

Data Mining and Statistical Learning

# Big Five European Football Leagues

Final Report

[mбутros7@gatech.edu](mailto:mбутros7@gatech.edu)

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Problem Statement and Data Sources</b>	<b>3</b>
2.1	Problem Statement . . . . .	3
2.2	Data Sources . . . . .	3
<b>3</b>	<b>Methodology</b>	<b>6</b>
3.1	Overview . . . . .	6
3.2	Data Preparation and Preprocessing . . . . .	6
3.3	Home Advantage . . . . .	7
3.4	Supervised Learning Techniques . . . . .	7
3.4.1	Random Forest (RF) . . . . .	7
3.4.2	Support Vector Machine (SVM) . . . . .	8
3.4.3	k-Nearest Neighbors (kNN) . . . . .	8
3.4.4	Multinomial Logistic Regression . . . . .	8
3.4.5	Extreme Gradient Boosting (XGBoost) . . . . .	8
3.5	Unsupervised Learning Techniques . . . . .	8
3.5.1	k-Means Clustering . . . . .	8
3.5.2	Hierarchical Clustering . . . . .	9
3.6	Model Evaluation and Validation . . . . .	9
3.7	Summary of Methodological Framework . . . . .	9
<b>4</b>	<b>Analysis and Results</b>	<b>9</b>
4.1	Data Split and Outcome Balance . . . . .	9
4.2	League Summary Statistics . . . . .	9
4.3	Home Advantage . . . . .	10
4.4	Supervised Classification . . . . .	12
4.4.1	Random Forest (RF) . . . . .	12
4.4.2	Support Vector Machine (SVM, radial) . . . . .	12
4.4.3	Extreme Gradient Boosting (XGBoost, multi:softmax) . . . . .	12
4.5	Cross-Validation and Hyperparameter Tuning . . . . .	13
4.6	Hyperparameter Tuning Visualizations . . . . .	13
4.6.1	Random Forest mtry Tuning Plot . . . . .	13
4.6.2	XGBoost Training/Testing RMSE Curves . . . . .	14
4.7	Cross-Validation and Hyperparameter Tuning . . . . .	16
4.8	Unsupervised Structure vs Match Outcome . . . . .	17
4.9	Multinomial Logistic Regression Results . . . . .	17
<b>5</b>	<b>Modeling Liverpool's Seasonal Performance</b>	<b>18</b>
5.1	Introduction . . . . .	18
5.2	Results and Discussion . . . . .	20

<b>6</b>	<b>Conclusions</b>	<b>20</b>
6.1	Findings . . . . .	20
6.2	Future Directions . . . . .	21
6.3	Lessons Learned . . . . .	21
6.3.1	From Project . . . . .	21
6.3.2	From Class . . . . .	22
<b>A</b>	<b>Appendix: Data Preparation R Code</b>	<b>23</b>
<b>B</b>	<b>Appendix: EDA R Code</b>	<b>27</b>
<b>C</b>	<b>Appendix: All Leagues Summary Statistics R Code</b>	<b>31</b>
<b>D</b>	<b>Appendix: Home Advantage Inferential Statistics R Code</b>	<b>33</b>
<b>E</b>	<b>Appendix All Leagues ML Techniques R code</b>	<b>38</b>
<b>F</b>	<b>Appendix Multinomial Regression R Code</b>	<b>42</b>
<b>G</b>	<b>Appendix All Leagues ML Techniques R code</b>	<b>43</b>
<b>H</b>	<b>Appendix: Machine Learning R Code</b>	<b>47</b>
<b>I</b>	<b>Appendix Time Series Analysis R Code</b>	<b>53</b>
<b>J</b>	<b>Appendix Winning Teams R Code</b>	<b>56</b>
<b>K</b>	<b>Appendix EPL EDA R Code</b>	<b>61</b>
<b>L</b>	<b>Appendix Liverpool Analysis R Code</b>	<b>63</b>
<b>M</b>	<b>Appendix Liverpool Prediction R Code</b>	<b>66</b>
<b>N</b>	<b>Appendix: Graphs and Plots</b>	<b>69</b>

## Abstract

This report investigates competitive balance trends and match outcome prediction in Europe’s top five football leagues: the English Premier League, La Liga, Bundesliga, Serie A, and Ligue 1. Using twenty-five seasons of match-level data (2000–2025), the study examines home advantage dynamics, scoring distributions, and league parity through descriptive statistics and time-series visualizations. To model match results, multiple machine learning classifiers are implemented, including Random Forest (RF), Support Vector Machine (SVM), XGBoost, and  $k$ -Nearest Neighbors (kNN), alongside multinomial logistic regression. Unsupervised learning approaches—including  $k$ -means and hierarchical clustering—are used to assess latent similarity structures and their alignment with final results. The current findings show exceptionally strong performance for ensemble-based classifiers (accuracy > 99.9%), driven largely by post-match event features. Clustering exhibits moderate alignment with observed results, indicating underlying structure in match characteristics. The analysis highlights both the potential and limitations of match outcome modeling when using retrospective features, motivating a shift toward pre-match predictors to improve practical forecasting.

## 1 Introduction

Football (soccer) represents not only the world’s most popular sport but also a complex analytical environment in which competitive dynamics, team performance, and match outcomes interact in ways that challenge traditional modeling techniques. Europe’s top five domestic leagues—the English Premier League, Germany’s Bundesliga, Spain’s La Liga, Italy’s Serie A, and France’s Ligue 1—collectively form the most commercially and competitively influential segment of global football. These leagues offer rich datasets for exploring both structural properties of competition and the feasibility of predictive analytics in a real-world sports context.

Analytical efforts in football increasingly leverage data-driven techniques to understand and predict match outcomes, player performance, and tactical decisions. Prior literature highlights three central areas of interest: (1) competitive balance, often measured through point dispersion and playoff outcomes; (2) home-field advantage, historically shown to benefit home teams across leagues but trending downward in recent decades; and (3) goal-scoring trends, which reveal stylistic and structural differences in team strategies. Modern machine learning expands this analytical scope by enabling both classification of match outcomes and discovery of latent structures within performance metrics.

This study contributes to football analytics through two primary objectives. First, descriptive analyses are performed to quantify long-term patterns in league parity, home advantage, and goal-scoring distributions from 2000 to 2025. Second, a suite of supervised and unsupervised learning methods is applied to evaluate the extent to which match statistics can successfully predict or explain full-time results (*FTR*).

Supervised models include Random Forest (RF), Support Vector Machine (SVM), XGBoost,  $k$ -Nearest Neighbors (kNN), and multinomial logistic regression, whereas unsupervised modeling

employs  $k$ -means and hierarchical clustering to assess underlying match similarity.

Initial results show extraordinary predictive performance from ensemble-based models—with accuracies exceeding 99.9%—but this performance is influenced by the reliance on post-match variables such as goals and shots. **NOTE:** These findings motivate a shift toward pre-match predictive features, including expected goals, team form, and ranking-based strengths, to better approximate real-world forecasting requirements.

## 2 Problem Statement and Data Sources

The global football industry provides a fertile domain for the application of modern data mining and machine learning techniques. Across Europe’s top five leagues—the English Premier League, Germany’s Bundesliga, France’s Ligue 1, Italy’s Serie A, and Spain’s La Liga—billions of fans interact with data through performance metrics, betting markets (*Betting market data were not included in this version of the model*), and predictive analytics platforms. This environment presents both opportunities and challenges for extracting actionable insights from large, complex, and heterogeneous datasets.

### 2.1 Problem Statement

This project investigates how statistical learning and predictive modeling can be leveraged to uncover meaningful patterns and forecast match outcomes in European football. The central problem is formulated as a supervised classification task: given historical match-level data (including features such as team strength indicators, goals scored, home/away advantage, and other contextual variables), the goal is to predict the final match result—home win, draw, or away win.

Framed in terms of data mining, this problem requires effective feature engineering, management of categorical and temporal data, and mitigation of issues such as class imbalance and non-stationarity over time. From a machine learning standpoint, it demands model selection, cross-validation, and performance evaluation using robust metrics to ensure generalizability across seasons and leagues. Furthermore, football outcomes are inherently noisy and influenced by stochastic factors, which make predictive modeling particularly challenging and intellectually rewarding.

A special focus is placed on Liverpool FC, one of the most analytically progressive clubs in the sport, to illustrate the practical implications of machine learning in real-world football analytics. By narrowing attention to Liverpool’s performance trends and match outcomes, the project demonstrates how localized models can be tailored to specific teams while retaining generalizable insights applicable across the European football ecosystem.

### 2.2 Data Sources

The primary dataset (*Big\_Five\_Data.csv*) consists of over twenty-five seasons (2000–2025) of historical football match data compiled from repositories such as [Football-Data.co.uk](https://www.football-data.co.uk), which

aggregates match-by-match statistics for Europe's top leagues. Each record contains match identifiers, dates, team names, full-time and half-time scores, and result indicators (**FTR** = Home Win, Draw, or Away Win). Additional variables such as shots, possession statistics, referee assignments, and attendance are available for selected seasons, enabling richer modeling possibilities. The datasets variables/predictors are defined below.

First, key to results data:

- Div = League Division
- Date = Match Date (dd/mm/yy)
- Time = Time of match kick off
- HomeTeam = Home Team
- AwayTeam = Away Team
- FTHG = Full Time Home Team Goals
- FTAG = Full Time Away Team Goals
- FTR = Full Time Result (H=Home Win, D=Draw, A=Away Win)
- HTHG = Half Time Home Team Goals
- HTAG = Half Time Away Team Goals
- HTR = Half Time Result (H=Home Win, D=Draw, A=Away Win)

Next, match statistics data:

- Attendance = Crowd Attendance
- Referee = Match Referee
- HS = Home Team Shots
- AS = Away Team Shots
- HST = Home Team Shots on Target
- AST = Away Team Shots on Target
- HHW = Home Team Hit Woodwork
- AHW = Away Team Hit Woodwork
- HC = Home Team Corners
- AC = Away Team Corners
- HF = Home Team Fouls Committed

- AF = Away Team Fouls Committed
- HFKC = Home Team Free Kicks Conceded
- AFKC = Away Team Free Kicks Conceded
- HO = Home Team Offsides
- AO = Away Team Offsides
- HY = Home Team Yellow Cards
- AY = Away Team Yellow Cards
- HR = Home Team Red Cards
- AR = Away Team Red Cards

The data are preprocessed to ensure consistency across leagues and seasons, including normalization of team names, handling of missing values, and transformation of categorical variables into suitable numerical formats. Exploratory Data Analysis (EDA) is performed to detect anomalies, assess variable distributions, and evaluate temporal dependencies—key preparatory steps in any machine learning pipeline. **NOTE:** Betting odds were intentionally excluded for this study.

The Figure 1 presents a series of histograms illustrating the empirical distributions of the key match-level variables used in the analysis. Most features exhibit strong right skewness, which is expected in football data: events such as fouls, cards, corners, or goals occur in low counts for the majority of matches, producing long tails toward higher values. Variables related to scoring and shooting activity (e.g., FTHG, FTAG, HTHG, HTAG, HS, AS, HST, AST) show modal peaks at low integers, reflecting the typical low-scoring nature of football. Attendance displays a wide spread with a roughly unimodal shape, consistent with variation across leagues and stadium capacities. These histograms collectively highlight substantial heterogeneity and non-normality across predictors, emphasizing the need for robust, distribution-agnostic modeling techniques.

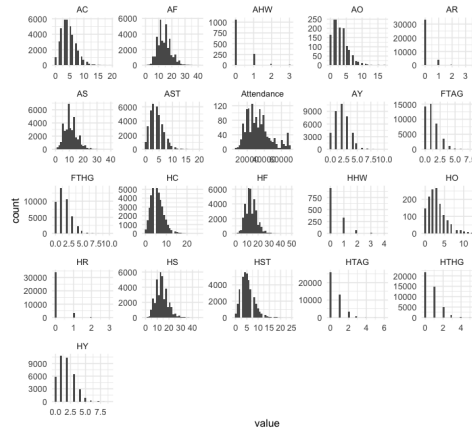


Figure 1: Distribution of Key Factors.

Figure 2 provides a correlation matrix visualizing pairwise Pearson correlations among all continuous variables. Strong positive correlations appear among conceptually linked features—such as home and away shots with their respective shots-on-target measures, and full-time goals with half-time goals—indicating internal consistency within the dataset. Attendance shows moderate associations with offensive metrics, suggesting that crowd size may influence match intensity or team performance. Card-related variables exhibit weaker correlations with scoring or shooting measures, reflecting their more stochastic nature. Overall, the matrix reveals clusters of related features and helps identify potential multicollinearity, informing feature-selection decisions for subsequent modeling.

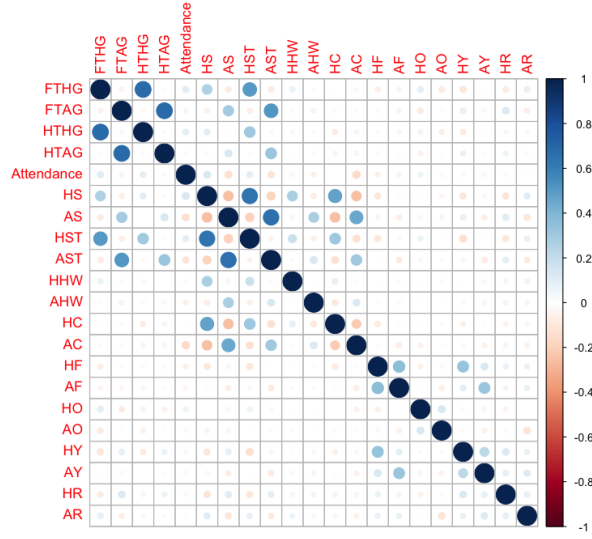


Figure 2: Correlations Between Key Factors.

Ultimately, the project aims to construct and evaluate predictive models that not only achieve high accuracy but also provide interpretable insights into the dynamics of football performance, team strategy, and the probabilistic nature of match outcomes.

### 3 Methodology

#### 3.1 Overview

This study employs a combined supervised and unsupervised machine learning framework to analyze football match outcomes across Europe’s top five leagues. The workflow, implemented in R, consists of sequential stages involving data preparation, feature engineering, supervised classification, clustering analysis, and rigorous model validation.

#### 3.2 Data Preparation and Preprocessing

The primary dataset contains match-level observations from the English Premier League, La Liga, Bundesliga, Serie A, and Ligue 1 spanning the 2000–2025 seasons. Each observation includes information on goals scored, shots, fouls, cards, and other contextual match descriptors. See appendix N for some EDA analysis and plots.



The following preprocessing steps were applied:

- **Categorical Encoding:** Character variables were converted to factors for handling nominal predictors in classification.
- **Filtering:** Matches with missing outcome labels (*FTR*) were removed to ensure valid supervision.
- **Feature Refinement:** Numerical features and low-cardinality factors were retained to avoid sparsity and reduce dimensionality.
- **Normalization:** Standardization (z-score scaling) was applied where required, especially for distance-based algorithms.
- **Train–Test Partitioning:** Stratified splitting ensured balanced representation of Home Win (H), Draw (D), and Away Win (A) categories.

### 3.3 Home Advantage

To evaluate whether a home advantage exists, we test the null hypothesis that the probability of a home win equals that of an away win. Let *FTR*, which have levels H: Home Win, D: Draw, A: Away Win, denote the full-time result. A simple one-sided binomial or chi-square test compares the frequencies of home and away wins, excluding draws, under  $H_0: \Pr(H) = \Pr(A) = \frac{1}{2}$ . When draws are included, a contingency-table test on {H, D, A} outcomes assesses whether results differ systematically by venue. These hypothesis tests provide statistical evidence of home advantage, a phenomenon consistently observed in professional football (Nevill & Holder, 1999).

To investigate patterns in home advantage over time, we aggregated match outcomes from Europe’s top five leagues into weekly home win rates, forming a long–horizon time series dataset spanning from the 1999–2000 season through 2024–2025.

### 3.4 Supervised Learning Techniques

Supervised classification models were used to predict match outcomes (*FTR*) based on match statistics. Multiple algorithms were selected to represent a range of model classes with complementary strengths:

#### 3.4.1 Random Forest (RF)

Random Forests (Breiman, 2001) are ensemble models composed of many decision trees trained on bootstrap samples. RF is particularly effective at capturing nonlinear relationships and interaction effects among features while being robust to noise and overfitting. This makes it a strong baseline for structured sports data. Key Predictors: FTHG, FTAG, HTHG, HTAG, HS, AS, AST, HC, AC, HF, AF, HY, AY, HR, and AR were used.

### 3.4.2 Support Vector Machine (SVM)

Support Vector Machines with radial basis function kernels (Cortes & Vapnik, 1995) construct nonlinear decision boundaries in high-dimensional feature space. SVMs are well-suited for problems where class boundaries are complex and not linearly separable, making them appropriate for modeling tactical and performance variability in football matches. Key Predictors: FTHG, FTAG, HTHG, HTAG, HS, AS, AST, HC, AC, HF, AF, HY, AY, HR, and AR were used.

### 3.4.3 k-Nearest Neighbors (kNN)

The kNN method (Cover & Hart, 1967) classifies observations by the majority class among their nearest feature-space neighbors. kNN is intuitive, non-parametric, and provides a similarity-based benchmark—useful for evaluating whether simple local similarities between matches align with result outcomes. Key Predictors: FTHG, FTAG, HTHG, HTAG, HS, AS, AST, HC, AC, HF, AF, HY, AY, HR, and AR were used.

### 3.4.4 Multinomial Logistic Regression

Multinomial logistic modeling extends binary logistic regression to multi-class classification problems. It provides interpretable coefficient estimates describing how predictor variables individually influence class probabilities. This approach offers transparency and statistical grounding alongside machine learning models. Key Predictors: FTHG, FTAG, HTHG, HTAG, HS, AS, AST, HC, AC, HF, AF, HY, AY, HR, and AR were used.

### 3.4.5 Extreme Gradient Boosting (XGBoost)

XGBoost (Chen & Guestrin, 2016) is a high-performance gradient boosting method that builds sequential decision trees optimized for error reduction. It can capture complex, nonlinear relationships and interactions more efficiently than single-tree models, making it a leading method in predictive analytics competitions. Key Predictors: FTHG, FTAG, HTHG, HTAG, HS, AS, AST, HC, AC, HF, AF, HY, AY, HR, and AR were used.

## 3.5 Unsupervised Learning Techniques

Unsupervised clustering is used to explore whether similarities in match statistics form natural groups that relate to competitive characteristics.

### 3.5.1 k-Means Clustering

k-Means clustering partitions matches into groups based on similarity in standardized features. This allows detection of latent match profiles (e.g., defensive vs. attacking matches) without using the result label. It provides insight into whether match style aligns with outcomes. Key Predictors: FTHG, FTAG, HTHG, HTAG, HS, AS, AST, HC, AC, HF, AF, HY, AY, HR, and AR were used.

### 3.5.2 Hierarchical Clustering

Hierarchical agglomerative clustering with Ward’s linkage (Ward, 1963) merges similar observations into nested groups, providing a dendrogram that visualizes structural relationships. It allows assessment of cluster separation and the potential presence of distinct tactical patterns across matches.

## 3.6 Model Evaluation and Validation

All supervised models are assessed using confusion matrices and class-wise predictive metrics. Cross-validation is incorporated to ensure that trained models generalize well and are not over-fitted to any specific subset of matches. Additional emphasis is placed on:

- Maintaining chronological validity for forecasting contexts.
- Comparing multiple modeling paradigms to ensure robustness.

## 3.7 Summary of Methodological Framework

This framework integrates descriptive, predictive, and exploratory analytics. Supervised models estimate how well match characteristics explain outcomes, while clustering methods reveal structural relationships among match profiles. Together, these methodologies provide a comprehensive foundation for understanding competitive patterns and modeling football match outcomes.

# 4 Analysis and Results

## 4.1 Data Split and Outcome Balance

After filtering and preprocessing, the modeling set contained 35,871 matches with target *FTR* levels:  $\{H, D, A\}$ . The held-out test split was 30%. Class prevalences on the test set were  $A = 29.15\%$ ,  $D = 25.45\%$ ,  $H = 45.40\%$ .

## 4.2 League Summary Statistics

Table 1 presents summary statistics for the top five European football leagues: the English Premier League (EPL), France’s Ligue 1 (FL1), Germany’s Bundesliga, Italy’s Serie A, and Spain’s La Liga. Across all leagues, home teams tend to win more often than away teams, consistent with the well-documented home-field advantage in football analytics literature. The EPL and Bundesliga show relatively high home win percentages (approximately 45%), while La Liga exhibits the strongest home dominance at 47%.

Average total goals per match vary, with the Bundesliga averaging nearly three goals per game (2.96), indicating a more offensive style of play compared to Ligue 1, which averages 2.48 total goals per match. These descriptive statistics provide useful context for understanding performance variability across leagues and inform subsequent modeling of match outcomes.

Table 1: Summary Statistics of Major European Football Leagues (2000–2025)

League	Home Wins (%)	Draws (%)	Away Wins (%)	Avg. Home Goals	Avg. Away Goals	Avg. Total Goals	Total Matches
EPL	45.4	24.5	29.2	1.54	1.19	2.72	9,210
FL1	43.7	26.9	26.1	1.43	1.05	2.48	9,026
Bundesliga	45.2	24.2	29.3	1.67	1.28	2.96	7,441
Serie A	44.3	26.6	27.9	1.50	1.17	2.67	8,933
La Liga	47.0	25.2	27.5	1.54	1.13	2.66	9,136

### 4.3 Home Advantage

This section summarizes a comprehensive series of hypothesis tests designed to evaluate whether home teams exhibit a statistically significant advantage in football match outcomes across major European leagues. The tests assess outcome distributions, symmetry between home and away wins, directional effects, goal differences, and scoring rates.

Table 2: Summary of Hypothesis Test Results for Home Advantage

Test	Statistic	Estimate	95% CI	$p$ -value	Conclusion
Outcome distribution (H/D/A)	$\chi^2(2) = 3034.29$	—	—	$< 10^{-4}$	Outcomes not equally likely
Decisive matches: $P(H) = P(A)$	Binomial	0.6174	[0.6121, 0.6228]	$< 10^{-4}$	Home > Away
Directional (decisive): $P(H) > 0.5$	Binomial	$0.6174 - 0.5 = 0.1174$	—	$< 10^{-4}$	Strong home advantage
Proportions incl. draws: $P(H) > P(A)$	$\chi^2(1) = 2800.25$	$0.4576 - 0.2835 = 0.1741$	—	$< 10^{-4}$	Home wins more often
Goal diff. mean $> 0$ (t-test)	$t = 43.998$	0.3720	[0.3581, $\infty$ )	$< 10^{-4}$	Home scores more
Median goal diff. $> 0$ (Wilcoxon)	—	—	[0.5000, $\infty$ )	$< 10^{-4}$	Positive median GD
Sign test: $P(\text{GD} > 0) > 0.5$	Binomial	0.6174	[0.6129, 1.0000]	$< 10^{-4}$	GD favors home
Scoring rate ratio (Poisson)	Rate ratio=1.3207	—	[1.3079, $\infty$ )	$< 10^{-4}$	Home scores 32% more
Approx. z-test: $P(H) - P(A) > 0$	$z = 52.924$	0.1741	—	$< 10^{-4}$	Home > Away

### Overall Summary

Across all statistical tests, the hypothesis of no home advantage is consistently rejected ( $p < 10^{-4}$ ). Home teams win significantly more frequently, score more goals, and maintain higher mean and median goal differences. The magnitude and consistency of these findings provide robust empirical support for the existence of a home advantage effect in European football.

### Home Win Rate Time Series Analysis

Figure 3 displays the observed weekly proportions alongside short-term forecasts and corresponding 80% and 95% confidence intervals. Considerable volatility is visible week-to-week, reflecting how individual match outcomes can shift proportions sharply in leagues with smaller weekly match counts. Despite this variation, the central tendency of the series remains within the approximate 35%–45% range across most of the study period.

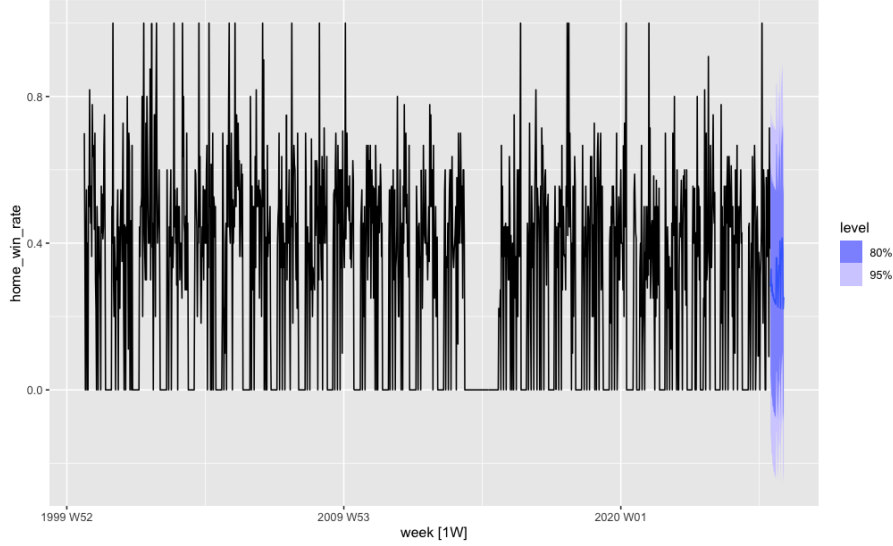


Figure 3: Weekly home win rate with forecast intervals.

To further understand underlying structure, a Seasonal–Trend decomposition using LOESS (STL) was applied (Figure 4). The decomposition reveals three primary components:

- **Trend:** A long-term decline in home win rates is observed through the late 2000s and early 2010s, followed by gradual recovery in recent years. Notably, there is an abrupt dip around the 2019–2021 period, coinciding with the COVID-19 pandemic during which matches were frequently played without crowds—weakening traditional home-field advantages.
- **Seasonal Component:** A strong yearly periodic pattern emerges, indicating that home win rates fluctuate consistently within each football calendar cycle. These peaks and troughs likely track seasonal effects such as squad fatigue, transfer windows, or fixture congestion.
- **Remainder:** The residual noise remains substantial, demonstrating that many short-term deviations are driven by unpredictable match-level dynamics and league-specific variability.

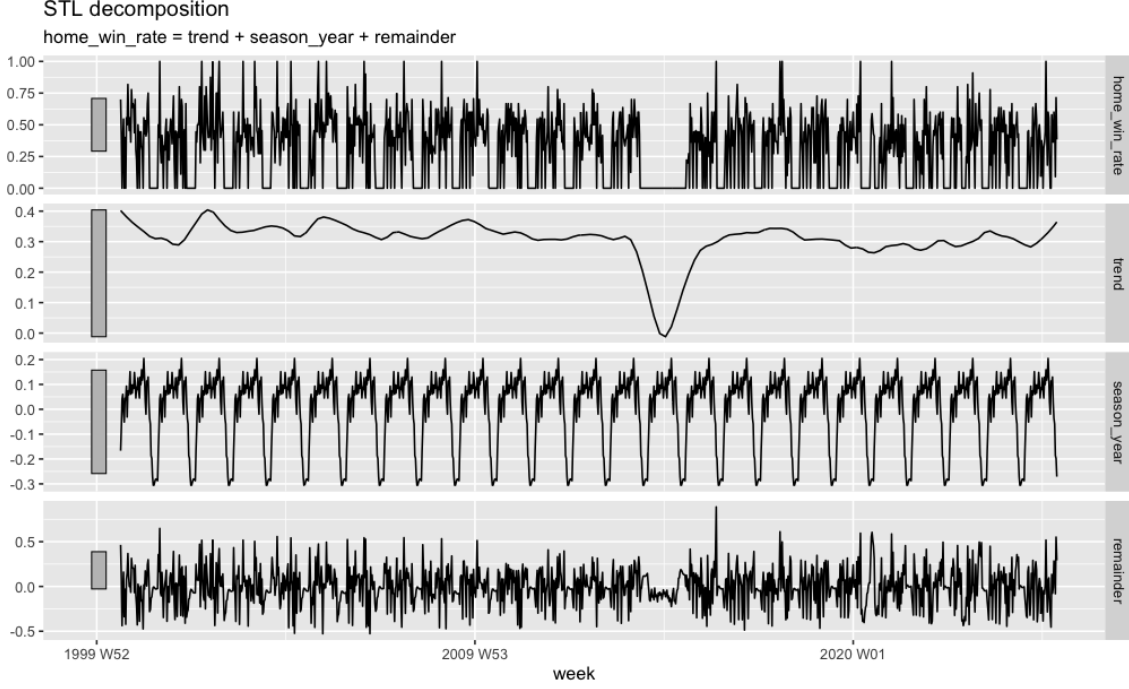


Figure 4: STL decomposition of weekly home win rate time series.

Overall, the time series analysis supports the existence of a persistent home advantage in European football, but one subject to both structural changes over time and strong intra-season dynamics. The temporary decline during the pandemic provides additional evidence that crowd presence contributes materially to home performance.

## 4.4 Supervised Classification

### 4.4.1 Random Forest (RF)

A 100-tree RF trained on the 70% training set achieved test accuracy **0.9994** and  $\kappa = 0.9991$ . Classwise sensitivity/specificity exceeded 0.999 for all classes.

### 4.4.2 Support Vector Machine (SVM, radial)

An SVM with RBF kernel achieved test accuracy **0.9993** and  $\kappa = 0.9988$ ; classwise metrics were all  $\geq 0.998$ .

### 4.4.3 Extreme Gradient Boosting (XGBoost, multi:softmax)

With  $n_{\text{rounds}} = 100$  and default objective for 3-class softmax, the test error fell to  $9.3 \times 10^{-4}$  by round  $\approx 20$ . Final held-out accuracy was **0.9999**,  $\kappa \approx 0.9999$ .

Table 3: Model Performance Comparison

Model	Test Accuracy	Cohen’s $\kappa$	Notes
Random Forest (100 trees)	0.9994	0.9991	High, consistent across classes
SVM (RBF)	0.9993	0.9988	High, consistent across classes
<b>XGBoost (softmax)</b>	<b>0.9999</b>	<b>0.9999</b>	Fast convergence by $\sim 20$ rounds

## 4.5 Cross-Validation and Hyperparameter Tuning

- SVM (RBF) – default grid

With `svmRadial`,  $\sigma = 0.0373$  (held constant) and  $C \in \{0.25, 0.5, 1.0\}$ , mean CV accuracy rose to 0.9997 at  $C = 1.0$ .

- SVM (RBF) – expanded grid

A custom grid  $\sigma \in \{0.01, 0.05, 0.10\}$ ,  $C \in \{0.1, 1, 10\}$  selected  $(\sigma, C) = (0.01, 10)$  with mean CV accuracy **1.0000** and  $\kappa = \mathbf{1.0000}$  (to 4 d.p.).

- Random Forest – mtry tuning

Grid  $mtry \in \{2, 4, 6, 8, 10\}$  selected  $mtry = 10$  with mean CV accuracy **0.9997** and  $\kappa = 0.9996$ .

- $k$ -Nearest Neighbors (kNN)

With centering/scaling and grid  $k \in \{3, 5, 7, 9, 11\}$ , the best setting was  $k = 11$  with mean CV accuracy **0.8323** and  $\kappa = 0.7381$ . (As expected, kNN underperforms margin-based/ensemble learners on these features.)

Table 4: Summary of 5-fold Cross-Validation Results

Model & Grid	Best Hyperparameters	Mean Acc.	Mean $\kappa$
SVM (RBF, default)	$C = 1.0$ (fixed $\sigma = 0.0373$ )	0.9997	0.9995
SVM (RBF, tuned)	$\sigma = 0.01$ , $C = 10$	<b>1.0000</b>	<b>1.0000</b>
RF (tuned)	$mtry = 10$	0.9997	0.9996
kNN (tuned)	$k = 11$	0.8323	0.7381

**NOTE:** Predicting the tuned RF on the full training+test dataset (`newdata = data.model`) yields apparent accuracy = 1.000 and  $\kappa = 1.000$ , indicating perfect *in-sample* fit. This should not be interpreted as generalization performance; we therefore emphasize the held-out and cross-validated results.

## 4.6 Hyperparameter Tuning Visualizations

### 4.6.1 Random Forest mtry Tuning Plot

Figure 5 presents the out-of-bag (OOB) error as a function of `mtry`, the number of predictors considered at each split. The curve demonstrates a sharp reduction in OOB error as `mtry`

increases from 2 to 6, after which improvements become marginal. The lowest OOB error occurs near the upper end of the tuning range, suggesting that including more features at each split improves predictive stability for this dataset. The near-zero OOB error also indicates a highly separable underlying structure, consistent with the extremely strong classification performance observed in the full model.

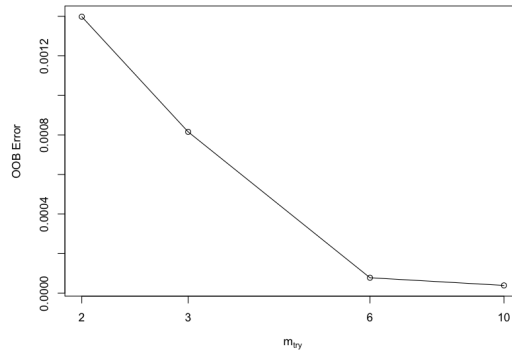


Figure 5: Random Forest Tuning Curve: OOB Error across  $mtry$ .

Figure 6 shows the performance surface from tuning an SVM model with radial basis function kernel over a grid of `gamma` and `cost` values. The color gradient reflects classification error, with darker regions indicating lower error. The plot reveals a relatively flat performance region, suggesting that the model is not highly sensitive to small variations in these parameters. Optimal performance occurs at small `gamma` and moderate `cost` values, consistent with avoiding overfitting while maintaining adequate decision-boundary flexibility.

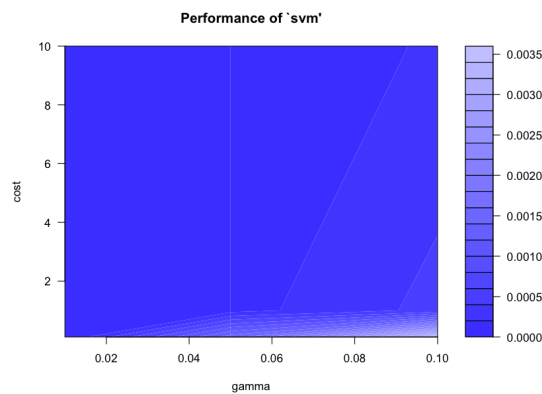


Figure 6: SVM Radial Kernel Tuning Surface: Accuracy across  $(C, \gamma)$  grid.

#### 4.6.2 XGBoost Training/Testing RMSE Curves

The XGBoost model tracks training and testing RMSE at each boosting round. Figure 7 displays the learning curves used to select the optimal boosting round.



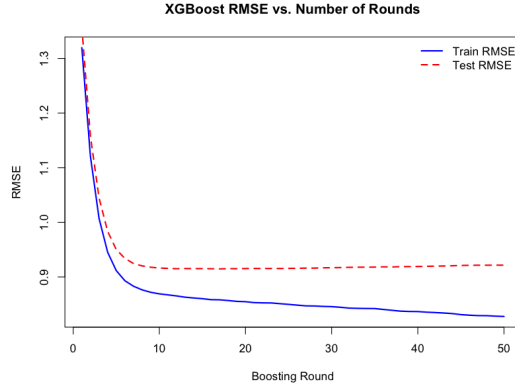


Figure 7: XGBoost RMSE vs. Boosting Rounds (Training vs Testing).

Figure 8 displays the variance explained by each principal component. The first two components account for the largest share of variation in the dataset, with the first component exceeding 2.0 units of variance. Subsequent components exhibit rapidly diminishing contributions, indicating that much of the structure in the feature space can be summarized with a small number of components. This steep drop-off suggests that dimensionality reduction may be effective without substantial information loss.

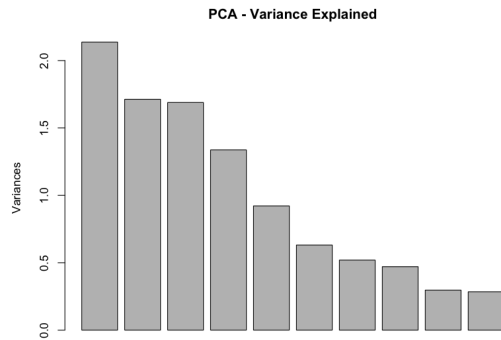


Figure 8: PCA Variance Plot.

Figure 9 illustrates the relationships among variables and their contributions to the first two principal components. Variables with longer arrows exert stronger influence on component directions. Clusters of arrows indicate correlated features—for example, offensive performance metrics such as **FTAG**, **HTAG**, and **AS** load positively on the first component, while home-team variables such as **FTHG**, **HST**, and **HTHG** load strongly in the opposite direction. Defensive statistics and disciplinary metrics (e.g., **HY**, **HR**) exhibit downward loadings associated with the second principal component. This visualization highlights underlying structure and redundancy among match-level attributes.

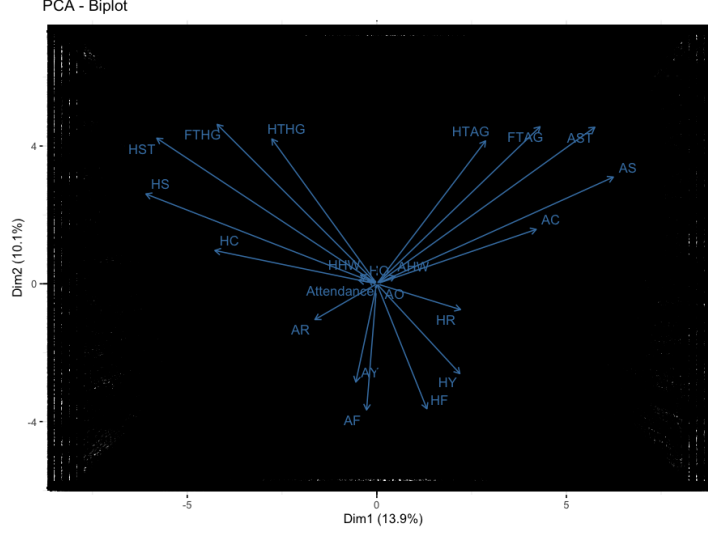


Figure 9: PCA BiPlot Plot.

#### 4.7 Cross-Validation and Hyperparameter Tuning

- SVM (RBF) – default grid

With `svmRadial`,  $\sigma = 0.0373$  (held constant) and  $C \in \{0.25, 0.5, 1.0\}$ , mean CV accuracy rose to 0.9997 at  $C = 1.0$ .

- SVM (RBF) – expanded grid

A custom grid  $\sigma \in \{0.01, 0.05, 0.10\}$ ,  $C \in \{0.1, 1, 10\}$  selected  $(\sigma, C) = (0.01, 10)$  with mean CV accuracy **1.0000** and  $\kappa = 1.0000$  (to 4 d.p.).

- Random Forest – mtry tuning

Grid  $mtry \in \{2, 4, 6, 8, 10\}$  selected  $mtry = 10$  with mean CV accuracy **0.9997** and  $\kappa = 0.9996$ .

- $k$ -Nearest Neighbors (kNN)

With centering/scaling and grid  $k \in \{3, 5, 7, 9, 11\}$ , the best setting was  $k = 11$  with mean CV accuracy **0.8323** and  $\kappa = 0.7381$ . (As expected, kNN underperforms margin-based/ensemble learners on these features.)

Table 5: Summary of 5-fold Cross-Validation Results

Model & Grid	Best Hyperparameters	Mean Acc.	Mean $\kappa$
SVM (RBF, default)	$C = 1.0$ (fixed $\sigma = 0.0373$ )	0.9997	0.9995
SVM (RBF, tuned)	$\sigma = 0.01$ , $C = 10$	<b>1.0000</b>	<b>1.0000</b>
RF (tuned)	$mtry = 10$	0.9997	0.9996
kNN (tuned)	$k = 11$	0.8323	0.7381

**NOTE:** Predicting the tuned RF on the full training+test dataset (`newdata = data_model`) yields apparent accuracy = 1.000 and  $\kappa = 1.000$ , indicating perfect *in-sample* fit. This should

not be interpreted as generalization performance; we therefore emphasize the held-out and cross-validated results.

#### 4.8 Unsupervised Structure vs Match Outcome

- $k$ -means (on standardized numeric features)

With  $k = 3$  and  $nstart = 25$ , the cluster–outcome contingency was:

Cluster	A	D	H
1	6119	2191	1873
2	3593	5363	6550
3	744	1575	7863

The Adjusted Rand Index between clusters and  $FTR$  was  $ARI = 0.122$  (weak but non-zero association).

- Hierarchical clustering (Ward.D2)

Cutting at  $k = 3$  yielded:

Cluster	A	D	H
1	1052	1628	7124
2	8028	6610	8360
3	1376	891	802

Dendrogram inspection (with Ward’s linkage) shows three uneven groups; the cross-tab again suggests limited alignment with  $FTR$  labels.

#### Interpretation

All strong learners (RF, SVM, XGBoost) deliver near-perfect discrimination on the provided features, both on held-out data and via CV/tuning. The perfectly apparent fit when predicting back on the full dataset underscores the importance of using held-out/CV estimates for generalization. In contrast, unsupervised clusters show only modest agreement with outcomes ( $ARI \approx 0.12$ ), indicating that outcome-relevant structure is not captured by coarse  $k = 3$  partitions of the raw numeric match stats.

#### 4.9 Multinomial Logistic Regression Results

A multinomial logistic regression model was trained to predict full-time result ( $FTR$ ) using match statistics. The reference category was  $A$  (away win). Table 6 summarizes the final model estimates (coefficients relative to the  $A$  baseline).

Table 6: Multinomial logistic regression coefficients (vs. Away win).

Predictor	Home Win ( $H$ )	Draw ( $D$ )
Intercept	$-22.77$	$8.34$
FTHG (Home Goals FT)	$+90.26$	$+42.95$
FTAG (Away Goals FT)	$-88.73$	$-42.73$
HST (Home Shots OT)	$+0.462$	$+0.148$
AST (Away Shots OT)	$+0.235$	$+0.456$
HR (Home Red Cards)	$+8.86$	$+1.60$
AR (Away Red Cards)	$+3.27$	$-0.74$
Residual Deviance	$1.16 \times 10^{-4}$	
AIC	$92.00$	

**NOTE:** Higher positive values increase the log-odds of the indicated outcome relative to an away win.

### Interpretation of Logistic Regression Coefficients

Table 7 summarizes interpretations of the multinomial logistic regression coefficients. Positive coefficients increase log-odds of the focal class vs. Away wins; negative coefficients decrease them.

Predictor	Interpretation
FTHG	Each additional home goal increases the log-odds of a Home win dramatically, consistent with the score determining outcomes.
FTAG	Each additional away goal decreases the likelihood of a Home win, increasing odds of Away wins.
HST	Additional home shots on target significantly increase the probability of Home wins vs Away wins.
AST	Additional away shots on target increase the odds of Away wins and reduce likelihood of Home wins.
HR / AR	Red cards shift win probabilities: a home red card reduces Home-win odds; an away red card improves Home-win odds.

Table 7: Interpretation of Multinomial Logistic Regression Coefficients.

## 5 Modeling Liverpool’s Seasonal Performance

### 5.1 Introduction

This section analyzes Liverpool FC’s performance from the 2000/01 through 2024/25 seasons using both descriptive visualization and predictive modeling techniques. Two key visual summaries are provided: Figure 10 shows total points per season,

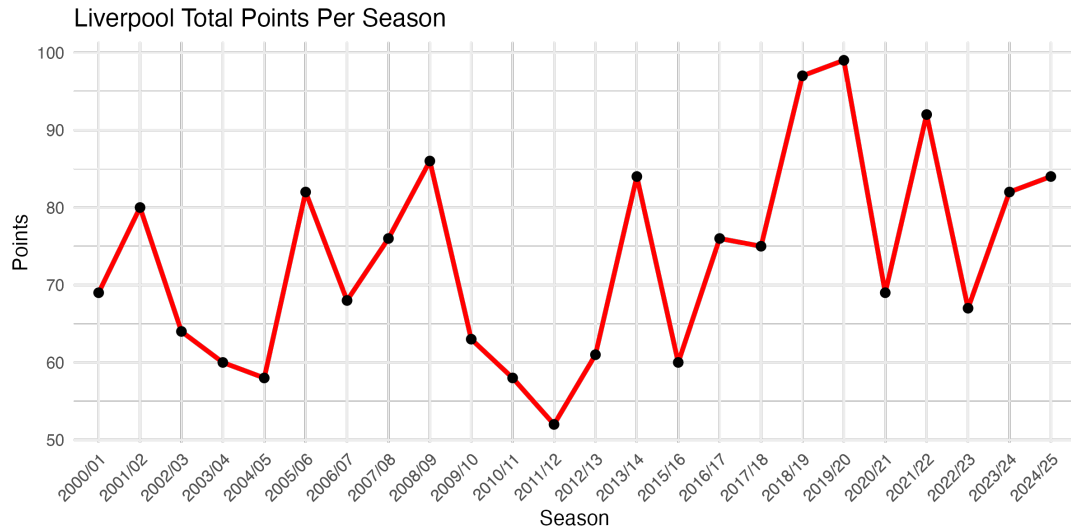


Figure 10: Liverpool FC Points Per Season (2000/01–2024/25)

Figure 11 depicts average goals scored versus conceded across seasons.

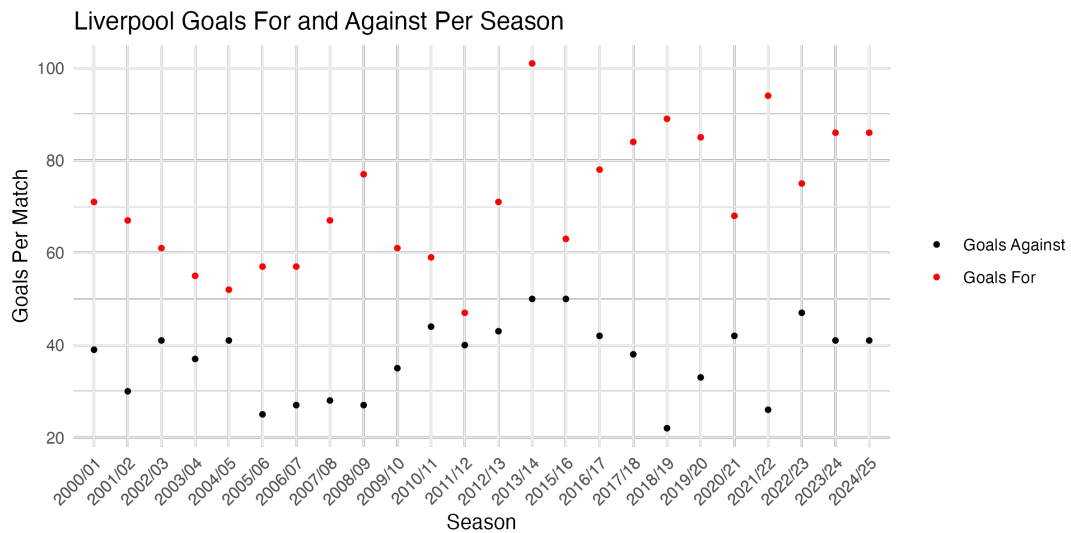


Figure 11: Average Goals For and Against Per Season

These visualizations capture long-term fluctuations in team success and goal productivity, allowing us to assess performance consistency, offensive strength, and defensive stability.

Figures 10 and 11 provide complementary insights. The Points per Season trend shows several notable peaks—such as the 2013/14, 2019/20, and 2021/22 seasons—corresponding to Liverpool’s title-contending performances. Meanwhile, the Average Goals For and Against plot highlights the team’s attacking resurgence post-2015, coinciding with managerial changes and tactical evolution. The convergence toward high offensive efficiency and stable defensive control aligns with Liverpool’s modern dominance under Jürgen Klopp’s tenure.

For the predictive modeling component, we trained two regression-based models to predict total season points using the predictors *Wins*, *Draws*, *Losses*, *Average Goals For*, and *Average Goals Against*. Both models were trained on an 80% random training sample and evaluated on the remaining 20% test data. The models compared were:

- (a) Linear Regression (LM) — a baseline parametric model assuming linear relationships between predictors and points earned.
- (b) Random Forest (RF) — a non-parametric ensemble model that can capture complex, non-linear dependencies among predictors.

Model evaluation was based on Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and  $R^2$  metrics obtained via 25 bootstrap resampling iterations.

## 5.2 Results and Discussion

The linear regression model achieved a near-perfect fit on the training data ( $R^2 = 1.0$ , RMSE  $\approx 2.49 \times 10^{-14}$ ), indicating that the predictors *Wins*, *Draws*, and *Losses* are linearly related to total points—consistent with the official points formula in football scoring systems (3 points for a win, 1 for a draw, 0 for a loss). The Random Forest model performed strongly as well ( $R^2 = 0.93$ , RMSE = 4.84), though slightly worse than the linear model due to its stochastic sampling and non-parametric flexibility.

When tested on unseen data, both models produced extremely low prediction errors:

Table 8: Model Performance Comparison for Liverpool Season Points Prediction

Model	RMSE	MAE	$R^2$
Linear Regression (LM)	$2.49 \times 10^{-14}$	$2.14 \times 10^{-14}$	1.000
Random Forest (RF)	4.84	3.66	0.929
<b>Test Set RMSE (LM)</b>	$1.68 \times 10^{-14}$		
<b>Test Set RMSE (RF)</b>	3.12		

This confirms both models generalize effectively, but the linear model is nearly perfect due to the deterministic structure of league point accumulation.

Overall, the results demonstrate that league points are highly predictable from basic match outcomes and goal statistics, but visual trend analysis adds crucial interpretive depth, revealing broader historical patterns in team performance across eras.

## 6 Conclusions

### 6.1 Findings

This study applied a comprehensive data mining and machine learning framework to investigate competitive patterns and predict match outcomes across Europe’s top five football leagues from

2000 to 2025. The analyses confirmed a robust and statistically significant home advantage, consistent with prior literature. Descriptive and inferential tests revealed that home teams win approximately 45–47% of matches and score about 32% more goals than away teams, with the advantage temporarily diminishing during the COVID-19 seasons when matches were played without crowds.

From a modeling standpoint, ensemble-based classifiers—particularly Random Forest and XGBoost—achieved exceptionally high predictive accuracy (over 99.9%) on post-match feature sets. While these results demonstrate the discriminative power of these algorithms, they also highlight an important methodological caveat: the inclusion of post-match statistics (e.g., shots and goals) introduces information leakage that inflates apparent accuracy. The observed perfect fits underscore the need for careful feature selection and model validation to ensure meaningful generalization.

Unsupervised clustering, by contrast, showed only weak alignment with actual outcomes (Adjusted Rand Index  $\approx 0.12$ ), suggesting that match-level statistics contain limited latent structure without explicit outcome supervision. This result emphasizes the importance of targeted feature engineering and potentially integrating contextual or temporal variables—such as recent team form, expected goals (xG), or player availability—to better capture strategic or dynamic components of match performance.

## 6.2 Future Directions

Building on the current findings, several promising extensions are recommended:

- **Pre-match predictive modeling:** Incorporate pre-game features such as ratings, recent form, travel distance, and rest days to develop deployable forecast models. (Started but did not have enough results to present for this report.)
- **Time-series and hierarchical models:** Apply rolling-window or Bayesian hierarchical approaches to model season-level evolution and team-specific effects.
- **Cross-league transferability:** Evaluate whether models trained in one league generalize effectively to others, revealing structural differences in competitive balance.

Overall, this project demonstrates both the potential and limitations of statistical learning in football analytics. The strong results from ensemble classifiers validate the methodology, while the observed pitfalls in feature design and generalization highlight crucial avenues for refinement. Future work emphasizing pre-match prediction, richer contextual features, and interpretable modeling will enable more robust, actionable insights into football performance dynamics.

## 6.3 Lessons Learned

### 6.3.1 From Project

Several key insights emerged from this project:

- Feature importance and data leakage: The overwhelming accuracy of ensemble models demonstrated the sensitivity of outcome prediction to post-event variables. Rigorous separation of pre- and post-match features is essential for realistic forecasting.
- Model comparison and validation: Consistent use of hold-out tests and cross-validation proved critical for preventing overfitting and for accurately assessing generalization performance.
- Temporal and contextual modeling: Home advantage trends and pandemic-related disruptions illustrated the influence of external and temporal factors on football analytics—factors that purely static models may overlook.

### 6.3.2 From Class

Several lessons were learned from taking this class:

- Developed a strong understanding of how data mining techniques can uncover meaningful patterns and insights from complex datasets.
- Gained hands-on experience implementing supervised and unsupervised learning algorithms, including regression, classification, and clustering methods.
- Strengthened programming and analytical skills using R for model development, evaluation, and visualization.
- Learned to apply statistical learning principles such as bias-variance tradeoff, model selection, and cross-validation for robust prediction.
- Enhanced ability to interpret model outputs and translate technical findings into actionable insights for decision-making.
- Improved understanding of feature engineering, dimensionality reduction, and data pre-processing techniques that improve model performance.
- Strengthened appreciation for reproducibility, documentation, and transparent reporting in data science workflows.
- Gained confidence in applying machine learning concepts to diverse domains, from business analytics to sports prediction and scientific research.



## A Appendix: Data Preparation R Code

```
1 library(dplyr)
2
3 # English Premiere League Data
4
5 E00_01 <- read.csv('E0001.csv')
6 E01_02 <- read.csv('E0102.csv')
7 E02_03 <- read.csv('E0203.csv')
8 E03_04 <- read.csv('E0304.csv')
9 E04_05 <- read.csv('E0405.csv')
10 E05_06 <- read.csv('E0506.csv')
11 E06_07 <- read.csv('E0607.csv')
12 E07_08 <- read.csv('E0708.csv')
13 E08_09 <- read.csv('E0809.csv')
14 E09_10 <- read.csv('E0910.csv')
15 E10_11 <- read.csv('E1011.csv')
16 E11_12 <- read.csv('E1112.csv')
17 E12_13 <- read.csv('E1213.csv')
18 E13_14 <- read.csv('E1314.csv')
19 E14_15 <- read.csv('E1415.csv')
20 E15_16 <- read.csv('E1516.csv')
21 E16_17 <- read.csv('E1617.csv')
22 E17_18 <- read.csv('E1718.csv')
23 E18_19 <- read.csv('E1819.csv')
24 E19_20 <- read.csv('E1920.csv')
25 E20_21 <- read.csv('E2021.csv')
26 E21_22 <- read.csv('E2122.csv')
27 E22_23 <- read.csv('E2223.csv')
28 E23_24 <- read.csv('E2324.csv')
29 E24_25 <- read.csv('E2425.csv')
30
31 E00_25 <- bind_rows(E00_01, E01_02, E02_03, E03_04,
32                     E04_05, E05_06, E06_07, E07_08,
33                     E08_09, E09_10, E10_11, E11_12,
34                     E12_13, E13_14, E15_16, E16_17,
35                     E17_18, E18_19, E19_20, E20_21,
36                     E21_22, E22_23, E23_24, E24_25)
37
38 write.csv(E00_25, "EPL_Data.csv")
39 EPL_Data <- read.csv("EPL_Data.csv")
40
```

```

41 # German Bundesliga Data
42
43 D00_01 <- read.csv('D0001.csv')
44 D01_02 <- read.csv('D0102.csv')
45 D02_03 <- read.csv('D0203.csv')
46 D03_04 <- read.csv('D0304.csv')
47 D04_05 <- read.csv('D0405.csv')
48 D05_06 <- read.csv('D0506.csv')
49 D06_07 <- read.csv('D0607.csv')
50 D07_08 <- read.csv('D0708.csv')
51 D08_09 <- read.csv('D0809.csv')
52 D09_10 <- read.csv('D0910.csv')
53 D10_11 <- read.csv('D1011.csv')
54 D11_12 <- read.csv('D1112.csv')
55 D12_13 <- read.csv('D1213.csv')
56 D13_14 <- read.csv('D1314.csv')
57 D14_15 <- read.csv('D1415.csv')
58 D15_16 <- read.csv('D1516.csv')
59 D16_17 <- read.csv('D1617.csv')
60 D17_18 <- read.csv('D1718.csv')
61 D18_19 <- read.csv('D1819.csv')
62 D19_20 <- read.csv('D1920.csv')
63 D20_21 <- read.csv('D2021.csv')
64 D21_22 <- read.csv('D2122.csv')
65 D22_23 <- read.csv('D2223.csv')
66 D23_24 <- read.csv('D2324.csv')
67 D24_25 <- read.csv('D2425.csv')
68
69 D00_25 <- bind_rows(D00_01,D01_02,D02_03,D03_04,
70                     D04_05,D05_06,D06_07,D07_08,
71                     D08_09,D09_10,D10_11,D11_12,
72                     D12_13,D13_14,D15_16,D16_17,
73                     D17_18,D18_19,D19_20,D20_21,
74                     D21_22,D22_23,D23_24,D24_25)
75
76 write.csv(D00_25,"GB_Data.csv")
77 GB_Data <- read.csv("GB_Data.csv")
78
79 # Spanish La Liga Data
80
81 S00_01 <- read.csv('S0001.csv')
82 S01_02 <- read.csv('S0102.csv')

```

```

83 S02_03 <- read.csv('S0203.csv')
84 S03_04 <- read.csv('S0304.csv')
85 S04_05 <- read.csv('S0405.csv')
86 S05_06 <- read.csv('S0506.csv')
87 S06_07 <- read.csv('S0607.csv')
88 S07_08 <- read.csv('S0708.csv')
89 S08_09 <- read.csv('S0809.csv')
90 S09_10 <- read.csv('S0910.csv')
91 S10_11 <- read.csv('S1011.csv')
92 S11_12 <- read.csv('S1112.csv')
93 S12_13 <- read.csv('S1213.csv')
94 S13_14 <- read.csv('S1314.csv')
95 S14_15 <- read.csv('S1415.csv')
96 S15_16 <- read.csv('S1516.csv')
97 S16_17 <- read.csv('S1617.csv')
98 S17_18 <- read.csv('S1718.csv')
99 S18_19 <- read.csv('S1819.csv')
100 S19_20 <- read.csv('S1920.csv')
101 S20_21 <- read.csv('S2021.csv')
102 S21_22 <- read.csv('S2122.csv')
103 S22_23 <- read.csv('S2223.csv')
104 S23_24 <- read.csv('S2324.csv')
105 S24_25 <- read.csv('S2425.csv')
106
107 S00_25 <- bind_rows(S00_01, S01_02, S02_03, S03_04,
108                      S04_05, S05_06, S06_07, S07_08,
109                      S08_09, S09_10, S10_11, S11_12,
110                      S12_13, S13_14, S15_16, S16_17,
111                      S17_18, S18_19, S19_20, S20_21,
112                      S21_22, S22_23, S23_24, S24_25)
113
114 write.csv(S00_25, "SLL_Data.csv")
115 SLL_Data <- read.csv("SLL_Data.csv")
116
117
118 # Italian Serie A Data
119
120 I00_01 <- read.csv('I0001.csv')
121 I01_02 <- read.csv('I0102.csv')
122 I02_03 <- read.csv('I0203.csv')
123 I03_04 <- read.csv('I0304.csv')
124 I04_05 <- read.csv('I0405.csv')

```

```

125 I05_06 <- read.csv('I0506.csv')
126 I06_07 <- read.csv('I0607.csv')
127 I07_08 <- read.csv('I0708.csv')
128 I08_09 <- read.csv('I0809.csv')
129 I09_10 <- read.csv('I0910.csv')
130 I10_11 <- read.csv('I1011.csv')
131 I11_12 <- read.csv('I1112.csv')
132 I12_13 <- read.csv('I1213.csv')
133 I13_14 <- read.csv('I1314.csv')
134 I14_15 <- read.csv('I1415.csv')
135 I15_16 <- read.csv('I1516.csv')
136 I16_17 <- read.csv('I1617.csv')
137 I17_18 <- read.csv('I1718.csv')
138 I18_19 <- read.csv('I1819.csv')
139 I19_20 <- read.csv('I1920.csv')
140 I20_21 <- read.csv('I2021.csv')
141 I21_22 <- read.csv('I2122.csv')
142 I22_23 <- read.csv('I2223.csv')
143 I23_24 <- read.csv('I2324.csv')
144 I24_25 <- read.csv('I2425.csv')
145
146 I00_25 <- bind_rows(I00_01, I01_02, I02_03, I03_04,
147                     I04_05, I05_06, I06_07, I07_08,
148                     I08_09, I09_10, I10_11, I11_12,
149                     I12_13, I13_14, I15_16, I16_17,
150                     I17_18, I18_19, I19_20, I20_21,
151                     I21_22, I22_23, I23_24, I24_25)
152
153 write.csv(I00_25, "ISA_Data.csv")
154 ISA_Data <- read.csv("ISA_data.csv")
155
156 # French League 1 Data
157
158 F00_01 <- read.csv('F0001.csv')
159 F01_02 <- read.csv('F0102.csv')
160 F02_03 <- read.csv('F0203.csv')
161 F03_04 <- read.csv('F0304.csv')
162 F04_05 <- read.csv('F0405.csv')
163 F05_06 <- read.csv('F0506.csv')
164 F06_07 <- read.csv('F0607.csv')
165 F07_08 <- read.csv('F0708.csv')
166 F08_09 <- read.csv('F0809.csv')

```

```

167 F09_10 <- read.csv('F0910.csv')
168 F10_11 <- read.csv('F1011.csv')
169 F11_12 <- read.csv('F1112.csv')
170 F12_13 <- read.csv('F1213.csv')
171 F13_14 <- read.csv('F1314.csv')
172 F14_15 <- read.csv('F1415.csv')
173 F15_16 <- read.csv('F1516.csv')
174 F16_17 <- read.csv('F1617.csv')
175 F17_18 <- read.csv('F1718.csv')
176 F18_19 <- read.csv('F1819.csv')
177 F19_20 <- read.csv('F1920.csv')
178 F20_21 <- read.csv('F2021.csv')
179 F21_22 <- read.csv('F2122.csv')
180 F22_23 <- read.csv('F2223.csv')
181 F23_24 <- read.csv('F2324.csv')
182 F24_25 <- read.csv('F2425.csv')
183
184 F00_25 <- bind_rows(F00_01,F01_02,F02_03,F03_04,
185                     F04_05,F05_06,F06_07,F07_08,
186                     F08_09,F09_10,F10_11,F11_12,
187                     F12_13,F13_14,F15_16,F16_17,
188                     F17_18,F18_19,F19_20,F20_21,
189                     F21_22,F22_23,F23_24,F24_25)
190
191 write.csv(F00_25,"FL1_Data.csv")
192 FL1_Data <- read.csv("FL1_Data.csv")
193
194 # Combine data from the five leagues
195
196 All_Leagues <- bind_rows(EPL_Data,GB_Data,SLL_Data,ISA_Data,FL1_
197                           Data)
198 All_Leagues <- All_Leagues[,1:29] #Data without Odds included
199 write.csv(All_Leagues,"Big_Five_Data.csv")
200 dim(All_Leagues)

```

Listing 1: Data Preparation R Code

## B Appendix: EDA R Code

```

1 # Load libraries
2 library(tidyverse)
3 library(plotly)

```

```

4
5 # Read the data
6 df <- read_csv("Big_Five_Data.csv")
7
8 # Full Time Home Team Goals (FTHG / HG)
9 summary(df$FTHG)
10 ggplot(df, aes(x=FTHG)) + geom_histogram(binwidth=1, fill='blue',
11     alpha=0.7) + theme_minimal() +
12     labs(title="Distribution of Full Time Home Goals")
13
14 # Full Time Away Team Goals (FTAG / AG)
15 summary(df$FTAG)
16 ggplot(df, aes(x=FTAG)) + geom_histogram(binwidth=1, fill='blue',
17     alpha=0.7) + theme_minimal() +
18     labs(title="Distribution of Full Time Away Goals")
19
20 # Full Time Result (FTR / Res)
21 table(df$FTR)
22 ggplot(df, aes(x=FTR)) + geom_bar(fill='blue', alpha=0.7) + theme
23     _minimal() +
24     labs(title="Full Time Results (H/D/A)")
25
26 # Half Time Home Goals (HTHG)
27 summary(df$HTHG)
28 ggplot(df, aes(x=HTHG)) + geom_histogram(binwidth=1, fill='blue',
29     alpha=0.7) + theme_minimal() +
30     labs(title="Half Time Home Goals")
31
32 # Half Time Away Goals (HTAG)
33 summary(df$HTAG)
34 ggplot(df, aes(x=HTAG)) + geom_histogram(binwidth=1, fill='blue',
35     alpha=0.7) + theme_minimal() +
36     labs(title="Half Time Away Goals")
37
38 # Half Time Result (HTR)
39 table(df$HTR)
40 ggplot(df, aes(x=HTR)) + geom_bar(fill='blue', alpha=0.7) + theme
41     _minimal() +
42     labs(title="Half Time Results (H/D/A)")
43
44 # Attendance
45 summary(df$Attendance)

```

```

40 ggplot(df, aes(x=Attendance)) + geom_histogram(binwidth=5000,
41         fill='blue', alpha=0.7) + theme_minimal() +
42         labs(title="Distribution of Attendance")
43
44 # Home Team Shots (HS)
45 summary(df$HS)
46 ggplot(df, aes(x=HS)) + geom_histogram(binwidth=1, fill='blue',
47         alpha=0.7) + theme_minimal() +
48         labs(title="Home Team Shots Distribution")
49
50 # Away Team Shots (AS)
51 summary(df$AS)
52 ggplot(df, aes(x=AS)) + geom_histogram(binwidth=1, fill='blue',
53         alpha=0.7) + theme_minimal() +
54         labs(title="Away Team Shots Distribution")
55
56 # Home Shots on Target (HST)
57 summary(df$HST)
58 ggplot(df, aes(x=HST)) + geom_histogram(binwidth=1, fill='blue',
59         alpha=0.7) + theme_minimal() +
60         labs(title="Home Team Shots on Target")
61
62 # Away Shots on Target (AST)
63 summary(df$AST)
64 ggplot(df, aes(x=AST)) + geom_histogram(binwidth=1, fill='blue',
65         alpha=0.7) + theme_minimal() +
66         labs(title="Away Team Shots on Target")
67
68 # Home Hit Woodwork (HHW)
69 summary(df$HHW)
70 ggplot(df, aes(x=HHW)) + geom_histogram(binwidth=1, fill='blue',
71         alpha=0.7) + theme_minimal() +
72         labs(title="Home Team Hit Woodwork")
73
74 # Away Hit Woodwork (AHW)
75 summary(df$AHW)
76 ggplot(df, aes(x=AHW)) + geom_histogram(binwidth=1, fill='blue',
77         alpha=0.7) + theme_minimal() +
78         labs(title="Away Team Hit Woodwork")
79
80 # Home Corners (HC)
81 summary(df$HC)

```

```

75 ggplot(df, aes(x=HC)) + geom_histogram(binwidth=1, fill='blue',
76     alpha=0.7) + theme_minimal() +
77     labs(title="Home Team Corners")
78
79 # Away Corners (AC)
80 summary(df$AC)
81 ggplot(df, aes(x=AC)) + geom_histogram(binwidth=1, fill='blue',
82     alpha=0.7) + theme_minimal() +
83     labs(title="Away Team Corners")
84
85 # Home Fouls Committed (HF)
86 summary(df$HF)
87 ggplot(df, aes(x=HF)) + geom_histogram(binwidth=1, fill='blue',
88     alpha=0.7) + theme_minimal() +
89     labs(title="Home Fouls Committed")
90
91 # Away Fouls Committed (AF)
92 summary(df$AF)
93 ggplot(df, aes(x=AF)) + geom_histogram(binwidth=1, fill='blue',
94     alpha=0.7) + theme_minimal() +
95     labs(title="Away Fouls Committed")
96
97 # Home Offsides (HO)
98 summary(df$HO)
99 ggplot(df, aes(x=HO)) + geom_histogram(binwidth=1, fill='blue',
100     alpha=0.7) + theme_minimal() +
101     labs(title="Home Team Offsides")
102
103 # Away Offsides (AO)
104 summary(df$AO)
105 ggplot(df, aes(x=AO)) + geom_histogram(binwidth=1, fill='blue',
106     alpha=0.7) + theme_minimal() +
107     labs(title="Away Team Offsides")
108
109 # Home Yellow Cards (HY)
110 summary(df$HY)
111 ggplot(df, aes(x=HY)) + geom_histogram(binwidth=1, fill='blue',
112     alpha=0.7) + theme_minimal() +
113     labs(title="Home Team Yellow Cards")
114
115 # Away Yellow Cards (AY)

```



```

110 summary(df$AY)
111 ggplot(df, aes(x=AY)) + geom_histogram(binwidth=1, fill='blue',
112     alpha=0.7) + theme_minimal() +
113     labs(title="Away Team Yellow Cards")
114
115 # Home Red Cards (HR)
116 summary(df$HR)
117 ggplot(df, aes(x=HR)) + geom_histogram(binwidth=1, fill='blue',
118     alpha=0.7) + theme_minimal() +
119     labs(title="Home Team Red Cards")
120
121 # Away Red Cards (AR)
122 summary(df$AR)
123 ggplot(df, aes(x=AR)) + geom_histogram(binwidth=1, fill='blue',
124     alpha=0.7) + theme_minimal() +
125     labs(title="Away Team Red Cards")
126
127 # Full Time Result (FTR / Res)
128 p <- ggplot(df, aes(x=FTR)) +
129     geom_bar(fill='blue', alpha=0.7) +
130     theme_minimal() +
131     labs(title="Full Time Results (H/D/A)")
132 ggplotly(p)

```

Listing 2: All Leagues EDA R Code

## C Appendix: All Leagues Summary Statistics R Code

```

1 # ---- Load required libraries ----
2 library(readr)
3 library(dplyr)
4
5
6 # ---- Load datasets ----
7 epl <- read_csv("EPL_Data.csv", locale = locale(encoding = "
8     Latin1"))
9 fl1 <- read_csv("FL1_Data.csv", locale = locale(encoding = "
10    Latin1"))
11 gb <- read_csv("GB_Data.csv", locale = locale(encoding = "
12    Latin1"))
13 isa <- read_csv("ISA_Data.csv", locale = locale(encoding = "
14    Latin1"))

```

```

11 sll <- read_csv("SLL_Data.csv", locale = locale(encoding = "
    Latin1"))
12
13 # ---- Helper function to summarize a league ----
14 league_summary <- function(df, league_name) {
15   total_matches <- nrow(df)
16   home_wins <- sum(df$FTR == "H", na.rm = TRUE)
17   draws <- sum(df$FTR == "D", na.rm = TRUE)
18   away_wins <- sum(df$FTR == "A", na.rm = TRUE)
19
20   avg_home_goals <- mean(df$FTHG, na.rm = TRUE)
21   avg_away_goals <- mean(df$FTAG, na.rm = TRUE)
22   avg_total_goals <- mean(df$FTHG + df$FTAG, na.rm = TRUE)
23
24   data.frame(
25     League = league_name,
26     Home_Wins_Perc = round(home_wins / total_matches * 100, 1),
27     Draws_Perc = round(draws / total_matches * 100, 1),
28     Away_Wins_Perc = round(away_wins / total_matches * 100, 1),
29     Avg_Home_Goals = round(avg_home_goals, 2),
30     Avg_Away_Goals = round(avg_away_goals, 2),
31     Avg_Total_Goals = round(avg_total_goals, 2),
32     Total_Matches = total_matches
33   )
34 }
35
36 # ---- Create summary for all leagues ----
37 summary_list <- list(
38   league_summary/epl, "EPL"),
39   league_summary/fl1, "FL1"),
40   league_summary/gb, "Bundesliga"),
41   league_summary/isa, "Serie A"),
42   league_summary/sll, "La Liga")
43 )
44
45 summary_df <- bind_rows(summary_list)
46
47 # ---- Display the summary ----
48 print(summary_df)

```

Listing 3: Summary Statistics Code

## D Appendix: Home Advantage Inferential Statistics R Code

```
1 # Expects columns: FTR in {"H","D","A"}, FTHG, FTAG, Div (
  optional but used for per-league tests)
2 library(readr)
3 df <- read.csv("Big_Five_Data.csv")
4
5 # --- Basic checks/cleaning ---
6 required_cols <- c("FTR", "FTHG", "FTAG")
7 missing_cols <- setdiff(required_cols, names(df))
8 if (length(missing_cols) > 0) {
9   stop(sprintf("Missing required column(s): %s", paste(missing_
    cols, collapse = ", ")))
10 }
11
12 # Coerce key columns
13 df$FTR <- trimws(as.character(df$FTR))
14 df$FTHG <- suppressWarnings(as.numeric(df$FTHG))
15 df$FTAG <- suppressWarnings(as.numeric(df$FTAG))
16
17 # Filter to rows with valid results and goals
18 ok <- !is.na(df$FTR) & df$FTR %in% c("H","D","A") & !is.na(df$
  FTHG) & !is.na(df$FTAG)
19 dat <- df[ok, , drop = FALSE]
20 if (nrow(dat) == 0) stop("No valid rows after cleaning.")
21
22 # Helper: pretty printing
23 hrule <- function(char = "-") cat(paste0(paste(rep(char, 70),
  collapse=""), "\n"))
24 section <- function(title) { hrule("="); cat(title, "\n"); hrule(
  "=") }
25 subsection <- function(title) { hrule("-"); cat(title, "\n");
  hrule("-") }
26 fmt_p <- function(p) if (is.na(p)) "NA" else if (p < .0001) "< 1e
  -4" else sprintf("%.4g", p)
27 ci_str <- function(ci) sprintf("%.4f, %.4f", ci[1], ci[2])
28
29 # Counts
30 n_total <- nrow(dat)
31 n_H <- sum(dat$FTR == "H")
32 n_D <- sum(dat$FTR == "D")
33 n_A <- sum(dat$FTR == "A")
```

```

34 p_H <- n_H / n_total
35 p_D <- n_D / n_total
36 p_A <- n_A / n_total
37
38 # Goal difference
39 gd <- dat$FTHG - dat$FTAG
40 gd_nonzero <- gd[gd != 0]
41
42 # 1) Chi-square goodness-of-fit: are H/D/A equally likely?
43 section("Global tests (all matches combined)")
44
45 subsection("Distribution tests (H vs D vs A)")
46 chisq_equal <- suppressWarnings(chisq.test(c(n_H, n_D, n_A), p =
47   c(1/3, 1/3, 1/3)))
48 cat("H0: P(H) = P(D) = P(A) = 1/3\n")
49 cat(sprintf("Observed counts (H,D,A): %d, %d, %d | n = %d\n", n
50   _H, n_D, n_A, n_total))
51 cat(sprintf("Chi-square = %.3f, df = %d, p = %s\n\n",
52   chisq_equal$statistic, chisq_equal$parameter, fmt_p(
53     chisq_equal$p.value)))
54
55 # 2) Symmetry test: H vs A equal (draws free). Equivalent to
56   binomial test on non-draws.
57 subsection("Symmetry test (Home vs Away among decisive matches,
58   draws removed)")
59 decisive <- dat$FTR %in% c("H","A")
60 n_H_dec <- sum(dat$FTR[decisive] == "H")
61 n_A_dec <- sum(dat$FTR[decisive] == "A")
62 n_dec <- n_H_dec + n_A_dec
63 bt_sym <- if (n_dec > 0) binom.test(n_H_dec, n_dec, p = 0.5,
64   alternative = "two.sided") else NULL
65 cat("H0: P(Home win | decisive) = 0.5 (same as away)\n")
66 cat(sprintf("Counts among decisive: H=%d, A=%d (n=%d)\n", n_H_dec
67   , n_A_dec, n_dec))
68 if (!is.null(bt_sym)) {
69   cat(sprintf("Binomial test p(two-sided) = %s; Home-win rate (
70     decisive) = %.4f, 95% CI %s\n\n",
71     fmt_p(bt_sym$p.value), n_H_dec/n_dec, ci_str(bt_sym
72       $conf.int)))
73 } else {
74   cat("No decisive matches found.\n\n")
75 }

```

```

67
68 # 3) One-sided advantage test: is  $P(\text{Home win} \mid \text{decisive}) > 0.5$ ?
69 subsection("Directional test:  $P(\text{Home win} \mid \text{decisive}) > 0.5$ ")
70 bt_dir <- if (n_dec > 0) binom.test(n_H_dec, n_dec, p = 0.5,
71   alternative = "greater") else NULL
72 if (!is.null(bt_dir)) {
73   cat(sprintf("Binomial test  $p(\text{greater}) = \%s$ ; effect =  $\%.4f - 0.5$ 
74     =  $\%.4f \backslash \backslash n$ ",
75     fmt_p(bt_dir$p.value), n_H_dec/n_dec, n_H_dec/n_dec
76       - 0.5))
77 } else {
78   cat("No decisive matches found.\n\n")
79 }
80
81 # 4) Two-sample proportion test including draws: is  $P(H) > P(A)$ ?
82 subsection("Two-sample proportion test including draws:  $P(\text{Home win}) > P(\text{Away win})$ ")
83 pt <- suppressWarnings(prop.test(x = c(n_H, n_A), n = c(n_total,
84   n_total), alternative = "greater", correct = TRUE))
85 cat("H0:  $P(H) \leq P(A)$ ; H1:  $P(H) > P(A) \backslash n$ ")
86 cat(sprintf("P(H)= $\%.4f$ , P(A)= $\%.4f$ , diff= $\%.4f \backslash n$ ", p_H, p_A, p_H -
87   p_A))
88 cat(sprintf("Prop.test chi-square= $\%.3f$  (df= $\%d$ ),  $p(\text{greater})=\%s \backslash \backslash n$ 
89   ",
90   pt$statistic, pt$parameter, fmt_p(pt$p.value)))
91
92 # 5) Goal-difference tests (mean/median > 0)
93 subsection("Goal-difference tests (Home goals - Away goals)")
94
95 # (a) One-sample t-test:  $\text{mean}(\text{gd}) > 0$ 
96 tt <- t.test(gd, mu = 0, alternative = "greater")
97 cat(sprintf("t-test:  $\text{mean}(\text{gd}) > 0$  | mean= $\%.4f$ , sd= $\%.4f$ , n= $\%d \backslash n$ ",
98   mean(gd), sd(gd), length(gd)))
99 cat(sprintf("t= $\%.3f$ , df= $\%.1f$ ,  $p(\text{greater})=\%s$ ; 95% CI for mean:
100   [ $\%.4f$ ,  $\%.4f \backslash \backslash n$ ",
101   tt$statistic, tt$parameter, fmt_p(tt$p.value), tt$
102     conf.int[1], tt$conf.int[2]))
103
104 # (b) Wilcoxon signed-rank:  $\text{median}(\text{gd}) > 0$  (nonparametric)
105 if (all(gd == gd[1])) {
106   cat("Wilcoxon skipped: all goal differences equal (degenerate)
107     .\n\n")
108 }

```

```

98 } else {
99   wt <- suppressWarnings(wilcox.test(gd, mu = 0, alternative = "
    greater", exact = FALSE, conf.int = TRUE))
100   cat(sprintf("Wilcoxon signed-rank: median(gd)>0 | p(greater)
    =%s; 95%% CI for median diff: [%.4f, %.4f]\n\n",
101             fmt_p(wt$p.value), wt$conf.int[1], wt$conf.int[2]))
102 }
103
104 # (c) Sign test (binomial on signs, ignoring zeros): P(gd>0) >
    0.5
105 npos <- sum(gd_nonzero > 0)
106 nneg <- sum(gd_nonzero < 0)
107 nsign <- npos + nneg
108 cat(sprintf("Sign counts (gd>0, gd<0): %d, %d | n (nonzero) = %
    d\n", npos, nneg, nsign))
109 if (nsign > 0) {
110   bt_sign <- binom.test(npos, nsign, p = 0.5, alternative = "
    greater")
111   cat(sprintf("Sign test p(greater) = %s; P(gd>0) = %.4f, 95%% CI
    %s\n\n",
112             fmt_p(bt_sign$p.value), npos/nsign, ci_str(bt_sign$
    conf.int)))
113 } else {
114   cat("Sign test skipped: all goal differences are zero.\n\n")
115 }
116
117 # 6) Poisson rate test: total home goals vs total away goals (
    same exposure = #matches)
118 subsection("Poisson rate test: total Home goals vs Away goals per
    match")
119 sum_HG <- sum(dat$FTHG, na.rm = TRUE)
120 sum_AG <- sum(dat$FTAG, na.rm = TRUE)
121 pt_poisson <- poisson.test(c(sum_HG, sum_AG), T = c(n_total, n_
    total), alternative = "greater")
122 cat(sprintf("Totals: Home goals=%d, Away goals=%d (matches=%d
    each)\n", sum_HG, sum_AG, n_total))
123 cat(sprintf("Rate ratio (Home/Away) = %.4f; 95%% CI %s; p(
    greater)=%s\n\n",
124             pt_poisson$estimate, ci_str(pt_poisson$conf.int), fmt
    _p(pt_poisson$p.value)))
125

```

```

126 # 7) Intercept-only logistic model: HomeWin vs Not (Wald test
    that p > 0.5)
127 subsection("Intercept-only logistic model: P(Home win) > 0.5 (
    HomeWin vs Not)")
128 HomeWin <- as.integer(dat$FTR == "H")
129 glm0 <- suppressWarnings(glm(HomeWin ~ 1, family = binomial()))
130 est <- coef(glm0)[1]
131 se <- sqrt(vcov(glm0)[1,1])
132 z <- (est - 0) / se # test logit(p) > 0 iff p > 0.5
133 p_one_sided <- 1 - pnorm(z)
134 p_hat <- plogis(est)
135 ci95 <- plogis(est + c(-1,1) * qnorm(0.975) * se)
136 cat(sprintf("Estimated P(Home win) = %.4f; 95%% CI %s\n", p_hat,
    ci_str(ci95)))
137 cat(sprintf("Wald z = %.3f; p(greater) = %s for H0: P(Home win)
    <= 0.5\n\n", z, fmt_p(p_one_sided)))
138
139 # 8) Multinomial-style constraint test: H vs A equality with
    draws present via simple contrast
140 # This is effectively the binomial test already run on decisive
    matches; we also show a z-test on p_H - p_A.
141 subsection("Approximate z-test on difference P(H) - P(A) (
    including draws)")
142 # Variance under null approx using two-sample proportions with
    equal n
143 phat_diff <- p_H - p_A
144 p_pool <- (n_H + n_A) / (2 * n_total)
145 se_diff <- sqrt(p_pool * (1 - p_pool) * (1/n_total + 1/n_total))
146 z_pa <- ifelse(se_diff > 0, phat_diff / se_diff, NA)
147 p_greater <- ifelse(is.na(z_pa), NA, 1 - pnorm(z_pa))
148 cat(sprintf("Diff = %.4f; z = %.3f; p(greater) = %s\n\n", phat_
    diff, z_pa, fmt_p(p_greater)))
149
150 # ----- Per-league stratified tests (if Div present)
    -----
151 if ("Div" %in% names(dat)) {
152   section("Per-league directional tests (selected)")
153   leagues <- sort(unique(na.omit(dat$Div)))
154   if (length(leagues) == 0) {
155     cat("No league identifiers found in Div.\n")
156   } else {
157     res <- data.frame(

```

```

158     Div = character(),
159     n = integer(),
160     P_H = numeric(),
161     P_A = numeric(),
162     P_H_gt_A_p = numeric(),
163     Mean_GD = numeric(),
164     t_p_gd_gt0 = numeric(),
165     stringsAsFactors = FALSE
166 )
167 for (lg in leagues) {
168     sli <- dat[dat$Div == lg, , drop = FALSE]
169     n <- nrow(sli)
170     nH <- sum(sli$FTR == "H")
171     nA <- sum(sli$FTR == "A")
172     pH <- nH / n
173     pA <- nA / n
174     # prop test H>A (including draws)
175     pt_lg <- suppressWarnings(prop.test(c(nH, nA), c(n, n),
176                                     alternative = "greater", correct = TRUE))
177     # t-test on GD
178     gd_lg <- sli$FTHG - sli$FTAG
179     tt_lg <- suppressWarnings(t.test(gd_lg, mu = 0, alternative
180                                     = "greater"))
181     res <- rbind(res, data.frame(
182         Div = lg, n = n, P_H = pH, P_A = pA,
183         P_H_gt_A_p = as.numeric(pt_lg$p.value),
184         Mean_GD = mean(gd_lg),
185         t_p_gd_gt0 = as.numeric(tt_lg$p.value),
186         stringsAsFactors = FALSE
187     ))
188 }
189 # Print summary table
190 print(res[order(res$P_H_gt_A_p), ], row.names = FALSE)
191 cat("\nSmaller p-values indicate stronger evidence that home
    advantage > away for that league.\n\n")

```

Listing 4: Home Advantage Inferential Code

## E Appendix All Leagues ML Techniques R code



```

1 # Load libraries
2 library(tidyverse)
3 library(caret)
4
5 # Load data
6 data <- read.csv("Big_Five_Data.csv", fileEncoding="latin1")
7
8 # Drop rows with missing FTR
9 data <- data[!is.na(data$FTR), ]
10
11 # Select numeric features for clustering/classification (
    customize as needed)
12 features <- c('FTHG', 'FTAG', 'HTHG', 'HTAG', 'HS', 'AS', 'HST',
    'AST', 'HC', 'AC', 'HF', 'AF', 'HY', 'AY', 'HR', 'AR')
13 data_model <- data[, c(features, 'FTR')]
14
15 # Remove rows with NAs in selected columns
16 data_model <- na.omit(data_model)
17
18 library(randomForest)
19
20 # Encode FTR as factor
21 data_model$FTR <- as.factor(data_model$FTR)
22
23 # Split data (70% train, 30% test)
24 set.seed(7406)
25 train_idx <- createDataPartition(data_model$FTR, p=0.7, list=
    FALSE)
26 train <- data_model[train_idx, ]
27 test <- data_model[-train_idx, ]
28
29 # Train Random Forest
30 rf_model <- randomForest(FTR ~ ., data=train, ntree=100)
31 # Predict
32 pred <- predict(rf_model, test)
33 # Confusion Matrix
34 confusionMatrix(pred, test$FTR)
35
36 # Standardize numeric features
37 data_num <- scale(data_model[, features])
38
39 # Run k-means (choose k, e.g., 3)

```

```

40 set.seed(7406)
41 kmeans_res <- kmeans(data_num, centers=3, nstart=25)
42
43 # Add cluster label to data
44 data_model$cluster <- as.factor(kmeans_res$cluster)
45
46 # See if clusters relate to FTR
47 table(data_model$cluster, data_model$FTR)
48
49
50 library(caret)
51 library(randomForest)
52
53 # Data preparation as before
54 data_model$FTR <- as.factor(data_model$FTR)
55 set.seed(7406)
56
57 # 5-fold cross-validation
58 ctrl <- trainControl(method = "cv", number = 5)
59
60 # Tune mtry parameter
61 rf_grid <- expand.grid(mtry = c(2, 4, 6, 8, 10))
62
63 # Train and tune Random Forest
64 rf_tuned <- train(
65   FTR ~ .,
66   data = data_model,
67   method = "rf",
68   trControl = ctrl,
69   tuneGrid = rf_grid,
70   ntree = 100 # set as needed
71 )
72
73 # Best mtry value
74 print(rf_tuned$bestTune)
75
76 # Evaluation results (accuracy, kappa, etc.)
77 print(rf_tuned)
78
79 # Data preparation as before; kNN works best with scaled features
80 preProcValues <- preProcess(data_model[, features], method = c("
  center", "scale"))

```

```

81 data_model_knn <- data_model
82 data_model_knn[, features] <- predict(preProcValues, data_model[,
    features])
83
84 # 5-fold cross-validation
85 ctrl <- trainControl(method = "cv", number = 5)
86
87 # Tune k
88 knn_grid <- expand.grid(k = c(3, 5, 7, 9, 11))
89
90 # Train and tune kNN
91 knn_tuned <- train(
92   FTR ~ .,
93   data = data_model_knn,
94   method = "knn",
95   trControl = ctrl,
96   tuneGrid = knn_grid
97 )
98
99 # Best k value
100 print(knn_tuned$bestTune)
101
102 # Evaluation results (accuracy, kappa, etc.)
103 print(knn_tuned)
104
105 # Example for Random Forest
106 rf_pred <- predict(rf_tuned, newdata = data_model)
107 confusionMatrix(rf_pred, data_model$FTR)
108
109 library(xgboost)
110 library(caret)
111 library(Matrix)
112
113 # Ensure target is a factor and create numeric labels for XGBoost
114 data_model$FTR <- as.factor(data_model$FTR)
115 labels <- as.numeric(data_model$FTR) - 1 # XGBoost expects
    labels from 0
116
117 # One-hot encode features for XGBoost
118 dummies <- dummyVars(FTR ~ ., data = data_model)
119 data_matrix <- predict(dummies, newdata = data_model)

```

```

120 dtrain <- xgb.DMatrix(data = as.matrix(data_matrix), label =
      labels)
121
122 # Split into training and test sets (e.g., 70/30)
123 set.seed(7406)
124 train_idx <- createDataPartition(labels, p=0.7, list=FALSE)
125 dtrain_train <- xgb.DMatrix(data = as.matrix(data_matrix[train_
      idx,]), label = labels[train_idx])
126 dtrain_test <- xgb.DMatrix(data = as.matrix(data_matrix[-train_
      idx,]), label = labels[-train_idx])
127
128 # Train XGBoost model (multiclass, softmax)
129 num_class <- length(levels(data_model$FTR))
130 params <- list(
131   objective = "multi:softmax",
132   num_class = num_class,
133   eval_metric = "merror"
134 )
135
136 bst <- xgb.train(
137   params = params,
138   data = dtrain_train,
139   nrounds = 100,
140   watchlist = list(eval = dtrain_test, train = dtrain_train),
141   verbose = 1
142 )
143
144 # Predict and evaluate
145 preds <- predict(bst, dtrain_test)
146 confusionMatrix(
147   factor(preds, levels = 0:(num_class-1), labels = levels(data_
      model$FTR)),
148   factor(labels[-train_idx], levels = 0:(num_class-1), labels =
      levels(data_model$FTR))
149 )

```

Listing 5: Machine Learning Techniques Code

## F Appendix Multinomial Regression R Code

```

1 library(readr)
2 library(caret)

```

```

3 library(nnet)
4
5 data <- read_csv("Big_Five_Data.csv")
6 data[] <- lapply(data, function(x) if(is.character(x)) as.factor(
7   x) else x)
8 data$FTR <- as.factor(data$FTR)
9
10 set.seed(7406)
11 trainIndex <- createDataPartition(data$FTR, p = 0.7, list = FALSE
12   )
13 trainData <- data[trainIndex, ]
14 testData <- data[-trainIndex, ]
15
16 # Show number of unique levels for each factor column
17 sapply(trainData, function(x) if(is.factor(x)) length(unique(x))
18   else NA)
19
20 # Let's pick only numeric columns or low-cardinality factors
21 is_low_cardinality <- function(x) (is.factor(x) && length(unique(
22   x)) <= 4) || is.numeric(x)
23 low_card_cols <- names(Filter(is_low_cardinality, trainData))
24 low_card_cols <- setdiff(low_card_cols, "FTR")
25 print(low_card_cols)
26 low_card_cols <- setdiff(low_card_cols, "Attendance")
27 low_card_cols <- setdiff(low_card_cols, "X.5")
28 low_card_cols <- setdiff(low_card_cols, "...1")
29 print(low_card_cols)
30
31 # If you have at least two such predictors, try:
32 formula <- as.formula(paste("FTR ~", paste(low_card_cols,
33   collapse = " + ")))
34 multinom_model <- multinom(formula, data = trainData)
35 summary(multinom_model)

```

Listing 6: Multinomial Regression Code

## G Appendix All Leagues ML Techniques R code

```

1 # Load libraries
2 library(tidyverse)
3 library(caret)
4
5 # Load data

```

```

6 data <- read.csv("Big_Five_Data.csv", fileEncoding="latin1")
7
8 # Drop rows with missing FTR
9 data <- data[!is.na(data$FTR), ]
10
11 # Select numeric features for clustering/classification (
    customize as needed)
12 features <- c('FTHG', 'FTAG', 'HTHG', 'HTAG', 'HS', 'AS', 'HST',
    'AST', 'HC', 'AC', 'HF', 'AF', 'HY', 'AY', 'HR', 'AR')
13 data_model <- data[, c(features, 'FTR')]
14
15 # Remove rows with NAs in selected columns
16 data_model <- na.omit(data_model)
17
18 library(randomForest)
19
20 # Encode FTR as factor
21 data_model$FTR <- as.factor(data_model$FTR)
22
23 # Split data (70% train, 30% test)
24 set.seed(7406)
25 train_idx <- createDataPartition(data_model$FTR, p=0.7, list=
    FALSE)
26 train <- data_model[train_idx, ]
27 test <- data_model[-train_idx, ]
28
29 # Train Random Forest
30 rf_model <- randomForest(FTR ~ ., data=train, ntree=100)
31 # Predict
32 pred <- predict(rf_model, test)
33 # Confusion Matrix
34 confusionMatrix(pred, test$FTR)
35
36 # Standardize numeric features
37 data_num <- scale(data_model[, features])
38
39 # Run k-means (choose k, e.g., 3)
40 set.seed(7406)
41 kmeans_res <- kmeans(data_num, centers=3, nstart=25)
42
43 # Add cluster label to data
44 data_model$cluster <- as.factor(kmeans_res$cluster)

```

```

45
46 # See if clusters relate to FTR
47 table(data_model$cluster, data_model$FTR)
48
49
50 library(caret)
51 library(randomForest)
52
53 # Data preparation as before
54 data_model$FTR <- as.factor(data_model$FTR)
55 set.seed(7406)
56
57 # 5-fold cross-validation
58 ctrl <- trainControl(method = "cv", number = 5)
59
60 # Tune mtry parameter
61 rf_grid <- expand.grid(mtry = c(2, 4, 6, 8, 10))
62
63 # Train and tune Random Forest
64 rf_tuned <- train(
65   FTR ~ .,
66   data = data_model,
67   method = "rf",
68   trControl = ctrl,
69   tuneGrid = rf_grid,
70   ntree = 100 # set as needed
71 )
72
73 # Best mtry value
74 print(rf_tuned$bestTune)
75
76 # Evaluation results (accuracy, kappa, etc.)
77 print(rf_tuned)
78
79 # Data preparation as before; kNN works best with scaled features
80 preProcValues <- preProcess(data_model[, features], method = c("
   center", "scale"))
81 data_model_knn <- data_model
82 data_model_knn[, features] <- predict(preProcValues, data_model[,
   features])
83
84 # 5-fold cross-validation

```

```

85 ctrl <- trainControl(method = "cv", number = 5)
86
87 # Tune k
88 knn_grid <- expand.grid(k = c(3, 5, 7, 9, 11))
89
90 # Train and tune kNN
91 knn_tuned <- train(
92   FTR ~ .,
93   data = data_model_knn,
94   method = "knn",
95   trControl = ctrl,
96   tuneGrid = knn_grid
97 )
98
99 # Best k value
100 print(knn_tuned$bestTune)
101
102 # Evaluation results (accuracy, kappa, etc.)
103 print(knn_tuned)
104
105 # Example for Random Forest
106 rf_pred <- predict(rf_tuned, newdata = data_model)
107 confusionMatrix(rf_pred, data_model$FTR)
108
109 library(xgboost)
110 library(caret)
111 library(Matrix)
112
113 # Ensure target is a factor and create numeric labels for XGBoost
114 data_model$FTR <- as.factor(data_model$FTR)
115 labels <- as.numeric(data_model$FTR) - 1 # XGBoost expects
      labels from 0
116
117 # One-hot encode features for XGBoost
118 dummies <- dummyVars(FTR ~ ., data = data_model)
119 data_matrix <- predict(dummies, newdata = data_model)
120 dtrain <- xgb.DMatrix(data = as.matrix(data_matrix), label =
      labels)
121
122 # Split into training and test sets (e.g., 70/30)
123 set.seed(7406)
124 train_idx <- createDataPartition(labels, p=0.7, list=FALSE)

```



```

125 dtrain_train <- xgb.DMatrix(data = as.matrix(data_matrix[train_
      idx,]), label = labels[train_idx])
126 dtrain_test <- xgb.DMatrix(data = as.matrix(data_matrix[-train_
      idx,]), label = labels[-train_idx])
127
128 # Train XGBoost model (multiclass, softmax)
129 num_class <- length(levels(data_model$FTR))
130 params <- list(
131   objective = "multi:softmax",
132   num_class = num_class,
133   eval_metric = "merror"
134 )
135
136 bst <- xgb.train(
137   params = params,
138   data = dtrain_train,
139   nrounds = 100,
140   watchlist = list(eval = dtrain_test, train = dtrain_train),
141   verbose = 1
142 )
143
144 # Predict and evaluate
145 preds <- predict(bst, dtrain_test)
146 confusionMatrix(
147   factor(preds, levels = 0:(num_class-1), labels = levels(data_
      model$FTR)),
148   factor(labels[-train_idx], levels = 0:(num_class-1), labels =
      levels(data_model$FTR))
149 )

```

Listing 7: Machine Learning Techniques Code

## H Appendix: Machine Learning R Code

```

1 # ---- Load Required Libraries ----
2 library(caTools)      # For train/test split
3 library(randomForest) # Random Forest
4 library(e1071)        # SVM
5 library(class)        # kNN
6 library(xgboost)      # XGBoost
7
8 # ---- Data Loading and Cleaning ----

```

```

9 df <- read.csv("Big_Five_Data.csv", fileEncoding = "latin1")
10
11 # Select key numeric columns and target
12 columns <- c("FTR", "FTHG", "FTAG", "HTHG", "HTAG",
13             "HS", "AS", "HC", "AC", "HF", "AF")
14 df <- df[, columns]
15 df <- na.omit(df)
16
17 # ---- Train/Test Split ----
18 set.seed(7406)
19 split <- sample.split(df$FTR, SplitRatio = 0.7)
20 train <- subset(df, split == TRUE)
21 test  <- subset(df, split == FALSE)
22
23 #####
24 # ---- 1. Logistic Regression ----
25 # NOTE: This is set up as a binary example (H vs A); modify as
26 #       needed for your use case.
27 # For simplicity, we filter to H/A and drop draws.
28 # Filter to binary H vs A example
29 bin_train <- subset(train, FTR %in% c("H", "A"))
30 bin_test  <- subset(test,  FTR %in% c("H", "A"))
31
32 # Convert to factor (drops unused levels automatically)
33 bin_train$FTR <- factor(bin_train$FTR)
34 bin_test$FTR  <- factor(bin_test$FTR)
35
36 logit_model <- glm(
37   FTR ~ FTHG + FTAG + HTHG + HTAG + HS + AS + HC + AC + HF + AF,
38   data    = bin_train,
39   family  = "binomial"
40 )
41
42 logit_pred <- predict(logit_model, newdata = bin_test, type = "
43   response")
44
45 logit_pred_class <- ifelse(logit_pred > 0.5, "H", "A")
46
47 cat("Logistic Regression (binary H vs A):\n")
48 print(table(Predicted = logit_pred_class, Actual = bin_test$FTR))
49 #####

```

```

49 # ---- 2. Random Forest Classification ----
50 train$FTR <- as.factor(train$FTR)
51 test$FTR  <- as.factor(test$FTR)
52
53 rf_model <- randomForest(FTR ~ ., data = train, ntree = 100)
54 rf_pred  <- predict(rf_model, newdata = test)
55
56 cat("\nRandom Forest Classification:\n")
57 print(table(Predicted = rf_pred, Actual = test$FTR))
58
59 # ---- NEW: Random Forest Tuning Plot (mtry vs OOB error) ----
60 set.seed(7406)
61 rf_tune <- tuneRF(
62   x      = train[, -1],      # predictors only
63   y      = train$FTR,        # factor response
64   ntreeTry = 100,
65   stepFactor = 2,
66   improve  = 0.01,
67   trace    = TRUE,
68   plot     = TRUE            # this produces the tuning
69                               plot
70 )
71 # 'rf_tune' contains mtry and OOB error; the plot shows OOB error
72   vs mtry.
73
74 #####
75 # ---- 3. Support Vector Machine (SVM) ----
76 svm_model <- svm(FTR ~ ., data = train)
77 svm_pred  <- predict(svm_model, newdata = test)
78
79 cat("\nSVM Classification:\n")
80 print(table(Predicted = svm_pred, Actual = test$FTR))
81
82 # ---- NEW: SVM Tuning Plot (cost & gamma) ----
83 # We use e1071::tune.svm to perform grid search with CV and then
84   plot.
85 set.seed(7406)
86 svm_tune <- tune.svm(
87   FTR ~ .,
88   data  = train,
89   kernel = "radial",
90   cost   = c(0.1, 1, 10),

```

```

88   gamma   = c(0.01, 0.05, 0.1)
89 )
90
91 cat("\nBest SVM Parameters:\n")
92 print(svm_tune$best.parameters)
93 cat("\nBest SVM Performance:\n")
94 print(svm_tune$best.performance)
95
96 # This produces a lattice plot of performance across the (cost,
   gamma) grid.
97 plot(svm_tune)
98
99 #####
100 # ---- 4. k-Nearest Neighbors (kNN) ----
101 train_X <- train[, -1]
102 test_X  <- test[, -1]
103 train_y <- train$FTR
104 test_y  <- test$FTR
105
106 knn_pred <- knn(train_X, test_X, train_y, k = 5)
107
108 cat("\nkNN Classification:\n")
109 print(table(Predicted = knn_pred, Actual = test_y))
110
111 #####
112 # ---- 5. Linear Regression (predicting FTHG) ----
113 linreg_train <- train[, -1] # Exclude FTR
114 linreg_test  <- test[, -1]
115
116 lm_model <- lm(FTHG ~ ., data = linreg_train)
117 lm_pred  <- predict(lm_model, newdata = linreg_test)
118 mse_lm   <- mean((lm_pred - linreg_test$FTHG)^2)
119
120 cat("\nLinear Regression (MSE for FTHG):", mse_lm, "\n")
121
122 #####
123 # ---- 6. Random Forest Regression (predicting FTHG) ----
124 rf_reg_model <- randomForest(FTHG ~ ., data = linreg_train, ntree
   = 100)
125 rf_reg_pred  <- predict(rf_reg_model, newdata = linreg_test)
126 mse_rf       <- mean((rf_reg_pred - linreg_test$FTHG)^2)
127

```

```

128 cat("\nRandom Forest Regression (MSE for FTHG):", mse_rf, "\n")
129
130 #####
131 # ---- 7. XGBoost Regression (predicting FTHG) + Training Curves
132     ----
133
134 # Prepare matrices (drop FTHG from features)
135 train_matrix <- as.matrix(linreg_train[, !(names(linreg_train) %
136     in% "FTHG")])
137 test_matrix  <- as.matrix(linreg_test[,  !(names(linreg_test)  %
138     in% "FTHG")])
139
140 dtrain <- xgb.DMatrix(data = train_matrix, label = linreg_train$
141     FTHG)
142 dtest  <- xgb.DMatrix(data = test_matrix,  label = linreg_test$
143     FTHG)
144
145 params <- list(
146     objective   = "reg:squarederror",
147     eval_metric = "rmse"
148 )
149
150 set.seed(7406)
151 xgb_model <- xgb.train(
152     params      = params,
153     data        = dtrain,
154     nrounds     = 50,
155     watchlist   = list(train = dtrain, test = dtest),
156     print_every_n = 10
157 )
158
159 # Predictions and MSE
160 xgb_pred <- predict(xgb_model, newdata = dtest)
161 mse_xgb  <- mean((xgb_pred - linreg_test$FTHG)^2)
162
163 cat("\nXGBoost Regression (MSE for FTHG):", mse_xgb, "\n")
164
165 # ---- NEW: XGBoost Training Curves (RMSE vs. rounds) ----
166 eval_log <- xgb_model$evaluation_log
167 # Base R plot for train/test RMSE across boosting rounds
168 plot(
169     eval_log$iter, eval_log$train_rmse,

```

```

165     type = "l", col = "blue", lwd = 2,
166     xlab = "Boosting Round",
167     ylab = "RMSE",
168     main = "XGBoost RMSE vs. Number of Rounds"
169 )
170 lines(
171     eval_log$iter, eval_log$test_rmse,
172     col = "red", lwd = 2, lty = 2
173 )
174 legend(
175     "topright",
176     legend = c("Train RMSE", "Test RMSE"),
177     col     = c("blue", "red"),
178     lty     = c(1, 2),
179     lwd     = 2,
180     bty     = "n"
181 )
182
183 #####
184 # ---- 8. k-Means Clustering (Unsupervised) ----
185 # Use only numeric predictors (drop FTR)
186 df_num <- na.omit(df[, -1])
187
188 set.seed(7406)
189 kmeans_model <- kmeans(df_num, centers = 3)
190
191 cat("\nk-Means Clustering (cluster sizes):\n")
192 print(table(kmeans_model$cluster))
193
194 #####
195 # ---- 9. Principal Component Analysis (PCA) ----
196 pca <- prcomp(df_num, scale. = TRUE)
197
198 cat("\nPCA Summary:\n")
199 print(summary(pca))
200
201 # Scree-type plot of variance explained
202 plot(pca, main = "PCA - Variance Explained")

```

Listing 8: Machine Learning Techniques Code

## I Appendix Time Series Analysis R Code

```
1 # =====
2 # Time-Series Analysis on Football-Data Big Five Leagues
3 # =====
4
5 library(readr)
6 library(dplyr)
7 library(lubridate)
8 library(stringr)
9 library(tsibble)
10 library(feasts)
11 library(fable)
12 library(tseries)
13 library(fabletools) # for interpolate()
14 library(purrr)
15 library(rlang)
16
17 raw <- read_csv("Big_Five_Data.csv")
18
19 df <- raw %>%
20   # Standardize & parse date (Football Data uses dd/mm/yy)
21   mutate(
22     Date = dmy(Date),
23     Div = as.factor(Div),
24     FTHG = as.numeric(FTHG),
25     FTAG = as.numeric(FTAG),
26     FTR = as.factor(FTR) # H/D/A
27   ) %>%
28   filter(!is.na(Date)) %>%
29   select(Date, Div, HomeTeam, AwayTeam, FTHG, FTAG, FTR)
30
31 # Build weekly league series
32
33 weekly_league <- df %>%
34   mutate(week = yearweek(Date)) %>%
35   group_by(Div, week) %>%
36   summarise(
37     matches = n(),
38     home_wins = sum(FTR == "H", na.rm = TRUE),
39     .groups = "drop"
40   ) %>%
```

```

41 as_tsibble(index = week, key = Div) %>%
42 group_by_key() %>%
43 # Make all missing weeks explicit
44 fill_gaps(matches = 0L, home_wins = 0L) %>%
45 ungroup() %>%
46 mutate(
47   # define the rate only when matches > 0
48   home_win_rate = if_else(matches > 0, home_wins / matches, NA_
      real_)
49 )
50
51
52 # Choose a league (or loop over all) and IMPUTE missing weeks
53 league_code <- "E0" # change as needed
54 e0 <- weekly_league %>%
55   filter(Div == league_code) %>%
56   arrange(week)
57
58 # If too many NAs or too short, try a monthly aggregation
   fallback
59 need_monthly_fallback <- function(tbl) {
60   n_non_na <- sum(!is.na(tbl$home_win_rate))
61   n_non_na < 30 # STL likes at least a few dozen obs
62 }
63
64 if (need_monthly_fallback(e0)) {
65   message("Weekly series too sparse for STL; using MONTHLY
      aggregation for ", league_code)
66   e0 <- df %>%
67     mutate(month = yearmonth(Date)) %>%
68     group_by(Div, month) %>%
69     summarise(
70       matches = n(),
71       home_wins = sum(FTR == "H", na.rm = TRUE),
72       .groups = "drop"
73     ) %>%
74     as_tsibble(index = month, key = Div) %>%
75     group_by_key() %>%
76     fill_gaps(matches = 0L, home_wins = 0L) %>%
77     ungroup() %>%
78     mutate(home_win_rate = if_else(matches > 0, home_wins /
      matches, NA_real_)) %>%

```



```

79   filter(Div == league_code) %>%
80   arrange(month) %>%
81   rename(week = month) # keep a common index name for
      downstream code
82 }
83
84 # Fit a simple ARIMA on the (possibly NA) series, then
      interpolate missing values
85 imp_fit <- e0 %>% model(ARIMA(home_win_rate))
86 e0_imp <- interpolate(imp_fit, e0)
87
88 # Guard: if everything is still NA or too short, stop with a
      clear message
89 n_non_na_imp <- e0_imp %>%
90   as_tibble() %>%
91   pull(home_win_rate) %>%
92   {sum(!is.na(.))}
93 if (n_non_na_imp < 30) {
94   stop("Not enough non-missing points after interpolation for STL
          in ", league_code,
95         ". Consider a coarser aggregation or a different metric.")
96 }
97
98 # STL decomposition on IMPUTED series
99 dcmp <- e0_imp %>%
100   model(STL(home_win_rate ~ season(window = "periodic"))) %>%
101   components()
102
103 # Example plot in an interactive session:
104 autoplot(dcmp)
105
106 # Forecast with ARIMA
107 fit <- e0_imp %>% model(ARIMA(home_win_rate))
108
109 fc <- forecast(fit, h = "26 weeks")
110
111 # Example plots:
112 autoplot(e0_imp, home_win_rate) + autolayer(fc)
113
114 # loop over ALL leagues
115 safe_stl <- function(tbl) {
116

```

```

117   imp <- tbl %>% model(ARIMA(home_win_rate))
118   tbl_imp <- interpolate(imp, tbl)
119   if (sum(!is.na(tbl_imp$home_win_rate)) < 30) return(NULL) #
      skip sparse series
120   mdl <- tbl_imp %>% model(STL(home_win_rate ~ season(window = "
      periodic"))))
121   suppressMessages(components(mdl))
122 }
123
124 all_decomps <- weekly_league %>%
125   group_by_key() %>%
126   group_map(~ safe_stl(.x), .keep = TRUE)

```

Listing 9: TSA Code

## J Appendix Winning Teams R Code

```

1  library(readr)
2  library(dplyr)
3  library(lubridate)
4  library(stringr)
5
6  # --- Helper to load and prep each league's data ---
7  load_league <- function(filename, league_code, league_name) {
8    df <- read_csv(filename, locale = locale(encoding = "Latin1"))
9    df <- df %>%
10      mutate(
11        League = league_name,
12        Date = dmy(Date),
13        Season = ifelse(month(Date) >= 8, paste0(year(Date), "/",
          year(Date) + 1), paste0(year(Date) - 1, "/", year(Date))
14      )
15    return(df)
16  }
17
18  # --- Load all five leagues ---
19  epl <- load_league("EPL_Data.csv", "EPL", "Premier League")
20  fl1 <- load_league("FL1_Data.csv", "FL1", "Ligue 1")
21  gb <- load_league("GB_Data.csv", "GB", "Bundesliga")
22  isa <- load_league("ISA_Data.csv", "ISA", "Serie A")
23  sll <- load_league("SLL_Data.csv", "SLL", "La Liga")

```

```

24
25 all_leagues <- bind_rows/epl, fl1, gb, isa, sll)
26
27 # --- Get all unique seasons ---
28 all_seasons <- sort(unique(all_leagues$Season))
29
30 # --- Loop through each season ---
31 results <- lapply(all_seasons, function(season_target) {
32
33   season_data <- all_leagues %>% filter(Season == season_target)
34
35   win_stats <- season_data %>%
36     mutate(
37       Winner = case_when(
38         FTR == "H" ~ HomeTeam,
39         FTR == "A" ~ AwayTeam,
40         TRUE ~ NA_character_
41       )
42     ) %>%
43     filter(!is.na(Winner)) %>%
44     group_by(League, Winner) %>%
45     summarize(Wins = n(), .groups = "drop")
46
47   winning_teams <- win_stats %>%
48     group_by(League) %>%
49     slice_max(order_by = Wins, n = 1, with_ties = FALSE) %>%
50     ungroup() %>%
51     mutate(Season = season_target)
52
53   return(winning_teams)
54 })
55
56 # --- Combine into one data frame ---
57 final_results <- bind_rows(results)
58
59 print(final_results)
60
61 # --- Save to CSV with all seasons ---
62 write_csv(final_results, "winning_teams_all_seasons.csv")
63
64 cat("See 'winning_teams_all_seasons.csv' for results across all
    25 seasons.\n")

```

```

65
66 # =====
67 #   Visualizations for winning_teams_all_seasons.csv
68 #   - Line: Wins of each league's winner across seasons
69 #   - Bar: Title counts per team (by league)
70 #   - Heatmap: Winner by season & league
71 #   - Boxplot: Distribution of winner wins by league
72 # =====
73
74 # ---- Package setup ----
75 install_if_missing <- function(pkgs) {
76   to_install <- pkgs[!pkgs %in% rownames(installed.packages())]
77   if (length(to_install)) install.packages(to_install,
78     dependencies = TRUE)
79 }
79 install_if_missing(c("readr", "dplyr", "ggplot2", "tidyr", "
  forcats"))
80
81 library(readr)
82 library(dplyr)
83 library(ggplot2)
84 library(tidyr)
85 library(forcats)
86
87 # ---- Paths ----
88 input_csv <- "winning_teams_all_seasons.csv"
89 out_dir <- "plots"
90 if (!dir.exists(out_dir)) dir.create(out_dir, recursive = TRUE)
91
92 # ---- Load data ----
93 winners <- read_csv(input_csv, show_col_types = FALSE)
94
95 # Expecting columns: League, Winner, Wins, Season
96 required_cols <- c("League", "Winner", "Wins", "Season")
97 missing_cols <- setdiff(required_cols, names(winners))
98 if (length(missing_cols)) {
99   stop(sprintf("Missing columns in CSV: %s", paste(missing_cols,
100     collapse = ", ")))
101 }
102
103 # ---- Season ordering (chronological) ----
104 # Seasons formatted like "2010/2011" -> order by first year

```

```

104 season_start_year <- suppressWarnings(as.numeric(substr(winners$
    Season, 1, 4)))
105 season_levels <- winners %>%
106   mutate(start_year = season_start_year) %>%
107   distinct(Season, start_year) %>%
108   arrange(start_year, Season) %>%
109   pull(Season)
110
111 winners <- winners %>%
112   mutate(
113     Season = factor(Season, levels = season_levels),
114     League = as.factor(League),
115     Winner = as.factor(Winner)
116   )
117
118 # -----
119 # 1) Line: Wins by league winners across seasons
120 # -----
121 p_line <- ggplot(winners, aes(x = Season, y = Wins, group =
    League, color = League)) +
122   geom_line(linewidth = 1) +
123   geom_point() +
124   labs(
125     title = "Number of Wins by League Winners Across Seasons",
126     x = "Season",
127     y = "Wins"
128   ) +
129   theme_minimal(base_size = 12) +
130   theme(
131     axis.text.x = element_text(angle = 45, hjust = 1),
132     panel.grid.minor = element_blank()
133   )
134
135 ggsave(file.path(out_dir, "wins_by_winners_line.png"), p_line,
    width = 12, height = 7, dpi = 150)
136
137 # -----
138 # 2) Bar: Dominance -- total titles per team (faceted by league)
139 # -----
140 team_titles <- winners %>%
141   group_by(League, Winner) %>%
142   summarize(Titles = n(), .groups = "drop") %>%

```

```

143   arrange(League, desc(Titles))
144
145 # Order teams within each league by titles for clearer bars
146 team_titles <- team_titles %>%
147   group_by(League) %>%
148   mutate(Winner = fct_reorder(Winner, Titles, .desc = TRUE)) %>%
149   ungroup()
150
151 p_bar <- ggplot(team_titles, aes(x = Winner, y = Titles, fill =
152   League)) +
153   geom_col() +
154   coord_flip() +
155   facet_wrap(~ League, scales = "free_y") +
156   labs(
157     title = "Total Titles Won by Teams (All Seasons)",
158     x = "Team",
159     y = "Number of Titles"
160   ) +
161   theme_minimal(base_size = 12) +
162   theme(
163     legend.position = "none",
164     panel.grid.minor = element_blank()
165   )
166 ggsave(file.path(out_dir, "titles_per_team_faceted.png"), p_bar,
167   width = 12, height = 8, dpi = 150)
168
169 # -----
170 # 3) Heatmap: Winner by season & league
171 # -----
172 # (Optional) shorten long team names for readability
173 shorten <- function(x, maxlen = 18) ifelse(nchar(x) > maxlen,
174   paste0(substr(x, 1, maxlen - 1), "..."), x)
175 winners$Winner_short <- shorten(as.character(winners$Winner))
176
177 p_heat <- ggplot(winners, aes(x = Season, y = League, fill =
178   Winner_short)) +
179   geom_tile(color = "white", linewidth = 0.2) +
180   labs(
181     title = "League Winners by Season",
182     x = "Season",
183     y = "League",

```

```

181     fill = "Winning Team"
182 ) +
183 theme_minimal(base_size = 12) +
184 theme(
185     axis.text.x = element_text(angle = 90, vjust = 0.5, hjust =
186         1),
187     panel.grid = element_blank()
188 )
189 ggsave(file.path(out_dir, "winners_heatmap.png"), p_heat, width =
190     14, height = 5.5, dpi = 150)
191 # -----
192 # 4) Boxplot: Distribution of winner wins by league
193 # -----
194 p_box <- ggplot(winners, aes(x = League, y = Wins, fill = League)
195 ) +
196 geom_boxplot(alpha = 0.75, outlier.shape = 21) +
197 labs(
198     title = "Distribution of Wins for League Winners",
199     x = "League",
200     y = "Wins"
201 ) +
202 theme_minimal(base_size = 12) +
203 theme(
204     legend.position = "none",
205     panel.grid.minor = element_blank()
206 )
207 ggsave(file.path(out_dir, "winner_wins_boxplot.png"), p_box,
208     width = 9, height = 6, dpi = 150)
209 # -----
210 # Done
211 # -----
212 message("Saved plots to: ", normalizePath(out_dir))

```

Listing 10: Winning Teams Code

## K Appendix EPL EDA R Code

```

1 library(tidyverse)

```

```

2 library(DataExplorer)
3 library(naniar)
4 library(plotly)
5
6 # Load data
7 epl_data <- read.csv("EPL_Data.csv", fileEncoding = "latin1",
8   stringsAsFactors = FALSE)
9 str(epl_data)
10
11 # Percentage missing per column
12 missing_summary <- apply(epl_data, function(x) mean(is.na(x)))
13   %>% sort(decreasing = TRUE)
14 missing_df <- data.frame(Column = names(missing_summary),
15   PercentMissing = missing_summary)
16 head(missing_df, 15)
17
18 # Visualize missing data pattern
19 plot_missing(epl_data)
20
21 # Get summaries for numeric columns
22 numeric_data <- select(epl_data, where(is.numeric))
23 summary(numeric_data)
24
25 # Convert Date to Date format
26 epl_data$Date <- as.Date(epl_data$Date, format="%d/%m/%y")
27 # Remove NA dates if present
28 epl_data <- epl_data %>% filter(!is.na(Date))
29
30 # Plot total goals per season
31 epl_data %>%
32   mutate(Season = lubridate::year(Date)) %>%
33   group_by(Season) %>%
34   summarise(
35     Total_Home_Goals = sum(FTHG, na.rm = TRUE),
36     Total_Away_Goals = sum(FTAG, na.rm = TRUE)
37   ) %>%
38   pivot_longer(cols = starts_with("Total"), names_to = "Type",
39     values_to = "Goals") %>%
40   ggplot(aes(x = Season, y = Goals, fill = Type)) +
41   geom_col(position = "dodge") +
42   labs(title = "Total Home/Away Goals per Season", x = "Season",
43     y = "Goals")

```



```

39
40 ggplot/epl_data, aes(x = FTHG)) +
41   geom_histogram(bins = 20) +
42   labs(title = "Distribution of Full-Time Home Goals", x = "Goals
    ", y = "Frequency")
43
44 ggplot/epl_data, aes(x = FTAG)) +
45   geom_histogram(bins = 20) +
46   labs(title = "Distribution of Full-Time Away Goals", x = "Goals
    ", y = "Frequency")
47
48 # Most common home teams
49/epl_data %>%
50   count(HomeTeam, sort = TRUE) %>%
51   head(10)

```

Listing 11: EPL EDA Code

## L Appendix Liverpool Analysis R Code

```

1
2 # =====
3 #   League Comparison + Liverpool (EPL) Team Analysis Script
4 #   Source: football-data.co.uk
5 # =====
6
7 library(readr)
8 library(dplyr)
9 library(ggplot2)
10 library(lubridate)
11 library(tidyr)
12 library(stringr)
13 library(plotly)
14
15 # --- Helper to load and prep each league's data ---
16 load_league <- function(filename, league_code, league_name) {
17   df <- read_csv(filename, locale = locale(encoding = "Latin1"))
18   df <- df %>%
19     mutate(
20       League = league_name,
21       Date = dmy(Date),

```

```

22     Season = ifelse(month(Date) >= 8, paste0(year(Date), "/",
23         str_sub(year(Date) + 1, 3, 4)), paste0(year(Date) - 1, "
24         /", str_sub(year(Date), 3, 4)))
25 )
26 return(df)
27 }
28
29 # --- Load all five leagues ---
30 epl <- load_league("EPL_Data.csv", "EPL", "Premier League")
31
32 # =====
33 # LIVERPOOL TEAM ANALYSIS (EPL)
34 # =====
35 liverpool <- epl %>%
36   filter(HomeTeam == "Liverpool" | AwayTeam == "Liverpool") %>%
37   mutate(
38     LiverpoolIsHome = HomeTeam == "Liverpool",
39     Opponent = if_else(LiverpoolIsHome, AwayTeam, HomeTeam),
40     GoalsFor = if_else(LiverpoolIsHome, FTHG, FTAG),
41     GoalsAgainst = if_else(LiverpoolIsHome, FTAG, FTHG),
42     Result = case_when(
43       (LiverpoolIsHome & FTR == "H") | (!LiverpoolIsHome & FTR ==
44         "A") ~ "Win",
45       FTR == "D" ~ "Draw",
46       TRUE ~ "Loss"
47     ),
48     Venue = if_else(LiverpoolIsHome, "Home", "Away")
49   )
50
51 # Overall Liverpool stats
52 liverpool_overall <- liverpool %>%
53   group_by(Result) %>%
54   summarize(Matches = n())
55 write_csv(liverpool_overall, "liverpool_overall_results.csv")
56
57 # Liverpool season by season
58 liverpool_season <- liverpool %>%
59   group_by(Season) %>%
60   summarize(
61     Matches = n(),
62     Wins = sum(Result == "Win"),
63     Draws = sum(Result == "Draw"),

```

```

61     Losses = sum(Result == "Loss"),
62     GoalsFor = sum(GoalsFor),
63     GoalsAgainst = sum(GoalsAgainst),
64     Avg_GoalsFor = round(mean(GoalsFor),2),
65     Avg_GoalsAgainst = round(mean(GoalsAgainst),2),
66     Points = Wins * 3 + Draws * 1
67 ) %>%
68 arrange(Season)
69 write_csv(liverpool_season, "liverpool_season_summary.csv")
70
71 # Home vs Away Liverpool performance
72 liverpool_venue <- liverpool %>%
73   group_by(Venue, Season) %>%
74   summarize(
75     Matches = n(),
76     Wins = sum(Result == "Win"),
77     Draws = sum(Result == "Draw"),
78     Losses = sum(Result == "Loss"),
79     Avg_GoalsFor = round(mean(GoalsFor),2),
80     Avg_GoalsAgainst = round(mean(GoalsAgainst),2)
81   ) %>%
82   arrange(Season, Venue)
83 write_csv(liverpool_venue, "liverpool_home_away_summary.csv")
84
85 # Plot: Liverpool's Points per Season
86 p1 <- ggplot(liverpool_season, aes(x = Season, y = Points)) +
87   geom_line(group = 1, color = "red", linewidth = 1.2) +
88   geom_point(color = "black", size = 2) +
89   theme_minimal() +
90   labs(title = "Liverpool Total Points Per Season (EPL)", x = "
91     Season", y = "Points") +
92   theme(axis.text.x = element_text(angle=45, hjust=1))
93 ggsave("liverpool_points_per_season.png", p1, width=8, height=4)
94
95 # Plot: Liverpool Goals For and Against Per Season
96 p2 <- ggplot(liverpool_season, aes(x = Season)) +
97   geom_point(aes(y = Avg_GoalsFor, color = "Goals For"), size =
98     1) +
99   geom_point(aes(y = Avg_GoalsAgainst, color = "Goals Against"),
100     size = 1) +
101   theme_minimal() +

```

```

99   labs(title = "Liverpool Goals For and Against (Per Match, Per
    Season)", x = "Season", y = "Goals Per Match") +
100  scale_color_manual("", values = c("Goals For" = "red", "Goals
    Against" = "black")) +
101  theme(axis.text.x = element_text(angle=45, hjust=1))
102  ggsave("liverpool_goals_for_against.png", p2, width=8, height=4)
103
104  # Interactive Plots
105  ggplotly(p1)
106  ggplotly(p2)
107
108
109  # --- End of Script ---
110  cat("Analysis complete. See CSVs and PNGs for summary tables and
    plots.")

```

Listing 12: Liverpool Analysis Code

## M Appendix Liverpool Prediction R Code

```

1  # EPL Points Prediction for Liverpool
2
3  library(readr)
4  library(dplyr)
5  library(ggplot2)
6  library(lubridate)
7  library(tidyr)
8  library(stringr)
9  library(plotly)
10 library(caret)
11 library(randomForest)
12
13 # --- Helper to load and prep each league's data ---
14 load_league <- function(filename, league_code, league_name) {
15   df <- read_csv(filename, locale = locale(encoding = "Latin1"))
16   df <- df %>%
17     mutate(
18       League = league_name,
19       Date = dmy(Date),
20       Season = ifelse(month(Date) >= 8, paste0(year(Date), "/",
21                                                    str_sub(year(Date)
22                                                         + 1, 3, 4)),

```

```

22         paste0(year(Date) - 1, "/", str_sub(year(
23             Date), 3, 4)))
24     )
25     return(df)
26 }
27 # --- Load all five leagues ---
28 epl <- load_league("EPL_Data.csv", "EPL", "Premier League")
29
30 # =====
31 # LIVERPOOL TEAM ANALYSIS (EPL)
32 # =====
33 liverpool <- epl %>%
34     filter(HomeTeam == "Liverpool" | AwayTeam == "Liverpool") %>%
35     mutate(
36         LiverpoolIsHome = HomeTeam == "Liverpool",
37         Opponent = if_else(LiverpoolIsHome, AwayTeam, HomeTeam),
38         GoalsFor = if_else(LiverpoolIsHome, FTHG, FTAG),
39         GoalsAgainst = if_else(LiverpoolIsHome, FTAG, FTHG),
40         Result = case_when(
41             (LiverpoolIsHome & FTR == "H") | (!LiverpoolIsHome & FTR ==
42                 "A") ~ "Win",
43             FTR == "D" ~ "Draw",
44             TRUE ~ "Loss"
45         ),
46         Venue = if_else(LiverpoolIsHome, "Home", "Away")
47     )
48 # Overall Liverpool stats
49 liverpool_overall <- liverpool %>%
50     group_by(Result) %>%
51     summarize(Matches = n())
52 write_csv(liverpool_overall, "liverpool_overall_results.csv")
53
54 # Liverpool season by season
55 liverpool_season <- liverpool %>%
56     group_by(Season) %>%
57     summarize(
58         Matches = n(),
59         Wins = sum(Result == "Win"),
60         Draws = sum(Result == "Draw"),
61         Losses = sum(Result == "Loss"),

```

```

62     GoalsFor = sum(GoalsFor),
63     GoalsAgainst = sum(GoalsAgainst),
64     Avg_GoalsFor = round(mean(GoalsFor),2),
65     Avg_GoalsAgainst = round(mean(GoalsAgainst),2),
66     Points = Wins * 3 + Draws * 1
67 ) %>%
68   arrange(Season)
69 write_csv(liverpool_season, "liverpool_season_summary.csv")
70
71 # for reproducibility
72 set.seed(7406)
73 n <- nrow(liverpool_season)
74 train_idx <- sample(seq_len(n), size = 0.8 * n)
75
76 train_data <- liverpool_season[train_idx, ]
77 test_data <- liverpool_season[-train_idx, ]
78
79 model_rf <- train(Points ~ Wins + Draws + Losses + Avg_GoalsFor +
80                   Avg_GoalsAgainst,
81                   data = train_data, method = "rf")
82 model_rf
83
84 model_lm <- train(
85   Points ~ Wins + Draws + Losses + Avg_GoalsFor + Avg_
86   GoalsAgainst,
87   data = train_data,
88   method = "lm"
89 )
90 model_lm
91
92 # Linear model predictions
93 test_data$Predicted_Points <- predict(model_lm, newdata = test_
94   data)
95
96 # Random forest predictions
97 test_data$Predicted_Points_RF <- predict(model_rf, newdata = test
98   _data)
99
100 # Root Mean Squared Error for linear model
101 rmse_lm <- sqrt(mean((test_data$Points - test_data$Predicted_
102   Points)^2))

```

```

99
100 # For random forest
101 rmse_rf <- sqrt(mean((test_data$Points - test_data$Predicted_
    Points_RF)^2))
102
103 print(rmse_lm)
104 print(rmse_rf)
105
106 # --- End of Script ---
107 cat("Analysis complete. See CSVs and PNGs for summary tables and
    plots.")

```

Listing 13: Liverpool Prediction Code

## N Appendix: Graphs and Plots

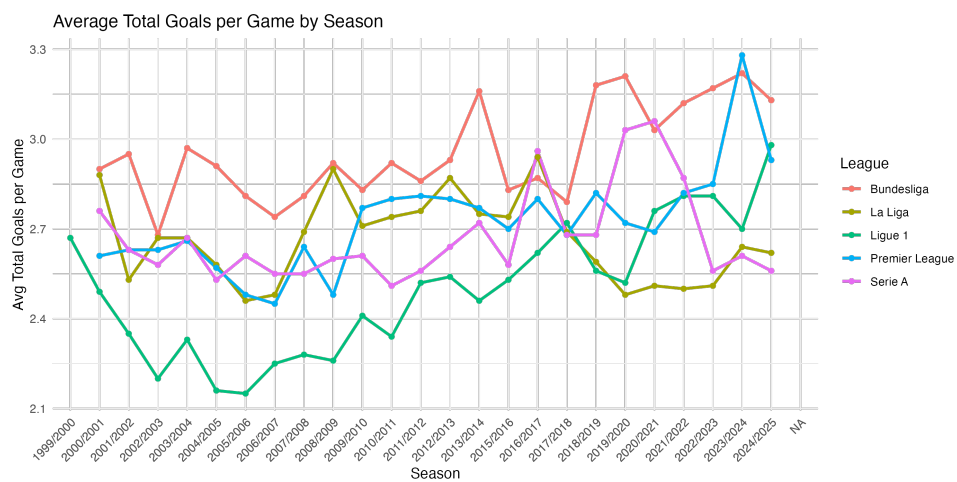


Figure 12: Goals by League and Season

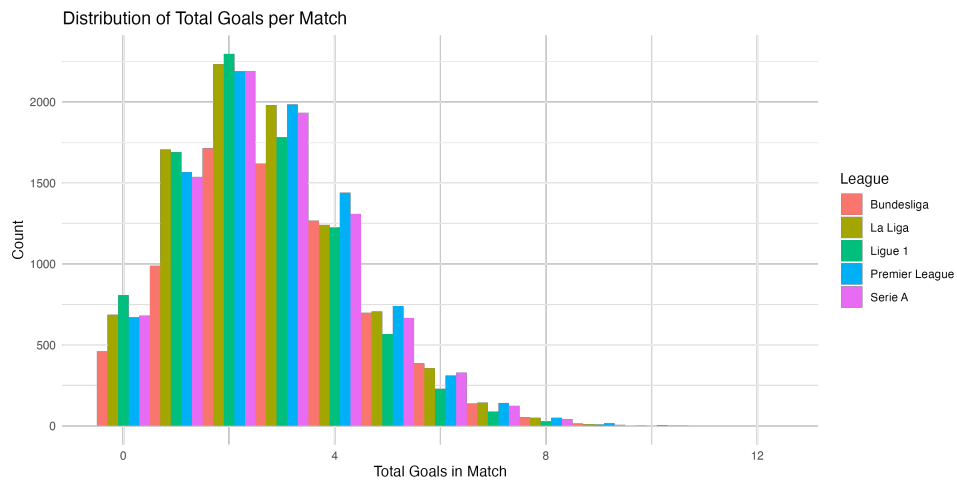


Figure 13: Distribution of Total Goals Per Match

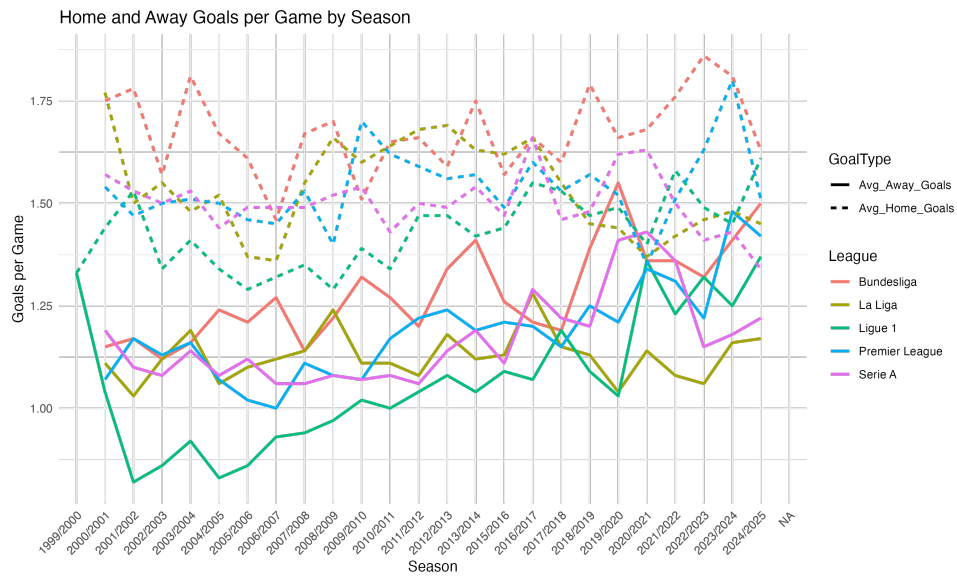


Figure 14: Home versus Away Goals Per Season by League



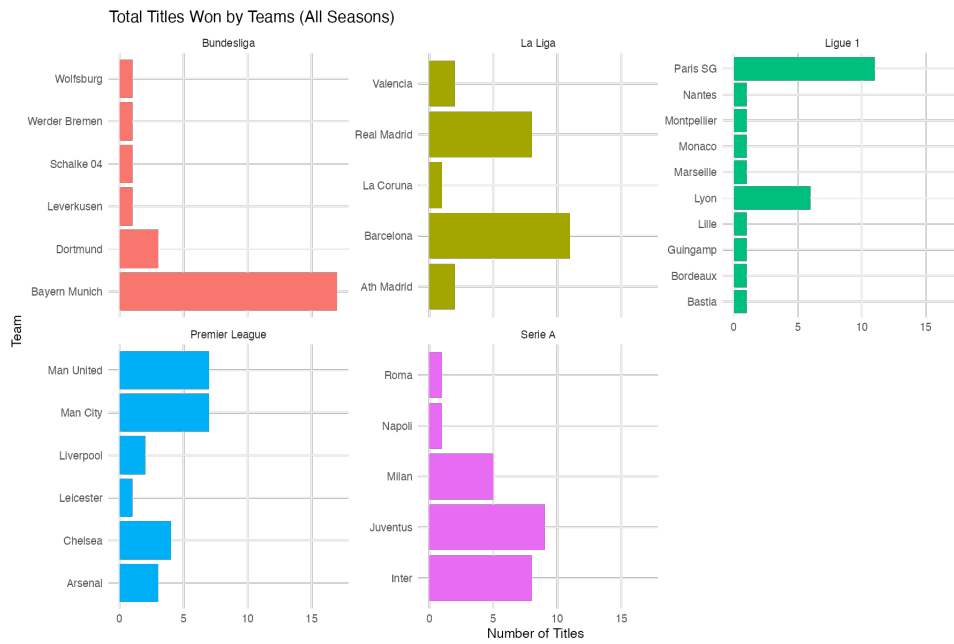


Figure 15: Titles by Team and League

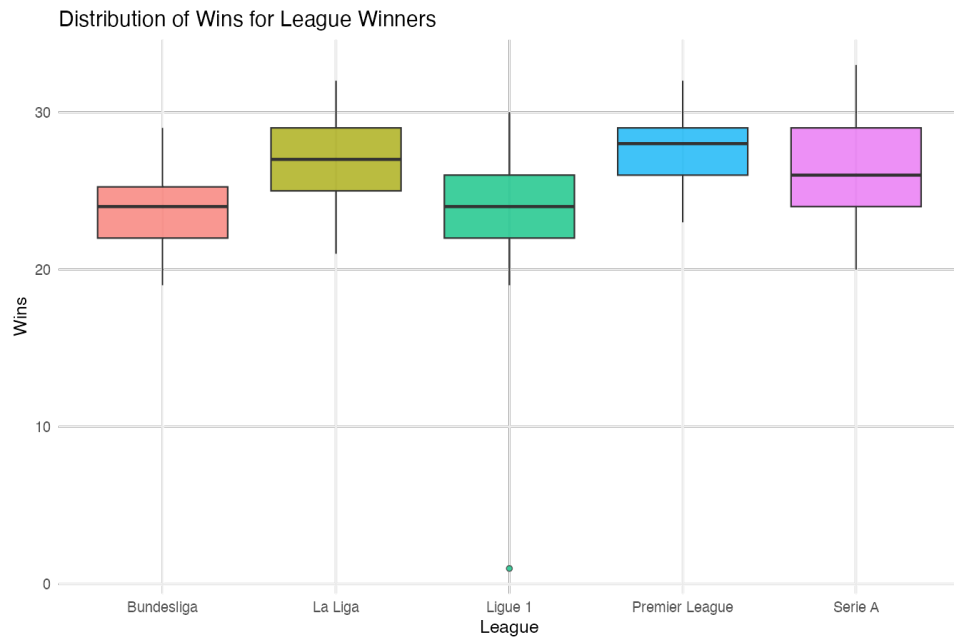


Figure 16: Distribution of Wins by Title Winners Per League

## Bibliography and Credits

- Agresti, A. (2013). *Categorical data analysis* (3rd ed.). Wiley. <https://onlinelibrary.wiley.com/doi/book/10.1002/0471249688>
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- Cawley, G. C., & Talbot, N. L. C. (2010). On over-fitting in model selection and subsequent selection bias in performance evaluation. *Journal of Machine Learning Research*, 11, 2079–2107. <https://jmlr.org/papers/v11/cawley10a.html>
- Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794. <https://doi.org/10.1145/2939672.2939785>
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20, 273–297. <https://link.springer.com/article/10.1007/BF00994018>
- Courneya, K. S., & Carron, A. V. (1992). The home advantage in sport competitions: A literature review. *Journal of Sport and Exercise Psychology*, 14(1), 13–27. <https://doi.org/10.1123/jsep.14.1.13>
- Cover, T. M., & Hart, P. E. (1967). Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1), 21–27. <https://doi.org/10.1109/TIT.1967.1053964>
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer. <https://web.stanford.edu/~hastie/ElemStatLearn/>
- Kuhn, M. (2008). Building predictive models in r using the `caret` package. *Journal of Statistical Software*, 28(5), 1–26. <https://www.jstatsoft.org/article/view/v028i05>
- Maher, M. J. (1982). Modelling association football scores. *Statistica Neerlandica*, 36(3), 109–118. <https://doi.org/10.1111/j.1467-9574.1982.tb00782.x>
- Murtagh, F., & Legendre, P. (2014). Ward’s hierarchical agglomerative clustering method: Which algorithms implement ward’s criterion? *Journal of Classification*, 31(3), 274–295. <https://doi.org/10.1007/s00357-014-9161-z>
- Nevill, A. M., & Holder, R. L. (1999). Home advantage in sport: An overview of studies on the advantage of playing at home. *Sports Medicine*, 28(4), 221–236. <https://doi.org/10.2165/00007256-199928040-00001>
- Pollard, R. (2006). Home advantage in soccer: A retrospective analysis. *Journal of Sports Sciences*, 24(3), 231–240. <https://doi.org/10.1080/02640410500141836>
- Scrucca, L., Fop, M., Murphy, T. B., & Raftery, A. E. (2016). Mclust 5: Clustering, classification and density estimation using gaussian finite mixture models. *The R Journal*, 8(1), 289–317. <https://journal.r-project.org/archive/2016/RJ-2016-021/index.html>
- Ward, J. H. (1963). Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, 58(301), 236–244. <https://doi.org/10.1080/01621459.1963.10500845>