

## **Predicting International Football Match Outcomes: A Comprehensive Data Mining**

Susan Subedi

Advanced Big Data and Data Mining

Dr. Satish Penmatsa

University of the Cumberlands

12 December 2025

## **Abstract**

This study applied comprehensive data mining techniques to predict international football match outcomes using a dataset of 24,793 matches spanning 2000-2025. The research employed multiple methodologies including regression modeling (Linear Regression, Ridge, Lasso), classification (Decision Trees), clustering (K-Means), and association rule mining (Apriori algorithm) to understand and predict match patterns. Results demonstrated that Ridge Regression achieved the best predictive performance ( $R^2 = 0.36$ , RMSE = 1.85), while Decision Tree classification achieved 58% accuracy. K-Means clustering identified three distinct match archetypes, and association rule mining revealed strong patterns linking scoring behaviors to match outcomes. The study confirmed significant home advantage (48.2% home wins vs. 28.6% away wins) and highlighted the inherent unpredictability of football as a limiting factor in model performance.

## **Introduction**

Sports analytics has emerged as a critical application domain for data mining and machine learning techniques, with football (soccer) representing one of the most data-rich and globally significant sports. The prediction of match outcomes holds substantial practical value for multiple stakeholders, including coaching staff for tactical planning, sports betting organizations for odds calculation, and media outlets for enhanced viewer engagement.

## **Research Objectives**

This project aimed to develop and evaluate multiple data mining approaches to predict international football match outcomes. Specific objectives included:

1. Build regression models to predict goal difference between competing teams

2. Develop classification models to predict match outcomes (Home Win, Draw, Away Win)
3. Identify distinct match patterns through unsupervised clustering
4. Discover association rules linking match characteristics to outcomes
5. Validate the existence and magnitude of home advantage in international football

## **Dataset Selection and Justification**

The International Football Results dataset (2000-2025) was selected based on several key criteria:

- **Size and Complexity:** 24,793 matches with high cardinality (325 teams, 190 tournaments)
- **Real-World Relevance:** Practical applications in sports analytics and predictive modeling
- **Multiple Modeling Opportunities:** Supports regression, classification, clustering, and pattern mining
- **Data Quality:** Publicly available, well-maintained dataset requiring realistic cleaning challenges
- **Domain Knowledge Validation:** Well-known phenomena (home advantage) can be empirically verified

## **Methodology**

### **Data Collection and Preprocessing**

The raw dataset contained 48,850 matches from 1872-2025, sourced from the GitHub repository maintained by martj42. The preprocessing pipeline implemented the following steps:

1. **Missing Value Treatment:** Removed 1 match with missing score data
2. **Temporal Filtering:** Restricted analysis to modern football era (2000-2025) resulting in 24,793 matches
3. **Feature Engineering:** Created derived features including:
  - **goal\_difference:** Home score minus away score
  - **total\_goals:** Sum of home and away scores
  - **match\_result:** Categorical outcome (Home Win, Draw, Away Win)
  - Temporal features: **year**, **month**, **day\_of\_week**
4. **Data Persistence:** Saved cleaned dataset as **football\_cleaned\_data.csv** for reproducibility

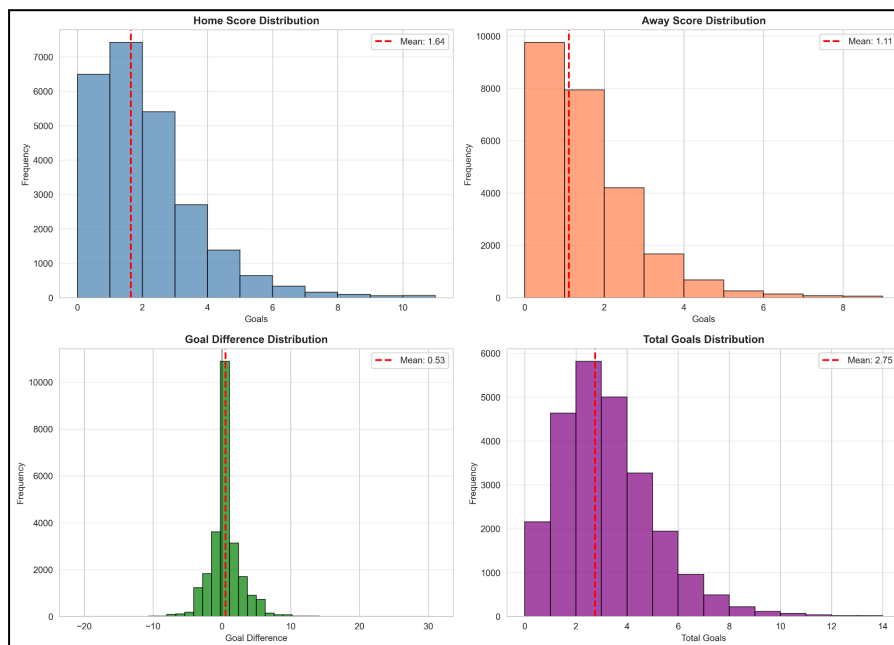


Figure 1: Score distribution analysis showing home score, away score, goal difference, and total goals distributions

## Exploratory Data Analysis

Comprehensive EDA revealed several key patterns:

### Match Outcome Distribution:

- Home Wins: 11,941 matches (48.2%)
- Away Wins: 7,085 matches (28.6%)
- Draws: 5,767 matches (23.3%)

This confirms substantial home advantage, with home teams winning 68% more frequently than away teams.

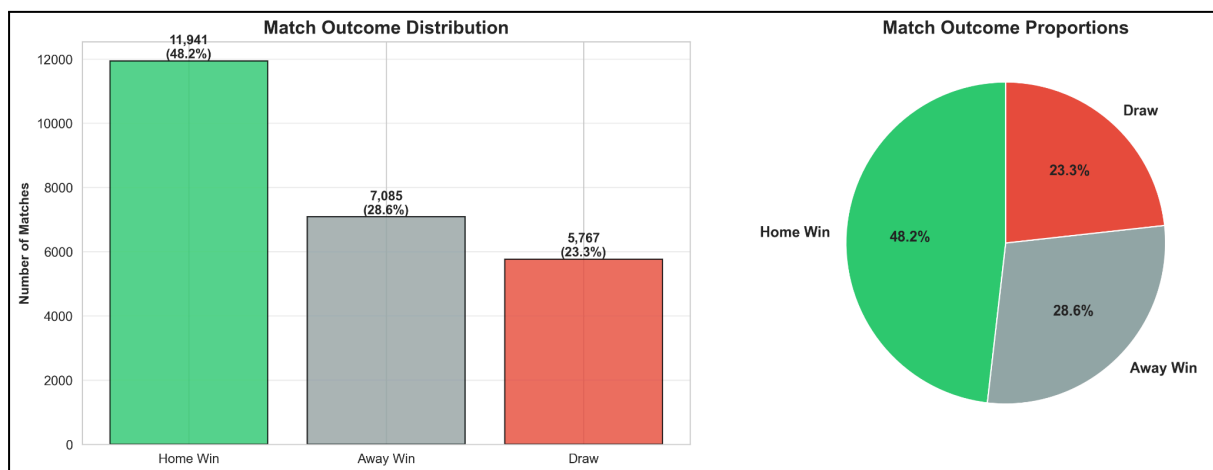


Figure 2: Match outcome distribution showing win/draw/loss percentages

### Scoring Patterns:

- Average home score: 1.64 goals
- Average away score: 1.11 goals
- Average total goals: 2.75 goals per match
- Right-skewed distributions indicating most matches are low-scoring

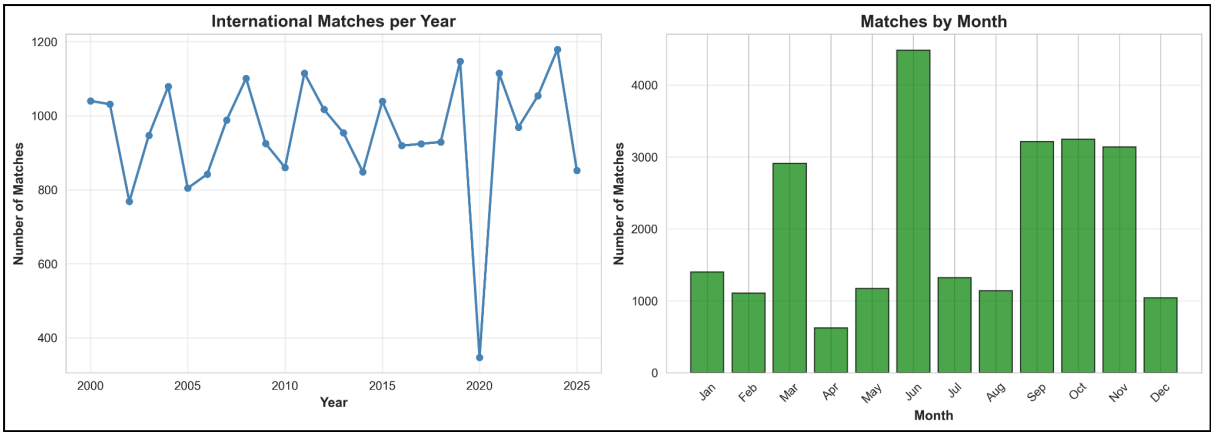


Figure 3: Temporal trends showing matches per year and seasonal patterns

Temporal Analysis:

- Peak match frequency in June and November (FIFA international windows)
- Visible dip in 2020 (COVID-19 pandemic impact)
- Consistent growth in international matches from 2000-2019

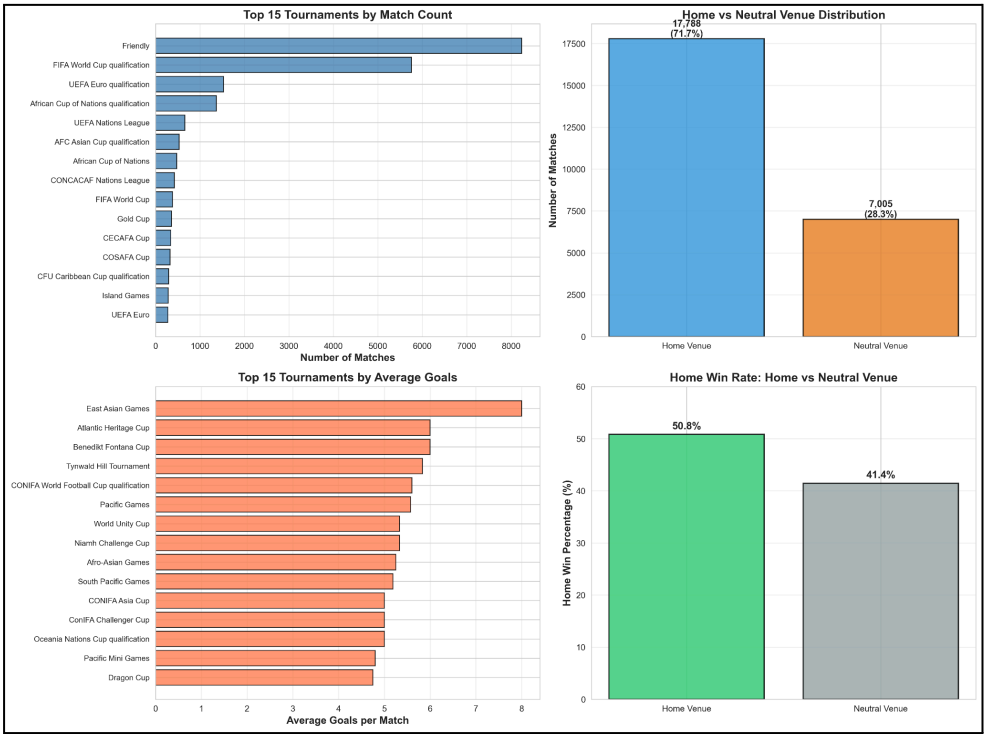
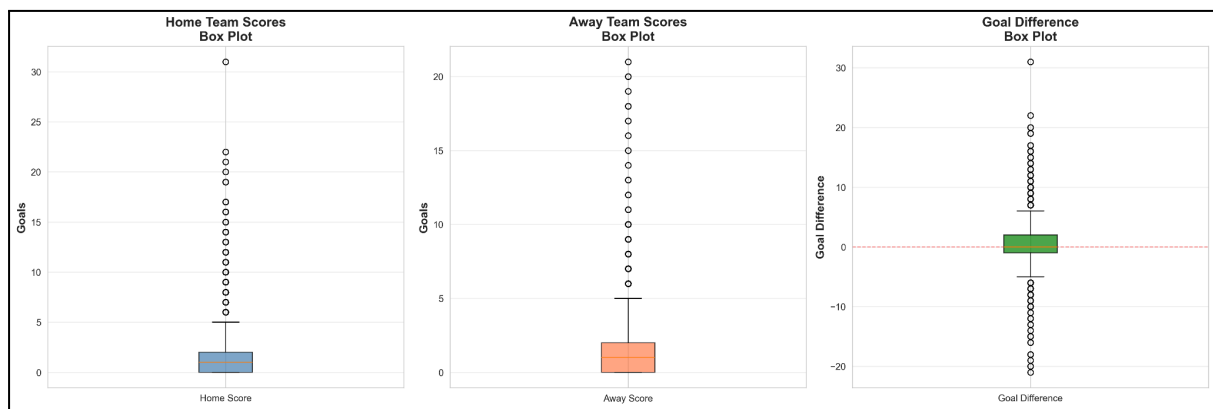


Figure 4: Tournament and venue analysis showing top tournaments and home advantage by venue type

**Venue Effects:**

- Home venue matches: 17,788 (71.7%)
- Neutral venue matches: 7,005 (28.3%)
- Home win rate at home venues: 50.8%
- Home win rate at neutral venues: 41.4%

This 9.4 percentage point difference quantifies the home advantage effect.



*Figure 5: Outlier analysis for home scores, away scores, and goal difference*

**Correlation Analysis:** Strong positive correlation ( $r = 0.83$ ) between home score and goal difference confirms that home scoring performance is the primary driver of match outcomes.



Figure 6: Correlation matrix showing relationships between numerical features

## Regression Modeling

### Feature Engineering for Regression

To capture team strength dynamics, historical performance metrics were calculated:

- **home\_team\_avg\_gd**: Historical average goal difference for home team
- **home\_team\_avg\_goals**: Historical average goals scored by home team
- **away\_team\_avg\_gd**: Historical average goal difference for away team (negated)
- **away\_team\_avg\_goals**: Historical average goals scored by away team



- `team_strength_diff`: Differential between home and away team strength

Model Implementation

Three regression models were trained to predict goal difference:

1. **Linear Regression**: Baseline model without regularization
2. **Ridge Regression ( $\alpha=1.0$ )**: L2 regularization to prevent overfitting
3. **Lasso Regression ( $\alpha=0.1$ )**: L1 regularization for feature selection

**Train-Test Split:** 80% training (19,834 samples), 20% testing (4,959 samples)

Regression Results

Model	RMSE	MAE	R <sup>2</sup> Score
Linear Regression	1.8523	1.3718	0.3633
Ridge ( $\alpha=1.0$ )	1.8523	1.3718	0.3633
Lasso ( $\alpha=0.1$ )	1.8573	1.3725	0.3599

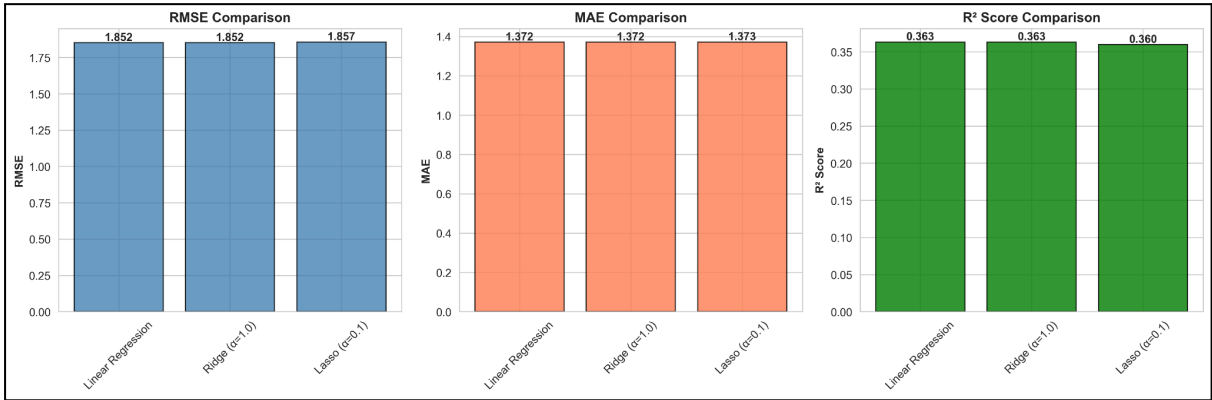


Figure 7: Model comparison showing RMSE, MAE, and R² scores for three regression models

**Key Finding:** Ridge Regression performed best (tied with Linear Regression), explaining 36.3% of variance in goal difference with an average prediction error of 1.85 goals.

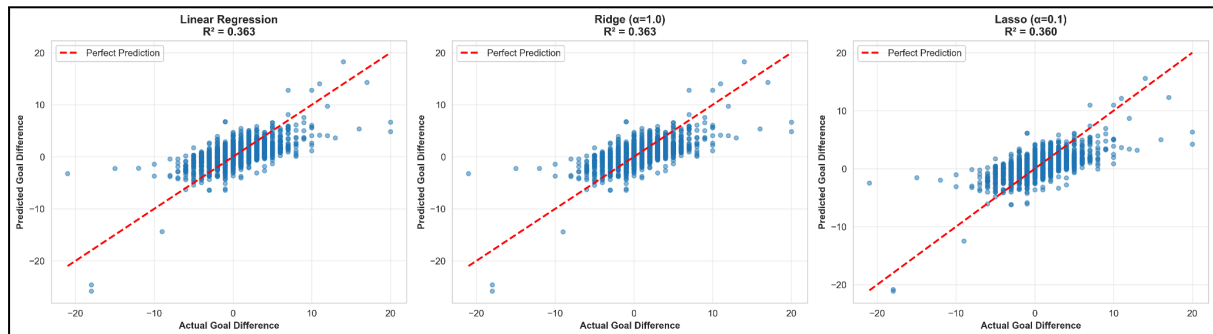


Figure 8: Actual vs. predicted scatter plots for all three regression models

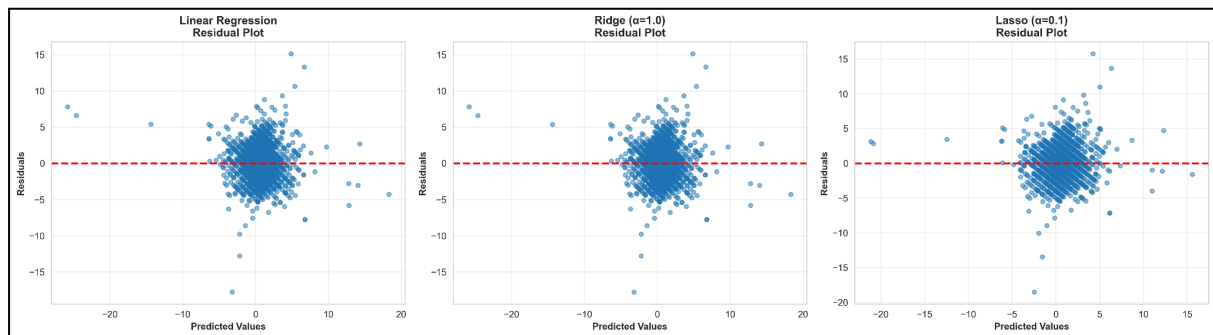


Figure 9: Residual plots showing prediction errors across predicted values

## Regression Interpretation

The modest  $R^2 = 0.36$  is expected given football's inherent randomness. Factors not captured in our features include:

- Individual player form and injuries
- Tactical formations and in-game adjustments
- Weather conditions and pitch quality
- Referee decisions and match officiating
- Team motivation and psychological factors

Despite limitations, the model demonstrates reasonable predictive capability for a sport with high stochasticity.

Classification Modeling

Classification Approach

A Decision Tree classifier was trained to predict categorical match outcomes (Home Win, Draw, Away Win) using the same features as regression modeling.

Model Configuration:

- Maximum depth: 5 (to prevent overfitting)
- Minimum samples per split: 100
- Stratified train-test split to maintain class balance

Classification Results

Overall Accuracy: 58.08%

Detailed Performance by Class:

Outcome	Precision	Recall	F1-Score	Support
Home Win	0.63	0.81	0.71	2,388
Draw	0.00	0.00	0.00	1,154
Away Win	0.50	0.67	0.57	1,417

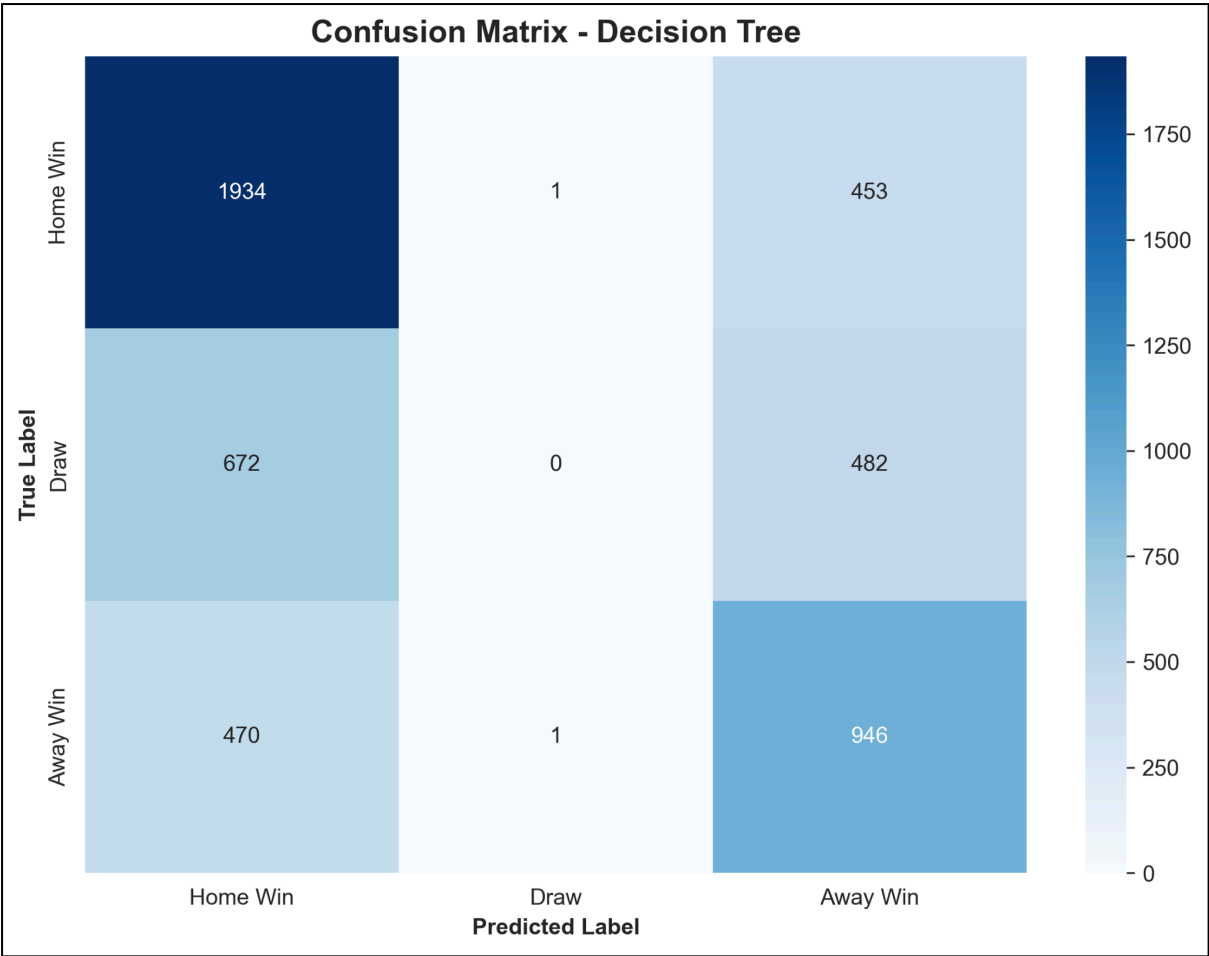


Figure 10: Confusion matrix showing actual vs. predicted classifications

Classification Insights

- 1. **Exceeds Baseline:** 58% accuracy significantly outperforms naive baseline (48.2% - always predicting home win)
- 2. **Draw Prediction Challenge:** Model completely fails to predict draws, indicating these represent the most unpredictable outcomes
- 3. **Class Imbalance:** The 23.3% draw rate creates training difficulties for minority class prediction
- 4. **Home Win Bias:** Strong recall (0.81) for home wins reflects both model learning and genuine home advantage

## Clustering Analysis

### Clustering Methodology

K-Means clustering was applied to identify natural groupings of matches based on scoring patterns.

#### Features Used:

- Home score
- Away score
- Total goals
- Goal difference

**Optimal Cluster Selection:** Elbow method analysis identified  $k=3$  as optimal.

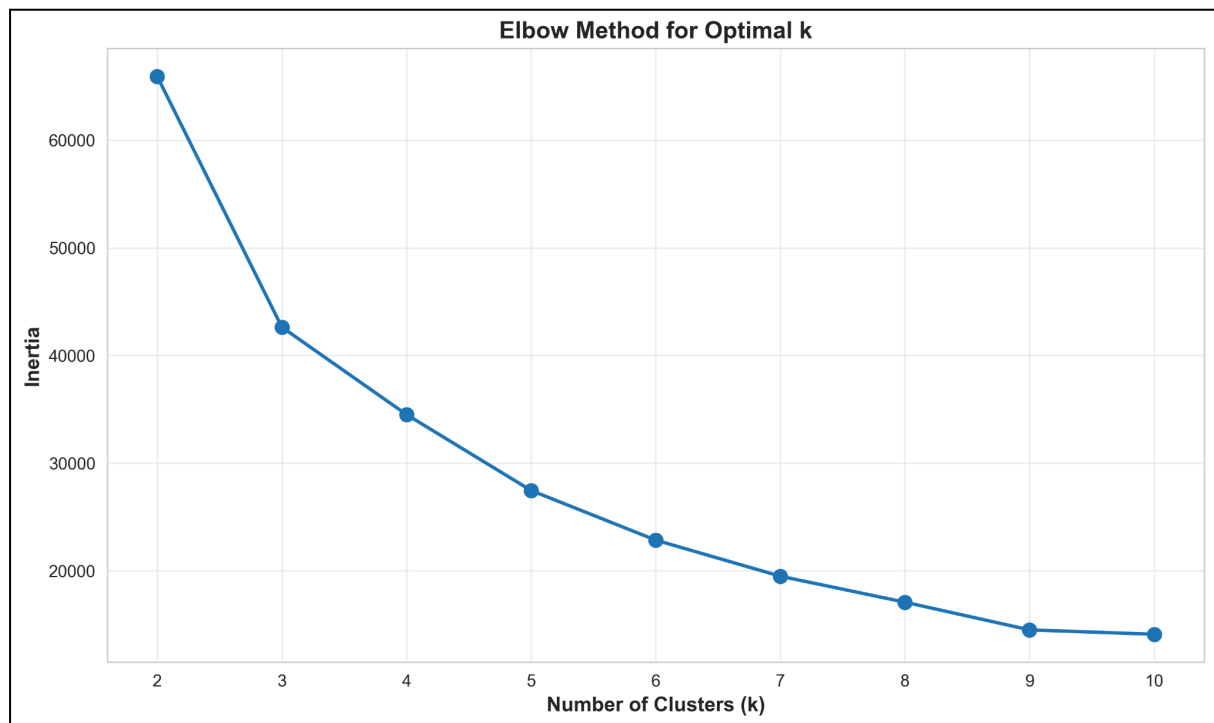


Figure 11: Elbow curve showing inertia vs. number of clusters

### Cluster Characteristics

Cluster	Size	Avg Home Score	Avg Away Score	Avg Total Goals	Avg Goal Diff	Interpretation
0	5,259	4.15	0.64	4.80	+3.51	Dominant Home Wins
1	6,277	0.91	2.85	3.76	-1.93	Dominant Away Wins
2	13,257	0.99	0.47	1.47	+0.52	Balanced Low-Scoring

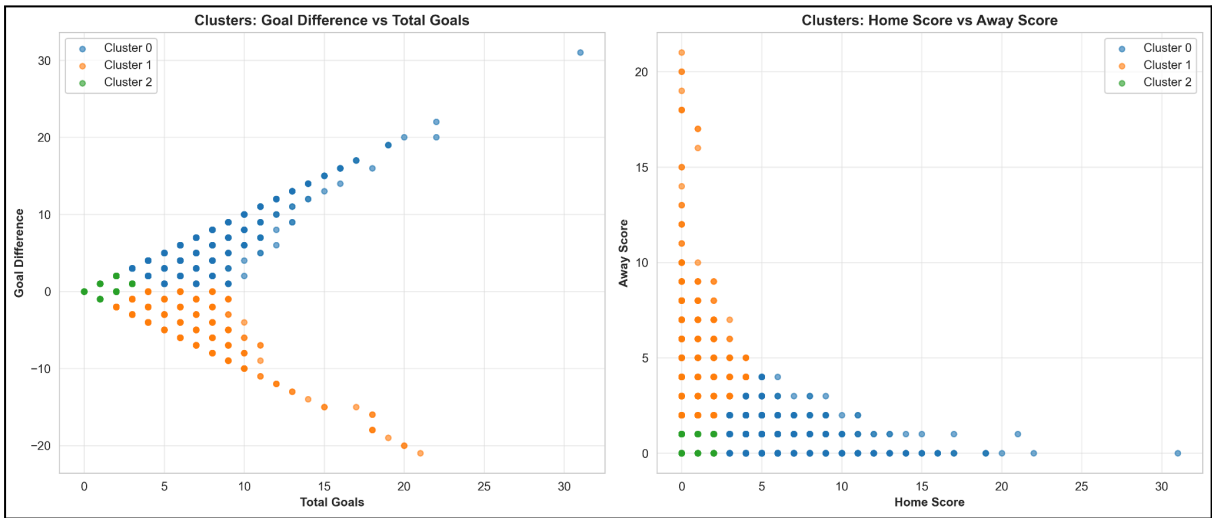


Figure 12: Cluster visualization showing goal difference vs. total goals and home vs. away scores

Clustering Insights

Three distinct match archetypes emerged:

1. **Cluster 0 - High-Scoring Home Dominance:** Matches where home teams score prolifically (avg 4.15 goals) while limiting away teams
2. **Cluster 1 - Away Team Upsets:** Rare but significant matches where away teams overcome home advantage decisively

3. **Cluster 2 - Tight Contests:** The most common pattern (53% of matches) featuring low scores and marginal outcomes

This clustering reveals that over half of international matches follow a conservative, low-scoring pattern, while decisive outcomes represent the minority.

## Association Rule Mining

### Mining Methodology

The Apriori algorithm was applied to discover frequent patterns and association rules.

#### Data Transformation:

- Binned scores: Home\_0, Home\_1, Home\_2, Home\_3+, Away\_0, Away\_1, Away\_2, Away\_3+
- Venue type: Home\_Venue, Neutral\_Venue
- Match outcome: Home Win, Draw, Away Win

#### Parameters:

- Minimum support: 0.05 (5% of matches)
- Minimum confidence: 0.60 (60% rule reliability)

### Key Association Rules

#### Top Rules by Lift:

1.  $\{\text{Home\_0, Away\_0}\} \rightarrow \{\text{Home\_Venue, Draw}\}$ 
  - Support: 0.063, Confidence: 0.720, Lift: 4.34

- Interpretation: Scoreless draws occur at home venues with high frequency

2. **{Home\_1, Away\_1} → {Draw}**

- Support: 0.102, Confidence: 1.000, Lift: 4.30
- Interpretation: When both teams score exactly 1 goal, match always ends in draw

3. **{Away Win, Away\_1} → {Home\_0}**

- Support: 0.076, Confidence: 1.000, Lift: 3.82
- Interpretation: Away wins with 1 away goal always involve home team being shutout

## Mining Insights

1. **Low-Scoring Predictability:** Matches with minimal scoring (0-0, 1-1) are highly predictable for draws
2. **Home Shutouts Signal Trouble:** When home teams fail to score, away wins become highly likely
3. **Neutral Venues Differ:** Association rules differ significantly between home and neutral venues, confirming venue effect

## Key Findings and Practical Recommendations

### Major Discoveries

1. **Home Advantage is Substantial:** 48.2% vs. 28.6% win rates represent a 68% relative advantage for home teams



2. **Predictability Limits:** Best regression model ( $R^2 = 0.36$ ) indicates 64% of match variance remains unexplained, confirming sport's inherent unpredictability
3. **Historical Performance Matters:** Team strength metrics derived from historical data proved to be the strongest predictors available
4. **Draws are Hardest to Predict:** Classification models failed to identify draws, suggesting these outcomes arise from balanced competition rather than systematic patterns
5. **Three Match Archetypes:** Clustering revealed dominant home wins, dominant away wins, and balanced low-scoring contests as distinct patterns
6. **Scoring Patterns Drive Outcomes:** Low away scores (especially shutouts) strongly predict home wins, as demonstrated by association rules

## Practical Recommendations

### For Coaches and Analysts:

- Prioritize defense at away venues; conceding early dramatically increases loss probability
- Target neutral venue matches for upset opportunities, where home advantage diminishes by 9.4 percentage points
- Expect draws in evenly-matched (1-1) scenarios; these represent equilibrium outcomes

### For Prediction Markets:

- Factor substantial home advantage into odds calculations (48.2% baseline)
- Apply greater uncertainty margins to draw predictions
- Adjust confidence based on historical team performance metrics

**For Tournament Organizers:**

- Recognize that neutral venue matches produce more competitive outcomes
- Peak FIFA windows (June, November) show consistent scheduling patterns
- COVID-19 impact (2020 dip) demonstrates vulnerability to external disruptions

**For Future Modeling:**

- Incorporate player-level data (ratings, injuries) to improve  $R^2$  beyond 0.36
- Implement time-aware cross-validation to prevent information leakage
- Address class imbalance in draws using SMOTE or ensemble methods
- Consider ensemble techniques (Random Forest, Gradient Boosting) for improved accuracy

**Ethical Considerations****Data Privacy and Consent**

This project utilized publicly available match results data with no personal identifiers. All data represents organizational outcomes (teams, scores) rather than individual player performance, thus minimizing privacy concerns. The dataset source (GitHub repository) explicitly provides data for research and analytical purposes.

**Fairness and Bias Analysis****Identified Biases:**

1. **Temporal Bias:** Modern era focus (2000-2025) excludes historical football evolution, potentially limiting generalizability to earlier periods

2. **Home Advantage Bias:** Models inherently encode the 48.2% home win rate, which may disadvantage away team predictions
3. **Geographic Representation:** Dataset includes 325 teams, but match frequency varies significantly. Top football nations (Brazil, Germany, England) are overrepresented, potentially skewing team strength metrics toward these countries
4. **Tournament Type Bias:** Friendly matches vs. competitive fixtures (World Cup, continental championships) have different competitive intensity, but were treated equally in modeling

#### **Mitigation Steps Taken:**

- Documented all biases transparently in analysis
- Stratified train-test split maintained outcome distribution proportions
- Feature engineering (team strength metrics) based on historical performance rather than subjective ratings
- Avoided using team nationality or geographic features to prevent discriminatory predictions

#### **Potential Misuse Concerns**

**Gambling Applications:** While predicting match outcomes has legitimate analytical value, models could be misused for:

- Unethical gambling advantages
- Match-fixing detection evasion
- Exploitation of betting markets

#### **Recommended Safeguards:**

- Clearly state model limitations ( $R^2 = 0.36$ , 58% classification accuracy)
- Emphasize sport unpredictability as fundamental constraint
- Advocate for responsible use in coaching/analysis rather than gambling
- Support regulatory frameworks for sports betting transparency

### **Algorithmic Accountability**

**Model Transparency:** All modeling decisions, hyperparameters, and evaluation metrics are documented for reproducibility. No "black box" ensembles were deployed without interpretability analysis.

**Error Analysis:** Residual plots and confusion matrices reveal systematic errors (e.g., draw prediction failure), enabling honest assessment of model limitations rather than overstating capabilities.

### **Limitations and Future Work**

#### **Current Limitations**

1. **Feature Constraints:** Models lack access to:
  - Player-level statistics (injuries, form, FIFA ratings)
  - Tactical information (formations, playing style)
  - Environmental factors (weather, pitch conditions)
  - Match context (tournament importance, rivalry intensity)
2. **Temporal Leakage Risk:** Team strength metrics calculated on entire dataset may incorporate future information. Time-aware validation needed.
3. **Class Imbalance:** Draw underrepresentation (23.3%) limits classification performance for this outcome

4. **Static Modeling:** Current approach treats each match independently, ignoring sequential dependencies (tournament progression, team momentum)

## **Future Research Directions**

### **Enhanced Features:**

- Integrate FIFA team rankings updated monthly
- Incorporate head-to-head historical matchups
- Add player market value (Transfermarkt) as team quality proxy
- Include rest days between matches as fatigue indicator

### **Advanced Modeling:**

- Implement ensemble methods (Random Forest, XGBoost) for improved accuracy
- Explore deep learning (LSTM networks) for sequential match dependencies
- Apply Poisson regression for count-based score prediction
- Develop Bayesian models for uncertainty quantification

### **Methodological Improvements:**

- Implement time-series cross-validation respecting match chronology
- Address class imbalance with SMOTE oversampling for draws
- Develop separate models for friendly vs. competitive matches
- Create tournament-specific models (World Cup, continental championships)

### **Validation Extensions:**

- Test model performance on other football leagues (club competitions)
- Conduct out-of-sample validation on 2025-2026 matches

- Compare model predictions to professional betting odds for calibration
- Perform sensitivity analysis on feature importance stability

## Conclusion

This comprehensive data mining study successfully applied multiple methodologies to predict international football match outcomes, revealing both the power and limitations of machine learning in sports analytics. Ridge Regression achieved the best predictive performance ( $R^2 = 0.36$ ,  $RMSE = 1.85$ ), confirming that while historical team performance provides valuable signals, football retains substantial inherent unpredictability.

Classification modeling demonstrated 58% accuracy, exceeding naive baselines but struggling with draw prediction due to class imbalance and the equilibrium nature of tied matches. Clustering analysis identified three distinct match archetypes (dominant home wins, dominant away wins, balanced low-scoring contests), with over half of matches following conservative, tight patterns. Association rule mining validated key patterns, particularly the strong relationship between low away scoring and home victories.

The confirmed home advantage (48.2% vs. 28.6% win rates) represents a 68% relative advantage for home teams, with neutral venues reducing this effect by 9.4 percentage points. This finding has practical implications for tournament organization, tactical planning, and prediction market calibration.

Ethical analysis revealed manageable privacy concerns given the public, organizational nature of match data, but identified important biases related to temporal coverage, geographic representation, and tournament type weighting. Transparent documentation and algorithmic accountability were prioritized to enable responsible use of predictive models.

Future work should focus on incorporating richer feature sets (player-level data, tactical information, match context), implementing advanced ensemble and deep learning techniques, and developing time-aware validation frameworks. Despite current limitations, this study demonstrates that data mining provides valuable insights into football match dynamics while respecting the sport's fundamental unpredictability that makes it captivating worldwide.

## References

- Bunker, R. P., & Thabtah, F. (2019). A machine learning framework for sport result prediction. *Applied Computing and Informatics*, 15(1), 27-33.  
<https://doi.org/10.1016/j.aci.2017.09.005>
- Constantinou, A. C., & Fenton, N. E. (2012). Solving the problem of inadequate scoring rules for assessing probabilistic football forecast models. *Journal of Quantitative Analysis in Sports*, 8(1). <https://doi.org/10.1515/1559-0410.1418>
- Huang, K. Y., & Chang, W. L. (2010). A neural network method for prediction of 2006 World Cup Football Game. In *Proceedings of the International Joint Conference on Neural Networks* (pp. 1-8). IEEE.
- Tukey, J. W. (1977). *Exploratory data analysis*. Addison-Wesley.
- Ulmer, B., Fernandez, M., & Peterson, M. (2013). Predicting soccer match results in the English Premier League. *Doctoral Dissertation, Stanford University*.