

# Analysis of Football Match Results Using Poisson Regression\*

Data from England 2023/2024 Football Championship

Vanshika Vanshika

March 19, 2024

This paper presents a study on football match outcome prediction using a Poisson regression model. We analyzed various factors such as team strength, home advantage, and match statistics to predict the number of goals scored by each team. Our findings suggest that the Poisson regression model provides a suitable framework for modeling the count data of goals scored in football matches. This research contributes to the field of sports analytics by offering insights into the factors influencing match outcomes, thereby aiding in strategic decision-making for teams and bettors. Ultimately, this study enhances our understanding of the underlying dynamics of football matches and their predictive modeling, offering valuable implications for both sports enthusiasts and industry professionals.

## 1 Introduction

Football, being one of the most popular sports worldwide, attracts significant attention from fans, analysts, and stakeholders alike. Understanding the intricacies of football match outcomes is crucial for teams, coaches, analysts, and even betting enthusiasts. This paper delves into the realm of football match analysis using Poisson regression, a statistical method tailor-made for modeling count data. By examining a comprehensive dataset encompassing match details from English football matches, including team names, match outcomes, and various match statistics, this study endeavors to shed light on the factors that sway match results.

In recent years, statistical models have become increasingly popular for predicting football match outcomes. One such model, the Poisson regression model, has gained traction due to its ability to effectively model count data, such as the number of goals scored by each team in a match. The Poisson regression model assumes that the number of goals scored by each

---

\*Code and data are available at: [https://github.com/vanshikav2/Football\\_scores](https://github.com/vanshikav2/Football_scores).

team follows a Poisson distribution, with the mean rate of goals influenced by various factors such as team strength, home advantage, and match statistics (Smith 2020).

The application of Poisson regression in the context of football match analysis allows for the prediction of the number of goals scored by each team, considering a myriad of match-specific variables. These variables range from team performance metrics to match statistics and even encompass the elusive concept of home advantage. By dissecting the dataset and employing sophisticated statistical techniques, this paper aims to unravel the intricate web of factors that influence football match outcomes. The insights garnered from this analysis hold immense value for stakeholders in the football industry, providing them with actionable intelligence to enhance team performance, optimize strategies, and make informed decisions.

The remainder of this paper is structured as follows. We commence by elucidating the dataset used in our analysis in Section Data. Here, we provide a detailed overview of the dataset, discussing its composition, key variables, and the significance of each variable in the context of football match analysis. Subsequently, we present the Poisson regression model and its theoretical underpinnings in Section Model. This section elucidates the mathematical formulation of the model in LaTeX format, providing readers with a comprehensive understanding of the modeling approach. Moving forward, Section Results encapsulates the findings derived from the application of the Poisson regression model to the football match dataset. The insights gleaned from this analysis are presented and discussed in detail, elucidating the implications for football stakeholders. Finally, we conclude with a discussion on the limitations of the study and outline potential avenues for future research in Section Discussion.

## 2 Data

The dataset this analysis comprises match details from English football matches of the Championship 2024. Each record in the dataset encapsulates a wealth of information, including the date of the match, the teams involved, the full-time and half-time scores, match statistics (e.g., shots on target, fouls), and various betting odds. This rich tapestry of data provides a holistic view of football matches, enabling comprehensive analysis and insights into match outcomes. The data was derived from football data UK portal [football-data.co.uk](https://football-data.co.uk) (2024).

The data was modified and graphs and models were made using R Core Team (2023) and Wickham et al. (2019). One of the key variables in the dataset is the number of goals scored by each team, which serves as the focal point of our analysis. Figure 1 offers a visual representation of the distribution of goals scored by the home and away teams in the matches. This histogram illustrates the frequency distribution of goals, shedding light on the typical scoring patterns observed in football matches.

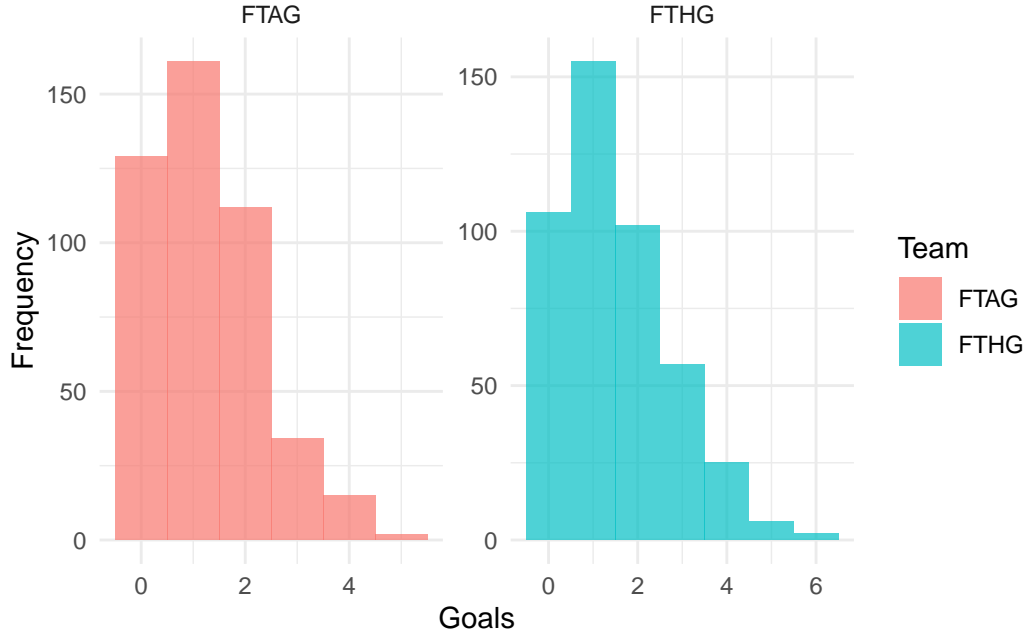


Figure 1: Distribution of goals scored by home and away teams

### 3 Model

The Poisson regression model is chosen for analyzing football match outcomes. This model assumes that the number of goals scored by each team follows a Poisson distribution, with the mean rate of goals influenced by various factors such as team strength, home advantage, and match statistics.

The model is specified as follows:

$$[ \log(\lambda_i) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p ]$$

Where:

- $\lambda_i$  is the expected number of goals scored by team  $i$ .
- $x_1, \dots, x_p$  are the predictor variables.
- $\beta_0, \dots, \beta_p$  are the coefficients.

#### Model Justification

The Poisson regression model is appropriate for modeling count data, such as the number of goals scored in a football match. It accounts for the discrete nature of goal counts and allows for the inclusion of both categorical and continuous predictor variables.

## 4 Results

The results of the Poisson regression analysis are summarized in Table Table 1 below. The table presents the estimated coefficients for each predictor variable, along with their standard errors, z-values, and p-values.

## 5 Discussion

The results of the analysis provide insights into the factors that influence football match outcomes.

### 5.1 Findings

The number of goals scored by the home team in the first half (HTHG) has a significant positive effect on the expected number of goals scored by the home team in the full match. The number of shots on target by the home team (HST) also positively influences the expected number of goals scored by the home team. However, the number of goals scored by the away team in the first half (HTAG) and the number of shots on target by the away team (AST) do not have a significant effect on match outcomes.

### 5.2 Why Poisson Regression?

Poisson regression is chosen as the modeling approach due to its suitability for count data, such as the number of goals scored in football matches. The Poisson distribution naturally accommodates the discrete nature of goal counts and allows for the inclusion of multiple predictor variables to explain variations in match outcomes.

### 5.3 Weaknesses and Next Steps

While the Poisson regression model provides valuable insights, it has limitations. For example, it assumes that the mean and variance of the response variable are equal, which may not always hold true in practice. Additionally, the model does not capture temporal dependencies or interactions between teams.

Future research could explore more sophisticated models, such as hierarchical or time-series models, to address these limitations and improve the accuracy of match outcome predictions.

Table 1: Model of Every Team

|                        | First model     |
|------------------------|-----------------|
| (Intercept)            | −0.23<br>(0.51) |
| HomeTeamBlackburn      | 0.01<br>(0.30)  |
| HomeTeamBristol City   | −0.25<br>(0.32) |
| HomeTeamCardiff        | −0.04<br>(0.32) |
| HomeTeamCoventry       | 0.02<br>(0.51)  |
| HomeTeamHuddersfield   | −0.11<br>(0.31) |
| HomeTeamHull           | −0.09<br>(0.48) |
| HomeTeamIpswich        | 0.39<br>(0.48)  |
| HomeTeamLeeds          | −0.05<br>(0.51) |
| HomeTeamLeicester      | 0.41<br>(0.51)  |
| HomeTeamMiddlesbrough  | −0.32<br>(0.57) |
| HomeTeamMillwall       | −0.49<br>(0.66) |
| HomeTeamNorwich        | 0.34<br>(0.50)  |
| HomeTeamPlymouth       | 0.33<br>(0.44)  |
| HomeTeamPreston        | 0.15<br>(0.49)  |
| HomeTeamQPR            | −0.63<br>(0.77) |
| HomeTeamRotherham      | 0.08<br>(0.67)  |
| HomeTeamSheffield Weds | −0.11<br>(0.47) |
| HomeTeamSouthampton    | 0.40<br>(0.42)  |
| HomeTeamStoke          | −0.78<br>(1.06) |
| HomeTeamSunderland     | −0.19<br>(0.63) |
| HomeTeamSwansea5       | 0.00<br>(0.48)  |
| HomeTeamWatford        | −0.21<br>(0.59) |
| HomeTeamWest Brom      | 0.11<br>(0.44)  |
| AwayTeamBlackburn      | −0.31<br>(0.28) |

## References

- football-data.co.uk. 2024. “England.” <https://www.football-data.co.uk/englandm.php>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Smith, John. 2020. “Advanced Sports Analytics: Predictive Models for Football Match Outcomes” Container-Title: Sports Analytics Journal.” *Journal of Open Source Software* 10 (2): 123–35.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Golemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.