# Predicting Football Match Outcomes with eXplainable Machine Learning and the Kelly Index

Yiming Ren[1] and Teo Susnjak[1]

[1]School of Mathematical and Computational Sciences, Massey University, Auckland, New Zealand

November 30, 2022

## Abstract

In this work, a machine learning approach is developed for predicting the outcomes of football matches. The novelty of this research lies in the utilisation of the Kelly Index to first classify matches into categories where each one denotes the different levels of predictive difficulty. Classification models using a wide suite of algorithms were developed for each category of matches in order to determine the efficacy of the approach. In conjunction to this, a set of previously unexplored features were engineering including Elo-based variables.

The dataset originated from the Premier League match data covering the 2019-2021 seasons. The findings indicate that the process of decomposing the predictive problem into sub-tasks was effective and produced competitive results with prior works, while the ensemble-based methods were the most effective.

The paper also devised an investment strategy in order to evaluate its effectiveness by benchmarking against bookmaker odds. An approach was developed that minimises risk by combining the Kelly Index with the predefined confidence thresholds of the predictive models. The experiments found that the proposed strategy can return a profit when following a conservative approach that focuses primarily on easy-to-predict matches where the predictive models display a high confidence level.

*Keywords*— Football match prediction; Machine learning;Kelly index ; eXplainable AI, Investment strategy

## 1 Introduction

Due to the worldwide appeal of football, the football industry has occupied an unassailable position in the sports business since its inception. The popularity of the sport is

1

continuing to increase and increasing sports fans are becoming involved in the football industry. At the same time, betting on sporting events, including football, has become a growth industry. The number of legal sports betting organisers and participants worldwide is increasing every year.

Meanwhile, the pre-match predictability of the sport is fairly low with some attribution placed on the sheer length of the matches as well as the number of players involved [29]. Others [10] have argued that the lower predictability of football can be explained by its low-scoring natures as well as the higher competitiveness relative to other sports and the fact that football can have more outcomes than merely wins and losses, but also draws, which contribute to the uncertainty of predicting the outcome of the matches[60]. Other reasons have also been suggested such as power struggles within the leadership of football clubs[30], match-fixing[30], referees with low ethical standards[13] and tacit understandings between teams on both sides of the game simultaneously reduce the predictability of football matches. The limited ability to counter match-fixing of local football leagues[46] and illegal internet gambling[25] are contributing components to the uncertainty of match outcomes.

The drive for increasing the predictability of football matches and profitability in gambling contexts, has encouraged researchers to work on devising effective strategies. These approaches have taken the form of statistical[4], machine learning[3], [21], natural language processing of football forums[5] as well as in-play image processing from camera footage of players[40] have been used in the attempt to create more accurate predictive systems of football matches. However, even in recent years, research efforts have realised limited and incremental successes at predicting the outcome of pre-match football matches.

In contrast to other sports, such as horse racing, there is no significant difference in the predictive accuracy between the opinions of the general public on internet forums, or the judgement of football experts and that of sophisticated machine learning methods[9] which underscores the challenge of the task. Even with access to rich sources of raw data and information, as well as the expert judgement of pundits, the accuracy of the best predictions from literature after converting odds set by bookies into probabilities of events reaches ~55%.

## Contribution

This study joins the growing body of literature that is attempting to improve the predictive accuracy of football matches using machine learning techniques, demonstrating the proposed methods on three seasons of Premier League matches in Europe sourced from publicly available datasets. In this work, we use a range of new features and algorithms that have not yet been used on this problem domain, and we employ tools from the explainable AI field in order to infuse the predictive models with interpretability and the ability to expose their reasoning behind the predictions. This study ultimately proposes a novel strategy for maximising the return on investment which relies on the Kelly Index in order to first categorise matches into levels of confidence with respect to their outcomes that are calculated on betting odds. We show how the decomposition of the problem with the aid of the Kelly Index and subsequent machine learning can indeed form an investment strategy that returns a profit.

# 2 Related work

This review divides previous studies in predicting the outcome of football matches into the target type of outcome (win/loss/draw) and those that have attempted at predicting the final scores of matches. However, because of the emergence of new betting options on betting platforms, predictions have emerged about whether or not both teams in the match will score goals[16] [33]. In fact, researchers now rarely predict the score of a match, as the uncertainty of the score and the existence of draws make score prediction approach challenging and the prediction accuracy low, even with the inclusion of betting odds data[52]. The methods used for predicting the result of football matches fall into 3 categories. In the literature, statistical models, machine learning algorithms and rating systems have been explored.

## Statistical and rating-system approaches

Inan [33] fitted the team's offensive and defensive capabilities calculated from the goals made and conceded by teams each week to a Poisson distribution in his study while also taking the home advantage into account. After fitting the team's offensive and defensive capabilities, home team advantage and betting parameters into a Poisson model, Egidi et al. [18] has developed a method of testing the prediction results to obtain more accurate predictions.

Robberechts and Davis [53] used ELO ratings in predicting the result of FIFA World Cup matches. this approach considers the team's earlier result (whether it won or not) as the team's performance, and then gradually adjusts for subsequent results to obtain the final team performance rating. In order to identify and adjust a team's score, Beal et al. [5] used goals by the team. They also considered the strength of the opponent of the match which is the source of the data to correct the team performance scores. In addition, Constantinou [15] used a hybrid Bayesian network to develop a dynamic rating system. The system corrects the team performance ratings by assigning greater weight to data on the results of the most recent matches that occurred. The researchers found that the study that uses scoring systems to predict the results of a football match usually uses databases with less content. These approaches need a small amount of data to extract features that can be used for prediction.

Koopman and Lit [38] used both a bivariate Poisson distribution model and a team performance rating system in their research. They held home team advantage constant when predicting National League outcomes but varied the team's offensive and defensive capabilities over time. The results proved that there was no major difference in prediction between the two methods and that both could make good predictions. However, the team performance rating system was ineffective at predicting draws. Most of the above research on the use of rating systems noted the near absence of ties in the predicted results. Berrar et al. [6] used team performance scores as a predictive feature for machine learning.

## Machine learning approaches

Efforts to improve the predictability of football have focused on two aspects. They have either attempted to use more sophisticated machine learning algorithms, or on feature engineering in order to develop more descriptive variables. Sometimes both aspects have been pursued. Table 1 summarises all the key studies and their attributes.

### Algorithms

A mixture of more simplistic approaches like Linear Regression[50], K-nearest neighbors[8] and Decision Trees[62] have been explored. Subsequently, the algorithms have increased in their sophistication with the use of Random Forests[23], Support Vector Machines[44], Artificial Neural Networks[11] and boosting Boosting[28]. More recently, Deep learning methods have emerged with the use of Convolutional Neural Networks [26] and LSTM Zhang et al. [64], Malini and Qureshi [42] where the authors concluded that the LSTM model displayed superior characteristics to traditional machine learning algorithms and artificial neural networks in predicting the results of matches [42].

### Features

Literature has shown that some of the effective features so far are half-time goal data as used by Yekhande et al. [61], first goal team data and individual technical behaviour data used by Parim et al. [45], ball possession and passing over data used by Bilek and Ulas [7], key player position data used by Joseph et al. [35].

Kınalıoğlu and Coskun[37] compared predictive models of six machine learning algorithms and evaluated models with different hyperparameters settings using a large number of model evaluation methods. Beal et al. [5] improved prediction accuracy by ~7% by extracting expert and personal opinions published on self-published media and combining environmental, player sentiment, competition, and external factors in their predictions.

Overall, studies using machine learning methods mention the need to use richer predictive features and validate them more robustly in future works. However, even newer studies have a tendency to reuse existing features and limited proposals have been made for innovating with novel features. Moreover, suggestions for new predictive features have primarily focused on ones that are not available in public datasets [60] such as player transfer data, injuries, expert advice, and psychological data.

Direct prediction of odds are rarely researched in this context, this is because the calculation of the odds itself contains the bookmaker's prediction of the results of the matches[63]. Using odds to predict the outcome of a match requires a combination of machine learning or statistical methods. Zeileis et al. [63] and Wunderlich and Memmert [60] both reverse the odds into the probability of a team winning and fit it into the ELO team performance scoring system to predict the outcome of the match. Štrumbelj [56] used traditional linear regression models and Shin models to predict the probability of winning directly from the odds, with the aim of inverting the bookmaker's method of calculating odds based on the betting.

## Summary of literature

The reported accuracies by existing studies as summarised in Table 1 are chiefly determined first by the nature of the features used, and specifically if these are generated in real-time during in-play match outcome predictions. Features generated in-play are more deterministic and thus result in considerably higher accuracies than more static pre-match features when predicting outcomes. Therefore these kinds of studies are different in nature than those that simply rely on pre-match prediction.

Next, the accuracy is strongly affected by the nature of the predictive problem in the manner in which it is formulated. Models that attempt to predict a simple win or loss scenario are more accurate than those that also attempt to predict draws. This is expected

**Table 1:** Prediction accuracy of different studies

| Study | Competition | Features | Non-pre-match features | Best Algorithm | Accuracy | Test set duration | Matches | Class |
|---|---|---|---|---|---|---|---|---|
| Joseph et al. [35]2006 | Tottenham Hotspur Football Club | 7 | None | Bayesian networks | 58% | 1995-1997(2 years) | 76 | 3 |
| Constantinou [15] 2019 | 52 football leagues | 4 | None | Hybrid bayesian networks | - | 2000-2017(17 years) | 216,743 | 3 |
| Hubáček et al. [27]2019 | 52 football leagues | 2 | None | Double Poisson | 48.97% | 2000-2017(17 years) | 218,916 | 3 |
| Mendes-Neves and Mendes-Moreira [43]2020 | 6 Leagues | 28 | None | Bagging | 51.31% | 2016-2019(3 years) | 1,656 | 3 |
| Danisik et al. [17]2018 | 5 Leagues | 139 | None | LSTM regression | 52.5% | 2011-2016 | 1520 | 3 |
| Berrar et al. [6]2019 | 52 football leagues | 8 | None | KNN | 53.88% | 2000-2017(17 years) | 216,743 | 3 |
| Herbinet [24]2018 | 5 leagues | 6 | None | XGBoost | 54% | 2014-2016(2 years) | 3,800 | 3 |
| Carloni et al. [12]2021 | 12 countries | 47 | None | SVM | 57% | 2008-2020(12 years) | 49,319 | 3 |
| Baboota and Kaur [3]2019 | English Premier League | 12 | None | Random Forest | 57% | 2017-2018(1 year) | 380 | 3 |
| Chen [14]2019 | La Liga | 34 | None | Convolution neural network | 57% | 2016-2017(1 year) | 380 | 3 |
| Esme and Kiran [20]2018 | Super League of Turkey | 17 | None | KNN | 57.52% | 2015-2016(half year) | 153 | 3 |
| [54]2022 | English Premier League | 31 | None | SVM | 61.32% | 2018-2019(1 year) | 380 | 3 |
| Alfredo and Isa [1]2019 | English Premier League | 14 | Half-time results | Random Forest | 68.16% | 2007-2017(10 years) | 3,800 | 3 |
| Pathak and Wadhwa [47]2016 | English Premier League | 4 | None | Logistic Regression | 69.5% | 2001-2015(14 years) | 2280 | 2 |
| Prasetio et al. [50]2016 | English Premier League | 13 | Ball possession, Distance run | Logistic Regression | 69.51% | 2015-2016(1 year) | 380 | 2 |
| Danisik et al. [17]2018 | 5 leagues | 139 | None | LSTM regression | 70.2% | 2011-2016 | 1520 | 2 |
| Ievoli et al. [32]2021 | UEFA Champions League | 29 | Number of passes and other features | BLR | 81% | 2016-2017(1 year) | 75 | 3 |
| Azeman et al. [2]2021 | English Premier League | 11 | Shooting, fouls and other 8 features | Multiclass decision forest | 88% | 2005-2006(1 year) | 380 | 3 |

since all predictive problems become more complex as the number of categories being predicted increases. For this reason, a number of researchers have opted for eliminating the prediction of tied matches.

In addition, the size of the datasets also affects the overall capability of the models, however, the number of features used in football match prediction has a small effect on prediction accuracy (Figure 1). The figure highlights that the accuracy of pre-match prediction is consistently below 60%. However, when predicting pre-match results over multiple seasons, the prediction accuracy was almost always below 55%. The researchers were unable to break the bottleneck in predicting football matches with more advanced algorithms and scoring systems.

When predictions are made over a large time span, for example, predicting the results of a game over 10 seasons, then these models will be less accurate than those predicting the results of a game over 1 season. Furthermore, the choice of machine learning method does not seem to have a significant impact on the accuracy of the model.
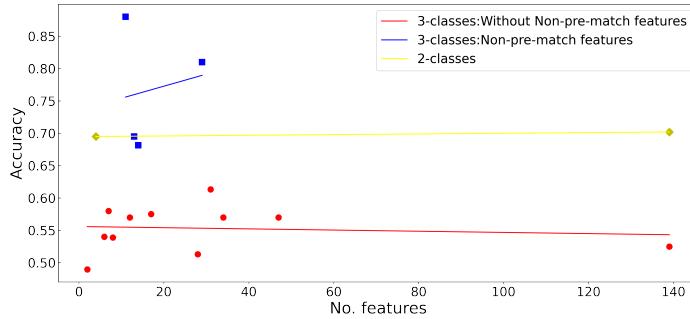


**Figure 1:** The plot of the number of features against prediction accuracy for all types of literature

**Research question**

In light of the literature, this study poses the following research questions:

1. Is it possible to use public data to find out which matches are difficult to predict and which results are already certain before the match?

2. Can prediction accuracy be improved by training models using hard-to-predict matches or easy-to-predict matches?

3. Which machine learning predictions perform best for hard-to-predict matches and easy-to-predict matches?

4. Is it possible to increase returns by applying different weights to hard-to-predict matches and easy-to-predict matches when investing?

# 3 Methodology

## 3.1 Machine learning algorithms

The machine learning methods used are drawn from open-source machine learning libraries. These include Scikit-learn[48], XGboost and CatBoost.

**Logistic regression**   Logistic regression predicts the probability on the basis of linear regression. It converts predictions to probabilities and then implements classification. The algorithm demonstrates a strong ability to deal with data with a strong linear relationship between features and labels. Logistic regression also has the advantage of being noise-resistant and performs well on small data sets. Prasetio et al. [50] used logistic regression in predicting the English Premier League. Using in-play features, the authors obtained a prediction accuracy of 69%.

**Decision Tree**   This algorithm builds a Decision Tree [51] by choosing features that have the greatest information gain. The advantage is that it does not require any domain knowledge or parameter assumptions and is suitable for high-dimensional data. However, this algorithm is prone to overfitting and tends to ignore correlations between attributes. The use of a small number of features in the Decision Tree methods can prevent overfitting when predicting the results of football matches. The Decision Tree approach may therefore produce a more explanatory and thus informative model.

**Random Forest**   This is an ensemble-based method which combines multiple Decision Trees. The Random Forest approach integrates multiple Decision Trees and builds each one using a non-backtracking approach. Each of these Decision Trees relies on independent sampling and has the same randomness in data selection as the other trees. When classifying the data, each tree is polled to return the category with the most votes, giving a score for the importance of each variable and assessing the importance of each variable in the classification. Random Forests tend to perform better than Decision Trees in preventing overfitting.

   The correlation between the model features and the test predictions in the Random Forest approach is not strong. For poorly predicted football match results, this may be an advantage of the Random Forest approach. Using a training set containing a large number of features can result in a generated model with strong predictive power on a small number of them and weak predictive power on others. The classification criteria are different for each Decision Tree in the Random Forest model. This can also be a good way to reduce the correlation between Decision Trees.

**k-Nearest Neighbor**   The k-Nearest Neighbor(KNN) method classifies samples by comparing a particular sample with its $k$ closest data points in a given dataset and therefore does not generate a classifier or a model [49].

   KNN methods have an excellent performance when using data from in-play contexts to predict the results of contests[20]. However, the method needs to consider all the data in the sample when generating the prediction model. Such an approach creates a strong risk of overfitting in the prediction of pre-match results in competitive sports. Nevertheless, it is still not known how well the KNN method can perform for hard-to-predict matches.

7

**Gradient Boosting**   The idea behind gradient boosting is to iteratively generate multiple weak models and then add up the predictions of each weak model. Each residual calculation increases the weight of the wrongly split samples, while the weights of the split pairs converge to zero, so the generalisation performance is better. The importance of the features of the trained model can be extracted.

In previous studies[1] of match prediction gradient boost methods have often failed to produce the best models. However, the predictive power of the method in hard-to-predict matches cannot be denied outright.

**XGBoost**   XGBoost extends and improves upon the Gradient Boosting. The XGBoost algorithm is faster. It takes advantage of the multi-threading of the CPU based on traditional Boosting and introduces regularisation to control the complexity of the model. Prior to the iteration, the features are pre-ranked for the nodes and by traversing them the best segmentation points are selected, which results in lower complexity of data segmentation.

The XGBoost method is popular amongst machine learning researchers. However, the method does not perform well in capturing high-dimensional data, such as images, audio, text, etc. For predicting football matches with based on lower-dimensional data, XGBoost may produce promising models.

**CatBoost**   CatBoost algorithm is a type of gradient boosting algorithm. In contrast to the gradient boosting algorithm, this algorithm can also achieve excellent results using data that has not been feature engineered. It does better than gradient boosting in preventing overfitting. The CatBoost method has not been explored in previous studies of football match result prediction.

**Voting and Stacking**   The voting classifier can vote on the results of different models, with the majority deciding the outcome. The Voting classifier is divided into Hard and Soft implementation. The Hard voting is votes on the results obtained using multiple machine learning methods, with the majority getting the result. Soft voting is the process of calculating the mean of the classification probabilities of the different methods, and finally selecting the one with the highest mean as the prediction result.

The stacking algorithm refers to a hybrid estimator approach in which a number of estimators are fitted separately on training data, while the final estimator is trained using the stacked predictions of these basic estimators. However, the stacking method uses cross-validation which is not favoured in the field of football result prediction when training the model.

Voting and stacking methods have produced the best prediction models in other studies of athletic match result prediction[39]. However, in football these two methods are rarely used. This study uses both methods in an attempt to find the most suitable model for prediction in difficult-to-predict matches.

## 3.2   Dataset characteristics and feature engineering

The Premier League data used in this study was sourced from the website https://www.footballdata.co.uk/englandm.php. The results from the 2001 to 2018 seasons were used to develop features based on ELO ratings, attack and defence ratings and home-team and away-team winning percentages. Data from the 2019 to 2021 seasons were used for

testing. The raw data was recorded using a total of 107 columns of data for the 2019 to 2021 seasons.

The Premier League system divides the season into 38 rounds and guarantees that each of the 20 teams will play once in each of the 10 matches in each round. After 380 matches per season, each team plays 19 times in the home position and 19 times in the away position. This study predicts the results of 3 seasons, implying a total of 1140 predicted matches.

The Premier League is based on a points system, with the winning team gaining three points after the match and the two drawing teams gaining one point each. 24 teams are entered for the 2019-2021 season, all of which appeared in the 2001-2018 season. However, due to differences in the raw data on each season on bookmaker odds records, in fact, only 21 of the 83 columns of data were used in the 2019-2021 season predictions. The records are for the six European bookmakers and the European average odds respectively.

The premise of this study for feature engineering is the assumption that all matches occur continuously.

A total of 52 features were generated and explored in this project. Their names and descriptions are given in Table 2.

## ELO Rating

The ELO rating was first introduced by ELO[19] to assess the ability of players in chess competitions and has been extended to the assessment of the performance of players or teams in many sports. The rationale is to achieve a dynamic evaluation of a particular player's ability by comparing the ability of the player in a recent match with that player's past performance. The time-series nature of football matches makes the ELO model suitable for scoring the performance of football teams[34].

The score of the participating teams in a football match is determined by their performance in past matches (Eq 1) and the rating of their performance in the current match of the competition (Eq'2) are as follows:

$$E^H = \frac{1}{1 + \dfrac{c^{R^H - R^A}}{d}} \qquad E^A = 1 - E^H \tag{1}$$

$$S^H = \begin{cases} 1 & Home \quad team \quad win \\ 0.5 & Draw \\ 0 & Away \quad team \quad win \end{cases} \tag{2}$$

$$S^A = 1 - S^H \tag{3}$$

where $E^H$ and $E^A$ are the performance ratings that teams are supposed to have in that match, $S^H$ and $S^A$ are the actual performance ratings of two teams. $R^H$ and $R^A$ are the ratings that the two teams have constantly revised to represent their own team's ability, and $c$ and $d$ are two constants based on the scale of scoring. In this study $c$ was set to 10 and $d$ to 400 to ensure that the effect of the new matches' results on the previous ELO ratings was kept in a balanced range.. The scoring of the two teams is appropriately corrected by the error between the ability of the two teams during the newly occurring match and previously assessed performance, using the following:

$$R^{'H} = R^H + k(S^H - S^E) \tag{4}$$

Predicting Football Match Outcomes with eXplainable ML

**Table 2:** Description of features

| Feature Name | Description |
| --- | --- |
| AvgGoalDiff | Average goal difference between the two teams in the previous six games |
| TotalGoalDiff | Goal difference between the two teams in the previous six games |
| HomeELO | ELO ratings for the home team |
| AwayELO | ELO ratings for away teams |
| ELOsta | Standard deviation of the ELO ratings of the two teams |
| ELOHomeW | Probability of home team winning converted by ELO ratings |
| ELOAwayW | Probability of away team winning converted by ELO ratings |
| ELODraw | Probability of a draw occurring converted by ELO ratings |
| one_ELO | Probability of conversion from ELO ratings after PCA dimensionality reduction |
| HomeHELO | ELO ratings for the home team's half-time result |
| AwayHELO | ELO ratings for the away team's half-time result |
| HELOSta | Standard deviation of the half-time ELO ratings of the two teams |
| ELOHHomeW | Probability of home team winning converted by half-time ELO ratings |
| ELOHAwayW | Probability of away team winning converted by half-time ELO ratings |
| ELOHDrawW | Probability of a draw occurring converted by half-time ELO ratings |
| one_HELO | Probability of conversion from half-time ELO ratings after PCA dimensionality reduction |
| HomeTeamPoint | Current home team points in the season |
| AwayTeamPoint | Current away team points in the season |
| PointDiff | Current point difference between home and away teams in the season |
| AvgHOddPro | Average of pre-match home team odds offered by all bookmakers in the available data |
| AvgAOddPro | Average of pre-match away team odds offered by all bookmakers in the available data |
| AvgDOddPro | Average of pre-match draw odds offered by all bookmakers in the available data |
| one_Odd_Pro | Average odds after PCA dimensionality reduction |
| HomeOff | Rating of the home team's offensive capabilities |
| HomeDef | Rating of the home team's defensive capabilities |
| AwayOff | Rating of the away team's offensive capabilities |
| AwayDef | Rating of the away team's defensive capabilities |
| Offsta | Difference in offensive capability rating |
| Defsta | Difference in defensive capability rating |
| AvgShotSta | Standard deviation of shots on goal for both sides in six matches |
| AvgTargetSta | Standard deviation of shots on target for both sides in six matches |
| ShotAccSta | Standard deviation of shot accuracy between the home team and the away team |
| AvgCornerSta | Standard deviation of corners for both sides in six matches |
| AvgFoulSta | Standard deviation of fouls for both sides in six matches |
| HomeHWin | Home team's winning percentage in home games |
| HomeHDraw | Home team's draw percentage in home games |
| AwayAWin | Away team's winning percentage in away games |
| AwayADraw | Away team's draw percentage in away games |
| HomeWin | Home team's win percentage in all previous matches |
| HomeDraw | Home team's draw percentage in all previous matches |
| AwayWin | Away team's win percentage in all previous matches |
| AwayDraw | Away team's draw percentage in all previous matches |
| LSHW | Home team win percentage last season |
| LSHD | Home team draw percentage last season |
| LSAW | Away team win percentage last season |
| LSAD | Away team draw percentage last season |
| Ysta | Difference in the number of yellow cards between the home team and the away team |
| Rsta | Difference in the number of red cards between the home team and the away team |
| StreakH | Home team winning streak index |
| StreakA | Away team winning streak index |
| WStreakH | Home team weighted winning streak index |
| WStreakA | Away team weighted winning streak index |

10

where $k$ indicates the magnitude of the correction. $k$ is used to determine the effect of the results of the new competition on the scoring. $k$ is calculated as follows:

$$k = k_0(1 + \delta)^\gamma \qquad (5)$$

where $\delta$ is the absolute goal difference and $k_0$ and $\gamma$ are set to 10 and 1.

The performance ability of teams was first assessed in this study using match results from the 2001-2005 English Premier League season. Teams were ranked and assigned ratings by comparing the total scores and wins of teams in the 2001-2002 season. The team ratings were corrected using match data from the 2002-2017 season. 2018-2019 season was formally characterised using the ratings generated by the ELO model, and each team's rating was subsequently corrected for the previous match results.

In a subsequent study on ELO ratings, a regression model was generated to predict wins and losses by the ELO rating profile of the teams from both sides before the match[31]. The model has a high degree of confidence among those involved in football gambling. The probability of different match outcomes occurring can be simply predicted using the ELO scores of the teams on both sides of the match using the following formula:

$$P_H = 0.448 + (0.0053 * (E^H - E^A)) \qquad (6)$$

$$P_A = 0.245 + (0.0039 * (E^A - E^H)) \qquad (7)$$

$$P_D = 1 - (P_H + P_A) \qquad (8)$$

where $P^H$ indicates the probability of the home team victory, $P^A$ indicates the probability of the away team victory occurring and $P^D$ indicates the probability of a draw situation occurring.

## Offensive and defensive capabilities

The Offensive and Defensive Model (ODM) [22] is used to generate a team's offensive and defensive abilities. The calculation formula (Eq 9) for scoring given in the introduction of the ODM model by Govan et al. [22] is as follows:

$$o_j = \sum_{i=1}^{n} \frac{A_{ij}}{d_i} \qquad dj = \sum_{i=1}^{n} \frac{A_{ji}}{o_i} \qquad (9)$$

where $A_{ij}$ is a goal scored by team $j$ in a match between teams from both sides of the match, $o_j$ is the offensive rating and $d_j$ is the defence rating. The equation shows that the two ratings affect each other. We directly assess the offensive and defensive ratings of a team one season and apply them to the next season. A team with high goal-scoring has a high offensive rating. Likewise, a team with a high number of goals conceded has an inefficient defensive rating.

For teams new to the season, we extrapolate their offensive and defensive ratings based on that team's ranking in goals scored and goals against during the season. This is calculated as follows:

$$o_j = \frac{XG_2 + X_2G_1 - X_1G_1 - X_1G_2}{X_2 - X_1} \qquad (10)$$

where $O_j$ represents the new team's offensive score. $G$ is the goals scored by the new team during the season. $G1$ and $R1$ are the goals scored and offensive rankings of the teams that are one place above the new team in the season goal rankings. $G2$ and $R2$ are the goals scored and rankings of the teams that are one place below. The defensive rating is calculated in the same way.

**Streak Index**

The Streak Index represents a team's capacity to go on a winning streak which captures a team's most recent form. The index was generated by Baboota and Kaur [3] in predicting the results of the Premier League for the two seasons 2014-2016. the Streak Index is further divided into two categories, the first of which directly measures the team's form over the previous $k$ games, with the following formula:

$$S = (\sum_{p=j-k}^{j-1} res_p)/3k \tag{11}$$

where $j$ indicates the number of matches to be predicted and $res_p$ indicates the team's score in a particular match, where a win and a draw correspond to 3 points and 1 point respectively. No points are awarded for losses. In the second category, the Streak Index, is weighted and normalised based on the first category. Relatively low weights are given to matches that occur further back in time. The formula for this is as follows:

$$\omega_S = \sum_{p=j-k}^{j-1} 2\frac{p - (j - k - 1)res_p}{3k(k-1)} \tag{12}$$

The Streak Index is logically problematic when calculating across seasons. This means that when calculating the Streak Index for the first few games of a season, the data used is from the last season's end. The changes in team form that occur during the intervening season can compromise the effectiveness of the Streak Index.

**Principal Component Analysis**

As the outcome of a football match includes a draw, bookmakers also offer odds in the event of a draw. This requires a feature reduction method to turn three-dimensional data into one-dimensional data.

This study uses the Principal Component Analysis (PCA) [59] to reduce the dimensionality of the data. This function linearly reduces the dimensionality on the original data by projecting the original data to an alternative dimension where the data loss is minimised. PCA is generally used to process and compress high-dimensional data sets. The main objective is to reduce the dimensionality of a large feature set into one that is most essential with the tradeoff that the new features lose interpretability.

**Kelly Index**

The Kelly index[57] was used as a tool to classify the matches in the dataset into categories that define different levels of predictability. In literature, the Kelly Index was originally used to calculate the flow-through rate of electronic bits[36]. However, due to its probabilistic nature and its similarity to the nature of betting, the Kelly Index is also widely used by bettors.

Each bookmaker calculates the Kelly Index before a match and makes it available to members who gamble. The Kelly Index can be calculated from the available data using the following formula:

$$K_H = (O_H/avgO_H)F_{99} \tag{13}$$
$$K_A = (O_A/avgO_A)F_{99} \tag{14}$$
$$K_D = (O_D/avgO_D)F_{99} \tag{15}$$

where $K_H$, $K_A$, and $K_D$ indicate the three results of matches in terms of the Kelly Index. $avgO_H$, $avgO_A$, and $avgO_D$ indicate in our dataset the average European betting market odds for each of the three results of the matches. $F_{99}$ indicates the return rate of the European average odds, which is calculated by the following formula:

$$F_{99} = \frac{1}{\frac{1}{avgO_H} + \frac{1}{avgO_A} + \frac{1}{avgO_D}} \tag{16}$$

If a bookmaker has set lower odds on an event occurring than any other bookmaker, this means that the bookmaker is confident enough to believe that a particular outcome will occur. This confidence may come from the bookmaker's perception of possessing more effective prediction models, having access to more extensive sources of information, or illegal black market trading. The bookmaker as a dealer is willing to set the lowest odds on a more likely event to maximise profits. Therefore a bookmaker's Kelly Index is the most likely event to occur. At the same time, if there is a low Kelly Index at one bookmaker, there will be a high Kelly Index for other events accordingly. However, there is no guaranteed correlation between the level of the Kelly Index and the result of a match as this is just a technique for understanding what odds have been determined to be in the best interest of the bookmaker for maximising their returns.
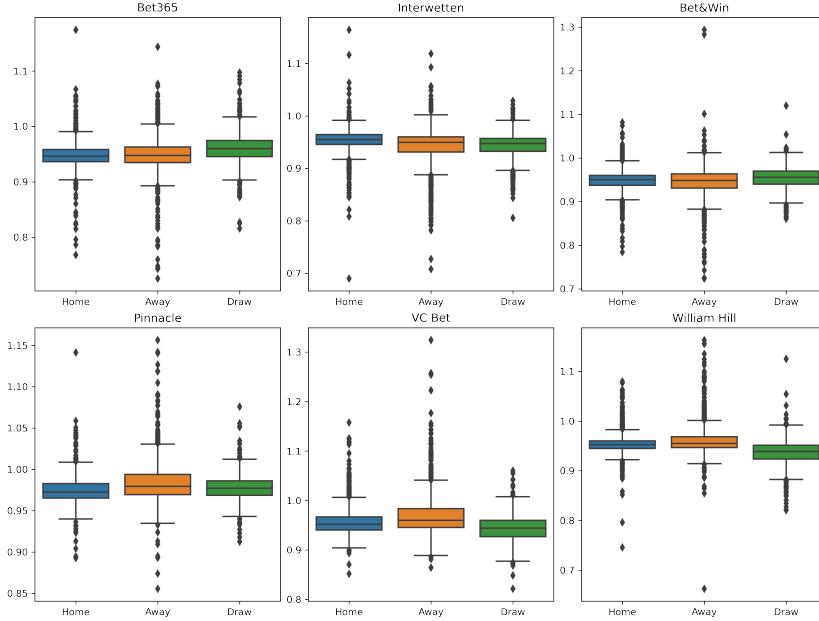


**Figure 2:** The distribution of the Kelly Index of the six European bookmakers for the 2019-2021 season.

As shown in Figure 2, the Kelly Indexes of these six bookmakers are mainly concentrated in the range of 0.9-1.0. In the setup of the predictive problem, the matches for the 2005-2021 season were divided into 3 categories with each one reflecting the confidence levels of their predictability. These were, matches with Kelly indexes greater than 1 (Type

1), matches with only one Kelly index greater than 1 (Type 2), and matches with no Kelly Index greater than 1 (Type 3).

## 3.3   Model Optimisation

**Hyperparameter tuning**

This study used the RandomizedSearchCV method from the scikit-learn library to select the most effective parameter values. All the parameter tuning values for the machine learning algorithms used in this study are shown in Table 3.

**Table 3:** Values for parameter tuning

| Algorithm | Parameter | Range of values |
|---|---|---|
| CatBoost | learning_rate | 0.01,0.02,0.03,0.04 |
| | iterations | 10,20,30,40,50,60,70,80,90 |
| | depth | 4,5,6,7,8,9,10 |
| Decision Tree | criterion | 'gini', 'entropy' |
| | max_depth | 2,4,6,8,10,12 |
| Gradient boost | min_samples_leaf | 2,5,8 |
| | min_samples_split | 3,5,7,9 |
| | max_features | Based on number of filtered features |
| | max_depth | 2,5,7,10 |
| | learning_rate | 0.1,1.0,2.0 |
| | subsample | 0.5,0.8,1 |
| k-Nearest Neighbor | n_neighbors | Based on sample size |
| Logistic regression | solver | liblinear |
| | penalty | 'l1','l2' |
| | class_weight | 1:2:1,3:3:4,4:3:3 |
| Random Forest | min_samples_leaf | 2,5,8 |
| | max_features | Based on number of filtered features |
| | max_depth | 2,5,7,10 |
| Stacking | Same as above algorithm | Same as above algorithm |
| Voting | Same as above algorithm | Same as above algorithm |

**Feature selection**

Models developed with feature selection tend to have greater explanatory power, run faster, and have a decreased risk of over-fitting. In some cases, feature selection will improve the predictive power of models. With the exception of the two ensemble algorithms Voting classifier and Stacking, feature selection was used to enhance the models.

Feature selection decisions were made using Shapley Additive exPlanations[41] (SHAP). SHAP is based on the Shapley value[58], a concept derived from cooperative games, and is a common metric for quantitatively assessing the marginal contribution of users in an equitable manner. SHAP has been implemented as a visualisation tool for model interpretation that is able to score every feature in terms of its effects on the final outcome. Efficient feature selection can be achieved by removing features that are indicated to have low impacts on predictions. For multiclassification gradient boosting models where the SHAP explainer is not available, the Recursive Feature Elimination method in the scikit-learn library was used to remove the least important features recursively selected in the dataset.

14

In the feature selection process, after generating a model that has been modelled using all features and tuned with hyperparameters, the features in the validation set with an average SHAP value of less than zero were removed and the dataset was remodelled after the low-impact features were. This loop continues until the number of features is less than 2 or no features with a SHAP value less than zero appear. After the iterative process is completed, the model with the highest prediction accuracy in the validation set is selected as the prediction model for the test set.

## 3.4 Model training and testing framework

### 3.4.1 Modelling process

This study uses an extended window approach to test the models in order to avoid the mistake of using information from future matches to predict past matches which is in contrast to a number of other studies on sports outcome prediction. Bunker and Thabtah [11] outline in detail why standard cross-validation is not suitable for evaluation of prediction models in this domain. The pattern of the training and testing sequences comprising windows of division is shown in Figure 3 while the modelling process is shown in Figure 4.

Training sets from the immediately preceding season are used to make predictions for the next season. This strategy was followed since it has been shown that it is more informative to choose the most recent matches for predictions of upcoming matches. The reason for this is that the game of football, being a high-intensity competitive sport, advances over time in terms of its macro tactics and the detail of the players' actions[55]. As the number of different types of matches varies from season to season, the size of the test set is set at 20 to 30 matches depending on the number of matches of a particular type in the current season.
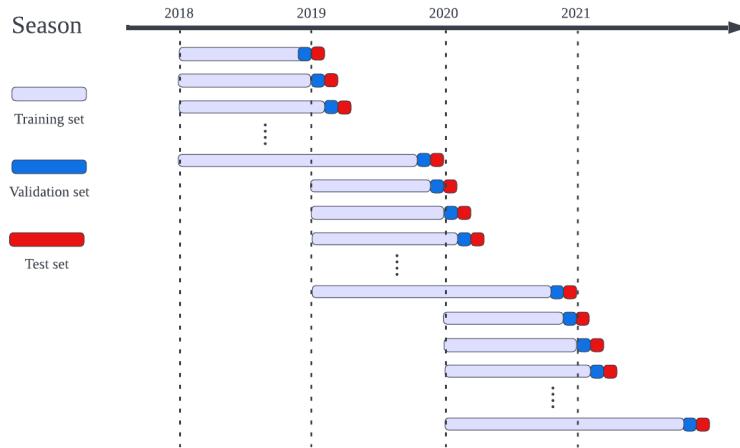


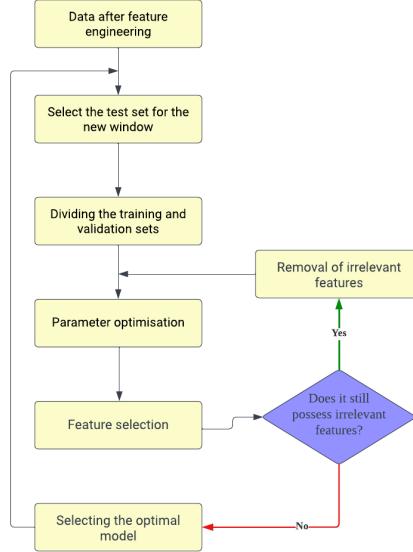**Figure 3:** Window extension strategies for predictive models

**Figure 4:** Modelling process

The validation set is used in order to select the best model during parameter tuning. When using models where parameter tuning has little impact on model quality (e.g. Voting Classifier) the validation set is not used and the part of the model that would otherwise be part of the validation set is merged with the training set. The validation set is usually set to the same size as the test set for that window.

### 3.4.2 Model evaluation

This study evaluates the predictive ability of each classifier by using four metrics. These four reported metrics are accuracy, precision, recall and F1 Score. Accuracy is used as a primary metric for ranking the performances of the algorithms which the remaining metrics are used as a reference for completion. These evaluation metrics are calculated as follows:

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN} \tag{17}$$

where the denominator part indicates the number of all samples and the numerator part indicates the number of all correctly predicted samples.

$$Precision = \frac{TP}{TP + FP} \tag{18}$$

where the denominator part indicates the number of samples for a particular result and the numerator part indicates the number of correctly predicted samples for that result species. The precision of the three classification questions is weighted by the proportion of each result in the sample using the precision of the three results.

16

**Table 4:** Number of matches according to different types and categories

| Category | Season | Result Number | | | | Draw proportion |
|----------|--------|------|------|------|-----|-----------------|
| | | Home | Away | Draw | All | |
| Type 1 | 2019 | 61 | 28 | 17 | 106 | 16.0% |
| | 2020 | 50 | 43 | 18 | 111 | 16.2% |
| | 2021 | 42 | 18 | 16 | 76 | 21.1% |
| Type 2 | 2019 | 35 | 28 | 18 | 81 | 22.2% |
| | 2020 | 36 | 29 | 22 | 87 | 25.3% |
| | 2021 | 55 | 42 | 27 | 124 | 21.8% |
| Type 3 | 2019 | 76 | 60 | 57 | 193 | 29.5% |
| | 2020 | 58 | 81 | 43 | 182 | 23.6% |
| | 2021 | 66 | 69 | 45 | 180 | 25.0% |
| All | 2019 | 172 | 116 | 92 | 380 | 24.2% |
| | 2020 | 144 | 153 | 83 | 380 | 21.9% |
| | 2021 | 163 | 129 | 88 | 380 | 23.2% |

$$Recall = \frac{TP}{TP + FN} \tag{19}$$

Theoretically, the three classification problems have the same recall and accuracy on a given result. The recall in this study uses macro-Recall without considering the weights of each result.

$$F_1 = \frac{2 \times Precision \times Recall}{Precison + Recall} \tag{20}$$

To evaluate the relative capability of the models, two baseline models were added for comparison. Both models are referenced from the Dummy Classifier class in the scikit-learn library. Both models ignore the data in the training set and make random predictions on the test set. Baseline model 1 is a random selection of one result from three results of matches as the predicted result. Baseline model 2 has a similar prediction process to baseline model 1, but predictions follow the distribution of each result in the test set.

## 4 Result and discussion

### 4.1 Match classification

This study classified 1,140 matches in the English Premier League for the three seasons of 2019-2021 into three possible categories, namely win, loss, draw. Matches with multiple bookmakers having a Kelly Index greater than 1 were categorised as Type 1. Matches with only one bookmaker having a Kelly Index greater than 1 were categorised as Type 2. Matches with no bookmaker having a Kelly Index greater than 1 were categorised as Type 3. Their precise numbers in each season are shown in Table 4. The number of Type 1 matches decreases yearly, while the number of Type 2 matches increases yearly. This may indicate that bookmakers are generally becoming progressively more conservative in setting their odds.

Nearly half of all 1140 matches belonged to the Type 3. In these matches, the bookmakers were not confident in setting odds that constituted an increased liability for them. Ultimately, Type 3 matches produced the highest percentage of draws, while Type 1 matches had the lowest percentage of draws confirming the efficacy of the Kelly Index

to reduce uncertainty to some degree. Although draws are difficult to predict, they do not account for more than a third of matches in any category. In addition, home-team wins in Type 1 matches were much higher than the number of away-team wins. However, away-team wins were higher in the Type 3 matches. Bookmakers consider both draws and home team advantages when formulating their predictions. Bookmakers are more likely to reduce the odds of a home win when considering setting odds or home teams.

Although bookmakers can adjust the odds based on their own predictions, upsets often happen in football matches. It is not always straightforward to precisely define which matches are upsets based on the available data. This study identifies matches as an upset when the study's prediction was correct and the bookmaker not only predicted incorrectly, but also set the highest odds on the actual outcome of the match. See the next sub-section for the predicted results. The percentage of upsets for each type is shown in Figure 5.
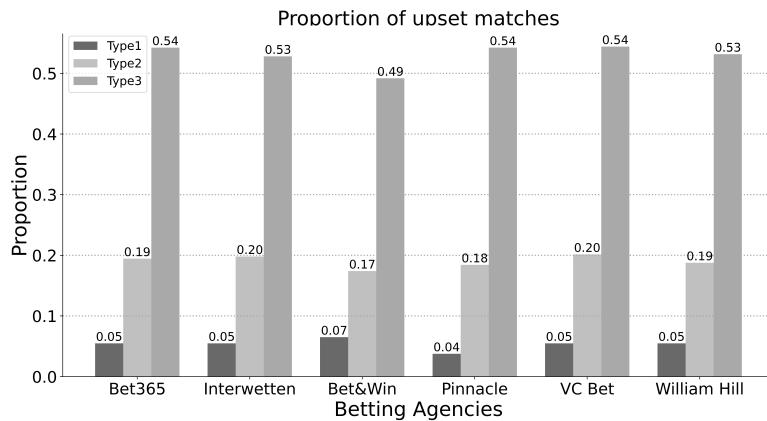


**Figure 5:** Percentage of upset in each type of match

Betting agencies are poor predictors of Type 3 matches compared to the other two types. Upsets are very unlikely in the Type 1 matches, but they become a regular occurrence in the Type 3 matches. Figure 6 aggregates the spread of the odds across all six betting agencies in this study. According to this figure, on average, bookmakers set significantly lower odds on results in Type 1 matches than in the other two categories.

## 4.2 Result prediction

The prediction accuracy, precision, recall and average ranking of each algorithm across the prediction process for Type 1, Type 2 and Type 3 matches are all shown in Table 5, Table 6 and Table 7 respectively. The algorithms are rank-ordered from best-performing onwards. The average ranking is likewise calculated by the accuracy of each algorithm across all different forecasting windows. Meanwhile, the accuracies of all the algorithms across the different Types are contrasted with the accuracies attained across all matches without taking the Kelly Index into account. These results can be seen in Table 8.

The model offering the greatest predictive performance for the Type 1 matches across all forecasting windows was produced by CatBoost, with an accuracy of 70%. This

**Table 5:** Prediction results from various machine learning algorithms in the Type 1 match

| Algorithm | Accuracy | Precision | Recall | F1_Score | Rank |
|---|---|---|---|---|---|
| CatBoost | 70.0% | 66.6% | 54.8% | 63.6% | 1.1 |
| Logistic Regression | 67.9% | 64.8% | 53.6% | 62.3% | 3.1 |
| Random Forest | 68.9% | 64.3% | 54.9% | 63.4% | 3.2 |
| Voting | 66.6% | 59.4% | 52.4% | 61.7% | 4.3 |
| KNN | 62.1% | 60.6% | 52.8% | 61.3% | 4.5 |
| Gradient Boosting | 62.5% | 59.6% | 51.7% | 60.8% | 4.8 |
| Stacking | 65.5% | 54.1% | 50.9% | 59.1% | 5.0 |
| Decision Tree | 60.4% | 56.7% | 48.9% | 58.2% | 5.1 |
| XGBoost | 63.8% | 57.1% | 51.0% | 59.6% | 5.3 |
| Baseline 1 | 47.1% | 33.5% | 31.5% | 37.4% | 8.8 |
| Baseline 2 | 30.4% | 36.9% | 29.9% | 32.1% | 9.8 |

**Table 6:** Prediction results from various machine learning algorithms in the Type 2 match

| Algorithm | Accuracy | Precision | Recall | F1_Score | Rank |
|---|---|---|---|---|---|
| CatBoost | 49.7% | 38.8% | 42.0% | 43.3% | 2.8 |
| Logistic Regression | 51.0% | 49.3% | 45.4% | 48.8% | 3.1 |
| Random Forest | 52.7% | 50.9% | 46.5% | 49.2% | 3.4 |
| Decision Tree | 45.9% | 45.3% | 42.6% | 45.6% | 3.7 |
| Stacking | 49.3% | 44.5% | 42.1% | 43.8% | 4.3 |
| Voting | 47.9% | 44.5% | 42.7% | 45.7% | 4.9 |
| XGBoost | 48.6% | 46.0% | 43.6% | 46.7% | 5.0 |
| KNN | 39.7% | 40.1% | 37.5% | 39.9% | 5.9 |
| Gradient Boosting | 40.8% | 39.7% | 37.6% | 40.0% | 6.2 |
| Baseline 1 | 40.1% | 29.0% | 32.3% | 31.7% | 6.9 |
| Baseline 2 | 31.5% | 33.5% | 31.3% | 32.1% | 8.7 |

**Table 7:** Prediction results from various machine learning algorithms in the Type 3 match

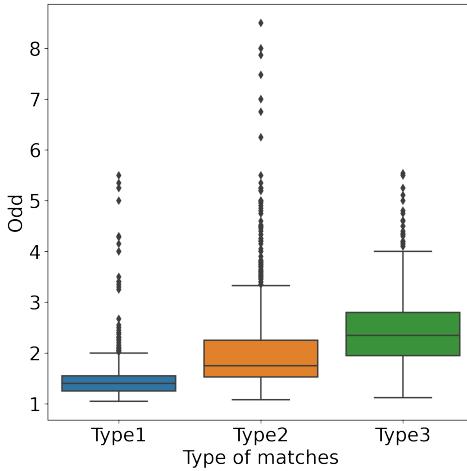| Algorithm | Accuracy | Precision | Recall | F1_Score | Rank |
|---|---|---|---|---|---|
| Decision Tree | 40.0% | 38.5% | 37.5% | 38.5% | 3.9 |
| CatBoost | 40.0% | 37.6% | 36.3% | 35.2% | 4.0 |
| Random Forest | 41.1% | 39.4% | 38.4% | 39.1% | 4.2 |
| Logistic Regression | 40.7% | 38.8% | 38.2% | 39.1% | 4.5 |
| Voting | 40.0% | 39.1% | 37.7% | 38.5% | 4.6 |
| Gradient Boosting | 37.7% | 38.0% | 37.0% | 37.8% | 5.0 |
| Stacking | 39.1% | 35.5% | 35.9% | 36.2% | 5.2 |
| XGBoost | 39.6% | 38.8% | 37.9% | 39.0% | 5.4 |
| KNN | 35.3% | 34.9% | 33.9% | 34.8% | 5.8 |
| Baseline 2 | 37.3% | 38.3% | 37.0% | 37.6% | 6.2 |
| Baseline 1 | 36.4% | 26.3% | 33.4% | 26.4% | 6.4 |

**Figure 6:** Distribution of odds set by the six bookmakers on the predicted results of this study

performance was significantly better than that of the baseline methods which scored well below 50%.

The best model in the Type 2 matches was generated by Random Forest, although Catboost ranked as the winner across all the forecasted windows. Random Forest achieved 53% which represents a significant reduction in accuracy compared to Type 1 models. However, the best-performing models under the Type 2 setting still clearly outperformed baseline predictions.
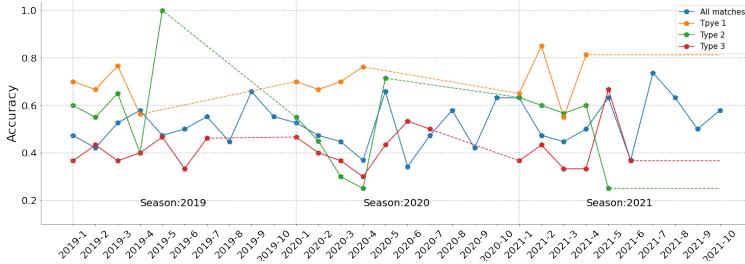
The accuracy of Type 3 accuracies continued the decrease down to 41%, which was the best accuracy registered by Random Forest. Though overall, both the Decision Tree and CatBoost performed more consistently across all the forecasting windows. While the best-performing models' efficacy significantly reduced for Type 3 matches, they still marginally contributed more value than baseline models. Unexpectedly, stacking and voting methods were consistently outperformed by other algorithms which is in contrast to the performance of these algorithms in other studies covering competitive sports.

The results from the previous tables are depicted in Figure 7 highlighting the accuracy of each of the best-performing algorithms across time for each season, and Type of matches. The dashed lines in the figure represent interpolations in cases where no matches were predicted for a given Type when matches meeting the Kelly Index criteria did not exist. The figure demonstrates that Type 1 predictions consistently outperformed other Types. . Taken altogether, the results demonstrate the efficacy of the Kelly Index to reduce an element of uncertainty from the predictive problem and thereby increase the overall predictive accuracy.

The next analysis considers the internals of the predictive moedls and uses SHAP in order to understand the impact of the key features on the best-performing algorithms based on their ranks, across each of the Types of matches. The feature importances and

**Table 8:** Prediction results from various machine learning algorithms in the all matches

| Algorithm | Accuracy | Precision | Recall | F1_Score | Rank |
|---|---|---|---|---|---|
| CatBoost | 51.9% | 45.7% | 44.7% | 45.3% | 2.6 |
| Logistic Regression | 50.6% | 44.0% | 43.5% | 45.2% | 3.2 |
| Random Forest | 52.0% | 47.7% | 44.7% | 45.7% | 3.4 |
| Stacking | 51.3% | 46.7% | 44.5% | 45.9% | 3.6 |
| Decision Tree | 50.2% | 46.6% | 44.6% | 46.2% | 3.8 |
| Voting | 49.9% | 46.7% | 44.5% | 47.0% | 4.6 |
| Gradient Boosting | 47.1% | 45.8% | 43.3% | 46.3% | 5.0 |
| XGBoost | 48.4% | 45.7% | 43.5% | 46.5% | 5.7 |
| KNN | 45.4% | 44.4% | 42.0% | 44.8% | 5.7 |
| Baseline 1 | 40.5% | 28.4% | 32.8% | 29.7% | 8.0 |
| Baseline 2 | 31.9% | 33.2% | 32.1% | 32.2% | 9.2 |



**Figure 7:** As the extended window progresses, the trend in prediction accuracy of the algorithms with the best predictive ability for different types of matches.

their effects are shown in Figure 8.

For Type 1 matches, the top most influential features are based on the team's historical match statistics based on their home-win record. These are matches where the home team's advantage exists and it is easier to determine whether the match will be one-sided based on the features of the home team. The rarity of draw results of such matches also leads to a higher prediction accuracy.

In the cases of Type 2 matches, we can see that as the prediction becomes more difficult, historical data on both teams becomes less important in and is replaced by odds data provided by bookmakers before the match.

Meanwhile, in predicting the Type 3 matches, historical data appears to be inconsequential. The features based on bookmakers' prediction estimates still have the advantage in predicting these matches. The already low accuracy accuracy of the Type 3 matches still depends on the results of the bookmakers' forecasts. The odds in this case have an almost direct influence on the prediction of uncategorised matches.
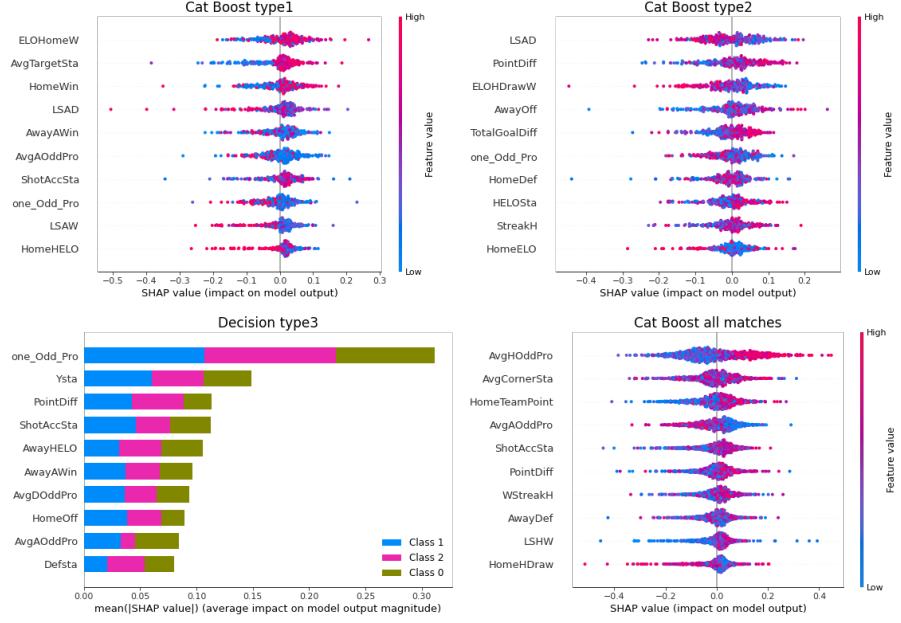
**Figure 8:** Scatterplot of SHAP values for the top 10 features of importance when being predicted for the three different types of matches and for all matches.

## 4.3    Betting returns

This section compares traditional betting methods in football betting with the returns of classified matches betting to examine if there is a way to maximise returns or minimise risk. We first begin by highlighting the distribution in confidence probabilities of the machine learning algorithms across different Types as this will form a key factor when comparing betting strategies.

Table 9 shows the distribution of confidence in the predicted outcomes for the three types of matches. Even for the most predictable matches, machine learning has difficulty predicting a match's result with confidence. Surprisingly, the proportion of these predictions that machine learning considered extremely confident was similar for the easy-to-predict matches and the hard-to-predict matches, and these matches were rarely present in both types of matches.

We next construct a simulated betting scenario in order to determine the return on investment(ROI) based on several strategies. We hypothetically invest $1 on each match and receive corresponding profits according to the odds for a correct prediction, and a loss of the investment of $1 for a wrong prediction. The ROI is calculated by dividing the investment profit by the total amount invested. A positive or negative ROI is used to determine whether an investment strategy is profitable or not.

Firstly, a set of benchmark ROIs are prepared for this study and are shown in Table 10. These are returns based on the model with the highest prediction accuracy across each of the Kelly index categories as well as the baseline which includes all matches, investing in the six bookmakers. A naive approach indicates marginal returns on two cases only,

**Table 9:** The distribution of confidence in the predicted results of each type across seasons in different ranges

| Type | Season | Range of confidence | | | | | Match number |
|------|--------|---------|---------|---------|---------|----------|--------------|
|      |        | 0.7-1.0 | 0.6-0.7 | 0.5-0.6 | 0.4-0.5 | 0.33-0.4 |              |
| Type1 | 2019 | 8.5% | 10.4% | 15.09% | 55.7% | 10.4% | 106 |
|       | 2020 | 0.0% | 0.0% | 4.50% | 42.3% | 53.2% | 111 |
|       | 2021 | 2.6% | 25.0% | 27.63% | 36.8% | 7.9% | 76 |
| Type2 | 2019 | 17.3% | 28.4% | 13.58% | 25.9% | 14.8% | 81 |
|       | 2020 | 5.8% | 10.3% | 33.33% | 41.4% | 9.2% | 87 |
|       | 2021 | 11.3% | 33.1% | 18.55% | 24.2% | 12.9% | 124 |
| Type3 | 2019 | 6.2% | 6.2% | 19.69% | 48.7% | 19.2% | 193 |
|       | 2020 | 2.2% | 8.2% | 21.98% | 47.8% | 19.8% | 182 |
|       | 2021 | 2.8% | 8.9% | 21.11% | 47.2% | 20.0% | 180 |

**Table 10:** Return on total investment by type of competition

|  | Type1 | Type2 | Type3 | Baseline |
|---|-------|-------|-------|----------|
| **Match number** | 292 | 293 | 555 | 1140 |
| **Bet365** | -1.1% | -5.9% | -3.8% | -6.7% |
| **Interwetten** | 0.9% | -4.8% | -3.4% | -5.4% |
| **Bet&Win** | -0.2% | -5.5% | -3.8% | -6.2% |
| **Pinnacle** | 1.4% | -3.5% | -1.1% | -4.1% |
| **VC Bet** | -0.7% | -6.0% | -3.3% | -6.3% |
| **William Hill** | -0.5% | -5.6% | -3.6% | -6.1% |

with higher losses across all other scenarios.

The next step in the experiments was to calibrate the investment strategy based on the confidence level of the predictive model in the outcome of each match. Thus, a bet would be placed only if the predictive model's confidence level met or exceeded a predefined threshold. Several Thresholds were chosen and examined for returns. For each Type level, the best-performing model was selected. Odds from Pinnacle were chosen.

The results of the experiments are shown in Figure 9 highlighting the ROI across time and seasons. Over time, all strategies based both on Type and models' confidence level conclude with a negative ROI. The exception to this is found in a single strategy based on Type 1 matches and a confidence threshold of 70%. The final ROI using this strategy is 17%.
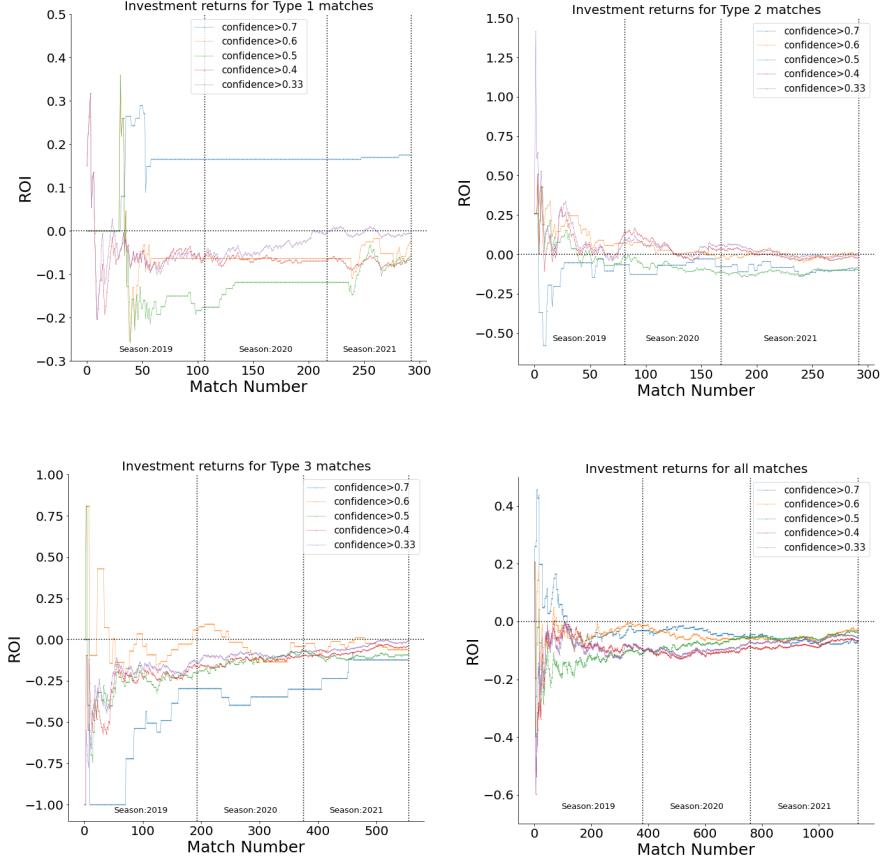
**Figure 9:** A graph of the return on investment as the match progresses, based on how confident the machine learning method is in predicting the result to determine whether to invest in the match.

## 5   Conclusions and future work

This study has considered the problem of predicting the outcomes of football matches using the Premier League match data from 2019-2021 seasons. The proposed strategy used the Kelly Index to first categorise the matches into three groups, with each one representing different levels of uncertainty, or predictability. A range of machine learning algorithms were explored for predicting the outcomes of the matches from each category in order to determine the utility of the process of decomposing the predictive problem into sub-tasks. This paper devised a range of new features previously unexplored as well as machine learning algorithms not investigated in this domain. The study found that ensemble-based algorithms outperformed all other approaches including the benchmark

approaches, while the models produced competitive results with prior works.

The paper also validates the proposed approaches by benchmarking them against bookmaker odds in order to determine with strategies are able to return a profit on investment. A method was developed that combines both the Kelly Index together with predictive confidence thresholds and investigated. The findings indicate that a strategy comprised of a combination of focusing on easy-to-predict matches that have a high predictive confidence level from machine learning models can return a profit over the long term. is a non-negligible part of the match data for machine learning. Moreover, the available features cannot explain the reasons for the match result.

# References

[1] Y. F. Alfredo and S. M. Isa. Football match prediction with tree based model classification. *International Journal of Intelligent Systems and Applications*, 11(7):20–28, 2019.

[2] A. A. Azeman, A. Mustapha, N. Razali, A. Nanthaamomphong, and M. H. Abd Wahab. Prediction of football matches results: Decision forest against neural networks. In *2021 18th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology (ECTI-CON)*, pages 1032–1035. IEEE, 2021.

[3] R. Baboota and H. Kaur. Predictive analysis and modelling football results using machine learning approach for english premier league. *International Journal of Forecasting*, 35(2):741–755, 2019.

[4] G. Baio and M. Blangiardo. Bayesian hierarchical model for the prediction of football results. *Journal of Applied Statistics*, 37(2):253–264, 2010.

[5] R. Beal, S. E. Middleton, T. J. Norman, and S. D. Ramchurn. Combining machine learning and human experts to predict match outcomes in football: A baseline model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 15447–15451, 2021.

[6] D. Berrar, P. Lopes, and W. Dubitzky. Incorporating domain knowledge in machine learning for soccer outcome prediction. *Machine learning*, 108(1):97–126, 2019.

[7] G. Bilek and E. Ulas. Predicting match outcome according to the quality of opponent in the english premier league using situational variables and team performance indicators. *International Journal of Performance Analysis in Sport*, 19(6):930–941, 2019.

[8] J. Brooks, M. Kerr, and J. Guttag. Using machine learning to draw inferences from pass location data in soccer. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 9(5):338–349, 2016.

[9] A. Brown and J. J. Reade. The wisdom of amateur crowds: Evidence from an online community of sports tipsters. *European Journal of Operational Research*, 272(3): 1073–1081, 2019.

[10] R. Bunker and T. Susnjak. The application of machine learning techniques for predicting match results in team sport: A review. *Journal of Artificial Intelligence Research*, 73:1285–1322, 2022.

[11] R. P. Bunker and F. Thabtah. A machine learning framework for sport result prediction. *Applied computing and informatics*, 15(1):27–33, 2019.

[12] L. Carloni, A. D. Angelis, G. Sansonetti, and A. Micarelli. A machine learning approach to football match result prediction. In *International Conference on Human-Computer Interaction*, pages 473–480. Springer, 2021.

[13] K. Carpenter. Match-fixing—the biggest threat to sport in the 21st century? *International sports law review*, 2(1):13–24, 2012.

[14] H. Chen. Neural network algorithm in predicting football match outcome based on player ability index. *Advances in Physical Education*, 9(4):215–222, 2019.

[15] A. C. Constantinou. Dolores: a model that predicts football match outcomes from all over the world. *Machine Learning*, 108(1):49–75, 2019.

[16] I. B. da Costa, L. B. Marinho, and C. E. S. Pires. Forecasting football results and exploiting betting markets: The case of "both teams to score". *International Journal of Forecasting*, 38(3):895–909, 2022.

[17] N. Danisik, P. Lacko, and M. Farkas. Football match prediction using players attributes. In *2018 World Symposium on Digital Intelligence for Systems and Machines (DISA)*, pages 201–206. IEEE, 2018.

[18] L. Egidi, F. Pauli, and N. Torelli. Combining historical data and bookmakers' odds in modelling football scores. *Statistical Modelling*, 18(5-6):436–459, 2018.

[19] A. ELO. The rating of chessplayers: past and present. *Arco Pub, New York*, 1978.

[20] E. Esme and M. S. Kiran. Prediction of football match outcomes based on bookmaker odds by using k-nearest neighbor algorithm. *International Journal of Machine Learning and Computing*, 8(1):26–32, 2018.

[21] G. Fialho, A. Manhães, and J. P. Teixeira. Predicting sports results with artificial intelligence–a proposal framework for soccer games. *Procedia Computer Science*, 164:131–136, 2019.

[22] A. Y. Govan, A. N. Langville, and C. D. Meyer. Offense-defense approach to ranking team sports. *Journal of Quantitative Analysis in Sports*, 5(1), 2009.

[23] A. Groll, C. Ley, G. Schauberger, and H. Van Eetvelde. A hybrid random forest to predict soccer matches in international tournaments. *Journal of quantitative analysis in sports*, 15(4):271–287, 2019.

[24] C. Herbinet. Predicting football results using machine learning techniques. *MEng thesis, Imperial College London*, 2018.

[25] D. Hill. A critical mass of corruption: Why some football leagues have more match-fixing than others. *International Journal of Sports Marketing and Sponsorship*, 11(3):38–52, 2010.

[26] Y.-C. Hsu. Using convolutional neural network and candlestick representation to predict sports match outcomes. *Applied Sciences*, 11(14):6594, 2021.

[27] O. Hubácek, G. Sourek, and F. Zelezny. Score-based soccer match outcome modeling–an experimental review. *MathSport International*, 2019.

[28] O. Hubáček, G. Šourek, and F. Železnỳ. Learning to predict soccer results from relational data with gradient boosted trees. *Machine Learning*, 108(1):29–47, 2019.

[29] J. Hucaljuk and A. Rakipović. Predicting football scores using machine learning techniques. In *2011 Proceedings of the 34th International Convention MIPRO*, pages 1623–1627. IEEE, 2011.

[30] M. Huggins. Match-fixing: A historical perspective. *The International Journal of the History of Sport*, 35(2-3):123–140, 2018.

[31] L. M. Hvattum and H. Arntzen. Using elo ratings for match result prediction in association football. *International Journal of forecasting*, 26(3):460–470, 2010.

[32] R. Ievoli, L. Palazzo, and G. Ragozini. On the use of passing network indicators to predict football outcomes. *Knowledge-Based Systems*, 222:106997, 2021.

[33] T. Inan. Using poisson model for goal prediction in european football. 2021.

[34] P. K. Jain, W. Quamer, and R. Pamula. Sports result prediction using data mining techniques in comparison with base line model. *Opsearch*, 58(1):54–70, 2021.

[35] A. Joseph, N. E. Fenton, and M. Neil. Predicting football results using bayesian nets and other machine learning techniques. *Knowledge-Based Systems*, 19(7):544–553, 2006.

[36] S. Kelly, G. Noone, and J. Perkins. Synchronization effects on probability of pulse train interception. *IEEE transactions on aerospace and electronic systems*, 32(1):213–220, 1996.

[37] İ. H. KINALIOĞLU and K. Coşkun. Prediction of uefa champions league elimination rounds winners using machine learning algorithms. *Cumhuriyet Science Journal*, 41(4):951–967, 2020.

[38] S. J. Koopman and R. Lit. Forecasting football match results in national league competitions using score-driven time series models. *International Journal of Forecasting*, 35(2):797–809, 2019.

[39] T. Li and H. Han. A high-performance basketball game forecast using magic feature extraction. In *International conference on Data Science, Medicine and Bioinformatics*, pages 35–50. Springer, 2019.

[40] J. E. Lopes, D. M. Jacobs, D. Travieso, and D. Araújo. Predicting the lateral direction of deceptive and non-deceptive penalty kicks in football from the kinematics of the kicker. *Human Movement Science*, 36:199–216, 2014.

[41] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.

[42] P. Malini and B. Qureshi. A deep learning framework for temperature forecasting. In *2022 7th International Conference on Data Science and Machine Learning Applications (CDMA)*, pages 67–72. IEEE, 2022.

[43] T. Mendes-Neves and J. Mendes-Moreira. Comparing state-of-the-art neural network ensemble methods in soccer predictions. In *International Symposium on Methodologies for Intelligent Systems*, pages 139–149. Springer, 2020.

[44] A. T. Oluwayomi, A. O. Olajide, A. A. Adetayo, A. O. Gabriel, O. A. Okunola, and O. T. Gabriel. Evaluation of team's false '9'for match winner prediction. 2022.

[45] C. Parim, M. Ş. Güneş, A. H. Büyüklü, and D. Yıldız. Prediction of match outcomes with multivariate statistical methods for the group stage in the uefa champions league. *Journal of Human Kinetics*, 79(1):197–209, 2021.

[46] J.-H. Park, C.-H. Choi, J. Yoon, and V. Girginov. How should sports match fixing be classified? *Cogent Social Sciences*, 2019.

[47] N. Pathak and H. Wadhwa. Applications of modern classification techniques to predict the outcome of odi cricket. *Procedia Computer Science*, 87:55–60, 2016.

[48] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.

[49] L. E. Peterson. K-nearest neighbor. *Scholarpedia*, 4(2):1883, 2009.

[50] D. Prasetio et al. Predicting football match results with logistic regression. In *2016 International Conference On Advanced Informatics: Concepts, Theory And Application (ICAICTA)*, pages 1–5. IEEE, 2016.

[51] J. R. Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.

[52] J. J. Reade, C. Singleton, and A. Brown. Evaluating strange forecasts: The curious case of football match scorelines. *Scottish Journal of Political Economy*, 68(2):261–285, 2021.

[53] P. Robberechts and J. Davis. Forecasting the fifa world cup–combining result-and goal-based team ability parameters. In *International Workshop on Machine Learning and Data Mining for Sports Analytics*, pages 16–30. Springer, 2018.

[54] F. Rodrigues and Â. Pinto. Prediction of football match results with machine learning. *Procedia Computer Science*, 204:463–470, 2022.

[55] N. Smith. Nature as accumulation strategy. *Socialist register*, 43, 2007.

[56] E. Štrumbelj. On determining probability forecasts from betting odds. *International journal of forecasting*, 30(4):934–943, 2014.

[57] E. O. Thorp. Portfolio choice and the kelly criterion. In *Stochastic optimization models in finance*, pages 599–619. Elsevier, 1975.

[58] E. Winter. The shapley value. *Handbook of game theory with economic applications*, 3: 2025–2054, 2002.

[59] S. Wold, K. Esbensen, and P. Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987.

[60] F. Wunderlich and D. Memmert. The betting odds rating system: Using soccer forecasts to forecast soccer. *PloS one*, 13(6):e0198668, 2018.

[61] A. Yekhande, Y. D. Kumar, K. Sanjay, and A. Phalke. Predicting premier league match odds using machine learning. *Journal homepage: www. ijrpr. com ISSN*, 2582: 7421.

[62] B. F. YILDIZ. Applying decision tree techniques to classify european football teams. *Journal of Soft Computing and Artificial Intelligence*, 1(2):86–91, 2020.

[63] A. Zeileis, C. Leitner, and K. Hornik. Probabilistic forecasts for the 2018 fifa world cup based on the bookmaker consensus model. Technical report, working papers in economics and statistics, 2018.

[64] Q. Zhang, X. Zhang, H. Hu, C. Li, Y. Lin, and R. Ma. Sports match prediction model for training and exercise using attention-based lstm network. *Digital Communications and Networks*, 2021.

# A  Additional Accuracy Results

The best three performing algorithms from the above tables are selected and their detailed accuracies are depicted in Figure 10 which shows their relative accuracies across time. A pattern can be detected where the divergence in accuracies becomes increasingly pronounced between different algorithms as the Type of matches being predicted increases, indicating a higher uncertainty.
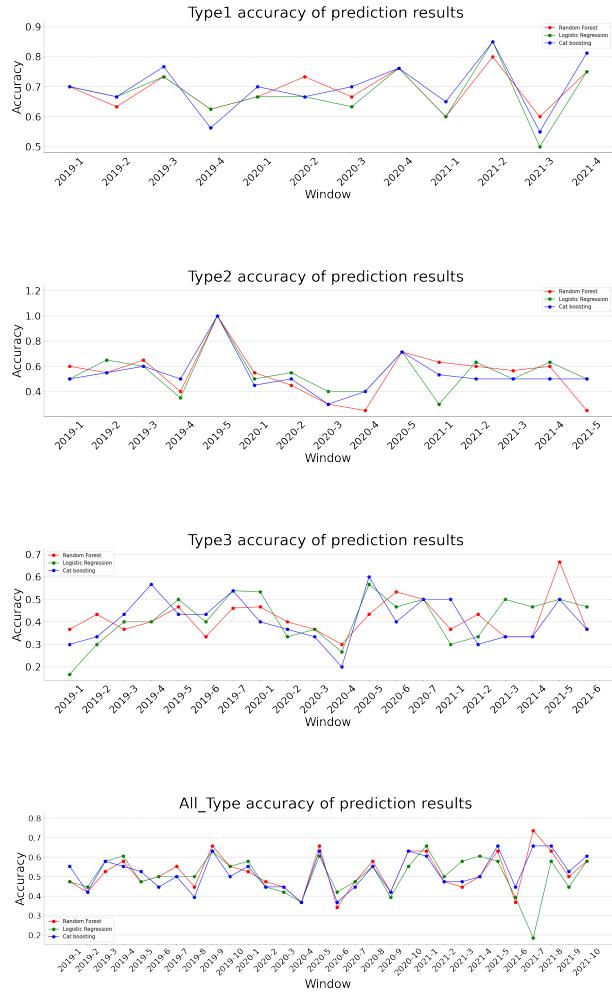
**Figure 10:** Prediction accuracy per window for three machine learning algorithms that produce excellent prediction results

Figure 11 shows whether the predictions of this study are the same as those of the bookmakers. The vertical axis shows the proportion of all matches where the two predictions are the same. Both are almost identical in their predictions for Type 1 matches.

However, as the difficulty of predicting matches increases, the difference between the two increases, although they are similar in terms of prediction accuracy. There is no way of knowing what algorithms and features the bookmakers used in their predictions. However, in terms of results, a model using only historical match result data can achieve the predictive power of a bookmaker's forecasting model. This also demonstrates that it is practical and meaningful to use the Kelly Index to classify football matches.

There is no one fixed way to predict the result of a match. There is no pattern between the features used when predicting from different windows. The available features do not explain how the results of football matches are generated. This is one of the reasons why pre-match prediction has hit an upper accuracy ceiling.

The high fluctuation in accuracy over a short period of time is what makes the game of football so attractive as a high-intensity competition. This explains the necessity of using the extended window method or the rolling window method when predicting competitive events that occur in the same chronological order as a football match. When predicting matches that occur over a short period of time, such as a single season, the changes in competitive strategy, team personnel changes and changes in player mentality that arise from advancing time are minimised. It is possible for a machine learning algorithm to produce a model that is overfitting for the entire league's history of matches but predicts a particular season's matches well. However, such models have no practical implications for future matches.
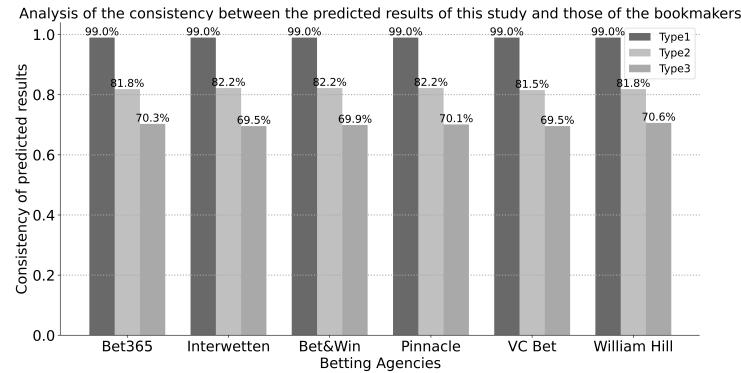


**Figure 11:** The consistency of the different bookmakers' predictions in terms of results with those of this study.