mu1 - mu2) / sqrt(s1<sup>2</sup> / n1 + s2<sup>2</sup> / n2), df, 'upper');
6
   Significance = 0.05;
\overline{7}
8
   if tValProb < Significance</pre>
9
        fprintf('Reject Null Hypothesis.\n');
10
   else
11
        fprintf('Do not reject Null Hypothesis.\n');
12
   end
13
```

And the result is that we do <u>not</u> reject the null hypothesis, which means we do not have enough evidence to show that the ground game before the Denver game is better than that after the Denver game. Thus we can put much less faith on that "improving the ground game will boost the offense".

5.5 Sack

Here we will look into another dominant value in the third PC: *Sack*. Sack is somewhat represented in negative plays, but in critical game situations, a sack can be a game changer. Like what we did in Section 5.2, we filter the data first and calculate the number of sacks per game.

```
% Filter out sack data
  SackIndex = find(OffRaw(:, Attr == 'IsSack') == '1');
2
  SackRaw = OffRaw(SackIndex, :);
3
   SackYards = OffYards(SackIndex);
4
\mathbf{5}
   DenverGameIndex = find(SackRaw(:, DefCol) == 'DEN', 1);
   sackLength = length(SackYards);
6
7
   gameLength = length(Games);
8
   preDenLength = find(Games == 'DEN') - 1;
9
   postDenLength = gameLength - preDenLength;
10
11
  % The mean of the number of sacks in a game
12
  SackPreAvg = length(2 : DenverGameIndex) / preDenLength
13
   SackPostAvg = length(DenverGameIndex : sackLength) / postDenLength
14
```

We see that before the Denver game, the average number of sacks is only 2.1667, while after the Denver game, this number bumps $\approx 75.38\%$ to 3.8. An increase in sacks can indicate the following:

- Lack of ground game caused too many passing attempts, which naturally generated more sacks. This is already declined in Section 5.4.1.
- Lack of protection from offensive line. The Packers' offensive line was plagued by injuries in the second half of the season, with which the data are consistent.
- Aaron Rodgers held the ball too long. The protection will always break down as time-aftersnap increases, which will result in more sacks. NFL Analysts have already found the cause: receives can't get open.

Since there is nothing novel coming out of the *sack* data, I will not do more data analysis on this part.

6 Conclusion

Because of my familiarity with American football and NFL, I decided to analyze NFL play-by-play data. In my humble opinion, except *ProFootballFocus*, no other media really pay attention to data mining, and analysts either only refer to simple statistics or rely solely on their "eye-test". For example, for Green Bay Packers' struggle in the second half of last season, I heard the following rhetoric for quite a while:

- 1. "Aaron Rodgers is not the same."
- 2. "There is no running game."
- 3. "There is no protection from the offensive line."
- 4. "The receivers can't get open."

For 1. and 2., often analysts do not refer to statistics at all. They use their *eye-test* to say what they feel, which is not very reliable.

In my analysis, for example, I found that there is no sufficient evidence that the ground game was worse in the second half of last season. NFL Analysts also did not notice an increase in serious penalties in the second half of the season, as shown in Section section::penalty, which created difficult tasks for the Packers' offense.

In this project, converting csv files / cells into string matrices costs me some time. I looked through several MATLAB documents to properly import the play-by-play data [3] into MATLAB. After my Principal Component Analysis, I spent quite some time getting histograms work. Since I need to use the average of yards per play (YPP), or YPP per game, I need to modify the counts in the histogram counts. In this course, we only saw hist, and I could not find a way to manipulate this function to get the average of counts. I looked several documents and decided to use a combination of histcounts and histogram. Then I used χ^2 Test and Hypothesis Test (mainly the latter) to confirm my intuitions and found out that one of them is false. The hardest part of this project will be converting cells / char vectors into string matrices, and then manipulating and extracting data from the raw data string matrix. For the report, the hardest part is to find a proper arrangement so that the graphs and relevant paragraphs are not very far apart. In conclusion, I found out that

- The ground game was not worse in the second half of the season.
- Packers' had noticeably more "serious" penalties, penalties with 10 or 15 yards.
- The number of sacks go up quite a bit, consistent with Analysts' "eye-test".

Although the dataset used in this project is already very large, it also lacks some attributes required by more advanced analyses. For example, if the receivers really cannot get open in press coverage, can we confirm that from the data; e.g. is the yards per play / completion percentage significantly down from non-press coverages? If I have this kinds of data, I would assign numeric labels to replace their string representations and use correlation matrix to do PCA again to confirm whether they play a huge role in offense.

References

- DEMOVSKY, R. Passing problem: Aaron rodgers, packers have nfl's slowest weapons. http: //www.espn.com/, November 2015.
- [2] HOROWITZ, M. Detailed nfl play-by-play data 2015. https://www.kaggle.com, October 2016.
- [3] NFL-SAVANT. Nfl play-by-play data 2015. http://nflsavant.com/about.php, 2016.