

Introduction

The Premier League is one of the most widely followed football competitions globally, captivating audiences with its intense matches and huge fan bases. With the huge following of the league, the ability to predict match outcomes has garnered considerable attention. Using Machine Learning methods, we will work to identify the factors that determine success in Premier League games.

This project revolves around the application of the Machine Learning technique known as the Random Forest Classifier. By utilising key metrics such as team performance, match scheduling, venue, opponents, and tactical formations, we aim to make informed predictions about match winners.

This report aims to provide a comprehensive overview of our methodology and the steps we go through along the way. We will begin by establishing a clear understanding of the problem, including exploration of the data points and their relevance. Then, in order to ensure the dataset's quality and applicability, we will explain where it came from and how it was created.

In this project, we are not only attempting to forecast Premier League match winners, but we also hope to learn more about the complex factors that influence these results. By mixing our passion for football with machine learning techniques, we hope to explain the subtle aspects of the "beautiful game."

Problem Formulation

To approach the task of predicting Premier League game winners as a Machine Learning problem, we first need to define the key components

We are formulating this as a supervised binary classification problem. Each instance in our dataset represents a Premier League match. The goal is to predict whether a specific team will win or not based on a set of features and compare it to the actual results.

This is a supervised learning task, since the model learns from labelled examples, where it's provided with both the input features and the corresponding target labels (in this case, whether a team wins or not).

Data Points, Features, and Labels:

Data Points: Each data point in our dataset corresponds to a single Premier League match. It contains information about both competing teams, such as their performance statistics, the date and time of the match, the round of play, the venue, the opponent, and the tactical

formation utilised. The data is categorical eg. team names and also numerical eg. match date and time and round numbers.

Features:

These are the attributes or characteristics of a match that we use to make predictions. Our features encompass a range of statistics and match details, including team performance metrics, match timing, round number, venue, opponent identity, and tactical formations employed by the teams.

Labels:

The labels indicate the outcome of each match. Specifically, it denotes whether the team we are interested in predicting (referred to as the "home" team) emerged as the winner (label = 1) or not (label = 0) for simplicity reasons a draw will also be labelled as 0.

Source of the Dataset:

The data we will use comes from FBref.com which is a reputable sports statistics source, including official Premier League records. It consists of a thorough collection of match data spanning numerous Premier League games, giving our machine learning model a large and varied set of cases to learn from. The data needs a lot of cleaning before being usable since we need different data from multiple datasets on the website. We have created our own file matches.csv containing the cleaned up data.

Methods

Number of Datapoints and Dataset Description:

The dataset contains a total of 1389 records. It encompasses a comprehensive collection of Premier League match data spanning two different seasons. Each record represents a single match and includes a range of features, such as team performance metrics, match details such as date, time, round, venue information, opponent identity, and tactical formations used by the teams.

Data Preprocessing:

Prior to utilising the dataset for modelling, we conducted necessary data preprocessing steps. We filtered out the needed features from different sets to form our "matches.csv" file. We also handled missing values, ensuring data consistency.

Feature Selection Process:

The feature selection process involved evaluating the relevance and importance of each feature for predicting Premier League match outcomes. As football fans we have selected features that have a visibly significant impact on match results, such as team performance metrics and venue details, and prioritised them.

Choice of ML Models:

We chose the Random Forest Classifier as the main machine learning model for our project. This is because of the model's capacity to handle complicated relationships in the data, robustness to noisy input, and ability to identify nonlinear patterns. The Random Forest Classifier was the most viable contender for the prediction job considering the variety and complexity of Premier League match data.

We chose Linear Regression as the secondary model to provide a different perspective on the problem. While it may seem counterintuitive to apply a linear model to a problem with complex, non-linear relationships, Linear Regression can still offer valuable insights as it provides a clear and interpretable view of the relationship between individual features and the outcome variable.

This can help identify which specific factors have a direct and linear impact on match outcomes. For instance, it may reveal if metrics like team performance or match timing have a more straightforward, linear influence.

By incorporating both Random Forest and Linear Regression, we aim to gain a comprehensive understanding of the predictive power of different models and potentially uncover unique insights into the factors that affect Premier League match outcomes.

Choice of Loss Function

For this particular project we have decided to use both Cross Entropy and Mean Squared Error as the loss functions. For our Random Forest Classifier, we utilized Cross Entropy as the primary loss function since it is well-suited for multi-class classification tasks like predicting Premier League match outcomes. It effectively penalises misclassifications which is important for accurately predicting results.

We chose Mean Squared Error as the principal loss function in the context of Linear Regression because while our primary goal is classification, including MSE is a useful indicator for evaluating regression performance. This option is especially useful for analysis using match statistics or performance measures, since it gives a clear metric of prediction errors, improving our ability to analyse and assess the regression-based approach's outcomes.

Model Validation Process

To evaluate the model's performance, we will use a splitting approach. The dataset will be divided into three subsets. The training set, used for training the model, which will be 60% of the data, the validation set, 20% of the data, to finetune the hyperparameters and model evaluation and the testing set also 20% of the data, reserved for final evaluation and simulating real scenarios with unseen data.

This approach achieves a balance between the necessity for sufficient training data and the need for different sets for fine-tuning and final evaluation, which results in a reliable and adaptable predictor of Premier League match outcomes.

Results

In this analysis, two machine learning models were evaluated: Random Forest and Linear Regression.

Random Forest:

- Validation Accuracy using Entropy: 0.64
- Test Accuracy using Entropy: 0.57
- Precision when predicting wins: 38%

Linear Regression:

- Training Set MSE: 0.24
- Test Set Accuracy: 0.75
- Test Set Precision: 0.33
- Test Set MSE: 0.21
- Validation Set MSE: 0.25

Final Chosen Method:

After comparing the results, it appears that the Linear Regression model performs better in terms of accuracy and precision. It has a higher test set accuracy (75%) compared to the Random Forest model (57%). Additionally, the precision of the Linear Regression model for predicting wins is 33%, which is higher than the Random Forest model's precision of 38%.

Test Error of the Final Chosen Method:

The test error for the final chosen method, which is Linear Regression, is measured by the Mean Squared Error (MSE) and is found to be 0.21.

In conclusion, based on the analysis and comparison of both models, the Linear Regression model is chosen as the final method for predicting match outcomes. It demonstrates better performance in terms of accuracy, precision, and MSE on the test set. The test error, measured by MSE, is 0.21 for the Linear Regression model. This indicates that, on average, the predicted outcomes are within a reasonable range of the actual outcomes.

Conclusion

In conclusion, this report evaluated how well two machine learning models—Linear Regression and Random Forest—performed at predicting match results for a particular dataset. To choose the most useful model, the analysis compared training and validation errors.

The Linear Regression model demonstrated promising results, with a lower test set Mean Squared Error (MSE) of 0.21, indicating that its predictions were closer to the actual outcomes. However, it's important to note that while the Linear Regression model showed good results, there may still be room for improvement.

The possibility of overfitting should be considered, as the training error was noticeably smaller than the validation error. This suggests that the model might be too complex for the given data, and regularisation techniques or alternative models could be explored to address this issue.

To enhance the performance of the ML method, several strategies can be considered like leveraging additional features of data points, experimenting with different models, and exploring alternative loss functions for training. Additionally, collecting a larger and more diverse dataset could provide the models with a broader range of patterns to learn from, potentially leading to more accurate predictions.

Overall, while the Linear Regression model shows promise, there is still room for refinement. By addressing potential overfitting and exploring alternative modelling techniques, we can work towards a more accurate prediction system for match outcomes. Additionally, the use of more data can be important for achieving even more accurate and reliable results in the future.

Bibliography/References

https://en.wikipedia.org/wiki/Random_forest

<https://www.geeksforgeeks.org/python-mean-squared-error/>

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

<https://towardsdatascience.com/cross-entropy-loss-function-f38c4ec8643e>

https://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html

<https://madewithml.com/courses/foundations/pandas/>