

AN ECONOMETRIC ANALYSIS OF THE RELATIONSHIP BETWEEN CORNER  
KICK NUMBERS AND FOOTBALL OUTCOMES

by

Toyosi Ashimolowo

A thesis submitted to the faculty of  
The University of North Carolina at Charlotte  
in partial fulfillment of the requirements  
for the degree of Master of Science in  
Economics

Charlotte

2018

Approved by:

---

Dr. Craig A. Depken, II

---

Dr. Rob Roy McGregor

---

Dr. L. Ted Amato

©2018  
Toyosi Ashimolowo  
ALL RIGHTS RESERVED

## ABSTRACT

TOYOSI ASHIMOLOWO. An Econometrics Analysis of the Relationship Between Corner Kick Numbers and Football Outcomes. (Under the direction of DR. CRAIG A. DEPKEN, II)

Corner kicks are arguably important during match play, but questions remain about the impact of corner kicks on match outcomes. This study analyzes the relevance of corner kicks numbers and their importance to match outcomes, then compares this result with that of shots on target and free kicks. In total, 7368 matches played in the English Premier League, French Ligue 1, German Bundesliga, Italian Serie A, Spanish La Liga, and the 2018 FIFA World Cup are analyzed using the Multinomial logistic regression (MLR). The study finds that both the number of corner kicks and shots on target are significantly associated with the outcome of football games while free kick numbers are not. Also, from the results of the MLR analysis, the study finds that the winning team tends to play more shots accurately while the losing team tends to be awarded more corners during the game. The findings thereby support Riquelme's (2012) contention that match status affects the number of corners in a match and also Collert's (2012) claim that there is a significant relationship between shooting efficiency and overall team success. The findings also disagree with Gordon et al's (2013) argument that corner kicks numbers are statistically useless in soccer games.

## ACKNOWLEDGMENT

Firstly, I would like to thank my encyclopedia of econometrics, Dr. Craig A. Depken, who inspired me to work on this topic and as well provided me with valuable help and direction through my work. His guidance helped me in all the time of research and writing of this thesis. Beside my advisor, I would like to thank the rest of my thesis committee: Dr. L. Ted Amato and Dr. Rob McGregor for their insightful comments.

Finally, I would like to thank my family: my parents and to my brothers and sisters for supporting me spiritually throughout writing this paper and my life in general.

## TABLE OF CONTENTS

List of Tables .....	vi
CHAPTER 1: INTRODUCTION .....	1
CHAPTER 2: LITERATURE REVIEW .....	4
CHAPTER 3: METHODOLOGY .....	9
3.1 Multinomial Logit Regression .....	9
3.2 Model Assumptions .....	11
3.3 Methodological Procedures.....	12
3.3.1. Likelihood Ratio and Wald Test .....	12
3.3.2. Multicollinearity.....	13
3.3.3. Independence of Irrelevant Alternatives .....	14
3.3.4. Average Marginal Effects .....	16
3.3.5. Goodness-of-Fits.....	17
3.3.6. Model Comparison.....	18
CHAPTER 4: RESULTS AND INTERPRETATION .....	20
4.1 Data Description .....	20
4.2 Multinomial Logistic Regression Results .....	22
4.2.1 Multinomial Logit Coefficients.....	22
4.2.2 Multinomial Logit Average Marginal Effects.....	25
4.3 Methodological Procedures.....	27
4.3.1 Likelihood Ratio and Wald Test .....	27
4.3.2 Multicollinearity.....	28
4.3.3 Independence of Irrelevant Alternatives .....	29
4.3.4 Goodness-of-fits(Hosmer-Lemeshow Statistic).....	30
4.3.5 Model Comparison (McFadden R-square, AIC, BIC) .....	31

CHAPTER 5: CONCLUSION AND FUTURE RESEARCH.....	33
REFERENCES .....	35
APPENDIX A: SUPPLEMENTAL TABLE AND FIGURES .....	37
APPENDIX B: AVERAGE MARGINAL EFFECTS at 95% CONFIDENCE INTERVAL .....	38

## LIST OF TABLES

Table 1: Description of Variables.....	20
Table 2: Summary of Independent Variables .....	21
Table 3: Description of Variables cont.....	22
Table 4: Multinomial Logistic Coefficient Estimates.....	23
Table 5: Multinomial Logit Marginal Effects Estimates.....	26
Table 6: Likelihood Ratio (LR) AND Wald Test results.....	27
Table 7: Multicollinearity Test; VIF results .....	29
Table 8: Independence of Irrelevant Alternatives (IIA) results .....	30
Table 9: Hosmer-Lemeshow test results.....	31
Table 10: Mc-Fadden R-squared, AIC and BIC test results .....	32

## CHAPTER 1: INTRODUCTON

With over 380 million viewers, the 2014 Champions League final between Atletico Madrid and Real Madrid in Lisbon will forever be remembered for Sergio Ramos' 93<sup>rd</sup> minute equalizer. Luka Modric found the center back on a "corner kick", and the Los Blancos hero buried his header into the far corner, just past the reach of the goalkeeper; it was not the first goal of the game, nor was it the last, but it was the one that changed history. Corner kicks are not particularly valuable events in soccer but when a player is away from action near the corner flag, fans are on their feet, knowing that a goal could come at any moment. Scoring goals is the ultimate determinant of a successful team and in order to be efficient in front of goal, the art of goal scoring ranks at the top of every manager, player, and teams' agenda.

Football's worldwide popularity and support has resulted in varieties of research investigating different aspects within the game. A significant number of researchers have analyzed the pattern of play and possession in relation to goal scoring (Hook and Hughes, 2001; Hughes and Frank, 2005). Others (Sousa and Gargantua, 2001 and Armatus et al., 2007) have focused on the impact of set-plays on scoring opportunities, focusing on free kicks, corners, and penalties. Set plays are important, but do their numbers really influence the outcome of a game? In this study, I intend to show that football outcomes are not only influenced by how efficient a team is in front of goal but that the number of corner kicks in one way or the other is also a significant determinant in match outcomes.

Set plays in football are crucial, providing the opportunity for a "free" shot towards goal. It offers teams, or individuals, the opportunity to perform tactical routines. Whether it is through the use of team strategies for corners or the individual technique of a penalty,



the anticipation of a goal increases as they are deemed to be a goal scoring scenario (Lowe, 2015). A corner kick is a type of set-play awarded when the whole of the ball goes out of play over the goal line, without a goal being scored, and having last been touched by a member of the defending team. Corners provide attacking teams with the option of crossing the ball into an advantageous attacking position. For this reason, some authors have decided to analyze corner kicks in detail, with the majority of findings originating from tournaments such as the FIFA World Cup and the European Championship. Most of these studies carried out descriptive analyses of how corner kicks can be more effective and the effects of match status on corner kicks. There is little research on how the number of corner kicks taken affects match results, which is quite surprising as corner analysis across the duration of a season provides comprehensive data and the opportunity for detailed team comparisons.

### **1.1 Aim**

The aim of this study is to broaden our understanding of corner kicks and their relevance to match outcomes by employing comparative data across five different leagues (English Premier League, Italian Serie A, French Ligue 1, German Bundesliga, Spanish La Liga) within the period 2014 to 2018 as well as the 2018 FIFA World Cup. This study also makes an effort to compare the relevance of corner kicks to shots on targets and free kicks.

### **1.2 Research Questions**

Drawing on the literature reviewed, the study will aim to answer the following research questions;

- I. Does the relative number of corners a team plays affect the match outcome across the seven leagues during the period of 2014/2015 – 2017/2018 season and the 2018 World Cup?

- II. Does the relative number of shots a team plays on target affect the match outcome across the seven leagues during the period of 2014/2015 – 2017/2018 season and the 2018 World Cup ?
- III. Does the relative number of free kicks a team gets affect the match outcome across the seven leagues during the period of 2014/2015 – 2017/2018 season and the 2018 World Cup?

### **1.3 Hypothesis**

The null hypothesis ( $H_0$ ) states that the relative number of corner kicks is not significant towards match outcomes. The following hypotheses are also stated:

- The relative number of shot on targets has no significant impacts on match outcomes.
- The relative number of free kicks has no significant impacts on match outcomes.

### **1.4 Limitations**

The following limitations of the study have the potential to influence the results obtained;

- a) Knowing fully well that an individual corner kick will lead to a goal around 0.03 (3%) of the time, the study does not account for the efficiency of corner kicks.
- b) The study ignores the tactical approach of teams towards attacking or defending corner kicks (i.e. man to man marking, zonal marking, combined marking, or aerial threats).
- c) Even though the aim is to analyze the impact of corner kicks numbers on match outcome, the research does not account for when home total corners are equal to away total corners.

- d) Lastly, the results do not account for the playing style implemented by the teams during games.

Having affirmed the rationale for the current research topic, chapter two presents a discussion of the literature. Chapter three considers the methodology employed. Chapter four presents' findings from the tests conducted and the final chapter reinstates key findings, and general conclusion of the thesis.

## CHAPTER 2: LITERATURE REVIEW

Corner kicks play an important role in football, providing teams with opportunities to deploy offensive strategy onto the opposition's defense. They represent the chance to maintain possession in an attacking area of the field and to hit an unopposed pass into the opposition's penalty area (Hill and Hughes, 2001). Research has been conducted (Olsen and Larsen, 1997; Hill and Hughes, 2001; Taylor et al., 2005, and Sainz de Baranda and Riquelme, 2012) that has provided information into why and how the corner kick is an important tool. All these studies investigate the effectiveness of corners during tournaments or league campaigns with Casal et al. (2015) providing the most in-depth and recent corner kick research, exploring corner characteristics across three competitions (FIFA World Cup 2010, UEFA European Championships 2012, and UEFA Champions League 2010-2011).

The general consensus is that between nine and thirteen corner kicks on average are taken per match (Hill and Hughes, 2001; Yamanaka et al., 2002; Borrás and Sainz de Baranda 2005; Carling et al., 2005; Taylor et al., 2005; Acar et al., 2009; Sainz de Baranda, and Lopez Riquelme 2012, and Casal et al., 2015). As a result of this finding, Casal et al. (2015), who analyzed 1139 corner kicks across multiple competitions, states that only 26 per cent of these resulted in a shot. From the total shots, 9.8 per cent hit the target and only 2.2 per cent ended in a goal being scored. This finding is supported by Marquez and Raya (1998) who argue that corner kicks are efficient in 2.28 per cent of all cases and Perez and Vicente (1996) who analyze the 1994 World Cup in the USA and find 1.6 per cent of shots from corners result in a goal.

Furthermore, common trends have been discovered in terms of corner type and the location that produces the highest number of shots on goal and goals scored. For example,

Olsen and Larsen (1997), Taylor et al. (2004) and Taylor et al. (2005) argue that a direct, swung corner provides a team with the best chance to create a scoring opportunity, while Hill and Hughes (2001) also state that corners with curl provide the most opportunities on goal. Both Olsen and Larsen (1997) and Taylor et al. (2005) support each other's findings, stating that one in five out-swing and one in three in-swing corners result in a chance on goal. Taylor et al. (2005) further support these findings by arguing that out-swing corners lead to the highest number of attempts on goal (60.7 per cent), whilst in-swing corners lead to more goals (66 per cent). The nature of the out-swing corner means the ball is directed further away from the goal-keeper, increasing the chances of the attacking team making first contact with the ball. However, as a result, the chance of a goal being scored is reduced due to the difficulty of generating enough power on a shot as the ball is travelling away from the goal. The consensus that in-swing corners produce more goals does not come as a surprise, due to its proximity to the goal and the fact that more emphasis is put on shot direction rather than power, creating a greater chance of a goal.

On the other hand, the execution of the short corner has been found to produce numerous goal-scoring opportunities, although they are less common (Hill and Hughes, 2001; Taylor et al., 2005). This could be due to the fact that these types of corners are more infrequent, and defenders are less familiar with the angle of delivery. As a consequence, this causes confusion as defenders become attracted to the ball, creating greater spaces behind the defense (Ali, 1998). Multiple discrepancies have been highlighted across research investigating corner kick delivery location. Some researchers (Hughes and Petit, 2001; Taylor et al., 2005) declare that the area between the six-yard box and penalty spot is the optimal area of delivery, as this has led to the greatest number of attempts and goals.

Similarly, Taylor et al. (2005) state that in-swing corners into that same area provide the most attempts and the highest chance of a shot on target. Arguably, the corner is being delivered into the “goalkeepers’ area of uncertainty”, not knowing whether to stay on the goal line or come to collect the cross. Hughes (1999) also emphasize the use of in-swing will lead to a higher number of goals.

A variation in results was found in regard to match status and the target area of corners. Sainz de Baranda and Lopez-Riquelme (2012) discovered that there was a tendency for teams during the 2006 World Cup to target the front post regardless of the score-line (36.6 per cent), agreeing with Marquez and Raya (1998) who previously concluded similar results at the 1998 World Cup in France. An interesting finding by Ali (1998) is a trend in teams winning a match to perform more short corners. A discovery complemented by Sainz de Baranda and Lopez-Riquelme (2012) who find when winning 29.2 per cent were played short compared to 10.1 per cent when drawing and 15.2 per cent when losing.

In addition, a wealth of research examines whether goal scoring is affected by time, (Jishan et al., 1993; Michaildis et al., 2004; Yiannakos and Armatas, 2006; Armatas et al., 2007) yet multiple discrepancies have been found. Jishan et al. (1993) conclude in a study of the 1990 World Cup that most goals were scored in the final 15 minutes. Yiannakos and Armatas (2006) support this claim by investigating set-plays at the European 2004 Championships and report more goals were scored in the second half. However, Michaildis et al. (2004) report that time has no effect on goal scoring when analyzing how, where, and when goals were scored during the 2002-03 Champions League. A recent study by Armatas et al. (2007) analyzes the effect that time has on goal-scoring from three World Cups (1998,

2002 and 2006) and conclude more goals are scored as time progresses during the second half, agreeing with Jishan et al. (1993) and Yiannakos and Armatas (2006). Specifically, across the three World Cups analyzed, no significant differences were found between goals being scored in 15 intervals. Evidently, some reports have supported that time-scale affects goal-scoring, however this is an area that requires further development.

There is abundance of literature concerning corner kicks in football however most of this research relies heavily on tournaments such as the FIFA World Cup and the European Championships; the nature of these tournaments can result in misleading findings as corner samples are very small and few matches are played per team. Also, none of the previous research conducted fully assesses how the relative number of corner kicks affects football outcomes.

This study aims to fully assess whether or not corner kick numbers affect football outcomes using data across six different domestic leagues and the FIFA world cup. The analysis will use Stata 15.

## CHAPTER 3: METHODOLOGY

Since the dependent variable (football match result) has two or more outcome categories (i.e., win, draw or loss), the logit and probit models should be used to model the impact of corner kicks on football results. Logit models can either be binary or multinomial; binary models consider two response outcomes (i.e. win vs lose or draw vs lose) while the multinomial models consider three or more response outcomes. This chapter further describes the method used in the study, the assumptions of multinomial logit regression (MLR), and several methodological procedures that should be used in testing the assumptions of the MLR.

### 3.1. Multinomial Logit Regression

The multinomial logit regression (MLR) is an extension of the binary logistic regression with multiple predictors. It is used to model the relationship between a “polytomous” dependent variables (with more than two outcomes) and a set of independent variables. MLR compares multiple groups through a combination of binary logistic regressions which allows each category of the dependent variable to be compared to a reference category. The reference category, also known as the base category, serves as a contrast point for all analyses, and the effects of the analysis are always in reference to the contrast category.

The general form of the multinomial logit model is;

$$Prob[y_i = j] = \exp(\beta_j X_i) / \sum_j \exp(\beta_j X_i),$$
$$j = 0, 1, \dots, m,$$



where  $j$  denotes the specific one of the “ $m + 1$ ” possible unordered choices,  $y_i$  is the indicator variable of choices,  $X_i$  denotes the vector of the independent variables, and  $\beta_j$  is the corresponding coefficient vector.

The MLR has many advantages in modelling football outcomes, such as:

- The results can be interpreted by both the regression coefficient estimates and the exponentiated coefficients.
- The estimates are asymptotically consistent with requirements of the nonlinear regression.
- It produces valid estimates as it applies transformation of the multinomial dependent variable to a continuous variable ranging from negative infinity to positive infinity.

The dependent variable (match outcome) in this paper consists of three outcomes categories (i.e. home win, draw, home loss), and is assumed to be unordered. Since the MLR works by choosing one outcome category as the base (reference) category for the other categories, hence home win is considered as the reference group, because it is the most frequent outcome of football results and the other outcome levels are estimated relative to home win.

The dependent variable  $Y$ , outcome, will then take on three values: a home win (1), draw (2), and home loss (3). This analysis will then compare draw (2) relative to home win (1) and home loss (3) relative to home win (1). The two equations of the MLR model are then given by;

$$\text{Log} \left[ \frac{\text{Pr}(Y=2)}{\text{Pr}(Y=1)} \right] = \beta_{20} + \beta_{21}x_1 + \beta_{22}x_2 + \cdots \beta_{2p}x_p$$

$$\text{Log} \left[ \frac{\text{Pr}(Y=3)}{\text{Pr}(Y=1)} \right] = \beta_{30} + \beta_{31}x_1 + \beta_{32}x_2 + \cdots + \beta_{3p}x_p ,$$

where  $p$  denotes the number of predictors for the binary response,  $Y$ , by  $x_1, x_2, \dots, x_p$ .

There are a few applications of the MLR in modelling match outcomes. For example, Christian Collet (2013) applied the ordered MLR to analyze the impact of possession on team success/ football outcomes. Jefferey Allan (2014) applied the multinomial logistic regression to predict the outcomes of 380 matches of the 2011/2012 Premier league season. Despite these few applications of the MNL, this paper seeks to introduce new variables in capturing the determinants of match outcome and new methods of presenting the results of the MLR applications that have not been reported in other football outcome research. This new method includes proper analysis of the assumption of the independence of irrelevant alternatives (IIA), which is very crucial in the MLR modeling; tests for multicollinearity among the independent variables; the Likelihood Ratio (LR) and Wald test are used to properly test for the effect of independent variables; the marginal effects of all independent variables upon the dependent variables are presented.

### **3.2. Model Assumptions**

To get a valid result using the MLR, the following assumptions needs to be met.

- a) The MLR assumes that the odds for any pair of outcomes are determined without referring to the other outcomes that might be available. This assumption is called the Independence of Irrelevant Alternatives (IIA) and it is very crucial in the modeling.
- b) There should be no multi collinearity. Multicollinearity occurs when you have two or more independent variables that are highly correlated with each other. This can lead to problems of understanding which variable contributes to the explanation of

the dependent variable and issues while calculating a multinomial logistic regression.

- c) The MLR assumes that there should be a linear relationship between any continuous independent variables and the logit transformation of the dependent variables.

### 3.3. Methodological Procedures/ Application

#### 3.3.1. Likelihood Ratio (LR) and Wald Test

If the dependent variable has  $M$  categories, there are  $M - 1$  non-redundant coefficients associated with each independent variable  $x_N$ . The hypothesis that  $x_N$  does not affect the dependent variable can be written as;

$$H_0: \beta_{N,1|Base} = \dots = \beta_{N,M|Base} = 0$$

where *Base* is the base category (home win) used in the comparison. Since  $\beta_{N,Base|Base}$  is necessarily zero, the hypothesis imposes constraints on  $M - 1$  parameters. This hypothesis can be tested with either the Wald or a LR test.

The LR test estimates the full model that contains all of the independent variables with the resulting LR statistic  $LR_F$ . It then estimates the restricted model formed by excluding the independent variable  $x_N$  with the resulting LR statistic  $LR_R$ . Finally, the LR test takes the difference between  $LR_F$  and  $LR_R$  which is distributed as chi-square with  $M - 1$  degrees of freedom if the hypothesis that  $x_N$  does not affect the outcome is true:

$$LR = LR_F - LR_R .$$

If the  $LR$  statistic for the overall model is significant then there is evidence that the independent variables have contributed to the prediction of the outcome.

The Wald test is an alternative to the LR test and can be computed without estimating additional models. The test is defined as follows:

$$W_N = \widehat{\beta}'_N \widehat{Var} \left( \frac{1}{\widehat{\beta}_N} \right) \widehat{\beta}_N,$$

where,  $\beta_N$  is the  $M - 1$  coefficients associated with  $x_N$  and,  $\widehat{Var}(\widehat{\beta}_N)$ , is the estimated covariance matrix.

If the null hypothesis is true, the  $W_N$  is distributed as chi-square with  $M - 1$  degrees of freedom.

In Stata, the command **mlogtest, lr** computes the likelihood ratio (LR) test, and the command **mlogtest, wald** computes the Wald test.

### 3.3.2. Multicollinearity

Multicollinearity is a state of very high inter-correlations or inter associations among the independent variables. It is mostly caused by inaccurate use of dummy variables, repetition of the same kind of variables, variables being highly correlated with each other, and inclusion of a variable which is computed from other variable in the data set.

Multicollinearity can cause inaccurate estimates of the regression coefficients, inflate the standard errors of the regression coefficients, give false and non-significant p-values, and degrade the predictability of the model.

Testing for multi-collinearity can be achieved by either using the Eigen values method or Variance Inflation Factors (VIF) test. VIF is the most widely used test to measure how much the variance of the estimated regression coefficients are inflated compared to when the predictors are linearly related. It helps to identify the severity of any multicollinearity issues so that the model can be adjusted. The VIF may be calculated for each predictor by linearly regressing that predictor on all other predictors and obtaining the

$R^2$  from that regression. The VIF obtained by the regression can be used in logistic regression models, because the concern is with the relationship among the independent variables included in the model and not with the functional form of the model. Thus, a VIF of 1.6 tells us that the variance of a particular coefficient is 60% larger than it would be if that predictor was completely uncorrelated with all other predictors. A VIF has a lower value of 1.0 but no upper bound. A VIF of more than 10.0, indicates high correlation which is a cause of concern. In a nutshell, the more a VIF increases, the less reliable the results are going to be.

Specifically,

$$VIF_j = \frac{1}{1-R_j^2}$$

where,  $R_j^2$  is the coefficient of determination of the regression model that includes all predictors except the  $j^{\text{th}}$  predictors. If  $R_j^2$  equals zero (no correlation between  $j$  and the remaining predictors), the VIF equals 1.0, which is the minimum value. In Stata, the command **collins** test is used to compute the VIF estimates.

### 3.3.3. Independence of Irrelevant Alternatives (IIA)

The MLR assumes that the odds for any pair of outcomes are determined without reference to the other outcomes that might be available. This is known as the independence of irrelevant alternatives or IIA. If the IIA holds, the MLR can be used. It can be tested by either the Hausman specification test or the Small and Hsiao test. The null hypothesis for both tests is that the IIA does not exist and estimators of the full and restricted models are consistent. On the other hand, under the alternative hypothesis the IIA does exist and only the estimator of the restricted model is consistent.

The Hausman specification test was proposed by Hausman and McFadden (1984) and involves the following steps;

- Estimate the error coefficients of the full model with all  $K$  categories of the dependent variable included; these coefficients are contained in  $\widehat{E}_f$ .
- Estimate the error coefficients of a restricted model by eliminating one or more outcomes categories; these coefficients are contained in  $\widehat{E}_r$ .
- Let  $\widehat{E}_f^*$  represents  $\widehat{E}_f$  after eliminating all coefficients not estimated in the restricted model. The Hausman test of IIA is defined as:

$$H_{IIA} = (\widehat{E}_r - \widehat{E}_f^*)' [Var(\widehat{E}_f) - Var(\widehat{E}_f^*)]^{-1} (\widehat{E}_r - \widehat{E}_f^*).$$

Hausman and McFadden (1984:1226) note that  $H_{IIA}$  can be negative when  $Var(\widehat{E}_f) - Var(\widehat{E}_f^*)$  is not positive semidefinite and suggest a negative  $H_{IIA}$  is evidence that IIA holds.  $H_{IIA}$  is asymptotically distributed as a chi-square with the degrees of freedom equal to the rows in  $\widehat{E}_r$  if IIA is true.

To compute Small and Hsiao's test, the sample is divided into two random subsamples of approximately equal size and the unrestricted MLR is estimated on both subsamples. The weighted average of the coefficients from the two samples is defined as follows:

$$\widehat{E}_u = \left(\frac{1}{\sqrt{2}}\right) \widehat{E}_u^{S1} + \left[1 - \left(\frac{1}{\sqrt{2}}\right)\right] \widehat{E}_u^{S2},$$

where  $\widehat{E}_u^{S1}$  is a vector of estimates from the unrestricted model on the first subsample and  $\widehat{E}_u^{S2}$  is its counterpart for the second subsample. The next step involves creating a restricted sample from the second subsample by eliminating all cases with a chosen value of the

dependent variable. MLR is estimated using the restricted sample yielding the estimates  $\hat{E}_r^{S2}$  and the likelihood  $L(\hat{E}_r^{S2})$ . The Small-Hsiao statistic is the difference:

$$SH = -2[L(\hat{E}_u^{S1S2}) - L(\hat{E}_r^{S2})].$$

SH is asymptotically distributed as a chi-square with the degrees of freedom equal  $N+1$ , where  $N$  is the number of independent variables.

This paper will rely more on the Hausman specification test. It will be applied on each outcome pair of the dependent variable (i.e. match outcome). Since the home win is assumed to be the base category, the test will be applied on basically draw and home loss. This will be done for each of the 6 different leagues chosen and compared to the model with all leagues. If the value of the  $H_{IIA}$  computed in any of the leagues is significant, then the IIA assumption is violated in that model which implies that MLR cannot be used in the modelling. On the other hand, if the values of the  $H_{IIA}$  are insignificant, then the IIA assumption holds and the MLR can be used in the modeling process.

#### **3.3.4. Average Marginal Effects**

Marginal effects are useful estimates of the impact of a one-unit change of a predictor on the dependent variable. Marginal effects (ME) or partial effects, most often measure the effect on the conditional mean of  $y$  of a change in one of the regressors, say,  $x_j$ . In the linear regression model, the ME equals the relevant slope coefficient but for nonlinear models it is quite different. In MLR, the marginal effect of an explanatory variable is the partial derivative of the event probability with respect to the predictor of interest (i.e. the change in the event probability of the dependent variable for a unit change in the predictor), and can be positive or negative. Positive values indicate that the explanatory variable positively contributes to the dependent outcome (i.e. would increase the degree of the predictor

affecting the match outcome), while negative values indicate that the predictor negatively contributes to the dependent outcome.

An average marginal effect is interpreted as the effect of a one-unit change in an independent variable (keeping all other independent variables constant at their mean values) on the dependent variable.

For a multinomial logistic regression model, the probability of response level is given by:

$$p_i = Pr[y_i = j | y_i = j \text{ or } 1] = \frac{e^{x_i' \beta_j}}{1 + e^{x_i' \beta_j}}$$

Where  $x'$  is the predictor of interest, and  $\beta_j$  is the regression coefficient (i.e. log odd) of  $x'$ .

The marginal effect of the  $j^{th}$  predictor is then given by:

$$\frac{\partial p_i}{\partial X_j} = p_i \left[ \frac{\partial X' \beta_i}{\partial X_j} - \sum_K \left[ p_k \frac{\partial X' \beta_k}{\partial X_j} \right] \right]$$

In Stata, the **margins** command computes the average marginal effect.

### 3.3.5. Goodness of Fit

Testing the goodness-of-fit is an important step in evaluating a statistical model. The goodness-of-fit test basically compares observed and estimated frequencies in groups of observations defined by the estimated probability of the reference outcome. It basically targets model misspecification and may help detect a poorly fitting model. Most of the goodness-of-fit tests for the logistic regression are designed for a binary outcome. As the multinomial logistic model can be considered a generalization of the binomial logistic regression with multiple possible outcomes, most authors basically just extended their test statistic from the binary case. Hosmer and Lemeshow (1980, 1989) proposed an extension to the Pearson's chi-square test using a grouping method based on estimated probabilities.



The generalized Hosmer-Lemeshow test is an important goodness of fit measure to assess whether or not the observed events match expected events by sorting data according to the probabilities estimated from the final fitted MLR model. The sorted dataset is partitioned into several equal-sized groups which inherently leads to a construction of a chi-square distribution based on the observed and predicted group frequencies.

The null hypothesis of the Hosmer-Lemeshow test is that the differences between the observed and predicted events are insignificant while the alternative hypothesis is that the differences are significant. An insignificant test statistic implies that the fitted model is a good fit while a significant test implies that the fitted model is not a good fit.

In Stata, the command **mlogitgof, table** is used to generate the Hosmer-Lemeshow test results.

### 3.3.6. Model Comparison

In as much as the goodness-of-fit test is useful in comparing observed and estimated frequencies in groups of observations, it is not something we use in the model building stage to compare different models. Several approaches such as the Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC), the McFadden R-squares, and Cox & Snell's R-squares have all being invented to compare various MLR models and select the best.

The McFadden R-squares result ranges from 0-1, with higher values indicating better model fit. It is defined as:

$$R_{MCF}^2 = 1 - \frac{\ln L_M}{\ln L_0},$$

Where  $L_0$  is the value of the likelihood function for a model with no predictors (i.e. with only a intercept), and  $L_M$  is the likelihood function for the model being estimated. The

model with a small ratio of likelihoods indicates that the full model is a better fit than the intercept model. Thus, when comparing two models on the same data, the McFadden R-square would be higher for the model with the greater likelihood.

The Akaike Information Criterion (AIC), proposed by Akaike (1973, 1974), is widely used for selecting the best model among the candidate models on the same data. The model with the smallest AIC among the candidate models is the best model. In the MLR model, the subset explanatory variables in the best model is the best subset.

AIC has the form:

$$AIC = -2\log L + 2k,$$

where  $k$  is the number of explanatory variables included in the model.

However, the AIC may perform poorly; that is, a model having too many parameters tends to be chosen as the best when the sample size is small, or the number of unknown parameters is large (Satoh, Kamo, and Imori 2012).

The Bayesian Information Criterion (BIC) on the other hand, proposed by Schwarz (1978), extends the AIC, arguing from a Bayesian point of view. The BIC has an advantage over the AIC since the BIC selects the correct model with a probability of 1 as the sample size increases or decreases. The BIC has the form:

$$BIC = -2\log L + k\log(n),$$

where  $n$  is the sample size,  $L$  is the maximized likelihood, and  $k$  is the number of regressors including the intercept. Both the BIC and AIC will be used in the study. In Stata, **fitstat**, **dif** command provides this model comparison methods.

## CHAPTER 4: RESULTS

### 4.1 Data Description

The study is based on an analysis of publicly archived match reports and uses no experimental data or information that involved human subjects. Data were collected for all club matches in the English Premier League, the Italian Serie A, French Ligue 1, German Bundesliga, and the Spanish La Liga from the 2014/2015 season to the 2017/2018 season. A separate dataset of a national team tournament based on the 2018 FIFA World Cup was constructed for comparison purposes. The club dataset consists of 7304 matches while the national dataset consists of 64 matches. Club data were obtained from Opta and ESPN Soccernet. While FIFA World Cup data are obtained from FIFA website. The descriptive statistics of variables are summarized in Table 1.

Table 1: *Description of Variables*

Variables	Description	Total	Percentage
Match Outcome			
1	Home Win	3373	45.78%
2	Draw	1827	24.8%
3	Home Loss	2168	29.42%
H/Corner	Home team total corners		
A/Corner	Away team total corners		
H/Sot	Home team shots on target		
A/Sot	Home team shots on target		
H/FK	Home team free kicks		
A/FK	Away team total free kicks		

**Note:** Total number of games= 7368

As seen in Table 1, football match outcomes (the dependent variable) is modelled using the following three categories:

- Home win: Occurs when the home team scores more goals than the away team.

- Draw: Occurs when the home team number of goals is equal to the away team number of goals.
- Home Loss: Occurs when the home team scored less goals than the away team.

Variables used as predictors of match outcomes include the number of corner kicks for both the home and away team, the number of shots on target for both the home and away team and the total number of free kicks for both the home and away team. Also, a new set of dummy variables were created to account for when a home predictor is greater or less than an away predictor.

The summary statistics of the independent variables hypothesized as affecting football outcomes across the different leagues and overall are reported in Table 2.

Table 2: Summary of Independent Variables

League	Statistic	H/Corner	A/Corner	H/SOT	A/SOT	H/FK	A/FK
<b>All</b>	<i>Mean</i>	5.58	4.43	4.78	3.86	13.07	13.54
	<i>S.D</i>	2.95	2.57	2.58	2.27	4.28	4.36
<b>England</b>	<i>Mean</i>	5.85	4.71	4.70	3.80	10.65	11.26
	<i>S.D</i>	3.14	2.64	2.65	2.23	3.42	3.49
<b>France</b>	<i>Mean</i>	5.35	4.22	4.53	3.74	12.85	13.51
	<i>S.D</i>	2.76	2.51	2.40	2.27	3.94	4.22
<b>Germany</b>	<i>Mean</i>	5.18	4.27	5.08	4.12	13.94	14.79
	<i>S.D</i>	2.84	2.43	2.66	2.31	4.20	4.37
<b>Italy</b>	<i>Mean</i>	5.82	4.70	4.84	3.92	14.17	14.38
	<i>S.D</i>	3.07	2.71	2.59	2.30	4.42	4.54
<b>Spain</b>	<i>Mean</i>	5.68	4.24	4.84	3.78	13.87	13.93
	<i>S.D</i>	2.88	4.24	2.56	2.23	4.30	4.25
<b>W/Cup</b>	<i>Mean</i>	4.84	4.53	4.55	3.30	14.11	14.56
	<i>S.D</i>	2.66	2.27	2.47	1.79	4.43	4.93

**Note:** “H” indicates Home team, “A” indicates away team, “SOT” indicates shots on target, “FK” indicates Free kicks awarded.

Since the research is aimed at analyzing if the relative number of each predictors affect football results, dummy variables were created to compare the number of each home team predictor versus the number of each away team predictor in every game.

Table 3: *Description of Variables contd.*

<b>Variables</b>	<b>Description</b>
Corner	
<b>0</b>	H/corners < A/corners
<b>1</b>	H/corners > A/corners
Shot on target	
<b>0</b>	H/SOT < A/SOT
<b>1</b>	H/SOT > A/SOT
Free kicks	
<b>0</b>	H/FK < A/FK
<b>1</b>	H/FK > A/FK

## 4.2. Multinomial Logistic Regression Results

The prediction results of the MLR are shown in the following sections;

### 4.2.1 Multinomial Logit Coefficients

The multinomial logistic regression model estimates (m-1) equations, where m is the number of outcome levels of the dependent variable, and the m<sup>th</sup> equation is relative to the reference group. In this research, home win is considered as the reference group (base outcome) because it is the most frequent outcome of football games, and the other outcome levels (i.e. home loss and draw) are estimated relative to home win. The standard interpretation of the multinomial logistic regression is that for a unit change in the predictor variable, the probability of one of the outcomes relative to the referent group is expected to change by its respective parameter estimate given the other predictors in the model are held constant.

If the coefficients are positive, then the predictors would increase the likelihood

of the match outcome and if the coefficients are negative, then the predictors would decrease the likelihood of the match outcome.

Table 4.  
*Multinomial Logistic Coefficient Estimates*

Independent Variable	ALL LEAGUES	England	France	Germany
	<i>DRAW</i>			
CORNER	0.443*** (0.0640)	0.325* (0.140)	0.686*** (0.140)	0.395* (0.154)
SOT	-1.333*** (0.0656)	-1.054*** (0.142)	-1.356*** (0.142)	-1.225*** (0.159)
F/K	-0.0594 (0.0612)	0.0642 (0.133)	-0.111 (0.135)	-0.295 (0.153)
CONS	0.00856 (0.0608)	-0.129 (0.130)	-0.0982 (0.131)	0.0685 (0.144)
<i>HOME LOSS</i>				
CORNER	0.618*** (0.0661)	0.640*** (0.147)	0.823*** (0.144)	0.596*** (0.164)
SOT	-2.521*** (0.0692)	-2.576*** (0.156)	-2.301*** (0.149)	-2.508*** (0.172)
F/K	0.129 (0.0625)	0.0775 (0.140)	0.222 (0.136)	0.234 (0.155)
CONS	0.428*** (0.0583)	0.387** (0.123)	0.143 (0.128)	0.326* (0.142)
<i>N</i>	7368	1520	1520	1224

Standard errors in parentheses

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

**Note:** “SOT” indicates shots on target, “FK” indicates Free kicks awarded.

Many of the estimated coefficients in Table 4 themselves render interesting plausible explanations, especially the shots target variable. For example, in all leagues, when home teams play more shots on target, they have a lower log likelihood ratio of losing over winning the game than away teams. This is consistent with previous findings such as Collet (2012) who claims that there is a significant relationship between shooting efficiency and overall team success. Collet’s conjecture is indirectly corroborated by Hacker’s (2013)

findings that even though possession is a key factor in relative team success, the ability to convert possessions into shots on goal distinguishes successful teams from unsuccessful teams.

For the corner predictor, across all leagues, when home teams play more corner kicks, they have a slightly higher log likelihood ratio of drawing or losing over winning the game. This is also consistent with previous findings such as Baranda (2010) who claims that match status plays a huge role on the number of corner kicks played. The losing team tends to intensify their attacks, which increases the probability of getting corner kicks. However, free kicks are found to be insignificant across all leagues.

Table 4.  
*Multinomial Logistic Coefficient Estimates contd.*

Independent Variables	Italy	Spain	World Cup
<b><i>DRAW</i></b>			
CORNER	0.563*** (0.147)	0.252 (0.142)	-0.298 (0.803)
SHOTS	-1.664*** (0.151)	-1.397*** (0.146)	-1.103 (0.809)
F/K	0.0342 (0.136)	-0.105 (0.136)	1.510 (0.803)
CONS	0.143 (0.136)	0.129 (0.147)	-0.729 (0.706)
<b><i>HOME LOSS</i></b>			
CONRER	0.545*** (0.149)	0.558*** (0.147)	-0.453 (0.622)
SHOTS	-2.709*** (0.156)	-2.652*** (0.152)	-0.701 (0.631)
F/K	0.128 (0.139)	0.0305 (0.139)	-0.458 (0.647)
CONS	0.677*** (0.128)	0.570*** (0.142)	0.800 (0.508)
<i>N</i>	1520	1520	64

Standard errors in parentheses  
\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

**Note:** “SOT” indicates shots on target, “FK” indicates free kicks awarded.

#### 4.2.2 Multinomial Logit Marginal Effects

Marginal effects are defined as the slope of the prediction function at a given value of the explanatory variable and thus inform us about the change in predicted probabilities due to a change in a particular predictor. The marginal effect reflects the impact of one-unit change of a predictor on the event probability of the predicted variable (keeping all other predictors constant). In MLR, the marginal effect of an explanatory variable is the partial derivative of the event probability with respect to the predictor of interest (i.e. the change in the event probability of the dependent variable for a unit change in the predictor), and the could be positive or negative. Positive values indicate that the predictor would positively contribute to the football outcome (i.e. would increase the probability of home win, draw, or home loss), while negative values indicate that the predictor would negatively contribute to the football outcome.

The signs of the estimated marginal effects reported in Table 5 are generally consistent across all leagues. For instance, in the English Premier League, the variable corner has a marginal effect of 9.81% (negative and significant), for home win. This implies that every game the home team play more corner kicks then the probability that the home team won decreases by 9.81%. Similarly, the corner has a marginal effect of 8.5% (positive and significant) for home loss, which implies that every game the away team play more corner kicks, the probability that the away team lost increases by 8.5%. This could be because the losing team tends to attack more when in search of goals, and corner kicks are product of attacking play.



For shots on targets, take for instance in Italian Serie A, the shots predictor has a marginal effect of 41.84% (positive and significant) for home win, which implies that when the home team plays more shots accurately than the away team the probability that the

Table 5.  
*Multinomial Logit Marginal Effects Estimates*

<i>LEAGUES</i>		<b>1</b>	<b>2</b>	<b>3</b>
		<b>Home Win</b>	<b>Draw</b>	<b>Home Lose</b>
<b>All Leagues</b>	CORNER	-0.1072*** (0.0110)	0.0343*** (0.0105)	0.0729*** (0.0101)
	SOT	0.3863*** (0.0074)	-0.0517*** (0.0088)	-0.3346*** (0.0077)
	F/K	-0.0060 (0.0107)	-0.0210 (0.0102)	0.0270 (0.0098)
<b>England</b>	CORNER	-0.0981*** (0.0249)	0.0131*** (0.0234)	0.0850*** (0.0221)
	SOT	0.3662*** (0.0175)	-0.0055*** (0.0198)	-0.3607*** (0.0169)
	F/K	-0.0147 (0.0240)	0.0062 (0.0227)	0.0084 (0.0216)
<b>France</b>	CORNER	-0.1554*** (0.0239)	0.0638* (0.0228)	0.0916*** (0.0222)
	SOT	0.3734*** (0.0172)	-0.0734* (0.0199)	-0.3000*** (0.0180)
	F/K	-0.0097 (0.0237)	-0.0380 (0.0226)	0.0478 (0.0215)
<b>Germany</b>	CORNER	-0.1011*** (0.0274)	0.0300* (0.0258)	0.0711** (0.0245)
	SOT	0.3760*** (0.0191)	-0.0443*** (0.0221)	-0.3316*** (0.0190)
	F/K	0.0113 (0.0268)	-0.0724 (0.0257)	0.0611 (0.0235)
<b>Italy</b>	CORNER	-0.1071*** (0.0241)	0.0581*** (0.0233)	0.0489*** (0.0227)
	SOT	0.4184*** (0.0146)	-0.0798** (0.0192)	-0.3329*** (0.0168)
	F/K	-0.0153* (0.0227)	-0.0044* (0.0222)	0.0197 (0.0216)
<b>Spain</b>	CORNER	-0.0793** (0.0241)	0.0037*** (0.0226)	0.0756** (0.0219)
	SOT	0.3978*** (0.0150)	-0.0529*** (0.0184)	-0.3449*** (0.0155)

	F/K	0.0080 (0.0232)	-0.0210 (0.0219)	0.0129 (0.0213)
--	-----	--------------------	---------------------	--------------------

Standard errors in parentheses. \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

**Note:** “SOT” indicates shots on target, “FK” indicates free kicks awarded.

home team wins the game increases by 41.84%. Similarly, the shots predictor has a marginal effect of 33.29% (negative and significant) for home loss, which implies that when the home team plays fewer shots accurately than the away team, the probability that the home team wins that game decreases by 33.29%. The marginal effect of the total free kicks predictor on match outcomes is insignificant across all leagues

### 4.3. Multinomial Logistic Regression Assumptions

#### 4.3.1. Likelihood Ratio and Wald test

Table 6. *Likelihood Ratio (LR) and Wald Test results*

Test	# Observations	Test Statistic	<i>p</i> -value
LR	7367	2933.20	0.0000
WALD	7367	1808.14	0.0000

The likelihood ratio (LR) statistic test is used for testing the effect of any independent variable on the outcome (dependent variable). The null hypothesis of this test is that the predictor variables do not affect the predicted variable. It is calculated by obtaining the log likelihood of the observations with just the outcome in the model with the intercept alone. The final fitted model is then calculated by obtaining the log likelihood of observations with all the predictors in the model. The difference between these two models yields a chi-square Likelihood ratio statistic which is a measure of how well the independent variables affect the dependent variable categories. A significant LR statistic indicates that there is evidence that the predictors are effective, and they have contributed to the prediction of the outcome. On the other hand, if the LR statistic is insignificant then there is evidence that

the predictors are not effective. The results in Table 6 indicates that the LR stat for match outcome is significant at 95% confidence level with its p-value less than 0.05. This implies that all independent/predictor variables included in the model are not equal to zero and that they are effectively contributing to modelling the impacts of corner kicks on match outcome for all categories. Thus, the overall chosen models have good fit.

The Wald test, on the other hand, approximates the LR test, but with the advantage that it only requires estimating one model. The Wald test works by testing the null hypothesis that a set of parameters is equal to some value. In the model being tested, the null hypothesis is that the four coefficients of interest (home shots on target, away shot on target, home free kicks, and away free kicks) are simultaneously equal to zero. If the test fails to reject the null hypothesis, this suggests that removing the variables from the model will not substantially harm the fit of the model. The result presented in Table 6 indicate that the p-value is less than the generally used criterion of 0.05 so the null hypothesis is rejected, indicating that the coefficients are not simultaneously equal to zero. Because including statistically significant predictors should lead to better prediction (i.e., better model fit) it can be concluded that including these four variables results in a statistically significant improvement in the fit of the model.

#### **4.3.2. Multicollinearity**

Multicollinearity is a common problem when estimating linear or generalized linear models, including Logistic regression. It occurs when there are high correlations among predictor variables, leading to unreliable and unstable estimates of regression coefficients. The Multinomial Logit Regression (MLR) model requires that multicollinearity be low between predictors in the model. To test this assumption, the Variance Inflation Factor

(VIF) is used to detect multicollinearity among all predictors in the MLR model as it is the most widely used test for multicollinearity. The VIF measures how much the variance of the estimated regression coefficients is inflated as compared to when the predictors are not linearly related; in other words, the VIF measures how much the behavior of an independent variable is influenced by its interaction with other independent variables.

Table 7. *Multicollinearity Test; VIF results*

<b>Variables</b>	<b>VIF</b>
Home Corner	1.16
Away Corner	1.14
Home total shot on target	1.12
Away total shot on target	1.11
Home total free kicks	1.07
Away total free kicks	1.06

**Mean VIF - 1.11**

VIF has a lower value of 1.0 but no upper bound. A value of 1 implies that the predictor is not correlated with other variables. The higher the value, the greater the correlation of the variable with other variables. Values of more than 4 or 5 are sometimes regarded as being moderate to high, with values of 10 or more regarded as very high. As seen in Table 7, the VIFs of all independent variables are much less than 4 and for this reason, it can be concluded that multicollinearity is not a problem in the data.

#### **4.3.3. Independence of Irrelevant Alternatives (IIA)**

Multinomial Logit models are valid under the Independence of Irrelevant Alternatives (IIA) assumption which states that characteristics of one particular choice alternative do

not impact the relative probabilities of choosing other alternatives. It implies that adding or deleting alternative outcome categories does not affect the odds ratios among the remaining outcomes. The Hausman specification test is used to test the IIA assumption for this football outcome model.

Table 8. *Independence of Irrelevant Alternatives (IIA) Assumption results*

Category	H <sub>IIA</sub>	p-value
Draw vs. Home loss (2 vs 3)	5.618	0.585

The Hausman test was run on each outcome pair of the dependent variable (i.e. football outcome) separately, excluding the base category which in this case is “home win”. Since there are just basically just two categories left after the exclusion of the base category, the test is performed on just the second and third categories, (i.e., draw vs home loss). The null hypothesis of the test is that IIA does not exist and under the alternative hypothesis the IIA does exist. The Hausman test statistic is asymptotically distributed as chi square and significant values indicates that the IIA assumption is violated. As seen in Table 8, the Hausman test statistic was insignificant at the 95% confidence level with its p-value greater than 0.05. This implies that the null hypothesis cannot be rejected, and it can be concluded that the IIA assumption has not been violated.

#### **4.3.4. Goodness of Fits (Hosmer-Lemeshow Statistic)**

The goodness-of-fit tests are based on a comparison of observed and estimated frequencies in groups of observations defined by the estimated probability statistic assesses whether or not the observed events match the predicted events. The Hosmer-Lemeshow

test works by sorting the data according to the probabilities estimated from the final fitted MLR model, the sorted dataset is partitioned into several equal-sized groups, in this research ten groups. The Hosmer Lemeshow (HL) test statistic follows a chi-square distribution that is constructed based on the observed and predicted group frequencies. The null hypothesis is that the difference between the observed and predicted events are insignificant, so the fitted model is correct, while the alternative hypothesis is that the differences are insignificant, so the fitted model has a deficiency. A significant HL test statistic implies that we reject the null hypothesis and conclude that the data do not fit the hypothesized fitted MLR model. On the other hand, an insignificant HL statistic implies that we fail to reject the null hypothesis and conclude that the fitted model is a good fit.

Table 9. Hosmer-Lemeshow test results

<b>Model</b>	<b>HL statistic</b>	<b><i>p</i>-value</b>
Corners	13.410	0.643
Shot on target	29.864	0.019
Free kicks	14.277	0.578
Overall model	28.640	0.026

**Note:** Total number of observations is 7367 and total number of groups is 10.

The results in Table 9 show that even though the HL test statistic for the corner and free kicks model are insignificant at the 95% level, the overall model is significant at the 95% confidence level with a *p*-value lower than 0.05. This implies that we fail to reject the null hypothesis and it can be concluded that the overall models are not a good fit as there is no good match between the predicted events and observed events for all categories of the dependent variable. Goodness-of-fit tests should be considered as just one of several

tools for assessing how fit a model is. Specifically, we cannot conclude that a model fits on the basis of a nonsignificant result from the goodness-of-fit test.

#### 4.3.5. Model Comparison

The McFadden R-square, Akaike Information Criterion (AIC), and Bayesian Information Criterion (BIC) are used to assess which model better predicts the outcome. They are applied to the intercept only model for each dataset and then they are applied to the full model with

Table 10. *McFadden R-squared, AIC and BIC test results*

<b>Criterion</b>	<b>Intercept</b>	<b>Full</b>	<b>Difference</b>
McFadden R-squared	0.001	0.188	-0.187
AIC	15667.270	12750.067	2917.203
BIC	15708.699	12846.734	2861.965

all predictors to capture any improvement in the fitted full model. The McFadden R-square treats the log likelihood of the intercept model as a total sum of squares and the log likelihood of the full model as the sum of squared errors. It ranges from 0-1, with higher value indicating better model fit. AIC and BIC on the other hand, assess the overall fit of a model and allows the comparison of both the full and intercept models. The model with the smaller AIC and BIC is preferred.

As indicated in Table 10, the improvement of the full model over the intercept model through these three approaches is clear. The McFadden R-squared is higher in the full model while the AIC and BIC values are smaller in the full model. This indicates that the fitted full model better predicts the outcomes of the dependent variable, and the predictions are effective in modelling the different outcomes of football matches.

## CHAPTER 5: CONCLUSION

### 5.1. Conclusion

This paper applies multinomial logistic regression (MLR) to model the relationship between relative number of corner kicks and football outcomes across five different domestic leagues and the 2018 FIFA World Cup. The categories of the dependent variables are home win, draw and home loss whereby home win was considered the base outcome (reference group). This paper investigated the influence of corner kicks numbers, shot on targets (SOT) numbers, and free kicks (F/K) numbers on match outcomes using a wider range of datasets, given that past research made use of limited number of independent variables.

The findings show the existence of a relationship between the relative number of corner kicks, shots on target and match outcome are significant. However, even though there exists some sort of relationship, the relative number of free kicks played by teams are insignificant towards match outcome. With reference to the results, this finding disagrees with those of Gordon et al (2013) who claim that corner kicks numbers are statistically useless in soccer games. However, the study supported findings by Baranda (2010) and Riquelme (2012) who found that match status affects the number of corners in a match and how corners are delivered. The results also support Collert (2012) and Hacker (2013) that the key to success in a football game is how well possessions are converted into shot on goals.

In addition, this paper introduces a variety of new procedures in presenting the results of the MLR applications that have not being reported in other football research including:

- 1) a focus on the assumption of Independent and Irrelevant Alternatives (IIA) that is very



important in in the MLR modelling, using the Hausman specification test; 2) testing multicollinearity among independent variables as precondition assumption; 3) presenting the marginal effects of all the predictors upon the dependent variable; 4) using the generalized Hosmer-Lemeshow test as an important goodness of fit measure; and 5) the use of McFadden R-squared, the Akaike Information Criteria (AIC) and the Bayesian Information Criteria (BIC) as potential goodness of fits and model comparison. Results showed the effective use of the MLR approach in analyzing the match outcomes.

## **5.2. Future Research**

With regard to future research, the study unveils some interesting as well as obvious findings when analyzing the effect of corner kick numbers on match outcome. The current study does not consider matches whereby the home team and away team play equal number of corner kicks, shots on goal, and free kicks. Therefore, a comparison of how the effectiveness of corner kicks regardless of their number affects match outcome should be undertaken. Also, the current study does not consider playing style implemented by teams such as counter attack, possession, and direct play when analyzing the effects of corner kicks on match outcomes. Future studies may devote more attention to how differing styles of play affect corner success which simultaneously affects both shot on goals and a team success.

## REFERENCES

- Arda, T., Maneiro, R., Rial, A., Losada, J.L. and Casal, C.A. (2014). Analysis of the effectiveness of corner in the football World Cup 2010. An attempt to identify explanatory variables. *Journal of Sport Psychology*, 23(1): 165-172.
- Azad, A. (2017). Incorporating the Multinomial Logistic Regression in vehicle crash severity modelling: A detailed overview. *Journal of Transportation Terminologies*, 7(3): 279-303.
- Bangsbo, J. and Peitersen, B. (2000). *Soccer systems & strategies*. Champaign, IL: Human Kinetics.
- Bartus, T. (2005). Estimation of marginal effects using `margins`. *Stata Journal*, 5, 309-329.
- Bloomfield, J.R., Polman, R. C. J., & O'Donoghue, P.G. (2005). Effects of score-line on team strategies in FA Premier League Soccer. *Journal of Sports Sciences*, 23(2), 192-193.
- Carling, C., Williams, A. & Reilly, T. (2005). *Handbook of soccer match analysis*. London: Routledge.
- Collet, C. (2013). The possession games? A comparative analysis of ball retention and team success in Europe and International football, 2007-2010. *Journal of Sports Sciences*, 31(2): 123-136.
- Depken, C.A. (2018). Multinomial Models. *Advanced Microeconometrics. Lecture notes*.
- Ender, P. B. (2010). Collin: Collinearity diagnostics for Stata. *Web*.
- Fagerland, M.W. and Hosmer, D.W.J. (2012). A Generalized Hosmer Lemeshow Goodness-of-Fit Test for Multinomial Logistic Regression Models. *Stata Journal*, 12(3): 447- 453
- FIFA. (2018). Statistics Overview of the World Cup, Russia 2018.
- Harry, L. (2015). A comparison of corner kicks between the top and bottom four teams during the 2014/2015 English Premier League season. *Cardiff Metropolitan University, England*
- Hill, A. and Hughes, M. (2001). Corner kicks in the European Championship for association football. 2000. *PASS.COM: computer science and sport III and performance analysis of sport V*, 284-294.
- Hook, C. and Hughes, M. (2001). Patterns of play leading to shots in Euro 2000. *PASS.COM: computer science and sport III and performance analysis of sport V*, 295-302.

- Hughes, C. (1990). *The winning formula*. Collins London
- Hughes, M. and Churchill, S. (2004). Attacking profiles of successful and unsuccessful teams in Copa America 2001. *Journal of Sports Sciences*, 22(6): 505-508
- Hughes, M. and Franks, I. (2005). Analysis of passing sequences, shots and goals in soccer. *Journal of Sports Sciences*, 23(5): 509-514.
- James, N. (2006). Notational analysis in soccer: past, present and future. *International journal of performance analysis in sport*, 6(2): 67-81.
- Lago-Penas C. (2009). The influence of match location, quality of opposition, and match status on possession strategies in professional association football. *Journal of Sports Sciences*, 27(13): 1463-1469.
- Lago-Penas, C., & Lago-Ballesteros, J., Dellal, A., & Gomez, M. (2010). Game related statistics that discriminated winning, drawing and losing teams from the Spanish soccer league. *Journal of Sports Science and Medicine*, 9: 288-293.
- Mara, J., Wheeler, K. and Lyons, K. (2012). Attacking Strategies That Lead to Goal Scoring Opportunities in High Level Women's Football. *International Journal of Sports Science & Coaching*, 7(3): 565-577.
- McFadde, D. (1974). Conditional Logit Analysis of Qualitative Choice Behavior. *Frontiers in Econometrics*, pp. 105-142.
- Sainz de Baranda, P. and Borrás, D. (2005). Analysis of the corner kicks in the World Cup Korea and Japan 2002. *Cultura, Ciencia y Deporte*, 1(2): 87-93.
- Saltas, P. and Ladis, S. (1992). Soccer and study in shots. *Thessaloniki, Greece*.
- Sousa, T. and Garganta, J. (2001). The importance of set-plays in soccer. *Proceedings of the IV Congress of Notational Analysis of Sport*, 53-7.
- Taylor, J., James, N. and Mellalieu, S. (2004). Notational analysis of corner kicks in English premier league soccer. *Journal of Sports Sciences*, 22(5): 518-519.
- Winker, W. (1996). Qualitative and quantitative match analysis in soccer. In Hughes, M. *Notational Analysis of Sport III*, 43-56.

APPENDIX A

Table 5.  
*Multinomial Logit Marginal Effects Estimates contd.*

<b>LEAGUES</b>		1	2	3
		<b>Win</b>	<b>Draw</b>	<b>Lose</b>
World Cup	CORNER	0.0921 (0.1286)	-0.0112 (0.1003)	-0.0809* (0.1252)
	SHOTS	0.1884*** (0.1239)	-0.1048* (0.0973)	-0.0836*** (0.1240)
	F/K	-0.0362* (0.1246)	0.2353 (0.0898)	-0.1991* (0.1172)

Standard errors in parentheses

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

## APPENDIX B

### Average Marginal Effects at 95% Confidence Interval

