# Research on predicting football matches based on handicap data and BPNN

**Jiahao Hu**

Dalian University of Technology, Dalian, Liaoning 116620, China

hu1124666055@163.com

**Abstract.** Football is one of the most influential sports in the world, and billions of people around the globe pay much attention to the football matches. With the growing popularity of football and the continuous development of the football betting industry, the prediction of the outcomes of football matches has become a hot topic in the commercial operations of sports especially footballs in recent years. It is also an important subject of academic research. In this paper, we develop a football match result prediction model based on the back propagation neural network. We take the German Bundesliga competitions as the research object in this paper. In addition to utilizing historical statistic data and team attributes from previous matches, we also incorporate a new dataset, known as handicap data, which refers to the odds data of the football matches, as the input layer of the BPNN (back propagation neural networks) for prediction. We also innovatively use varying numbers of hidden nodes, which greatly improves the prediction accuracy and stability of the model. Experimental results indicate that the average prediction accuracy of this football match prediction model is around 57.2%, with the highest prediction accuracy reaching 59.8% and the lowest prediction accuracy at 53.8%. The prediction model demonstrates relative stability, with no significant fluctuations in prediction accuracy.

**Keywords:** Football Match Win-Draw-Lose Prediction, Back Propagation Neural Networks, Handicap Data, Machine Learning.

## 1. Introduction

Football has a far-reaching influence worldwide, engaging numerous individuals involved in football-related work. The football industry has experienced rapid development and is currently the outdoor sport with the highest production value, greatest influence, and highest public attention in the world. In the field of sports, football has achieved a leading position far ahead of other sports. According to relevant data statistics, football, hailed as the "17th largest economy in the world," has an annual industry production value accounting for 43% of the total output value of the sports industry, reaching $500 billion, surpassing the Gross Domestic Product (GDP) of many developed countries and regions. Undoubtedly, it is the world's number one sport, far surpassing other sports such as basketball, golf, baseball, and F1 racing. According to Fédération Internationale de Football Association (FIFA) statistics, as of July 2016, there were over 200 million athletes from 1.5 million football teams worldwide participating in various football competitions. In addition, the number of individuals working in football-related jobs has reached 30 million. It is estimated that those involved in football and related occupations account for approximately 3% of the world's total population.[1] Since the 21st century, people's living

standards have greatly improved, and sports-related industries have also advanced rapidly and steadily. Sports lottery, as one of the sports-related industries, has also gained wide support and attention from various countries. Over the past few decades, this sport has experienced steady development, and as a result, the business revenue in the related betting industry has significantly grown. However, due to the high unpredictability of football matches, it is a highly complex task to achieve stable and accurate predictions of football match outcomes. Without analyzing match data and statistical information, the success rate of betting on football matches is typically low.

However, with the rapid development of computer technology, the emergence of machine learning has opened up new opportunities for addressing this issue. Machine learning, particularly the widespread adoption of neural networks, has recently reached its pinnacle. The results of football matches can be predicted based on various machine learning algorithms and various characteristic data, such as the strength of both teams, their rankings, and their performance in home and away matches. In this paper, we utilize BP neural networks, multiple match feature data, and handicap data to predict football match results. And we also innovatively use varying numbers of hidden nodes of the BP neural network, which greatly improves the prediction accuracy and stability of the model.

The remainder of this paper is organized as follows. Section 2 provides a brief introduction of football betting rules. And we also provide a brief summary of previous related work in this section. Section 3 mainly elaborates on the work we have conducted, including data acquisition, data processing, and data mining and predicting. It mainly describes the methods used in our research and the process of building the BP neural network. Section 4 discusses and interprets the prediction results of the BP neural network model. We also explore the limitations of the prediction model and the areas where our work can be improved in this section. Finally, in Section 5, we summarize the manuscript and provide prospects for future work.

## 2. Background and literature review

### 2.1. Background

Back Propagation (BP) neural network is one of the most popular neural network architectures, not only because of its implementation complexity, but also because of its high productivity and efficiency. It is a multi-layer feedforward neural network trained by the error backpropagation algorithm. The BP neural network generally consists of three parts: input layer, hidden layer, and output layer, each of which is composed of several neurons, that is, several feature data. The hidden layer may also include multiple layers of neurons, and its quantity is one of the decisive factors affecting the efficiency and accuracy of the neural network. The learning of the BP neural network is a supervised learning process. The model continuously adjusts the parameters in the neural network by learning and backpropagating the error to correct the error, ultimately achieving or approaching the expected mapping relationship between input and output.

Sports betting is a practice that involves predicting the outcome of a sporting event, placing a certain amount of money, and potentially earning profits if the prediction is correct. However, if the prediction is incorrect, the individual stands to lose the entire amount of money placed on the bet. In the case of football matches, there are typically three possible outcomes: a win, a draw, or a loss. This makes the probability of randomly guessing the correct result only 1/3, rendering random betting highly likely to result in financial losses. Therefore, it is important to conduct thorough analysis and research to improve the accuracy of predictions and increase the chances of making a profit in sports betting. By analyzing match data and statistical information, individuals can make informed decisions and potentially increase their chances of success in sports betting.

### 2.2. Literature review

In recent years, the application of machine learning in the sports field has significantly increased to help predict sports match results and make decisions. In 2019, Rahul Baboota, Harleen Kaur, and others used data from the English Premier League to create a feature set for determining the most important factors

for predicting football match results. They also used machine learning to create a highly accurate predictive system. Their best model achieved a performance of 0.2156 on the Ranked Probability Score (RPS) metric for game weeks in the English Premier League, but this performance still fell short of bookmaker predictions.[2] In 2022, Fátima Rodrigues, Ângelo Pinto, and others focused on the profitability of football betting and used data from 1900 matches in five seasons (2013/2014 to 2018/2019) of the English Premier League. They conducted correlation analysis on 31 variables and compared the predictive accuracy of several machine learning algorithms, such as K-nearest Neighbors (KNN), Random Forest (RF), and Support Vector Machine (SVM), using the selected highly correlated variables. The RF model achieved the highest prediction accuracy of 65.26%.[3]

Research on using BP neural networks for predicting football match results is also a hot topic in the academic community. For example, in 2018, Fu Yu utilized neural networks and feature data such as shots on target, free kicks, and corner kicks from 24 matches of FC Barcelona during the 2015-2016 and 2016-2017 seasons to predict match outcomes. The dataset was divided into training and testing sets in a 3:1 ratio. His Multi-Layer Perceptron neural network model consisted of 4 input layers and 3 output layers, with a fixed hidden layer of 6 neurons. He implemented the research using the R programming language and assessed the prediction accuracy through a confusion matrix, which indicated reliable overall prediction accuracy.[4] Another example is the research conducted by Yiqi Wang, Hongrun Zhao, and others in 2021. They employed social network analysis and BP neural networks to predict match outcomes based on player behaviors on the football field. Their study considered passing success rate, participation in passing, player interaction, and core player dependence as passing factors. They also included attacking-defending strategy, ball possession rate, successful interception rate, and defensive interception rate as attacking-defending factors. Furthermore, they incorporated team injury occurrences, coach performance, and home/away advantage as player status factors. By utilizing these three sets of factors, their BPNN achieved a prediction accuracy of 81.8%.[5] In the same year, Yi Huang utilized BP neural networks to identify the most influential factors on match results and used these selected feature factors for predicting football match outcomes. The research achieved a prediction accuracy of 75%.[6]

There is also a case study that utilized handicap data for football match result prediction. In 2016, Xichen Ao and others used Python tools to scrape odds data from some matches in the English League One, and they built a model using principal component analysis and logistic regression. The model's prediction accuracy was verified to be above 70.1% [7].

It is worth mentioning that the researchers mentioned above used different feature data and different datasets, resulting in significant variations in prediction accuracy.

## 3. Method

### 3.1. Data Description

The data used in this study consists of 612 records from a single season (2021/2022) of the German Bundesliga, which includes a total of 306 football matches. We extracted eight types of feature data related to football matches from the "All Football" app:

(1) Possession rate of both teams.
(2) Number of shots by both teams.
(3) Number of shots on target by both teams.
(4) Number of attacks by both teams.
(5) Number of dangerous attacks by both teams.
(6) Number of corner kicks by both teams.
(7) Home and away match indicators.
(8) Team strength index of both teams.

The team strength index is a comprehensive indicator calculated by the "All Football" app, taking into account recent team performance, goals scored and conceded, and market value.

For obtaining handicap data, similar to [7], we used a web scraping tool, Python, to extract odds data from the Bet365 company for the 2021/2022 German Bundesliga season, sourced from the website https://odds.500.com/.

(9) Win-draw-loss odds data from Bet365 company.

For missing values, we manually supplemented the data. In the end, we obtained a total of 612 original data records, as shown in Fig. 1.



**Figure 1.** Football match raw data.

### 3.2. Data Processing

In Section 3.1, we normalized the nine types of feature data. Among them, (1), (7), and (8) were already normalized and can be directly used as input to the neural network. As for (2), (3), (4), (5), (6), and (9), we grouped the data pairwise for both teams in a football match and applied the following formula for normalization:

$$p = \frac{data1}{data1 + data2} \quad (data1 \neq 0 \text{ 或} data2 \neq 0)$$

$$p = 0.5 \quad (data1 = data2 = 0)$$

Among the feature data, data1 represents the features of the home team, and data2 represents the features of the opposing team. Both data1 and data2 are normalized to the range [0, 1], which facilitates their use as input for the neural network.

Regarding the representation of football match results, similar to [4], we adopted a categorical label approach. We represent the outcomes of "win," "draw," and "loss" using category labels as follows: "win" is represented as "1 0 0," "draw" is represented as "0 1 0," and "loss" is represented as "0 0 1." These labels serve as the output of the neural network.

In the end, we obtained 612 records of normalized data, as shown in Fig. 2.



**Figure 2.** Football match processed data.

*3.3. Data mining and predicting*

In Section 3.2, we used the normalized feature data obtained in Section 3.1 as the input layer of the neural network. The corresponding category labels for different football match results were used as the output layer of the neural network. The dataset was divided into a training set and a test set in a 7:3 ratio, with 428 records for training and 184 records for testing.

As for the activation function of the BP neural network, we chose the "tansig" function. The "tansig" function allows the output to maintain a nonlinear monotonic relationship with the input, which is suitable for the gradient calculation in the BP network and exhibits good fault tolerance. The function expression is as follows:

$$y = \frac{2}{1+e^{-2x}} - 1$$

In the case of selecting the number of neurons in the hidden layer, there is no standard method so far. However, a rough range can be determined using the following formula:

$$h = \sqrt{m+n} + a \quad (1 \leq a \leq 10, a\epsilon Z)$$

Based on the range of hidden neurons mentioned earlier, we iterate from 4 to 13 hidden nodes and calculate the Mean Squared Error (MSE) for each training set. The hidden layer node count that yields the minimum MSE is selected as the optimal configuration for testing the model on the test set.

We utilize the trained model to make predictions, similar to the approach in [4]. We generate a confusion matrix using the confusionMat() function to visualize the experimental results and calculate the prediction accuracy based on the confusion matrix.

The pseudocode introduction of my MATLAB code is provided below:
1. Read data from 'input.xlsx' file and split it into input and output data.
2. Define the number of training and testing samples.
3. Split the data into training and testing sets.
4. Initialize the mean square error (MSE) to a high value.
5. Define the activation and training functions for the neural network.
6. For each hidden layer neuron number between sqrt(inputnum+outputnum)+1 and sqrt(inputnum+outputnum)+10:
      a. Create a neural network with the defined architecture.
      b. Set the training parameters.
      c. Train the neural network.
      d. Compute the mean square error of the training set.
      e. If the MSE is lower than the previous value, keep the current number of hidden neurons as the best.
7. Create the final neural network with the best number of hidden neurons.
8. Set the training parameters.
9. Train the neural network using the training set.
10. Use the trained neural network to predict the output of the testing set.
11. Convert the predicted output to class labels.
12. Convert the true output to class labels.
13. Compute the confusion matrix.
14. Compute the accuracy of the model.
15. Print the confusion matrix and the accuracy.

## 4. Discussion

The proposed BP neural network prediction model in this study achieves an average prediction accuracy of approximately 57.2%. The highest prediction accuracy is around 59.8%, while the lowest prediction accuracy is around 53.8%. The prediction model is relatively stable, with no significant fluctuations in accuracy. Furthermore, the prediction accuracy significantly surpasses the 33.3% probability of randomly predicting football match results without considering match and analysis data. The main

contribution of this study lies in incorporating a new dataset - odds data - as a feature for predicting football match results. Compared to [4], this study also innovatively uses varying numbers of hidden nodes, which greatly improves the prediction accuracy and stability of the model.

During practical predictions, since we cannot obtain data such as ball possession, shot attempts, and corner kicks for upcoming matches in advance, we can substitute these missing data with the average values of the team's previous 10 matches. This allows us to make predictions for actual football match results.

However, this study also has several limitations. Firstly, due to limited data collection capabilities, our dataset is still relatively small, which directly affects the final prediction accuracy. Secondly, we have a limited number of feature variables and have not performed correlation testing on the feature data. Future work can involve increasing the variety of feature data and conducting correlation testing to select relevant features and improve prediction accuracy. Lastly, in building the neural network model, we did not conduct in-depth research and exploration of parameters such as learning rate. We only used commonly used values for parameter tuning and did not depict the learning rate's variation curve with respect to loss to select the optimal parameters for tuning.

## 5. Conclusion

This paper describes the development of a football match result prediction model using BP neural networks. It utilizes a dataset of 612 records from 306 matches of the 2021/2022 German Bundesliga. The dataset consists of two different sources: football match-related data and team strength indices extracted from the "All Football" app, and odds data obtained from Bet365. A total of 9 feature variables are collected as inputs for the neural network, and categorical labels representing "win," "draw," and "loss" are used as the network's outputs. By varying the number of hidden layers, the best configuration is determined through validation on the test set, achieving an average prediction accuracy of approximately 57.2%, with a highest accuracy of around 59.8% and a lowest accuracy of around 53.8%. These results significantly outperform the 33.3% prediction probability achieved without considering match and analysis data.

The model exhibits great potential for further development. As future work, expanding the dataset and increasing the variety of feature variables can be considered. Additional factors such as passing attempts, passing accuracy, tackles, clearances, and the physiological and psychological states of players can be included to enhance the model's prediction accuracy for football match results. Moreover, integrating the prediction model with other machine learning models into a decision support system can help evaluate betting risks and provide gamblers with insights into the risks involved. This integration can create more value for the sports betting industry and have a greater impact on the world of football.

## References

[1]     Xia Fei. (2017). Application research of BP neural network in predicting football match outcomes. Chongqing Normal University Press, Chongqing.

[2]     Rahul Baboota, Harleen Kaur. (2019). Predictive analysis and modelling football results using machine learning approach for English premier league. International Journal of Forecasting, 35:741-755.

[3]     Fátima Rodrigues, Ângelo Pinto. (2022). Prediction of football match results with machine learning. Procedia Computer Science, 204:463-470.

[4]     Fu Yu. (2018). Neural network for predicting football match outcomes. Science and Technology Review, 23:238.

[5]     Wang Yiqi, Zhao Hongrun, Zhao Hongwen, Xu Xichen. (2021). Research on football match prediction based on social network analysis and BP neural network. Journal of Weifang Engineering Vocational College, 34(3):104-108.

[6]     Huang Yi. (2021). Prediction of football match results by using neural network. Microcomputer Applications, 37(11):137-140.

[7]   Ao Xiqin, Gong Yujie, Li Jian. (2016). Prediction of football match outcomes based on odds data. Journal of Chongqing Technol Business Univ (Nat Sci Ed), 33(6):85-89.