

# **EVALUATION OF THE CURRENT STATE OF FOOTBALL MATCH OUTCOME PREDICTION MODELS**

Thiebe SLEEUWAERT

Student ID: 01302061

Promotor: Prof. Dr. Christophe Ley

Tutor(s): Prof. Dr. Christophe Ley

A dissertation submitted to Ghent University in partial fulfilment of the requirements for the degree of Master of Science in Statistical Data Analysis.

Academic year: 2019 - 2020

The author and promoter give permission to consult this master dissertation and to copy it or parts of it for personal use. Every other use falls under the restrictions of the copyright, in particular concerning the obligation to mention explicitly the source when using results of this master dissertation.

Gent, September 4, 2020

The promotor,

The author,

Prof. Dr. Christophe Ley

Thiebe SLEEUWAERT

# **ACKNOWLEDGEMENTS**

First and foremost, I would like to thank the promotor of this master dissertation, Prof. Dr. Christophe Ley, for allowing me to work on this exciting subject. In these unusual times of the coronavirus pandemic, it was not always easy to discuss some approaches and results with your fellow students or professors. Still, Prof. Dr. Christophe Ley adapted to this situation and granted me the guidance I needed.

Secondly, I would like to thank both Prof. Dr. Lars Magnus Hvattum and Dr. Hans Van Eetvelde. Prof. Dr. Lars Magnus Hvattum provided the necessary data to include the plus-minus ratings in this thesis, and Dr. Hans Van Eetvelde shaped this data in the right format to work with.

Last but not least, I would like to thank everyone that made my unusual educational path throughout the university of Ghent possible. Starting out in 2013 as a bachelor student in biology and finishing in 2020 as a MSc in Statistical Data Analysis.



# CONTENTS

<b>Acknowledgements</b>	<b>i</b>
<b>Contents</b>	<b>v</b>
<b>Abstract</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 A brief history of sports betting . . . . .	1
1.2 Literature review . . . . .	2
1.2.1 Goal-based models . . . . .	2
Independent Poisson . . . . .	3
Dependent Poisson . . . . .	3
Skellam distribution . . . . .	4
1.2.2 Result-based models . . . . .	4
Regression approaches . . . . .	5
Thurnstone-Mosteller and Bradley-Terry models . . . . .	5
Machine learning techniques . . . . .	5
1.3 Football leagues . . . . .	6
1.3.1 English Premier League . . . . .	7
1.4 Scoring . . . . .	7
1.4.1 Ranked probability score . . . . .	8
1.4.2 Brier score . . . . .	9
1.4.3 Ignorance score . . . . .	9
1.5 Problem statement . . . . .	10
1.6 Objectives . . . . .	10
<b>2 Methods</b>	<b>11</b>
2.1 Protocol . . . . .	11
2.2 Data . . . . .	12

2.3	Models . . . . .	12
2.3.1	Naive models . . . . .	12
	Uniform . . . . .	12
	Frequency . . . . .	13
2.3.2	Logit regression models . . . . .	13
	ELO based models . . . . .	13
	Plus-minus based models . . . . .	15
2.3.3	Poisson based models . . . . .	18
	Independent Poisson . . . . .	18
	Bivariate Poisson . . . . .	19
2.3.4	Weibull count model . . . . .	20
2.3.5	Machine learning models . . . . .	23
	Data . . . . .	23
	Random Forest . . . . .	25
	Gradient boosting . . . . .	27
2.3.6	Hybrid random forest model . . . . .	28
	Data . . . . .	29
<b>3</b>	<b>Results</b>	<b>31</b>
3.1	General results . . . . .	31
3.2	ELO based models . . . . .	32
3.3	Plus minus ratings . . . . .	33
3.4	Poisson based models . . . . .	33
3.5	Weibull count . . . . .	34
3.6	Machine learning models . . . . .	34
3.7	Hybrid model . . . . .	34
<b>4</b>	<b>Discussion</b>	<b>39</b>
4.1	ELO based models . . . . .	39
4.2	Plus minus ratings . . . . .	40
4.3	Weibull count . . . . .	41
4.4	Machine learning models . . . . .	42
4.5	Hybrid models . . . . .	43
4.6	Practical implications . . . . .	43

<b>References</b>	<b>46</b>
<b>Appendix A Appendix</b>	<b>51</b>
A.1 Code for the scoring rules . . . . .	51
A.2 Code for the ELO based models . . . . .	53
A.2.1 basic ELO . . . . .	53
A.2.2 Goal-based ELO . . . . .	59
A.3 Code for the plus-minus models . . . . .	65
A.3.1 Plus-minus . . . . .	65
A.4 Code for the Poisson based models . . . . .	79
A.4.1 Independent Poisson . . . . .	79
A.4.2 Bivariate Poisson . . . . .	84
A.5 Code for the Weibull count model . . . . .	89
A.5.1 Weibull count . . . . .	89
A.6 Code for the machine learning models . . . . .	96
A.6.1 Random forest Baboota . . . . .	96
A.6.2 Gradient boosting . . . . .	99
A.7 Code for the hybrid models . . . . .	104
A.7.1 Random forest Groll . . . . .	104
A.7.2 Hybrid random forest . . . . .	110



# **ABSTRACT**

Sports betting knows a long history and has always excited the fans. In recent times, football is one of the most popular sports worldwide. Thus the forecasting of football matches is prevalent. Aside from betting, forecasting the outcome of football matches is also relevant for sports journalists and decision-makers within the sport. Many statistical models have been proposed to predict football match outcomes. These models usually incorporate or estimate the strengths of the opposing teams relative to each other. Two main types of models that predict football matches are recognised, namely the result-based and goal-based approaches.

The result-based models predict the outcome classes (homewin/draw/awaywin) of the football matches directly. A simple example of these result-based models is the ordered logistic regression. At the same time, more advanced methods include machine learning techniques, such as random forest classification.

The goal-based alternatives first estimate or assume a suitable distribution for the goals scored by the opposing teams. From those distributions, the outcome classes are then derived. Simple examples include the independent Poisson model, and more advanced examples use machine learning techniques, such as the random forest regression model.

Many models are conducted on different scales. The first scale is how various scholars use different scoring rules to evaluate the models, which makes it challenging to interpret the results over multiple articles. A second scale is that various scholars evaluate their models on diverse leagues. Diverse leagues have distinct competitive levels, which influence the randomness and predictability of the matches. The last scale is how various scholars use vastly different training and testing protocols.

We aim to compare the performance of the current literature, by using a mixture of both simple and more complex models, from both the result-based and goal-based approaches. To circumvent the problems of different scales, we first bring all the models on a uniform training and testing procedure. Secondly, three essential scoring rules (ranked probability score, Brier score and ignorance score) are used to evaluate and rank the performance of the included models. The data for all the models originates

from the English premier league, which circumvents the problem of using diverse leagues.

In total, we compared 12 models to each other: two naive methods, named the uniform and frequency model, three order logit regression approaches, from which two were based on the ELO rating system and one based on the plus-minus rating system, three goal-based models, namely the independent Poisson, the bivariate Poisson and Weibull count model and four machine learning models, namely two random forest approaches, a gradient boosting approach and a hybrid random forest approach. From the two random forest models, one used a result-based approach and the other a goal-based approach.

Our results show that more sophisticated machine learning models have better predictions compared to more simple alternatives. In particular, the gradient boosting model from Baboota and Kaur (2019) had the best performance overall the scoring rules. For this model, we report an ignorance score of 0.933, a Brier score of 0.353 and a ranked probability score of 0.132. We were not able to distinguish between the performances of the result-based and goal-based models. However, goal-based models hold the most promising results. Our results also illustrate the importance of informative and qualitative features, e.g. the models by Baboota and Kaur (2019) had a mixture of both historical features and features that captured the recent performance of the team.

Although our included models cover some key contributions and recent additions to this field of statistical modelling, there are still a lot more models remaining. We suggest that instead of creating new and highly sophisticated models and features, the scholars in this field should focus on ensembling established models and use comparative studies to rank the current literature.

## CHAPTER 1

# INTRODUCTION

In this chapter, the first section gives a brief introduction to the history of sports betting. The following section presents a summary of the current literature concerning football match prediction models. After that, a section dedicated to the most popular football leagues and scoring rules is given. Finally, the last section specifies the problem statement and objectives.

### **1.1 A brief history of sports betting**

Sports betting and sports forecasting have been around for as long as sports events have existed. Records of sports betting go back as far as 2000 years ago. It originates from the ancient Greeks who used to bet on athletic competitions during the Olympics and from the ancient Roman empire where betting on gladiator fights occurred. In those empires, sports betting eventually became legal and consequently spread to neighbouring kingdoms. During medieval times, religious leaders abolished sports betting, which forced it to stay underground, where it continued to grow in popularity (Milton, 2017).

Later in the United Kingdom, sports betting became again increasingly popular and mainstream in the form of outcome betting on horse races. People could place their bets through bookmakers, or so-called "bookies". Presumably, the first bookmaker in the United Kingdom opened in the 1790s (Munting, 1996). Throughout the 20<sup>th</sup> century, multiple countries legalised sports betting again, e.g. 1931 in Nevada and 1961 in The United Kingdom. In 1994, the country Antigua and Barbuda passed laws that allowed enterprises to apply for online betting licenses and so the modern form of online sports betting was born (RightCasino, 2014).

In recent times, football (or soccer) has become one of the most popular sports in the world, so sports betting on football matches is also prevalent. Betting on football matches is estimated to take up 70% of the sports betting market (Keogh and Rose, 2013). The most popular form of betting is outcome betting, where the aim is to predict which of the two competing teams will win the match.

The prediction of which team will win the match excites football fans all around the world. Online bookmakers have created a business out of this excitement and made it into a vast and competitive industry. The prediction of football match outcomes also gives valuable insights for sports journalists and decision-makers within the sport. The following section gives a review of the current literature concerning the prediction of football match outcomes.

## **1.2 Literature review**

There are numerous studies regarding the prediction of football match outcomes. Stefani (1977) predicted football matches by using a least-squares model that rated the strengths of both competing teams relative to each other. The trend of first estimating the relative strengths of the competing teams is still present in numerous recent studies. Estimating the relative strength of a team can be done by a multitude of methods and is usually conducted on a team-based level. Recently, with the increase in available data, researchers have rated players on an individual-level and derived the ratings of a team from those individual player ratings (Arntzen and Hvattum, 2020).

Most statistical models that aim to predict the results of football matches are categorised under two main types. The first type of models aims to estimate the distribution of the goals for both competing teams. From those distributions, the outcome classes (win/draw/loss) are then indirectly derived. We will refer to these types of models as *goal-based* models. The second type of statistical models aims to predict the outcome classes directly. We will refer to these types of models as *result-based* models. The following subsections give an overview of the methods that hold notoriety in this field.

### **1.2.1 Goal-based models**

Goal-based models aim at estimating the goal distribution of the opposing teams. These models first assume a suitable count distribution for the football goals scored by each team. One of the most frequently used approaches assumes that the goals follow a Poisson distribution.

**Independent Poisson**

The independent Poisson model, first proposed by Maher (1982), assumes that the number of goals scored by each team follows a Poisson distribution and that these distributions are independent from each other (Katti and Rao, 1968). The following formula gives the Poisson density function:

$$P(X = x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda} \quad (1.1)$$

This formula shows the probability of observing  $x$  number of goals given  $\lambda$ , with  $\lambda$  equal to the expected number of goals for a given team. We will refer to this parameter  $\lambda$  as the rate intensity parameter.

Most statistical models first estimate  $\lambda$  by constructing it to represent the relative strengths of the competing teams. For example, Ley et al. (2019) used  $\lambda_{i,m} = \exp(c + (r_i + h) - r_j)$  to estimate this parameter, where  $r_i$  and  $r_j$  represent the relative strengths of respectively the home and away team, and  $h$  represents the effect of playing as the home team. Chapter 2 gives a more thorough explanation of this model.

**Dependent Poisson**

Experts in the field mostly agree that the assumption of independence between the goal distributions of both competing teams is flawed. These goal distributions have some apparent dependence between them, since the estimation of the  $\lambda$  parameters usually involves the relative strengths of both competing teams (Ley et al., 2019). Additionally, in team sports, it is reasonable to assume that the goals made by each team are dependent since both teams interact during a match (Karlis and Ntzoufras, 2003).

Dixon and Coles (1997) extended the basic independent Poisson model by including an indirect correlation term between the goal distributions of the competing teams. Dixon and Coles (1997) identified this correlation to be slightly negative. The correlation term is indirect since it ignores the direct correlation between the intensity parameters  $\lambda$  of the opposing teams.

Karlis and Ntzoufras (2003) used a bivariate Poisson model with a direct dependence term between the goal distributions of the competing teams. This bivariate Poisson model has the advantage that each goal distribution still follows a Poisson distribution marginally. A deeper explanation is given in section 2.3.3 - Bivariate Poisson.

There is a multitude of other methods available which introduce some form of dependence to the basic independent Poisson model. Some examples include McHale and Scarf (2011) who extended the independent Poisson with copula dependent structures. Boshnakov et al. (2017) presented a Weibull interval-arrival-time-based count process with a copula to model the number of goals for each team. We will discuss this model in more detail in section 2.3.4 - Weibull count model. The Weibull count model, and many other models, assume different distributions to model the goals made by the opposing teams. Some other examples include negative-binomial distributions, gamma-Poisson distributions and zero-inflated-Poisson distributions.

### **Skellam distribution**

One of the advantages of using Poisson based models is that the Skellam distribution can be derived from it (Skellam, 1946). The Skellam distribution, or Poisson difference distribution, is the discrete probability distribution of the difference between two Poisson random variables (Ley et al., 2019). If we consider  $GD_m = G_{i,m} - G_{j,m}$  as the difference between the expected goals scored by team  $i$  and team  $j$  during match  $m$ , then the probability of a win of team  $i$  over team  $j$  is calculated as  $p(GD_m) > 0$ . The probability of a draw and loss is calculated as respectively  $p(GD_m) = 0$  and  $p(GD_m) < 0$ .

### **1.2.2 Result-based models**

Although modelling the goal distributions is the most frequently used approach to predict football match outcomes, other types of statistical models with this aim exist. Namely, result-based models that directly predict the outcome classes. These models usually apply some form of ordered logit or probit regression. Recently, with the increase in available data and advanced computational algorithms, other types of result-based models have been proposed.

Directly modelling the outcomes does not indicate anything significant about the estimated goal difference between teams, which makes these result-based models milder in their assumptions and usually they have fewer parameters to estimate (Egidi and Torelli, 2020). A potential downside of these models is the overestimation or underestimation of the relative strengths of the competing teams since these models do not include the actual goal difference to derive these strengths (Egidi and Torelli, 2020).

### **Regression approaches**

One of the earliest articles concerning result-based models, comes from Koning (2000), where a probit regression model was used to estimate the outcome classes directly. Goddard (2005) also used a probit regression model and compared it to the bivariate Poisson regression models. The article reports that the best performance was achieved by a hybrid model that combined a result-based dependent variable with goal-based lagged performance covariates. However, the differences among the models were small, and thus both approaches are considered relevant.

Hvattum and Arntzen (2010) used an ordered logit regression on the difference of the opposing teams ELO ratings and reported that the predictive performance of this model worked better than the result-based models of Goddard (2005). Hvattum (2017) reports that these logit regression models had difficulties predicting draws. To circumvent this problem, Egidi and Torelli (2020) used multinomial regression models with subtracted factors to inflate the probability of draws. They also compared goal-based and result-based approaches and found that the multinomial regression models were slightly lower in predictive performance when looking at the Brier score as a scoring rule, compared to the goal-based alternatives. However, the differences were again insignificant.

Recently, Arntzen and Hvattum (2020) also used an ordered logit regression model based on the plus-minus ratings of the players. The vast increase in data availability has made it possible to estimate such individual player rating. The authors report that the ordered logistic regression model based on the plus-minus ratings outperforms the ordered logistic regression model based on the team-based ELO ratings. However, when both covariates are combined, the predictive performance is significantly enhanced.

### **Thurnstone-Mosteller and Bradley-Terry models**

Both Thurnstone-Mosteller (Mosteller, 2006; Thurstone, 1927) and Bradley-Terry (Bradley and Terry, 1952) type models were used successfully by Ley et al. (2019) to directly model and predict football outcome classes. These models predict the outcomes of pairwise comparisons by using latent variables.

### **Machine learning techniques**

Recently, machine learning techniques have been used to predict the outcome classes of football matches directly. Joseph et al. (2006) showed that Bayesian nets out-

performed other supervised machine learning classification models such as decision trees, naive Bayes and k nearest neighbours. Constantinou et al. (2012) proposed a Bayesian network to predict the outcome classes of a match. In a follow-up study, Constantinou and Fenton (2013) illustrated that a new ranking system, called the pi-ratings, incorporated in their Bayesian network model, significantly outperforms the ELO ratings. Groll et al. (2018) used a random forest with multiple informative covariates to predict the goals scored by the opposing teams. These estimates were then used as the intensity parameters  $\lambda$  for the Poisson distributions of the opposing teams, and used to derive the outcome class probabilities. Groll et al. (2019) expanded this random forest model to a hybrid random forest model, where first the relative strengths of the competing teams were estimated based on a Poisson maximum likelihood approach. Then secondly these relative strengths were used in a random forest model, combined with other covariates from Groll et al. (2018), to estimate the goals of the opposing teams, from which the outcome class probabilities were again indirectly derived. Baboota and Kaur (2019) used the gradient boosting, naive Bayes, linear support vector machine, RBF support vector machine and random forest algorithm to estimate the outcome distributions directly. The article reports that the gradient boosting algorithm outperformed all other models.

The performance of statistical models is usually only evaluated on a single football league. The following section gives a brief overview of the most popular and important football leagues. Another issue is that different authors use diverse scoring rules to evaluate the performance of their statistical models. After the section about the football leagues, a section dedicated to the various scoring rules is given.

## **1.3 Football leagues**

There are many football leagues throughout the world. Some of the most competitive and popular include the English - Premier League, the German - Bundesliga and the Spanish - La Liga. There are also matches between national teams organised by international football federations, such as the International Federation of Association Football (FIFA), and FIFA confederation such as the Union of European Football Associations (UEFA). These international federations organise football tournaments between nations, usually yearly or over multiple years, e.g. the FIFA World Cup tournament takes place every four years and is the most famous international football championship (Suzuki et al., 2010).

Due to the immense popularity and highly competitive level of the English Premier League, we will include the data from this domestic league in this master dissertation.

### 1.3.1 English Premier League

The English Premier League is the most famous football league worldwide and is considered to be of the highest competitive level. It consists of 20 teams that play against each other twice a season, for a total of 380 matches. The English Premier League is broadcasted worldwide in 212 countries and reaching 4.7 billion people (Kundu et al., 2019). The revenues generated are therefore enormous and estimated at 2.2 billion euro in television rights and 5.8 billion euro from other sources such as merchandise and ticket sale. These numbers illustrate the magnitude of how successful the English Premier League is.

The highly competitive nature of the English Premier League gives the outcome distribution of matches much randomness. Therefore it is rather challenging to come up with accurate prediction models. One way to measure the randomness of a dataset is by looking at its entropy. An entropy score of 1 means complete randomness. Kundu et al. (2019) reports that for the historical English Premier League data between the seasons of 2005 and 2016, the entropy of the dataset was 0.96, which again supports the assumption that the English Premier League is highly competitive.

## 1.4 Scoring rules

Scoring rules are functions used to evaluate the performance of predictive models. There is a wide range of scoring rules available, each developed for different purposes and situations. For football match outcomes, which is considered to be an ordinal outcome parameter with three classes, namely home win, draw and away win, there are also numerous choices available. Debate exists over which scoring rule is the most appropriate (Wheatcroft, 2019).

Wheatcroft (2019) considered three properties in scoring rules to be relevant for the evaluations of models that aim to predict football match outcomes. Firstly, a scoring rule must be proper, meaning that it favours predictions that consist of distributions drawn from the actual outcome distribution. Another central property in scoring rules that aim to evaluate football match prediction models, is locality. A scoring rule can either be local or non-local. It is considered local if it only takes the probability from the predicted class into account. Therefore, a non-local scoring rule takes the probability from multiple classes into account. If a score is non-local, a final property to consider is sensitivity to distance. This sensitivity to distance follows the rationale that the predictive outcome classes are ordinal. A scoring rule should, therefore, penalise a model more in case of an observed home win and predicted away win, compared to

a predicted draw (Ley et al., 2019). A scoring rule that is insensitive to distance does not follow the rationale that the predictive outcome classes are ordinal.

Most articles focus on the ranked probability score (RPS), the Brier score (BS) and ignorance score (IGN). Some papers use accuracy to evaluate the performance of their models. However, it is often challenging to predict draws, hence scoring rules that focus on the probability placed on each outcome class, are preferred.

### 1.4.1 Ranked probability score

Ranked probability score, proposed by Epstein (1969), is considered to be the most appropriate scoring rule by Constantinou and Fenton (2012) and has since gained more recognition. It is currently the most popular and widely used scoring rule for evaluation of football match outcome models. The ranked probability score is a non-local scoring rule that is sensitive to distance. The following function gives the ranked probability score:

$$RPS = \frac{1}{2N} \sum_{n=1}^N ((P_{Hn} - y_{Hn})^2 + (P_{An} - y_{An})^2) \quad (1.2)$$

Here the parameters  $P_{Hm}$  and  $P_{Am}$  are the predicted probabilities of a home win or away win.  $y_{Hm}$  and  $y_{Am}$  are the observed outcomes, so either 1 or 0.  $N$  is the total number of predicted matches.

The ranked probability score is a topic of present debate in the literature. Ever since Constantinou and Fenton (2012) proposed it to be the most appropriate scoring rule, it gained notoriety. However, Wheatcroft (2019) recently published an article where he argued against the importance of the non-locality and sensitivity to distance properties in a scoring rule for the evaluation of models that predict football match outcomes.

The main argument in favour of the ranked probability score is that probabilities placed on outcomes close to the observed outcome should receive a higher reward (Constantinou and Fenton, 2012). If a team is winning by one goal, it takes the opposing team one goal to end the match in a draw and two goals to end in a win for the opposing team. In this regards, the outcome classes are considered ordinal, and thus sensitivity to distance is deemed to be essential in its scoring rule, which makes the ranked probability score an obvious choice. Constantinou and Fenton (2012) then gives hypothetical examples of football matches, from which they conclude that the ranked probability score is favoured, because it assigns the best score to the favoured forecast in each case.

The main counter-argument of Wheatcroft (2019) is that the examples provided by Constantinou and Fenton (2012) are flawed since it compares the performance of the scores under specific outcomes. Instead, the underlying probability distribution of the match must be taken into consideration because the observed outcome gives no information about the actual underlying distribution. Wheatcroft (2019) then reproduced the examples given by Constantinou and Fenton (2012) and shows that the Brier score and in particular the ignorance score are more appropriate scoring rules for football match prediction models.

For a deeper understanding of these examples, we refer the readers to both articles. In this master dissertation, however, we will use a combination of different scoring rules to circumvent this debate.

### 1.4.2 Brier score

The Brier score (Brier, 1950), or the squared loss function, is quite similar to the ranked probability score but it is insensitive to distance. This insensitivity to distance means that it does not penalize a model more according to the ordinal structure of the outcome classes. The following function gives the Brier score:

$$BS = \frac{1}{N} \sum_{n=1}^N \sum_{r=1}^R ((P_{nr} - y_{nr})^2) \quad (1.3)$$

Here  $N$  is the number of predicted matches. The parameter  $R$  stands for all of the possible outcome classes.  $P_{nr}$  is the probability placed on instance  $n$  for outcome class  $r$ .  $y_{nr}$  is the observed outcome, for instance  $n$  and outcome class  $r$ , so either 0 or 1.

### 1.4.3 Ignorance score

The ignorance score (Gneiting and Raftery, 2007), or the logarithmic loss function, is a scoring rule that is both local and insensitive to distance. The following function gives the ignorance score:

$$IGN = \frac{1}{N} \sum_{n=1}^N [-\log_2(p(y_n))] \quad (1.4)$$

Here  $p(y_n)$  is the probability placed on the correct outcome class  $y$  of match  $n$ . Wheatcroft (2019) found practical evidence that the ignorance score is the most proper scoring rule (Bröcker and Smith, 2007).

## **1.5 Problem statement**

The main problem is the multitude of models that aim to predict football match outcomes on incomparable scales. The first layer of the problem is that scholars use different scoring rules to evaluate the performance of their statistical models. Therefore the comparison between them is challenging to interpret. A second layer is that scholars use data of different leagues. Comparison of models trained on different data is often difficult to understand since different leagues have different competitive levels. A third layer and final layer is that scholars use different training and test protocols for their statistical models.

## **1.6 Objectives**

In this master dissertation, we aim to evaluate the current literature of football outcome prediction models and circumvent the problems mentioned above in section 1.5. The remainder of the master dissertation is built up as follows: section 2 explains the methods used, section 3 shows the results and finally in section 4 the results are discussed.

## CHAPTER 2

# METHODS

This chapter explains the methods used to create a detailed overview of the models on comparable scales, in detail. First, we will go over the general protocol and data. The following sections explain each evaluated model separately. For each model, first, the method is conceptualised, and if necessary, the data specific to the model are explained.

### **2.1 Protocol**

A uniform protocol brings the evaluated models on comparable scales. The evaluation of each model will use all the scoring rules, mentioned in section 1.4 (ranked probability score, brier score and ignorance score), in conjunction. The models are ranked for each scoring rule, and the best performing model is defined by having the highest average rank. After that, the models are estimated on data originating from the English Premier League, which circumvents the problem of using different leagues.

Most models use vastly different predicting procedures. We aim to bring them all together on a comparable scale. All English Premier League matches in the seasons of 2008 to 2015 are predicted. During the prediction of a season, the previous two seasons of matches are used as data, combined with the first five weeks of the current season, to estimate the model coefficients. The reason for the burn-in period of the first five weeks is to get reliable information on the possible new teams entering the league.

The weeks of a season are predicted in a stepwise manner. After the prediction of each week, the information is added to the data. A football season consists out of 38 weeks, so in total, for each season, 810 matches are used as initial data, and 330 matches are predicted. In total, for the eight seasons, 2640 matches are predicted

Many statistical models from the current literature have some hyperparameters or variables that require estimation. These hyperparameters are not estimated under the new protocol, but merely the same values from the articles are used. If these val-

ues are absent, the parameters will take arbitrary reasonable values as recommended by the literature. We will explicitly mention our reasoning where this occurs.

## 2.2 Data

The English Premier League data used is from the *engsoccerdata* R package (Curley, 2016). This package is mainly a repository that contains different European football datasets, such as the three English ones (Premier League, FA Cup, Playoff) and also other European leagues (Spain - La Liga, Germany - Bundesliga, Italy - Seria A, Netherlands - Eredivisie).

Date	Season	home	visitor	hgoal	vgoal	result	round
2011-08-13	2011	Fulham	Aston Villa	0	0	D	1
2011-08-13	2011	Liverpool	Sunderland	1	1	D	1
2011-08-13	2011	Newcastle	Arsenal	0	0	D	1
2011-08-13	2011	Wigan	Norwich City	1	1	D	1
...	...	...	...	...	...	...	...
...	...	...	...	...	...	...	...

Table 2.1: Example of the data used for the English Premier League (Curley, 2016)

## 2.3 Models

This section explains the statistical models used in detail. The additional data required by some models is also specified here. The models are ranked to go from more simple approaches to more complex ones.

### 2.3.1 Naive models

The first two models utilise the available information naively. We will refer to them as the *uniform* and *frequency* models. These models will serve as benchmark instruments since any model that does not outperform them, is rendered ineffective.

#### Uniform

The uniform model ignores all available data on past football matches and assumes that the probabilities for a home win, away win or draw are uniformly distributed. Effectively this means that for every match, every outcome class gets a probability of  $\frac{1}{3}$  assigned to it.

**Frequency**

Unlike the uniform model, the frequency model does incorporate some information about past matches played. The frequency model estimates that the probabilities for a home win, away win or draw are distributed as the observed frequencies of the home wins, away wins and draws in the past  $k$  matches. So effectively the win probability is equal to  $p(\text{homewin}) = \frac{1}{N} \sum_{k=1}^N I(\text{homewin}_k)$ , analogous we find  $p(\text{awaywin}) = \frac{1}{N} \sum_{k=1}^N I(\text{awaywin}_k)$  and  $p(\text{draw}) = \frac{1}{N} \sum_{k=1}^N I(\text{draw}_k)$ .

**2.3.2 Logit regression models**

This section explains two articles that both utilise a logit regression approach to predict the outcome classes.

The first article focuses on the ELO rating system to represent the strengths of the opposing teams. The second article uses the plus-minus rating system to calculate the strengths of the individual players relative to the players within the same team and players between opposing teams.

**ELO based models**

The ELO based models come from the article by Hvattum and Arntzen (2010). Here the ELO rating system, adapted for football matches, is used to estimate the current strengths of the competing teams. An ordered logit regression model then incorporates the ELO ratings as the single covariate.

Hvattum and Arntzen (2010) compared the predictive performance of this model to six other methods, namely the two naive methods equal to the uniform and frequency model mentioned above, two probit regression models derived from Goddard (2005) and two models based on the odds offered by bookmakers. The article reports that the logit regression based on the ELO ratings outperforms the naive and probit regression models from Goddard (2005), but does not outperform the models based on the bookmaker odds.

The ELO rating system, modified for football matches, estimates the current strengths of a team based on historical data. A scoring system is first defined to derive this rating. In this scoring system, a win grants a score of 1, a draw a score of 0.5 and a loss a score of 0. If  $l_0^i$  and  $l_0^j$  are the current ELO ratings of teams  $i$ , the home team, and team  $j$ , the away team, respectively, then the ELO rating framework assumes

that, on average, each team should score  $\gamma^i$  and  $\gamma^j$  against each other in a given match. The following formulas give the functions for  $\gamma^i$  and  $\gamma^j$ :

$$\begin{aligned}\gamma^i &= \frac{1}{1 + c^{-(l_0^i - l_0^j)/d}} \\ \gamma^j &= \frac{1}{1 + c^{-(l_0^j - l_0^i)/d}} = 1 - \gamma^i\end{aligned}\tag{2.1}$$

The formulas above are dependent on two parameters, namely  $c$  and  $d$ , with  $c > 1$  and  $d > 0$ . These parameters are only used to scale the expected scores. Hvattum and Arntzen (2010) reported that  $c = 10$  and  $d = 400$  is sufficient. Alternative values can give identical rating systems. In 2018, the FIFA rating system was changed to the ELO-based system, with  $d = 600$ . Alternative methods extend the ELO ratings by incorporating a home advantage effect  $h$ . The formula above is then formulated as  $\gamma^i = \frac{1}{1 + c^{-(l_0^i + h - l_0^j)/d}}$ . For example the website <https://eloratings.net/> uses this formula with  $h = 100$ .

The scoring system gives the observed scores  $\alpha^i$  and  $\alpha^j$  for both teams, with  $\alpha^i = 1$  if team  $i$  won, 0.5 if the match was a draw and 0 otherwise.  $\alpha^j$  is calculated as  $1 - \alpha^i$ . The ELO ratings are updated after every match from  $l_0$  to  $l_1$  by the following formulas:

$$\begin{aligned}l_1^i &= l_0^i + k(\alpha^i - \gamma^i) \\ l_1^j &= l_0^j + k(\alpha^j - \gamma^j)\end{aligned}\tag{2.2}$$

The formulas above are dependent on the parameter  $k$ . Unlike the values for the parameters  $c$  and  $d$ , a more careful value consideration is needed for  $k$ . If  $k$  is set too low, team ratings take too long to stabilise, and if  $k$  is set too high, team ratings will be too volatile. Hvattum and Arntzen (2010) reported that  $k = 20$  is reasonable and that after 30 matches the ELO ratings are usually stabilised.

Note that the ELO rating system requires some initial data in order to indicate the current strengths of the opposing teams reliably. For this reason, the protocol mentioned in 2.1 is extended with an initial period of one year, e.g. for the prediction of the season for the year 2008, the season of 2005 is used to get reliable initial estimates of the ELO ratings before starting the training protocol.

Two prediction models are created from this ELO rating system, namely the *basic ELO* and *goal-based ELO*. The basic ELO uses the methodology described above, with  $c = 10$ ,  $d = 400$  and  $k = 20$ . The goal-based ELO extends the basic ELO, by improving the parameter  $k$  to depend on the absolute goal difference between the teams. The

incorporation of this goal difference ensures that a higher difference translates to a higher gain in ELO. The parameter  $k$  is now extended to  $k = k_0(1 + \delta)^\lambda$ , with  $\delta$  equal to the absolute goal difference. The goal-based ELO thus requires the estimation of four hyperparameters. Hvattum and Arntzen (2010) reported that  $c = 10$ ,  $d = 400$ ,  $k_0 = 10$  and  $\lambda = 1$  are reasonable.

The ordered logit regression model is then used, with the difference in the ELO ratings of the opposing teams as a single covariate, to make predictions for the match results. The difference in ELO ratings is defined as  $x = l_0^i - l_0^j$ , with  $l_0^i$  equal to the current ELO rating of team  $i$ , the home team, and  $l_0^j$  for team  $j$ , the away team. After each prediction, the ELO ratings are updated. Figure 2.1 shows how these ELO ratings have evolved from 2000 to 2015. The ordered logit regression model is defined by:

$$p(y = i|x) = F(-\Theta_i - \beta x) - F(-\Theta_{i-1} - \beta x) \quad (2.3)$$

$$F(z) = \frac{1}{1 + e^{-z}}$$

In the formula above  $j \in \{1, 2, 3\}$  represents the outcome classes, with  $j = 1$ , for an away win,  $j = 2$  for a draw and  $j = 3$  for a home win. The parameters  $\Theta$  are used to differentiate between the ordinal values of the dependent variable. Only  $\Theta_1$  and  $\Theta_2$  need to be estimated, since  $\Theta_0 = \infty$  and  $\Theta_3 = -\infty$ . The function  $F(z)$  represents the logit link, with  $F(-\infty) = 0$  and  $F(\infty) = 1$ . So in total the the logit regression needs to estimate three parameters, namely  $\Theta_1$ ,  $\Theta_2$  and  $\beta_{ELOdiff}$ .

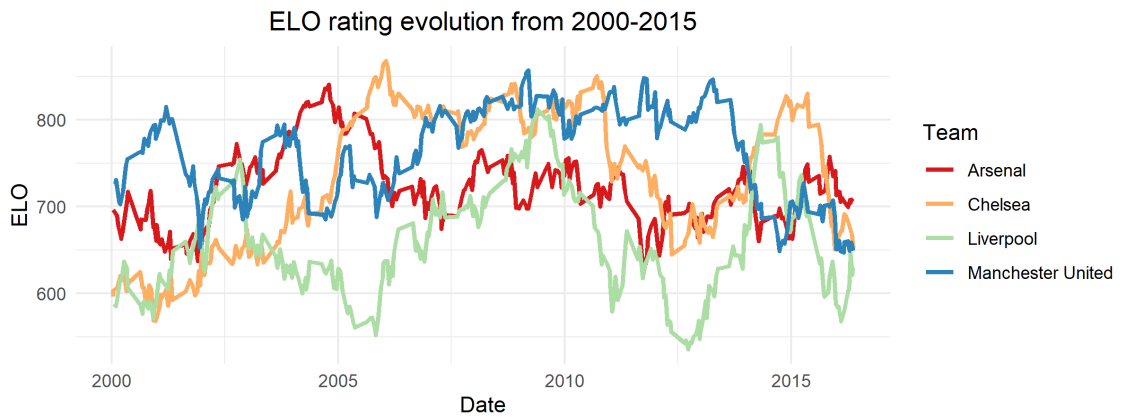


Figure 2.1: Evolution of ELO ratings from 2000-2015

### Plus-minus based models

The plus-minus based models are derived from the article by Sæbø and Hvattum (2015). Here the plus-minus rating system, adapted for football matches, is used to estimate the strengths of the players relative to its own teammates and the opposing

team's players. The plus-minus ratings of a team are then calculated as the average plus-minus ratings of the players within a team. Again, an ordered logit regression model, with the difference in the plus-minus ratings of the opposing teams, is used to estimate the outcome class probabilities.

Plus-minus ratings attempt to distribute credit for the goals of a team onto the players responsible for it (Hvattum, 2019). In its simplest form, the plus-minus rating system calculates the goals scored minus the goals conceded for every player during a match. Specifically for football matches, Sæbø and Hvattum (2015), came up with a regularised adjusted plus-minus rating system.

For the regularised adjusted plus-minus ratings, we first need to define segments, within a match, of constant players on the field. Every time a team changes a player on the field, or a player receives a red card, the match is split into two separate segments.

For each segment  $i$ , we define an appearance of a player  $j$ , by the parameter  $\alpha_{ij}$ . The value of  $\alpha_{ij}$  is given by the following formula:

$$\alpha_{ij} = \begin{cases} e^{-kt} & \text{If player } j \text{ plays for the home team in segment } i \\ 0 & \text{If player } j \text{ does not play in segment } i \\ -e^{-kt} & \text{If player } j \text{ plays for the away team in segment } i \end{cases} \quad (2.4)$$

In the formulas above, the factor  $e^{-kt}$  is dependent on two parameters, namely  $k$  and  $t$ , and represents a time depreciation effect. The parameter  $t$  is the difference between the current time, when the plus-minus ratings need to be estimated, and the time when a match is played, expressed in years. The parameter  $k \in [0, 1]$  represents the magnitude of the time depreciation effect. Sæbø and Hvattum (2015) report that  $k = 0.2$  is within reason. This value implies that the weight of a match played five years ago contributes only for  $\frac{1}{e} (\approx 0.37)$  to the estimation of the plus-minus rating, compared to matches played in the present time.

Since the plus-minus rating system is based on the goals scored minus the goals conceded, we define a parameter  $\beta_i$ . This parameter represents a scaled version of the goals scored minus the goals conceded, in favor of the home team, during a segment  $i$ . The following function defines  $\beta_i$ :

$$\beta_i = \frac{90(H_i - A_i)e^{-kt}}{D_i} \quad (2.5)$$

In the formula above, the parameter  $D_i$  represents the duration of segment  $i$ ,  $H_i - A_i$  represents goals scored minus the goals conceded in favour of the home team during segment  $i$  and  $e^{-kt}$  is again the time depreciation effect.

A twelfth dummy player, that is included in every home team's starting lineup, accounts for the home advantage effect. The appearance of the home advantage is thus always equal to  $e^{-kt}$ . For the calculation of the plus-minus ratings of a home team, the average of the players and home team advantage must be taken.

Four dismissal dummy variables account for the effect of red cards. When a player  $j$  in the home team gets a red card, then the first dismissal dummy gets a value of  $e^{-kt}$ , and the player's appearance is changed to 0. A second red card for the home team would lead to a transfer of appearance from the player getting the red card to the second dismissal dummy. Red cards for the opposing teams can nullify the dismissal dummy variables, e.g. if the home team has two red cards and the away team one, then the first dismissal dummy has a value of 0 and the second a value of  $e^{-kt}$ .

Finally the plus-minus ratings  $x$  are calculated by:

$$x = (a^T a + \lambda I)^{-1} a^T \beta \quad (2.6)$$

The formula above corresponds to the estimation of the coefficients in a penalised linear regression by using the least-squares criterion. The penalisation term is crucial, since many players are joint for most of their playtime. The rating system then struggles to differentiate between them, which causes collinearity issues and inflates the errors. Also, the ratings for players with little playing time are prone to large errors. The penalisation term in equation 2.6 is equivalent to the ridge regression approach and was found by Macdonald (2012) to be the most appropriate penalisation term when estimating plus-minus ratings. Sæbø and Hvattum (2015) report that  $\lambda = 3000$  is sufficient.

Note that the plus-minus ratings, just like the ELO ratings, requires some initial data in order to indicate the current strengths of the opposing teams reliably. Arntzen and Hvattum (2020) used the previous five seasons to get those initial ratings, so we will extend the protocol, mentioned in 2.1, with an initial period of five years, e.g., for the prediction of the season of 2015, the seasons from 2008 till 2012 are used to get the initial ratings. Then the seasons 2013 and 2014, in combination with the first five weeks of season 2015, are used to estimate the coefficients of the ordered logit regression.

In this master dissertation, we were not able to incorporate the data for the red cards due to time constraints. Also important to note is that Sæbø and Hvattum (2015) did not use the plus-minus ratings initially to predict football match outcomes, but instead, they used it to evaluate the efficiency of the transfer market. However, in a follow-up study by Arntzen and Hvattum (2020) the difference in plus-minus ratings between the opposing teams, in favour of the home team, is taken as a single covariate in an ordered logit regression, see equation 2.4. Arntzen and Hvattum (2020) also improved the plus-minus rating system to depend on an age effect, the similarity of the players and the competition type.

### 2.3.3 Poisson based models

The third set of models comes from the article by Ley et al. (2019). Here eight different statistical models are compared in their performance to predict football match outcomes. From those models, the independent and bivariate Poisson had the highest predictive power. Additionally, these models create a ranking method based on a maximum likelihood approach. We will refer to this ranking parameter as the ability of a team. Figure 2.2 shows how these have evolved between 2000 and 2015.

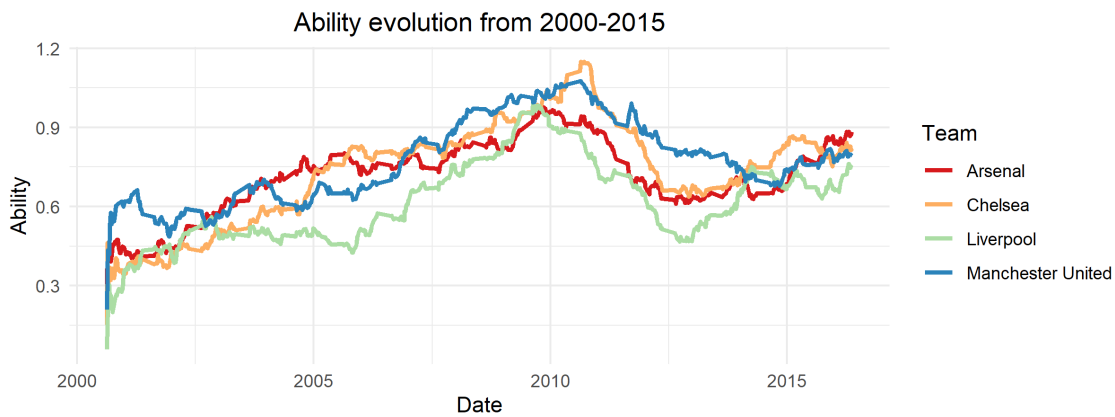


Figure 2.2: Evolution of the abilities from 2000-2015

#### Independent Poisson

The independent Poisson models the goals of the opposing teams  $i$  and  $j$  by the random variables  $G_{i,m}$  and  $G_{j,m}$  for match  $m$ . These random variables follow a Poisson distribution. The following formulas give the joint density distribution of observing  $x$  goals for team  $i$  and  $y$  goals for team  $j$ , under the assumption that both random variables are independent from each other:

$$P(G_{i,m} = x, G_{j,m} = y) = \frac{\lambda_{i,m}^x}{x!} \exp(-\lambda_{i,m}) \cdot \frac{\lambda_{j,m}^y}{y!} \exp(-\lambda_{j,m}) \quad (2.7)$$

The parameters  $\lambda_{i,m}$  and  $\lambda_{j,m}$  are the expected goals scored by team i and j, respectively. These  $\lambda$ 's are estimated to represent the abilities of the opposing teams. The following formulas demonstrate this:

$$\begin{aligned} \lambda_{i,m} &= \exp(c + (r_i + h) - r_j) \\ \lambda_{j,m} &= \exp(c + r_j - (r_i + h)) \end{aligned} \quad (2.8)$$

Here h stands for the home effect, c is the intercept,  $r_i$  and  $r_j$  are the relative abilities of team i, the home team, and team j, the away team, respectively. The ability parameters, intercept and home team advantage are estimated by a maximum likelihood approach which takes the following formula:

$$L = \prod_{m=1}^M \prod_{i,j \in \{1, \dots, T\}} \left( \frac{\lambda_{i,m}^{g_{i,m}}}{g_{i,m}!} \exp(-\lambda_{i,m}) \cdot \frac{\lambda_{j,m}^{g_{j,m}}}{g_{j,m}!} \exp(-\lambda_{j,m}) \right)^{y_{ijm} \cdot W_{time,m}} \quad (2.9)$$

In the formula above, the parameter  $m \in \{0, \dots, M\}$  is an index for the match, and the parameter T encompasses all the different teams. The random variable  $y_{ijm}$  is equal to 1, if i and j stand for the home and away team, respectively, in match m, else  $y_{ijm}$  is equal to 0. The parameters  $g_{i,m}$  and  $g_{j,m}$  are the actual observed goals scored by each team.  $W_{time,m}$  serves as a decay function to reflect a smooth time depreciation effect. This function is given by  $W_{time,m}(x_m) = \frac{1}{2} \frac{x_m}{Half\ period}$ , which implies that a match played *Half period* days ago contributes half as much to the likelihood function compared to matches played in the present time.

Ley et al. (2019) reports that the optimal *Half period* was 360 days for the *Independent Poisson* model and 390 days for the *Bivariate Poisson* model for the data coming from the English Premier League.

### Bivariate Poisson

Ley et al. (2019) extended the basic independent Poisson by adding a direct correlation coefficient between the scores, based on the bivariate Poisson model suggested by Karlis and Ntzoufras (2003). The goals of the opposing teams i and j are now modelled by the random variables  $G_{i,m} = X_{i,m} + X_C$  and  $G_{j,m} = X_{j,m} + X_C$  for match m, where  $X_{j,m}$ ,  $X_{i,m}$  and  $X_C$  follow a Poisson distributed with the respective intensity parameters  $\lambda_{j,m}$ ,  $\lambda_{i,m}$  and  $\lambda_C$ . The expected goals,  $E(G_{i,m}) = \lambda_{i,m} + \lambda_C$  and  $E(G_{j,m}) = \lambda_{j,m} + \lambda_C$

, scored by team  $i$  and  $j$  respectively, now take the correlation term into account. The parameter  $\lambda_C$  is the correlation between the scores of both opposing teams. The  $\lambda_{i,m}$  and  $\lambda_{j,m}$  parameters are still estimated to represent the abilities of the opposing teams, with the same equation as 2.8. The following formula gives the joint density distribution of the bivariate Poisson:

$$P(G_{i,m} = x, G_{j,m} = y) = \frac{\lambda_{i,m}^x \lambda_{j,m}^y}{x!y!} \exp(-(\lambda_{i,m} + \lambda_{j,m} + \lambda_C)) \sum_{k=0}^{\min(x,y)} \binom{x}{k} \binom{y}{k} k! \left( \frac{\lambda_C}{\lambda_{i,m} \lambda_{j,m}} \right)^k \quad (2.10)$$

The next formula gives the appropriate likelihood function:

$$L = \prod_{m=1}^M \prod_{i,j \in (1, \dots, T)} \left( \frac{\lambda_{i,m}^{g_{i,m}} \lambda_{j,m}^{g_{j,m}}}{g_{i,m}! g_{j,m}!} \exp(-(\lambda_{i,m} + \lambda_{j,m} + \lambda_C)) \sum_{k=0}^{\min(g_{i,m}, g_{j,m})} \binom{g_{i,m}}{k} \binom{g_{j,m}}{k} k! \left( \frac{\lambda_C}{\lambda_{i,m} \lambda_{j,m}} \right)^k \right)^{y_{ijm} \cdot W_{time,m}} \quad (2.11)$$

The parameters of the likelihood function are interpreted in the same way as in equation 2.9. If  $i$  and  $j$  stand for the home and away team,  $y_{ijm}$  is still equal to 1 and else to 0. The observed goals are denoted as  $g_{i,m}$  and  $g_{j,m}$ .  $W_{time,m}$  is still the weight function for the time decay effect.

### 2.3.4 Weibull count model

The fourth model comes from the article by Boshnakov et al. (2017). Just like the Poisson based models, the fourth approach is also goal-based, which means that it models the goal distributions of the opposing teams. However, unlike the Poisson based models, the fourth approach assumes that the goals now follow a Weibull count distribution. The following formula gives the Weibull count density distribution:

$$p(X(t) = x) = \sum_{j=x}^{\infty} \frac{(-1)^{x+j} (\lambda t^c)^j \alpha_j^x}{\Gamma(cj+1)} \quad (2.12)$$

$$\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt$$

In the function above  $\alpha_j^0 = \Gamma(cj+1)/\Gamma(j+1)$ , for  $j = 0, 1, 2, \dots$ , and  $\alpha_j^{x+1} = \sum_{m=x}^{j-1} \alpha_m^x \Gamma(cj - cm + 1)/\Gamma(j - m + 1)$ , for  $x = 0, 1, 2, \dots$  and  $j = x + 1, x + 2, x + 3, \dots$ . The parameter  $\lambda$  can be seen as the scoring rate per match, which is comparable to the  $\lambda$  used in the Poisson distribution. The  $c$  parameter is the shape of the distribution and allows

for extra flexibility. The dispersion of the Weibull count is denoted by the hazard  $h(t) = \lambda c t^{c-1}$ . Note that if  $c = 1$ , the Weibull count distribution is equal to a Poisson distribution, since the dispersion is then equal to  $\lambda$ . If  $c > 1$  the distribution is over-dispersed and if  $c < 1$  the distribution is under-dispersed. For the estimation of this density function, we used the R package *Countr* (Baker et al., 2016). Figure 2.3 illustrates the differences between the Poisson and Weibull count distribution, fitted to the goals of the home and away teams in the English Premier League, and clearly shows that the Weibull count distribution has a better fit, in particular for lower values of the observed goal values.

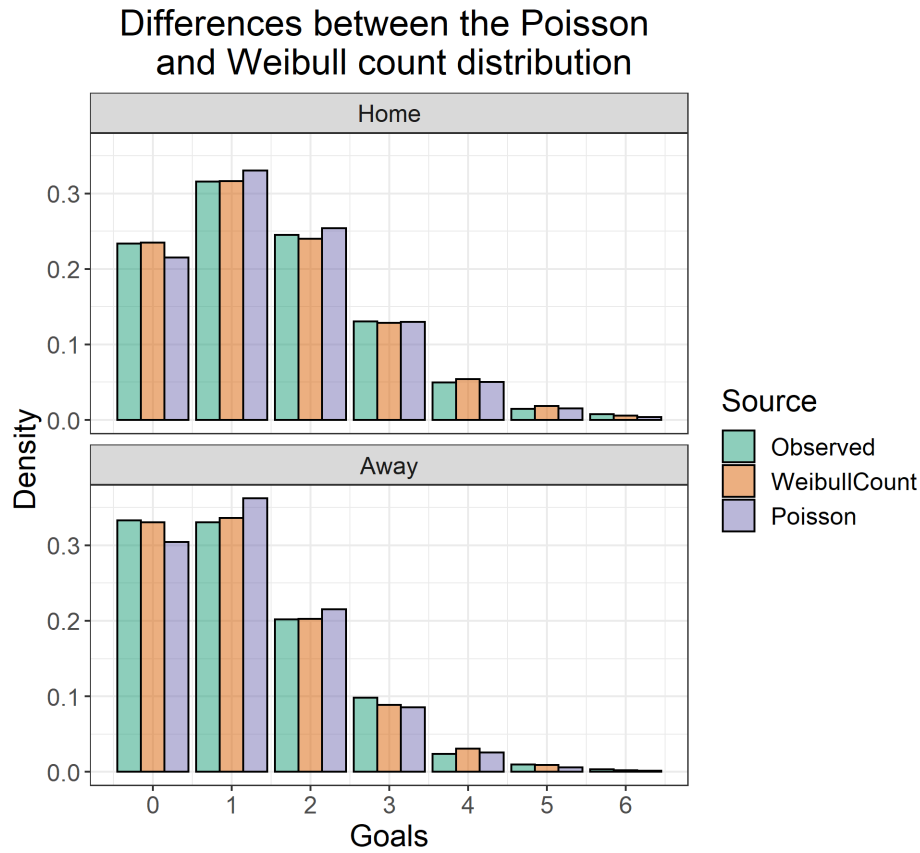


Figure 2.3: Difference between Weibull Count and Poisson distribution

Boshnakov et al. (2017) used the Weibull count distribution with a copula dependence between the goal distributions of the home and away teams to create a bivariate prediction model for the outcomes of football matches. A copula  $C$  is a multivariate distribution, for which the marginal distributions are uniform between  $[0, 1]$ . Here a copula  $C$  is used to glue the marginal cumulative distributions, of the home and away goals, together. The following formula illustrates this:

$$C(F_1(y_1), F_2(y_2)) = F(y_1, y_2) \quad (2.13)$$

Equation 2.14 shows how a copula  $C$  can be used to glue the marginal cumulative distributions,  $F_1(y_1)$  and  $F_2(y_2)$ , together from the joint density distribution  $F(y_1, y_2)$ . Boshnakov et al. (2017) reports that a Frank's copula provided the best fit to the data, and will thus be the copula of choice. The formula below gives the equation for Frank's copula:

$$C(u, v) = -\frac{1}{k} \log \left( 1 + \frac{(e^{-ku} - 1)(e^{-kv} - 1)}{e^{-k} - 1} \right) \quad (2.14)$$

In this formula, the parameter  $k$  is the dependence between  $u$  and  $v$ , which in our case are the marginal cumulative distributions. The coefficients of the bivariate Weibull count model are then estimated by using a maximum likelihood approach. The maximum likelihood takes the following formula:

$$L(k, \alpha, \beta, c) = \prod_{k \in (k: t_k < t)}^M (C(F_1(y_{1i}), F_2(y_{2i})) - C(F_1(y_{1i} - 1), F_2(y_{2i})) - C(F_1(y_{1i}), F_2(y_{2i} - 1)) + C(F_1(y_{1i} - 1), F_2(y_{2i} - 1))) * e^{-\epsilon(t - t_k)} \quad (2.15)$$

Here  $F_1$  and  $F_2$  are again the cumulative Weibull count distributions for the home and away teams respectively. The cumulative Weibull count distributions require two variables, the rate parameter  $\lambda$  and shape parameter  $c$ . For every home team the parameter  $\lambda$  is defined as:  $\log(\lambda_i) = \alpha_i + \beta_i + \gamma$ , with  $\alpha_i$  equal to the attack strength of team  $i$ ,  $\beta_i$  equal to the defence strength and  $\gamma$  the effect of playing as the home team. For every away team the parameter  $\lambda$  is defined as:  $\log(\lambda_j) = \alpha_j + \beta_j$ . The shape parameters  $c_h$  and  $c_a$ , for the home and away team respectively, are expected to be constant. The parameters  $y_{1i}$  and  $y_{2i}$  are the observed goals.

The maximum likelihood procedure thus has to estimate  $2T + 4$  parameters, namely  $\alpha$ 's and  $\beta$ 's for all the teams  $T$ ,  $c_h$  and  $c_a$  which are the shape parameters of the home and away team,  $h$  the home team effect and  $k$  the dependence term from Frank's copula.

Finally, the last unknown parameter  $\epsilon$  models the time depreciation effect. In formula 2.16, the term  $t - t_k$  stands for the difference between the time  $t_k$ , when a historical match  $m$  was played, and the current time  $t$ , expressed as number of days. The value for  $\epsilon$  is obtained by maximizing the next function:

$$T(\epsilon) = \sum_{k=1}^N (\delta_k^H \log p_k^H + \delta_k^A \log p_k^A + \delta_k^D \log p_k^D + \gamma_k^{02.5} \log p_k^{02.5} + \gamma_k^{U2.5} \log p_k^{U2.5}) \quad (2.16)$$

In this formula,  $\delta_k^H = 1$  if the home team won the match  $k$ ,  $\delta_k^A$  and  $\delta_k^D$  are interpreted analogously.  $p_k^H$ ,  $p_k^A$  and  $p_k^D$  are the maximum likelihood estimates for the probability of a home win, home loss and draw, respectively, in match  $k$ . The parameter  $\gamma^{O2.5} = 1$ , if there are more than 2.5 goals in match  $k$ , and  $\gamma^{U2.5} = 1$  if there are fewer than 2.5 goals.  $p_k^{O2.5}$  and  $p_k^{U2.5}$  are the maximum likelihood estimates for the probability of observing more or fewer than 2.5 goals in match  $k$ . Boshnakov et al. (2017) reports that a value of  $\epsilon = 0.002$  is reasonable. This value implies that the weight of a match played 500 days ago contributes only for  $\frac{1}{e} (\approx 0.37)$  to the maximum likelihood estimation, compared to matches played in the present time.

### 2.3.5 Machine learning models

The next models come from the article by Baboota and Kaur (2019). In this article, Baboota and Kaur (2019) compared five different machine learning approaches in their ability to predict football match outcomes. The considered models are the naive Bayes, linear support vector machine, RBF support vector machine, random forest and gradient boosting. From those five models, the random forest and gradient boosting outperform the other models and will thus be included in this master dissertation. Firstly, we will go over the specific data used for these models.

#### Data

The models used by Baboota and Kaur (2019) use a set of highly informative engineered features. The first set feature involves the different FIFA-ratings of the opposing teams, namely the attack, midfield, defence and overall rating. The EA Sports company constructed these ratings to be used in the FIFA game series. The ratings can be scraped from the FIFA index database (<https://www.fifaindex.com/>). The statistical model then uses the difference between the respective home and away team ratings, given by the following formula,  $R_i = R_i^H - R_i^A$ , for  $R_i \in$  (attack, midfield, defence, overall), and  $R^H$  and  $R^A$ , standing for the home and away team respectively. This trend, of using the difference between the home and away team, can be assumed for all the following continuous features.

The next feature is the goal difference. This feature is a sum of the goal differences from the preceding matches. The goal differences for the  $k^{th}$  match is given by  $GD = \sum_{j=1}^{k-1} GS_j - \sum_{j=1}^{k-1} GC_j$ , where  $GS$  and  $GC$  stand respectively for the goals scored and conceded.

The third set of features incorporate information of a team's recent performance. Specifically, the feature contains the average the number of corners, shots on target

and goals of the past  $k$  matches. the formula for the  $j^{th}$  match is given by,  $\mu_j^i = (\sum_{p=j-k}^{j-1} \mu_p^i)/k$ , with  $\mu^i \in (Corners, Shotsontarget, Goals)$ . The hyperparameter  $k$  requires tuning. Baboota and Kaur (2019) report that  $k = 6$  is optimal for the random forest and gradient boosted model. The data for this feature can be obtained from the Football UK website (<https://www.bbc.com/sport/football/>).

The fourth set of features are also engineered to represent the recent performance of a team. Namely two indicators, the streak and weighted streak, were used for this purpose. The streak is meant to capture the recent increase/decrease in performance of a team. The streak is calculated by giving a score to each match result, and then the mean score of the  $k$  preceding matches is taken. Note that  $k$  is the same hyperparameter as mentioned above. The scores follow the 3-1-0 rule. A win grants a score of 3, a draw a score of 1 and a loss a score of 0. The weighted streak is calculated by updating the streak with a time depreciation effect. The oldest observation ( $j-k$ ) gets a value 1, and the most recent observation ( $j-1$ ) gets a value  $k$ . The following formulas give the streak ( $\delta$ ) and weighted streak ( $\omega$ ) for the  $j^{th}$  match:

$$\delta_j = \left( \sum_{p=j-k}^{j-1} res_p \right) / 3k \quad (2.17)$$

$$\omega_j = \sum_{p=j-k}^{j-1} 2^{\frac{p-(j-k-1)}{k}} \frac{res_p}{3k(k+1)}$$

Here  $res_p \in (0, 1, 3)$  stands for the score given to the outcome of the match. The term  $3k$  ensures that the streak is normalised between 0 and 1. In the weighted streak formula, the term  $p - (j - k - 1)$  ensures that the oldest observation ( $j-k$ ) gets a value 1, and the most recent observation ( $j-1$ ) gets a value  $k$ , and the term  $3k(k+1)/2$  now ensures the normalisation.

The last feature, named form, aims to display the performance of a team during an individual match. Just like the streak feature, a team's form gives information about the recent performance of a team. However, contrary to the streak feature, the form feature displays a team's performance relative to the opposing team. After every match, the form values of a team are updated. The form feature rewards a team with low form, defeating a team with high form, more, and vice-versa. If the match ends in a draw, the form of the weaker team increases, while the form of the stronger team decreases. The initial value of the form of each team is 1, and updated after every match by the following formula:

In the case of a win or loss:

$$\begin{aligned}\epsilon_j^\alpha &= \epsilon_{j-1}^\alpha + \gamma \epsilon_{j-1}^\beta \\ \epsilon_j^\beta &= \epsilon_{j-1}^\beta - \gamma \epsilon_{j-1}^\alpha\end{aligned}\tag{2.18}$$

In the case of draw:

$$\begin{aligned}\epsilon_j^\alpha &= \epsilon_{j-1}^\alpha - \gamma (\epsilon_{j-1}^\alpha - \epsilon_{j-1}^\beta) \\ \epsilon_j^\beta &= \epsilon_{j-1}^\beta - \gamma (\epsilon_{j-1}^\beta - \epsilon_{j-1}^\alpha)\end{aligned}\tag{2.19}$$

With  $\epsilon_j^\alpha$  and  $\epsilon_j^\beta$  equal to the form of team  $\alpha$  and team  $\beta$  respectively, in the  $j^{th}$  match. The parameter  $\gamma$  is referred to as the stealing fraction, e.g. if team  $\alpha$  wins over team  $\beta$ , it steals a fraction  $\gamma$  of team  $\beta$ 's form. If the match ends in a draw, the weaker team gets a positive update and the stronger team a negative update, proportional to the difference in their respective forms. Baboota and Kaur (2019) report that the value of the stealing fraction  $\gamma$  is equal to 0.33.

Feature	Equation
Form	$\epsilon_j^\alpha - \epsilon_j^\beta$
Streak	$\delta_j^\alpha - \delta_j^\beta$
Weighted streak	$\omega_j^\alpha - \omega_j^\beta$
Corners	$\mu_{corners,j}^\alpha - \mu_{corners,j}^\beta$
Goals	$\mu_{goals,j}^\alpha - \mu_{goals,j}^\beta$
Shots on target	$\mu_{shots,j}^\alpha - \mu_{shots,j}^\beta$
Goal difference	$GD_j^\alpha - GD_j^\beta$
FIFA <sub>attack</sub>	$FIFA_{attack,j}^\alpha - FIFA_{attack,j}^\beta$
FIFA <sub>midfield</sub>	$FIFA_{midfield,j}^\alpha - FIFA_{midfield,j}^\beta$
FIFA <sub>defence</sub>	$FIFA_{defence,j}^\alpha - FIFA_{defence,j}^\beta$
FIFA <sub>overall</sub>	$FIFA_{overall,j}^\alpha - FIFA_{overall,j}^\beta$

Table 2.2: Features used in (Baboota and Kaur, 2019)

Note that many of the variables are only calculated after  $k = 6$  weeks. Thus a burn-in period of 6 weeks is needed to get reliable estimates of the features. Baboota and Kaur (2019) reports that the highest performance is found by the random forest and gradient boosting algorithm. The following subsections give some detail about these models.

### Random Forest

Random forest models, first proposed by Breiman (2001), is a method that ensembles decision trees. These decision trees are used for regression or classification

problems. The decision trees work by partitioning the predictor space ( $X_1, \dots, X_p$ ) repeatedly into multiple, non-overlapping, distinct regions ( $R_1, \dots, R_j$ ). Usually, this is done with binary splits, intending to find homogeneous response values within the same region and heterogeneous response values between them (Groll et al., 2019). Figure 2.4 illustrates how this progress works with a dendrogram. The decision trees predict values by averaging over the response values, for regression trees, or using the majority vote, for classification approaches. On figure 2.4, we observe four different splits, that define five distinct regions. However, many more splits could be used, which makes decision trees prone to overfitting.

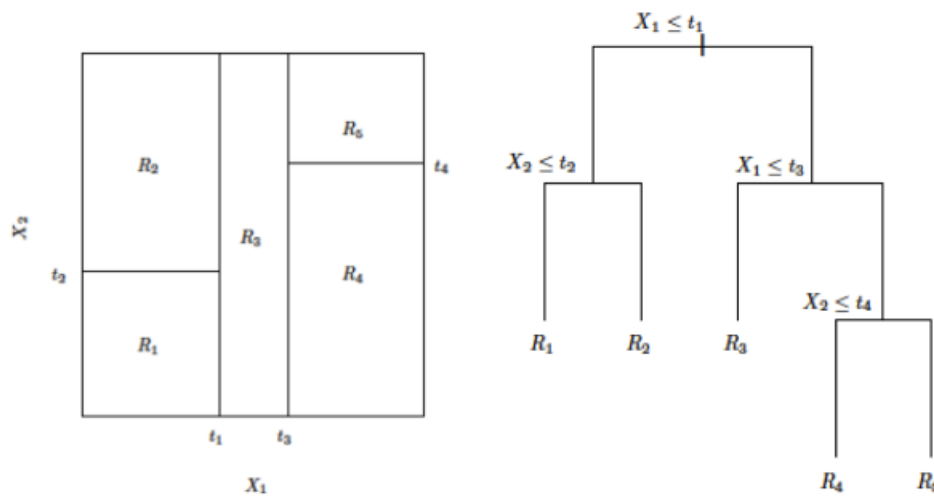


Figure 2.4: Example of decision trees. source: Ley, C. (2020). Big data science, Chapter 4 - Tree-based method [PowerPoint slides], course material University Ghent.

Random forests circumvent this issue of overfitting by using a bagging approach. This bagging approach first creates multiple bootstrapped datasets, and for each such dataset, a decision tree is created. The random forest model aggregates the predictions over the multiple individual decision trees. Combining all of these individual decision trees has the advantage of making the predictions unbiased and reducing the variance among them. Another improvement of random forests over decision trees is that random forest only uses a random subsample of the original predictor space. This reduces the correlation between the multiple decision trees over different bootstrapped datasets. The selection of the random subset of predictors is usually done in two steps. First, the decision trees only use a random subset of original predictor space, and secondly, at every node, a random subset of the predictor space of the decision tree is used to find the best split.

We use the programming language Python version 3.7.6 with the machine learning software sklearn to create the random forest classifier. The table below shows the specific hyperparameters used.

The first hyperparameter, name *criterion*, measures the quality of a split. It measures the amount of entropy or impurity that each feature removes when it is used in a node. The second hyperparameter, named *max features*, gives the maximum number of features that can be randomly selected at every node to find the best split. A value of log2 indicates that if the predictor space of the decision tree has eight features, three randomly selected features can be used in every node. The third hyperparameter, *min samples leaf*, indicate the minimum number of samples needed to form a distinct region (or leaf). A value of two would mean that a split will only occur if the distinct regions, that are formed after the splitting, each have at least two samples in them. The fourth hyperparameter, *min sample split*, is similar. A value of 100 indicates that 100 samples are needed to split an internal node. Both the min samples leaf and min samples split are used to smooth out the data and attempt to reduce the probability of learning noise. The last hyperparameter, *N estimators*, indicates the number of bootstrapped datasets, with their respective decision trees, that are aggregated into a final random forest. Usually, the value for this hyperparameter is set high enough, so that the predictions obtain the feature of being unbiased, and low enough to reduce computational demands.

Hyperparameter	Description	Value used
Criterion	The function used to measure the quality of a split	gini
Max features	The maximum number of features for the best split	log2
Min samples leaf	The minimum number of samples at region	2
Min samples split	The minimum number of samples needed to split an internal node	100
N estimators	The number of independent trees in the forest	150

Table 2.3: Hyperparameters used in the random forest of (Baboota and Kaur, 2019)

## Gradient boosting

Just like the random forest, gradient boosting is a technique that ensembles decision trees. However, unlike the random forest, it does not use a bagging approach, but alternatively, it uses the boosting principle. The difference between a bagging and a boosting approach is that in a boosting approach, the decision trees are not trained independently from each other. Instead, the decision trees are trained in sequence on the entire data set. The model's accuracy increases with each iteration.

We again use the programming language Python version 3.7.6 with the open-source software XgBoost to create the gradient boosted classifier (Chen and Guestrin, 2016). The table below shows the specific hyperparameters used.

The first hyperparameter, *gamma*, is the minimum loss reduction required to make a further partition on an internal node of the tree. A value of 0.5 means that the loss needs to be reduced by 0.5 before splitting a node. A large value of *gamma* means that the algorithm will be more conservative. The second hyperparameter is the *learning rate*. In gradient boosting algorithms, trees are added to the model sequentially. The newly added trees aim to correct the residual errors in the older trees, which can cause the model to overfit. The learning rate aims to shrink the weight of these new trees. Lower values for the learning rate means that the boosting algorithm stays more conservative. The third hyperparameter, *max depth*, indicates the maximum depth of a tree. A value of three means that only three consecutive splits can be used for each tree. Larger values for this hyperparameter would make the model more complex and likely to overfit. The final hyperparameter, *lambda*, is a ridge regularization term on weights. It is used to reduce the probability of overfitting, and a large value will make the model more conservative.

Hyperparameter	Description	Value used
Gamma	The minimum loss reduction required to make a split	0.5
Learning rate	Boosting learning rate	0.025
Max depth	Maximum depth per tree	3
Lambda	Controls the regularization strength	0.6

Table 2.4: Hyperparameters used in the random forest of (Baboota and Kaur, 2019)

### 2.3.6 Hybrid random forest model

The final model, proposed by Groll et al. (2019), uses a hybrid modelling approach. This hybrid approach combines both the bivariate Poisson, from section 2.3.3, and a random forest regression model.

Many covariates are gathered and concatenated into a single, highly informative dataset. One of those covariates are the ability parameters estimated by the maximum likelihood approach of the bivariate Poisson, see section 2.3.3 and figure 2.2. The random forest approach then uses all these covariates to model the goals scored by both the home and away team. The Skellam distribution, see section 1.2.1, uses these estimated goals as intensity parameters,  $\lambda_1$  and  $\lambda_2$ , to derive the probabilities for each outcome class.

Figure 3.2, which comes from the original article by Groll et al. (2019), shows that the ability parameter is highly informative and essential for the performance of the model. The implementation of this model was done in the programming language R version 3.6.1 by using the *ctree* function from the *party* package. The package *mlr* is used to specify the hyperparameters of the hybrid random forest model. As mentioned in section 2.3.5, the random forest model uses a bagging approach, which bootstraps

the data and creates independent decision trees on each of them. If the number of trees used is large enough, the predictions are considered to be unbiased, and the variance is greatly reduced. Predictions are made by averaging over the predictions from the individual trees. The hyperparameter  $B$ , for the number of trees, does not require tuning as long as it is sufficiently large, so it gets the feature of unbiasedness (Probst and Boulesteix, 2017). Groll et al. (2019) reports that the number of trees,  $B = 500$ , is sufficient.

## Data

We will now go over all the covariates used in the dataset from the original article by Groll et al. (2019). The dataset used in the original article was created for the national teams. Unfortunately, due to time constraints, we were not able to reproduce the complete dataset for the English Premier League.

Variable	Description
<b>Sportive variables</b>	
ODDSET probability	The bookmaker odds converted to winning probabilities
FIFA rank	The FIFA rank of the national teams based on the last four years
ELO rating	The ELO rating of the national teams
ABILITIES	The ability of a team estimated by the bivariate Poisson model, see 2.3.3
<b>Economic variable</b>	
GDP per capita	The increase in GDP between 2002-2004 for every nation
Population	The nation's population size
<b>Team structure variables</b>	
Host	Dummy indicating if a national team is the hosting country
Continent	Dummy indicating if the national team is on the same continent as the host
Confederation	A categorical variable indicating the confederation of a team
<b>Home advantage variables</b>	
Maximum teammates	The maximum number of team mates for each squad
Second maximum teammates	The second maximum number of team mates for each squad
Age	The average age of each squad
Champion league players	Number of Champion league players in a squad
Legionnaires	Number of players in a squad, playing in clubs abroad
<b>Coach variables</b>	
Age	The age of a team's coach
Tenure	The tenure of a team's coach
Nationality	Dummy indicating if a coach has the same nationality as the team

Table 2.5: Variables used in the hybrid random forest of Groll et al. (2019)

Since we were not able to replicate all the essential features of the dataset during the timeframe of this master dissertation, we decided to focus on a subset of the most informative ones based on figure 3.2, which comes from the original article by Groll et al. (2019).

First, all of the sportive variables were included since those were the most informative for the hybrid random forest model. The ability variable is by far the most important feature. This variable was replicated by using the method described in the bivariate Poisson, section 2.3.3. The second most informative feature are the ELO ratings of

the teams. These were obtained through the goal-based ELO model from Hvattum and Arntzen (2010), section 2.3.2. For the FIFA rank data, we used the overall FIFA rank from the data set of Baboota and Kaur (2019), section 2.3.5. The odds data were obtained from the website <http://www.football-data.co.uk/> and specifically, we used the odds from the bet365 company.

From the home advantage variables, only the age variable was used, since all others were deemed irrelevant for the English Premier League. The age variable, which is the average age of the players in a team, was scraped from the website <https://www.worldfootball.net/>. On that website, the dates of birth of all the English Premier League players, given a season, were listed. Through their dates of birth and the current season, their age could be reverse-engineered.

Most of the economic and team structure variables were not relevant for the English Premier League, e.g. the GDP per capita would be the same for every team and were thus not included. The data concerning the coaches variables were not included due to time constraints. These variables were only mildly informative for the hybrid random forest model, so no significant performance loss is expected.

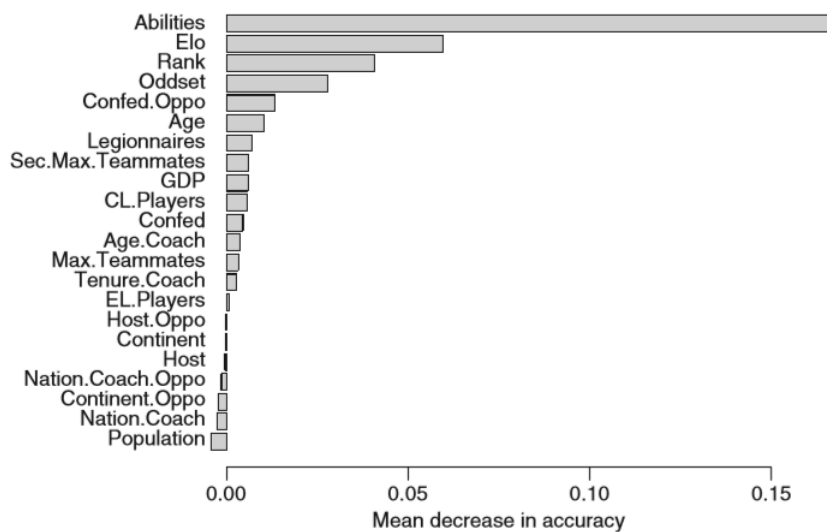


Figure 2.5: Barplot of the variable importance in the original hybrid random forest of Groll et al. (2019)

We also incorporate the random forest model, discussed in Groll et al. (2018), which is identical to the hybrid random forest described above, but without incorporating the information of the abilities of the teams. In order to differentiate between the random forest of Baboota and Kaur (2019) and Groll et al. (2018), we will refer to both models as Random forest<sub>1</sub> and Random forest<sub>2</sub> respectively. The reason for incorporating the basic random forest is to illustrate how much predictive power is improved by using the hybrid approach.

## CHAPTER 3

# RESULTS

This chapter will discuss the results of this master dissertation. The scoring rules, mentioned in section 1.4, evaluate and rank the models per scoring rule. The average rank then defines the overall best performing model. After that, we will go over the specific results for each model separately and reference the values for the scoring rules reported in the original articles.

### **3.1 General results**

Table 3.1 shows the average scores (standard deviation) for all the models per scoring rule. Lower scores translate to a better predictive performance of the model. This table shows that models with more sophisticated techniques performed better on average compared to more simplistic models. The models from Baboota and Kaur (2019) significantly outperform the other models, with the gradient boosting model the best among them. Nonetheless, the most complex model, the hybrid random forest, only mildly outperforms the goal-based ELO model, which is considered to use a more simplistic approach.

Another example of this is the Weibull count model, which, despite having a better fit to the observed goals, performs worse than the Poisson-based alternatives. The Poisson-based and ELO models are very similar in their performance. More parsimonious models and models with fewer assumptions are seemingly better at predicting football match outcomes.

Model	Ignorance score	Brier score	RPS Score	Mean rank
Gradient boost	0.933 (0.361)	0.353 (0.175)	0.132 (0.073)	1.000
Random forest <sub>1</sub>	0.964 (0.350)	0.363 (0.172)	0.141 (0.072)	2.000
Hybrid random forest	1.387 (0.729)	0.571 (0.353)	0.190 (0.135)	3.000
Random forest <sub>2</sub>	1.408 (0.723)	0.581 (0.350)	0.195 (0.136)	4.000
Goal-based ELO	1.409 (0.752)	0.581 (0.359)	0.197 (0.147)	5.000
Bivariate Poisson	1.417 (0.801)	0.583 (0.355)	0.198 (0.145)	6.333
Basic ELO	1.415 (0.749)	0.584 (0.357)	0.198 (0.147)	6.667
Independent Poisson	1.421 (0.762)	0.585 (0.355)	0.198 (0.141)	8.000
Plus-minus	1.460 (0.759)	0.605 (0.356)	0.207 (0.152)	9.333
Weibull count	1.464 (0.526)	0.610 (0.261)	0.203 (0.101)	9.667
Frequency	1.533 (0.395)	0.641 (0.193)	0.225 (0.089)	11.000
Uniform	1.585 (0.000)	0.667 (0.000)	0.235 (0.073)	12.000

Table 3.1: The mean (sd) score for each model per scoring rule and the average rank.  
Random forest<sub>1</sub> from Baboota and Kaur (2019) and Random forest<sub>2</sub> from Groll et al. (2018)

## 3.2 ELO based models

For the basic ELO, Hvattum and Arntzen (2010) reported a Brier score with a mean (standard deviation) of 0.627(0.248) and for the goal-based ELO a Brier score of 0.626(0.249). For these models, we respectively report brier scores of 0.584(0.357) and 0.581(0.359). Hvattum and Arntzen (2010) reports an ignorance score of 1.502(0.503) and 1.499(0.506) for the basic ELO and goal-based ELO, respectively. We report ignorance scores of 1.415(0.749) and 1.409(0.752). Our findings are likely to be better due to the burn-in of the first five weeks, which excluded these observations from the testing procedure. The reported scores are still within each other's standard deviations. Note that the deviations we report are from the individual observations to the mean, not the standard deviation of the mean, which we would specifically denote as the standard error.

The standard deviations reported in the original paper by Hvattum and Arntzen (2010) are lower compared to our findings. This might be due to the differences in the training protocol. Hvattum and Arntzen (2010) used two seasons to get initial ELO ratings, then five seasons were used to estimate the coefficients of the ordered logit regression, and finally, eight seasons were used to test on. The coefficients of the ordered logit regression from Hvattum and Arntzen (2010) were also updated after the prediction of a match.

Table 3.2 shows the coefficients of the ordered logit regression, estimated in this master dissertation, by the end of the season 2015.

Parameter	Value	Standard error	t-value
$\beta_{ELO}$	0.006	0.001	12.310
$\Theta_{A D}$	-0.912	0.068	-13.328
$\Theta_{D H}$	0.238	0.063	3.784

Table 3.2: Coefficients for the ELO ordered logit regression by the end of season 2015

### 3.3 Plus minus ratings

For the ordered logit regression based on the plus-minus ratings, Arntzen and Hvattum (2020) report a brier score of 0.614 and an ignorance score of 1.476. Our results report a brier score of 0.605(0.356) and an ignorance score of 1.460(0.759), which are thus within each other's magnitude. Surprisingly, we only included a more simplistic version of the plus-minus ratings, but the predictive performance does not seem to suffer. We have to note, however, that the original paper does not use a burn-in period.

Table 3.3 shows the coefficients of the ordered logit regression, estimated in this master dissertation, by the end of the season 2015.

Parameter	Value	Std. Error	t-value
$\beta_{PM}$	45.68	4.04	11.31
$\Theta_{A D}$	0.43	0.14	3.20
$\Theta_{D H}$	1.63	0.15	10.89

Table 3.3: Coefficients for the plus-minus ordered logit regression by the end of season 2015

### 3.4 Poisson based models

For the independent and bivariate Poisson, Ley et al. (2019) reports a ranked probability score of 0.1954 and 0.1953, respectively. Our findings are very similar, 0.1983(0.1407) and 0.1975(0.1452) for the independent and bivariate Poisson, respectively. The original article does not report any standard deviations, but the scores are well within the range of our observed findings. The protocol used by Ley et al. (2019) is identical to the protocol used in this master dissertation. First, the past two years and the first five weeks of the current season were used as initial data, and the data was updated after the prediction of a week. The only difference is that Ley et al. (2019) predicted from season 2008 till 2017, while we predicted from the season of 2008 till 2015.

### 3.5 Weibull count

The article for the Weibull count by Boshnakov et al. (2017) does not incorporate any of the scoring rules from section 1.4. Instead, it uses the Akaike information criterion (AIC) to compare the copula-based Weibull count model to a copula Poisson model, an independent Weibull count model and independent Poisson model. The lowest AIC was reported for the copula Weibull count model. Unfortunately, we were not able to confirm these findings, since the Weibull count did not outperform the independent Poisson.

### 3.6 Machine learning models

For the machine learning models, by Baboota and Kaur (2019), the article reports a ranked probability score of 0.2100 and 0.2137 for the gradient boosting model and random forest model respectively. However, Baboota and Kaur (2019) also include figure 3.3 in their article and from that figure, we can indicate that the ranked probability scores of both the gradient boosting and random forest are around 0.173. This information seems inconsistent in the original article. Our finding reports a much lower ranked probability score of 0.132(0.073) and 0.141(0.072), for the gradient boosting model and random forest model, respectively. In the original article, the matches are not estimated in a stepwise manner. The matches from 2005 to 2014 were used as training data, and the matches from 2015 to 2016 were used as test data. The stepwise addition of information might explain why our findings report lower ranked probability scores. Figure 3.1 shows the feature importance for both models.

### 3.7 Hybrid model

Finally, the article of the hybrid random forest, by Groll et al. (2019), reports a ranked probability score of 0.187. The original article does not report any standard deviations. Our findings are just a little higher 0.190(0.135), but again within the range of the standard deviations. For the random forest model, by Groll et al. (2018), a ranked probability score of 0.192 was reported. This within the magnitude of our findings of 0.195(0.136).

Groll et al. (2018, 2019) used multiple extra covariates, to which we had no access to and they had better data quality for them. This could explain the slightly higher ranked probability scores found in this master dissertation. For example, the FIFA

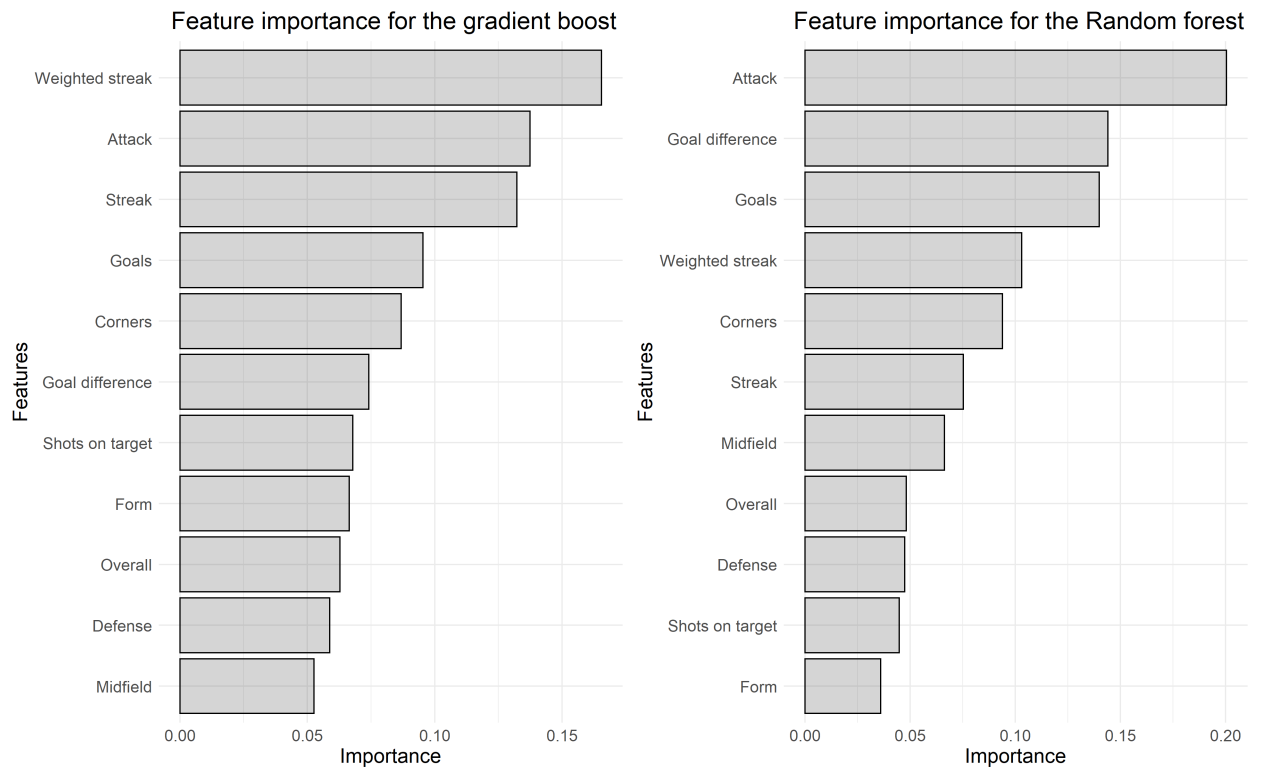


Figure 3.1: Feature importance of the gradient boost from Baboota and Kaur (2019)

ratings, obtained from Baboota and Kaur (2019), and the age of the players, which were obtained from <https://www.worldfootball.net>, proved challenging to merge with the datasets that contained other the covariates and the match outcomes. This might explain there lower feature importance of these features in our findings, compared to the original article (see figure 3.2).

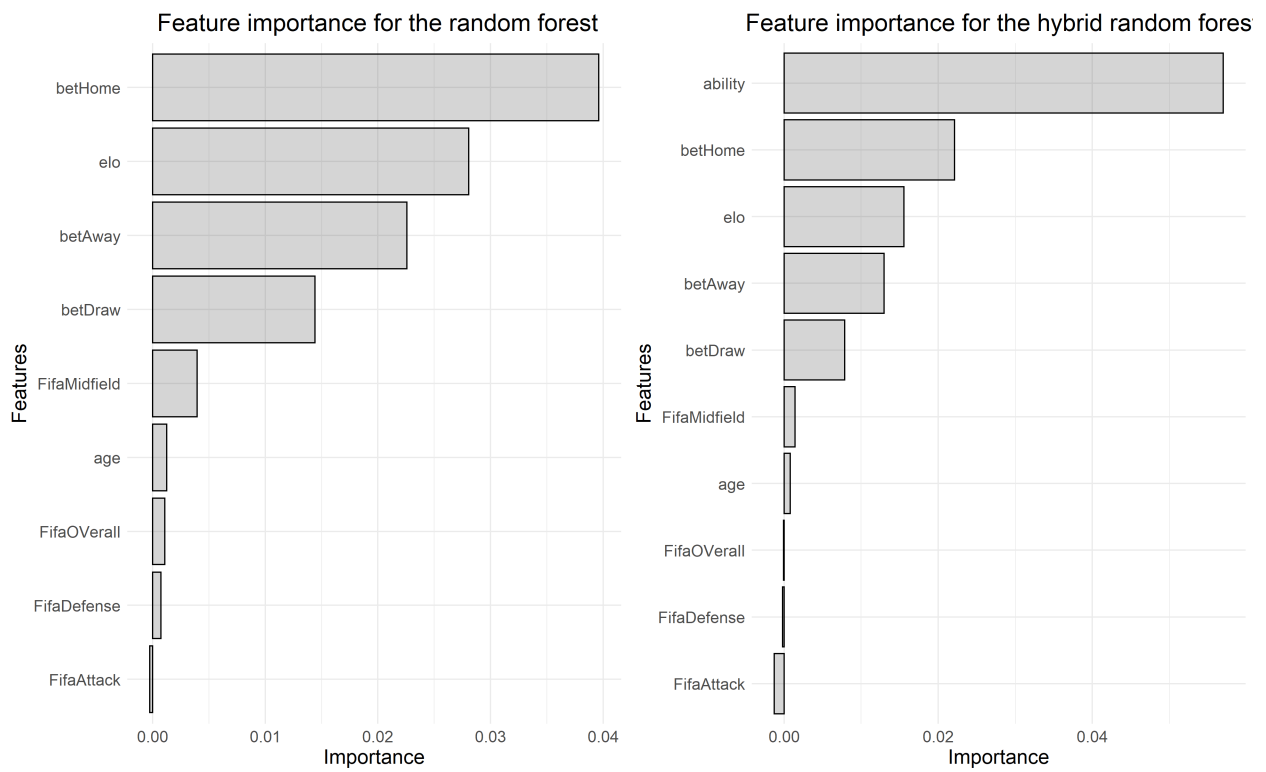


Figure 3.2: Feature importance of the random forest and hybrid random forest from Groll et al. (2018, 2019)

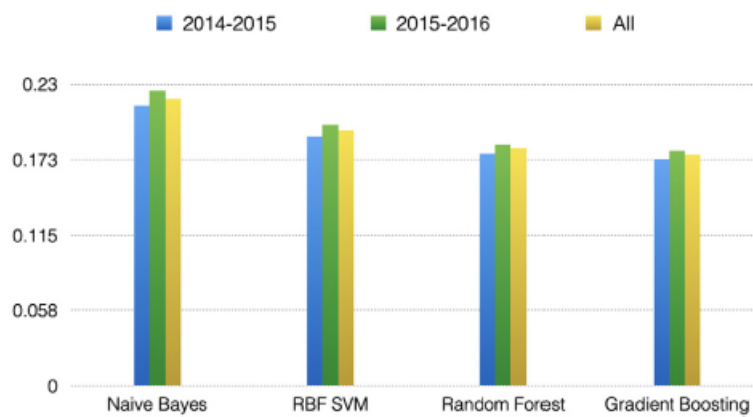


Figure 3.3: Ranked probability scores reported by Baboota and Kaur (2019)

### Comparison of different models over different scoring rules

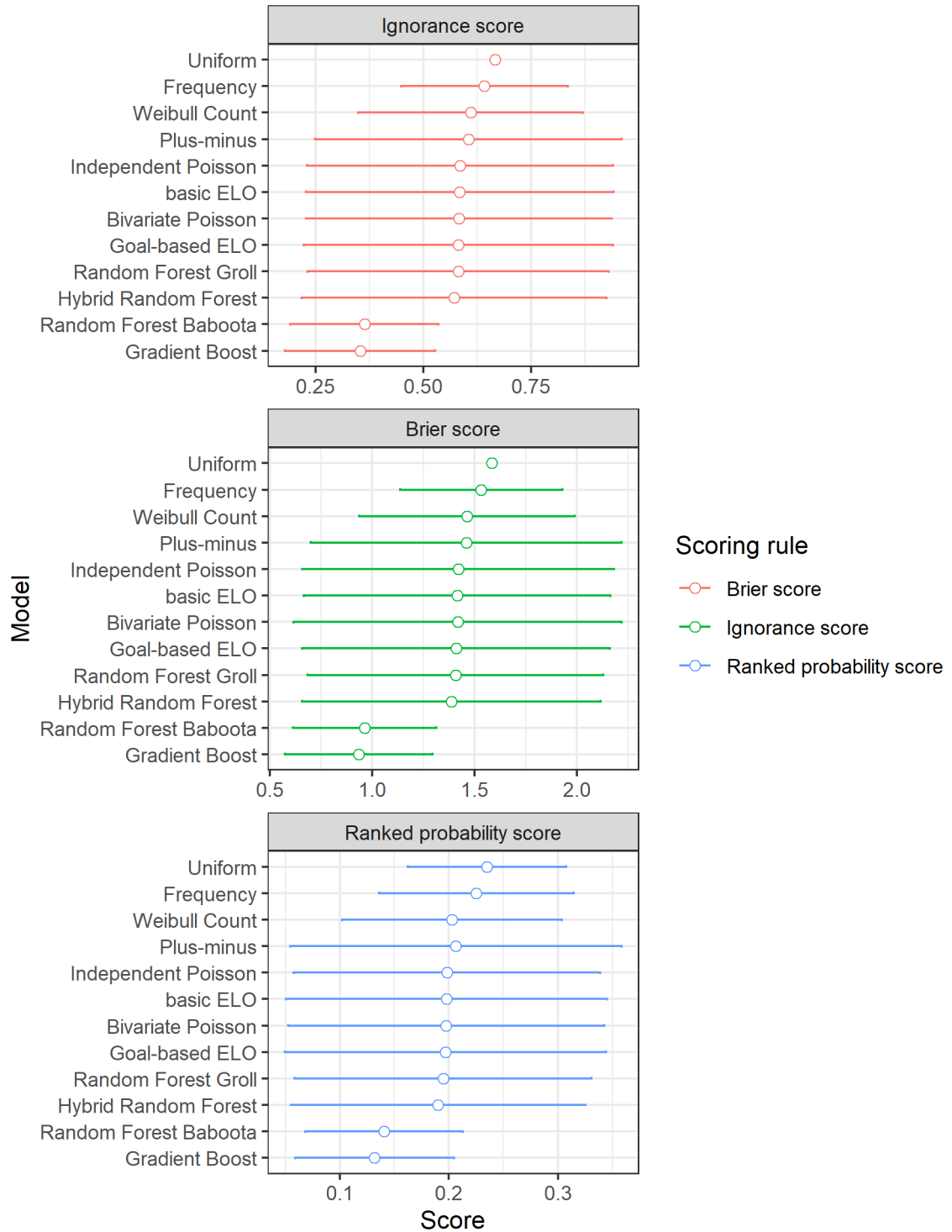


Figure 3.4: Comparison of the models over different scoring rules



## CHAPTER 4

# DISCUSSION

This chapter discusses the results found in this master dissertation. First, the controversial findings are discussed on a per model basis. After that, we give practical implications about the models and insights for future research.

### 4.1 ELO based models

For the ELO based models, our findings exceeded the expectations. We report better scores for both the Brier and ignorance score, compared to the original article of Hvattum and Arntzen (2010). As mentioned in 3.2, this might be due to the initial burn-in of five weeks. This burn-in procedure ensures that the models get some information on the new teams entering the league for that season. It is reasonable to assume that the estimations for the relative strength of the teams stabilise throughout a season. The predictions should thus increase in precision in the course of a season. Figure 4.1 shows the linear trends of the average of the scoring rule per week of a season. On this figure, the expected decrease is only observed for the ELO based model, but not for the models by Baboota and Kaur (2019).



Figure 4.1: Evolution of the scoring rules during the season of 2012

Another possible reason for our improved findings is the scale on which the ELO ratings are calculated. Hvattum and Arntzen (2010) calculated the ELO ratings over all the domestic leagues from England, namely the English Premier League, the Championship, League One and League Two. We specifically calculated only the ratings for the English Premier League on a per season basis, as described in section 2.3.2.

The ELO based models were very similar in performance compared to the Poisson based models. The goal-based ELO mildly outperforms the bivariate Poisson, and the basic ELO mildly outperforms the basic ELO. One possible reason for this might be the fact that the ELO based models indirectly got an additional season as data since we extended the initial protocol with one season to get initial reliable ELO ratings for the competing teams.

Although comparative studies are relatively absent in the current literature, Egidi and Torelli (2020) recently published an article where goal-based and result-based approaches were compared. The article compared two multinomial regression approaches to an independent and bivariate Poisson model, which were based on the models from Karlis and Ntzoufras (2003). Egidi and Torelli (2020) report that the result-based approaches, namely the two multinomial regressions, had higher Brier scores compared to the Poisson based alternatives. However, when predicting the World Cup tournament, Egidi and Torelli (2020) showed that the multinomial models outperformed the Poisson based models, by using the LOOIC (Leave-one-out cross-validation). The LOOIC is an approach for estimating the pointwise out-of-sample prediction accuracy from a fitted Bayesian model, by using the log-likelihood evaluated at the posterior simulations of the parameter value (Vehtari et al., 2017).

Egidi and Torelli (2020) point out that the multinomial models are less complex and that the outcome classes are easier to predict compared to the exact number of goals scored by each team. On the contrary, Ley et al. (2019) argued that goal-based models have a better performance since they incorporate additional information as the goal difference.

## **4.2 Plus minus ratings**

The plus-minus ratings are a top-down rating system, that aims to divide the success of a team's performance onto the players responsible for it, which results in ratings for the individual players. To put this into perspective, for the prediction of the season 2015 a total of 1415 unique players were observed over 2889 matches, which are further divided into 17820 segments. For every segment, there are only 23 non-zero

data points over the features. Datasets with such high dimensions often have to deal with a problem called the curse of dimensionality (Bellman, 2003).

This curse of dimensionality occurs when the number of predictors increases, because then the distance between the data points increases exponentially. A regularised approach is thus a necessity for the plus-minus ratings. A potential downside of these regularised approaches is that they trade the reduction in variance for an increase in bias, for the parameter estimations. The high amount of parameters to estimate, in combination with the introduced bias, might explain why the model based on the plus-minus ratings is underperforming compared to the other models. Player ratings may thus be noisier compared to direct team ratings and not enhance predictions.

In a more recent study, Arntzen and Hvattum (2020) used an improved version of the plus-minus ratings, described in Pantuso and Hvattum (2019), as a covariate in an ordered logit regression. The improved plus-minus rating system also includes the effect of red cards and the age of the players. Interestingly, it also takes the differences in competitive leagues and player similarities into account, which allows for parameter estimation on a much broader scale. The plus-minus ratings might also become more reliable with an increase in data.

Arntzen and Hvattum (2020) also report that when both the ELO and plus-minus rating are combined in an ordered logit regression, both the Brier and ignorance score are significantly reduced, and thus the predictive performance is enhanced.

Finally, figure 4.2 shows that the time depreciation effect used by Arntzen and Hvattum (2020) pays more attention to older observations compared to the models from Ley et al. (2019) and Boshnakov et al. (2017). These higher weights for older observations might not be optimal for the small scale comparative study of this master dissertation.

### **4.3 Weibull count**

Boshnakov et al. (2017) found that the Weibull count model had a better fit to the observed goal distribution compared to the Poisson based alternatives, according to the AIC criterion. However, this improved fit does not seem to be reflected in its predictive performance. A possible explanation for this could be that more parsimonious models are favoured over more complex ones. Ley et al. (2019) found that the independent and bivariate Poisson with one strength parameter outperformed their counterparts with two strength parameters, by using the RPS scoring rule.

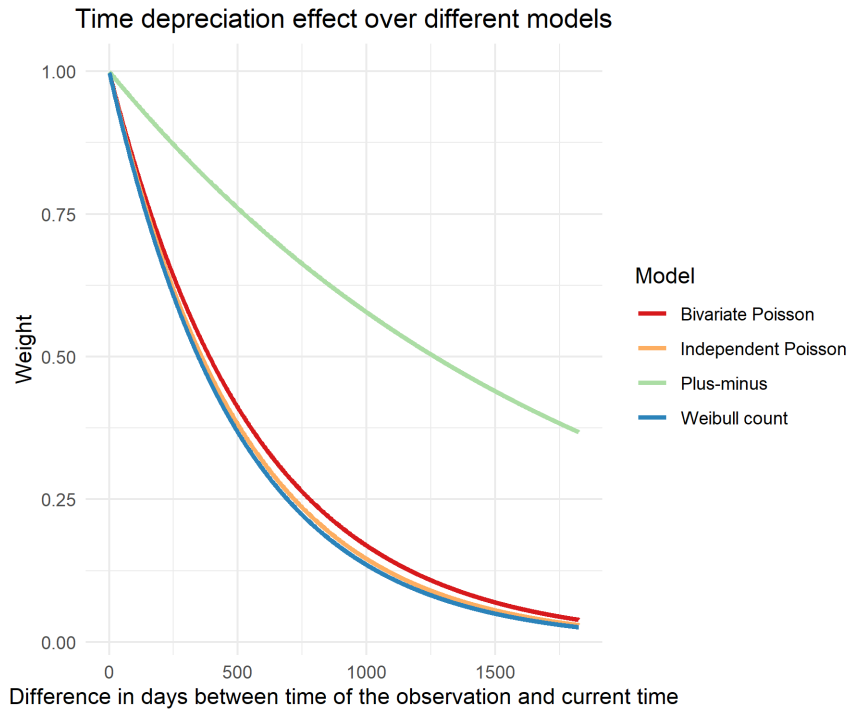


Figure 4.2: Time delay function for different models

Another possible reason for the lower performance of the Weibull count model is that it aims to estimate the parameters to have an optimal fit to the observed goal distributions, while the Poisson based models used by Ley et al. (2019) were mainly used to come up with an effective system to rank the opposing teams relative to each other. Ley et al. (2019) used some constraints in order to achieve this, e.g. the strengths of all the opposing teams were set to sum to zero. Boshnakov et al. (2017) did not use any constraints on the parameter estimations, which might make them unreliable and prone to learn noise.

## 4.4 Machine learning models

All the included models in this master dissertation that used more sophisticated machine learning techniques obtained better predictions. In particular, the models from Baboota and Kaur (2019) vastly outperformed the other models. The features for the models by Baboota and Kaur (2019) had better quality compared to the features used in the models by Groll et al. (2019), because for the features in the model by Groll et al. (2019) multiple datasets had to be merged by a unique key, which proved to be difficult. The features used in the models by Baboota and Kaur (2019) are published on the author's personal GitHub and were of high quality.

While most models are solely based on historical data, e.g. for the Poisson based models two years of data is needed to estimate the current abilities of the opposing teams, the features of Baboota and Kaur (2019) incorporate a mixture of features based on historical data (FIFA ranks) and features that capture the recent performance of a team. The importance of the features that capture the recent performances of teams might be exciting to investigate further since the scoring rules stayed relatively stable over the weeks during a season, see figure 4.1, and the variation is much smaller compared to the other models. The feature importance plots (see figure 3.1) also show that many of the features that capture the recent performance are essential for the accuracy of the model.

Interestingly, Baboota and Kaur (2019) give the suggestion to explore the gradient boosting and random forest regression approaches, which aim to estimate the goal distributions of each team. Groll et al. (2018) compared the random forest goal-based model, which predicts the intensity parameter  $\lambda$  for both opposing teams, to a result-based random forest model, which directly predicts the outcome classes, and found that the goal-based random forest models performed better than the result-based alternatives.

## **4.5 Hybrid models**

The feature importance plots on figure 3.2 show the importance of a key engineered feature, namely the abilities of the teams. This ability parameter significantly increased the accuracy of the model when included, which shows the importance of incorporating highly informative features.

Both comparative studies and studies that aim at combining different models are relatively absent from the current literature. The hybrid random forest illustrates that rather than focusing on unique techniques or new engineered features to enhance predictions, scholars should focus on ensembled models.

## **4.6 Practical implications**

We consider both the computing time and predictive performance to be valuable when considering the practical implications of the models. Football match prediction models are of great interest to decision-makers in the sport, football fans, news reporters and experts in the field.

Although single predictions from all models did not take longer than 45 seconds, when multiple matches or tournaments have to be predicted, the computation times can increase rapidly for specific models. The lowest computational time was obtained with the ELO-based models. The ELO-based models are built upon relatively simple calculations, which caused predictions of a season to take only a couple of seconds. Their standard results and swift computational times make them ideal for running quick simulations. The ELO parameters also proved to be of importance when used as a covariate in the random forest models from Groll et al. (2018, 2019) (see figure 3.2).

The computational times for the estimation of the ability parameters took around 20 seconds, depending on the amount of historical data incorporated. When making multiple predictions, these computational times can increase drastically and take up several hours. When choosing between the Poisson and ELO-based models, the ELO-based models are favoured, since they have equal performance with faster computations. However, the ability parameter incorporated in the hybrid random forest model proved to be the most important covariate. The advantage of the Poisson based models lies primarily in the creation of this highly informative ranking feature.

We do not recommend using the Weibull count model for predictions, since the computations were demanding and the performance was unsatisfactory. The attack and defence parameters, estimated as a result of the Weibull count model, could be further investigated in a random forest approach, to evaluate if they hold some value.

The random forest and boosting approaches had the best performances. The computational times were mostly dependent on the number of independent trees that had to be grown. The models of Baboota and Kaur (2019) only used 100 independent trees, while Groll et al. (2019) used 5000. Consequently, the individual predictions by the models of Baboota and Kaur (2019) took less than a second to complete, while for the models used by Groll et al. (2019) took around 30 seconds.

The plus-minus ratings were very promising, but due to a lack of data availability, the estimates were unreliable. It seems that for the prediction of football match outcomes, the player-based ratings can not yet compete with the team based ratings, and perhaps they are more useful when trying to evaluate the market value of players, as done in Sæbø and Hvattum (2015). Arntzen and Hvattum (2020) found that the combination of both the ELO ratings and plus-minus ratings significantly improved the predictive performance of the models, so it might be worth further investigating the plus-minus ratings in ensembled techniques.

To conclude, we give table 4.1, which ranks the models based on the performance and computational time. As mentioned above the more advanced machine learning

methods outperform the other models based on predictive power. The ELO-based models have an equal predictive performance as the Poisson-based models but obtain the results with much lower computational demands. The Weibull count and plus-minus rating model have the lowest predictive performance. From those two, we gave the lowest rank to the plus-minus ratings, since it requires much more computational time. Scholars and researchers in this field mainly try to develop new sophisticated features to rank the opposing teams or come up with advanced algorithms to predict football match results. We want to emphasise on the importance of combining various models and engineered features from the existing literature, which proves, in this master dissertation, to enhance the predictive power for most models.

Model	Performance	Computational time	Rank
Gradient boost	Very good	Medium	1
Random forest <sub>1</sub>	Very good	Medium	2
Hybrid random forest	Good	Medium	3
Random forest <sub>2</sub>	Good	Medium	4
Goal-based ELO	Average	Fast	5
Basic ELO	Average	Fast	6
Bivariate Poisson	Average	Medium	7
Independent Poisson	Average	Medium	8
Weibull count	Bad	Slow	9
Plus-minus	Bad	Very slow	10

Table 4.1: The final ranking of the models.  
Random forest<sub>1</sub> from Baboota and Kaur (2019) and Random forest<sub>2</sub> from Groll et al. (2018)



# **BIBLIOGRAPHY**

- Arntzen, H. and Hvattum, L. M. (2020). Predicting match outcomes in association football using team ratings and player ratings. *Statistical Modelling*, page 1471082X20929881.
- Baboota, R. and Kaur, H. (2019). Predictive analysis and modelling football results using machine learning approach for english premier league. *International Journal of Forecasting*, 35(2):741–755.
- Baker, R., Boshnakov, G., Kharrat, T., and McHale, I. (2016). Countr: an r package to generate flexible count models. *Journal of Statistical Software*.
- Bellman, R. (2003). Dynamic programming, 1957. *A very comprehensive reference with many economic examples is*.
- Boshnakov, G., Kharrat, T., and McHale, I. G. (2017). A bivariate weibull count model for forecasting association football scores. *International Journal of Forecasting*, 33(2):458–466.
- Bradley, R. A. and Terry, M. E. (1952). Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.
- Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3.
- Bröcker, J. and Smith, L. A. (2007). Scoring probabilistic forecasts: The importance of being proper. *Weather and Forecasting*, 22(2):382–388.
- Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794.
- Constantinou, A. C. and Fenton, N. E. (2012). Solving the problem of inadequate scoring rules for assessing probabilistic football forecast models. *Journal of Quantitative Analysis in Sports*, 8(1).

- Constantinou, A. C. and Fenton, N. E. (2013). Determining the level of ability of football teams by dynamic ratings based on the relative discrepancies in scores between adversaries. *Journal of Quantitative Analysis in Sports*, 9(1):37–50.
- Constantinou, A. C., Fenton, N. E., and Neil, M. (2012). pi-football: A bayesian network model for forecasting association football match outcomes. *Knowledge-Based Systems*, 36:322–339.
- Curley, J. P. (2016). engsoccerdata. English Soccer Data 1871-2016. R package version 0.1.5.
- Dixon, M. J. and Coles, S. G. (1997). Modelling association football scores and inefficiencies in the football betting market. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 46(2):265–280.
- Egidi, L. and Torelli, N. (2020). Comparing goal-based and result-based approaches in modelling football outcomes. *Social Indicators Research*, pages 1–13.
- Epstein, E. S. (1969). A scoring system for probability forecasts of ranked categories. *Journal of Applied Meteorology*, 8(6):985–987.
- Gneiting, T. and Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American statistical Association*, 102(477):359–378.
- Goddard, J. (2005). Regression models for forecasting goals and match results in association football. *International Journal of forecasting*, 21(2):331–340.
- Groll, A., Ley, C., Schauburger, G., and Van Eetvelde, H. (2018). Prediction of the fifa world cup 2018-a random forest approach with an emphasis on estimated team ability parameters. *arXiv preprint arXiv:1806.03208*.
- Groll, A., Ley, C., Schauburger, G., and Van Eetvelde, H. (2019). A hybrid random forest to predict soccer matches in international tournaments. *Journal of Quantitative Analysis in Sports*, 15(4):271–287.
- Hvattum, L. M. (2017). Ordinal versus nominal regression models and the problem of correctly predicting draws in soccer. *International Journal of Computer Science in Sport*, 16(1):50–64.
- Hvattum, L. M. (2019). A comprehensive review of plus-minus ratings for evaluating individual players in team sports. *International Journal of Computer Science in Sport*, 18(1):1–23.
- Hvattum, L. M. and Arntzen, H. (2010). Using elo ratings for match result prediction in association football. *International Journal of forecasting*, 26(3):460–470.

- Joseph, A., Fenton, N. E., and Neil, M. (2006). Predicting football results using bayesian nets and other machine learning techniques. *Knowledge-Based Systems*, 19(7):544–553.
- Karlis, D. and Ntzoufras, I. (2003). Analysis of sports data by using bivariate poisson models. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 52(3):381–393.
- Katti, S. and Rao, A. V. (1968). Handbook of the poisson distribution.
- Keogh, F. and Rose, G. (2013). Football betting-the global gambling industry worth billions. *BBC*, 3rd October, available at: <http://www.bbc.com/sport/football/24354124> (accessed 6th July, 2016).
- Koning, R. H. (2000). Balance in competition in dutch soccer. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 49(3):419–431.
- Kundu, T., Choudhury, A. R., and Rai, S. (2019). Predicting english premier league matches using classification and regression.
- Ley, C., Wiele, T. V. d., and Eetvelde, H. V. (2019). Ranking soccer teams on the basis of their current strength: A comparison of maximum likelihood approaches. *Statistical Modelling*, 19(1):55–73.
- Macdonald, B. (2012). Adjusted plus-minus for nhl players using ridge regression with goals, shots, fenwick, and corsi. *Journal of Quantitative Analysis in Sports*, 8(3).
- Maher, M. J. (1982). Modelling association football scores. *Statistica Neerlandica*, 36(3):109–118.
- McHale, I. and Scarf, P. (2011). Modelling the dependence of goals scored by opposing teams in international soccer matches. *Statistical Modelling*, 11(3):219–236.
- Milton, J. (2017). History of sports betting.
- Mosteller, F. (2006). Remarks on the method of paired comparisons: I. the least squares solution assuming equal standard deviations and equal correlations. In *Selected Papers of Frederick Mosteller*, pages 157–162. Springer.
- Munting, R. (1996). *An economic and social history of gambling in Britain and the USA*. Manchester University Press.
- Pantuso, G. and Hvattum, L. M. (2019). Maximizing performance with an eye on the finances: a chance-constrained model for football transfer market decisions. *arXiv preprint arXiv:1911.04689*.
- Probst, P. and Boulesteix, A.-L. (2017). To tune or not to tune the number of trees in random forest. *The Journal of Machine Learning Research*, 18(1):6673–6690.

- RightCasino (2014). History of online casinos: infographic.
- Sæbø, O. D. and Hvattum, L. M. (2015). Evaluating the efficiency of the association football transfer market using regression based player ratings. In *Norsk IKT-konferanse for forskning og utdanning*.
- Skellam, J. G. (1946). The frequency distribution of the difference between two poisson variates belonging to different populations. *Journal of the Royal Statistical Society. Series A (General)*, 109(Pt 3):296–296.
- Stefani, R. T. (1977). Football and basketball predictions using least squares. *IEEE Transactions on systems, man, and cybernetics*, 7(2):117–21.
- Suzuki, A. K., Salasar, L. E., Leite, J., and Louzada-Neto, F. (2010). A bayesian approach for predicting match outcomes: the 2006 (association) football world cup. *Journal of the Operational Research Society*, 61(10):1530–1539.
- Thurstone, L. L. (1927). Psychophysical analysis. *The American journal of psychology*, 38(3):368–389.
- Vehtari, A., Gelman, A., and Gabry, J. (2017). Practical bayesian model evaluation using leave-one-out cross-validation and waic. *Statistics and computing*, 27(5):1413–1432.
- Wheatcroft, E. (2019). Evaluating probabilistic forecasts of football matches: The case against the ranked probability score. *arXiv preprint arXiv:1908.08980*.

## APPENDIX A

# APPENDIX

[Appendix]

### A.1 Code for the scoring rules

```
1 BrierScore <- function(score){
2
3   pwin <- score[1]
4   pdraw <- score[2]
5   ploss <- score[3]
6
7   owin <- score[4]
8   odraw <- score[5]
9   oloss <- score[6]
10
11   brier<- sum((pwin-owin)**2,
12              (pdraw-odraw)**2,
13              (ploss-oloss)**2)
14   return(brier)
15 }
16
17
18 RPS2 <- function(score){
19
20   pwin <- score[1]
21   pdraw <- score[2]
22   ploss <- score[3]
23
24   owin <- score[4]
25   odraw <- score[5]
26   oloss <- score[6]
27
28   rps2<- sum((pwin-owin)**2,
29             (ploss-oloss)**2)
30   return(rps2)
31 }
32
33 IGN <- function(score){
```

```
34
35 pwin <- score[1]
36 pdraw <- score[2]
37 ploss <- score[3]
38
39 owin <- score[4]
40 odraw <- score[5]
41 oloss <- score[6]
42
43 if (owin == 1) {
44   ignorance <- -log(base = 2,x = pwin)
45 } else
46   if (odraw == 1) {
47     ignorance <- -log(base = 2,x = pdraw)
48   } else ignorance <- -log(base = 2,x = ploss)
49
50 return(ignorance)
51 }
52
53 IGN2 <- function(score){
54
55   pwin <- score[1]
56   pdraw <- score[2]
57   ploss <- score[3]
58
59   owin <- score[4]
60   odraw <- score[5]
61   oloss <- score[6]
62
63   if (owin == 1) {
64     ignorance <- -log(base = 2,x = pwin)
65   } else
66     if (odraw == 1) {
67       ignorance <- -log(x = pdraw)
68     } else ignorance <- -log(base = 2,x = ploss)
69
70   return(ignorance)
71 }
72
73
74
75
76
77
78 RPS <- function(score, outcomes = 3){
79
80   pwin <- score[1]
81   pdraw <- score[2]
82   ploss <- score[3]
```

```
83
84 p <- c(pwin,pdraw,ploss)
85 p <- unlist(p)
86 p <- unname(p)
87
88 owin <- score[4]
89 odraw <- score[5]
90 oloss <- score[6]
91
92 o <- c(owin,odraw,oloss)
93 o <- unlist(o)
94 o <- unname(o)
95
96 cumulative <- c()
97 for (i in 1:outcomes-1){
98   for (j in 1:i){
99     cumulative <- c(cumulative , (p[j]-o[j])^2)
100   }
101
102 }
103
104 return(sum(cumulative))
105 }
```

## A.2 Code for the ELO based models

### A.2.1 basic ELO

```
1 #####
2 # Using ELO ratings for match result prediction in association football
3 # M. hvattum
4 #####
5
6 # packages
7 {
8   setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/'
9         )
9   library(engsoccerdata)
10  library(tidyverse)
11  library(skellam)
12  library(caret)
13  library(MASS)
14  source('BrierRpsIgnorance.R')
15 }
16
17 #data
18 {
```

```
19 # Premier League
20 PL_data <- england
21
22 PL_data%>% filter(tier == 1)
23
24
25 data <- PL_data %>%
26   filter(tier == 1,
27          Season %in% 2005:2017)
28
29 }
30
31 # Save output
32 Prob_home <- c()
33 Prob_draw <- c()
34 Prob_away <- c()
35 observed_home <- c()
36 observed_draw <- c()
37 observed_away <- c()
38
39 # algorithm
40 for (season_id in 2008:2015){
41
42   data %>% group_by(Season) %>%
43     summarise(matches = length(Date),
44               team = length(levels(factor(home))))
45
46
47   #make round parameter
48   id <- order(data$Season, data$Date)
49   data <- data[id,]
50   data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
51
52
53
54   # 1. initial elo ranking
55   {
56     Init_data <- data %>% filter(Season == season_id-3)
57
58
59     teams <- sort(levels(factor(Init_data$home)))
60     elo <- rep(0, length(teams))
61
62     elo_hash <- data.frame(team = teams, elo = elo)
63
64     for (i in 1:nrow(Init_data)){
65
66       # 1 . Hyperparameters
67       c <- 10
```

```
68     d <- 400
69     k <- 20
70
71     # 2. observation
72     obs <- Init_data[i,]
73
74     elo_home <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)]
75     elo_away <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)]
76     x <- elo_home - elo_away
77
78     E_score_home <- (1+c^(-x/d))^-1
79     E_score_away <- 1 - E_score_home
80
81     actual_score_home <- ifelse(obs$result == 'H' , 1 ,
82                                ifelse(obs$result == 'A', 0, 0.5 ))
83     actual_score_away <- 1 - actual_score_home
84
85     # 3. Update Elo
86     elo_home_update <- elo_home + k*(actual_score_home-E_score_home)
87     elo_away_update <- elo_away + k*(actual_score_away-E_score_away)
88
89     elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)] <- elo_home_
90         update
91     elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)] <- elo_away_
92         update
93
94 }
95
96
97 # 2. Training data set
98 {
99     train_data <- data %>%
100         filter(Season %in% (season_id-2):(season_id-1))
101
102     train_data <- rbind( train_data, data %>% filter(Season == season_id,round <= 5)
103         )
104
105     print('training dataset')
106     print(train_data %>% group_by(Season) %>% summarise(rounds = length(levels(
107         factor(round)))))
108
109
110     teams <- levels(factor(train_data$home))[!levels(factor(train_data$home)) %in%
111         elo_hash$team]
112
113     elo_hash <- rbind(elo_hash,
114         data.frame(team = teams,
```

```

112         elo = rep(0,length(teams)))
113
114
115     x_vect <- c()
116     y_vect <- c()
117
118     for (i in 1:nrow(train_data)){
119
120         # 1. Hyperparameters
121         c <- 10
122         d <- 400
123         k <- 20
124
125         # 2. Draw observation and save information
126         obs <- train_data[i,]
127
128         y_vect <- c(y_vect,as.character(obs$result))
129
130         elo_home <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)]
131         elo_away <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)]
132
133         x <- elo_home - elo_away
134         x_vect <- c(x_vect,x)
135
136         E_score_home <- (1+c^(-x/d))^-1
137         E_score_away <- 1 - E_score_home
138
139         actual_score_home <- ifelse(obs$result == 'H' , 1 ,
140                                     ifelse(obs$result == 'A', 0, 0.5 ))
141         actual_score_away <- 1 - actual_score_home
142
143         # 3. Update Elo
144         elo_home_update <- elo_home + k*(actual_score_home-E_score_home)
145         elo_away_update <- elo_away + k*(actual_score_away-E_score_away)
146
147         elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)] <- elo_home_
            update
148         elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)] <- elo_away_
            update
149
150     }
151
152     # Use information for OLR
153     OLR <- polr(formula = 'y~x',
154                 data = data.frame(y = factor(y_vect, levels = c('A', 'D', 'H'),ordered
155                                     = TRUE),
156                                     x = x_vect))
157     summary(OLR)

```

```
158   }
159
160
161   # 3. Test dataset
162   {
163     test_data <- data %>% filter(Season == season_id, round > 5)
164     print(test_data %>% group_by(Season) %>% summarise(rounds = length(levels(factor
165       (round)))))
166
167     for (i in 1:nrow(test_data)){
168
169       # 1. Hyper parameters
170       c <- 10
171       d <- 400
172       k <- 20
173
174       # 2. Match to predict
175       obs <- test_data[i,]
176
177       y_vect <- c(y_vect, as.character(obs$result))
178
179       elo_home <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)]
180       elo_away <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)]
181
182       x <- elo_home - elo_away
183       x_vect <- c(x_vect, x)
184
185       # 3. Predict
186       prob <- predict(OLR,
187         newdata = data.frame(y = factor(obs$result, levels = c('A', 'D'
188           , 'H'), ordered = TRUE),
189           x = x),
190         type = 'p')
191
192       Prob_home <- c(Prob_home, unname(prob['H']))
193       Prob_draw <- c(Prob_draw, unname(prob['D']))
194       Prob_away <- c(Prob_away, unname(prob['A']))
195
196       # 4. observed
197       if (as.character(obs$result) == 'H'){
198         observed_home <- c(observed_home, 1)
199         observed_draw <- c(observed_draw, 0)
200         observed_away <- c(observed_away, 0)
201       }
202       else {
203         if (as.character(obs$result) == 'A'){
204           observed_home <- c(observed_home, 0)
205           observed_draw <- c(observed_draw, 0)
```

```

205     observed_away <- c(observed_away,1)
206
207   }
208   else {
209     observed_home <- c(observed_home,0)
210     observed_draw <- c(observed_draw,1)
211     observed_away <- c(observed_away,0)
212
213   }
214 }
215
216 # 5. Update
217 E_score_home <- (1+c^(-x/d))^-1
218 E_score_away <- 1 - E_score_home
219
220 actual_score_home <- ifelse(obs$result == 'H' , 1 ,
221                             ifelse(obs$result == 'A', 0, 0.5 ))
222 actual_score_away <- 1 - actual_score_home
223
224 elo_home_update <- elo_home + k*(actual_score_home-E_score_home)
225 elo_away_update <- elo_away + k*(actual_score_away-E_score_away)
226
227 elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)] <- elo_home_
    update
228 elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)] <- elo_away_
    update
229
230 OLR <- polr(formula = 'y~x',
231             data = data.frame(y = factor(y_vect, levels = c('A','D','H'),
232                                     ordered = TRUE),
233                               x = x_vect))
234
235 }
236 }
237
238 out <- data.frame(pwin = Prob_home,
239                  pdraw = Prob_draw,
240                  ploss = Prob_away,
241                  owin = observed_home,
242                  odraw = observed_draw,
243                  oloss = observed_away)
244
245 out$ignorance <- apply(out[,1:6], 1, IGN)
246 out$brierScore <- apply(out[,1:6], 1, BrierScore)
247 out$RPS <- apply(out[,1:6], 1, RPS2)
248
249
250 mean(out$ignorance)

```

```
251 sd(out$ignorance)
252
253 mean(out$brierScore)
254 sd(out$brierScore)
255
256 mean(out$RPS/2)
257 sd(out$RPS/2)
258
259 saveRDS(object = out, file = 'ELOb_Hvattum2010.rds')
```

### A.2.2 Goal-based ELO

```
1 #####
2 # Using ELO ratings for match result prediction in association football
3 # M. hvattum
4 #####
5
6 # packages
7 setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/')
8 library(engsoccerdata)
9 library(tidyverse)
10 library(skeLLAM)
11 library(caret)
12 library(MASS)
13 source('BrierRpsIgnorance.R')
14
15
16 # Premier League
17 PL_data <- england
18
19 PL_data %>% filter(tier == 1)
20
21
22 data <- PL_data %>%
23   filter(tier == 1,
24          Season %in% 2005:2017)
25
26
27
28 # Save output
29 Prob_home <- c()
30 Prob_draw <- c()
31 Prob_away <- c()
32 observed_home <- c()
33 observed_draw <- c()
34 observed_away <- c()
35 week <- c()
36
37 # algorithm
```

```

38 for (season_id in 2008:2015){
39
40   data %>% group_by(Season) %>%
41     summarise(matches = length(Date),
42               team = length(levels(factor(home))))
43
44
45   #make round parameter
46   id <- order(data$Season, data$Date)
47   data <- data[id,]
48   data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
49
50
51
52   # 1. initial elo ranking
53   {
54     Init_data <- data %>% filter(Season == season_id-3)
55
56
57     teams <- sort(levels(factor(Init_data$home)))
58     elo <- rep(0, length(teams))
59
60     elo_hash <- data.frame(team = teams, elo = elo)
61
62     for (i in 1:nrow(Init_data)){
63
64       # 1. observation
65       obs <- Init_data[i,]
66
67       # 2 . Hyperparameters
68       c <- 10
69       d <- 400
70       k0 <- 10
71       lambda <- 1
72       gamma <- abs(obs$hgoal-obs$vgoal)
73       k <- k0*((1+gamma)**lambda)
74
75       # 3. variables
76       elo_home <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)]
77       elo_away <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)]
78       x <- elo_home - elo_away
79
80       E_score_home <- (1+c^(-x/d))^-1
81       E_score_away <- 1 - E_score_home
82
83       actual_score_home <- ifelse(obs$result == 'H' , 1 ,
84                                   ifelse(obs$result == 'A', 0, 0.5 ))
85       actual_score_away <- 1 - actual_score_home
86

```

```
87     # 4. Update Elo
88     elo_home_update <- elo_home + k*(actual_score_home-E_score_home)
89     elo_away_update <- elo_away + k*(actual_score_away-E_score_away)
90
91     elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)] <- elo_home_
      update
92     elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)] <- elo_away_
      update
93
94   }
95
96 }
97
98
99 # 2. Training data set
100 {
101   train_data <- data %>%
102     filter(Season %in% (season_id-2):(season_id-1))
103
104   train_data <- rbind( train_data, data %>% filter(Season == season_id,round <= 5)
      )
105
106   print('training dataset')
107   print(train_data %>% group_by(Season) %>% summarise(rounds = length(levels(
      factor(round)))))
108
109
110   teams <- levels(factor(train_data$home))[!levels(factor(train_data$home)) %in%
      elo_hash$team]
111
112   elo_hash <- rbind(elo_hash,
113     data.frame(team = teams,
114       elo = rep(0,length(teams))))
115
116
117   x_vect <- c()
118   y_vect <- c()
119
120   for (i in 1:nrow(train_data)){
121
122     # 1. Draw observation and save information
123     obs <- train_data[i,]
124
125     y_vect <- c(y_vect,as.character(obs$result))
126
127     elo_home <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)]
128     elo_away <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)]
129
130     x <- elo_home - elo_away
```

```

131     x_vect <- c(x_vect,x)
132
133     E_score_home <- (1+c^(-x/d))^-1
134     E_score_away <- 1 - E_score_home
135
136     actual_score_home <- ifelse(obs$result == 'H' , 1 ,
137                                ifelse(obs$result == 'A', 0, 0.5 ))
138     actual_score_away <- 1 - actual_score_home
139
140
141     # 2. Hyperparameters
142     c <- 10
143     d <- 400
144     k0 <- 10
145     lambda <- 1
146     gamma <- abs(obs$hgoal-obs$vgoal)
147     k <- k0*((1+gamma)**lambda)
148
149
150     # 3. Update Elo
151     elo_home_update <- elo_home + k*(actual_score_home-E_score_home)
152     elo_away_update <- elo_away + k*(actual_score_away-E_score_away)
153
154     elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)] <- elo_home_
155       update
156     elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)] <- elo_away_
157       update
158
159     # Use information for OLR
160     OLR <- polr(formula = 'y~x',
161                 data = data.frame(y = factor(y_vect, levels = c('A', 'D', 'H')),ordered
162                                   = TRUE),
163                                   x = x_vect))
164
165     summary(OLR)
166   }
167
168   # 3. Test dataset
169   {
170     test_data <- data %>% filter(Season == season_id,round > 5)
171     print(test_data %>% group_by(Season) %>% summarise(rounds = length(levels(factor
172       (round)))))
173
174     for (i in 1:nrow(test_data)){
175       week <- c(week,test_data$round)

```

```
176
177   # 1. Match to predict
178   obs <- test_data[i,]
179
180   y_vect <- c(y_vect,as.character(obs$result))
181
182   elo_home <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)]
183   elo_away <- elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)]
184
185   x <- elo_home - elo_away
186   x_vect <- c(x_vect,x)
187
188   E_score_home <- (1+c^(-x/d))^(-1)
189   E_score_away <- 1 - E_score_home
190
191   actual_score_home <- ifelse(obs$result == 'H' , 1 ,
192                               ifelse(obs$result == 'A', 0, 0.5 ))
193   actual_score_away <- 1 - actual_score_home
194
195
196
197   # 2. Hyper parameters
198   c <- 10
199   d <- 400
200   k0 <- 10
201   lambda <- 1
202   gamma <- abs(obs$hgoal-obs$vgoal)
203   k <- k0*((1+gamma)**lambda)
204
205
206   # 3. Predict
207   prob <- predict(OLR,
208                   newdata = data.frame(y = factor(obs$result, levels = c('A','D'
209                                     , 'H')),ordered = TRUE),
209                                     x = x),
210                   type = 'p')
211
212   Prob_home <- c(Prob_home,unname(prob['H']))
213   Prob_draw <- c(Prob_draw,unname(prob['D']))
214   Prob_away <- c(Prob_away,unname(prob['A']))
215
216   # 4. observed
217   if (as.character(obs$result) == 'H'){
218     observed_home <- c(observed_home,1)
219     observed_draw <- c(observed_draw,0)
220     observed_away <- c(observed_away,0)
221
222   }
223   else {
```

```

224     if (as.character(obs$result) == 'A'){
225         observed_home <- c(observed_home,0)
226         observed_draw <- c(observed_draw,0)
227         observed_away <- c(observed_away,1)
228
229     }
230     else {
231         observed_home <- c(observed_home,0)
232         observed_draw <- c(observed_draw,1)
233         observed_away <- c(observed_away,0)
234
235     }
236 }
237
238 # 5. Update
239
240 elo_home_update <- elo_home + k*(actual_score_home-E_score_home)
241 elo_away_update <- elo_away + k*(actual_score_away-E_score_away)
242
243 elo_hash$elo[levels(elo_hash$team) == as.character(obs$home)] <- elo_home_
    update
244 elo_hash$elo[levels(elo_hash$team) == as.character(obs$visitor)] <- elo_away_
    update
245
246 OLR <- polr(formula = 'y~x',
247             data = data.frame(y = factor(y_vect, levels = c('A', 'D', 'H'),
248                                     ordered = TRUE),
249                               x = x_vect))
250
251 }
252 }
253
254 out <- data.frame(pwin = Prob_home,
255                  pdraw = Prob_draw,
256                  ploss = Prob_away,
257                  owin = observed_home,
258                  odraw = observed_draw,
259                  oloss = observed_away)
260
261 out$ignorance <- apply(out[,1:6], 1, IGN)
262 out$brierScore <- apply(out[,1:6], 1, BrierScore)
263 out$RPS <- apply(out[,1:6], 1, RPS2)
264
265
266 mean(out$ignorance)
267 sd(out$ignorance)
268
269 mean(out$brierScore)

```

```
270 sd(out$brierScore)
271
272 mean(out$RPS/2)
273 sd(out$RPS/2)
274
275 saveRDS(object = out, file = 'EL0g_Hvattum2010.rds')
276
277
278
279 #####
280 #plot
281
282 out <- data.frame(pwin = Prob_home,
283                   pdraw = Prob_draw,
284                   ploss = Prob_away,
285                   owin = observed_home,
286                   odraw = observed_draw,
287                   oloss = observed_away,
288                   week = week)
289
290 out$IgnoranceScore <- apply(out[,1:6], 1, IGN)
291 out$BrierScore <- apply(out[,1:6], 1, BrierScore)
292 out$RankedProbScore <- apply(out[,1:6], 1, RPS2)
293
294
295 out$Model <- 'out'
296 out <- out[,c('IgnoranceScore', 'BrierScore', 'RankedProbScore', 'week')] %>% gather(
297   key = 'ScoringRule', value = 'Score', IgnoranceScore:RankedProbScore )
298
299 p <- out %>%
300   ggplot(aes(x= week, y=Score, color = ScoringRule, fill = ScoringRule)) +
301   #geom_point(alpha = 0.25) +
302   geom_smooth(method='glm', se = F)+
303   theme_bw() +
304   labs(x = 'Week',
305        y = 'Score',
306        title = 'Evolution of the scoring rules over season 2012',
307        color = 'Scoring rule',
308        fill = 'Scoring rule') +
309   theme(plot.title = element_text(hjust = 0.5))
310
311 p
312 ggsave(plot = p, 'scorerulee.png', height = 4, width = 5.2)
```

## A.3 Code for the plus-minus models

### A.3.1 Plus-minus

```

1 #####
2 # Player Ratings
3 # M. hvattum
4 #####
5
6 # packages
7 {
8   setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/'
9         )
10  library(engsoccerdata)
11  library(tidyverse)
12  library(skeLLam)
13  library(caret)
14  library(MASS)
15  source('BrierRpsIgnorance.R')
16 }
17 # data
18 {
19
20   load('eplhvattum.Rdata')
21   data <- data1
22   str(data)
23   data <- data[data$date <= as.Date('2016-05-16'),]
24
25
26   # current times
27   # Premier League
28   PL_data <- england
29
30   PL_data <- PL_data %>% filter(tier == 1,
31                                Season %in% 2000:2017)
32
33   id <- order(PL_data$Season, PL_data$date)
34   PL_data <- PL_data[id,]
35   PL_data$round <- rep(sort(rep(1:38,10)), max(PL_data$Season) - min(PL_data$Season)
36                        +1 )
37   DataSeasonRound <- unique(PL_data[,c('date', 'Season', 'round')])
38   colnames(DataSeasonRound) <- c('date', 'season', 'round')
39   DataSeasonRound$date <- as.Date(DataSeasonRound$date)
40
41   data$Season <- DataSeasonRound$season[match(data$date, DataSeasonRound$date)]
42   data$round <- DataSeasonRound$round[match(data$date, DataSeasonRound$date)]
43 }
44
45 # functions
46 {
47   SplitSegments <- function(obs){

```

```
48
49   # information
50   {
51     # player changes
52     {
53       playerChanges <- data.frame(time = c(obs$home_change_minute1, obs$home_change_
54         _minute2, obs$home_change_minute3,
55         obs$away_change_minute1, obs$away_change_
56         _minute2, obs$away_change_minute3),
57       athleteIn = c(obs$home_change_in1, obs$home_
58         change_in2, obs$home_change_in3,
59         obs$away_change_in1, obs$away_
60         change_in2, obs$away_change_
61         minute3),
62       athleteOut = c(obs$home_change_out1, obs$home_
63         change_out2, obs$home_change_out3,
64         obs$away_change_out1, obs$away_
65         change_out2, obs$away_change_
66         out3))
67
68       playerChanges <- playerChanges[order(playerChanges$time),]
69       playerChanges <- na.omit(playerChanges)
70     }
71
72     # goal changes
73     {
74       Goalchanges <- data.frame(goals = c(obs$home_goal_minute1, obs$home_goal_
75         minute2, obs$home_goal_minute3,
76         obs$home_goal_minute4, obs$home_goal_
77         minute5, obs$home_goal_minute6,
78         obs$home_goal_minute7, obs$home_goal_
79         minute8, obs$home_goal_minute9,
80         obs$away_goal_minute1, obs$away_goal_
81         minute2),
82       team = c(rep('home', 9), rep('away', 2)))
83       Goalchanges <- Goalchanges[order(Goalchanges$goals),]
84       Goalchanges <- na.omit(Goalchanges)
85     }
86
87     # params
88     {
89       segments <- list()
90       start <- 0
91       time <- 0
92     }
93
94     #players
```

```

85   {
86     players <- unique(c(levels(obs$home_team_player1), levels(obs$home_team_
      player2), levels(obs$home_team_player3), levels(obs$home_team_player4),
87       levels(obs$home_team_player5), levels(obs$home_team_
      player6), levels(obs$home_team_player7), levels(obs$
      home_team_player8),
88       levels(obs$home_team_player9), levels(obs$home_team_
      player10), levels(obs$home_team_player11)))
89
90     matrix_players <- t(matrix(c(players, rep(0, length(players))), ncol = 2))
91
92     home_players <- c(as.character(obs$home_team_player1), as.character(obs$home_
      team_player2), as.character(obs$home_team_player3),
93       as.character(obs$home_team_player4), as.character(obs$home_
      team_player5), as.character(obs$home_team_player6),
94       as.character(obs$home_team_player7), as.character(obs$home_
      team_player8), as.character(obs$home_team_player9),
95       as.character(obs$home_team_player10), as.character(obs$home_
      _team_player11))
96
97     away_players <- c(as.character(obs$away_team_player1), as.character(obs$away_
      team_player2), as.character(obs$away_team_player3),
98       as.character(obs$away_team_player4), as.character(obs$away_
      team_player5), as.character(obs$away_team_player6),
99       as.character(obs$away_team_player7), as.character(obs$away_
      team_player8), as.character(obs$away_team_player9),
100      as.character(obs$away_team_player10), as.character(obs$away_
      _team_player11))
101
102     matrix_players[2, match(home_players, matrix_players[1,])] <- 1
103     matrix_players[2, match(away_players, matrix_players[1,])] <- -1
104   }
105
106 }
107
108
109 if (nrow(playerChanges) > 0){
110   for (i in 1:(nrow(playerChanges))){
111
112     #segment duration
113     end <- playerChanges$time[i]
114     duration <- end - start
115     start <- end
116     if (duration == 0) duration <- 1
117
118     #goals
119     starttime <- time
120     time <- time + duration
121     endtime <- time

```

```
122     goals_start <- sum(ifelse(Goalchanges$team[Goalchanges$goals<starttime] == '
      home' , 1,-1))
123     goals_end <- sum(ifelse(Goalchanges$team[Goalchanges$goals<endtime] == 'home
      ' , 1,-1))
124     gd <- goals_end - goals_start
125
126     #beta
127     beta <- 90*gd/duration
128
129     #players
130     a <- as.numeric(matrix_players[2,])
131
132
133     # red cards
134
135
136     # age
137
138
139     #home
140     homeeffect <- 1
141
142     # add to segment list
143     segment <- data.frame(matrix(c(beta,a,homeeffect,as.Date(obs$date),obs$
      Season,obs$round),nrow=1))
144     colnames(segment) <- c('beta',matrix_players[1,],'homeadv','date','Season',
      'round')
145     segments[[paste0('segment',i)]] <- segment
146
147     # update players
148     if (as.character(playerChanges$athleteOut[i]) %in% away_players){
149       away_players[match(as.character(playerChanges$athleteOut[i]),away_players)
        ] <- as.character(playerChanges$athleteIn[i])
150     }
151     if (as.character(playerChanges$athleteOut[i]) %in% home_players){
152       home_players[match(as.character(playerChanges$athleteOut[i]),home_players)
        ] <- as.character(playerChanges$athleteIn[i])
153     }
154
155     matrix_players <- t(matrix(c(players, rep(0,length(players))),ncol = 2))
156     matrix_players[2,match(home_players,matrix_players[1,])] <- 1
157     matrix_players[2,match(away_players,matrix_players[1,])] <- -1
158
159   }
160 }
161 else { i <-1 }
162
163 # last segment
164 {
```

```

165     #segment duration
166     end <- 90
167     duration <- end - start
168     start <- end
169     if (duration == 0) duration <- 1
170     #goals
171     starttime <- time
172     time <- time + duration
173     endtime <- time
174     goals_start <- sum(ifelse(Goalchanges$team[Goalchanges$goals<starttime] == '
        home' , 1,-1))
175     goals_end <- sum(ifelse(Goalchanges$team[Goalchanges$goals<endtime] == 'home'
        , 1,-1))
176     gd <- goals_end - goals_start
177
178     #beta
179     beta <- 90*gd/duration
180
181     #players
182     a <- as.numeric(matrix_players[2,])
183
184     # red cards
185
186
187     # age
188
189
190     #home
191     homeeffect <- 1
192
193     segment <- data.frame(matrix(c(beta,a,homeeffect,obs$date,obs$Season,obs$round
        ),nrow=1))
194     colnames(segment) <- c('beta',matrix_players[1,],'homeadv','date','Season','
        round')
195     segments[[paste0('segment',i+1)]] <- segment
196
197 }
198
199     return(segments)
200 }
201 }
202
203 #split data in segments
204 {
205     df <- data
206     df$home_team_player1 <- factor(df$home_team_player1)
207     df$home_team_player2 <- factor(df$home_team_player2)
208     df$home_team_player3 <- factor(df$home_team_player3)
209     df$home_team_player4 <- factor(df$home_team_player4)

```

```
210 df$home_team_player5 <- factor(df$home_team_player5)
211 df$home_team_player6 <- factor(df$home_team_player6)
212 df$home_team_player7 <- factor(df$home_team_player7)
213 df$home_team_player8 <- factor(df$home_team_player8)
214 df$home_team_player9 <- factor(df$home_team_player9)
215 df$home_team_player10 <- factor(df$home_team_player10)
216 df$home_team_player11 <- factor(df$home_team_player11)
217
218 df$away_team_player1 <- factor(df$away_team_player1)
219 df$away_team_player2 <- factor(df$away_team_player2)
220 df$away_team_player3 <- factor(df$away_team_player3)
221 df$away_team_player4 <- factor(df$away_team_player4)
222 df$away_team_player5 <- factor(df$away_team_player5)
223 df$away_team_player6 <- factor(df$away_team_player6)
224 df$away_team_player7 <- factor(df$away_team_player7)
225 df$away_team_player8 <- factor(df$away_team_player8)
226 df$away_team_player9 <- factor(df$away_team_player9)
227 df$away_team_player10 <- factor(df$away_team_player10)
228 df$away_team_player11 <- factor(df$away_team_player11)
229 #str(df)
230
231 seglist <- list()
232 for (id in 1:nrow(df)){
233   obs <- df[id,]
234   segments <- SplitSegments(obs)
235   segdf <- do.call(rbind, segments)
236   seglist[[as.character(id)]] <- segdf
237 }
238
239 df <- do.call(rbind, seglist)
240 }
241
242
243 df_2015 <- df %>% filter(Season %in% 2008:2015)
244 data2015 <- data %>% filter(Season %in% 2008:2015)
245
246 # Save output
247 Prob_home <- c()
248 Prob_draw <- c()
249 Prob_away <- c()
250 observed <- c()
251
252 # algorithm
253 yrsago <- 4
254 for (season_id in 2008:2015){
255
256   # train data
257   train_data <- data %>% filter(Season %in% (season_id-2):(season_id-1))
258
```

```

259 train_data <- rbind( train_data, data %>% filter(Season == season_id,round <= 5))
260 train_data$seasonround <- paste0(train_data$Season,train_data$round)
261
262 print('training dataset')
263 print(train_data %>% group_by(Season) %>% summarise(rounds = length(levels(factor(
    round)))))
264
265 ## get results and pm
266 train_results <- c()
267 train_pm <- c()
268 for (i in 1:length(unique(train_data$seasonround))){
269
270   #train obs
271   train_obs <- train_data[train_data$seasonround==unique(train_data$seasonround)[
    i],]
272
273   print(paste('train',train_obs$Season[1],i))
274   currenttime <- max(train_obs$date)
275
276   # df of all segments before current time
277   {
278     df_seg_train <- df %>% filter((date <= currenttime)&((date >= currenttime -
    yrsago*365)))
279     k <- 0.2
280     time_depre <- exp(-k*(as.numeric(currenttime) - df_seg_train$date)/365)
281   }
282
283   # get coefficients
284   {
285     beta <- df_seg_train$beta*time_depre
286
287     a <- as.matrix(df_seg_train[,-c(match('beta',colnames(df_seg_train)),
288                                     match('Season',colnames(df_seg_train)),
289                                     match('date',colnames(df_seg_train)),
290                                     match('round',colnames(df_seg_train))]))
291
292     a <- a * time_depre
293
294     lambda <- 3000
295     x <- solve(t(a) %*% a + diag(lambda,ncol(a))) %*% t(a) %*% beta
296   }
297
298
299   # for every obs of in currenttime get result and difference in pm of teams
300   {
301     gd <- train_obs$home_goals - train_obs$away_goals
302     result <- ifelse(gd < 0, 'A', ifelse(gd == 0, 'D', 'H'))
303     train_results <- c(train_results,result)
304

```

```
305     x_home_player1 <- x[match(as.character(train_obs$home_team_player1), rownames(x
    ))]
306     x_home_player2 <- x[match(as.character(train_obs$home_team_player2), rownames(x
    ))]
307     x_home_player3 <- x[match(as.character(train_obs$home_team_player3), rownames(x
    ))]
308     x_home_player4 <- x[match(as.character(train_obs$home_team_player4), rownames(x
    ))]
309     x_home_player5 <- x[match(as.character(train_obs$home_team_player5), rownames(x
    ))]
310     x_home_player6 <- x[match(as.character(train_obs$home_team_player6), rownames(x
    ))]
311     x_home_player7 <- x[match(as.character(train_obs$home_team_player7), rownames(x
    ))]
312     x_home_player8 <- x[match(as.character(train_obs$home_team_player8), rownames(x
    ))]
313     x_home_player9 <- x[match(as.character(train_obs$home_team_player9), rownames(x
    ))]
314     x_home_player10 <- x[match(as.character(train_obs$home_team_player10), rownames
    (x))]
315     x_home_player11 <- x[match(as.character(train_obs$home_team_player11), rownames
    (x))]
316     x_home_player12 <- x[match(rep('homeadv', nrow(train_obs)), rownames(x))]
317     x_home <- apply(matrix(c(x_home_player1, x_home_player2, x_home_player3,
318                             x_home_player4, x_home_player5, x_home_player6,
319                             x_home_player7, x_home_player8, x_home_player9,
320                             x_home_player10, x_home_player11, x_home_player12),
321                             ncol = 12),
322                        1, mean, na.rm = TRUE)
323     x_home[is.na(x_home)] <- 0
324
325     x_away_player1 <- x[match(as.character(train_obs$away_team_player1), rownames(x
    ))]
326     x_away_player2 <- x[match(as.character(train_obs$away_team_player2), rownames(x
    ))]
327     x_away_player3 <- x[match(as.character(train_obs$away_team_player3), rownames(x
    ))]
328     x_away_player4 <- x[match(as.character(train_obs$away_team_player4), rownames(x
    ))]
329     x_away_player5 <- x[match(as.character(train_obs$away_team_player5), rownames(x
    ))]
330     x_away_player6 <- x[match(as.character(train_obs$away_team_player6), rownames(x
    ))]
331     x_away_player7 <- x[match(as.character(train_obs$away_team_player7), rownames(x
    ))]
332     x_away_player8 <- x[match(as.character(train_obs$away_team_player8), rownames(x
    ))]
333     x_away_player9 <- x[match(as.character(train_obs$away_team_player9), rownames(x
    ))]
```

```

334     x_away_player10 <- x[match(as.character(train_obs$away_team_player10),rownames
      (x))]
335     x_away_player11 <- x[match(as.character(train_obs$away_team_player11),rownames
      (x))]
336     x_away <- apply(matrix(c(x_away_player1,x_away_player2,x_away_player3,
337                             x_away_player4,x_away_player5,x_away_player6,
338                             x_away_player7,x_away_player8,x_away_player9,
339                             x_away_player10,x_away_player11),
340                             ncol = 11),
341                             1, mean,na.rm = TRUE)
342     x_away[is.na(x_away)] <- 0
343
344     pmdiff <- x_home - x_away
345     train_pm <- c(train_pm,pmdiff)
346   }
347
348 }
349
350 # Use information for OLR
351 OLR <- polr(formula = 'result~pmdiff',
352             data = data.frame(result = factor(train_results, levels = c('A','D','H
353                                     '),ordered = TRUE),
354                                 pmdiff = train_pm))
355
356 # test data
357 test_data <- data %>% filter(Season == season_id,round > 5)
358 print(test_data %>% group_by(Season) %>% summarise(rounds = length(levels(factor(
359     round)))))
360
361 for(round_id in 6:max(test_data$round)){
362
363   #obs
364   obs_test <- test_data %>% filter(round == round_id)
365   print(paste('test',obs_test$Season[1],round_id))
366
367   # for every obs of in currenttime get result and prob
368   {
369     gd <- obs_test$home_goals - obs_test$away_goals
370     result <- ifelse(gd < 0, 'A', ifelse(gd == 0, 'D', 'H'))
371     #train_results <- c(train_results,result)
372
373     x_home_player1 <- x[match(as.character(obs_test$home_team_player1),rownames(x)
374                             )]
375     x_home_player2 <- x[match(as.character(obs_test$home_team_player2),rownames(x)
376                             )]
377     x_home_player3 <- x[match(as.character(obs_test$home_team_player3),rownames(x)
378                             )]

```

```
376     x_home_player4 <- x[match(as.character(obs_test$home_team_player4), rownames(x)
377                               )]
378     x_home_player5 <- x[match(as.character(obs_test$home_team_player5), rownames(x)
379                               )]
380     x_home_player6 <- x[match(as.character(obs_test$home_team_player6), rownames(x)
381                               )]
382     x_home_player7 <- x[match(as.character(obs_test$home_team_player7), rownames(x)
383                               )]
384     x_home_player8 <- x[match(as.character(obs_test$home_team_player8), rownames(x)
385                               )]
386     x_home_player9 <- x[match(as.character(obs_test$home_team_player9), rownames(x)
387                               )]
388     x_home_player10 <- x[match(as.character(obs_test$home_team_player10), rownames(
389                                x))]
390     x_home_player11 <- x[match(as.character(obs_test$home_team_player11), rownames(
391                                x))]
392     x_home_player12 <- x[match(rep('homeadv', nrow(obs_test)), rownames(x))]
393     x_home <- apply(matrix(c(x_home_player1, x_home_player2, x_home_player3,
394                               x_home_player4, x_home_player5, x_home_player6,
395                               x_home_player7, x_home_player8, x_home_player9,
396                               x_home_player10, x_home_player11, x_home_player12),
397                               ncol = 12),
398                      1, mean, na.rm = TRUE)
399     x_home[is.na(x_home)] <- 0
400
401     x_away_player1 <- x[match(as.character(obs_test$away_team_player1), rownames(x)
402                               )]
403     x_away_player2 <- x[match(as.character(obs_test$away_team_player2), rownames(x)
404                               )]
405     x_away_player3 <- x[match(as.character(obs_test$away_team_player3), rownames(x)
406                               )]
407     x_away_player4 <- x[match(as.character(obs_test$away_team_player4), rownames(x)
408                               )]
409     x_away_player5 <- x[match(as.character(obs_test$away_team_player5), rownames(x)
410                               )]
411     x_away_player6 <- x[match(as.character(obs_test$away_team_player6), rownames(x)
412                               )]
413     x_away_player7 <- x[match(as.character(obs_test$away_team_player7), rownames(x)
414                               )]
415     x_away_player8 <- x[match(as.character(obs_test$away_team_player8), rownames(x)
416                               )]
417     x_away_player9 <- x[match(as.character(obs_test$away_team_player9), rownames(x)
418                               )]
419     x_away_player10 <- x[match(as.character(obs_test$away_team_player10), rownames(
420                                x))]
421     x_away_player11 <- x[match(as.character(obs_test$away_team_player11), rownames(
422                                x))]
423     x_away <- apply(matrix(c(x_away_player1, x_away_player2, x_away_player3,
424                               x_away_player4, x_away_player5, x_away_player6,
```

```

406             x_away_player7,x_away_player8,x_away_player9,
407             x_away_player10,x_away_player11),
408             ncol = 11),
409             1, mean,na.rm = TRUE)
410
411     pmdiff <- x_home - x_away
412     prob <- predict(OLR,
413                   newdata = data.frame(result = factor(result, levels = c('A','D',
414                                     'H'),ordered = TRUE),
415                                     pmdiff = pmdiff),
416                                     type = 'p')
417
418     Prob_home <- c(Prob_home, as.numeric(prob[,3]))
419     Prob_draw <- c(Prob_draw, as.numeric(prob[,2]))
420     Prob_away <- c(Prob_away, as.numeric(prob[,1]))
421     observed <- c(observed,result)
422 }
423
424 # add test to train and update coefficients
425 train_data$seasonround <- NULL
426 train_data <- rbind(train_data,obs_test)
427
428 currenttime <- max(train_data$date)
429
430 {
431     df_seg_train <- df %>% filter((date <= currenttime)&((date >= currenttime -
432                                     yrsago*365)))
433     k <- 0.2
434     time_depre <- exp(-k*(as.numeric(currenttime) - df_seg_train$date)/365)
435 }
436
437 # get coefficients
438 {
439     beta <- df_seg_train$beta*time_depre
440
441     a <- as.matrix(df_seg_train[,-c(match('beta',colnames(df_seg_train)),
442                                     match('Season',colnames(df_seg_train)),
443                                     match('date',colnames(df_seg_train)),
444                                     match('round',colnames(df_seg_train)))])
445
446     a <- a * time_depre
447
448     lambda <- 3000
449     x <- solve(t(a) %*% a + diag(lambda,ncol(a))) %*% t(a) %*% beta
450 }
451
452 # for every obs of in currenttime get result and difference in pm of teams

```

```
453 {
454   gd <- obs_test$home_goals - obs_test$away_goals
455   result <- ifelse(gd < 0, 'A', ifelse(gd == 0, 'D', 'H'))
456   train_results <- c(train_results, result)
457
458   x_home_player1 <- x[match(as.character(obs_test$home_team_player1), rownames(x)
459     )]
460   x_home_player2 <- x[match(as.character(obs_test$home_team_player2), rownames(x)
461     )]
462   x_home_player3 <- x[match(as.character(obs_test$home_team_player3), rownames(x)
463     )]
464   x_home_player4 <- x[match(as.character(obs_test$home_team_player4), rownames(x)
465     )]
466   x_home_player5 <- x[match(as.character(obs_test$home_team_player5), rownames(x)
467     )]
468   x_home_player6 <- x[match(as.character(obs_test$home_team_player6), rownames(x)
469     )]
470   x_home_player7 <- x[match(as.character(obs_test$home_team_player7), rownames(x)
471     )]
472   x_home_player8 <- x[match(as.character(obs_test$home_team_player8), rownames(x)
473     )]
474   x_home_player9 <- x[match(as.character(obs_test$home_team_player9), rownames(x)
475     )]
476   x_home_player10 <- x[match(as.character(obs_test$home_team_player10), rownames(
477     x))]
478   x_home_player11 <- x[match(as.character(obs_test$home_team_player11), rownames(
479     x))]
480   x_home_player12 <- x[match(rep('homeadv', nrow(obs_test)), rownames(x))]
481   x_home <- apply(matrix(c(x_home_player1, x_home_player2, x_home_player3,
482     x_home_player4, x_home_player5, x_home_player6,
483     x_home_player7, x_home_player8, x_home_player9,
484     x_home_player10, x_home_player11, x_home_player12),
485     ncol = 12),
486     1, mean, na.rm = TRUE)
487   x_home[is.na(x_home)] <- 0
488
489   x_away_player1 <- x[match(as.character(obs_test$away_team_player1), rownames(x)
490     )]
491   x_away_player2 <- x[match(as.character(obs_test$away_team_player2), rownames(x)
492     )]
493   x_away_player3 <- x[match(as.character(obs_test$away_team_player3), rownames(x)
494     )]
495   x_away_player4 <- x[match(as.character(obs_test$away_team_player4), rownames(x)
496     )]
497   x_away_player5 <- x[match(as.character(obs_test$away_team_player5), rownames(x)
498     )]
499   x_away_player6 <- x[match(as.character(obs_test$away_team_player6), rownames(x)
500     )]
```

```

484     x_away_player7 <- x[match(as.character(obs_test$away_team_player7), rownames(x)
485                               )]
486     x_away_player8 <- x[match(as.character(obs_test$away_team_player8), rownames(x)
487                               )]
488     x_away_player9 <- x[match(as.character(obs_test$away_team_player9), rownames(x)
489                               )]
490     x_away_player10 <- x[match(as.character(obs_test$away_team_player10), rownames(
491                               x))]
492     x_away_player11 <- x[match(as.character(obs_test$away_team_player11), rownames(
493                               x))]
494     x_away <- apply(matrix(c(x_away_player1,x_away_player2,x_away_player3,
495                               x_away_player4,x_away_player5,x_away_player6,
496                               x_away_player7,x_away_player8,x_away_player9,
497                               x_away_player10,x_away_player11),
498                               ncol = 11),
499                               1, mean, na.rm = TRUE)
500     pmdiff <- x_home - x_away
501     train_pm <- c(train_pm, pmdiff)
502 }
503
504 # Use information for OLR
505 OLR <- polr(formula = 'result~pmdiff',
506             data = data.frame(result = factor(train_results, levels = c('A','D',
507                               'H'), ordered = TRUE),
508                               pmdiff = train_pm))
509
510
511 score <- data.frame(observed = factor(observed, levels = c('H', 'D', 'A')),
512                    win = Prob_home,
513                    draw = Prob_draw,
514                    loss = Prob_away)
515
516
517 saveRDS(object = score, file = 'PM_hvattum4y_new.rds')
518
519
520 #get scores
521
522 score <- readRDS('PM_hvattum4y_new.rds')
523 score$observed_win <- ifelse(score$observed == 'H', 1, 0)
524 score$observed_draw <- ifelse(score$observed == 'D', 1, 0)
525 score$observed_loss <- ifelse(score$observed == 'A', 1, 0)
526

```

```
527 score$IgnoranceScore <- apply(score[,2:7], 1, IGN)
528 score$BrierScore <- apply(score[,2:7], 1, BrierScore)
529 score$RPS <- apply(score[,2:7], 1, RPS2)
530
531
532 mean(score$IgnoranceScore)
533 sd(score$IgnoranceScore)
534
535 mean(score$BrierScore)
536 sd(score$BrierScore)
537
538 mean(score$RPS/2)
539 sd(score$RPS/2)
540
541 nrow(score)
542
543 s <- score[-c(2405:2466),]
544
545 mean(s$IgnoranceScore)
546 sd(s$IgnoranceScore)
547
548 mean(s$BrierScore)
549 sd(s$BrierScore)
550
551 mean(s$RPS/2)
552 sd(s$RPS/2)
553
554
555 s<- summary(OLR)
556 xtable(s$coefficients)
```

## A.4 Code for the Poisson based models

### A.4.1 Independent Poisson

```
1 #####
2 # Ranking soccer teams on current strength: MLE approach
3 # C. Ley
4 #####
5
6 # packages
7 {
8 setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/')
9 library(engsoccerdata)
10 library(tidyverse)
11 library(skellam)
12 library(caret)
```

```

13 source('BrierRpsIgnorance.R')
14 }
15
16 #DATA
17 {
18 # Premier League
19 PL_data <- england
20
21
22 # Data cleaning
23
24 # filter for right dataset
25 data <- PL_data %>%
26   filter(tier == 1,
27          Season %in% 2006:2017)
28
29 data %>% group_by(Season) %>%
30   summarise(matches = length(Date),
31             team = length(levels(factor(home))))
32
33
34 #make round parameter
35 id <- order(data$Season, data$Date)
36 data <- data[id,]
37 data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
38 }
39
40 #likelihood function
41
42 indeppoiss.mle <- function(par, data){
43
44   #parameters from data
45   teams <- sort(levels(factor(data$home)))
46   gi <- data$hgoal
47   gj <- data$vgoal
48   hteam <- as.character(data$home)
49   vteam <- as.character(data$visitor)
50
51   #restrictions
52   c <- abs(par[1])
53   h <- abs(par[2])
54   strength <- par[3:(length(par))]
55   last_strength <- -sum(strength)
56   strength <- c(strength, last_strength )
57
58   #features for likelihood
59   ri <- strength[match(hteam, teams)]
60   rj <- strength[match(vteam, teams)]
61

```

```
62 lambda_i <- exp(c + (ri+h)-(rj))
63 lambda_j <- exp(c + (rj)-(ri+h))
64
65 date <- as.Date(data$Date)
66 x <- as.numeric(max(as.Date(data$Date)) - date)
67 halfperiod <- 360
68 wtime <- 0.5 ^ (x/halfperiod)
69
70
71 lambda_i_exp_gi <- apply(matrix(c(lambda_i,gi),ncol = 2), 1,function(a){return(a
    [1]**a[2])})
72 rev_fact_gi <- 1/(factorial(gi))
73 exp_neg_lambda_i <- exp(-lambda_i)
74
75 lambda_j_exp_gj <- apply(matrix(c(lambda_j,gj),ncol = 2), 1,function(a){return(a
    [1]**a[2])})
76 rev_fact_gj <- 1/(factorial(gj))
77 exp_neg_lambda_j <- exp(-lambda_j)
78
79 matrix <- matrix(c(lambda_i_exp_gi,
80                     rev_fact_gi,
81                     exp_neg_lambda_i,
82                     lambda_j_exp_gj,
83                     rev_fact_gj,
84                     exp_neg_lambda_j),
85                  ncol=6)
86
87 densities <- apply(matrix, MARGIN = 1, FUN = prod)
88 densities <- apply(matrix(c(densities,wtime),ncol = 2), 1,function(a){return(a[1]*
    *a[2])})
89 return(-sum(log(densities)))
90
91 }
92
93
94 # 2008:max(season)
95
96 predict <- c()
97 observed <- c()
98
99 p_win <- p_draw <- p_loss <- c()
100 HAD <- c('H', 'D', 'A')
101
102 for (season_id in 2008:max(data$Season)){
103
104   # datasets
105
106   train_data <- data %>%
107     filter(Season %in% (season_id-2):(season_id-1))
```

```

108
109 train_data <- rbind( train_data, data %>% filter(Season == season_id,round <= 5))
110
111 print('training dataset')
112 print(train_data %>% group_by(as.character(Season)) %>% summarise(rounds = length(
    levels(factor(round)))))
113
114 test <- data %>% filter(Season == season_id,round > 5)
115
116 print(test %>% group_by(Season) %>% summarise(rounds = length(levels(factor(round)
    ))))
117
118 teams <- sort(unique(c(levels(factor(train_data$home)),levels(factor(train_data$
    visitor)))))
119
120 # initialize prediction
121
122 par <- c(1,1,rep(0,length(levels(factor(train_data$home)))-1))
123 par_optim <- optim(par = par,fn = indeppoiss.mle, data = train_data,gr = "BFGS")
124 par_optim
125
126 # predict
127
128 for (i in 6:max(test$round)){
129
130     # 1. MLE based on updated train_data
131     par_optim <- optim(par = par,fn = indeppoiss.mle, data = train_data,gr = "BFGS")
132
133     # 2. Parameter extraction
134     par <- par_optim$par
135
136     c <- abs(par[1])
137     h <- abs(par[2])
138
139     strength <- par[3:(length(par))]
140     last_strength <- -sum(strength)
141     strength <- c(strength, last_strength)
142
143     teams <- sort(unique(c(levels(factor(train_data$home)),levels(factor(train_data$
        visitor)))))
144
145
146     # 3. Parameters from test_data (observed)
147     t <- test %>% filter(round == i)
148     gi <- t$hgoal
149     gj <- t$vggoal
150     hteam <- as.character(t$home)
151     vteam <- as.character(t$visitor)
152

```

```
153 # 4. Features
154 ri <- strength[match(hteam, teams)]
155 rj <- strength[match(vteam, teams)]
156
157 lambda_i <- exp(c + (ri+h)-(rj))
158 lambda_j <- exp(c + (rj)-(ri+h))
159
160
161 # 5. Predict loss, draw, win
162 loss <- pskellam(q = -1,
163                 lambda1 = lambda_i,
164                 lambda2 = lambda_j)
165
166
167 draw <- dskellam(x = 0,
168                 lambda1 = lambda_i,
169                 lambda2 = lambda_j)
170
171
172 win <- 1-pskellam(q = 0,
173                 lambda1 = lambda_i,
174                 lambda2 = lambda_j)
175
176
177 p_win <- c(p_win, win)
178 p_draw <- c(p_draw, draw)
179 p_loss <- c(p_loss, loss)
180 # 6. Save predictions
181 observed <- c(observed, as.character(t$result))
182
183 # 7. Update train data
184 train_data <- rbind(train_data, t)
185 print(paste(season_id,i, ' '))
186 }
187
188
189
190 }
191
192 score <- data.frame(observed = factor(observed, levels = c('H', 'D', 'A')),
193                    win = p_win,
194                    draw = p_draw,
195                    loss = p_loss)
196 score
197
198 saveRDS(object = score, file = 'IndependentPoisson_CL2019_MLE.rds')
199
200
201 #get scores
```

```

202
203 score <- readRDS('IndependentPoisson_CL2019_MLE.rds')
204 score$observed_win <- ifelse(score$observed == 'H', 1, 0)
205 score$observed_draw <- ifelse(score$observed == 'D', 1, 0)
206 score$observed_loss <- ifelse(score$observed == 'A', 1, 0)
207
208 score$IgnoranceScore <- apply(score[,2:7], 1, IGN)
209 score$BrierScore <- apply(score[,2:7], 1, BrierScore)
210 score$RPS <- apply(score[,2:7], 1, RPS2)
211
212 mean(score$IgnoranceScore)
213 mean(score$BrierScore)
214 mean(score$RPS)/2
215 sd(score$RPS/2)

```

### A.4.2 Bivariate Poisson

```

1 #####
2 # Ranking soccer teams on current strength: MLE approach
3 # C. Ley
4 #####
5
6 # packages
7 {
8 setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/')
9 library(engsoccerdata)
10 library(tidyverse)
11 library(skellam)
12 library(caret)
13 source('BrierRpsIgnorance.R')
14 }
15
16 #Data
17 {
18   # Premier League
19   PL_data <- england
20
21
22   # Data cleaning
23
24   # filter for right dataset
25   data <- PL_data %>%
26     filter(tier == 1,
27            Season %in% 2006:2017)
28
29   data %>% group_by(Season) %>%
30     summarise(matches = length(Date),
31               team = length(levels(factor(home))))
32

```

```
33
34 #make round parameter
35 id <- order(data$Season, data$Date)
36 data <- data[id,]
37 data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
38 }
39
40 #functions
41 {
42   Sumation_function <- function(row, lambda_c){
43
44     lambda_i <- row[1]
45     lambda_j <- row[2]
46     gi <- row[3]
47     gj <- row[4]
48     min_goals <- row[5]
49
50     sumation <- 0
51     for (k in 0:min_goals){
52       sumation <- sumation + choose(gi,k)*choose(gj,k)*factorial(k)*((lambda_c/(
53         lambda_i*lambda_j))**k)
54     }
55     return(sumation)
56   }
57   #likelihood function
58
59   bivpoiss.mle <- function(par, data){
60
61     #parameters from data
62     teams <- sort(unique(levels(factor(data$home)),levels(factor(data$visitor))))
63     gi <- data$hgoal
64     gj <- data$vgoal
65     hteam <- as.character(data$home)
66     vteam <- as.character(data$visitor)
67
68     #restrictions
69     c <- abs(par[1])
70     h <- abs(par[2])
71     lambda_c <- abs(par[3])
72     strength <- par[4:(length(par))]
73     last_strength <- -sum(strength)
74     strength <- c(strength, last_strength )
75
76     #features for likelihood
77     ri <- strength[match(hteam, teams)]
78     rj <- strength[match(vteam, teams)]
79
80     lambda_i <- exp(c + (ri+h)-(rj))
```

```

81   lambda_j <- exp(c + (rj)-(ri+h))
82
83   date <- as.Date(data$Date)
84   x <- as.numeric(max(as.Date(data$Date)) - date)
85   halfperiod <- 390
86   wtime <- 0.5 ^ (x/halfperiod)
87
88
89   lambda_i_exp_gi <- apply(matrix(c(lambda_i,gi),ncol = 2), 1,function(a){return(a
    [1]**a[2])})
90   rev_fact_gi <- 1/(factorial(gi))
91
92   lambda_j_exp_gj <- apply(matrix(c(lambda_j,gj),ncol = 2), 1,function(a){return(a
    [1]**a[2])})
93   rev_fact_gj <- 1/(factorial(gj))
94   exp_ijc <- exp(-(lambda_i+lambda_j+lambda_c))
95
96
97   min_goals <- apply(matrix(c(gi,gj),ncol = 2), MARGIN = 1 , min)
98   sumation_matrix <- matrix(c(lambda_i, lambda_j, gi , gj, min_goals), ncol = 5)
99   sumations <- apply(sumation_matrix, MARGIN = 1, FUN = Sumation_function, lambda_
    c)
100
101
102   matrix <- matrix(c(lambda_i_exp_gi,
103                      rev_fact_gi,
104                      lambda_j_exp_gj,
105                      rev_fact_gj,
106                      exp_ijc,
107                      sumations),
108                      ncol=6)
109
110   densities <- apply(matrix, MARGIN = 1, FUN = prod)
111
112   densities <- apply(matrix(c(densities,wtime),ncol = 2), 1,function(a){return(a
    [1]**a[2])})
113
114
115   return(-sum(log(densities)))
116
117 }
118
119
120 }
121
122
123
124
125 # 2008:max(season)

```

```
126
127 predict <- c()
128 observed <- c()
129
130 p_win <- p_draw <- p_loss <- c()
131 HAD <- c('H', 'D', 'A')
132
133 for (season_id in 2008:max(data$Season)){
134   # datasets
135   train_data <- data %>%
136     filter(Season %in% (season_id-2):(season_id-1))
137
138   train_data <- rbind( train_data, data %>% filter(Season == season_id, round <= 5))
139
140   print('training dataset')
141   print(train_data %>% group_by(as.character(Season)) %>% summarise(rounds = length(
142     levels(factor(round)))))
143
144   test <- data %>% filter(Season == season_id, round > 5)
145
146   print(test %>% group_by(Season) %>% summarise(rounds = length(levels(factor(round)
147     ))))
148
149   # initialize prediction
150   teams <- sort(unique(c(levels(factor(train_data$home)), levels(factor(train_data$
151     visitor)))))
152
153   par <- c(1,1,1,rep(0,length(levels(factor(train_data$home)))-1))
154   par_optim <- optim(par = par,fn = bivpoiss.mle, data = train_data,gr = "BFGS")
155   par <- par_optim$par
156   par
157   # predict
158   for (i in 6:max(test$round)){
159     # 1. MLE based on updated train_data
160     par_optim <- optim(par = par,fn = bivpoiss.mle, data = train_data,gr = "BFGS")
161
162     # 2. Parameter extraction
163     par <- par_optim$par
164
165     c <- abs(par[1])
166     h <- abs(par[2])
167     lambda_c <- par[3]
168     strength <- par[4:(length(par))]
169     last_strength <- -sum(strength)
170     strength <- c(strength, last_strength)
171     teams <- sort(unique(c(levels(factor(train_data$home)), levels(factor(train_data$
172       visitor)))))
```

```

171
172
173   # 3. Parameters from test_data (observed)
174   t <- test %>% filter(round == i)
175   gi <- t$hgoal
176   gj <- t$vgoal
177   hteam <- as.character(t$home)
178   vteam <- as.character(t$visitor)
179
180   # 4. Features
181   ri <- strength[match(hteam, teams)]
182   rj <- strength[match(vteam, teams)]
183
184   lambda_i <- exp(c + (ri+h)-(rj))
185   lambda_j <- exp(c + (rj)-(ri+h))
186
187
188   # 5. Predict loss, draw, win
189   loss <- pskellam(q = -1,
190                   lambda1 = lambda_i,
191                   lambda2 = lambda_j)
192
193
194   draw <- dskellam(x = 0,
195                   lambda1 = lambda_i,
196                   lambda2 = lambda_j)
197
198
199   win <- 1-pskellam(q = 0,
200                   lambda1 = lambda_i,
201                   lambda2 = lambda_j)
202
203
204   p_win <- c(p_win, win)
205   p_draw <- c(p_draw, draw)
206   p_loss <- c(p_loss, loss)
207   # 6. Save prediction
208   observed <- c(observed, as.character(t$result))
209
210   # 7. Update train data
211   train_data <- rbind(train_data, t)
212   print(paste(season_id,i, ' '))
213 }
214
215
216
217 }
218
219 score <- data.frame(observed = factor(observed, levels = c('H', 'D', 'A')),

```

```
220             win = p_win,
221             draw = p_draw,
222             loss = p_loss)
223 score
224
225
226 saveRDS(object = score,file = 'BivariatePoisson_CL2019_MLE.rds')
227
228
229 #get scores
230
231 score <- readRDS('BivariatePoisson_CL2019_MLE.rds')
232 score$observed_win <- ifelse(score$observed == 'H', 1 , 0)
233 score$observed_draw <- ifelse(score$observed == 'D', 1 , 0)
234 score$observed_loss <- ifelse(score$observed == 'A', 1 , 0)
235
236 score$IgnoranceScore <- apply(score[,2:7], 1, IGN)
237 score$BrierScore <- apply(score[,2:7], 1, BrierScore)
238 score$RPS <- apply(score[,2:7], 1, RPS2)
239
240
241 mean(score$IgnoranceScore)
242 sd(score$IgnoranceScore)
243
244 mean(score$BrierScore)
245 sd(score$BrierScore)
246
247 mean(score$RPS/2)
248 sd(score$RPS/2)
```

## A.5 Code for the Weibull count model

### A.5.1 Weibull count

```
1 #####
2 # A bivariate Weibull count model for forecasting association football scores
3 # G. Boshnakov
4 #####
5
6
7 # packages
8 {
9   setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/'
10   )
11   library(engsoccerdata)
12   library(Countr)
13   library(skellam)
```

```

13 library(caret)
14 library(MASS)
15 library(tidyverse)
16 source('BrierRpsIgnorance.R')
17 }
18
19 # data
20 {
21   # Premier League
22   PL_data <- england
23   PL_data%>% filter(tier == 1)
24
25   data <- PL_data %>%
26     filter(tier == 1,
27            Season %in% 2006:2015)
28
29   #make round parameter
30   id <- order(data$Season, data$Date)
31   data <- data[id,]
32   data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
33
34 }
35
36 # Frank copula function
37 {
38   Frank.copula <- function(u,v,k){
39     teller <- (exp(-k*u)-1)*(exp(-k*v)-1)
40     noemer <- (exp(-k)-1)
41     fraction <- teller/noemer
42     copula <- -1/k*log(1+fraction)
43     return(copula)
44   }
45 }
46 }
47
48 # Cumulative density weibull function
49 {
50   CDWeibull <- function(yi,shape, scale){
51
52     if (yi < 0) {densities <- 0}
53     else {
54       seq <- c(0, seq_len(yi))
55       densities <- dWeibullCount(seq,shape,scale)
56       return(sum(densities))
57     }
58
59   }
60 }
61

```

```
62 # Likelihood function
63 {
64   WeibullLL <- function(par, data){
65
66     # parameters
67     teams <- sort(unique(levels(factor(data$home)), levels(factor(data$visitor))))
68     atk_teams <- par[1:length(teams)]# alpha's
69     def_teams <- par[(length(teams)+1):(2*length(teams))] # beta's
70     h <- par[(2*length(teams))+1] # home team advantage
71     k <- par[(2*length(teams))+2] # used for copula
72     Ch <- par[(2*length(teams))+3]
73     Ca <- par[(2*length(teams))+4]
74
75     #parameters from data
76     gi <- data$hgoal
77     gj <- data$vgoal
78     hteam <- as.character(data$home)
79     vteam <- as.character(data$visitor)
80
81     # restrictions ?
82
83     # derivations
84     atk_i <- atk_teams[match(hteam, teams)]
85     def_i <- def_teams[match(hteam, teams)]
86     lambda_i <- exp(atk_i+def_i+h)
87
88     atk_j <- atk_teams[match(vteam, teams)]
89     def_j <- def_teams[match(vteam, teams)]
90     lambda_j <- exp(atk_j+def_j)
91
92     # Time variety
93     epsilon <- 1/500
94     date <- as.Date(data$Date)
95     x <- as.numeric(max(as.Date(data$Date)) - date)
96     wtime <- exp(-epsilon*x)
97
98     # Weibull densities
99     MarginalDensHome <- apply(data.frame(x=gi,
100                                         shape=Ch,
101                                         scale=lambda_i),
102                               MARGIN = 1,
103                               FUN = function(a){return(dWeibullCount(x = a[1],shape
104                               = a[2],scale = a[3]))})
104
105
106
107     MarginalDensAway <- apply(data.frame(x=gj,
108                                         shape=Ca,
109                                         scale=lambda_j),
```

```

110             MARGIN = 1,
111             FUN = function(a){return(dWeibullCount(x = a[1],shape
              = a[2],scale = a[3]))})
112
113     # Cumulative densities
114
115     F1 <- apply(data.frame(x=gi,
116                           shape=Ch,
117                           scale=lambda_i),
118                MARGIN = 1,
119                FUN = function(a){return(CDWeibull(yi = a[1],shape = a[2],scale
              = a[3]))})
120
121     F2 <- apply(data.frame(x=gj,
122                           shape=Ca,
123                           scale=lambda_j),
124                MARGIN = 1,
125                FUN = function(a){return(CDWeibull(yi = a[1],shape = a[2],scale = a
              [3]))})
126
127     Fm1 <- apply(data.frame(x=gi-1,
128                            shape=Ch,
129                            scale=lambda_i),
130                 MARGIN = 1,
131                 FUN = function(a){return(CDWeibull(yi = a[1],shape = a[2],scale =
              a[3]))})
132
133     Fm2 <- apply(data.frame(x=gj-1,
134                            shape=Ca,
135                            scale=lambda_j),
136                 MARGIN = 1,
137                 FUN = function(a){return(CDWeibull(yi = a[1],shape = a[2],scale
              = a[3]))})
138
139
140     L <- Frank.copula(u = F1,v = F2, k = k) - Frank.copula(u = Fm1,v = F2, k = k) -
      Frank.copula(u = F1,v = Fm2, k = k) + Frank.copula(u = Fm1,v = Fm2, k = k)
141
142     # own made up
143     densities <- apply(matrix(c(L,wtime),ncol = 3), 1,prod)
144
145     return(-sum(log(densities)))
146
147 }
148
149
150 }
151
152

```

```
153 # Algorithm
154 Prob_home <- c()
155 Prob_draw <- c()
156 Prob_away <- c()
157 observed_home <- c()
158 observed_draw <- c()
159 observed_away <- c()
160 for (season_id in 2008:2015){
161
162   # Training data
163   train_data <- data %>%
164     filter(Season %in% (season_id-2):(season_id-1))
165   train_data <- rbind( train_data, data %>% filter(Season == season_id,round <= 5))
166   print('training dataset')
167   print(train_data %>% group_by(as.character(Season)) %>% summarise(rounds = length(
168     levels(factor(round))))))
169
170   # Testing data
171   test <- data %>% filter(Season == season_id,round > 5)
172   print(test %>% group_by(Season) %>% summarise(rounds = length(levels(factor(round)
173     ))))
174
175   # parameters
176   teams <- sort(unique(levels(factor(train_data$home)),levels(factor(train_data$
177     visitor))))
178   atk_teams <- rep(0,length(teams)) # alpha's
179   def_teams <- rep(0,length(teams)) # beta's
180   h <- 0.2948 # home team advantage
181   k <- -0.4561 # used for copula
182   Ch <- 1.050
183   Ca <- 0.9831
184
185   par <- c(atk_teams,def_teams,h,k,Ch,Ca)
186
187   par_optim <- optim(par = par,fn = WeibullLL, data = train_data,gr = "BFGS")
188
189   for (round_id in 6:(max(test$round))){
190
191     # Observation
192     obs <- test %>% filter(round == round_id)
193
194     # 1. MLE based on updated train_data
195     par_optim <- optim(par = par,fn = WeibullLL, data = train_data,gr = "BFGS")
196
197     # 2. Parameter extraction
198     par <- par_optim$par
199     teams <- sort(unique(levels(factor(train_data$home)),levels(factor(train_data$
200       visitor))))
201     atk_teams <- par[1:length(teams)]# alpha's
```

```

198   def_teams <- par[(length(teams)+1):(2*length(teams))] # beta's
199   h <- par[(2*length(teams))+1] # home team advantage
200   k <- par[(2*length(teams))+2] # used for copula
201   Ch <- par[(2*length(teams))+3]
202   Ca <- par[(2*length(teams))+4]
203
204   # 3. Parameters from test_data (observed)
205   gi <- obs$hgoal
206   gj <- obs$vgoal
207   hteam <- as.character(obs$home)
208   vteam <- as.character(obs$visitor)
209
210   # 4. Features
211   atk_i <- atk_teams[match(hteam, teams)]
212   def_i <- def_teams[match(hteam, teams)]
213   lambda_i <- exp(atk_i+def_i+h)
214
215   atk_j <- atk_teams[match(vteam, teams)]
216   def_j <- def_teams[match(vteam, teams)]
217   lambda_j <- exp(atk_j+def_j)
218
219
220   # 5. Predict loss, draw, win
221
222   # prob draw
223   for (i in 1:nrow(obs)){
224     draw <- c()
225     for (goal in 0:60){
226       ProbForGoal <- dWeibullCount(x=goal,shape = Ch, scale = lambda_i[i])*
227         dWeibullCount(x=goal,shape = Ca, scale = lambda_j[i])
228       draw <- c(draw, ProbForGoal)
229     }
230     p_draw <- sum(draw)
231
232     win <- c()
233     for (goal in 0:60){
234       ProbForGoal <- dWeibullCount(x=goal,shape = Ca, scale = lambda_j[i])*(1-
235         CDWeibull(yi = goal, shape= Ch, scale = lambda_i[i]))
236       win <- c(win, ProbForGoal)
237     }
238     p_win <- sum(win)
239
240     loss <- c()
241     for (goal in 0:60){
242       ProbForGoal <- dWeibullCount(x=goal,shape = Ch, scale = lambda_i[i])*(1-
243         CDWeibull(yi = goal, shape= Ca, scale = lambda_j[i]))

```

```
244     p_loss <- sum(loss)
245
246     if(sum(p_draw,p_win,p_loss)!=1){print('Sum probabilities is not 1')}
247
248     # 6. Save predictions
249     if (as.character(obs$result[i]) == 'H'){
250         observed_home <- c(observed_home,1)
251         observed_draw <- c(observed_draw,0)
252         observed_away <- c(observed_away,0)
253
254     }
255     else {
256         if (as.character(obs$result[i]) == 'A'){
257             observed_home <- c(observed_home,0)
258             observed_draw <- c(observed_draw,0)
259             observed_away <- c(observed_away,1)
260
261         }
262         else {
263             observed_home <- c(observed_home,0)
264             observed_draw <- c(observed_draw,1)
265             observed_away <- c(observed_away,0)
266
267         }
268     }
269
270
271     Prob_home <- c(Prob_home,p_win)
272     Prob_draw <- c(Prob_draw,p_draw)
273     Prob_away <- c(Prob_away,p_loss)
274 }
275 # 7. Update train data
276 train_data <- rbind(train_data, obs)
277 print(paste(season_id,round_id, ' '))
278 }
279 }
280
281
282 out <- data.frame(pwin = Prob_home,
283                  pdraw = Prob_draw,
284                  ploss = Prob_away,
285                  owin = observed_home,
286                  odraw = observed_draw,
287                  oloss = observed_away)
288
289 all.equal(apply(out[,1:3],1,sum),rep(1,nrow(out)))
290
291 out$ignorance <- apply(out[,1:6], 1, IGN)
292 out$brierScore <-apply(out[,1:6], 1, BrierScore)
```

```
293 out$RPS <- apply(out[,1:6], 1, RPS2)
294
295
296 mean(out$ignorance)
297 sd(out$ignorance)
298
299 mean(out$brierScore)
300 sd(out$brierScore)
301
302 mean(out$RPS/2)
303 sd(out$RPS/2)
304
305 saveRDS(object = out, file = 'BivariateWeibullCountCopula.rds')
```

## A.6 Code for the machine learning models

### A.6.1 Random forest Baboota

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # In[1]:
5
6
7 get_ipython().run_line_magic('matplotlib', 'inline')
8 import itertools
9 import numpy as np
10 import pandas as pd
11 import seaborn as sns
12 from scipy import interp
13 from pprint import pprint
14 from itertools import cycle
15 import matplotlib.pyplot as plt
16 from collections import OrderedDict
17 from sklearn.externals import joblib
18 from xgboost.sklearn import XGBClassifier
19 #from TrainTestSplit import TrainTestSplit
20 from sklearn.metrics import roc_curve, auc
21 from sklearn.metrics import confusion_matrix
22 from sklearn.model_selection import GridSearchCV
23 from sklearn.preprocessing import label_binarize
24 from sklearn.model_selection import ShuffleSplit
25 from sklearn.metrics import classification_report
26 from sklearn.multiclass import OneVsRestClassifier
27 from sklearn.model_selection import learning_curve
28 from sklearn.model_selection import cross_val_score
29 from sklearn.ensemble import RandomForestClassifier
```

```
30 from sklearn.impute import SimpleImputer
31
32
33 np.set_printoptions(precision = 10)
34
35 # load data
36 path = "C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/Predicting-English-
Premier-Results-master/EngineeredData.csv"
37 data = pd.read_csv(path)
38
39 data.head(5)
40
41 # Change Season
42
43 vector = []
44 for i in range(len(data['Season'])):
45     vector.append(data['Season'][i][0:4])
46 data['Season'] = vector
47
48 # Make Round Feature
49 data = data.sort_values(by=['Date'])
50
51 roundsperseason = np.repeat([i for i in range(1,39)],10)
52
53 listofdata = []
54 for season_id in sorted(set(data['Season'])):
55     test = data.copy()
56     test = test[test['Season'] == season_id]
57     test['Round'] = roundsperseason
58     listofdata.append(test)
59
60 data = pd.concat(listofdata)
61 data = data.sort_values(by=['Date'])
62
63 # Features needed
64 A = ["ATGD", "HTGD", "ACKPP", "HCKPP", "AGKPP", "HGKPP", "ASTKPP", "HSTKPP", "ASt",
      "HSt", "AStWeighted", "HStWeighted", "AForm", "HForm", "AOverall", "HOverall", "
      AAttack", "HAttack", "AMidfield", "HMidfield", "ADefense", "HDefense"]
65 B = ["GD", "CKPP", "GKPP", "STKPP", "Streak", "WeightedStreak", "Form", "Overall",
      "Attack", "Midfield", "Defense"]
66
67 frames = []
68 for season_id in range(2008,2016):
69
70     #1. make train and test
71     train = data[(data['Season'] == str(season_id-2)) | (data['Season'] == str(
        season_id-1))].copy()
72     firstrounds = data[(data['Season'] == str(season_id)) & (data['Round'] <= 5)].
        copy()
```

```

73     train = pd.concat([train,firstrounds])
74     test = data[(data['Season'] == str(season_id)) & (data['Round'] > 5)].copy()
75
76     for match_id in test.index:
77
78         #1. define match to predict
79         obs = test[test.index == match_id].copy()
80
81         #2. Split in x and y
82         train_x = train[A+B]
83         train_y = train['FTR']
84         test_x = obs[A+B]
85         test_y = obs['FTR']
86
87         train_x = train_x.fillna(0)
88         test_x = test_x.fillna(0)
89
90         #3. define classifier
91         Classifier = RandomForestClassifier(bootstrap=True,
92                                           class_weight=None,
93                                           criterion='gini',
94                                           max_depth=None,
95                                           max_features='log2',
96                                           max_leaf_nodes=None,
97                                           min_impurity_split=1e-07,
98                                           min_samples_leaf=2,
99                                           min_samples_split=100,
100                                          min_weight_fraction_leaf=0.0,
101                                          n_estimators=150,
102                                          n_jobs=-1,
103                                          oob_score=True,
104                                          random_state=None,
105                                          verbose=0,
106                                          warm_start=False)
107
108         #4. Train and test model
109         Classifier.fit(train_x,train_y)
110
111         #5. Predict
112         YPred = Classifier.predict(test_x)
113         classLabels = ['H','A','D']
114         Yprob = Classifier.predict_proba(test_x)
115
116         #6.output
117         probs = pd.DataFrame(Yprob, columns=['p_loss','p_draw','p_win'])
118         probs['Observed'] = YPred
119         frames.append(probs)
120
121         #7. update train

```

```
122         train = train.append(obs)
123 out = pd.concat(frames)
124 out.to_csv("C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/RFbaboota.csv",
            index=False)
```

## A.6.2 Gradient boosting

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # In[1]:
5
6
7 get_ipython().run_line_magic('matplotlib', 'inline')
8 import itertools
9 import numpy as np
10 import pandas as pd
11 import seaborn as sns
12 from scipy import interp
13 from pprint import pprint
14 from itertools import cycle
15 import matplotlib.pyplot as plt
16 from collections import OrderedDict
17 from sklearn.externals import joblib
18 from xgboost.sklearn import XGBClassifier
19 #from TrainTestSplit import TrainTestSplit
20 from sklearn.metrics import roc_curve, auc
21 from sklearn.metrics import confusion_matrix
22 from sklearn.model_selection import GridSearchCV
23 from sklearn.preprocessing import label_binarize
24 from sklearn.model_selection import ShuffleSplit
25 from sklearn.metrics import classification_report
26 from sklearn.multiclass import OneVsRestClassifier
27 from sklearn.model_selection import learning_curve
28 from sklearn.model_selection import cross_val_score
29
30
31 np.set_printoptions(precision = 10)
32
33
34 # In[2]:
35
36
37 # load data
38 path = "C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/Predicting-English-
Premier-Results-master/EngineeredData.csv"
39 data = pd.read_csv(path)
40
41 data.head(5)
```

```
42
43
44 # ## Data engineering
45 #
46
47 # In[3]:
48
49
50 # Change Season
51
52 vector = []
53 for i in range(len(data['Season'])):
54     vector.append(data['Season'][i][0:4])
55 data['Season'] = vector
56
57
58 # In[4]:
59
60
61 # Make Round Feature
62 data = data.sort_values(by=['Date'])
63
64 roundsperseason = np.repeat([i for i in range(1,39)],10)
65
66 listofdata = []
67 for season_id in sorted(set(data['Season'])):
68     test = data.copy()
69     test = test[test['Season'] == season_id]
70     test['Round'] = roundsperseason
71     listofdata.append(test)
72
73 data = pd.concat(listofdata)
74 data = data.sort_values(by=['Date'])
75
76
77 # In[5]:
78
79
80 # Features needed
81 A = ["ATGD", "HTGD", "ACKPP", "HCKPP", "AGKPP", "HGKPP", "ASTKPP", "HSTKPP", "ASt",
      "HSt", "AStWeighted", "HStWeighted", "AForm", "HForm", "AOverall", "HOverall", "
      AAttack", "HAttack", "AMidfield", "HMidfield", "ADefense", "HDefense"]
82 B = ["GD", "CKPP", "GKPP", "STKPP", "Streak", "WeightedStreak", "Form", "Overall",
      "Attack", "Midfield", "Defense"]
83
84
85 # In[7]:
86
87
```

```
88 season_id = 2008
89
90 #1. make train and test
91 train = data[(data['Season'] == str(season_id-2)) | (data['Season'] == str(season_id
    -1))].copy()
92 firstrounds = data[(data['Season'] == str(season_id)) & (data['Round'] <= 5)].copy()
93 train = pd.concat([train,firstrounds])
94 test = data[(data['Season'] == str(season_id)) & (data['Round'] > 5)].copy()
95
96
97 # In[11]:
98
99
100 test
101
102
103 # In[10]:
104
105
106 test['Round']
107
108
109 # In[20]:
110
111
112 frames = []
113 rounds = []
114 seasons = []
115
116 for season_id in range(2008,2016):
117
118     #1. make train and test
119     train = data[(data['Season'] == str(season_id-2)) | (data['Season'] == str(
        season_id-1))].copy()
120     firstrounds = data[(data['Season'] == str(season_id)) & (data['Round'] <= 5)].
        copy()
121     train = pd.concat([train,firstrounds])
122     test = data[(data['Season'] == str(season_id)) & (data['Round'] > 5)].copy()
123
124     #train = train[['FTR']+B].dropna(axis=0)
125     #test = test[['FTR']+B].dropna(axis=0)
126
127     for match_id in test.index:
128
129         #1. define match to predict
130         obs = test[test.index == match_id].copy()
131         rounds.append(int(obs['Round']))
132         seasons.append(int(obs['Season']))
133
```

```

134
135     #2. Split in x and y
136     train_x = train[B]
137     train_y = train['FTR']
138     test_x = obs[B]
139     test_y = obs['FTR']
140
141
142     #3. define classiffier
143     Classifier = XGBClassifier(base_score=0.5,
144                               booster='gbtree',
145                               colsample_bylevel=1,
146                               colsample_bytree=1,
147                               gamma=0.2,
148                               learning_rate=0.05,
149                               max_delta_step=0,
150                               max_depth=3,
151                               min_child_weight=1,
152                               n_estimators=100,
153                               n_jobs=-1,
154                               nthread=None,
155                               objective='multi:softprob',
156                               random_state=0,
157                               reg_alpha=0,
158                               reg_lambda=0.6,
159                               scale_pos_weight=1,
160                               seed=None,
161                               silent=True,
162                               subsample=1)
163
164     #4. Train and test model
165     Classifier.fit(train_x,train_y)
166
167     #5. Predict
168     YPred = Classifier.predict(test_x)
169     classLabels = ['H', 'A', 'D']
170     Yprob = Classifier.predict_proba(test_x)
171
172     #6.output
173     probs = pd.DataFrame(Yprob, columns=['p_loss', 'p_draw', 'p_win'])
174     probs['Observed'] = YPred
175     frames.append(probs)
176
177     #7. update train
178     train = train.append(obs)
179 out = pd.concat(frames)
180 out.to_csv("C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/Xgboost_dropna.
181           csv",index=False)

```

```
182
183 # In[22]:
184
185
186 out['season'] = seasons
187 out['round'] = rounds
188
189 out.to_csv("C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/
           Xgboostevolution.csv",index=False)
190
191
192 # In[40]:
193
194
195 importances = Classifier.feature_importances_
196 indices = np.argsort(importances[::-1])
197 names = [B[i] for i in indices]
198
199 # Create plot
200 plt.figure()
201
202 # Create plot title
203 plt.title("Feature Importance")
204
205 # Add bars
206 plt.barh(range(train_x.shape[1]), importances[indices])
207
208 # Add feature names as x-axis labels
209 plt.yticks(range(train_x.shape[1]), names)
210
211 # Show plot
212 plt.show()
213
214
215 # In[44]:
216
217
218 names
219
220
221 # In[43]:
222
223
224 pd.DataFrame(zip(names,importances), columns=['feature','importance']).to_csv("C:/
           Users/Thiebe/Documents/academie jaar 2019-2020/thesis/Xgboostfeatureimportance.
           csv",index=False)
```

## A.7 Code for the hybrid models

### A.7.1 Random forest Groll

```

1 # hybrid random forest data
2
3 # packages
4 {
5   setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/'
6         )
7   library(engsoccerdata)
8   library(tidyverse)
9   library(mlr)
10  library(skellam)
11  library(Hmisc)
12  library(party)
13  source('BrierRpsIgnorance.R')
14
15  substrRight <- function(x, n){
16    substr(x, nchar(x)-n+1, nchar(x))
17  }
18
19 # Data
20 {
21   # basic data (outcomes)
22   {
23     # Premier League
24     PL_data <- england
25
26     data <- PL_data %>%
27       filter(tier == 1,
28              Season %in% 2006:2015)
29
30     #make round parameter
31     id <- order(data$Season, data$Date)
32     data <- data[id,]
33     data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
34
35     # data
36     data$Date <- as.Date(data$Date)
37     data$key <- paste0(data$Date, substr(as.character(data$home),1,1), substr(as.
38                                     character(data$visitor),1,1), data$hgoal, data$vggoal)
39     droplevels(data$home)
40     droplevels(data$visitor)
41   }
42 }
43
44 # elo

```

```
43 {
44   elo_data <- readRDS('elo_data.rds')
45   elo_data$Date <- as.Date(elo_data$Date)
46   elo_data$Season <- NULL
47   merged <- merge(data, elo_data, by = c('Date', 'home', 'visitor'))
48
49 }
50
51
52 # age check why this doesnt work!
53 {
54   agedata <- read.csv('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/
      models/EPL_player_age_data.csv')
55   agedata$dob <- as.Date(agedata$dob, format = '%d/%m/%Y')
56   agedata$startseason <- as.Date(agedata$startseason)
57   agedata$age <- agedata$startseason - agedata$dob
58   agedata$age <- as.numeric(agedata$age)/365.25
59   agedata$season <- as.numeric(agedata$season)
60   agedata <- agedata %>% group_by(season, club) %>% summarise(mean_age = mean(age))
61   agedata$club <- str_remove(agedata$club, ' FC')
62   agedata$club <- str_remove(agedata$club, ' AFC')
63   colnames(agedata) <- c('Season', 'home', 'homeage')
64   merged <- merge(merged, agedata, by = c('Season', 'home'))
65   merged$homeage <- with(merged, impute(homeage, mean))
66   colnames(agedata) <- c('Season', 'visitor', 'awayage')
67   merged <- merge(merged, agedata, by = c('Season', 'visitor'))
68   merged$awayage <- with(merged, impute(awayage, mean))
69
70
71 }
72
73 # fifa
74 {
75   fifa <- read.csv("C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/
      Predicting-English-Premier-Results-master/EngineeredData.csv")
76   fifa <- fifa[, c('Date', 'HomeTeam', 'AwayTeam',
77                   'HAttack', 'HDefense', 'HMidfield', 'HOverall' ,
78                   'AAttack', 'ADefense', 'AMidfield', 'AOverall' )]
79   #fifa$Season <- substr(fifa$Season, start = 6, stop = 9)
80   colnames(fifa) <- c('Date', 'home', 'visitor',
81                       'HAttack', 'HDefense', 'HMidfield', 'HOverall' ,
82                       'AAttack', 'ADefense', 'AMidfield', 'AOverall' )
83
84
85   fifa$home <- str_remove(fifa$home, ' FC')
86   fifa$home <- str_remove(fifa$home, ' AFC')
87   fifa$home[fifa$home == 'Spurs'] <- 'Tottenham Hotspur'
88
89   fifa$visitor <- str_remove(fifa$visitor, ' FC')
```

```

90   fifa$visitor <- str_remove(fifa$visitor, ' AFC')
91   fifa$visitor[fifa$visitor== 'Spurs'] <- 'Tottenham Hotspur'
92
93
94   merged$k <- paste0(merged$Date,
95                     substr(str_remove(merged$home, 'AFC '),1,1),
96                     substr(str_remove(merged$visitor, 'AFC '),1,1))
97
98   fifa$k <- paste0(fifa$Date,
99                   substr(str_remove(fifa$home, 'AFC '),1,1),
100                   substr(str_remove(fifa$visitor, 'AFC '),1,1))
101
102   fifa$Date <- fifa$home <- fifa$visitor <- NULL
103
104
105   merged <- merge(merged, fifa, by=c('k'))
106
107   #paste0(merged$Date,merged$home,merged$visitor)[! merged$k %in% fifa$k]
108   merged$k<- NULL
109
110 }
111
112
113 # oddset
114 {
115   Oddset <- readRDS('Oddset.rds')
116   Oddset$home <- Oddset$visitor <- Oddset$Date <- NULL
117   merged <- merge(merged, Oddset, by = c('key'))
118   merged$key <- NULL
119
120 }
121
122
123
124
125
126 # final Data set
127 {
128   data <- data.frame(result = merged$result,
129                      hgoal = merged$hgoal,
130                      vgoal = merged$vgoal,
131                      elo = merged$elo_home-merged$elo_away,
132                      age = merged$homeage - merged$awayage,
133                      betHome = merged$B365H,
134                      betDraw = merged$B365D,
135                      betAway = merged$B365A,
136                      round = merged$round,
137                      Season = merged$Season,
138                      FifaAttack = merged$HAttack - merged$AAttack,

```

```
139         FifaMidfield = merged$HMidfield - merged$AMidfield,
140         FifaDefense = merged$HDefense - merged$ADefense,
141         FifaOverall = merged$HOverall - merged$AOverall,
142         date = as.Date(merged$Date)
143     )
144     print(data %>% group_by(as.character(Season)) %>% summarise(rounds = length(
145         round)))
146
147     id <- order(data$date)
148     data <- data[id,]
149 }
150
151 #Algorithm
152 Prob_home <- c()
153 Prob_draw <- c()
154 Prob_away <- c()
155 observed_home <- c()
156 observed_draw <- c()
157 observed_away <- c()
158 for (season_id in 2008:2015){
159
160     # Training data
161     {
162         train_data <- data %>%
163             filter(Season %in% (season_id-2):(season_id-1))
164         train_data <- rbind( train_data, data %>% filter(Season == season_id, round <= 5)
165             )
166         print('training dataset')
167         print(train_data %>% group_by(as.character(Season)) %>% summarise(rounds =
168             length(levels(factor(round)))))
169     }
170
171     # Testing data
172     {
173         test <- data %>% filter(Season == season_id, round > 5)
174         print(test %>% group_by(Season) %>% summarise(rounds = length(levels(factor(
175             round)))))
176     }
177
178     # Hyperparameters
179     {
180         B <- 5000
181     }
182
183 }
```

```

184 for (round_id in 6:(max(test$round))){
185
186
187   print(paste(season_id, ' ', round_id))
188
189
190   # Observation
191   obs <- test %>% filter(round == round_id)
192   train_data$age <- as.numeric(train_data$age)
193   obs$age <- as.numeric(obs$age)
194
195   # 1. Train Random Forest
196
197   cf <- cforest(formula = hgoal + vgoal ~ elo + age +betHome + betDraw + betAway
198                 +FifaAttack + FifaMidfield + FifaDefense + Fifa0Verall,
199                 data = train_data,
200                 controls = cforest_control(ntree = B))
201
202   # 3. Predict goal difference test
203   observed_goals <- obs[,c('hgoal', 'vgoal')]
204   colnames(observed_goals) <- c('hgoal_observed', 'vgoal_observed')
205   predicted_goals <- predict(cf, newdata = obs, type = "response")
206   predicted_goals <- matrix(unlist(predicted_goals), ncol = 2, byrow = TRUE)
207   colnames(predicted_goals) <- c('hgoal_predicted', 'vgoal_predicted')
208   cbind(observed_goals, predicted_goals)
209
210   # 4. Predict win draww loss with skellam
211   p_loss <- pskellam(q = -1,
212                     lambda1 = predicted_goals[,1],
213                     lambda2 = predicted_goals[,2])
214
215
216   p_draw <- dskellam(x = 0,
217                     lambda1 = predicted_goals[,1],
218                     lambda2 = predicted_goals[,2])
219
220
221   p_win <- 1-pskellam(q = 0,
222                     lambda1 = predicted_goals[,1],
223                     lambda2 = predicted_goals[,2])
224
225
226   Prob_home <- c(Prob_home, p_win)
227   Prob_draw <- c(Prob_draw, p_draw)
228   Prob_away <- c(Prob_away, p_loss)
229
230   for (i in 1:nrow(obs)){
231     if (as.character(obs$result[i]) == 'H'){

```

```
232     observed_home <- c(observed_home,1)
233     observed_draw <- c(observed_draw,0)
234     observed_away <- c(observed_away,0)
235
236   }
237   else {
238     if (as.character(obs$result[i]) == 'A'){
239       observed_home <- c(observed_home,0)
240       observed_draw <- c(observed_draw,0)
241       observed_away <- c(observed_away,1)
242     }
243   }
244   else {
245     observed_home <- c(observed_home,0)
246     observed_draw <- c(observed_draw,1)
247     observed_away <- c(observed_away,0)
248   }
249 }
250 }
251 }
252
253 # 7. Update train data
254 train_data <- rbind(train_data, obs)
255 }
256
257
258 }
259
260 out <- data.frame(pwin = Prob_home,
261                  pdraw = Prob_draw,
262                  ploss = Prob_away,
263                  owin = observed_home,
264                  odraw = observed_draw,
265                  oloss = observed_away)
266
267 all.equal(apply(out[,1:3],1,sum),rep(1,nrow(out)))
268
269 out$ignorance <- apply(out[,1:6], 1, IGN)
270 out$brierScore <- apply(out[,1:6], 1, BrierScore)
271 out$RPS <- apply(out[,1:6], 1, RPS2)
272
273
274 mean(out$ignorance)
275 sd(out$ignorance)
276
277 mean(out$brierScore)
278 sd(out$brierScore)
279
280 mean(out$RPS/2)
```

```

281 sd(out$RPS/2)
282
283 saveRDS(object = out, file = 'RFgroll.rds')

```

## A.7.2 Hybrid random forest

```

1  # hybrid random forest data
2
3  # packages
4  {
5    setwd('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/models/R scripts/'
6          )
7    library(engsoccerdata)
8    library(tidyverse)
9    library(mlr)
10   library(skellam)
11   library(Hmisc)
12   library(party)
13   source('BrierRpsIgnorance.R')
14
15   substrRight <- function(x, n){
16     substr(x, nchar(x)-n+1, nchar(x))
17   }
18
19  # Data
20  {
21    # basic data (outcomes)
22    {
23      # Premier League
24      PL_data <- england
25
26      data <- PL_data %>%
27        filter(tier == 1,
28               Season %in% 2006:2015)
29
30      #make round parameter
31      id <- order(data$Season, data$Date)
32      data <- data[id,]
33      data$round <- rep(sort(rep(1:38,10)), max(data$Season)- min(data$Season) +1 )
34
35      # data
36      data$Date <- as.Date(data$Date)
37      data$key <- paste0(data$Date, substr(as.character(data$home),1,1), substr(as.
38                                         character(data$visitor),1,1), data$hgoal, data$vggoal)
39
40      droplevels(data$home)
41      droplevels(data$visitor)
42    }
43  }
44
45  }

```

```
42 # elo
43 {
44   elo_data <- readRDS('elo_data.rds')
45   elo_data$Date <- as.Date(elo_data$Date)
46   elo_data$Season <- NULL
47   merged <- merge(data, elo_data, by = c('Date', 'home', 'visitor'))
48
49 }
50
51 # abilities
52 {
53   ability_data <- readRDS('ability_data_final.rds')
54   colnames(ability_data) <- c("Date", "Season", "home", "visitor", "ability_home", "
      ability_away")
55   ability_data$Date <- as.Date(ability_data$Date)
56   ability_data$Season <- NULL
57   merged <- merge(merged, ability_data, by = c('Date', 'home', 'visitor'))
58
59 }
60
61 # age check why this doesnt work!
62 {
63   agedata <- read.csv('c:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/
      models/EPL_player_age_data.csv')
64   agedata$dob <- as.Date(agedata$dob, format = '%d/%m/%Y')
65   agedata$startseason <- as.Date(agedata$startseason)
66   agedata$age <- agedata$startseason - agedata$dob
67   agedata$age <- as.numeric(agedata$age)/365.25
68   agedata$season <- as.numeric(agedata$season)
69   agedata <- agedata %>% group_by(season, club) %>% summarise(mean_age = mean(age))
70   agedata$club <- str_remove(agedata$club, ' FC')
71   agedata$club <- str_remove(agedata$club, ' AFC')
72   colnames(agedata) <- c('Season', 'home', 'homeage')
73   merged <- merge(merged, agedata, by = c('Season', 'home'))
74   merged$homeage <- with(merged, impute(homeage, mean))
75   colnames(agedata) <- c('Season', 'visitor', 'awayage')
76   merged <- merge(merged, agedata, by = c('Season', 'visitor'))
77   merged$awayage <- with(merged, impute(awayage, mean))
78
79
80 }
81
82 # fifa
83 {
84   fifa <- read.csv("C:/Users/Thiebe/Documents/academie jaar 2019-2020/thesis/
      Predicting-English-Premier-Results-master/EngineeredData.csv")
85   fifa <- fifa[, c('Date', 'HomeTeam', 'AwayTeam',
86                   'HAttack', 'HDefense', 'HMidfield', 'HOverall',
87                   'AAttack', 'ADefense', 'AMidfield', 'AOverall' )]
```

```

88   #fifa$Season <- substr(fifa$Season,start = 6,stop = 9)
89   colnames(fifa)<- c('Date','home','visitor',
90                     'HAttack','HDefense','HMidfield','HOverall' ,
91                     'AAttack','ADefense','AMidfield','AOverall' )
92
93
94   fifa$home <- str_remove(fifa$home,' FC')
95   fifa$home <- str_remove(fifa$home,' AFC')
96   fifa$home[fifa$home == 'Spurs'] <- 'Tottenham Hotspur'
97
98   fifa$visitor <- str_remove(fifa$visitor,' FC')
99   fifa$visitor <- str_remove(fifa$visitor,' AFC')
100  fifa$visitor[fifa$visitor== 'Spurs'] <- 'Tottenham Hotspur'
101
102
103  merged$k <- paste0(merged$Date,
104                    substr(str_remove(merged$home,'AFC '),1,1),
105                    substr(str_remove(merged$visitor,'AFC '),1,1))
106
107  fifa$k <- paste0(fifa$Date,
108                  substr(str_remove(fifa$home,'AFC '),1,1),
109                  substr(str_remove(fifa$visitor,'AFC '),1,1))
110
111  fifa$Date <- fifa$home <- fifa$visitor <- NULL
112
113
114  merged <- merge(merged,fifa,by=c('k'))
115
116  #paste0(merged$Date,merged$home,merged$visitor)[! merged$k %in% fifa$k]
117  merged$k<- NULL
118
119 }
120
121
122 # oddset
123 {
124   Oddset <- readRDS('Oddset.rds')
125   Oddset$home <- Oddset$visitor <- Oddset$Date <- NULL
126   merged <- merge(merged, Oddset,by = c('key'))
127   merged$key <- NULL
128
129 }
130
131
132
133
134
135 # final Data set
136 {

```

```
137   data <- data.frame(result = merged$result,
138                       hgoal = merged$hgoal,
139                       vgoal = merged$vgoal,
140                       elo = merged$elo_home-merged$elo_away,
141                       ability = merged$ability_home-merged$ability_away,
142                       age = merged$homeage - merged$awayage,
143                       betHome = merged$B365H,
144                       betDraw = merged$B365D,
145                       betAway = merged$B365A,
146                       round = merged$round,
147                       Season = merged$Season,
148                       FifaAttack = merged$HAttack - merged$AAttack,
149                       FifaMidfield = merged$HMidfield - merged$AMidfield,
150                       FifaDefense = merged$HDefense - merged$ADefense,
151                       FifaOverall = merged$HOverall - merged$AOverall,
152                       date = as.Date(merged$Date)
153   )
154   print(data %>% group_by(as.character(Season)) %>% summarise(rounds = length(
155       round)))
156   id <- order(data$date)
157   data <- data[id,]
158 }
159 }
160
161 #Algorithm
162 Prob_home <- c()
163 Prob_draw <- c()
164 Prob_away <- c()
165 observed_home <- c()
166 observed_draw <- c()
167 observed_away <- c()
168 for (season_id in 2008:2015){
169
170   # Training data
171   {
172     train_data <- data %>%
173       filter(Season %in% (season_id-2):(season_id-1))
174     train_data <- rbind( train_data, data %>% filter(Season == season_id, round <= 5)
175       )
176     print('training dataset')
177     print(train_data %>% group_by(as.character(Season)) %>% summarise(rounds =
178       length(levels(factor(round)))))
179   }
180
181   # Testing data
182   {
183     test <- data %>% filter(Season == season_id, round > 5)
```

```

182   print(test %>% group_by(Season) %>% summarise(rounds = length(levels(factor(
183     round)))))
184 }
185 # Hyperparameters
186 {
187   B <- 5000
188
189
190
191 }
192
193
194 for (round_id in 6:(max(test$round))){
195
196
197   print(paste(season_id, ' ', round_id))
198
199
200   # Observation
201   obs <- test %>% filter(round == round_id)
202   train_data$age <- as.numeric(train_data$age)
203   obs$age <- as.numeric(obs$age)
204
205   # 1. Train Random Forest
206
207   cf <- cforest(formula = hgoal + vgoal ~ elo + ability + age +betHome + betDraw
208     + betAway +FifaAttack + FifaMidfield + FifaDefense + FifaOverall,
209     data = train_data,
210     controls = cforest_control(ntree = B))
211
212   # 3. Predict goal difference test
213   observed_goals <- obs[,c('hgoal', 'vgoal')]
214   colnames(observed_goals) <- c('hgoal_observed', 'vgoal_observed')
215   predicted_goals <- predict(cf, newdata = obs, type = "response")
216   predicted_goals <- matrix(unlist(predicted_goals), ncol = 2, byrow = TRUE)
217   colnames(predicted_goals) <- c('hgoal_predicted', 'vgoal_predicted')
218   cbind(observed_goals, predicted_goals)
219
220   # 4. Predict win draww loss with skellam
221   p_loss <- pskellam(q = -1,
222     lambda1 = predicted_goals[,1],
223     lambda2 = predicted_goals[,2])
224
225
226   p_draw <- dskellam(x = 0,
227     lambda1 = predicted_goals[,1],
228     lambda2 = predicted_goals[,2])

```

```
229
230
231   p_win <- 1-pskellam(q = 0,
232                     lambda1 = predicted_goals[,1],
233                     lambda2 = predicted_goals[,2])
234
235
236   Prob_home <- c(Prob_home,p_win)
237   Prob_draw <- c(Prob_draw,p_draw)
238   Prob_away <- c(Prob_away,p_loss)
239
240   for (i in 1:nrow(obs)){
241     if (as.character(obs$result[i]) == 'H'){
242       observed_home <- c(observed_home,1)
243       observed_draw <- c(observed_draw,0)
244       observed_away <- c(observed_away,0)
245
246     }
247     else {
248       if (as.character(obs$result[i]) == 'A'){
249         observed_home <- c(observed_home,0)
250         observed_draw <- c(observed_draw,0)
251         observed_away <- c(observed_away,1)
252
253       }
254       else {
255         observed_home <- c(observed_home,0)
256         observed_draw <- c(observed_draw,1)
257         observed_away <- c(observed_away,0)
258
259       }
260     }
261   }
262
263   # 7. Update train data
264   train_data <- rbind(train_data, obs)
265 }
266
267
268 }
269
270 out <- data.frame(pwin = Prob_home,
271                  pdraw = Prob_draw,
272                  ploss = Prob_away,
273                  owin = observed_home,
274                  odraw = observed_draw,
275                  oloss = observed_away)
276
277 all.equal(apply(out[,1:3],1,sum),rep(1,nrow(out)))
```

```
278
279 out$ignorance <- apply(out[,1:6], 1, IGN)
280 out$brierScore <- apply(out[,1:6], 1, BrierScore)
281 out$RPS <- apply(out[,1:6], 1, RPS2)
282
283
284 mean(out$ignorance)
285 sd(out$ignorance)
286
287 mean(out$brierScore)
288 sd(out$brierScore)
289
290 mean(out$RPS/2)
291 sd(out$RPS/2)
292
293 saveRDS(object = out, file = 'hybridRF.rds')
```