

Determining the Contribution of Performance Variables to Game Outcomes in Elite Male Soccer

Barış KARAKOÇ^{1*}  Erdem ŞANLI²  Alper AŞÇI¹ 

¹Faculty of Sport Sciences, Halic University, Istanbul, Türkiye; ²Vocational Schools, Fatih Sultan Mehmet Vakif University, Istanbul, Türkiye.

*Corresponding Author: bariskarakoc@halic.edu.tr

Abstract

The increased amount of soccer data challenges performance evaluation using classical methods. Therefore, recent machine learning approaches can help address complex datasets in sports. The present study aims to determine which performance variables most strongly contribute to game outcomes in elite male soccer by using machine learning models trained under different venue conditions. Technical, tactical, and physical variables obtained from 542 matches played over two consecutive seasons were used to predict results for three venue conditions (all, home, and away). Variance Inflation Factor analysis and BorutaPy were applied before extreme gradient boost (XG_{Boost}) modeling. Feature importance rankings and SHAP analysis were used to identify variables affecting model performance and outputs across different conditions. The models showed high accuracy on game outcome predictions, especially in the win and loss conditions (between 93.38-95.93%), while lower results in the draw (between 68.99-88.46%). The variables that most impacted the model's predictions were the Conversion rate, the Opponent's xG per goal, the Opponent's xG, and xG Conversion. The teams' performance predictions for game outcomes differ, and draws are difficult to predict in this study's competition. The technical variables contributed the most to the models and outputs. Coaches should consider the structure and needs of their competitions while evaluating the data they possess. Future research could develop models for different tournaments, especially using time-related variables, if applicable.

Keywords: Game Analysis, Extreme Gradient Boost, Machine Learning, SHAP Analysis

INTRODUCTION

In sports performance, whether a team or individual sport, coaches and players need as much information as possible to make necessary and appropriate decisions (Sampaio, 2013). It is also clear that players need effective feedback, grounded in informed and accurate measures, to improve performance. Thus, game analysis plays a crucial role in achieving success in sports. Studies conducted by researchers and coaches over the years have focused on several performance aspects, including technical, tactical, physical, and behavioral components (Carling

et al., 2007). In sports like soccer, from the researcher's perspective, data is collected to inform the game's needs. From the coach's strategic perspective, the needs are essential, as are the teams' (own and opponents') strengths and weaknesses, to prepare for future games. Regardless of the reason, both researchers and coaches require performance indicators to obtain information about gameplay. Performance indicators are the action variables defining some or all aspects of sports performance (Hughes & Bartlett, 2002). Studies in soccer focused on technical (Dellal et al., 2011; Rampinini et al., 2009), tactical (Perl et al., 2013; Taylor et al., 2005), and physical (Dellal et al., 2011; Mohr et al., 2003) parts of the game, similar to other sports. Researchers used different groups (age, league levels, competitions, etc.) based on their needs to obtain useful information from the data. A critical issue in soccer studies is associating the indicators with game results or rankings in the competitions to determine how successful teams play (Kubayi & Larkin, 2020; Moreira Praça et al., 2023; Yang et al., 2018). There have been many studies that have attempted to measure success in soccer. Meanwhile, different statistical methods were applied to learn which indicators discriminate winning/losing better (Castellano et al., 2012; Lago-Peñas et al., 2010, 2011) or which are highly correlated with better rankings (Rampinini et al., 2009). These methods helped understand the requirements of games across different competitions and build a template for defining how successful teams play.

The increase in the volume and variety of data in soccer has led to a similar increase in interest and attention from researchers outside the field of sports science. Their involvement in research on sports data analysis has employed a different approach to addressing scientific problems and statistical methods. In a study that investigated 576 games across three seasons, the researchers aimed to identify which variables strongly predict winning and losing in Belgian professional soccer (Geurkink et al., 2021). A total of 100 variables, including technical (shots, dribbles, duels, etc.) and physical (total distance, distance at different speed zones, number of accelerations and decelerations, etc.) performance indicators, were reduced to 13 variables by the Variance Inflation Factor method. Then, a tree-based machine learning technique was applied to select features and predict game outcomes with 90% model accuracy. In a similar study, 19 of 75 variables showed a strong effect on game outcomes, and the model's accuracy was 83.3% and 72.7% for winning and losing, respectively (Hassan et al., 2020). Similar to the studies above, researchers also selected the passing distribution as a criterion to predict match outcomes (win or lose) using different machine learning (ML) techniques (Cho et al., 2018). Despite the researcher's aim to predict match outcomes using various indicators, their tendency was to produce binary results. Nevertheless, it should be considered that drawn games are not uncommon in football.

Defining and comparing performances in professional male soccer leagues have been studied before. Still, the findings of those studies have shown limited insight into the performance and success of games, due to the use of a small number of performance indicators. Another issue noted in a previous study was the reliance on small samples and the univariate analyses conducted for the observed variables (Lago-Peñas et al., 2011). As mentioned, new statistical methods and models can help analyze and predict game outcomes from large amounts of soccer data (Bunker et al., 2024). The importance of collaboration between sports and computer/data scientists has been stated in the literature (Rein & Memmert, 2016). Therefore, this study can provide new insights for sports scientists interested in game analysis and for coaches through their decision-making process (Settembre et al., 2004). This study aims to determine which technical, tactical, and physical variables, as the match performance indicators, strongly contributed to game outcomes (win, draw, and loss) using three different machine learning models for each venue condition in the Turkish Super League.

METHODS

Research Model

The study is a cross-sectional study.

Data Sample

The data for the current research were collected from 542 games over two seasons (2021-2022 and 2022-2023) of the Turkish Super League, the highest level of male soccer in Türkiye. Although there was a total of 693 games played in two consecutive seasons, the games in which any players had red cards were excluded to avoid the imbalances of performance data (Badiella et al., 2023; Gomez et al., 2016). A total of 106 variables (61 technical, 42 tactical, and three physical) were included in the analysis. The variable descriptions are provided in Appendix 1. The data was obtained from the Instat Scout Platform, which collects technical, tactical, and physical data on individual and team performance using a video-based system. The reliability of their data has been established in previous studies (Kubayi, 2020; Kubayi & Larkin, 2020; Silva & Marcelino, 2023). The study received approval from the local university's Social Sciences and Humanities Research Ethics Committee (10.07.2024, 06).

Procedure

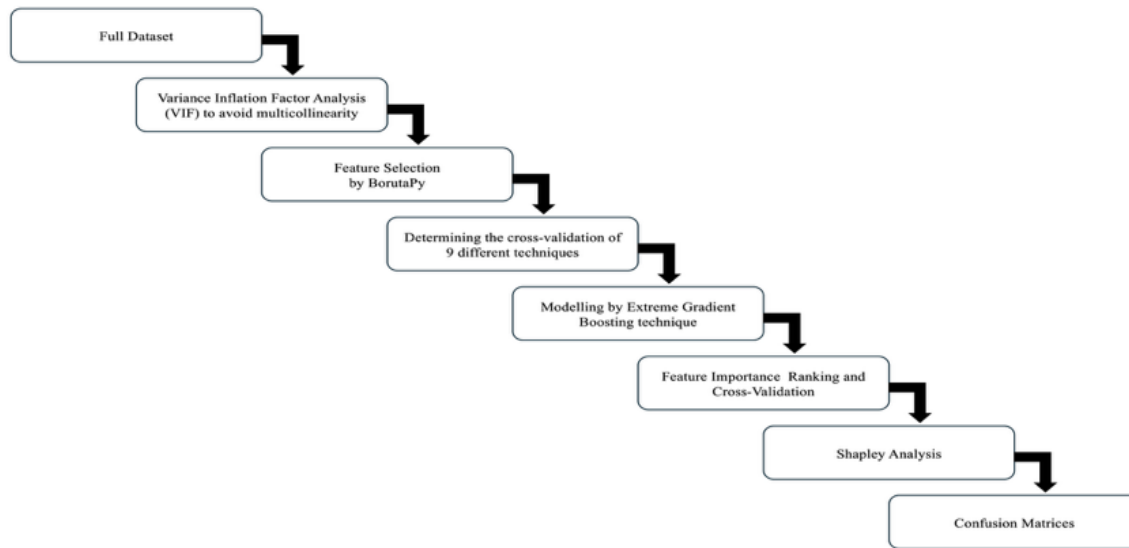
All data were collected from the Instat Scout Platform (<https://football.instatscout.com>) and downloaded into MS Excel. Columns indicating whether the team played at home or away, along

with the team's match result (win, draw, or loss), were subsequently added to the dataset. The dataset was then loaded into Visual Studio Code, a Python editor. The game-related performance variables were examined separately for three venue conditions: home (venue_home), away (venue_away), and all venues (venue_all), with home and away games evaluated together. Venue_home refers to the games where the team plays at its stadium, whereas venue_away refers to the games at the opponent team's stadium. Each venue's match outcomes (win, loss, and draw) were accepted as different conditions. In each match, the performance variables of the teams were evaluated based on the match result they achieved (e.g., if the match ended like Team A [W] - Team B [L]; Team A was evaluated under the Home_Win and All_Win conditions, while Team B was evaluated under the Away_Loss and All_Loss conditions). Totally 9 different conditions (All_Win, All_Loss, All_Draw, Home_Win, Home_Loss, Home_Draw, Away_Win, Away_Loss, Away_Draw) were observed for this study.

Data Analysis

Game outcomes were the dependent variables, and the other variables were assessed as independent. The analysis and study design are shown in Figure 1. To avoid multicollinearity among independent variables, a variance inflation factor (VIF) analysis was conducted using the statsmodels package, with a threshold of 5. Any variable equal to or greater than five was not used for further analysis. After VIF analysis, BorutaPy was applied to the data to select the variables for model building. Then, different machine learning techniques (including Random Forest, Gradient Boosting Machine, CatBoost, Extreme Gradient Boost, Logistic Regression, Decision Trees, Support Vector Machine, KNN + Univariate Feature Selection) were used and compared based on their cross-validation results. The results revealed that the Extreme Gradient Boost (XGboost) technique has the highest prediction scores at every venue and game outcome condition. Therefore, the XGboost technique was selected to build models for the games at home (model_home), away (model_away) and all venues (model_all). Before model construction, the dataset was randomly split into training (70%) and test (30%) sets using stratified sampling to preserve the proportions of match outcomes (win, draw, loss) across both sets. The scikit-learn library applied 5-fold stratified cross-validation to validate the results. Average accuracy, average precision, average recall, average F1-scores, and standard deviation results were recorded for each condition. The Feature Importance Rankings showed the variables' impact on the models. Additionally, Shapley (SHAP) analysis was applied to reflect the impact of the variables on the model output and their average impact. Finally, confusion matrices displayed the models' predictions for each class of game outcomes.

Figure 1
Study Design



RESULTS

The VIF analysis showed that no variables scored higher than five at any venue or game outcome condition. Therefore, the BorutaPy feature selection technique was applied to all variables in the dataset across all conditions. According to the results, 34 out of 106 dataset variables were selected for use in different models and conditions (Table 1). The chosen variables included 22 out of 61 technical, 12 out of 42 tactical, and no physical variables.

After BorutaPy feature selection, 16 variables for the win at venue_all condition (Low pressing %, Freekick attacks, Corner attacks, Attacks – left flank, Chances, Conversion rate, Shots on target, Key passes accurate, Air challenges, Ball interceptions, xG, Opponent's xG, xG Conversion, Opponent xG per shot, xG per goal, Opponent's xG per goal), 13 variables for loss at venue_all condition (Building ups, Chances, Conversion rate, Offsides, Shots on target, Shots on target %, Crosses, xG, Opponent's xG, xG Conversion, Opponent xG per goal, Lost balls in own half, Ball recoveries) and ten variables for draw at venue_all condition (Conversion rate, Shots on target, xG, Opponent's xG, xG Conversion, Opponent xG per shot, xG per goal, Opponent xG per goal, Lost balls, Penalties scored %) were selected. For venue_home condition, eight variables for win (Chances, Conversion rate, Shots on target, xG, Opponent's xG, xG Conversion, Opponent xG per shot, Opponent xG per goal), eight variables for loss (Attacks with shots Set pieces attacks, Chances, Conversion rate, Shots on target %, xG, Opponent's xG, xG Conversion, Opponent xG per goal), and five variables for draw (Conversion rate, Key passes accurate, xG, xG Conversion,

Opponent xG per goal) were selected. Nonetheless, for venue_away condition, BorutaPy selected 12 variables for the win (Team pressing successful, Set pieces attacks, Efficiency for attacks through the left flank %, Attacks with shots - right flank, Chances, Conversion rate, Shots on target, xG, Opponent's xG, Net xG, xG Conversion, Opponent xG per goal), ten variables for loss (Corner attacks with shots, Attacks with shots - center, Chances, Conversion rate, xG, Opponent's xG, Net xG, xG Conversion, xG per goal, Opponent xG per goal) and four variables for the draw (Shots on target, Defensive challenges, xG per goal, Opponent's passes per defensive action).

Table 1 shows the cross-validation results for three different models, built for venue_all (model_all), venue_home (model_home), and venue_away (model_away). The results were evaluated based on average accuracy, precision, recall, and F1-scores. When the cross-validation results for model_all were examined, the highest accuracy was in the win condition (95.93%), while the lowest was in the draw condition (88.46%). For precision results, the highest score was 95.87% in loss, and the lowest was in the draw condition (83.08%). Recall results showed the highest score of 95.04% in the win and the lowest score of 65.09% in the draw, similar to the accuracy results. In the last cross-validation run on model_all, the f1 score was highest in the win at $95.69 \pm 0.01\%$ and lowest in the draw at $82.76 \pm 0.01\%$. The feature importance ranking results show the impact of each variable on the model (Table 1). In model_all, 16 variables showed different impacts in the win condition; the highest score was for Conversion rate (0.1979), and the lowest was for Air challenges (0.0124). In the loss condition, 13 variables impacted the model. xG Conversion had the highest score with 0.3232, and Offsides had the lowest score, 0.0170. In the draw condition, ten variables had an impact on the model. The highest impact was found by Conversion rate (0.2016), and the lowest by Lost balls (0.0636).

The cross-validation results for model_home showed that the highest accuracy was in the win condition (95.84%) and the lowest in the draw condition (79.96%). Precision scores were the highest in the win (95.62%) and the lowest in the draw condition (60.20%). When the recall results were examined, the highest score was 95.59% in the condition win, and the lowest was 53.35% in the draw condition. The highest f1 score for model_home was found in the win condition ($95.82 \pm 0.01\%$) and the lowest in the draw condition ($71.07 \pm 0.04\%$). To model_home, eight variables showed different impacts in the win condition, according to the Feature Importance Ranking results. Conversion rate showed the highest impact with 0.2795, and Opponent's xG per shot the lowest (0.0406). In the loss condition, eight variables had impacts, where xG Conversion had the highest (0.3611) and the Attacks with shots, Set pieces attacks had the lowest (0.0370). Five

variables impacted the model in the draw condition. The Conversion rate had the highest (0.3547), and Key passes accurate showed the lowest impact on model_home.

For model_away, the highest accuracy was observed in the win condition (94.51%), while the lowest was in the draw condition (68.99%). For precision, the highest score was 93.62% in the loss condition, while the lowest was 29.41% in the draw condition. Recall results were the highest in the loss condition (93.97%) and the lowest in the draw condition (22.18%). The F1 score was also highest in the loss condition (94.11 ± 0.03%) and lowest in the draw condition (52.83 ± 0.03%), similar to the recall results. The Feature Importance Ranking results showed different impacts on model_away, similar to the other models. In the win condition, 12 variables impacted the model, where the Conversion rate has the highest with 0.2541, and Efficiency for attacks through the left flank% is the lowest with 0.0197. Among the ten variables that affected the model in the loss condition, xG per goal had the highest impact (0.4511), while Attacks with shots-center had the lowest (0.0225). In the draw condition, four variables showed different effects on the model, where Shots on target had the highest (0.2804), and Opponent's passes per defensive action were the lowest (0.2273).

Table 1
Feature Importance Rankings & Cross-Validation Results

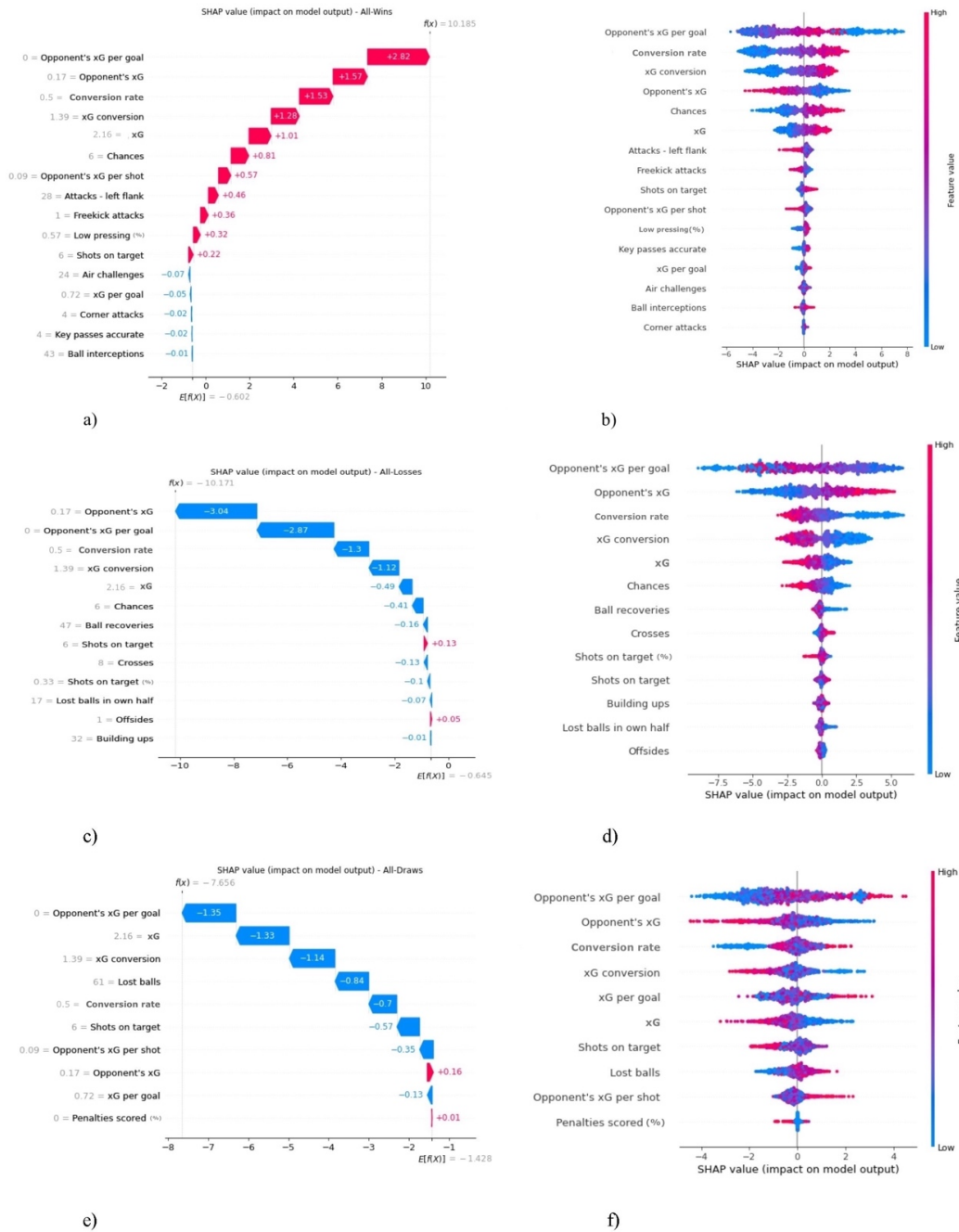
Model Game Outcome	All Win	All Loss	All Draw	Home Win	Home Loss	Home Draw	Away Win	Away Loss	Away Draw
Variables	n = 16	n = 13	n = 10	n = 8	n = 8	n = 5	n = 12	n = 10	n = 4
Conversion rate	0.1979	0.0737	0.2016	0.2795	0.0701	0.3547	0.2541	0.0463	-
xG Conversion	0.1334	0.3232	0.0975	0.1426	0.3611	0.1813	0.1883	0.0703	-
Shots on target	0.1314	0.0363	0.0867	0.1229	-	-	0.0303	-	0.2804
Chances	0.1122	0.0475	-	0.1693	0.0813	-	0.1152	0.0423	-
Opponent's xG per goal	0.0997	0.2092	0.1523	0.1389	0.2158	0.2068	0.1256	0.1554	-
xG per goal	0.0465	-	0.0763	-	-	-	-	0.4511	0.2579
Opponent's xG	0.0440	0.1064	0.0852	0.0638	0.1107	-	0.0604	0.0916	-
xG	0.0409	0.0658	0.0886	0.0424	0.0575	0.1454	0.0626	0.0354	-
Corner Attacks	0.0399	-	-	-	-	-	-	-	-
Opponent's xG per shot	0.0322	-	0.0698	0.0406	-	-	-	-	-
Freekick attacks	0.0279	-	-	-	-	-	-	-	-
Key passes accurate	0.0260	-	-	-	-	0.1118	-	-	-
Ball interceptions	0.0219	-	-	-	-	-	-	-	-
Low pressing (%)	0.0178	-	-	-	-	-	-	-	-
Attacks - left flanks	0.0159	-	-	-	-	-	-	-	-
Air challenges	0.0124	-	-	-	-	-	-	-	-

Table 1 (Continued)

Model Game Outcome	All Win	All Loss	All Draw	Home Win	Home Loss	Home Draw	Away Win	Away Loss	Away Draw
Variables	n = 16	n = 13	n = 10	n = 8	n = 8	n = 5	n = 12	n = 10	n = 4
Shots on target (%)	-	0.0346	-	-	0.0665	-	-	-	-
Ball recoveries	-	0.0251	-	-	-	-	-	-	-
Building ups	-	0.0220	-	-	-	-	-	-	-
Crosses	-	0.0213	-	-	-	-	-	-	-
Lost balls in own half	-	0.0180	-	-	-	-	-	-	-
Offsides	-	0.0170	-	-	-	-	-	-	-
Penalties scored (%)	-	-	0.0783	-	-	-	-	-	-
Lost balls	-	-	0.0636	-	-	-	-	-	-
Att. with shots - Set pieces	-	-	-	-	0.0370	-	-	-	-
Att. with shots - right flank	-	-	-	-	-	-	0.0475	-	-
Net xG	-	-	-	-	-	-	0.0378	0.0281	-
Set pieces attacks	-	-	-	-	-	-	0.0321	-	-
Team pressing successful	-	-	-	-	-	-	0.0265	-	-
Eff. for attacks - left flank (%)	-	-	-	-	-	-	0.0197	-	-
Corner attacks with shots	-	-	-	-	-	-	-	0.0570	-
Att. with shots - center	-	-	-	-	-	-	-	0.0225	-
Defensive challenges	-	-	-	-	-	-	-	-	0.2343
Opp. Pass per Def. Action	-	-	-	-	-	-	-	-	0.2273
Average Accuracy	95.93%	95.36%	88.46%	95.84%	93.38%	79.96%	94.51%	94.13%	68.99%
Average Precision	94.38%	95.87%	83.08%	95.62%	93.64%	60.20%	92.45%	93.62%	29.41%
Average Recall	95.04%	91.80%	65.09%	95.59%	83.13%	53.35%	88.90%	93.97%	22.18%
F1 Score	%95.69 ± 0.01	%95.03 ± 0.02	%82.76 ± 0.01	%95.82 ± 0.01	%91.64 ± 0.02	%71.07 ± 0.04	%93.28 ± 0.02	%94.11 ± 0.03	%52.83 ± 0.03

The results of the SHAP analysis for each condition are shown in Figures 2, 3, and 4. The graphs on the left side (a, c, and e for each Figure) reflect the average impact of each variable on the model output. On the other hand, the graphs on the right side (b, d, f) show the effect of the values for each sample, scaled vertically by different colors on the right, on the model output.

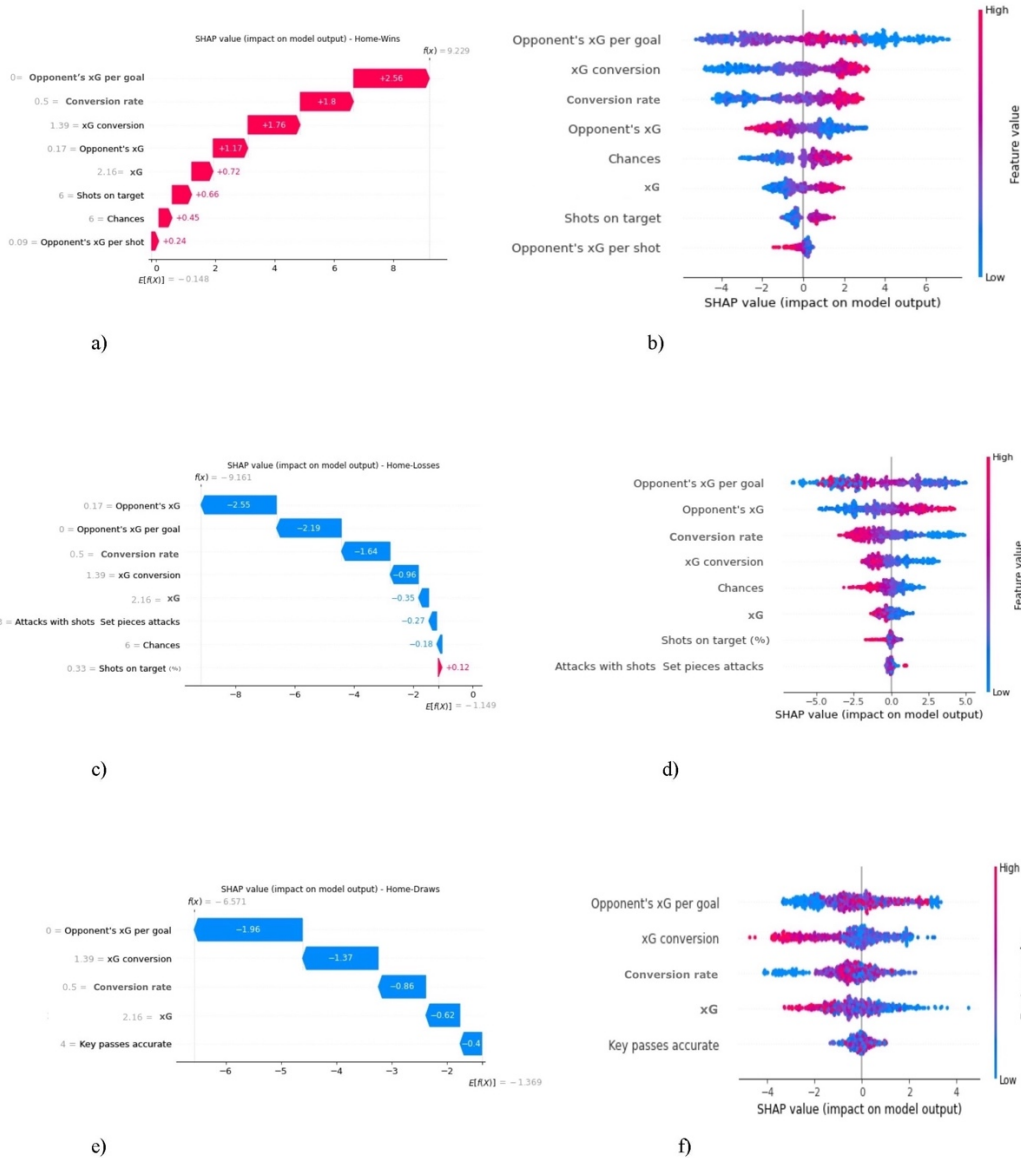
Figure 2
SHAP Results for model_all



The highest impact on the outputs of the model_all in the win condition was found positively for the Opponent's xG per goal with the value 2.82, and besides that, Opponent's xG (1.57) and Conversion rate (1.53) variables also had positive effects on the outputs of the model (Figure 2. a). When the results for the loss condition were examined, the Opponent's xG (-3.04) and Opponent's xG per goal (-2.87) showed the highest impacts on the model output negatively (Figure 2. c). For

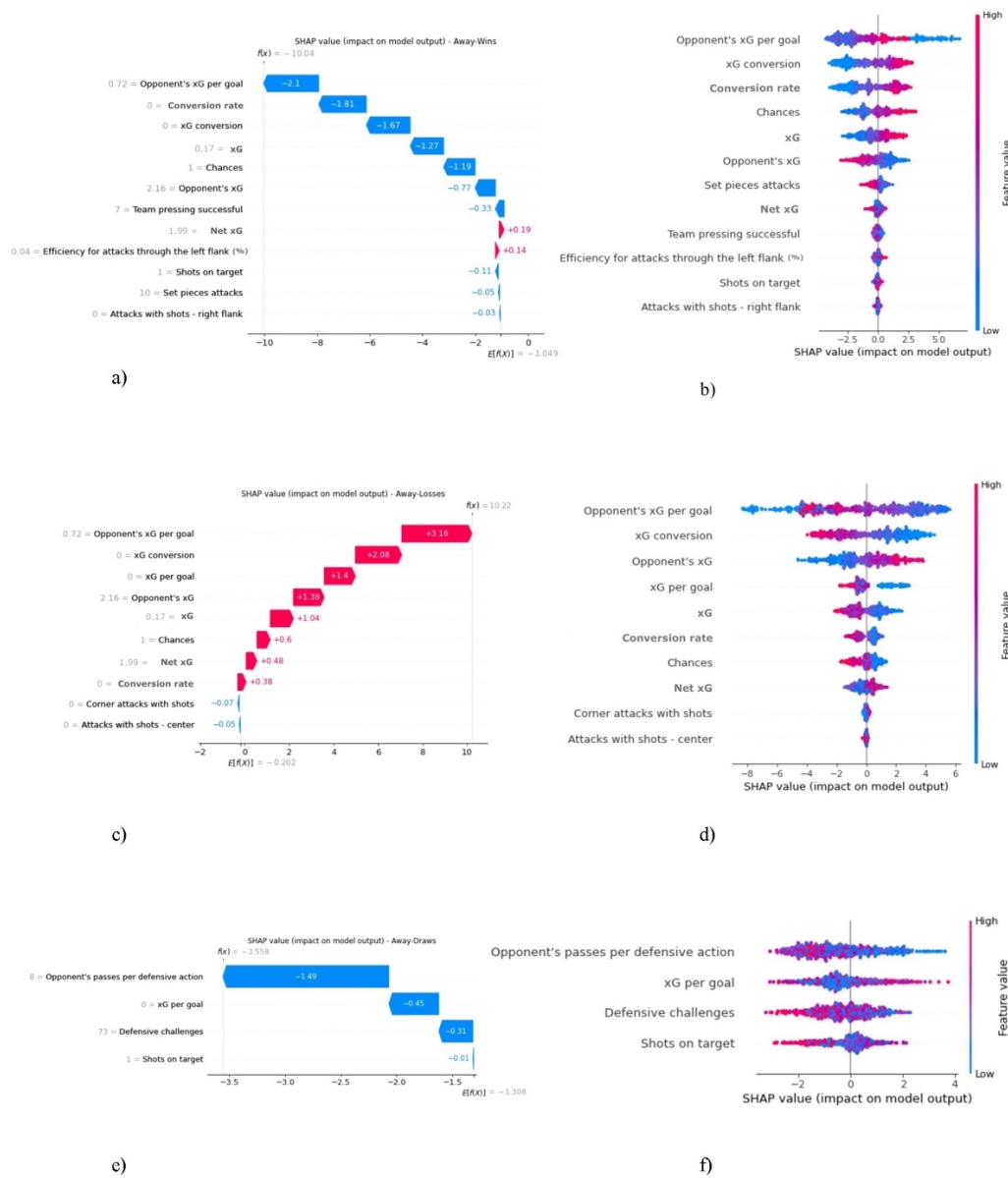
the effects on the outputs of model_all in the draw condition, the highest impacts were found in Opponent's xG per goal (-1.35), xG (-1.33), and xG Conversion (-1.14) negatively (Figure 2. e).

Figure 3
SHAP Results for model_home



In the win condition, the highest impacts on the outputs for model_home were observed positively in Opponent's xG per goal with 2.56, Conversion rate with 1.80, and xG Conversion with 1.76 (Figure 3. a). For the loss condition, the highest impacts were found negatively in Opponent's xG (-2.55), Opponent's xG per goal (-2.19), and Conversion rate (-1.64) (Figure 3. c). When the results for the draw condition were observed, the highest impacts were found negatively in the Opponent's xG per goal (-1.96) and xG Conversion (-1.37) (Figure 3. e).

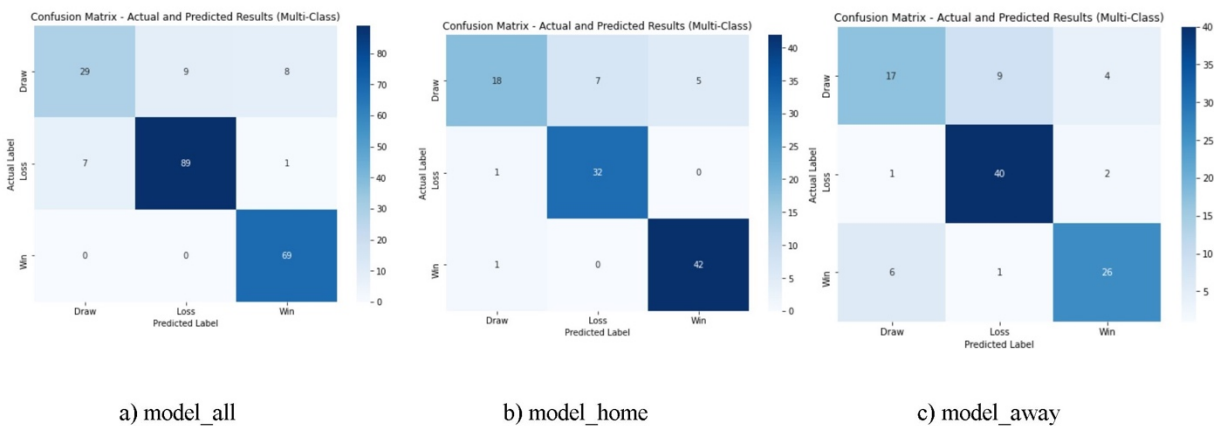
Figure 4
SHAP Results for model_away



The highest impacts on the output for model_away in win condition were examined negatively in Opponent's xG per goal (-2.10), in Conversion rate (-1.81), and in xG Conversion (-1.67) (Figure 4. a). In the loss condition, the highest impacts were found positively in the Opponent's xG per goal at 3.16, xG Conversion at 2.08, xG per goal at 1.40, and Opponent's xG at 1.38 (Figure 4. c). When the impacts for the outputs on model_away in draw condition were observed, the highest impact was found negatively in the Opponent's passes per defensive action (-1.49) (Figure 4. e).

The results of the confusion matrices for three different models (model_all, model_home, and model_away) are shown in Figure 5. For each model, a random sample dataset was selected, corresponding to 30% of the entire dataset, and used to test its ability to predict game results. Figure 5.a reflects that 69 out of 69 games (100%) were correctly classified for the win condition by the model_all, whereas for the loss condition, 89 out of 97 games (91.8%), and the draw condition, 29 out of 46 games (63.0%) were classified correctly. For model_home, the results revealed that 42 out of 43 games (97.7%) for the win, 32 out of 33 games (97.0%) for loss, and 18 out of 30 games (60.0%) for draw conditions were correctly classified (Figure 5.b). The results for model_away are given in Figure 5.c. For the win condition, 26 out of 33 (78.8%) games, for the loss condition, 40 out of 43 (93.0%) games, and for the draw, 17 out of 30 (56.7%) games were correctly classified by the model.

Figure 5
Confusion Matrices for Three Models



DISCUSSIONS

The present study aimed to determine which performance indicators strongly predict game outcomes (win, loss, and draw) using a machine learning model in the Turkish Super League over two consecutive seasons.

The BorutaPy variable selection was primarily applied to determine which variables would create a more accurate model for that condition before modeling. The results of BorutaPy revealed that 22 out of 61 technical variables, 12 out of 42 tactical variables, and none of the physical variables were selected for the models. A total of 34 out of 106 variables in the dataset were used in different models and conditions. The considerable difference in the number of selected variables across match venues likely reflects variations in data distribution and the relative influence of

performance indicators under distinct contextual conditions. While XGBoost is robust to multicollinearity, BorutaPy improved interpretability by identifying the most consistent and context-dependent predictors.

According to the Feature Importance Ranking analysis results, Conversion rate, and xG Conversion are the variables with the highest impact on the models in seven conditions (Table 1). Significantly, the Conversion rate had the greatest effect on the models in five conditions (except in all-loss conditions and the draw condition in model_away). The variable with the highest impact on the models among all conditions is xG per goal in the model_away loss condition. The Conversion rate variable is the percentage of chances created by the teams in the matches that are converted into goals. The observation that these values affect all models, regardless of whether they are high or low, underscores the variable's importance, particularly in win conditions. In soccer, the nature of the game is to score goals and change the result, thereby creating opportunities to achieve that purpose. Therefore, the importance of this variable is logical and expected.

While cross-validation results are evaluated using accuracy, precision, recall, and F1 scores, one of the most important findings in these evaluations is the accuracy of the models selected in previous studies (Cho et al., 2018; Joseph et al., 2006; Kvasnytsya et al., 2024). When the study's accuracy values are examined, model_all generally performs better than model_home and model_away. Accordingly, the correct prediction rate of model_away is lower than the other models. When examining the in-model results across the match outcome conditions in all three models, it is observed that draw predictions are less accurate than those for wins and losses (Table 1). In soccer matches, a draw can be achieved by teams with lower performance and lower frequencies of performance variables, especially when they are down by one goal. Alternatively, it can result from the inability of more powerful teams to find the winning goal, despite having high frequencies of the variables they possess. Especially considering the draw condition in model_away, which has the lowest accuracy among all models, this may be because home teams cannot gain the score advantage to win, away teams get points with the score equality, and performance variables cannot explain this. The impact of goal difference was not considered in the present study. In a study where goal difference was analyzed and inferences were made on match results, it was stated that the effect of a goal, even if scored by 'luck,' on the match result was exceptionally high in low-scoring sports like soccer (Geurkink et al., 2021). This conclusion also supports the low accuracy of the predictions. In addition, variables such as xG Conversion and Conversion rate, which have the highest impact on the models across seven conditions in

total, were not selected for the model_away draw condition. This may also be related to the low accuracy of the results. The highest accuracy values are observed in the win condition for all models (Table 1). In a study on predicting match results, a single model was created independent of venue effects; only the win or loss condition was evaluated, and the accuracy was 89.6% (Geurkink et al., 2021). Another study using a machine-learning model to predict match results showed 83.3% and 72.7% correct predictions for wins and losses, respectively (Hassan et al., 2020). In another study in which opponent quality was assessed, the prediction rate varied between 67.9% and 78.4% depending on the opponent. When all match results were evaluated in the same study, the rate was 64.8% (Bilek & Ulas, 2019). A recent study reported an effective approach achieving a 0.1 error rate and 87% accuracy in predicting offensive and defensive actions (Rahimian et al., 2024). It can be clearly stated that the accuracy values obtained in the present study are higher than those reported in other studies.

According to the accuracy results, the needs and variables required to reach the winning condition seem more predictable than those for the draw condition. When the SHAP analyses, revealing the effect of each variable on the outputs of the model, are examined, the variables that most affect the win condition for model_all are Opponent's xG per goal, Opponents' xG and Conversion rate in the positive direction (Figure 2.a). It is seen that high frequencies (red color) of the Opponent's xG per goal variable have a low impact on the model output in the model_all win condition. In contrast, low frequencies (blue color) have a high effect (positive and negative) on the model output (Figure 2.b). The Opponent's xG variable, which shows the opponent's goal expectation, indicates that opponents' low xGs (blue) positively affect the win condition. In contrast, high xGs (red color) have a negative effect (Figure 2.b).

According to the results of the SHAP analysis, the most influential variables in the model outputs in win condition of 3 models are the Opponent's xG per goal, Conversion rate, Opponent's xG, and xG Conversion. These variables positively affected the win condition in model_all and model_home (Figure 2.a, Figure 3.a) and a negative effect in model_away (Figure 4.a). For the Opponent's xG per goal variable, the positive impact on model_all and model_home stemmed from its low values. Similarly, for model_away, low values of the variable had a positive effect, but the effect range was less in this model (Figure 4.b). For Conversion rate and xG Conversion variables, high values positively affect the win condition in all three models. In contrast, for Opponent's xG, low values positively affect win conditions in all models. When the loss condition was analyzed, the variables with the greatest impact on the models' results were Opponent's xG, Opponent's xG per goal, Conversion rate, xG Conversion, and xG per goal. All of these variables

had a negative effect on the loss condition for model_all and model_home, while they had a positive impact for model_away. While low values of the Conversion rate and xG Conversion variables positively affected the loss condition across all three models, high values negatively affected it. In the Opponent's xG, the contrary situation was observed. In the draw condition, the highest impact on model_all and model_home outputs is seen in Opponent's xG per goal, xG, and xG Conversion variables, while in model_away, the highest impact is seen in Opponent's passes per defensive action variable. All of these variables negatively affected draw conditions across all models. The effects of high and low values of the variables on the model outputs vary as positive and negative (Figure 2.f, Figure 3.f, Figure 4.f). This variability explains why the cross-validation accuracy for the draw condition is lower across all three models.

Previous studies have indicated that total shots and shots on target variables are related to match results, especially wins (Castellano et al., 2012; Lago-Peñas et al., 2010; Moreira Praça et al., 2023). As a result of this study, xG-related variables were found to be more effective than other variables in different match outcome conditions. xG considers parameters such as the location of the shot, the distance to the goal, the angle of the shot, the number of opposing players at the moment of the shot, the position of the goalkeeper, etc (Rathke, 2017). The results revealed that this variable is a reliable and effective metric. As much as the xG metrics of the team's offensive performance, the opponent's xG values also impact the win, loss, and draw conditions as a defensive performance indicator.

In this study, 542 matches were analyzed. The correct classification rates for the three models, trained under different venue conditions, ranged from 56.7% to 100% (Figure 5). The highest correct classifications were generally in win and loss conditions, while draw conditions gave the lowest values. As similarly stated in another study, it is challenging to discriminate the draws. However, technical performance can effectively explain wins and losses (Pappalardo & Cintia, 2018). The classification results obtained in the draw condition agree with the cross-validation results for the present study. The unexpected outcome is that the selected sample data exhibits low accuracy in classifying match results in the model_away win condition. In the win condition, regardless of the venue, 100% accuracy was observed in the classification made by model_all, while this rate was 78.8% for model_away. While the cross-validation results give high values for the win condition in model_away, this is likely due to the sample data selected for model_away in the classification results (Figure 5).

Limitations

As this study is among the earliest to provide insights into machine learning models for predicting match results in the Turkish Super League, it has its limitations. As the physical variables considered as performance indicators in this study are few, they may need to be more numerous to adequately explain performance and assess their effects on the models. Another area for improvement for the present study is that the time-related data were not considered due to the lack of dataset diversity of the data provider.

CONCLUSION

It is essential to understand what factors affect match results on the path to winning in soccer. As a result of this study, it was found that technical variables, particularly xG-related variables (Opponent's xG per goal, Opponent's xG, and xG Conversion), are effective in predicting match results for the Turkish Super League in models and model outputs. The models showed high cross-validation scores for both wins and losses under whole-venue conditions, but lower scores for draws. Similarly, in the classification of match results, the confusion matrix shows that draw matches are less accurately classified than wins and losses across all venue conditions.

PRACTICAL IMPLICATIONS

The increase in the quantity and quality of performance indicators in elite football has made the data more complex. New analytical approaches are needed to evaluate match performance across different factors, such as venue (home or away), match result (win, lose, or draw), player position, and tournament characteristics. These approaches are expected to provide coaches with a useful new perspective on determining their teams' strategies for winning home and away matches. In this context, the results of this study show that xG-related variables (*Opponent's xG per goal*, *Opponent's xG*, and *xG conversion*) are the most effective in predicting wins for teams in the Turkish Super League, regardless of venue. Similar effects of these variables have been reported in other studies (Forcher et al., 2025). When developing defensive strategies to limit opponents' shooting opportunities and lower the expected goals from their shots, coaches should prioritize monitoring these variables and developing tactical organization, such as field positioning and timing of pressure. In terms of offensive strategies, a high *Conversion rate* has been found to be effective. Therefore, coaches should focus on tactical drills that increase the number of goal-scoring opportunities in the opponent's half, as well as the frequency and accuracy of shots that successfully convert these opportunities into goals.

Acknowledgments

The authors reported that no funding was associated with the work featured in this article.

Authors' Contribution

First and Third authors have made substantial contributions to the conception or the design of the manuscript, and second author has contributed to the analysis and interpretation of the data. All authors have participated in drafting the manuscript; first author revised it critically. All authors read and approved the final version of the manuscript.

Declaration of Conflict Interest

The authors declare that there is no conflict of interest regarding the publication of this paper.

Ethics Statement

This study is approved by the Social Sciences and Humanities Research Ethics Committee of Halic University (10.07.2024, 06).

Article History

Received 21 April 2025/Revised 24 October 2025/

Accepted 12 December 2025 / Available Online 7 March 2026

REFERENCES

- Badiella, L., Puig, P., Lago-Peñas, C., & Casals, M. (2023). Influence of Red and Yellow cards on team performance in elite soccer. *Annals of Operations Research*, 325(1), 149-165. <https://doi.org/10.1007/s10479-022-04733-0>
- Bilek, G., & Ulas, E. (2019). Predicting match outcome according to the quality of opponent in the English premier league using situational variables and team performance indicators. *International Journal of Performance Analysis in Sport*, 19(6), 930-941. <https://doi.org/10.1080/24748668.2019.1684773>
- Bunker, R., Yeung, C., & Fujii, K. (2024). *Machine learning for soccer match result prediction*. In Artificial Intelligence, Optimization, and Data Sciences in Sports (pp. 7-49). Cham: Springer Nature Switzerland. https://doi.org/10.1007/978-3-031-76047-1_2
- Carling, C., Williams, A. M., & Reilly, T. (2007). *Handbook of Soccer Match Analysis: A systematic approach to improving performance*. Routledge.
- Castellano, J., Casamichana, D., & Lago, C. (2012). The use of match statistics that discriminate between successful and unsuccessful soccer teams. *Journal of Human Kinetics*, 31(1), 139-147. <https://doi.org/10.2478/v10078-012-0015-7>

- Cho, Y., Yoon, J., & Lee, S. (2018). Using social network analysis and gradient boosting to develop a soccer win-lose prediction model. *Engineering Applications of Artificial Intelligence*, 72, 228–240. <https://doi.org/10.1016/j.engappai.2018.04.010>
- Dellal, A., Chamari, K., Wong, D. P., Ahmaidi, S., Keller, D., Barros, R., Bisciotti, G. N., & Carling, C. (2011). Comparison of physical and technical performance in European soccer match-play: Fa Premier League and La Liga. *European Journal of Sport Science*, 11(1), 51–59. <https://doi.org/10.1080/17461391.2010.481334>
- Forcher, L., Forcher, L., Woll, A., & Altmann, S. (2025). AI in Bundesliga match analysis – expected possession value (EPV) vs. expected goals (xG) to predict match outcomes in soccer. *Frontiers in Sports and Active Living*, 7, 1713852. <https://doi.org/10.3389/fspor.2025.1713852>
- Geurkink, Y., Boone, J., Verstockt, S., & Bourgois, J. G. (2021). Machine learning-based identification of the strongest predictive variables of winning and losing in Belgian professional soccer. *Applied Sciences (Switzerland)*, 11(5), 1–11. <https://doi.org/10.3390/app11052378>
- Gomez, M. A., Lago-Peñas, C., & Owen, L. A. (2016). The influence of substitutions on elite soccer teams' performance. *International Journal of Performance Analysis in Sport*, 16(2), 553–568. <https://doi.org/10.1080/24748668.2016.11868908>
- Hassan, A., Akl, A. R., Hassan, I., & Sunderland, C. (2020). Predicting wins, losses and attributes' sensitivities in the soccer world cup 2018 using neural network analysis. *Sensors (Switzerland)*, 20(11). <https://doi.org/10.3390/s20113213>
- Hughes, M. D., & Bartlett, R. M. (2002). The use of performance indicators in performance analysis. In *Journal of Sports Sciences (Vol. 20, Issue 10, pp. 739–754)*. <https://doi.org/10.1080/026404102320675602>
- Joseph, A., Fenton, N. E., & Neil, M. (2006). Predicting football results using Bayesian nets and other machine learning techniques. *Knowledge-Based Systems*, 19(7), 544–553. <https://doi.org/10.1016/j.knosys.2006.04.011>
- Kubayi, A. (2020). Analysis of Goal Scoring Patterns in the 2018 FIFA World Cup. *Journal of Human Kinetics*, 71(1), 205–210. <https://doi.org/10.2478/hukin-2019-0084>
- Kubayi, A., & Larkin, P. (2020). Technical performance of soccer teams according to match outcome at the 2019 FIFA Women's World Cup. *International Journal of Performance Analysis in Sport*, 20(5), 908–916. <https://doi.org/10.1080/24748668.2020.1809320>
- Kvasnytsya, O., Tyshchenko, V., Latyshev, M., Kvasnytsya, I., Kirsanov, M., Plakhotniuk, O., & Buhaiiov, M. (2024). Team Performance Indicators That Predict Match Outcome in Rugby Union. *Pamukkale Journal of Sport Sciences*, 15(1), 203–216. <https://doi.org/10.54141/psbd.1342340>
- Lago-Peñas, C., Lago-Ballesteros, J., Dellal, A., & Gómez, M. (2010). Game-related statistics that discriminated winning, drawing and losing teams from the Spanish soccer league. *Journal of Sports Science and Medicine*, 9, 288–293. <http://www.jssm.org>
- Lago-Peñas, C., Lago-Ballesteros, J., & Rey, E. (2011). Differences in performance indicators between winning and losing teams in the UEFA Champions League. *Journal of Human Kinetics*, 27(1), 135–146. <https://doi.org/10.2478/v10078-011-0011-3>

- Mohr, M., Krstrup, P., & Bangsbo, J. (2003). Match performance of high-standard soccer players with special reference to development of fatigue. *Journal of Sports Sciences*, 21(7), 519–528. <https://doi.org/10.1080/0264041031000071182>
- Moreira Praça, G., Braga Jacinto, A. L., de Sousa Pinheiro, G., de Oliveira Abreu, C., & Teoldo da Costa, V. (2023). What are the key performance indicators related to winning matches in the German Bundesliga? *International Journal of Performance Analysis in Sport*, 23(4), 284–295. <https://doi.org/10.1080/24748668.2023.2227923>
- Pappalardo, L., & Cintia, P. (2018). Quantifying the relation between performance and success in soccer. *Advances in Complex Systems*, 21(3–4). <https://doi.org/10.1142/S021952591750014X>
- Perl, J., Grunz, A., & Memmert, D. (2013). Tactics Analysis in Soccer-An Advanced Approach. In *International Journal of Computer Science in Sport* (Vol. 12, Issue 1).
- Rahimian, P., Mihalyi, B. M., & Toka, L. (2024). In-game soccer outcome prediction with offline reinforcement learning. *Machine Learning*, 113(10), 7393-7419. <https://doi.org/10.1007/s10994-024-06611-1>
- Rampinini, E., Impellizzeri, F. M., Castagna, C., Coutts, A. J., & Wisløff, U. (2009). Technical performance during soccer matches of the Italian Serie A league: Effect of fatigue and competitive level. *Journal of Science and Medicine in Sport*, 12(1), 227–233. <https://doi.org/10.1016/j.jsams.2007.10.002>
- Rathke, A. (2017). An examination of expected goals and shot efficiency in soccer. *Journal of Human Sport and Exercise*, 12(Proc2). <https://doi.org/10.14198/jhse.2017.12.proc2.05>
- Rein, R., & Memmert, D. (2016). Big data and tactical analysis in elite soccer: future challenges and opportunities for sports science. *SpringerPlus*, 5(1). <https://doi.org/10.1186/s40064-016-3108-2>
- Sampaio, J. (2013). *Routledge Handbook of Sports Performance Analysis*. In T. McGarry, P. O'Donoghue, & J. Sampaio (Eds.), *Routledge Handbook of Sports Performance Analysis*. Taylor&Francis Group. <https://doi.org/10.4324/9780203806913>
- Settembre, M., Buchheit, M., Hader, K., Hamill, R., Tarascon, A., Verheijen, R., & McHugh, D. (2024). Factors associated with match outcomes in elite European football—insights from machine learning models. *Journal of Sports Analytics*, 10(1), 1-16. <https://doi.org/10.3233/JSA-240745>
- Silva, H., & Marcelino, R. (2023). Inter-operator reliability of InStat Scout in female football games. *Science and Sports*, 38(1), 42–46. <https://doi.org/10.1016/j.scispo.2021.07.015>
- Taylor, B. J., Mellalieu, D. S., & James, N. (2005). A Comparison of Individual and Unit Tactical Behaviour and Team Strategy in Professional Soccer. *International Journal of Performance Analysis in Sport*, 5(2), 87–101. <https://doi.org/10.1080/24748668.2005.11868329>
- Yang, G., Leicht, A. S., Lago, C., & Gómez, M. Á. (2018). Key team physical and technical performance indicators indicative of team quality in the soccer Chinese super league. *Research in Sports Medicine*, 26(2), 158–167. <https://doi.org/10.1080/15438627.2018.1431539>

APPENDICES

Appendix 1

The Descriptions of the Technical, Tactical and Physical Variables

Technical variables	Description
% scored free kick shots	Percentage share of goals scored by shots in freekicks.
Accurate crosses (%)	Percentage share of successful crosses in the total number of crosses.
Accurate passes	Total number of accurate passes.
Accurate passes (%)	Percentage share of accurate passes in the total number of passes.
Air challenges	Two players of the opposing teams challenging for the ball in the air, at least above shoulder height, the rivals play or try to play with their heads.
Air challenges won	Successful attempt of air challenge that leads to a touch made by own team player.
Air challenges won (%)	Percentage share of air challenges won in the total number of air challenges.
Attacking challenges	Challenges involving a player of the team that currently possesses the ball.
Attacking challenges won	Successful attempts of attacking challenges that lead to the ball remaining in possession of own team.
Ball interceptions	Player's active, targeted and successful action to either prevent a potentially accurate pass or to change the ball trajectory.
Ball recoveries	First player's action in a team's ball possession after the team started possessing the ball, except for the cases when Ball Possession starts from a set piece (including a throw-in).
Ball recoveries in opponent's half	Ball recoveries occurred in team's opponent's half of the pitch.
Blocked shots	-
Challenges	The summary type of a parameter, includes duels for the neutral balls, air duels for the neutral balls, dribbles, tackles and losing the ball during opponent tackling attempts; the total amount of attacking and defensive challenges.
Challenges in attack won (%)	Percentage share of attacking challenges won in the total number of challenges.
Challenges in defence won (%)	Percentage share of defensive challenges won in the total number of challenges.
Challenges won	Successful challenge is registered for a player of a team that keeps possession of a ball after such challenge; lost challenge is simultaneously registered for a player's opponent.
Challenges won (%)	Percentage share of challenges won in the total number of challenges.
Chances	A goal-scoring opportunity, when the attacking team gets a clear-cut chance to score a goal.
Conversion rate	Percentage of chances created by the teams in the games that are converted into goals
Corners	Awarded after a ball being sent across the sideline of the own half of the field by a defending team player.
Crosses	A pass into the box from the flanks in the opponent's half of the field; strong and directed pass. It can be performed both in the air and on the ground, and it cannot be an action performed from a set piece.
Crosses accurate	Total number of accurate crosses.
Defensive challenges	Challenges involving a player of the team that does not currently possess the ball; the number of defensive challenges of the team is always equal to the number of attacking challenges of their opponents.
Defensive challenges won	Successful attempts of defensive challenges that lead to a touch made by own team player.
Dribbles	Is an active action performed by a player in order to get through an opponent; can be performed as a trick or fake movement, as a ball poked at speed, no-touch ball etc.
Dribbles successful	Successful attempt of a dribble. as a result a player committing a dribble always keeps the ball and improves his position, leaving the opponent behind.
Expected points	-
Fouls	Action that impedes the progress and success of the opposing team and obtaining an advantage by breaking the rules of the game.
Free ball pick ups	Recovering a neutral ball after an opponent lost it.
Freekick shots	-

Appendix 1 (Continued)

Technical variables	Description
Key passes	A pass to a partner who is in a goal scoring position (one-on-one situation. empty net etc.) or a pass to a partner that “cuts off” the whole defensive line of the opponent’s team (3 and more players) in the attacking phase.
Key passes accurate	Successful attempt of a key pass, when a teammate touches a ball; if a challenge was registered after a key pass, this pass is still considered as a “key pass accurate”.
Lost balls	It is registered when a player loses the ball by a poor trapping of the ball, errant pass, unsuccessful attempt to shoot or an unsuccessful dribble.
Lost balls in own half	Lost balls occurred in team’s own half of the pitch.
Offsides	A player is in an offside position if: any part of the head, body or feet is in the opponents’ half (excluding the halfway line) and any part of the head, body or feet is nearer to the opponents’ goal line than both the ball or the second-last opponent.
Opponent's passes per defensive action	-
Passes	An attempt to transfer a ball from one teammate to another with the purpose of attack build-up or keeping the possession.
Penalties	Total number of penalties taken in the game.
Penalties scored	-
Penalties scored (%)	Percentage share of goals scored by penalty in the total number of penalties.
Shots	Total number of all shots made during the course of a game; includes shots on target. shots wide, blocked shots and shots on post / bar.
Shots on post / bar	-
Shots on target	Shots going inside the goal, might end in a goal or be deflected by the goalkeeper or by a field player from the GK zone.
Shots on target (%)	Percentage share of shots on target in the total number of shots.
Shots wide	-
Successful actions	Successfully completed actions out of total actions.
Successful actions (%)	Percentage share of successfully completed actions in total actions.
Successful dribbles (%)	Percentage share of successful dribbles in the total number of dribbles.
Tackles	This parameter is registered automatically for own team player in case an opponent is making a dribbling attempt; successful or unsuccessful tackle depends on the success of a dribble.
Tackles successful	Successful attempt of a tackle, as a result an opponent’s player loses the ball while performing a dribble.
Tackles won (%)	Percentage share of successful tackles in the total number of tackles.
Total actions	Total number of all types of passes (including crosses and set pieces passes). challenges, interceptions, picking up free balls, dribbling, bad ball controls and all kinds of shots (including goals), shots saved and goals conceded. Fouls are not included in total actions.
xG	Expected Goal; defines the goal probability of each shot for the team. It is depending on the situations like; the distance from the goal, shot angle, the number of opponents between the player and the goal etc.
xG conversion	Sum of the xG values of the goals.
xG per goal	The ratio of dividing the team's xG value to team's scored goals.
xG per shot	The ratio of dividing the team's xG value to team's total shot attempts.
Opponent's xG	-
Opponent's xG per goal	-
Opponent's xG per shot	-
Net xG (xG - Opponent's xG)	The difference between xG values of the team and the opponent in the same game.
Efficiency for corner attacks (%)	Percentage share of corner-attacks with a shot in the total number of corner-attacks.
Efficiency for counterattacks (%)	Percentage share of counter-attacks with a shot in the total number of counter-attacks.
Efficiency for freekick attacks (%)	Percentage share of freekick-attacks with a shot in the total number of freekick-attacks.
Efficiency for positional attacks (%)	Percentage share of positional attacks with a shot in the total number of positional attacks.

Appendix 1 (Continued)

Tactical variables	Description
Efficiency for set piece attacks (%)	Percentage share of set-piece attacks with a shot in the total number of set-piece attacks.
Efficiency for throwin attacks (%)	Percentage share of throwin attacks with a shot in the total number of throwin attacks.
Attacks - center	Attacks occurred between the space of left-side and right-side attacks, or central zone; the attack is determined by the last action of an attack which isn't a shot or a goal and which didn't occur inside the penalty area; determined for positional attacks and counter-attacks only.
Attacks - left flank	Attacks occurred on the width of 20 meters from the left sideline, whole length of the sideline is considered; the attack is determined by the last action of an attack which isn't a shot or a goal and which didn't occur inside the penalty area.
Attacks - right flank	Attacks occurred on the width of 20 meters from the right sideline, whole length of the sideline is considered; the attack is determined by the last action of an attack which isn't a shot or a goal and which didn't occur inside the penalty area; determined for positional attacks and counter-attacks only.
Attacks with shots Set pieces attacks	Set pieces attacks included at least one shot of any type from the attacking side.
Attacks with shots - center	Central zone attacks included at least one shot of any type from the attacking side.
Attacks with shots - left flank	Left-side attacks that included at least one shot of any type from the attacking side.
Attacks with shots - right flank	Right-side attacks included at least one shot of any type from the attacking side.
Building ups	Number of team possessions during the preparation of carrying the ball from own half to opponent's half.
Building ups without pressing	Number of building ups without opponent teams' pressing.
Corner attacks	-
Corner attacks with shots	A corner finished with a shot.
Counter attacks	Attack from the open play that starts with winning the ball from a defensive position and then quickly transitioning to offense while the prior attacking team is caught in an offensive formation; the length of possession during the attack cannot exceed 8 seconds before the possession transition or end; alternatively the length of possession can last between 8 and 30 sec., but the speed of attack cannot be less than 2.6 m/s. A counterattack cannot begin with a pass from a goalkeeper if he controlled the ball for more than 4 seconds before the action.
Counter attacks with a shot	Counter-attacks that included at least one shot of any type from the attacking side.
Efficiency for attacks through the central zone (%)	Percentage share of central zone attacks with shots in the total number of central zone attacks.
Efficiency for attacks through the left flank (%)	Percentage share of left-side attacks with shots in the total number of left-side attacks.
Efficiency for attacks through the right flank (%)	Percentage share of right-side attacks with shots in the total number of right-side attacks.
Entrances on final third of opponent's half	Number of team possessions during which at least one entrance into the opponent's final third was made
Entrances on opponent's half	Number of team possessions during which at least one entrance into the opponent's half was made. Entrance is counted in as a result of one of the following actions: pass, challenge, tackle, dribble, ball recovery, ball loss, foul, YC, RC, all kinds of shots, interception, free ball pick up, GK interception, cross.
Entrances to the opponent's box	Number of team possessions during which at least one entrance into the opponent's penalty box was made.
Freekick attacks	-
Freekick attacks with shots	A freekick finished with a shot.
Goals - Freekick attack	-
High pressing	Total number of collective attempts to force the opponents to lose the ball or to stop the development of an attack on the opponent's half of the pitch.
High pressing successful	Number of successfully performed high pressings.
Low pressing	Total number of collective attempts to force the opponents to lose the ball or to stop the development of an attack on one's own half of the pitch.

Appendix 1 (Continued)

Tactical variables	Description
Low pressing successful	Number of successfully performed low pressings
Positional attacks	All attacks from the open play that do not fit into counter attacks.
Positional attacks with shots	Positional attacks included at least one shot of any type from the attacking side.
Set pieces attacks	Total number of free-kick attacks, corner attacks, throw-in attacks and penalties.
Team pressing	Total number of collective attempts to force the opponents to lose the ball or to stop the development of an attack.
Team pressing successful	Number of successfully performed team pressings.
Throwin attacks	Attack started with a throwin.
Throwin attacks with shots	A throwin finished with a shot.
Physical variables	Description
High Speed Distance (m)	High speed distance covered in the game, which is the 20 km/h
Sprint Distance (m)	Sprint distance covered in the game, which is above 25 km/h
Total Distance (m)	The total distance covered in the game by whole team.