

Imperial College London

Business Analytics Report

---

## Data in Football

---

**A machine learning-driven approach to predict game outcomes**

---

Author:

Cyrill Studer  
CID: 01456269

September 2, 2019

(Word Count: 5,163)

## Executive Summary

Over the recent years, an arms race to quantify the working mechanism of football games by making sense of the exploding amount of information and data available between opposing teams, analytics companies and bookmakers is visible. This development exposes the number-heavy sports betting market to newly arising opportunities and treats for bettors and bookmakers. The paper at hand develops a public data-based, machine learning-driven sports betting strategy for the English Premier League season 2018/19 by using team characteristics input features as well as time-series performance data in combination with dimensionality reduction methods and optimisation techniques. The final validated strategy achieved a return of investment of 38.38 % by betting on 8% of all possible games.

# Table of Contents

<b>1</b>	<b>Introduction and Motivation .....</b>	<b>1</b>
1.1	Research Subject.....	1
1.2	Literature Review.....	1
1.3	Structure.....	3
<b>2</b>	<b>Methodology .....</b>	<b>3</b>
<b>3</b>	<b>Data .....</b>	<b>6</b>
3.1	Sourcing.....	6
3.2	Exploratory Data Analysis.....	6
<b>4</b>	<b>Modeling .....</b>	<b>8</b>
4.1	Training Data .....	8
4.2	Validation Data .....	12
<b>5</b>	<b>Discussion .....</b>	<b>13</b>
<b>6</b>	<b>Conclusion .....</b>	<b>14</b>
6.1	Results.....	14
6.2	Limitations .....	14
6.3	Outlook .....	14
<b>References.....</b>		<b>16</b>
<b>Appendix .....</b>		<b>19</b>

## List of Abbreviations

Bn	Billions
EDA	Exploratory Data Analysis
KNN	K-Nearest Neighbor
LDA	Linear Discriminant Analysis
PCA	Principal Component Analysis
ROI	Return of Investment

# 1 Introduction and Motivation

## 1.1 Research Subject

*“The ball is round, the game lasts ninety minutes, and everything else is theory.”*

- *Sepp Herberger (1898-1977), Former Head Coach of the German National Football Team*

Football was long seen as a chaotic, unstructured and complex game where a few decisive actions that are accompanied by a high amount of variance decide the outcomes of games. The common opinion was that the critical, controllable elements for success were mainly intangible factors such as talent, passion or team spirit. Room and application areas for structured information processing, analytics and forecasting were limited for a long time. However, with the technological progress for gathering and analysing data, the tighter margins between opposing teams and the increasing economic significance of the football business, data science departments within football clubs and football analytics companies such as StatsBomb or Opta Sports are ever since on the rise. (Biermann, 2019)

Football analytics companies collect up to 2,000 geospatial and action-based data points per game from an individual player-based perspective as well as from a team perspective (Burn-Murdoch, 2018). The collected data can then be used for the individual player and team performance assessment, scouting and transfers, training and game plan design or the development of sports betting strategies. This paper will focus on modeling and forecasting Premier League game outcomes for the season 2018/19 with publicly available data in order to design a profitable betting strategy. To put the economic relevance of the topic into perspective, the English Premier League is estimated to be a € 5.44 bn industry on which bets worth € 67 bn are placed every season (Deloitte, 2019; Statista, 2019a).

## 1.2 Literature Review

In the following, the goal is to provide an outline of the advancement of scientific approaches to modeling and forecasting football results. The review will start with traditional statistical approaches, followed by the upcoming use of modern machine learning techniques and end with today's state-of-the-art use of individual player data.

The first statistical approaches to model football results were based on using results only as input variables. One of the earliest approaches of showing that football results are not occurring purely by chance goes back to Hill (1974) in 1974. The researcher found a strong correlation between a team's scored goals in the previous season and a team's scored goals in the current season. However, it was only in 1997 when researchers first came up with a complete statistical model to forecast English League game outcomes to systematically challenge the odds of bookmakers. Dixon & Coles (1997) used maximum likelihood to develop comprehensive parametric Poisson regression models to forecast goals in order to assign probabilities to potential scores. This model lied the foundation for a profitable betting strategy. Crowder et al.'s (2002) work builds upon Dixon & Coles's (1997) approach of applying Poisson regression models. The researchers improved the accuracy of Dixon & Coles model by enhancing the input data layer with a proxy variable for team strength based on past performances. In order to create this proxy, Markov Chains and Monte Carlo Simulations were used.

In the year 2005, Goddard (2005) made one of the first approaches of not only using previous results as input variables but to supplement the model with information such as the geographical distance between the two opposing teams or the teams' performances in other competitions. Moreover, the researcher further enhanced Dixon & Coles's approach with an ordered probit model. However, only slightly better forecasting power could be achieved. Some years later, Hucaljuk & Rakipović (2011) made a first approach of applying modern machine learning techniques such as KNN, random forest or neural networks while using input features such as the recent results, the teams' current rankings, the number of injured players, the average number of scored goals or the outcome of previous meetings between the two opposing teams. The researchers achieved a forecasting accuracy between 50-65% for predicting wins by the home team, draws or wins by the away team. However, the sample size was very small, and the variance of the accuracy was high depending on the statistical method used. Tax. et al. (2015) used a very similar approach to Hucaljuk & Rakipović (2011) by combining the previous work with dimensionality reduction techniques. The team achieved a stable accuracy of 49% using public data from the Dutch league.

In terms of input features, a recent popular trend is to model player and team strength. The basic concept goes back to Elo (1978) who modeled player strength to predict chess games. For football, Constantinou et al. (2013) and Lasek (2016) developed dynamic team rating systems based on features such as recent performances, previous meetings between the teams or the location of the match. Another hot trend is to use individual player data such as average number

of shots, average distance of shots, average angle of shots, defender proximity for shots, or the number of passes to create significant goal-scoring opportunities, also called packing, in order to predict the number of expected goals by a team (Lucey, 2015; Impect, 2019). However, to design these input features, geospatial and detailed player individual data is necessary. This information is not publicly available.

### 1.3 Structure

This paper is willing to combine different, in the previous chapter outlined modeling techniques and approaches from credit default forecasting methods to lay the foundation for the design of a profitable strategy in today's betting market. The paper's goal will be approached in four steps. First of all, based on the insights of the literature review, the methodology of the paper's modeling approach will be outlined. Subsequently, the sourcing of the data will be explained before an exploratory data analysis (EDA) will help to find relevant patterns in the data, check the validity of the study and guide the engineering of the most meaningful input features. Thirdly, different models will be assessed before the practical relevance of the results will be reflected. Finally, a conclusion chapter will summarise the results, outline limitations and highlight the further potential of the study.

## 2 Methodology

Based on the literature review, three findings for accurate football game outcome forecasting are apparent. First of all, the highest prediction accuracy was achieved by using state-of-the-art machine learning techniques in combination with dimensionality reduction techniques (Tax. et al., 2015). Secondly, team strength input features, as well as time-series performance data, are crucial to achieving high-scoring results (Constantinou et al., 2013; Lasek, 2016). Thirdly, player individual performance data enhances models further. As player individual performance data is not publicly available, the focus will lie on finding the best possible data sources, proxies and modeling techniques for findings one and two. Additionally, betting odds for the Premier League season 2018/19 will be sourced.

Graph 1 illustrates the planned usage of data. The input layer data consists of match-based features such as the number of collected points, scored goals and different kinds of shots taken for rolling windows of the last four to eight matches. Rolling windows are defined as averaging the number of actions for a certain metric over a defined period of past games. The modeling

of team characteristic is considered from three different perspectives. Firstly, the team characteristics are modeled as team strength. The strength will be based on the ratings of the video game Fifa 19. Each year, a group of experts systematically evaluates the abilities of each player within a team to subsequently assign a strength score to each team. Secondly, the Premier League teams' squad net worth is considered. The third perspective of team characteristics is the average squad salary as it is assumed that there is, as in every competitive market, a strong correlation between performance and salary. Considering all the above-mentioned data points, a data input layer of over sixty independent variables will be provided. In terms of the target variable, there are three different target variables modeled. The first target variable will be *home win/ away win/ draw*. The second target variable will be *home win/ no home win* and the third target variable will be *away win/ no away win*. The prediction of draws will be disregarded for target variables two and three as draws are the least frequently occurring outcome. Therefore, draws are harder to predict and associated with more variance. As a next step, the data is split into three parts. The first part, the first 80 matches of the season (the first eight gamedays), are used to initialise the rolling window time-series data. The next 200 matches (gamedays 9-28) will constitute the training data set. Lastly, the last 100 games of the season (gamedays 29-38) will be used as a validation data set.

**Graph 1: Planned Data Usage**

DATA SETUP		DATA SPLIT			
	Input Variables (X)	Target Variables (Y)	Initiation Data	Training Data	Validation Data
Match-Based	<ul style="list-style-type: none"> <li>• Shots</li> <li>• Shots on Target</li> <li>• Goals Conceded</li> <li>• Points</li> </ul>	<ul style="list-style-type: none"> <li>a. Home Win/ Draw/ Away Win</li> <li>b. Home Win/ No Home Win</li> <li>c. Away Win/ No Away Win</li> </ul>	Gameday 0-8	Gameday 9-28	Gameday 29-38
	<ul style="list-style-type: none"> <li>• Wage</li> <li>• Strength</li> <li>• Value</li> </ul>		Initialise Rolling Windows (4-8 Gamedays) for Match-Based Features	Cross-Validate, Optimise and Train Different Prediction Models	Validate and Test Developed Betting Strategies On New Data

To start with the modeling of the training data, four machine learning techniques are used to predict the target variables. The reason for evaluating multiple machine learning techniques is based on the fact that certain models lead to better results given specific characteristics of the data. As a next step, each method is twice applied for each of the three target variables, once

with the full set of input variables and once with a principal component-reduced (PCA) set of input variables. Dimensionality reduction techniques such as PCA allow reducing the number of input variables while conserving all critical information and therefore decreasing noise and creating sparse models (James et al., 2013). The number of principal components is determined by cross-validation. Cross-validation is a technique to further split the training data into smaller subsets in order to train and test a method on differently composed subsets. This allows to reduce variance and lower bias by addressing underfitting and overfitting tendencies as well as the low stability of results given the dependency on the composition of the particular training data set (James et al., 2013). Subsequently, the accuracy of each machine learning technique is evaluated by means of ten-fold cross-validation. Next, all the different techniques will be used to predict the probability for each possible match outcome.

To convert the predicted probabilities into a promising betting strategy, a recursion optimisation function will be developed to determine the optimal probability threshold (cutoff) of above which probabilistic certainty to invest into a predicted outcome in order to maximise the winnings. This idea comes from the area of credit default forecasting (Kürüm et al., 2012). In a first step, given the probabilities of the different outcomes for each match, the function places a bet of one on outcomes that exceed a given probability and does not bet on all other predicted match outcomes. As a second step, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. In case of mismatch, the loss of the betting amount of one is noted. As a third step, the sum of wins and losses as well as the total betting amount are accumulated for each probability threshold between 0.5 and 1 in steps of 0.01. Finally, the probability threshold (cutoff value) with the highest return of investment (ROI), defined as the sum of wins and losses plus the betting amount divided by the betting amount, will be evaluated. Subsequently, based on the ROI-maximisation criteria and the stability of the achieved results, the most promising machine learning technique and its corresponding probability threshold will then be applied to evaluate the investment strategy's out of sample performance on the validation data set. The stability of the results will be tested by the graphical inspection of the ROI development for all possible probability thresholds. The above-described function will be run for each of the three target variables and the details of the functions working mechanism can be found in Appendix 11.

## 3 Data

### 3.1 Sourcing

The match data for all 360 Premier League matches of the season 2018/19 was obtained from the website football-data.co.uk (2019). The dataset contains for each match 61 features. However, the majority of the features are betting odds from different betting companies. Therefore, a first reduced dataset with match prediction relevant data was created. The obtained data set includes features such as goals, shorts or shorts on target for both the home and the away team. For all features, time-series data as described in the methodology chapter was created. In terms of modeling team characteristics, the first feature squad net worth was based on the values assigned by the community of football enthusiasts and experts from transfermarkt.com (2019). The wage bills of each team for the season 2018/19 were obtained from the Statista (2019b). Lastly, as already mentioned, the strengths of the teams are approximated by the assigned values of the video game Fifa 19 and were scraped from fifaindex.com (2019).

The betting odds for home wins, away wins and draws for all 380 matches were considered from the betting company Bet365. These odds were included in the dataset from football-data.co.uk (2019) already. It is important to see that the betting odds of the major players in the British sports betting market only deviate marginally. Therefore, the odds of Bet365, one of the largest betting companies in Great Britain, constitute a representative market to develop a relevant betting strategy.

### 3.2 Exploratory Data Analysis

As a first step of the EDA, the goal is to evaluate whether there are significant differences between the characteristics of the different data sets. Table 1 shows the summary statistics of the match features for each of the three data sets.

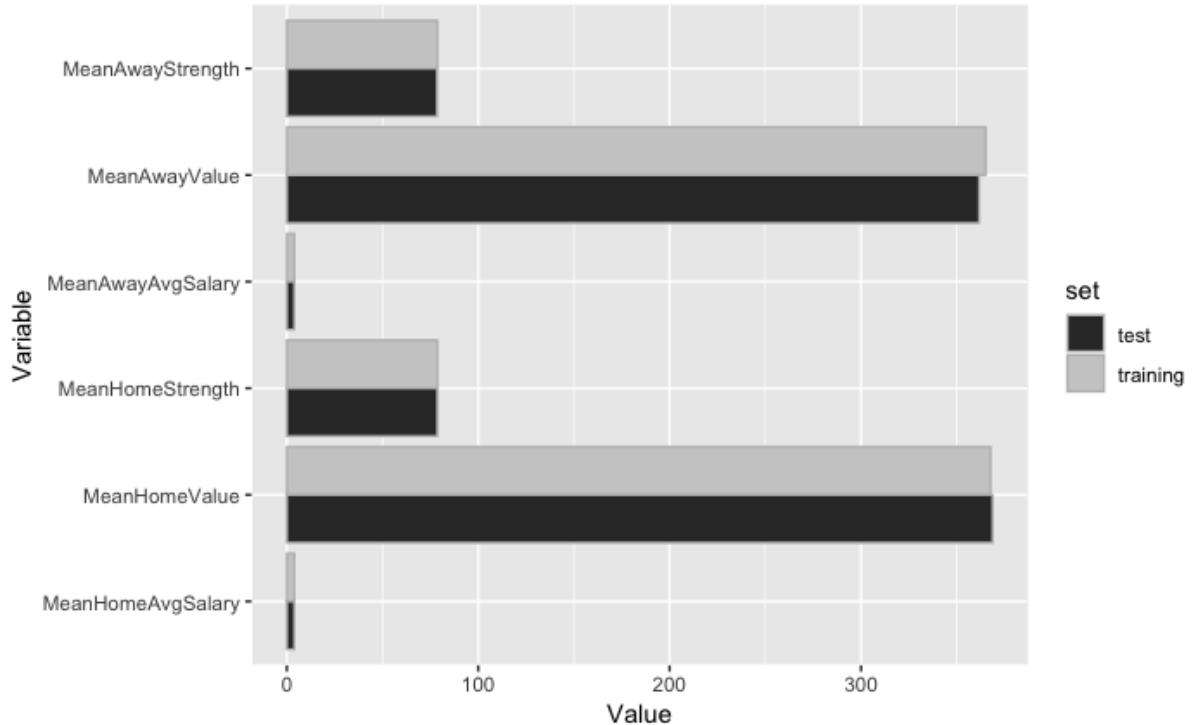
**Table 1: Summary Statistics Match Features**

	<b>Initialisation</b>	<b>Training</b>	<b>Validation</b>	<b>Overall</b>
<b>HomeGoals (<math>\emptyset</math>)</b>	1.475	1.605	1.570	1.57
<b>AwayGoals (<math>\emptyset</math>)</b>	1.338	1.230	1.230	1.25
<b>HomeShots (<math>\emptyset</math>)</b>	13.863	14.270	14.080	14.13
<b>AwayShots (<math>\emptyset</math>)</b>	11.438	10.710	11.780	11.15
<b>HomeShotsTarget (<math>\emptyset</math>)</b>	4.725	4.875	4.630	4.78
<b>AwayShotsTarget (<math>\emptyset</math>)</b>	4.525	3.705	3.900	3.93
<b>HomeWins (%)</b>	42.50	49.00	49.00	47.63
<b>AwayWins (%)</b>	37.50	32.00	34.00	33.68
<b>Draws (%)</b>	20.00	19.00	17.00	18.69
<b>Count</b>	80	200	100	380

The table provides two major insights. Firstly, it is important to notice that the outcomes of the target variables are not perfectly balanced. Home wins occur almost three times as likely as draws and approximately 50% more often than away wins. This is important to take into consideration when choosing appropriate machine learning prediction techniques at a later stage. Secondly, the characteristics of the training and validation data sets are highly similar while there is a slight discrepancy between the initialisation data set and the training and validation data sets. This is likely due to the small sample size and potentially also to the particular fixtures within each data set. For example, the higher amount of away goals and the lower amount of home goals in the initialisation data set could be caused by a higher proportion of weaker teams playing stronger opponents at home during that particular period of the season. Nevertheless, as no algorithm is trained, and no predictions are made with the initialisation data set, the differences have no impact on the results of the study. Therefore, in the following, the focus will lie on identifying patterns concerning the training and validation data set.

Graph 2 shows that also in terms of team-based features, the differences between the observations in the training data set and the observations in the validation data set are marginal. Therefore, the assumption of a homogenous test and validation data set holds. Consequently, the results of the study should be meaningful and representative.

**Graph 2: Summary Statistics Team-Based Features**



As a next step, a closer look at the correlation tables shows that all team-based features show relevant correlations with the *home win/ draw/ away win* target variable in the range of |0.21-0.26|. With regard to the time-series data, it appears that the rolling windows of the last six to eight games show the highest correlations with the target variables. The average shots on target, the average number of points and the average number of goals from both the away and the home team show the highest correlations with the *home win/ draw/ away win* target variable (range |0.14-0.19|) while other variables show lower correlations. Consequently, the results of the correlation tables indicate that dimensionality reduction techniques could be important tools to create models with a reduced amount of noise and therefore a higher predicting power. All the details for the correlation tables can be found in Appendix 1-10.

## 4 Modeling

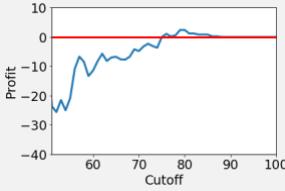
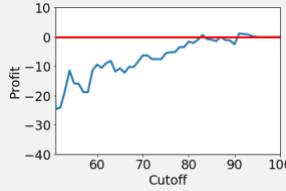
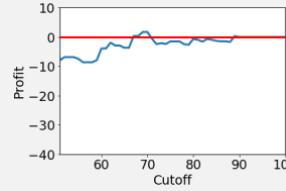
### 4.1 Training Data

As outlined before, the analysis will be based on the application of four different machine learning models to different input variables and output variables. In the following, the working

pattern of each of the four methods as well as their suitability for the particular problem at hand will be shortly introduced. Subsequently, the results for each of the modeling techniques for the training data set will be depicted.

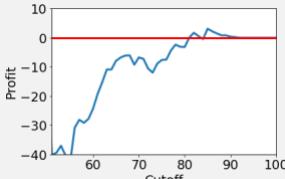
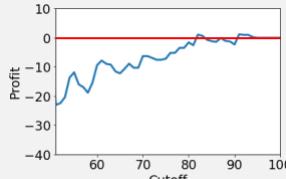
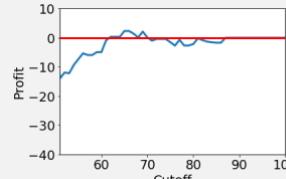
The first considered machine learning algorithm to predict the Premier League games outcomes of the season 2018/19 is linear discriminate analysis (LDA). “LDA is a linear, multivariate function that models conditional distributions of the target variable in order to allocate observations to different groups or assigning a probability of belonging” (James et al., 2013; Studer, 2019). A summary of the results for the LDA method can be found in Table 2. LDA results in an accuracy of 46.35% for the full model and 58.00 % accuracy for the PCA-reduced set of input variables when predicting *Home Win/ Away Win/ Draw*. When predicting the dummy variable *Home Win/ No Home Win* the full model achieves an accuracy of 57.31% and an accuracy of 65.35% for the reduced model. For the dummy variable *Away Win/ No Away Win* the full model’s accuracy is 68.32% while the accuracy for the reduced model is again significantly higher with 76.48%. The ROI maximising cutoff value for the *Home Win/ Away Win/ Draw* prediction is 0.9 and leads to an ROI of 29.88% while for the *Home Win/ No Home Win* dummy the cutoff is 0.9 as well which leads to an ROI of 23.60%. For the *Away Win/ No Away Win* prediction the cutoff is 0.68 with an ROI of 5.70%. The betting volume is distinct higher for the *Away Win/ No Away Win* prediction compared to the prediction of the other two variables. However, considering the graphs, the results seem to be stable for the *Home Win/ Away Win/ Draw* and the *Home Win/ No Home Win* target variables while for *Away Win/ No Away Win* in some cases losses are realised even after the optimal cutoff value. This indicates less stable predictions and higher variance.

**Table 2: Summary Results Training Data Set for LDA**

	Home Win/ Away Win/ Draw	Home Win/ No Home Win	Away Win/ No Away Win
<b>Accuracy Full Model (%)</b>	46.35	57.31	68.32
<b>Accuracy PCA-Reduced Model (%)</b>	58.00	65.54	76.48
<b>Graph Best Performing Model (Full or PCA-Reduced)</b>			
<b>Max Cutoff (%)</b>	90	90	68
<b>Max ROI (%)</b>	29.88	23.60	5.7

The second machine learning technique applied to the problem at hand is logistic regression. “Logistic regression is a linear modeling technique based on maximum likelihood specifically designed for classification and probability predictions” (James et al., 2013; Studer, 2019). This machine learning technique is particularly suitable as all classifications are based on the probability predictions of the method. The test data set results for logistic regression are displayed in Table 3. One can notice that also for the logistic regression method, higher accuracies can be achieved with the PCA-reduced input layer. Overall, the logistic regression performs similar but slightly worse than LDA except for the target variable *Home Win/ No Home Win* where the accuracy score could be slightly increased to 66.04%. Also in terms of betting volume, ROI and result stability for the optimal strategies, the obtained outcomes are highly similar to the results of LDA.

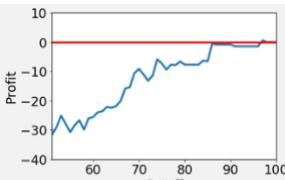
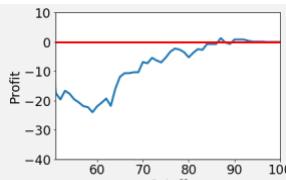
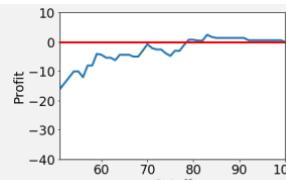
**Table 3: Summary Results Training Data Set for Logistic Regression**

	Home Win/ Away Win/ Draw	Home Win/ No Home Win	Away Win/ No Away Win
<b>Accuracy Full Model (%)</b>	52.44	61.42	73.95
<b>Accuracy PCA- Reduced Model (%)</b>	57.03	66.04	75.03
<b>Graph Best Performing Model (Full or PCA- Reduced)</b>			
<b>Max Cutoff (%)</b>	84	90	64
<b>Max ROI (%)</b>	31.10	23.60	7.48
<b>Betting Volume (%)</b>	5	2.5	15.5

As a third machine learning prediction method, random forest will be tested. “Random forest is a machine learning technique suitable for probability prediction that reduces the overfitting tendency of decision trees by building a multitude of trees with different subsets of predictors” (Studer, 2019). Random forest is considered a promising approach for the task at hand due to the superior handling of unbalanced target variable distributions, insignificant input variables as well as lower sample sizes (James et al., 2013). However, the results in Table 4 show that random forest has a strictly lower accuracy compared to LDA and logistic regression. Nevertheless, when it comes to the optimal investment strategy, the model performs well and

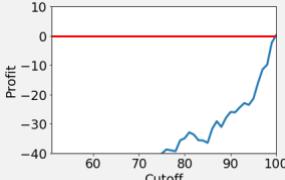
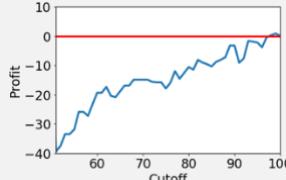
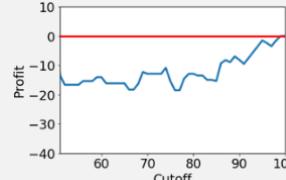
reliable for the variables *Home Win/ No Home Win* and *Away Win/ No Away Win*. However, the cutoff values for both strategies are high. The flip side of the high cutoff values and the high result stabilities are the low betting volumes.

**Table 4: Summary Results Training Data Set for Random Forest**

	<b>Home Win/ Away Win/ Draw</b>	<b>Home Win/ No Home Win</b>	<b>Away Win/ No Away Win</b>
<b>Accuracy Full Model (%)</b>	55.34	64.50	67.06
<b>Accuracy PCA- Reduced Model (%)</b>	57.55	58.96	71.99
<b>Graph Best Performing Model</b>			
<b>Max Cutoff (%)</b>	96	86	82
<b>Max ROI (%)</b>	56.99	16.25	49.99
<b>Betting Volume (%)</b>	0.5	8	2.5

Lastly, the gradient boosting algorithm will be applied. “Gradient boosting step-wise builds a decision tree to minimalise a defined loss function in order to correctly classify/assign probabilities to different classes. It can be particularly suitable for unbalanced data sets” (Studer, 2019; James et al., 2013). However, as one can see in Table 5, also the gradient boosting accuracies are significantly lower compared to the results obtained with LDA and logistic regression.

**Table 5: Summary Results Training Data Set for Gradient Boosting**

	<b>Home Win/ Away Win/ Draw</b>	<b>Home Win/ No Home Win</b>	<b>Away Win/ No Away Win</b>
<b>Accuracy Full Model (%)</b>	49.95	58.98	71.12
<b>Accuracy PCA- Reduced Model (%)</b>	53.15	56.51	73.12
<b>Graph Best Performing Model (Full or PCA- Reduced)</b>			
<b>Max Cutoff (%)</b>	99	98	-
<b>Max ROI (%)</b>	20.50	29.30	-

<b>Betting Volume (%)</b>	1	1.5	-
---------------------------	---	-----	---

Overall, the results show that LDA and logistic regression lead to the highest accuracy scores. However, this does not mean these techniques are automatically the best foundations for the most rewarding betting strategy. The main criteria to define the best betting strategy are the total winnings, the betting amount multiplied with the ROI, of the profit optimised cutoff value. Moreover, it is important to check the graph to see whether the results above the given cutoff threshold are strictly positive and therefore stable. To get more concrete, it can be seen that for the prediction of *Home Win/ Away Win/ Draw*, logistic regression with a probability threshold equal or above 0.84 leads to the highest winnings. Additionally, the strategy is also stable as the winnings never turn negative above this threshold. Therefore, *Home Win/ Away Win/ Draw* will be predicted with logistic regression and the corresponding 0.84 cutoff value on the validation data set. For the target variable *Home Win/ No Home Win* random forest with a probability threshold of 0.86 and above promises the most profitable and stable strategy. Lastly, *Away Win/ No Away Win* will be predicted with random forest as well. There, the threshold of 0.82 and above will be used. As for all three target variables a profitable and stable betting strategy could be found, the final out of sample strategy will combine all three different approaches.

## 4.2 Validation Data

The validation results for the three different strategies as well as for the overall strategy are displayed in Table 6. There are three major findings. First of all, each sub-strategy results in a positive ROI. Secondly, by combining the three strategies, the achieved ROI of 38.38% is highly promising. Thirdly, the betting volume for the overall strategy is with 8%, out of all possible games, fairly high and therefore indicates stable and reliable results.

**Table 6: Results Validation Data**

	Machine Learning Technique	Cutoff Threshold (%)	Betting Volume (%)	ROI (%)
<b>Home Win/ Away Win/ Draw</b>	Logistic Regression	0.84	3	17.33
<b>Home Win/ No Home Win</b>	Random Forest	0.86	3	9.66
<b>Away Win/ No Away Win</b>	Random Forest	0.82	2	113.00
<b>Overall Strategy</b>	mixed	mixed	8	38.38

## 5 Discussion

As a first step, the obtained results should be put into perspective and hence be compared to other forecasting approaches. Overall, the obtained accuracy of 58% for *the home win/ draw/ away win* target variable is higher than or similar to the achieved accuracies by other researchers such as Hucaljuk & Rakipović (2011) or Tax et al. (2015). However, accuracy is highly dependent on the particular data set. Therefore, it makes sense to compare the results to benchmarks for the same data set. Thereby, the random benchmark of 49% was outperformed. With regard to betting companies, it has to be acknowledged that the bookmaker Bet365 has with an accuracy of 59.5% still a slightly higher predictive power for the *home win/ draw/ away win* target variable compared to the best performing model of this paper. However, this is not surprising. The bookmakers have more data available with regard to features per game, the time frame of past data as well as the experience and human capital knowledge to develop superior models. In addition to that, the betting companies offer odds which do not add up, the companies introduce a margin, in order to establish a further edge and guarantee their profitability. Having a look at the financial reports of Bet365 over recent years confirms the common belief that overall, the bookmaker always wins (Ahmed & Bounds, 2018).

However, there are two distinct advantages that allow bettors to create successful betting strategies despite the lower resources and the disfavorable odds. Contrary to the bookmakers, the bettor does not have to act in every match but can focus on events where the certainty of his model is high. This lowers the variance and allows focusing on events that the bettor's model understands best. The betting company does not have this advantage as it is expected to make the market and assign probabilities to every possible outcome. Secondly, the bettor can lower the complexity of his models and gain a different perspective by grouping outcomes as done in this paper with the target variables *home win/ no home win* and *away win/ no away win*. This again introduces flexibility the bookmakers do not have.

As indicated by the above results, these two advantages combined with modern statistical, optimisation and machine learning techniques as well as creative feature engineering still allows bettors to create a profitable, public data-based strategy in today's football betting market.

## **6 Conclusion**

### **6.1 Results**

The present paper modeled a machine learning-driven sports betting strategy for the English Premier League season 2018/19 by using team characteristics input features as well as time-series performance data in combination with dimensionality reduction and optimisation techniques. The final validated strategy achieved a return of investment of 38.38 % by betting on 8% of all possible games.

### **6.2 Limitations**

As every statistical model, also this study does not come without any limitations. The major drawback of this paper's model is the small sample size of the different data sets. Even though it's comparable to other scientific research and a lot of attention was paid to reducing variance with statistical methods such as cross-validation and rolling windows, the betting volume of the final strategy is still exposed to the variance of a few events in the validation data set. Therefore, to prove the absolute validity and sustainability of the strategy over multiple seasons and leagues, the model would need to be applied to many other leagues and seasons. However, this would go beyond the goal and scope of this paper as it would be connected to a significant increase in data sourcing and data manipulation work.

With regard to the implementation of the betting strategy, it is important to notice that there are betting companies which simply exclude strong winning players from their service (Economist, 2017). Consequently, the opportunity to apply the developed strategy on such platforms is limited. Hence, alternatives such as implementing the strategy on betting exchange sites where the bet is matched with a counterposition of another bettor have to be considered. An example of such a platform would be betfair.com.

### **6.3 Outlook**

There are three major suggestions on how the study can be enhanced to increase its value further. First of all, an optimisation function can not only be applied to evaluate the optimal probability cutoff in order to maximalise the wins but also to optimise the gap between betting odds and modeled probabilities. However, this would increase the model's and its computational

complexity significantly. Moreover, it would also premise superior prediction accuracy in at least some identifiable cases, which directly leads to the second point.

As the publicly available data and its forecasting power are limited, there is a lot of room to create additional value by enhancing the developed modeling approach with data that is not publicly available. In order to create an even larger edge in the betting market, detailed team performance data, as well as player individual match data from companies such as OPTA, could be considered. With that, modeling will get much more complex but the potential rewards in the betting market will increase as well, as shown by the recent examples of Tony Bloom who bought with the winnings of his quant-based sports betting company the Premier League football club Brighton & Hove Albion (Biermann, 2019).

Lastly, moving away from the betting market, there is a lot of value for the football coaching staff by better understanding the relationship between detailed team-based features as well as player individual data and game outcomes. Examples of these input features are ball circulation speed, passing accuracy, room covering, individual player movement or defender proximity to opponents. Understanding the relationship of these features with game outcomes can help to develop better game plans, squad planning and transfer policies. For all these topics, this paper's model can be insightful by having a more detailed look at the feature importance of the random forest and the composition of the different principal components. However, once more, there would be a need for alternative, private data sources in order to get these insights. The Danish football champion of the season 2017/18 and former underdog FC Midtjylland already showed the power of analytics-based transfer strategies and game plan design (Biermann, 2019). In the future, with the ever-increasing amount of data, the most effective navigation through the flood of newly available information and its subsequent implantation onto the football pitch and into the football ball clubs board rooms', will likely be decisive tools for success in the increasingly competitive football world.

## References

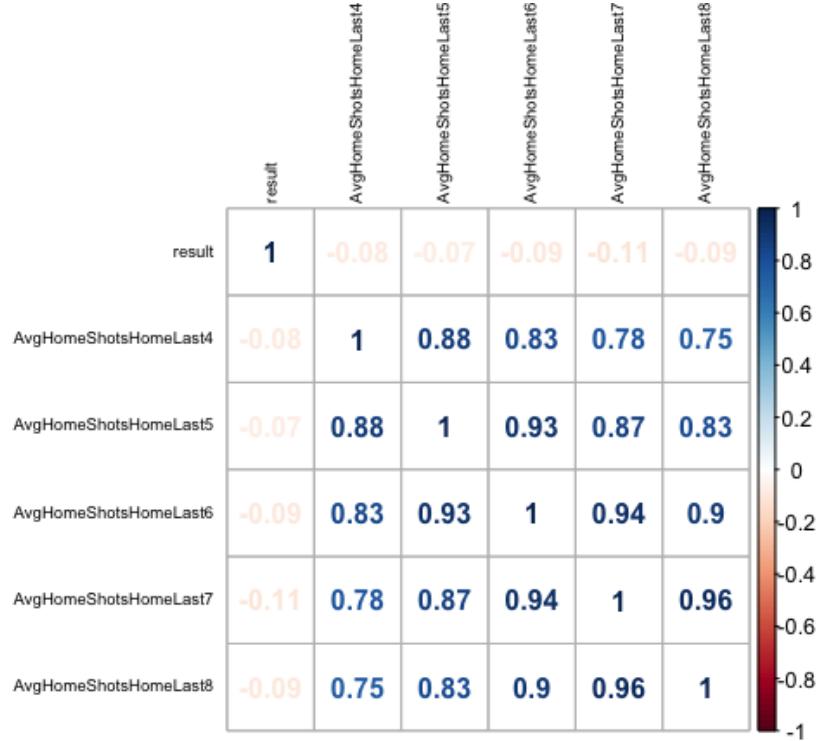
- Ahmed, M. and Bounds, A. (2018). Bet365 stands out among rivals for more than pay. *FT* [online]. Available at: <https://www.ft.com/content/f40e1b44-ee60-11e8-89c8-d36339d835c0> [Accessed 29 August 2019]
- Barnard, M., Boor, S., Winn, C., Wood, C., and Wray, I. (2019). World in motion – Annual review of football finance. *Deloitte* [online]. Available at: <https://www2.deloitte.com/content/dam/Deloitte/uk/Documents/sports-business-group/deloitte-uk-annual-review-of-football-finance-2019.pdf> [Accessed 19 August 2019]
- Biermann, C. (2019). *Football hackers: The science and art of data revolution*. Blink: London.
- Burn-Murdoch, J. (2018). How data analysis helps football clubs make better signings. *FT* [online]. Available at: [ft.com/content/84aa8b5e-c1a9-11e8-84cd-9e601db069b8](https://ft.com/content/84aa8b5e-c1a9-11e8-84cd-9e601db069b8) [Accessed 19 August 2019]
- Conn, D. (2018). How data analysis helps football clubs make better signings. *FT* [online]. Available at: [ft.com/content/84aa8b5e-c1a9-11e8-84cd-9e601db069b8](https://ft.com/content/84aa8b5e-c1a9-11e8-84cd-9e601db069b8) [Accessed 19 August 2019]
- 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), pp.37-50. doi: <https://doi.org/10.1515/jqas-2012-0036>
- 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*, 46(2), pp.265-280. doi: <https://doi.org/10.1111/1467-9876.00065>
- Economist (2017). *How bookmakers deal with winning customers* [online]. Available at: <https://www.economist.com/the-economist-explains/2017/10/04/how-bookmakers-deal-with-winning-customers> [Accessed 29 August 2019]
- Elo, A.E., 1978. *The rating of chessplayers, past and present*. Arco Pub.
- Fifindex.com (2019). Available at: <https://www.fifaindex.com/> [Accessed 10 August 2019]
- Football-data.co.uk (2019). *Data files: England*. [online]. Available at: <https://www.football-data.co.uk/englandm.php> [Accessed 10 August 2019]
- Goddard, J. (2005). Regression models for forecasting goals and match results in association football. *International Journal of forecasting*, 21(2), pp.331-340. doi: <https://doi.org/10.1016/j.ijforecast.2004.08.002>

- Hill, I.D. (1974). Association football and statistical inference. *Journal of the Royal Statistical Society, 23*(2), pp.203-208. doi: <https://doi.org/10.2307/2347001>
- Hucaljuk, J. and Rakipović, A. (2011). Predicting football scores using machine learning techniques [online]. In: *2011 Proceedings of the 34th International Convention MIPRO*. New Jersey: IEEE, pp. 1623-1627. Available at: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=5967321> [Accessed 19 August 2019]
- Impect (2019). *Impect – Making success in soccer measurable* [online]. Available at: <https://www.impect.com/en/> [Accessed 19 August 2019]
- James, G., Witten, D., Hastie, T. and Tibshirani, R. (2013). *An introduction to statistical learning*. New York: Springer.
- Kürüm, E., Yıldırak, K. and Weber, G.W. (2012). A classification problem of credit risk rating investigated and solved by optimisation of the ROC curve. *Central European Journal of Operations Research, 20*(3), pp.529-557. doi: <https://doi.org/10.1007/s10100-011-0224-5>
- Lasek, J. (2016). EURO 2016: Predictions using team rating systems. In: *MLSA* [online]. Dublin: PKDD/ECML. Available at: <https://pdfs.semanticscholar.org/378d/f9520dec6fc2792adf377b47baa5009b495c.pdf> [Accessed 19 August 2019]
- Lucey, P., Bialkowski, A., Monfort, M., Carr, P. and Matthews, I. (2014). Quality vs quantity: Improved shot prediction in soccer using strategic features from spatiotemporal data. In: *Proc. 8th annual MIT Sloan sports analytics conference* [online]. Cambridge: SSAC, pp. 1-9. Available at: <http://www.sloansportsconference.com/wp-content/uploads/2015/02/SSAC15-RP-Finalist-Quality-vs-Quantity.pdf> [Accessed 29 August 2019]
- Statista (2019a). *Total collection from sport betting at an international level for season 2014/15, by league (in billion euros)* [online]. Available at: <https://www.statista.com/statistics/620308/collection-from-betting-on-international-football-europe/> [Accessed 19 August 2019]
- Statista (2019b). *Average annual player salary in the English Premier League in 2018/19, by team (in million U.S dollars)* [online]. Available at: <https://www.statista.com/statistics/675303/average-epl-salary-by-team/> [Accessed 10 August 2019]
- Studer, C. (2019). *The 10% annual ROI in the peer to peer lending market in 2018? – A machine learning-driven investment strategy*. MSc. Imperial College London.
- Tax, N. and Joustra, Y., 2015. Predicting the Dutch football competition using public data: A machine learning approach. *Transactions on knowledge and data engineering, 10*(10), pp.1-13.

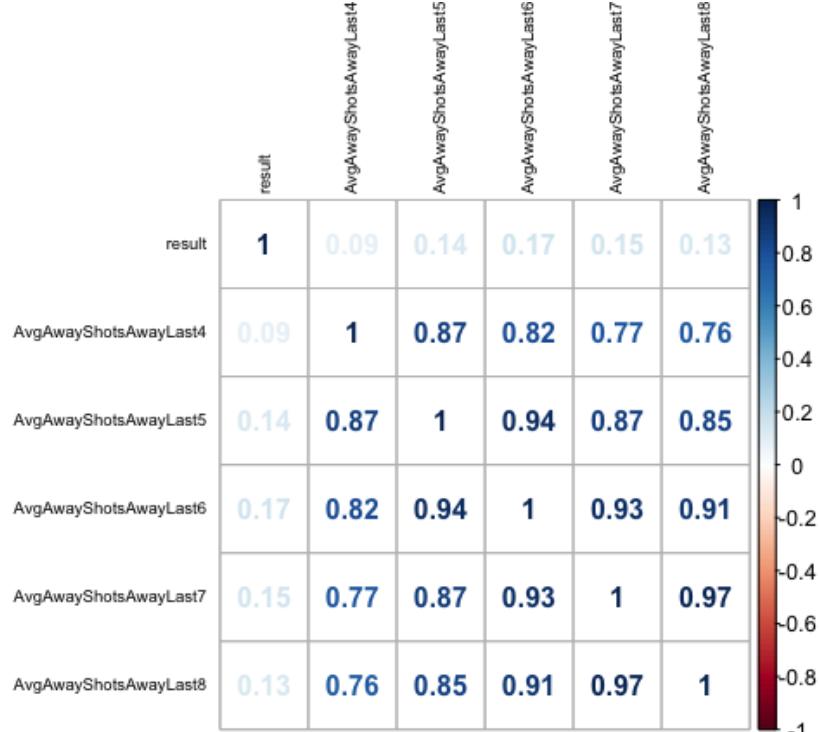
Transfermarkt.co.uk (2019). Clubs Premier League, 18/19 [*online*]. Available at:  
[https://www.transfermarkt.co.uk/premier-league/startseite/wettbewerb/GB1/saison\\_id/2018/plus/](https://www.transfermarkt.co.uk/premier-league/startseite/wettbewerb/GB1/saison_id/2018/plus/) [Accessed 10 August, 2019]

## Appendix

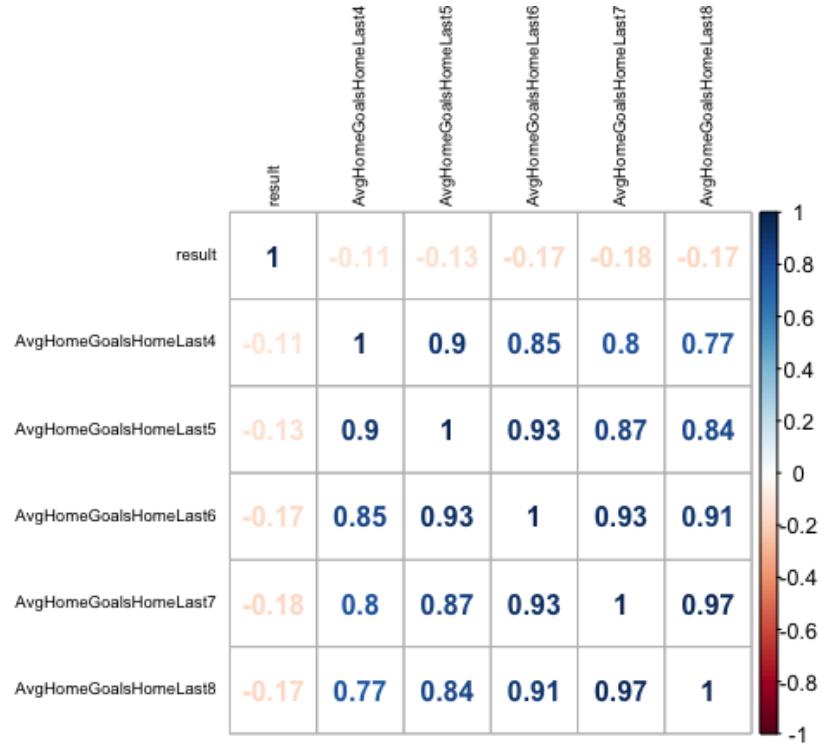
**Appendix 1: Correlation Table Result and Different Rolling Windows for Average Number of Shots Taken by the Home Team in Home Games**



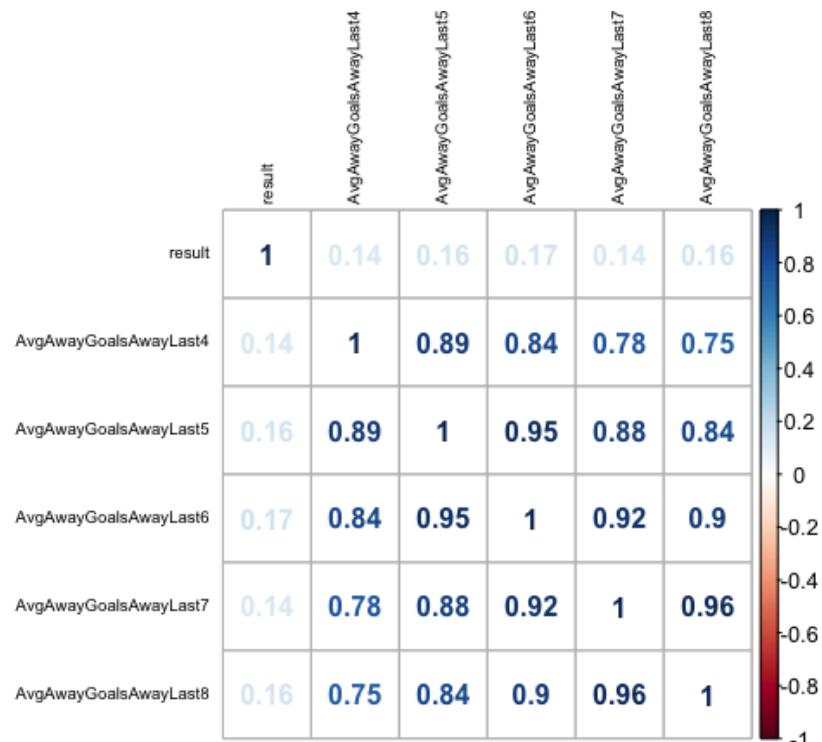
**Appendix 2: Correlation Table Result and Different Rolling Windows for Average Number of Shots Taken by the Away Team in Away Games**



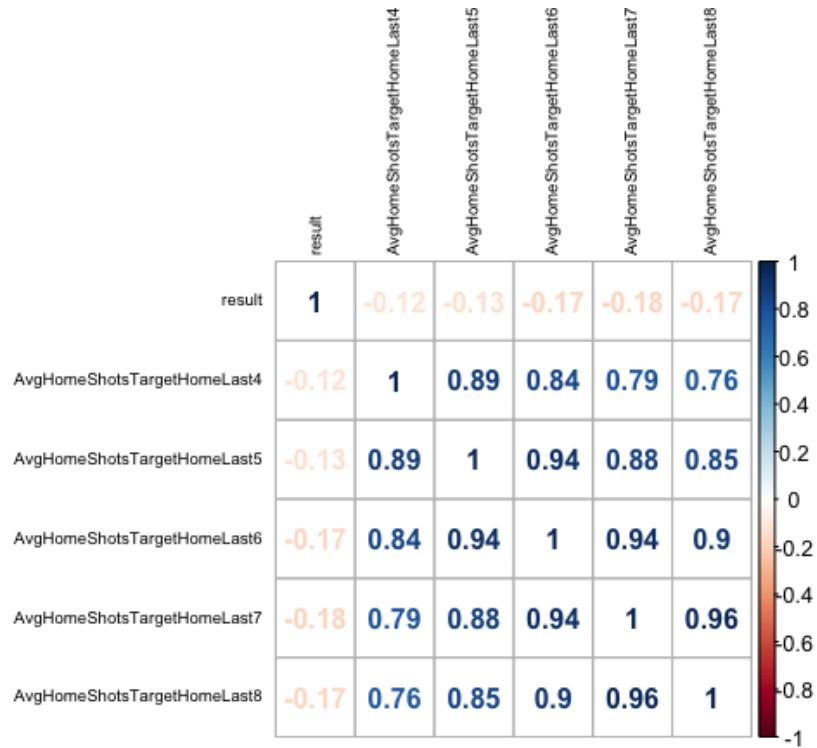
**Appendix 3: Correlation Table Result and Different Rolling Windows for Average Number of Goals Scored by the Home Team at Home**



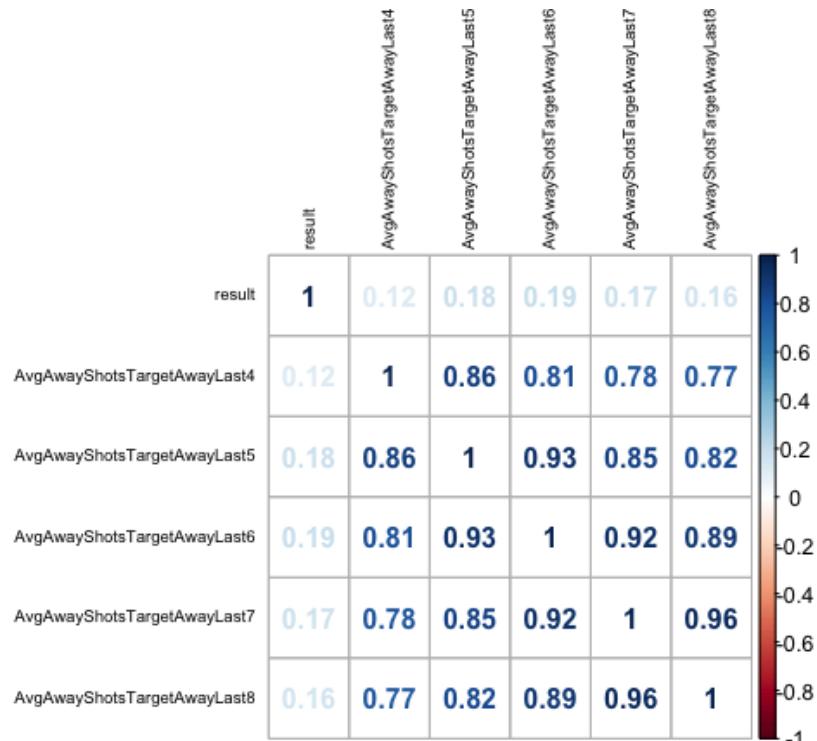
**Appendix 4: Correlation Table Result and Different Rolling Windows for Average Number of Goals Scored by the Away Team in Away Games**



**Appendix 5: Correlation Table Result and Different Rolling Windows for Average Number of Shots on Target by the Home Team at Home**



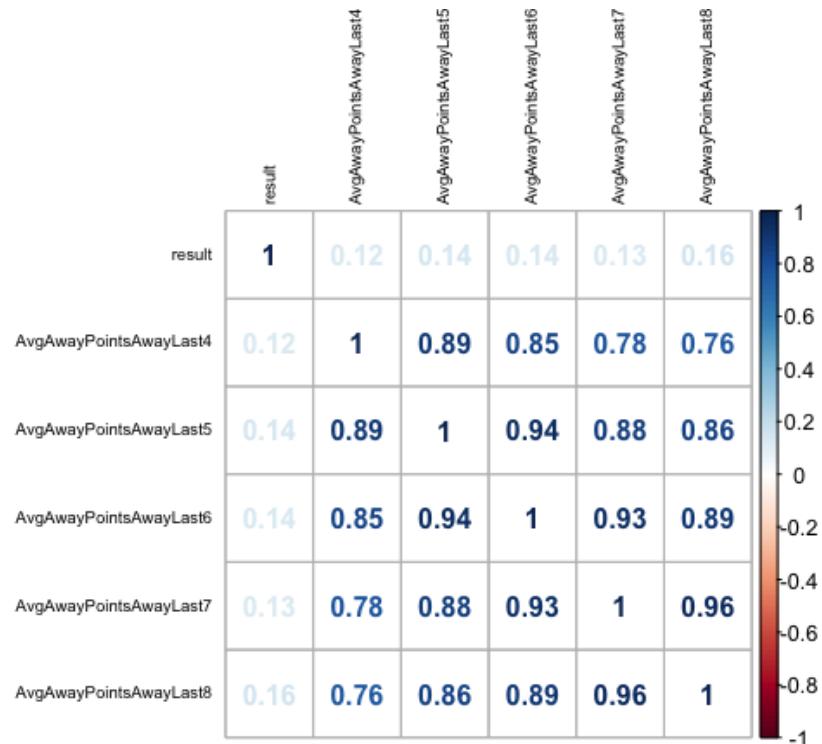
**Appendix 6: Correlation Table Result and Different Rolling Windows for Average Shots on Target by the Away Team in Away Games**



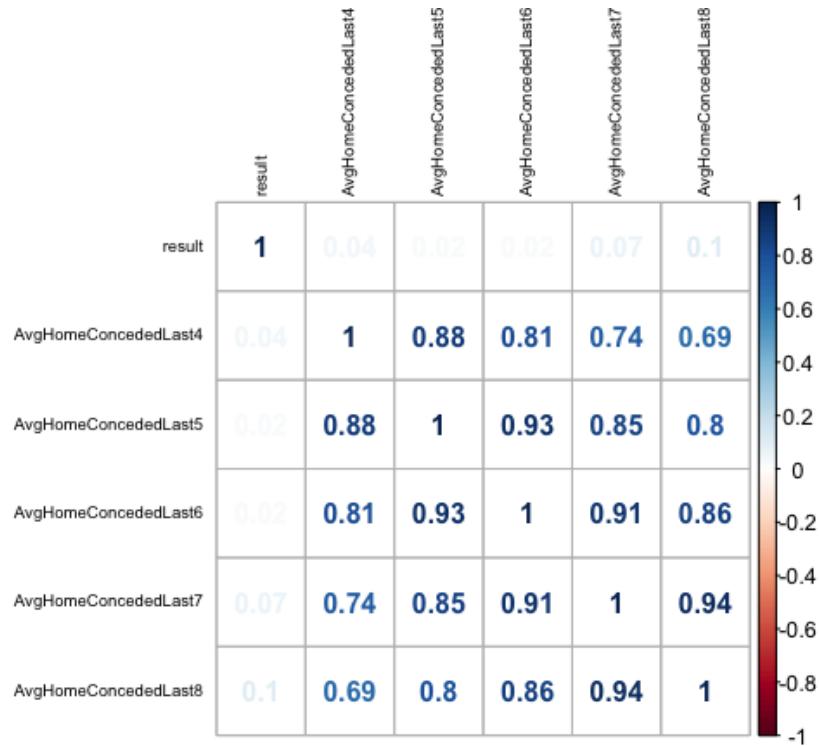
### Appendix 7: Correlation Table Result and Different Rolling Windows for Average Number of Points in by the Home Team at Home



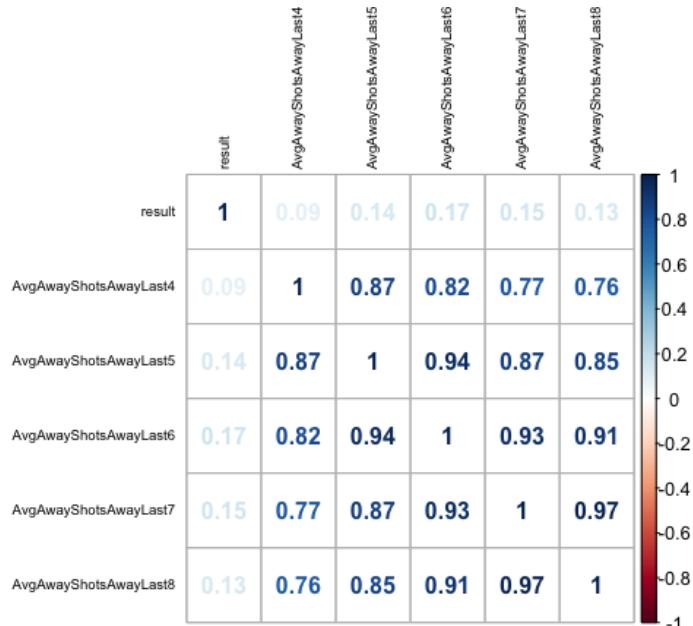
### Appendix 8: Correlation Table Result and Different Rolling Windows for Average Number of Away Points by the Away Team in Away Games



**Appendix 9: Correlation Table Result and Different Rolling Windows for Average Number of Conceded Goals by the Home Team**



**Appendix 10: Correlation Table Result and Different Rolling Windows for Average Number of Conceded Goals by the Away Team**



## Appendix 11

In [ ]:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import cross_val_predict
from sklearn.model_selection import cross_val_score
from sklearn.metrics import confusion_matrix
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn import linear_model
from sklearn.metrics import accuracy_score
from sklearn.decomposition import PCA
```

In [5]:

```
#load data
df = pd.read_csv('data_cleaned2.csv')
odds = pd.read_csv('odds365.csv')
```

In [6]:

```
#summary stats of data
df.describe()
```

Out[6]:

	stage	HomeGoal	AwayGoal	HomeShots	AwayShots	HomeShotsTarget
count	380.000000	380.000000	380.000000	380.000000	380.000000	380.000000
mean	19.507895	1.568421	1.252632	14.134211	11.144737	4.778947
std	10.989678	1.312836	1.180031	5.855371	4.654002	2.677686
min	1.000000	0.000000	0.000000	0.000000	2.000000	0.000000
25%	10.000000	1.000000	0.000000	10.000000	8.000000	3.000000
50%	19.500000	1.000000	1.000000	14.000000	11.000000	5.000000
75%	29.000000	2.000000	2.000000	18.000000	14.000000	6.000000
max	38.000000	6.000000	6.000000	36.000000	25.000000	14.000000

In [7]:

```
#create initialisation data set
ini=df.iloc[:80,:]
ini.describe()
```

Out[7]:

	stage	HomeGoal	AwayGoal	HomeShots	AwayShots	HomeShotsTarget
<b>count</b>	80.000000	80.000000	80.000000	80.000000	80.000000	80.000000
<b>mean</b>	4.500000	1.475000	1.337500	13.862500	11.437500	4.725000
<b>std</b>	2.305744	1.292383	1.157734	6.482193	4.431886	2.751179
<b>min</b>	1.000000	0.000000	0.000000	2.000000	3.000000	1.000000
<b>25%</b>	2.750000	0.000000	0.000000	9.000000	9.000000	3.000000
<b>50%</b>	4.500000	1.000000	1.000000	13.000000	11.000000	4.500000
<b>75%</b>	6.250000	2.000000	2.000000	18.000000	15.000000	6.000000
<b>max</b>	8.000000	6.000000	5.000000	32.000000	21.000000	14.000000

In [8]:

```
training=df.iloc[80:280,:]
training.describe()
```

Out[8]:

	stage	HomeGoal	AwayGoal	HomeShots	AwayShots	HomeShotsTarget
<b>count</b>	200.000000	200.000000	200.000000	200.00000	200.000000	200.000000
<b>mean</b>	18.510000	1.605000	1.230000	14.27000	10.710000	4.875000
<b>std</b>	5.797236	1.377849	1.172165	5.74737	4.237414	2.715926
<b>min</b>	9.000000	0.000000	0.000000	0.00000	2.000000	0.000000
<b>25%</b>	13.750000	1.000000	0.000000	10.00000	8.000000	3.000000
<b>50%</b>	18.500000	1.000000	1.000000	14.00000	10.000000	5.000000
<b>75%</b>	23.250000	2.000000	2.000000	17.00000	13.000000	7.000000
<b>max</b>	29.000000	6.000000	6.000000	36.00000	25.000000	12.000000

In [9]:

```
test=df.iloc[280:,:]
test.describe()
```

Out[9]:

	<b>stage</b>	<b>HomeGoal</b>	<b>AwayGoal</b>	<b>HomeShots</b>	<b>AwayShots</b>	<b>HomeShotsTarge</b>
<b>count</b>	100.000000	100.00000	100.000000	100.00000	100.000000	100.000000
<b>mean</b>	33.510000	1.57000	1.230000	14.08000	11.780000	4.630000
<b>std</b>	2.914592	1.19979	1.221483	5.58801	5.504048	2.557046
<b>min</b>	27.000000	0.00000	0.000000	2.00000	2.000000	0.000000
<b>25%</b>	31.000000	1.00000	0.000000	10.00000	7.000000	3.000000
<b>50%</b>	34.000000	1.00000	1.000000	14.00000	11.000000	5.000000
<b>75%</b>	36.000000	2.00000	2.000000	18.00000	16.000000	6.000000
<b>max</b>	38.000000	5.00000	5.000000	29.00000	25.000000	12.000000

## Pivot Tables

In [11]:

```
#create pivot tables for each of the match-based variables to prepare the outline of times series data based on rolling windows
HomeShots = df.pivot_table(index=['HomeTeam'],columns=['stage'], values =['HomeShots'], aggfunc='sum')
AwayShots = df.pivot_table(index=['AwayTeam'],columns=['stage'], values =['AwayShots'], aggfunc='sum')
HomePoints = df.pivot_table(index=['HomeTeam'],columns=['stage'], values =['HomePoints'], aggfunc='sum')
AwayPoints = df.pivot_table(index=['AwayTeam'],columns=['stage'], values =['AwayPoints'], aggfunc='sum')
HomeGoals = df.pivot_table(index=['HomeTeam'],columns=['stage'], values =['HomeGoals'], aggfunc='sum')
AwayGoals = df.pivot_table(index=['AwayTeam'],columns=['stage'], values =['AwayGoals'], aggfunc='sum')
HomeShotsTarget = df.pivot_table(index=['HomeTeam'],columns=['stage'], values =['HomeShotsTarget'], aggfunc='sum')
AwayShotsTarget = df.pivot_table(index=['AwayTeam'],columns=['stage'], values =['AwayShotsTarget'], aggfunc='sum')
HomeConceded = df.pivot_table(index=['HomeTeam'],columns=['stage'], values =['AwayGoal'], aggfunc='sum')
```

In [12]:

```
AwayConceded = df.pivot_table(index=[ 'AwayTeam' ],columns=[ 'stage' ], values =[ 'HomeGoal' ], aggfunc='sum')
AwayConceded
```

Out[12]:

	HomeGoal													
stage	1	2	3	4	5	6	7	8	9	10	...	29	30	31
<b>AwayTeam</b>														
<b>Arsenal</b>	NaN	3.0	NaN	2.0	1.0	NaN	NaN	1.0	NaN	2.0	...	1.0	NaN	3.0
<b>Bournemouth</b>	NaN	1.0	NaN	2.0	NaN	4.0	NaN	0.0	NaN	0.0	...	NaN	0.0	NaN
<b>Brighton</b>	2.0	NaN	1.0	NaN	2.0	NaN	2.0	NaN	0.0	NaN	...	NaN	1.0	1.0
<b>Burnley</b>	0.0	NaN	4.0	NaN	1.0	NaN	1.0	NaN	5.0	NaN	...	NaN	4.0	NaN
<b>Cardiff</b>	2.0	NaN	0.0	NaN	4.0	NaN	NaN	1.0	NaN	4.0	...	2.0	NaN	NaN
<b>Chelsea</b>	0.0	NaN	1.0	NaN	NaN	0.0	NaN	0.0	NaN	0.0	...	1.0	NaN	2.0
<b>Crystal Palace</b>	0.0	NaN	2.0	NaN	0.0	NaN	2.0	NaN	2.0	NaN	...	1.0	NaN	NaN
<b>Everton</b>	2.0	NaN	2.0	NaN	NaN	2.0	NaN	1.0	NaN	2.0	...	NaN	3.0	NaN
<b>Fulham</b>	NaN	3.0	NaN	2.0	3.0	NaN	3.0	NaN	4.0	NaN	...	NaN	3.0	NaN
<b>Huddersfield</b>	NaN	6.0	NaN	1.0	NaN	3.0	NaN	1.0	NaN	3.0	...	1.0	NaN	4.0
<b>Leicester</b>	2.0	NaN	1.0	NaN	4.0	NaN	0.0	NaN	3.0	NaN	...	2.0	NaN	1.0
<b>Liverpool</b>	NaN	0.0	NaN	1.0	1.0	NaN	1.0	NaN	0.0	NaN	...	0.0	NaN	1.0
<b>Man City</b>	0.0	NaN	1.0	NaN	NaN	0.0	NaN	0.0	NaN	0.0	...	0.0	NaN	0.0
<b>Man United</b>	NaN	3.0	NaN	0.0	1.0	NaN	3.0	NaN	2.0	NaN	...	NaN	2.0	NaN
<b>Newcastle</b>	NaN	0.0	NaN	2.0	NaN	0.0	NaN	3.0	NaN	0.0	...	2.0	NaN	2.0
<b>Southampton</b>	NaN	2.0	NaN	0.0	NaN	3.0	2.0	NaN	0.0	NaN	...	3.0	NaN	1.0
<b>Tottenham</b>	1.0	NaN	0.0	2.0	NaN	1.0	0.0	NaN	0.0	NaN	...	NaN	2.0	NaN
<b>Watford</b>	NaN	1.0	NaN	NaN	NaN	1.0	2.0	NaN	0.0	NaN	...	NaN	3.0	NaN
<b>West Ham</b>	4.0	NaN	3.0	NaN	1.0	NaN	NaN	1.0	NaN	1.0	...	NaN	2.0	NaN
<b>Wolves</b>	NaN	2.0	NaN	0.0	NaN	1.0	NaN	0.0	NaN	1.0	...	NaN	1.0	NaN

20 rows × 38 columns

## data base Loading

## Prepare Data with Time Series

In [13]:

```
#import only relevant data
df2 = pd.read_csv('data_cleaned3.csv')
```

In [14]:

```
df2.head()
```

Out[14]:

	homewin	draw	awaywin	stage	Date	HomeTeam	AwayTeam	FTHG	FTAG
0	3	0	0	1	01/01/00	Man United	Leicester	2	1
1	3	0	0	1	11/08/18	Bournemouth	Cardiff	2	0
2	0	0	3	1	11/08/18	Fulham	Crystal Palace	0	2
3	0	0	3	1	11/08/18	Huddersfield	Chelsea	0	3
4	0	0	3	1	11/08/18	Newcastle	Tottenham	1	2

In [15]:

```
#percentage of home wins in initialisation data set
df_i=df2.iloc[:80,:]
homewins_i=np.count_nonzero(df_i['homewin'] == 3)/(len(df_i))
homewins_i
```

Out[15]:

0.425

In [16]:

```
#percentage of away wins in initialisation data set
awaywins_i=np.count_nonzero(df_i['awaywin'] == 3)/(len(df_i))
awaywins_i
```

Out[16]:

0.375

In [17]:

```
#percentage of draws in initialisation data set
draws_i=1-homewins_i-awaywins_i
draws_i
```

Out[17]:

0.1999999999999996

In [18]:

```
#percentage of home wins over all data sets
homewins=np.count_nonzero(df2['homewin'] == 3)/(len(df2))
homewins
```

Out[18]:

0.4763157894736842

In [19]:

```
#percentage of draws over all data sets
draws=np.count_nonzero(df2['draw'] == 1)/(len(df2))
draws
```

Out[19]:

0.1868421052631579

In [20]:

```
#percentage of away wins over all data sets
awaywin=1-homewins-draws
awaywin
```

Out[20]:

0.33684210526315783

In [21]:

```
#creation of training set and validation data set
df3 = df2.iloc[80:,]
```

In [22]:

df3.head(11)

Out[22]:

	homewin	draw	awaywin	stage	Date	HomeTeam	AwayTeam	FTHG	FTA
80	0	1	0	9	20/10/18	Bournemouth	Southampton	0	0
81	3	0	0	9	20/10/18	Cardiff	Fulham	4	2
82	0	1	0	9	20/10/18	Chelsea	Man United	2	2
83	0	0	3	9	20/10/18	Huddersfield	Liverpool	0	1
84	3	0	0	9	20/10/18	Man City	Burnley	5	0
85	0	0	3	9	20/10/18	Newcastle	Brighton	0	1
86	0	0	3	9	20/10/18	West Ham	Tottenham	0	1
87	0	0	3	9	20/10/18	Wolves	Watford	0	2
88	3	0	0	9	21/10/18	Everton	Crystal Palace	2	0
89	3	0	0	9	22/10/18	Arsenal	Leicester	3	1
90	3	0	0	10	27/10/18	Brighton	Wolves	1	0

In [23]:

```
# function to append all rolling window times series data to the training and validation data set
def new_col2(lag,table, xx, colname, df):
    a=df3['stage'].values.tolist()
    b=df3[xx].values.tolist()
    new=[ ]
    for e1, e2 in zip(a,b):
        c=str(e2)
        d=int(e1)
        new.append(round(table.loc[c][d-lag-1:d-1].sum()/np.count_nonzero(~np.isnan(table.loc[c][d-lag-1:d-1])),2))
    df[colname] = new
```

In [24]:

df4=df2.iloc[80:, ]

In [25]:

df4.head(10)

Out[25]:

	homewin	draw	awaywin	stage	Date	HomeTeam	AwayTeam	FTHG	FTA
80	0	1	0	9	20/10/18	Bournemouth	Southampton	0	0
81	3	0	0	9	20/10/18	Cardiff	Fulham	4	2
82	0	1	0	9	20/10/18	Chelsea	Man United	2	2
83	0	0	3	9	20/10/18	Huddersfield	Liverpool	0	1
84	3	0	0	9	20/10/18	Man City	Burnley	5	0
85	0	0	3	9	20/10/18	Newcastle	Brighton	0	1
86	0	0	3	9	20/10/18	West Ham	Tottenham	0	1
87	0	0	3	9	20/10/18	Wolves	Watford	0	2
88	3	0	0	9	21/10/18	Everton	Crystal Palace	2	0
89	3	0	0	9	22/10/18	Arsenal	Leicester	3	1

In [26]:

```
# data for all team characteristics based features
clubs = {
    "Man United" : [8.6,769.5,84,84,81,83],
    "Man City" : [7.89,953.55,85,87,83,85],
    "Chelsea" : [6.6,787.28,84,85,82,83],
    "Liverpool" : [6.4,801.9,85,82,81,83],
    "Arsenal" : [6.39,498.6,85,82,81,82],
    "Tottenham" : [4.63,751.05,89,83,82,83],
    "Everton" : [4.28,389.70,80,79,78,79],
    "West Ham" : [4.2,257.85,79,78,76,79],
    "Crystal Palace" : [3.63,193.79,79,76,75,77],
    "Leicester" : [3.57,300.15,79,77,77,78],
    "Southampton" : [3.51,239.04,76,77,77,77],
    "Bournemouth" : [2.62,147.83,77,75,76,76],
    "Watford" : [2.34,147.78,77,78,76,77],
    "Wolves" : [2.31,177.93,75,80,75,76],
    "Fulham" : [2.27,223.43,76,77,74,76],
    "Brighton" : [2.23,161.73,76,76,75,76],
    "Newcastle" : [2.17,159.03,76,76,76,76],
    "Burnley" : [2.11,169.43,77,77,78,78],
    "Huddersfield" : [1.63,104.49,75,75,75,75],
    "Cardiff" : [1.26,85.64,73,74,73,73],
}
```

In [27]:

```
#append team characteristic-based features to training and validation data set
a=df3['HomeTeam'].values.tolist()
b=df3['AwayTeam'].values.tolist()
HomeSalary=[ ]
HomeValue=[ ]
HomeStrength=[ ]
AwaySalary=[ ]
AwayValue=[ ]
AwayStrength=[ ]
for e1, e2 in zip(a,b):
    HomeSalary.append(clubs[e1][0])
    HomeValue.append(clubs[e1][1])
    HomeStrength.append(clubs[e1][5])
    AwaySalary.append(clubs[e2][0])
    AwayValue.append(clubs[e2][1])
    AwayStrength.append(clubs[e2][5])
df4[ 'HomeAvgSalary' ]=HomeSalary
df4[ 'HomeValue' ]=HomeValue
df4[ 'HomeStrength' ]=HomeStrength
df4[ 'AwayAvgSalary' ]=AwaySalary
df4[ 'AwayValue' ]=AwayValue
df4[ 'AwayStrength' ]=AwayStrength
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:16: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
app.launch_new_instance()
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:17: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:18: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:19: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:20: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:21: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

In [28]:

```
#append average number of shots by the home team in home games to the training and validation data set for the rolling windows of the last 4-8 matches
for e in range(4,9):
    a='AvgHomeShotsHomeLast'+str(e)
    new_col2(e,HomeShots, 'HomeTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:
```

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [29]:

```
#append average number of shots by the away team in away games to the training a  
nd validation data set for the rolling windows of the last 4-8 matches  
for e in range(4,9):  
    a='AvgAwayShotsAwayLast'+str(e)  
    new_col2(e,AwayShots, 'AwayTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [30]:

```
#append average number of points by the home team in home games to the training  
and validation data set for the rolling windows of the last 4-8 matches  
for e in range(4,9):  
    a='AvgHomePointsHomeLast'+str(e)  
    new_col2(e,HomePoints, 'HomeTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [31]:

```
#append average number of points by the away team in away games to the training  
and validation data set for the rolling windows of the last 4-8 matches  
for e in range(4,9):  
    a='AvgAwayPointsAwayLast'+str(e)  
    new_col2(e,AwayPoints, 'AwayTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [32]:

```
#append average number of goals by the home team in home games to the training a  
nd validation data set for the rolling windows of the last 4-8 matches  
for e in range(4,9):  
    a='AvgHomeGoalsHomeLast'+str(e)  
    new_col2(e,HomeGoals, 'HomeTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [33]:

```
#append average number of goals by the away team in away games to the training a  
nd validation data set for the rolling windows of the last 4-8 matches  
for e in range(4,9):  
    a='AvgAwayGoalsAwayLast'+str(e)  
    new_col2(e,AwayGoals, 'AwayTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [34]:

```
#append average number of shots on target by the home team in home games to the  
training and validation data set for the rolling windows of the last 4-8 matches  
for e in range(4,9):  
    a='AvgHomeShotsTargetHomeLast'+str(e)  
    new_col2(e,HomeShotsTarget, 'HomeTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [35]:

```
#append average number of shots on target by the away team in away games to the
training and validation data set for the rolling windows of the last 4-8 matches
for e in range(4,9):
    a='AvgAwayShotsTargetAwayLast'+str(e)
    new_col2(e,AwayShotsTarget, 'AwayTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [36]:

```
#append average number of conceded by the home team to the training and validation
data set for the rolling windows of the last 4-8 matches
for e in range(4,9):
    a='AvgHomeConcededLast'+str(e)
    new_col2(e,HomeConceded, 'HomeTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [37]:

```
#append average number of conceded by the away team to the training and validation
data set for the rolling windows of the last 4-8 matches
for e in range(4,9):
    a='AvgAwayConcededLast'+str(e)
    new_col2(e,AwayConceded, 'AwayTeam',a, df4 )
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:9: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>

```
if __name__ == '__main__':
```

In [38]:

```
df4.shape
```

Out[38]:

```
(300, 67)
```

In [39]:

```
#add variable for home win / draw /away win to the data set
l1=df4[ 'FTHG' ].values.tolist()
l2=df4[ 'FTAG' ].values.tolist()
new=[ ]
for e1, e2 in zip (l1,l2):
    if e1>e2:
        new.append(1)
    elif e2>e1:
        new.append(2)
    else:
        new.append(0)
df4[ 'outcome' ]= new
```

```
/anaconda3/lib/python3.6/site-packages/ipykernel_launcher.py:11: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
# This is added back by InteractiveShellApp.init_path()
```

## x and y for 3 target variables

**! encoded training data sets are accidentally called test and test is called vali later on!**

In [46]:

```
#set up y variable
y = df4.outcome
```

In [47]:

```
#set up x variables
x= df4.drop(['homewin', 'draw', 'awaywin', 'stage', 'Date', 'HomeTeam', 'AwayTeam', 'FTHG', 'FTAG', 'outcome', 'HomePoints', 'AwayPoints'], axis=1)
```

In [42]:

```
#x2=x
```

In [43]:

```
'''ll=[]
for e in range (0,300):
    if e<200:
        ll.append('training')
    else:
        ll.append('test')
x2['set']=ll
x2['result']=y'''
```

In [44]:

```
#x2.to_csv('x2.csv', index=False)
```

In [48]:

```
#spliting into target variable for training data set(Y_fs), target variable for
the validation data set(y_test), input layer for training data set(x_fs), input
layer for validation data set(x-test)

y_fs = y.iloc[:200]
y_test = y.iloc[200:]
x_fs = x.iloc[:200]
x_test = x.iloc[200:]

#setting up PCA-reduced data set
pca = PCA(n_components=8)
x_fs_pca=pca.fit_transform(x_fs)
x_test_pca=pca.fit_transform(x_test)

x_fs.head()
```

Out[48]:

	<b>HomeAvgSalary</b>	<b>HomeValue</b>	<b>HomeStrength</b>	<b>AwayAvgSalary</b>	<b>AwayValue</b>	<b>AwayStr</b>
<b>80</b>	2.62	147.83	76	3.51	239.04	77
<b>81</b>	1.26	85.64	73	2.27	223.43	76
<b>82</b>	6.60	787.28	83	8.60	769.50	83
<b>83</b>	1.63	104.49	75	6.40	801.90	83
<b>84</b>	7.89	953.55	85	2.11	169.43	78

5 rows × 56 columns

In [49]:

```
pca.components_
```

Out[49]:

```
array([[ 4.65930877e-03,   6.69513053e-01,   7.85485805e-03,
       -5.40845281e-03,  -7.42471950e-01,  -8.35740377e-03,
       5.68064536e-03,   4.80736896e-03,   5.02224255e-03,
      4.18932387e-03,   4.17653998e-03,  -5.14371896e-03,
     -4.72863107e-03,  -5.05843482e-03,  -5.47577940e-03,
     -5.20311316e-03,   1.32998727e-03,   1.31296506e-03,
      1.40294685e-03,   1.36934516e-03,   1.34499420e-03,
     -6.83304153e-04,  -8.60745311e-04,  -8.37817746e-04,
     -8.27375248e-04,  -7.52303319e-04,   1.19339279e-03,
     1.01547776e-03,   1.08121963e-03,   1.04743977e-03,
     1.09655641e-03,  -2.08519520e-05,  -3.92210192e-05,
    -7.86237779e-06,  -5.37888529e-05,   1.86934566e-05,
     2.13620127e-03,   1.98673976e-03,   2.06390569e-03,
     1.82896481e-03,   2.01437520e-03,  -1.48518097e-03,
    -1.42595900e-03,  -1.23832953e-03,  -1.59518520e-03,
    -1.60426925e-03,  -1.08640376e-03,  -1.11022569e-03,
    -1.10342837e-03,  -1.09023301e-03,  -1.08156664e-03,
     8.26965143e-04,   7.65133242e-04,   7.68041369e-04,
     7.03057762e-04,   6.85530039e-04],
 [ 5.21671290e-03,   7.42446876e-01,   7.99402582e-03,
     4.52059306e-03,   6.69499621e-01,   8.09733224e-03,
     4.82227626e-03,   5.06558279e-03,   4.81245440e-03,
     5.59130082e-03,   5.67736884e-03,   6.32729349e-03,
     5.87259991e-03,   5.91918880e-03,   5.01147649e-03,
     4.70299753e-03,   1.44806387e-03,   1.10755992e-03,
     1.10094962e-03,   8.80562664e-04,   9.84722492e-04,
     6.93670351e-04,   6.67127594e-04,   8.69261571e-04,
     9.08853679e-04,   6.80008044e-04,   1.17709221e-03,
     1.08302851e-03,   1.06600712e-03,   1.00160905e-03,
     1.15673781e-03,   9.06047921e-05,   2.34060012e-04,
     2.98438936e-04,   2.44077556e-04,   1.48598852e-04,
     3.24599708e-03,   2.96031767e-03,   2.79185647e-03,
     2.85289747e-03,   2.78152169e-03,   1.10609258e-03,
     1.53346479e-03,   1.53133795e-03,   1.35806132e-03,
     1.24769386e-03,  -1.09688224e-03,  -7.09590681e-04,
    -7.37759823e-04,  -4.81885780e-04,  -5.16412841e-04,
    -6.70201955e-04,  -5.96739925e-04,  -8.67523725e-04,
    -8.49647991e-04,  -7.54634574e-04],
 [ 2.95808159e-02,   1.72485797e-02,   5.86811056e-03,
     7.34177563e-03,  -6.11868946e-04,   2.36935510e-03,
    -4.90099563e-01,  -4.62747164e-01,  -4.32917273e-01,
    -3.88839367e-01,  -3.51779494e-01,   6.39903337e-02,
    -5.56171039e-03,  -1.57726539e-02,   3.29470141e-03,
    -8.06539772e-03,  -2.64368532e-02,  -2.91722253e-02,
    -2.52763772e-02,  -2.43619912e-02,  -2.18114275e-02,
     1.29504228e-03,   8.68758803e-03,   6.78004441e-03,
     1.05710495e-02,   1.05794554e-02,  -2.69844900e-02,
    -2.82329946e-02,  -2.63168804e-02,  -2.14867558e-02,
    -2.26916386e-02,  -4.20669706e-03,   1.05671404e-04,
    -8.49675222e-04,   1.75711501e-03,   3.18095308e-03,
    -1.26506599e-01,  -1.27479840e-01,  -1.19849080e-01,
    -1.09797434e-01,  -1.03978667e-01,   2.02465737e-02,
     4.63019133e-03,  -1.24850141e-03,   1.29835216e-02,
     9.22515272e-03,   1.07780474e-02,   1.09248279e-02,
     8.45893918e-03,   9.73831991e-03,   7.11465660e-03,
     7.46348428e-03,  -9.90121133e-06,  -3.71368177e-03,
    -4.68824234e-03,  -6.34781467e-03],
 [-7.44538374e-03,  -2.10360126e-03,  -4.65756551e-03,
    -9.24568773e-03,  -1.67433592e-02,  -3.77573522e-02,
```

1.65376982e-02, 1.08386337e-02, -1.20756093e-03,  
 2.23326164e-02, 3.91742436e-02, 5.36538265e-01,  
 4.40041297e-01, 4.26039581e-01, 3.31215280e-01,  
 3.15649982e-01, 1.52979445e-02, -2.94558287e-04,  
 -2.14024834e-03, -6.19019896e-03, -2.35267741e-04,  
 5.37548026e-02, 4.23492798e-02, 4.43154159e-02,  
 3.18782746e-02, 3.04951699e-02, 2.81451339e-03,  
 -5.78503527e-03, -4.05480962e-03, -4.76527937e-03,  
 -3.55576190e-03, 5.13097508e-02, 4.29800303e-02,  
 4.03207358e-02, 3.19085249e-02, 3.19124219e-02,  
 -1.36428354e-02, -1.09922291e-02, -6.94614963e-03,  
 -4.86463459e-04, 4.53581604e-03, 1.82933689e-01,  
 1.48680604e-01, 1.46317938e-01, 1.23918601e-01,  
 1.15468747e-01, -8.55794908e-03, -2.19333474e-04,  
 -3.63085231e-03, 3.53069124e-03, 2.73700334e-03,  
 -2.59605096e-03, 3.19131754e-03, -1.76589097e-03,  
 6.94725077e-03, 4.30675074e-03],  
[ 2.33600013e-02, -4.04428446e-03, 4.53622054e-02,  
 8.30319762e-02, -2.09585830e-03, 7.19312877e-02,  
 -2.31638576e-01, -9.31773381e-02, -1.81399434e-02,  
 7.27506626e-02, 1.33823039e-01, -8.66926821e-02,  
 -2.62428995e-02, -1.00171690e-01, -9.28654790e-02,  
 -9.46642031e-02, 1.34801432e-01, 1.24017859e-01,  
 1.30095154e-01, 1.18481388e-01, 1.20286432e-01,  
 2.43506316e-01, 2.44455012e-01, 2.21967519e-01,  
 2.06762370e-01, 2.00865732e-01, 8.63328364e-02,  
 7.90614593e-02, 7.24350464e-02, 8.54149868e-02,  
 7.92486974e-02, 2.02428305e-01, 1.93146043e-01,  
 1.85137572e-01, 1.82041935e-01, 1.65304897e-01,  
 9.61673586e-02, 1.01831509e-01, 1.24355665e-01,  
 1.47760966e-01, 1.44531322e-01, 1.40903114e-01,  
 1.30898451e-01, 1.07086029e-01, 9.90886075e-02,  
 8.94438587e-02, -1.01640155e-01, -1.06141981e-01,  
 -9.89626258e-02, -7.82969655e-02, -7.30704749e-02,  
 -1.93869240e-01, -1.83923476e-01, -1.62744770e-01,  
 -1.46584980e-01, -1.46064030e-01],  
[-1.19144780e-02, 4.57524425e-03, -5.27713328e-02,  
 8.60823010e-03, -1.20773375e-04, -2.00294837e-02,  
 -2.17230392e-01, -8.39363681e-03, 1.88668166e-01,  
 2.38673605e-01, 3.53789401e-01, -1.19101265e-01,  
 -5.75140228e-02, -2.56839201e-02, 8.29403909e-03,  
 -1.52176786e-02, -2.02933647e-01, -1.72633292e-01,  
 -1.48569691e-01, -1.21494798e-01, -1.18973567e-01,  
 9.87567822e-02, 9.60003489e-02, 9.17880191e-02,  
 8.00041999e-02, 6.78252343e-02, -1.78300867e-01,  
 -1.36879297e-01, -1.20642161e-01, -1.02403416e-01,  
 -9.34136545e-02, 3.20929992e-02, 4.55178685e-02,  
 4.77261834e-02, 3.47417768e-02, 3.55563597e-02,  
 -3.77948626e-01, -2.88186509e-01, -2.33804079e-01,  
 -1.85492880e-01, -1.39536590e-01, 6.55071853e-02,  
 7.95634456e-02, 7.94294089e-02, 8.03807533e-02,  
 6.27367523e-02, 1.49441341e-01, 1.45152588e-01,  
 1.30544034e-01, 1.00938507e-01, 1.04695984e-01,  
 -7.45338414e-02, -6.23629923e-02, -5.69210140e-02,  
 -5.08975076e-02, -4.45941121e-02],  
[-2.53910521e-02, 1.82864367e-03, -5.38104292e-02,  
 6.26832762e-02, 6.04453241e-04, 1.40985124e-01,  
 3.87056683e-01, 3.40491175e-02, -3.88857848e-02,  
 -1.78287271e-01, -1.27934735e-01, 4.18038948e-01,  
 -1.87623848e-01, -2.61749168e-01, -3.00188205e-01,  
 -2.45353031e-01, -4.59467146e-02, -6.71975977e-02,

```

-8.67822488e-02, -6.74429190e-02, -5.42860873e-02,
2.58358609e-02, 2.56579711e-02, 1.40330191e-02,
5.35869182e-03, 2.97005672e-03, -4.99511171e-02,
-8.13446943e-02, -1.04707291e-01, -1.02226782e-01,
-8.34005592e-02, 1.04517855e-01, 8.78561292e-02,
8.06150458e-02, 5.41433746e-02, 5.04702792e-02,
1.08347626e-01, -2.64768889e-02, -3.95056368e-02,
-7.32732538e-02, -7.56970537e-02, 3.83180981e-01,
1.82566022e-01, 1.31796733e-01, 7.38615697e-02,
8.57430726e-02, 8.29556914e-02, 6.21057675e-02,
5.24554691e-02, 2.56144696e-02, 1.49894743e-02,
2.36444208e-02, 3.91165073e-03, 2.19794928e-02,
2.71866064e-02, 1.76116981e-02],
[ 1.86530603e-02, -2.99972945e-03, 2.14275002e-02,
7.71764004e-02, -8.46586416e-04, 6.87616452e-02,
-3.76966137e-01, -8.44507147e-03, 8.73266090e-02,
2.52426208e-01, 2.14926569e-01, 5.91839176e-01,
-3.44797778e-02, -1.17393746e-01, -1.59183370e-01,
-2.79525304e-01, 4.23311640e-02, 4.22128571e-02,
5.14206416e-02, 5.55042094e-02, 6.99792384e-02,
-1.41446039e-01, -7.80435249e-02, -7.64299453e-02,
-7.40040364e-02, -6.08346332e-02, 2.17650343e-03,
1.58434738e-02, 3.71762540e-02, 4.54621407e-02,
5.31453429e-02, -1.70977048e-01, -1.26445855e-01,
-1.18191050e-01, -1.13558987e-01, -9.39931936e-02,
-7.06224492e-02, -3.44985813e-02, -8.29956985e-04,
6.58213951e-02, 6.95692028e-02, -1.22449395e-02,
-1.40017840e-01, -1.30651634e-01, -1.31235117e-01,
-1.22387694e-01, -6.19217741e-02, -6.24799183e-02,
-6.33147514e-02, -7.04556754e-02, -6.58919977e-02,
-5.21566647e-02, -5.43559629e-02, -6.02802653e-02,
-3.22001792e-02, -3.91222028e-02]])

```

In [50]:

```
#percentage of home wins in training data set
homewins=np.count_nonzero(y_test == 1)/(len(y_test))
homewins
```

Out[50]:

0.49

In [51]:

```
#percentage of away wins in training data set
awaywins=np.count_nonzero(y_test == 2)/(len(y_test))
awaywins
```

Out[51]:

0.34

In [52]:

```
#percentage of home wins in validation data set
homewins=np.count_nonzero(y_fs == 1)/(len(y_fs))
homewins
```

Out[52]:

0.49

In [53]:

```
#percentage of away wins in validation data set
awaywins=np.count_nonzero(y_fs == 2)/(len(y_fs))
awaywins
```

Out[53]:

0.32

In [54]:

```
#percentage of draws in validation data set
draws=1-homewins-awaywins
draws
```

Out[54]:

0.19

**models target variable home win/ draw/ away win**

In [55]:

```
#prediction

# linear discriminate analysis full model
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis()
lda_p = cross_val_score(lda, x_fs, y_fs, cv=10)
print('lda: ' + str(lda_p.sum()/10))

# logit full model
from sklearn.linear_model import LogisticRegression
logit = LogisticRegression()
logit_p = cross_val_score(logit, x_fs, y_fs, cv=10)
print('logit: ' + str(logit_p.sum()/10))

#gba full model
from sklearn import ensemble
gbr = ensemble.GradientBoostingClassifier()
gbr_p = cross_val_score(gbr, x_fs, y_fs, cv=10)
print('gbr: ' + str(gbr_p.sum()/10))

#random forest full model
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier( max_depth=5,random_state=10)
rfc_p = cross_val_score(rfc, x_fs, y_fs, cv=10)
print('rfc: ' + str(rfc_p.sum()/10))
```

```
lda: 0.46349206349206346
logit: 0.5244444444444444
gbr: 0.503968253968254
rfc: 0.5534126984126984
```

In [56]:

```
# linear discriminate analysis pca-reduced model
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis()
lda_p= cross_val_score(lda, x_fs_pca, y_fs, cv=10)
print('lda_pca: ' + str(lda_p.sum()/10))

# logit full model pca-reduced model
from sklearn.linear_model import LogisticRegression
logit = LogisticRegression()
logit_p = cross_val_score(logit, x_fs_pca, y_fs, cv=10)
print('logit_pca: ' + str(logit_p.sum()/10))

#gba pca-reduced model
from sklearn import ensemble
gbr = ensemble.GradientBoostingClassifier()
gbr_p = cross_val_score(gbr, x_fs_pca, y_fs, cv=10)
print('gbr_pca: ' + str(gbr_p.sum()/10))

#random forest pca-reduced model
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier( max_depth=5,random_state=10)
rfc_p = cross_val_score(rfc, x_fs_pca, y_fs, cv=10)
print('rfc_pca: ' + str(rfc_p.sum()/10))
```

```
lda_pca: 0.5800793650793651
logit_pca: 0.5703174603174602
gbr_pca: 0.5367460317460317
rfc_pca: 0.5755555555555555
```

In [57]:

```
#feature ranking random forest
rfc2=rfc.fit(x_fs, y_fs)
importances = rfc2.feature_importances_
std = np.std([rfc2.feature_importances_ for tree in rfc2.estimators_],
             axis=0)
indices = np.argsort(importances)[::-1]

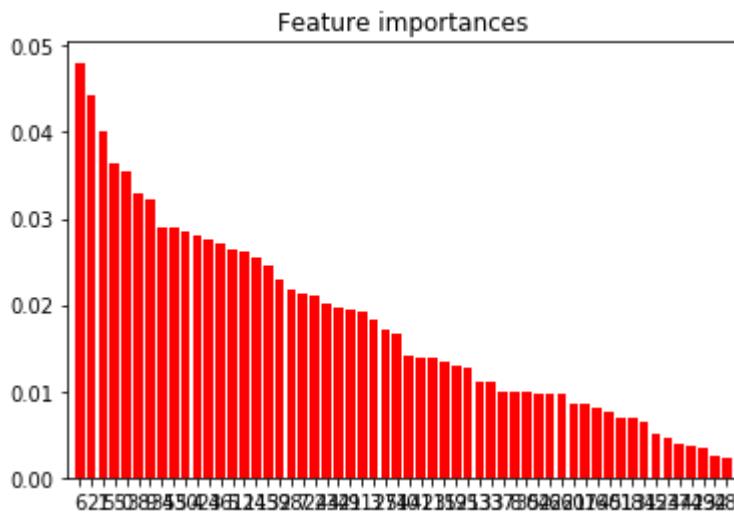
# Print the feature ranking
print("Feature ranking:")

for f in range(x_fs.shape[1]):
    print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))

# Plot the feature importances of the forest
plt.figure()
plt.title("Feature importances")
plt.bar(range(x_fs.shape[1]), importances[indices],
        color="r", yerr=std[indices], align="center")
plt.xticks(range(x_fs.shape[1]), indices)
plt.xlim([-1, x_fs.shape[1]])
plt.show()
```

**Feature ranking:**

1. feature 6 (0.048074)
2. feature 2 (0.044389)
3. feature 1 (0.040176)
4. feature 55 (0.036305)
5. feature 0 (0.035402)
6. feature 38 (0.032901)
7. feature 9 (0.032335)
8. feature 34 (0.028932)
9. feature 53 (0.028894)
10. feature 50 (0.028447)
11. feature 4 (0.028138)
12. feature 24 (0.027475)
13. feature 36 (0.027033)
14. feature 5 (0.026382)
15. feature 12 (0.026186)
16. feature 14 (0.025565)
17. feature 15 (0.024564)
18. feature 39 (0.023007)
19. feature 28 (0.021728)
20. feature 7 (0.021403)
21. feature 22 (0.021028)
22. feature 43 (0.020215)
23. feature 42 (0.019617)
24. feature 49 (0.019487)
25. feature 11 (0.019288)
26. feature 3 (0.018365)
27. feature 27 (0.017244)
28. feature 54 (0.016646)
29. feature 10 (0.014141)
30. feature 41 (0.013874)
31. feature 21 (0.013815)
32. feature 35 (0.013379)
33. feature 19 (0.012955)
34. feature 25 (0.012824)
35. feature 13 (0.011169)
36. feature 33 (0.011000)
37. feature 37 (0.009967)
38. feature 8 (0.009913)
39. feature 30 (0.009857)
40. feature 52 (0.009759)
41. feature 46 (0.009714)
42. feature 26 (0.009711)
43. feature 20 (0.008587)
44. feature 17 (0.008516)
45. feature 16 (0.007972)
46. feature 40 (0.007722)
47. feature 51 (0.006933)
48. feature 18 (0.006908)
49. feature 31 (0.006534)
50. feature 45 (0.005064)
51. feature 23 (0.004705)
52. feature 47 (0.003838)
53. feature 44 (0.003699)
54. feature 29 (0.003395)
55. feature 32 (0.002548)
56. feature 48 (0.002276)



In [58]:

```
importances
```

Out[58]:

```
array([0.03540232, 0.04017582, 0.04438912, 0.01836524, 0.02813773,
       0.02638201, 0.04807363, 0.02140286, 0.00991295, 0.03233549,
       0.01414105, 0.01928775, 0.02618611, 0.01116927, 0.02556468,
       0.02456381, 0.00797152, 0.00851589, 0.00690777, 0.01295529,
       0.008587 , 0.01381537, 0.02102805, 0.00470476, 0.02747482,
       0.01282379, 0.00971113, 0.01724406, 0.02172796, 0.00339495,
       0.0098573 , 0.0065338 , 0.00254804, 0.01099952, 0.02893193,
       0.01337883, 0.02703324, 0.00996749, 0.03290067, 0.02300673,
       0.00772189, 0.01387445, 0.01961673, 0.02021532, 0.00369891,
       0.00506439, 0.0097135 , 0.00383764, 0.00227567, 0.01948747,
       0.02844748, 0.00693292, 0.00975912, 0.02889357, 0.01664577,
       0.03630539])
```

## prepare betting odds of bet365

In [59]:

```
#test accuracy of their odds
prediction365=[]
for e in range(0,380):
    if odds.iloc[e][0]<odds.iloc[e][1] and odds.iloc[e][0]<odds.iloc[e][2]:
        prediction365.append(1)
    elif odds.iloc[e][2]<odds.iloc[e][1] and odds.iloc[e][2]<odds.iloc[e][0]:
        prediction365.append(2)
    else:
        prediction365.append(0)
prediction365_2=prediction365[80:280]
prediction365_3=prediction365[280:]
y_fs = y.iloc[:200]
print('Accuracy bet365: '+ str(accuracy_score(prediction365_2, y_fs)))
```

Accuracy bet365: 0.595

In [60]:

```
#predict probabilities for each match's outcome with all four different techniques
#lda with pca probability prediction
lda_pca_p= cross_val_predict(lda, x_fs_pca, y_fs, cv=10,method='predict_proba')
#rfc with pca probability prediction
rfc_pca_p= cross_val_predict(rfc, x_fs_pca, y_fs, cv=10,method='predict_proba')
#logit with pca probability prediction
logit_pca_p= cross_val_predict(logit, x_fs_pca, y_fs, cv=10,method='predict_proba')
#gbr probability prediction
gbr_pca_p= cross_val_predict(gbr, x_fs_pca, y_fs, cv=10,method='predict_proba')
```

In [61]:

```
#divide odds into the odds for the training data set and the odds for the validation data set
odds2=odds.iloc[80:280,:]
odds3=odds.iloc[280,:,:]
```

In [63]:

```
#the function places a bet of one on outcomes that exceed a given probability and does not bet on all other predicted match outcomes
def strat(t,m):
    bettingstrategy=[ ]
    for e in m:
        if e[1]>t:
            bettingstrategy.append(1)
        elif e[2]>t:
            bettingstrategy.append(2)
        else:
            bettingstrategy.append( '-' )
    return(bettingstrategy)

#in this function, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. returns the total amount lost or won
def res(bettingstrategy, odds,y):
    results=[]
    invest=0
    for i,e in enumerate(bettingstrategy):
        if e==1 :
            invest=invest+1
            if y.iloc[i]==1:
                results.append(1*odds.iloc[i][0])
            else:
                results.append(-1)
        if e==2:
            invest=invest+1
            if y.iloc[i]==2:
                results.append(1*odds.iloc[i][2])
            else:
                results.append(-1)
        if e==0 :
            invest=invest+1
            if y.iloc[i]==0:
                results.append(1*odds.iloc[i][1])
            else:
                results.append(-1)
        else:
            results.append(0)
    return(results)

#in this function, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. returns the total betting amount
def res2(bettingstrategy, odds,y):
    results=[]
    invest=0
    for i,e in enumerate(bettingstrategy):
        if e==1 :
            invest=invest+1
            if y.iloc[i]==1:
                results.append(1*odds.iloc[i][0])
            else:
                results.append(-1)
        if e==2 :
            invest=invest+1
            if y.iloc[i]==2:
                results.append(1*odds.iloc[i][2])
```

```

    else:
        results.append(-1)
if e==0 :
    invest=invest+1
    if y.iloc[i]==0:
        results.append(1*odds.iloc[i][1])
    else:
        results.append(-1)
else:
    results.append(0)
return(invest)

```

In [64]:

```

#this function collects the sum of wins and losses as well as the total betting
amount for each probability threshold between 0.5 to 1 in steps of 0.01
def output1(model, odds, target):
    roi=[]
    for e in range(50,101):
        e=e/100
        if res2(strat(e,model), odds,target) !=0:
            roi.append(((sum(res(strat(e,model), odds, target))))-(res2(strat(e,
model), odds,target)))
        else:
            roi.append(0)
    return(roi)

```

In [65]:

```

#this function graphs the return on investment for each probability threshold be
tween 0.5 to 1 in steps of 0.01
def output_graph(model, odds, target):
    roi=[]
    bets=[]
    for e in range(50,101):
        e=e/100
        if res2(strat(e,model), odds,target) !=0:
            roi.append(((sum(res(strat(e,model), odds, target))))-(res2(strat(e,
model), odds,target)))
            bets.append(res2(strat(e,model), odds,target))
        else:
            roi.append(0)
            bets.append(0)
    t=[ ]
    for e in range (0,51):
        t.append(0)

    roi=t+roi

    plt.plot(roi, lw=3)
    plt.axis([51,100,-40,10])
    plt.axhline(y=0, color='r', linestyle='-', lw=3)
    plt.ylabel('Profit')
    plt.xlabel('Cutoff')
    plt.show()
    plt.rc('font', size=19)          # controls default text sizes
    plt.rc('axes', titlesize=19)      # fontsize of the axes title
    plt.rc('axes', labelsize=19)

```

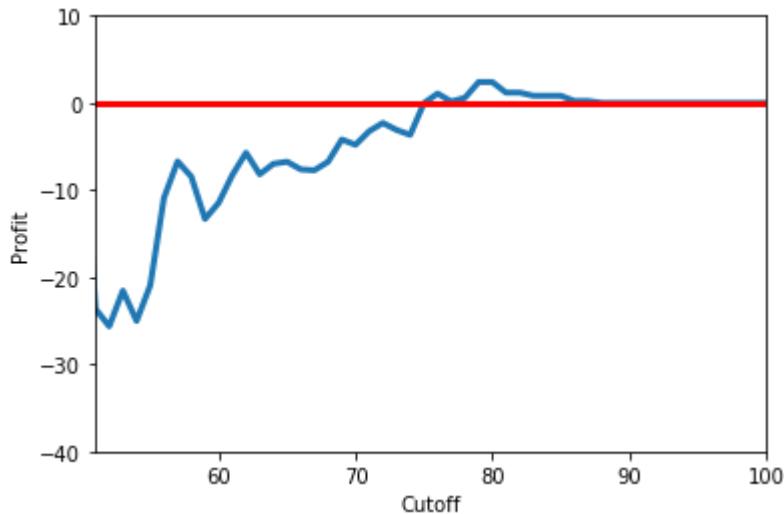
In [66]:

```
#this function evaluates the probability threshold (cutoff value) with the highest ROI, defined as the sum of wins and losses plus the betting amount divided by the betting amount and shows the corresponding number of bets,roi and cutoff value
def output_new(model, odds, target):
    roi=[ ]
    bets=[ ]
    for e in range(50,101):
        e=e/100
        if res2(strat(e,model), odds,target) !=0:
            roi.append(((sum(res(strat(e,model), odds, target)))-(res2(strat(e,
model), odds,target))))
            bets.append(res2(strat(e,model), odds,target)))
        else:
            roi.append(0)
            bets.append(0)

    a=round(max(roi),2)
    b=bets[roi.index(max(roi))]
    print('max bets profit: ' + str(a), '           cutoff: ' + str((roi.index(max(ro
i)))+50), '           bets: ' + str(b), '           roi: ' + str(((a+b)*100)/b)-100
) + '%' )
```

In [67]:

```
#output Lda for lasso reduced data set
output_graph(lda_pca_p, odds2, y_fs)
```



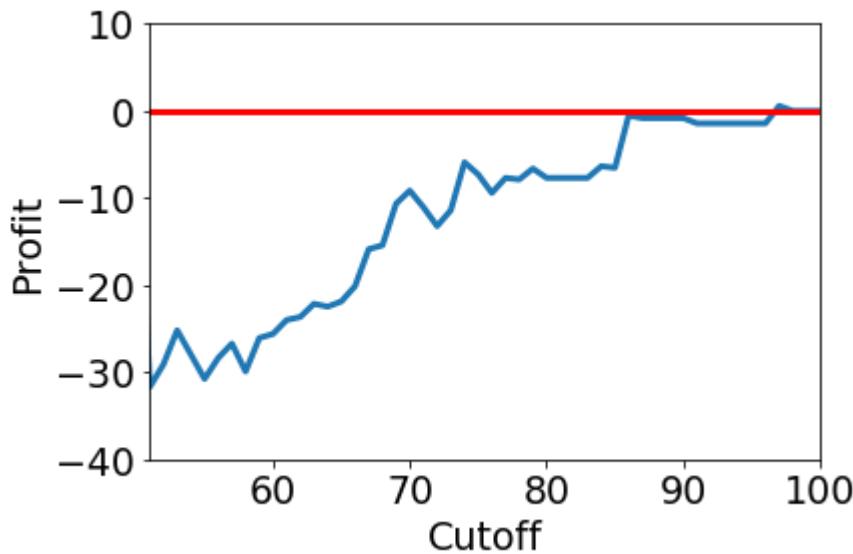
In [68]:

```
output_new(lda_pca_p, odds2, y_fs)
```

```
max bets profit: 2.39           cutoff: 78           bets: 8           roi:
29.875%
```

In [69]:

```
#output random forest for lasso reduced data set
output_graph(rfc_pca_p, odds2, y_fs)
```



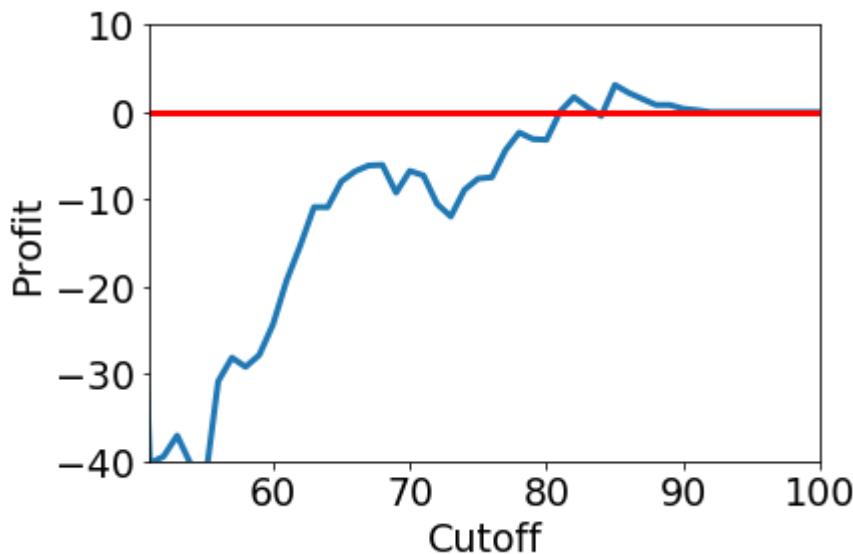
In [70]:

```
output_new(rfc_pca_p, odds2, y_fs)
```

max bets profit: 0.57                cutoff: 96                bets: 1                roi:  
56.9999999999997%

In [71]:

```
#output logit for lasso reduced data set
output_graph(logit_pca_p, odds2, y_fs)
```



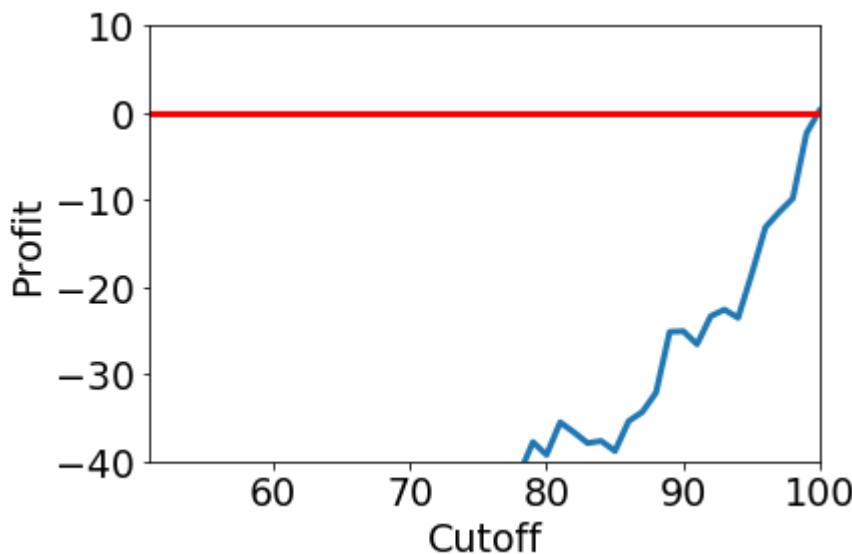
In [72]:

```
output_new(logit_pca_p, odds2, y_fs)
```

max bets profit: 3.11                cutoff: 84                bets: 10                roi:  
31.09999999999994%

In [73]:

```
##output gradient boosting for lasso reduced data set
output_graph(gbr_pca_p, odds2, y_fs)
```



In [74]:

```
output_new(gbr_pca_p, odds2, y_fs)
```

```
max bets profit: 0.41           cutoff: 99           bets: 2           roi:
20.5%
```

In [76]:

```
#this function gets out of sample performance (validation data set) with trainin
g set determined machine learning technique its is corresponding cutoff threshol
d
def output_vali(model, odds, target, e):
    e=e/100
    if res2(strat(e,model), odds,target) !=0:
        roi =((sum(res(strat(e,model), odds, target))))-(res2(strat(e,model), o
dds,target)))
        bets=(res2(strat(e,model), odds,target))
    else:
        roi= 0
        bets=0

    a=roi
    b=bets
    print('max bets profit: ' + str(a), '           cutoff: ' + str(e), '           bet
s: ' + str(b), '           roi: ' + str(((a+b)*100)/b)-100)+ '%')
```

In [78]:

```
#logit out of sample test (validation dat set)
logit.fit(x_fs_pca,y_fs)
logit_pca_vali = logit.predict_proba(x_test_pca)
```

In [135]:

```
output_vali(logit_pca_vali, odds3, y_test,84 )

max bets profit: 0.5199999999999996      cutoff: 0.84      bets:
3          roi: 17.33333333333314%
```

## home win/ no home win as target variable

In [81]:

```
#create y variable with new outcome
y_fs_h=[]
for e in y_fs:
    if e == 1:
        y_fs_h.append(1)
    else:
        y_fs_h.append(0)
```

In [82]:

```
#accuracies
# linear discriminate analysis full model
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis()
lda_p_h = cross_val_score(lda, x_fs, y_fs_h, cv=10)
print('lda: ' + str(lda_p_h.sum()/10))

# logit full model
from sklearn.linear_model import LogisticRegression
logit = LogisticRegression()
logit_p_h = cross_val_score(logit, x_fs, y_fs_h, cv=10)
print('logit: ' + str(logit_p_h.sum()/10))

#gba full model
from sklearn import ensemble
gbr = ensemble.GradientBoostingClassifier()
gbr_p_h = cross_val_score(gbr, x_fs, y_fs_h, cv=10)
print('gbr: ' + str(gbr_p_h.sum()/10))

#random forest full model
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier( max_depth=5,random_state=10)
rfc_p_h = cross_val_score(rfc, x_fs, y_fs_h, cv=10)
print('rfc: ' + str(rfc_p_h.sum()/10))

lda: 0.5730701754385964
logit: 0.614172932330827
gbr: 0.579812030075188
rfc: 0.6449874686716792
```

In [83]:

```
#accuracies
# linear discriminate analysis pca-reduced model
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis()
lda_p = cross_val_score(lda, x_fs_pca, y_fs_h, cv=10)
print('lda_pca: ' + str(lda_p.sum()/10))

# logit full model pca-reduced model
from sklearn.linear_model import LogisticRegression
logit = LogisticRegression()
logit_p = cross_val_score(logit, x_fs_pca, y_fs_h, cv=10)
print('logit_pca: ' + str(logit_p.sum()/10))

#gbrpca-reduced model
from sklearn import ensemble
gbr = ensemble.GradientBoostingClassifier()
gbr_p = cross_val_score(gbr, x_fs_pca, y_fs_h, cv=10)
print('gbr_pca: ' + str(gbr_p.sum()/10))

#random forest pca-reduced model
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier( max_depth=5,random_state=10)
rfc_p = cross_val_score(rfc, x_fs_pca, y_fs_h, cv=10)
print('rfc_pca: ' + str(rfc_p.sum()/10))
```

```
lda_pca: 0.655438596491228
logit_pca: 0.660438596491228
gbr_pca: 0.5801002506265664
rfc_pca: 0.5896240601503759
```

In [84]:

```
#accuracy bet365
prediction365_2=[]
for e in range(0,380):
    if odds.iloc[e][0]>odds.iloc[e][1] and odds.iloc[e][0]>odds.iloc[e][2]:
        prediction365_2.append(1)
    else:
        prediction365_2.append(0)
prediction365_2_2=prediction365_2[80:280]
prediction365_3_2=prediction365_2[280:]
accuracy_score(prediction365_2_2, y_fs_h)
```

Out[84]:

0.38

In [85]:

```
#strategy1
logit_pca_p_h=cross_val_predict(logit, x_fs_pca, y_fs_h, cv=10, method='predict_proba')
#strategy2
lda_pca_p_h=cross_val_predict(lda, x_fs_pca, y_fs_h, cv=10, method='predict_proba')
#strategy3
rfc_p_h=cross_val_predict(rfc, x_fs, y_fs_h, cv=10, method='predict_proba')
#gbr_p_h=cross_val_predict(gbr, x_fs, y_fs_h, cv=10, method='predict_proba')
```

In [86]:

```
#the function places a bet of one on outcomes that exceed a given probability and does not bet on all other predicted match outcomes
def strat_2(t,m):
    bettingstrategy=[ ]
    for e in m:
        if e[1]>t:
            bettingstrategy.append(1)
        else:
            bettingstrategy.append(' ')
    return(bettingstrategy)

#in this function, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. returns the total amount lost or won

def res_2(bettingstrategy, odds,y):
    results=[]
    invest=0
    for i,e in enumerate(bettingstrategy):
        if e==1:
            invest=invest+1
            if y[i]==1:
                results.append(1*odds.iloc[i][0])
            else:
                results.append(-1)
        else:
            results.append(0)
    return(results)

#in this function, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. returns the total betting amount
def res2_2(bettingstrategy,odds,y):
    results=[]
    invest=0
    for i,e in enumerate(bettingstrategy):
        if e==1:
            invest=invest+1
            if y[i]==1:
                results.append(1*odds.iloc[i][0])
            else:
                results.append(-1)
        else:
            results.append(0)
    return(invest)
```

In [87]:

```
#this function collects the sum of wins and losses as well as the total betting
amount for each probability threshold between 0.5 to 1 in steps of 0.01
def output_h(model, odds, target):
    roi=[]
    for e in range(50,101):
        e=e/100
        if res2_2(strat_2(e,model), odds,target) !=0:
            roi.append(((sum(res_2(strat_2(e,model), odds, target))))-(res2_2(st
rat_2(e,model), odds,target)))
        else:
            roi.append(0)
    return(roi)
```

In [88]:

```
#this function evaluates the probability threshold (cutoff value) with the highest
ROI, defined as the sum of wins and losses plus the betting
def output_new_h(model, odds, target):
    roi=[]
    bets=[]
    for e in range(50,101):
        e=e/100
        if res2_2(strat_2(e,model), odds,target) !=0:
            roi.append(((sum(res_2(strat_2(e,model), odds, target))))-(res2_2(st
rat_2(e,model), odds,target)))
            bets.append(res2_2(strat_2(e,model), odds,target))
        else:
            roi.append(0)
            bets.append(0)

    a=round(max(roi),2)
    b=bets[roi.index(max(roi))]
    print('max bets profit: ' + str(a), '           cutoff: ' + str((roi.index(max(ro
i)))+50), '           bets: ' + str(b), '           roi: ' + str(((a+b)*100)/b)-100
) + '%' )
```

In [89]:

```
#this function graphs the return on investment for each probability threshold between 0.5 to 1 in steps of 0.01
def output_graph_h(model, odds, target):
    roi=[ ]
    bets=[ ]
    for e in range(50,101):
        e=e/100
        if res2_2(strat_2(e,model), odds,target) !=0:
            roi.append(((sum(res_2(strat_2(e,model), odds, target))))-(res2_2(strat_2(e,model), odds,target)))
            bets.append(res2_2(strat_2(e,model), odds,target))
        else:
            roi.append(0)
            bets.append(0)

    t=[ ]
    for e in range (0,51):
        t.append(0)

    roi=t+roi

    plt.plot(roi, lw=3)
    plt.axis([51,100,-40,10])
    plt.axhline(y=0, color='r', linestyle='-', lw=3)
    plt.ylabel('Profit')
    plt.xlabel('Cutoff')
    plt.show()
    plt.rc('font', size=19)           # controls default text sizes
    plt.rc('axes', titlesize=19)       # fontsize of the axes title
    plt.rc('axes', labelsize=19)
```

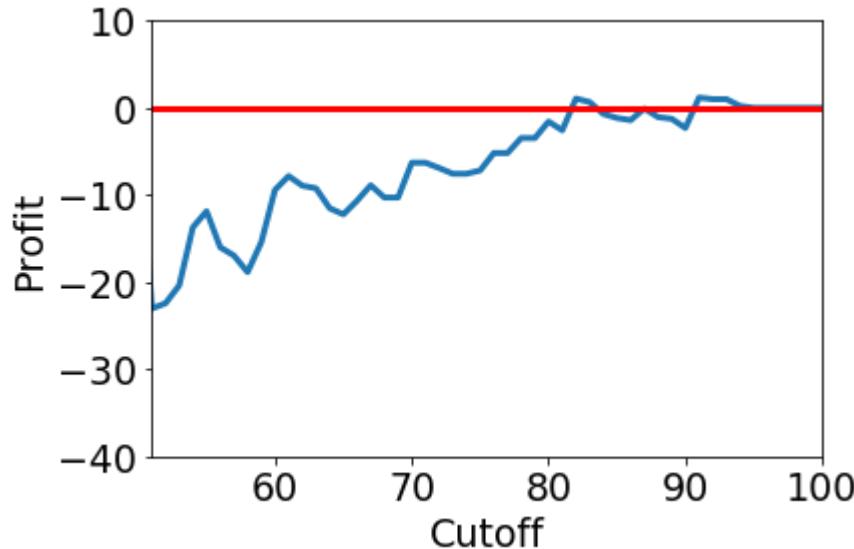
In [90]:

```
#output logit for lasso reduced data set
output_new_h(logit_pca_p_h, odds2, y_fs_h)

max bets profit: 1.18      cutoff: 90      bets: 5      roi:
23.599999999999994%
```

In [91]:

```
output_graph_h(logit_pca_p_h, odds2, y_fs_h)
```



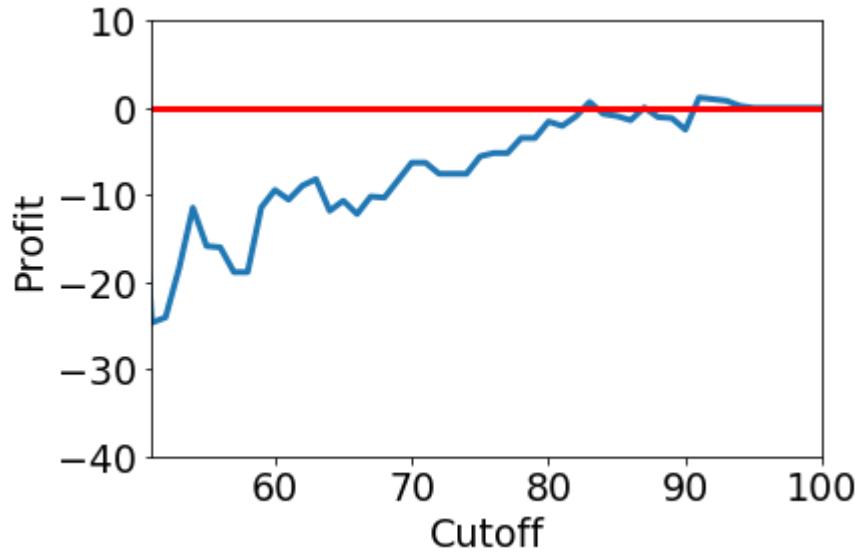
In [92]:

```
#output lda for lasso reduced data set  
output_new_h(lda_pca_p_h, odds2, y_fs_h)
```

```
max bets profit: 1.18          cutoff: 90          bets: 5          roi:  
23.599999999999994%
```

In [93]:

```
output_graph_h(lda_pca_p_h, odds2, y_fs_h)
```



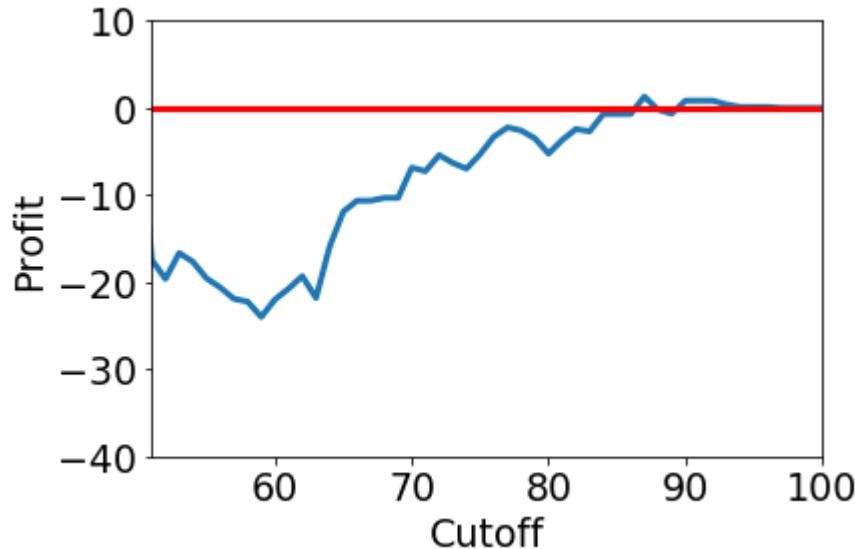
In [94]:

```
#output random forest for the full data set
output_new_h(rfc_p_h, odds2, y_fs_h)
```

```
max bets profit: 1.3          cutoff: 86          bets: 8          roi: 1
6.250000000000014%
```

In [95]:

```
output_graph_h(rfc_p_h, odds2, y_fs_h)
```



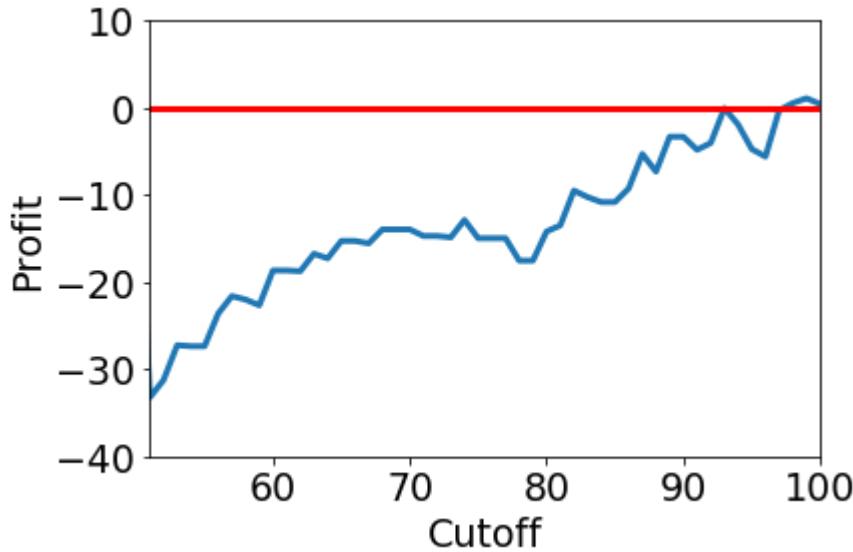
In [96]:

```
#output gradient boosting for the full data set
output_new_h(gbr_p_h, odds2, y_fs_h)
```

```
max bets profit: 1.1           cutoff: 98           bets: 4           roi: 2
7.499999999999986%
```

In [97]:

```
output_graph_h(gbr_p_h, odds2, y_fs_h)
```



In [98]:

```
#create y variable with new outcome for validation data set
y_test_h=[ ]
for e in y_test:
    if e == 1:
        y_test_h.append(1)
    else:
        y_test_h.append(0)
```

In [100]:

```
#this function gets out of sample performance (validation data set) with trainin
g set determined machine learning technique its is corresponding cutoff threshol
d
def output_vali_h(model, odds, target, e):
    e=e/100
    if res2_2(strat_2(e,model), odds,target) !=0:
        roi =((sum(res2_2(strat_2(e,model), odds, target))))-(res2_2(strat_2(e,m
odel), odds,target)))
        bets=(res2_2(strat_2(e,model), odds,target))
    else:
        roi= 0
        bets=0

    a=roi
    b=bets
    print('max bets profit: ' + str(a), '         cutoff: ' + str(e), '         bet
s: ' + str(b), '         roi: ' + str(((a+b)*100)/b)-100)+ '%')
```

In [104]:

```
#random forest out of sample test (validation dat set)
rfc.fit(x_fs,y_fs_h)
rfc_vali_h = rfc.predict_proba(x_test)
```

In [105]:

```
output_vali_h(rfc_vali_h, odds3, y_test_h,86)
```

```
max bets profit: 0.29000000000000004      cutoff: 0.86      bet
s: 3          roi: 9.66666666666671%
```

## away win/ no away win as target variable

In [106]:

```
#set target variable for new outcome
y_fs_a=[]
for e in y_fs:
    if e == 2:
        y_fs_a.append(2)
    else:
        y_fs_a.append(0)
```

In [107]:

```
#accuracies
# linear discriminate analysis full model
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis()
lda_p_a = cross_val_score(lda, x_fs, y_fs_a, cv=10)
print('lda: ' + str(lda_p_a.sum()/10))

# logit full model
from sklearn.linear_model import LogisticRegression
logit = LogisticRegression()
logit_p_a = cross_val_score(logit, x_fs, y_fs_a, cv=10)
print('logit: ' + str(logit_p_a.sum()/10))

#gba full model
from sklearn import ensemble
gbr = ensemble.GradientBoostingClassifier()
gbr_p_a = cross_val_score(gbr, x_fs, y_fs_a, cv=10)
print('gbr: ' + str(gbr_p_a.sum()/10))

#random forest full model
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier( max_depth=5,random_state=10)
rfc_p_a = cross_val_score(rfc, x_fs, y_fs_a, cv=10)
print('rfc: ' + str(rfc_p_a.sum()/10))
```

```
lda: 0.6831578947368422
logit: 0.7395238095238095
gbr: 0.6964035087719298
rfc: 0.6706390977443609
```

In [108]:

```
# linear discriminate analysis pca-reduced model
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
lda = LinearDiscriminantAnalysis()
lda_p_a = cross_val_score(lda, x_fs_pca, y_fs_a, cv=10)
print('lda: ' + str(lda_p_a.sum()/10))

# logit pca-reduced model
from sklearn.linear_model import LogisticRegression
logit = LogisticRegression()
logit_p_a = cross_val_score(logit, x_fs_pca, y_fs_a, cv=10)
print('logit: ' + str(logit_p_a.sum()/10))

#gba pca-reduced model
from sklearn import ensemble
gbr = ensemble.GradientBoostingClassifier()
gbr_p_a = cross_val_score(gbr, x_fs_pca, y_fs_a, cv=10)
print('gbr: ' + str(gbr_p_a.sum()/10))

#random forest pca-reduced model
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier( max_depth=5,random_state=10)
rfc_p_a = cross_val_score(rfc, x_fs_pca, y_fs_a, cv=10)
print('rfc: ' + str(rfc_p_a.sum()/10))
```

lda: 0.7647869674185463  
logit: 0.7502631578947369  
gbr: 0.7259774436090225  
rfc: 0.7199498746867168

In [109]:

```
#startegy1
logit_pca_p_a=cross_val_predict(logit, x_fs_pca, y_fs_a, cv=10, method='predict_proba')
#startegy2
lda_pca_p_a=cross_val_predict(lda, x_fs_pca, y_fs_a, cv=10, method='predict_proba')
#startegy3
rfc_pca_p_a=cross_val_predict(rfc, x_fs_pca, y_fs_a, cv=10, method='predict_proba')
#startegy3
gbr_pca_p_a=cross_val_predict(gbr, x_fs_pca, y_fs_a, cv=10, method='predict_proba')
```

In [110]:

```
#the function places a bet of one on outcomes that exceed a given probability and does not bet on all other predicted match outcomes
def strat_a(t,m):
    bettingstrategy=[ ]
    for e in m:
        if e[1]>t:
            bettingstrategy.append(2)
        else:
            bettingstrategy.append( '-')
    return(bettingstrategy)

#in this function, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. returns the total amount lost or won
def res_a(bettingstrategy, odds,y):
    results=[]
    invest=0
    for i,e in enumerate(bettingstrategy):
        if e==2:
            invest=invest+1
            if y[i]==2 :
                results.append(1*odds.iloc[i][2])
            else:
                results.append(-1)
        else:
            results.append(0)
    return(results)
#in this function, the placed bets are reconciled with the effective outcome. In case of a match between the bet and the outcome, the odds for the outcome are multiplied by the betting amount of one and noted as a win. returns the total betting amount
def res2_a(bettingstrategy,odds,y):
    results=[]
    invest=0
    for i,e in enumerate(bettingstrategy):
        if e==2:
            invest=invest+1
            if y[i]==2:
                results.append(1*odds.iloc[i][2])
            else:
                results.append(-1)
        else:
            results.append(0)
    return(invest)
```

In [111]:

```
#this function collects the sum of wins and losses as well as the total betting
amount for each probability threshold between 0.5 to 1 in steps of 0.01
def output_new_a(model, odds, target):
    roi=[]
    bets=[]
    for e in range(50,101):
        e=e/100
        if res2_a(strat_a(e,model), odds,target) !=0:
            roi.append(((sum(res_a(strat_a(e,model), odds, target))))-(res2_a(st
rat_a(e,model), odds,target)))
            bets.append(res2_a(strat_a(e,model), odds,target))
        else:
            roi.append(0)
            bets.append(0)

    a=round(max(roi),2)
    b=bets[roi.index(max(roi))]
    print('max bets profit: ' + str(a), '           cutoff: ' + str((roi.index(max(ro
i))+50), '           bets: ' + str(b), '           roi: ' + str(((a+b)*100)/b)-100
)+ '%')
```

In [112]:

```
#this function graphs the return on investment for each probability threshold be
tween 0.5 to 1 in steps of 0.01
def output_graph_a(model, odds, target):
    roi=[]
    bets=[]
    for e in range(50,101):
        e=e/100
        if res2_a(strat_a(e,model), odds,target) !=0:
            roi.append(((sum(res_a(strat_a(e,model), odds, target))))-(res2_a(st
rat_a(e,model), odds,target)))
            bets.append(res2_a(strat_a(e,model), odds,target))
        else:
            roi.append(0)
            bets.append(0)

    t=[]
    for e in range (0,51):
        t.append(0)

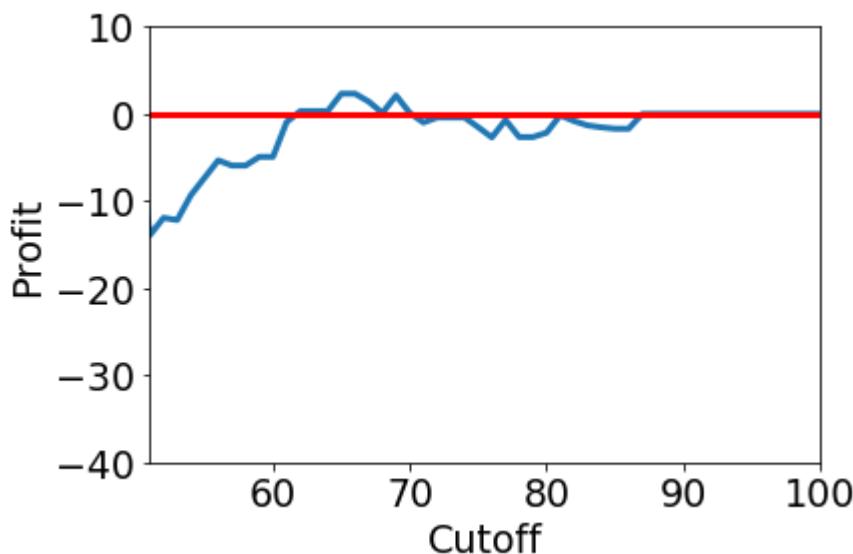
    roi=t+roi

    plt.plot(roi, lw=3)
    plt.axis([51,100,-40,10])
    plt.axhline(y=0, color='r', linestyle='-', lw=3)
    plt.ylabel('Profit')
    plt.xlabel('Cutoff')
    plt.show()
    plt.rc('font', size=19)          # controls default text sizes
    plt.rc('axes', titlesize=19)      # fontsize of the axes title
    plt.rc('axes', labelsize=19)
```

In [113]:

```
#output logit for the pca-reduced model  
output_new_a(logit_pca_p_a, odds2,y_fs_a )  
output_graph_a(logit_pca_p_a, odds2,y_fs_a )
```

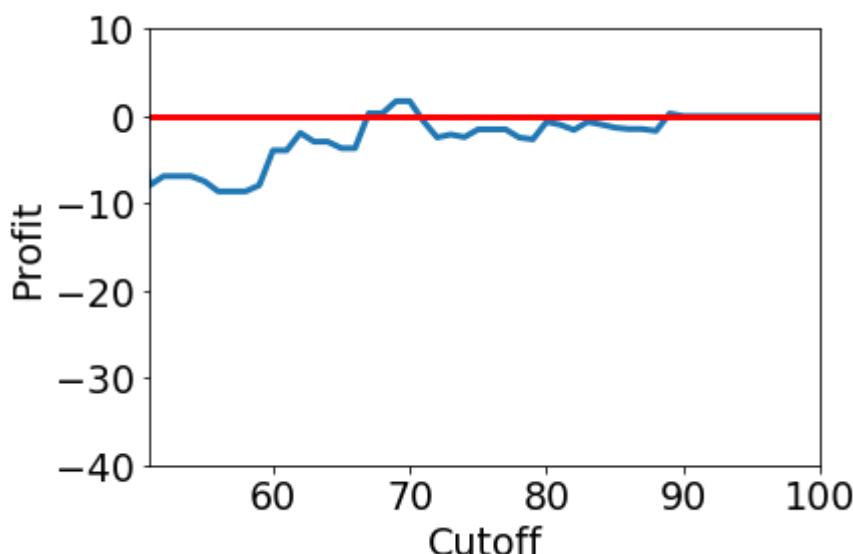
max bets profit: 2.32                  cutoff: 64                  bets: 31                  roi:  
7.483870967741936%



In [114]:

```
#output lda for the pca-reduced model  
output_new_a(lda_pca_p_a, odds2,y_fs_a )  
output_graph_a(lda_pca_p_a, odds2,y_fs_a )
```

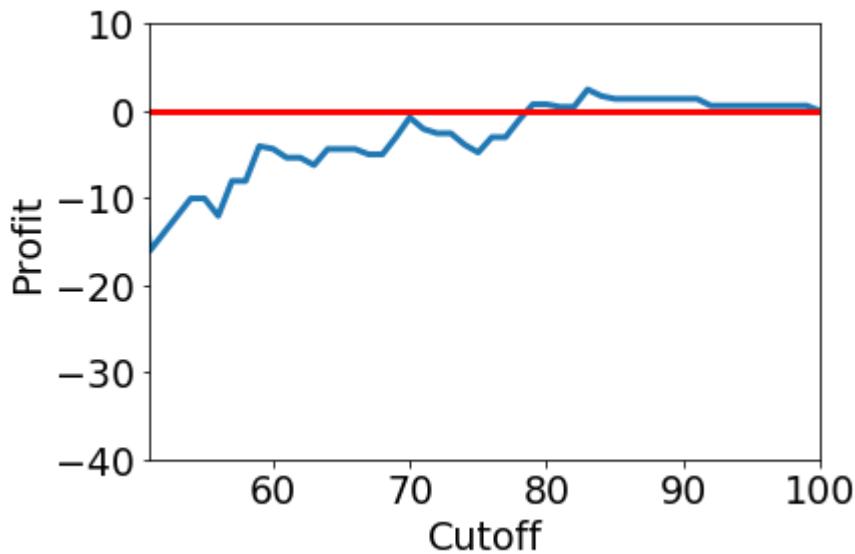
max bets profit: 1.71                  cutoff: 68                  bets: 30                  roi:  
5.700000000000003%



In [115]:

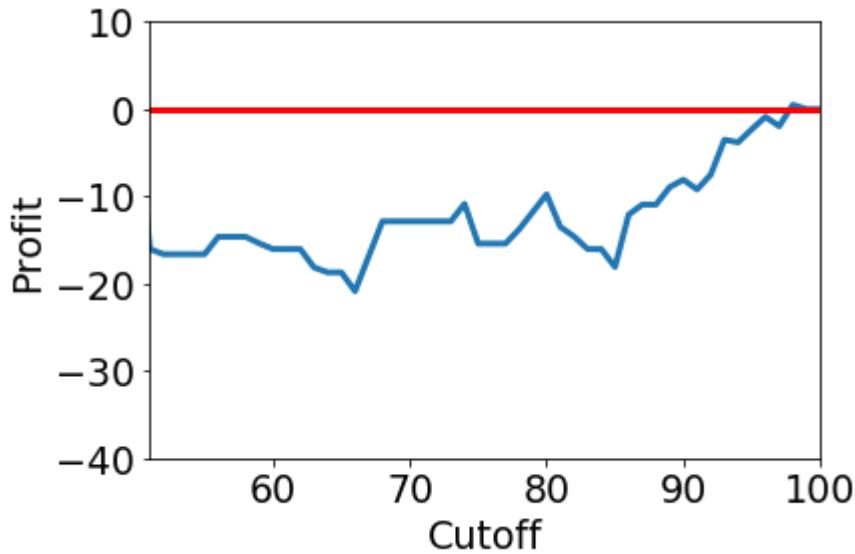
```
#output random forest for the pca-reduced model  
output_new_a(rfc_pca_p_a, odds2,y_fs_a )  
output_graph_a(rfc_pca_p_a, odds2,y_fs_a )
```

max bets profit: 2.46                    cutoff: 82                    bets: 5                    roi:  
49.1999999999999%



In [116]:

```
#output gradient boosting for the pca-reduced model  
output_graph_a(gbr_pca_p_a, odds2,y_fs_a )
```



In [118]:

```
#this function gets out of sample performance (validation data set) with trainin
g set determined machine learning technique its is corresponding cutoff threshol
d
def output_vali_a(model, odds, target, e):
    e=e/100
    if res2_h(strat_h(e,model), odds,target) !=0:
        roi =((sum(res_h(strat_h(e,model), odds, target)))-(res2_h(strat_2(e,m
odel), odds,target)))
        bets=(res2_h(strat_h(e,model), odds,target))
    else:
        roi= 0
        bets=0

    a=roi
    b=bets
    print('max bets profit: ' + str(a), '         cutoff: ' + str(e), '         bet
s: ' + str(b), '         roi: ' + str(((a+b)*100)/b)-100)+ '%')
```

In [120]:

```
#create y variable with new outcome for validation data set
y_test_a=[]
for e in y_test:
    if e == 2:
        y_test_a.append(2)
    else:
        y_test_a.append(0)
```

In [133]:

```
#random forest out of sample test (validation dat set)
rfc.fit(x_fs,y_fs_a)
rfc_vali_a = rfc.predict_proba(x_test)
```

In [134]:

```
output_vali_a(rfc_vali_a, odds3, y_test_a, 82)
```

```
max bets profit: 2.26         cutoff: 0.82         bets: 2         ro
i: 113.0%
```