

UCLA

UCLA Electronic Theses and Dissertations

Title

Goal-Oriented Forecasting: Predicting Soccer Match Outcomes with Deep Learning

Permalink

<https://escholarship.org/uc/item/8917g70n>

Author

Chen, Sheng

Publication Date

2024

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA

Los Angeles

Goal-Oriented Forecasting: Predicting Soccer Match Outcomes with Deep Learning

A thesis submitted in partial satisfaction
of the requirements for the degree
Master of Applied Statistics & Data Science

by

Sheng Chen

2024

© Copyright by
Sheng Chen
2024

ABSTRACT OF THE THESIS

Goal-Oriented Forecasting: Predicting Soccer Match Outcomes with Deep Learning

by

Sheng Chen

Master of Applied Statistics & Data Science

University of California, Los Angeles, 2024

Professor Guido Montúfar, Chair

Soccer, often referred to as “football” in the heart of Europe, has a deep-rooted cultural significance that transcends national boundaries. The sport’s appeal extends far beyond the pitch, encompassing a wide array of enthusiasts, from die-hard fans to data-driven strategists. In recent years, the fusion of deep learning models with the captivating world of football has taken center stage, revolutionizing our approach to predicting match outcomes. We delve into the fusion of state-of-the-art artificial intelligence and machine learning techniques with the intricacies of a sport that inspires fervent devotion. Our aim is straightforward: to unearth the potential of deep learning models in enriching our capacity to anticipate which team will emerge victorious on the hallowed turf.

We will delineate the primary objectives of this research essay, elaborate on the methodology employed, and elucidate our anticipated contributions to the existing body of knowledge in both the realms of deep learning and sports analytics. Finally, we will underscore the significance of precise football match outcome prediction as an evolving and multi-dimensional research domain that holds great promise for aficionados, professionals, and researchers throughout Europe and beyond.

The thesis of Sheng Chen is approved.

Nicolas Christou

Frederic R. Schoenberg

Guido Montúfar, Committee Chair

University of California, Los Angeles

2024

*To my grandpa . . .
who dedicated years fostering
my mathematical education
and patience to overcome any difficulties*

TABLE OF CONTENTS

List of Figures	vii
List of Tables	viii
1 Introduction	1
2 Prior Work	2
2.1 Expected Goals (xG)	2
2.2 Video and Image Analysis	4
3 Data	6
3.1 Data Collection	6
3.2 Objective	6
3.3 Organizing Data	7
3.4 Exploratory Data Analysis	7
3.4.1 PCA	8
4 Methodology	12
4.1 Data Reduction	12
4.1.1 Tomek links	12
4.2 Data Augmentation	13
4.2.1 SMOTE and its variations (borderline, ADASYN)	14
4.2.2 Tabular GAN	14
4.3 Feature Selection	14
4.3.1 Exhaustive Search	15

4.3.2	Lasso	15
4.3.3	Mutual Information	16
4.4	Models	19
4.4.1	Linear Regression	19
4.4.2	XGBoost	19
4.4.3	MLP	20
4.4.4	Tabular Transformer	20
4.4.5	Graph Convolution Network	20
5	Result	22
5.1	Features Selected	22
5.2	Experiments	25
6	Future Work	27

LIST OF FIGURES

2.1	Comparison of two players in their shooting style. Neves clearly has a lower conversion rate; however, since his shots are mostly outside of the box (penalty area), the actual goals scored exceeds his XG (Expected Goals), while Sánchez is quite the opposite, who, as a striker, likes to finish more in the penalty area. Despite the difference in shot preferences due to positions, both are elite shooters since their resulting goals are above expected goals.	3
2.2	Each player's trajectory is first processed by a recurrent neural network (RNN) with shared parameters to capture sequential patterns. The results are then refined by a graph neural network (GNN) to account for the relational dynamics among players (Figures taken from Bauer and Anzer[10].)	5
3.1	Scatterplot of Straight Corners Against Goal Differences	9
4.1	Histogram of Goal Differences	13
4.2	Bar Plot of MI gain of added Input Features shows the top 7 features added sequentially are shown above.	18
5.1	Force Plot of the Features' Shapley Value	25

LIST OF TABLES

3.1	The first row indicates the outcome classes, and the second row are the counts of observations that falls into each category. Outcome variable comes with skewed class proportion, with more games end up around with low absolute goal differences, and only 7 games out of 8975 games have results greater than 7 goal differences.	8
3.2	XGBoost Classifier Accuracy after kernel PCA	11
5.1	Glossary of Feature Selected by added MI	23
5.2	Glossary of Additional Features Selected by Lasso	23
5.3	Model Result with Combination of Methods	26

CHAPTER 1

Introduction

In the world of soccer, predicting game outcomes and player performances has always been a challenging task. With the increasing availability of data, including statistics, team lineups, and player profiles, there is a growing interest in using data analysis and machine learning techniques to gain insights and make predictions. This document provides a set of strategies for predicting soccer game outcomes, specifically focusing on matches played in the German leagues spanning from the 2017-2018 season to the 2021-2022 season. We aim to understand which team features play the most important role in goal difference, predict future game outcomes, and explore the possibility of predicting games with different starting lineups.

Modern data analysis methods, including the utilization of Expected Goals (xG), a statistical tool developed by Opta [1], are becoming increasingly prevalent. xG evaluates a team's likelihood of scoring by considering various factors such as the quality and quantity of scoring opportunities generated. These methods and their applications will be elaborated upon in Chapter 2 of the study. In the field of betting, where more advanced models have become popular, among which are Machine Learning algorithms like Random Forests, Gradient Boosting, and deep neural networks. All these methods share the same idea: use past match data, player performance info, and even factors like weather and injuries to make accurate game predictions. In-play betting has made things more exciting by allowing changes to betting odds while a game is currently underway. Crucial for accurately predicting the outcome, specific measurements are taken to quantify how well both teams and individual players perform. This includes things like where players are on the field, how accurate their passes are, and factors related to their mindset. Furthermore, the industry closely monitors market performance, enhancing their ability to comprehend the game's evolution over time.

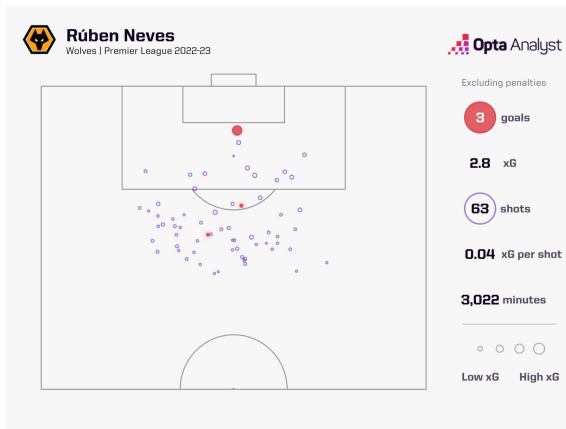
CHAPTER 2

Prior Work

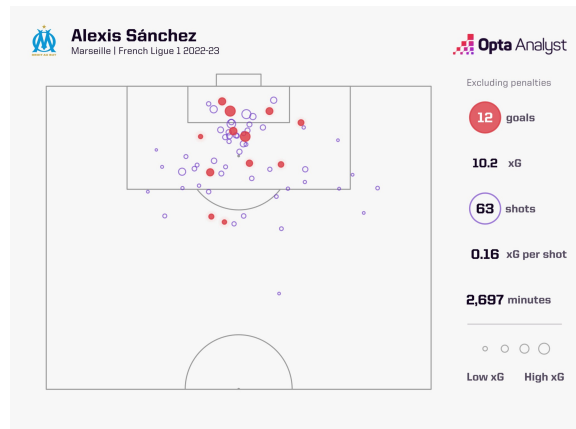
2.1 Expected Goals (xG)

StatsBomb Expected Goals (xG)[1] is a metric created by Opta, used to assess the likelihood of a soccer shot resulting in a goal. Calculated based on historical data from shots with similar characteristics, an xG model factors in elements such as distance to goal, angle, body part used for the shot, and type of assist. This model assigns a continuous value between 0 and 1 to each shot, representing the probability of it leading to a goal. For a more intuitive understanding, the 2022 update of the StatsBomb xG model gives all penalty kicks a shared static value of 0.78 xG, meaning out of 100 penalty shots, 78 are expected to be scored. This example is simple and intuitive since every penalty kick are taken at precisely the same spot every time. We can also focus on a single player's xg across his season. For example, when we look at Neves and Sanchez's xG shot maps in Fig. 2.1, despite they have both made exactly 63 shots that season, the comparison with xG emphasizes the difference in their preferred shooting position and conversion rate.

Opta's extensive data, which includes over 4.5 million shots with xG values for more than 100,000 players, enables us to compare and analyze the performances of players and teams worldwide. According to Lars Maurath's independent research[7], the correlation between expected goals (xG) and actual goals for a single season is projected to be in the range of 79% to 93%, within a 95% confidence interval, contingent on the model's quality.



(a) Rúben Neves's shot map



(b) Alexis Sánchez's shot map

Figure 2.1: Comparison of two players in their shooting style. Neves clearly has a lower conversion rate; however, since his shots are mostly outside of the box (penalty area), the actual goals scored exceeds his xG (Expected Goals), while Sánchez is quite the opposite, who, as a striker, likes to finish more in the penalty area. Despite the difference in shot preferences due to positions, both are elite shooters since their resulting goals are above expected goals.

2.2 Video and Image Analysis

In the convergence of computer vision and sports analytics, the combined force of Tracking Data Analysis and Video/Image Analysis emerges as a potent synergy. This integrated approach, propelled by high-resolution videos and computer vision techniques, combines precise player movements and object tracking techniques. By leveraging the strengths of both, these types of research not only advance soccer analytics but also underscore the potential of a unified model in the broader landscape of machine learning-driven computer vision research and real-time sports analysis.

Pascal Bauer and Gabriel Anzer[10] has made progress in the realm of soccer analytics, and notably implemented a Graph Neural Network to automatically detect tactical patterns with semi-supervised learning. As shown in Fig. 2.2, their work relies on exact positions of players on the field as their input, which means they either need high-resolution game recordings or a player tracker that sends real-time location to the data hub. The findings highlight a dependence on the exact positions of players on the soccer field as crucial input data. This reliance implies a necessity for either high-resolution game recordings or a real-time player tracking system that can relay precise location information to a central data hub. This technological requirement underscores a potential area for future investigation within the domain of soccer analytics. Subsequent research endeavors could explore the feasibility and impact of integrating sophisticated player tracking technologies into predictive models, potentially incorporating advanced deep learning methods such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and other methodologies. Furthermore, as datasets continue to grow in size and scope over time, the cumulative effect of accumulating substantial data can contribute to the refinement and sophistication of predictive models. The interplay of advanced technologies and the utilization of extensive datasets is promising for the enhancement of the precision and depth of understanding in match predictions. Future researchers may find value in exploring these avenues to further advance the sophistication of analytics in the field[10].

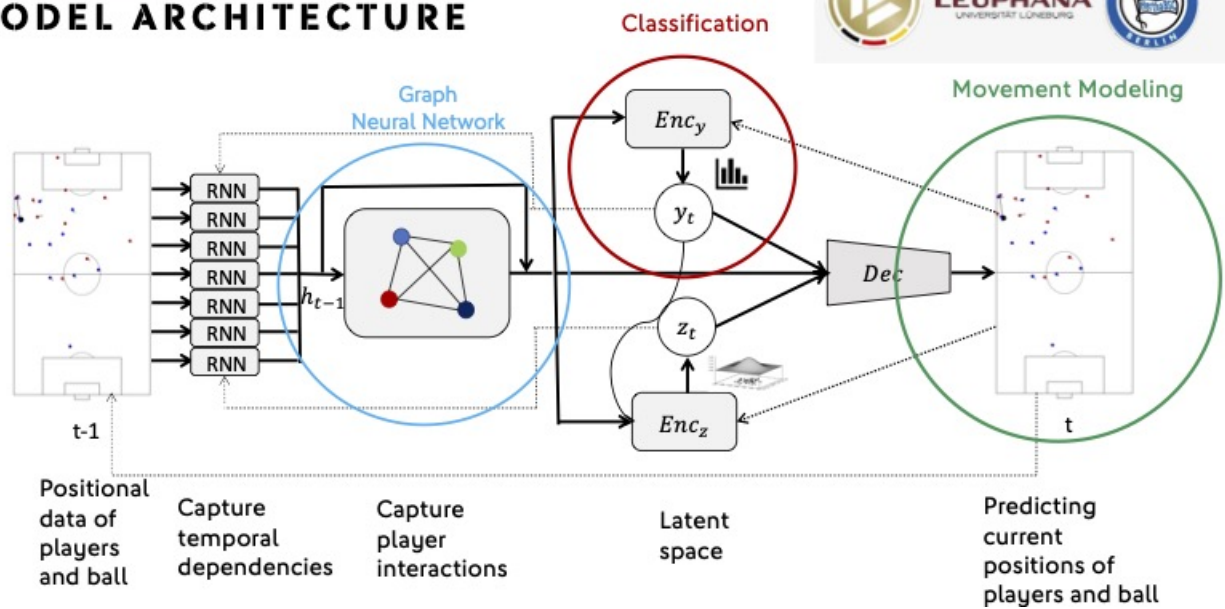


Figure 2.2: Each player’s trajectory is first processed by a recurrent neural network (RNN) with shared parameters to capture sequential patterns. The results are then refined by a graph neural network (GNN) to account for the relational dynamics among players (Figures taken from Bauer and Anzer[10].)

CHAPTER 3

Data

3.1 Data Collection

The data utilized in this study was sourced from fbref.com and encompassed statistical information from the five major European football leagues. This dataset covers a span of five seasons, specifically from the 2017-2018 season to the 2021-2022 season. To compile this dataset, information was extracted from a total of 20 league tables. Notably, the dataset comprises detailed information on 380 (or 306 for Bundesliga) matches for each league per year. This comprehensive dataset serves as a fundamental component of the research conducted in this thesis. We employed BeautifulSoup (bs4) to automatically extract links from each league page, which allows us to access detailed information about the games. The data contains 8975 observations, with each game being the unit of each row, and 234 initial columns, with different game features(eg. successful tackles, possession rate) representing each column. Each game feature contains the summation of individual performances in each team and among the 234 columns, the first half(117 columns) represents the home team and rest represents the away team. The output variable is the single variable goal difference in each team that we are interested in predicting.

3.2 Objective

We have two main objectives, namely:

1. Which (team) feature plays the most important role in goal difference?

2. Can we predict the goal difference of future games? Both questions can be considered as a supervised learning problem since we are using historical data on teams. Supervised learning is a type of machine learning where the model learns to predict labeled output data from input data.

Idea for objective 1: This is a classic regression problem with an emphasis on the interpretation of the coefficients, which can give us information on each feature's impact on the output variable.

Idea for objective 2: we can calculate a weighted average of team performance over the past N games to substitute for each row of the model matrix and predict future results. We can also construct some graph structure to represent team's relationship and game performances, and consequently convert the prediction problem into a node classification task.

3.3 Organizing Data

We organize raw data into tabular format, with rows being each game and columns being the team features. To investigate the questions of interest, we need two different matrices. Both matrices should have the same dimensions. The first matrix we need is the team-performance matrix, with rows indicating each game, and columns that include all team features. The columns can also be separated by home and away team features. The second matrix we need is the past-team-performance matrix, with each row replaced by the weighted average of the past n performances from the two teams. N can be a hyperparameter to tune later on.

3.4 Exploratory Data Analysis

We investigate the proportion of our output variable, the goal difference, which contains possible values from -9 to $+9$ excluding -8 . If we look at the value counts in Table 3.1, it is easy to notice that there are few cases of $+9$, -9 , and other classes of large absolute values, making it an extremely imbalanced classification problem. We have several ways of dealing

-9	-8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	9
1	0	3	4	44	134	351	762	1481	2229	1815	1072	550	236	99	30	6	4	2

Table 3.1: The first row indicates the outcome classes, and the second row are the counts of observations that falls into each category. Outcome variable comes with skewed class proportion, with more games end up around with low absolute goal differences, and only 7 games out of 8975 games have results greater than 7 goal differences.

with such problems, including removing outliers, generating data, or using a weighted loss function.

Examination of specific scatterplots depicting the relationship between the input feature and the output variable of home goals reveals a conspicuous trend. From Fig 3.1, we can notice that teams winning straight(as opposed to curved) corner kicks tend to exhibit lower values in terms of goals. This observation aligns with the judgement that seasoned players possess the proficiency to manipulate the trajectory of the ball, curving it either inward or outward. Such nuanced ball control strategies are evidently superior to the straightforward execution of corner kicks, which, in turn, diminishes the likelihood of scoring. This delineation underscores a noteworthy correlation, indicative of the team’s proficiency level, whereas a propensity for straight corner kicks correlates with suboptimal goal-scoring performance, thereby shedding light on specific facets of the game. The analytical framework employed herein substantiates the premise that certain intricacies of the game can be elucidated through systematic analysis, thereby contributing to a more nuanced understanding of soccer dynamics.

3.4.1 PCA

Principal Component Analysis utilizes some of the basic facts in linear algebra to apply data compression and in statistics are often used to extract main features or linear combination of features. Here, we use PCA to observe the explained variance before we start the modeling. Explained variance indicates the proportion of the dataset’s total variance that is captured

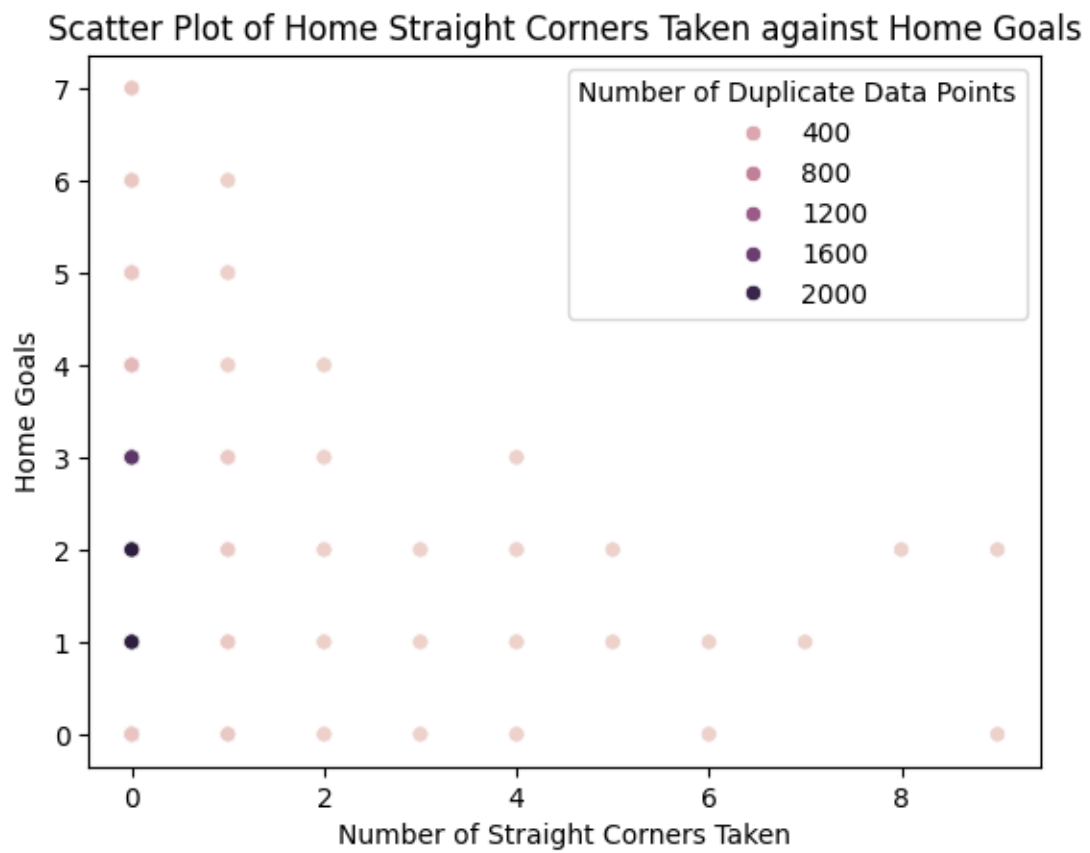


Figure 3.1: Scatterplot of Straight Corners Against Goal Differences

by each principal component. It helps in understanding how much information or variability in the data is retained by the principal components. By examining the explained variance, we can determine the number of principal components needed to adequately represent the data, ensuring that the essential patterns and structures are preserved while reducing dimensions. This step is crucial for identifying the most informative features and improving the efficiency and performance of subsequent modeling efforts. Assume the data consists of points $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\} \in \mathbb{R}^n$. We need some encoding function \mathbf{f} and decoding function \mathbf{g} , where $\mathbf{c}^{(i)} = \mathbf{f}(\mathbf{x}^{(i)})$ lies in the lower-dimensional space \mathbb{R}^l . For basic PCA, we use the simplest matrix product for the decoder $\mathbf{g}(\mathbf{c}^{(i)}) = \mathbf{D}\mathbf{c}^{(i)}$ with $\mathbf{D} \in \mathbb{R}^{m \times l}$. The goal is to minimize the reconstruction loss:

$$\mathbf{c}^* = \arg \min_{\mathbf{c} \in \mathbb{R}^l} \|\mathbf{x} - \mathbf{g}(\mathbf{c})\|_2^2$$

After solving the above, we get $\mathbf{c} = \mathbf{D}^T \mathbf{x}$, then reconstruction loss becomes

$$\arg \min_D \|\mathbf{X} - \mathbf{D}\mathbf{D}^T \mathbf{X}\|_F$$

where $\mathbf{X} = [\mathbf{x}^{(1)} \mathbf{x}^{(2)} \dots \mathbf{x}^{(m)}]^T$. This optimization problem is solved through eigendecomposition with the optimal D obtained through vertically combining the n eigenvectors of $\mathbf{X}^T \mathbf{X}$ that corresponds to the n largest eigenvalue.

In kernel PCA, we replace \mathbf{X} with function $\phi(\mathbf{X})$ that is able to operate in high dimensions, without ever computing in that space, but rather the dot product in a low dimension. If we use PCA with different kernel functions, we get the following logistic regression accuracy in Table 3.2.

Despite the application of diverse kernel transformations in the Principal Component Analysis (PCA) process, the restricted use of 5 principal components (5 is chosen for interpretability, and can be altered in future experiments) for computational efficiency yields unsatisfactory results, even when employing the best-performing model, which, as will be detailed subsequently, is XGBoost. This indicates that the limited dimensionality represented by 5 principal components falls short in adequately encapsulating the richness of information within the data distribution. In the intricate domain of soccer, where numerous features equally impact the game, each action contributing uniquely to the outcome, it

Kernel Function	Accuracy
linear	0.26096866
poly	0.25811966
rbf	0.25754986
sigmoid	0.26267806
cosine	0.25868945

Table 3.2: XGBoost Classifier Accuracy after kernel PCA

becomes apparent that a higher-dimensional representation is imperative for a more comprehensive understanding of the underlying dynamics. The inadequacy of 5 principal components underscores the necessity for a more expansive representation, emphasizing the need for additional dimensions to capture the nuanced relationships within the dataset. However, increasing the number of principal components will not be computationally efficient due to exponentially increased time of exhaustively searching through the possible feature space, which would significantly hinder the efficiency of the analysis.

CHAPTER 4

Methodology

The outcome variable, goal difference, is a discrete variable. Thus, the task at hand turns into a classification problem. Upon obtaining our training data, the initial imperative lies in addressing the challenge posed by imbalanced classification. Traditionally, there exist two primary approaches to tackle class imbalance: the first involves employing data augmentation techniques, while the second entails utilizing loss functions with re-weighting terms. While some resampling techniques offer advantages, it is essential to acknowledge their inherent limitations (there is no free lunch). The most basic approach to over-sampling involves duplicating random records from the minority class, potentially leading to overfitting. Conversely, in under-sampling, the simplest technique entails removing random records from the majority class, a practice that may result in a loss of valuable information. Here, we choose to apply several data augmentation techniques, and later compare the results on modeling testing data with the method without augmentation but with a re-weighting loss function.

4.1 Data Reduction

4.1.1 Tomek links

We start by implementing Tomek links[9], a simple undersampling technique, referring to pairs of instances that are in close proximity but belong to opposing classes. By eliminating instances from the majority class within each pair, the separation between the two classes is enhanced, thereby streamlining the classification process. A Tomek link is identified when two samples emerge as the nearest neighbors of one another. In leveraging this method for

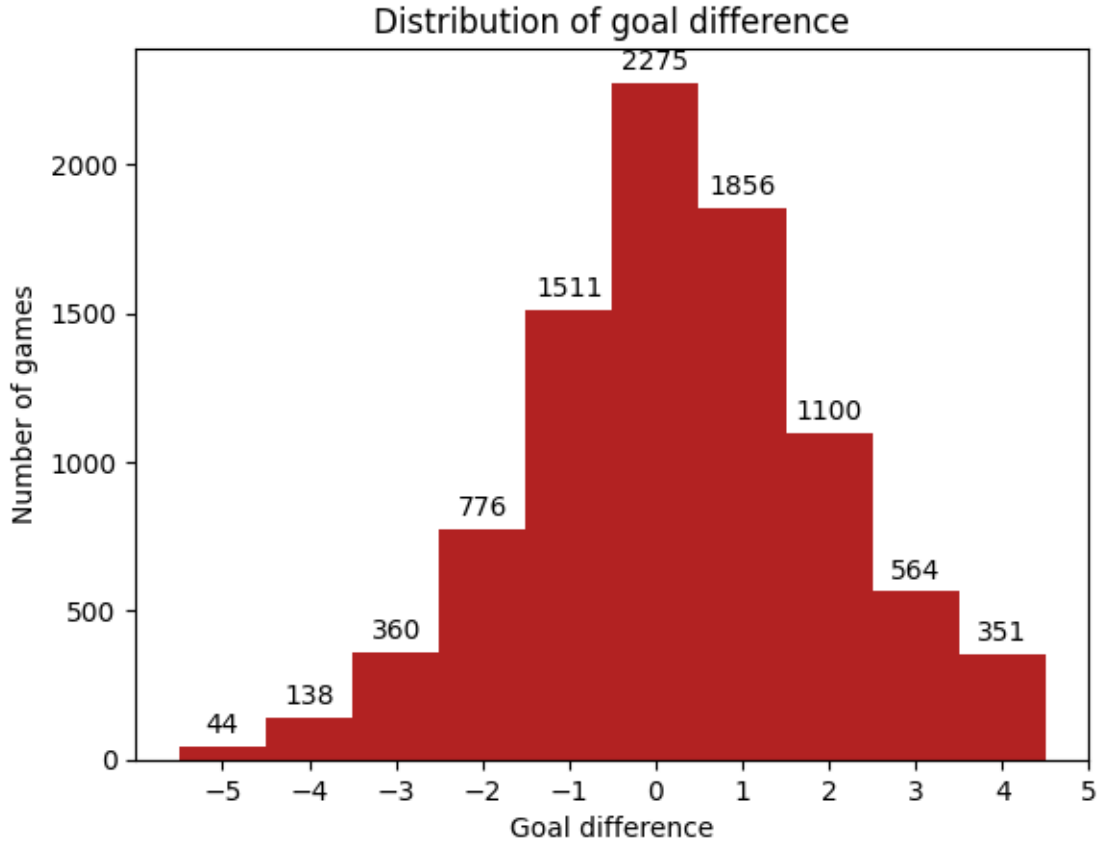


Figure 4.1: Histogram of Goal Differences

undersampling, the objective is to strategically reduce the dominance of the majority class, fostering a more balanced and discriminative dataset for classification tasks.

4.2 Data Augmentation

To apply some of the techniques below, there is a minimum number of points required (theoretically we can do augmentation even with only 2 points for a given class, but with such a lower number of observations, newly generated points would all be near the line connecting the original two points). For our dataset, which contains some classes with numbers of observations less than 10, it is reasonable to remove these extreme outliers in our classification modeling and left us with the outcome distribution shown in Fig. 4.1.

4.2.1 SMOTE and its variations (borderline, ADASYN)

The Synthetic Minority Over-sampling Technique (SMOTE) is a technique to generate synthetic tabular data. All variations of this oversampling technique create points by using a convex combination of points that belong to the selected minority classes, and by repeatedly applying this method to expand all classes to the same number of counts, which, in our case, is 2229.

4.2.2 Tabular GAN

For tabular Generative Adversarial Network (GAN)[12] and Conditional GAN[11], both methods emulate the distribution of the training matrix and outcome vector. Again, we are attempting to generate data using TGAN so that each outcome class has an equal count. Therefore, we split the data one by one for each class and generated the required number for each class.

4.3 Feature Selection

In our case, feature selection is used as a preprocessing step in conjunction with machine learning models for regression/classification purposes. The purpose is to find the best feature subset with a determined number of features. Feature selection methods can be classified into wrapper, embedded, filter methods, or the hybrid of the previous three methods. Wrapper methods treat the classifier as a black box and select features based on the classification result using some metric (e.g. forward, backward selection); embedded methods carry out selection through the classifier algorithm itself (e.g. regularization, tree-based algorithms); filter methods have the selection process independent of classifier algorithms (e.g. univariate ANOVA). [4] Typically, filter methods are the fastest, and wrapper the slowest.

4.3.1 Exhaustive Search

In general, exhaustive search over feature space is often impractical due to combinatorial explosion. The feature space refers to all the subset of the features, and the number of such combinations grows exponentially with the number of features. For example, with 20 features, there are 2^{20} possible subsets to consider. The computational cost increases exponentially with each additional feature, making an exhaustive search over the feature space quickly unmanageable. Some feasible approaches to get around is to assume a fixed number of features to select, denote it as k . We can implement a greedy search of $\binom{n}{k}$ number of possible subsets of size k to find the optimal subset based on the metrics we choose.

The combinatorial explosion is a fundamental challenge in feature selection, and it highlights the need for efficient and heuristic methods to identify relevant features without exploring the entire feature space. Techniques such as greedy algorithms, recursive feature elimination, or optimization methods are often employed to strike a balance between computational complexity and obtaining a good subset of features for predictive modeling. In this case, to implement a simple exhaustive search, we set the subset size fixed as $k = 5$ for the sake of computational efficiency.

4.3.2 Lasso

Lasso[13], or L1 regression, is a common regularization technique based on linear regression. It adds a weighted L1 norm of the coefficient vector to the least squares function. It is not hard to notice, once the weight was set large enough, some elements in the coefficient vector will diminish to zero, and thus performing feature selection while doing regression (embedded feature selection). However, the downside of using lasso is the its lack of stability and its natural assumption of a linear relationship, which stems from its fundamental nature of being a linear regression.

4.3.3 Mutual Information

Mutual information (MI)[6] is a measure of the amount of information that one random variable has about another variable, and in our case, it is used to quantify and compare the relevance of different feature subsets concerning the output variable. It is defined to be

$$I(X; Y) = D_{\text{KL}}(P_{(X,Y)} \parallel P_X \otimes P_Y) \quad (4.1)$$

where (X, Y) are pairs from the space $\mathcal{X} \times \mathcal{Y}$, $P_{(X,Y)}$ is the joint probability mass function, and \otimes indicates outer product, and for an individual point \mathbf{x} , simply the scalar product of the marginal distribution $P_X(x) \times P_Y(y)$.

The Kullback-Leibler (KL) divergence measures the difference between two probability distributions. For discrete probability distributions P and Q , the KL divergence ($D_{\text{KL}}(P \parallel Q)$) is defined as:

$$D_{\text{KL}}(P \parallel Q) = \sum_{\mathcal{X}} \sum_{\mathcal{Y}} P(x) \log \left(\frac{P(x)}{Q(x)} \right) dx dy \quad (4.2)$$

and for continuous variable, integrals instead of summations.

Thus, the mutual information can be written explicitly:

$$I(X; Y) = \int_{\mathcal{X}} \int_{\mathcal{Y}} P_{(X,Y)}(x, y) \log \left(\frac{P_{(X,Y)}(x, y)}{P_X(x)P_Y(y)} \right) dx dy \quad (4.3)$$

where $P_{(X,Y)}(x, y)$ is the joint probability distribution of X and Y , and $P_X(x)$ and $P_Y(y)$ are the marginal distributions.

Since all of our dataset's features are continuous, we need nonparametric methods based on entropy estimation from k-nearest neighbors distances as described in [8]. The result is implemented by the `mutual_info_classif` function in `scikitlearn.feature_selection`. We apply Mutual Information between each feature and output variable to perform this filter method.

In order to empirically elucidate the significance of mutual information among input features concerning home and away goals, our study adopts home goals as the output variable. As delineated in Figure 4.2, while the mutual information (MI) values across features do not

exhibit conspicuous prominence, an analysis reveals that pivotal features, notably attacking maneuvers from the home side encompassing key passes, touches in the penalty area, and total shots, alongside errors committed by opponents and penalties awarded to the away team, exert a discernible influence on home goal outcomes. The observed results align with our a priori expectations, thereby reinforcing our intuitive understanding. Furthermore, our findings establish a robust statistical foundation elucidating the salience of these factors in influencing home goals. This not only enhances our comprehension but also furnishes statistical validation for the pivotal role played by these factors. Subsequently, scouts or analysts may leverage these metrics to make informed decisions regarding the identification of key players on the field or determining the corresponding player's salary based on their contribution to the team's results from a statistical standpoint.

In the context of implementing feature selection using mutual information, the approach to selecting features, whether sequentially adding them or considering mutual information with respect to the target variable Y in one go, can yield significant differences in outcomes. Whether or not to augment the data before feature selection can also lead to a different result, as shown in Fig. 4.2.

When features are sequentially added, similar to the adapted method shown in greedy search in Sec. 4.3.1, the selection process involves iteratively adding features based on their individual mutual information with Y . This method allows for a step-by-step refinement of feature selection, where each added feature is chosen based on its incremental contribution to mutual information. On the other hand, considering mutual information with Y in one go involves evaluating the mutual information for all features collectively. This holistic approach considers the combined information provided by the entire set of features in relation to the target variable. To illustrate this process, we generate a barplot where each bar represents the information gain contributed by an added feature. This step-by-step addition of features results in a non-decreasing barplot, showcasing the cumulative improvement in information gain as we progressively incorporate more features into the model. This approach provides a clear visual representation of the utility of each feature, allowing us to discern trends and

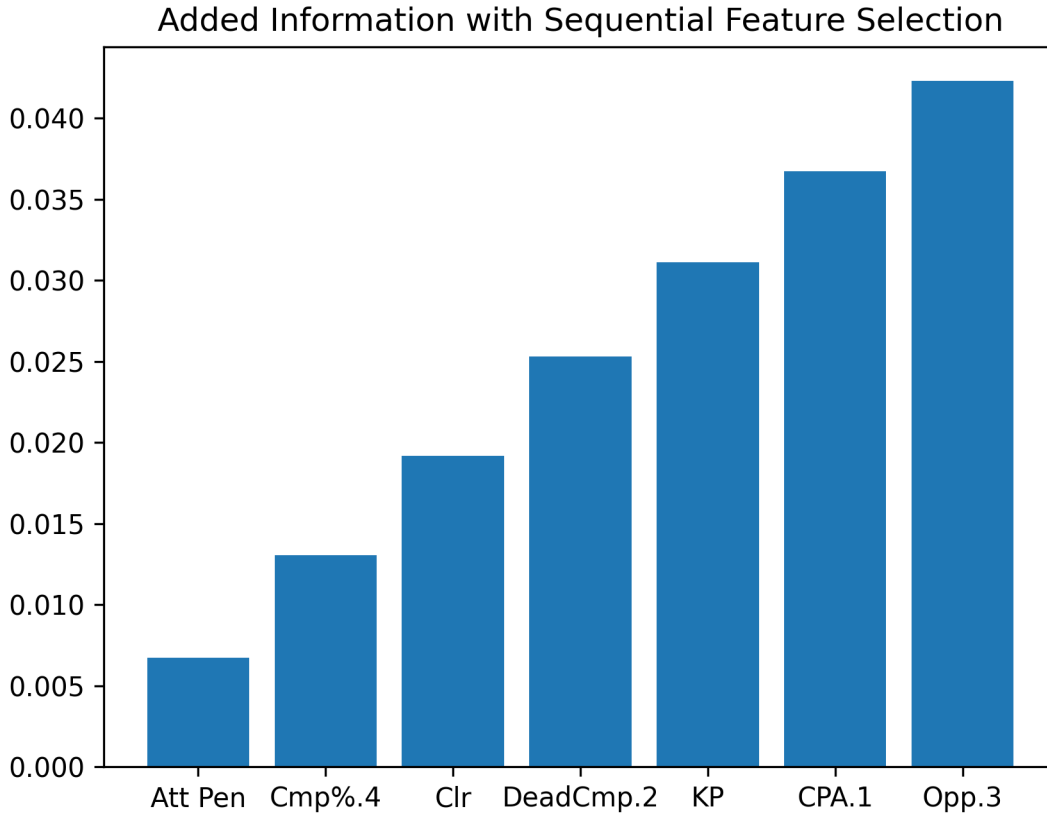


Figure 4.2: Bar Plot of MI gain of added Input Features shows the top 7 features added sequentially are shown above.

identify the point at which additional features contribute less substantially.

As we observe the barplot, we eventually reach a saturation point, where the incremental information gain diminishes, indicating that further addition of features does not significantly improve predictive performance. At this juncture, we strategically decide on a threshold, beyond which we cease adding features. This threshold is a crucial determinant in streamlining the feature set for optimal performance without unnecessary complexity. The chosen metrics for evaluating the effectiveness of this feature selection process is the direct result from XGBoost (XGB) modeling over possible feature set.

The choice between the above two strategies can impact the set of selected features and,

consequently, the model's performance. A sequential addition may capture nuanced relationships incrementally, while a simultaneous evaluation may reveal synergies or dependencies among features that contribute collectively to mutual information with Y .

4.4 Models

4.4.1 Linear Regression

Linear regression is the fundamental statistical method used to model the linear relationship between an outcome variable with independent variables. Geometrically speaking, the model helps us find the best fitted line that explains the data points in a way that minimizes the distance between the predicted and the actual data points. The equation of the line is commonly expressed as $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where \mathbf{y} is the outcome variable, \mathbf{X} is the model matrix, β is the slope of the line, and ϵ denotes the random noise.

4.4.2 XGBoost

Extreme Gradient Boosting, or XGBoost, is the state-of-the-art ensemble tree algorithm. It utilizes the concept where weak learners (in this case usually simple decision trees) are used sequentially to construct a strong learner. XGBoost enhances this process by incorporating ridge and lasso regularization techniques, enables parallel and distributed computing, and combines Newton's method and gradient descent to accelerate optimization speed. This method is known for its exceptional speed and ability to handle large datasets on tabular data for both regression and classification tasks.

Softmax Objective Function:

$$\text{Objective} = - \sum_{i=1}^N \sum_{k=1}^K [y_{ik} \log(p_{ik}) + (1 - y_{ik}) \log(1 - p_{ik})] + \sum_{k=1}^K \Omega(f_k) \quad (4.4)$$

4.4.3 MLP

Multilayer Perceptron (MLP), involves stacking layers of linear regression with activation function in order for it to learn patterns from data. While its basic structure is straightforward, the model’s simplicity conceals its capability to capture complex relationships in diverse datasets, making it a foundational and accessible entry point into deep learning.

4.4.4 Tabular Transformer

TabTransformer (Huang et al. 2020), utilizes self-attention based Transformers to enhance predictive accuracy on tabular data. The architecture includes a column embedding layer, a series of Transformer layers transforming categorical features into contextual embeddings, and a multi-layer perceptron (MLP). The concatenated contextual embeddings and continuous features are fed into the MLP, and the model is trained end-to-end by minimizing the loss function to learn all parameters.

4.4.5 Graph Convolution Network

we formulate the problem as a node classification task employing Graph Convolutional Networks (GCNs). Each node, denoted as h_i , is associated with the label of Goal Differences and contains relevant feature information. The graph is constructed based on recent games involving teams and the current opponent of the ongoing match, where edges connect nodes representing these entities. The Convolutional Layer at the $l + 1$ iteration updates the representation of each node using the following formula:

$$h_i^{(l+1)} = \sigma\left(\sum_{j \in N_i} \frac{1}{c_{ij}} W^{(l)} h_j^{(l)}\right)$$

Here, \mathbf{W} is the weight matrix of dimensions $\in \mathbf{R}^{100 \times 7}$, h_i represents the feature vector of node $\in \mathbf{R}^7$, c_{ij} is a normalizing constant, and N_i refers to the neighboring nodes of node i . The convolutional operation is designed to capture and propagate information through the graph structure. The diagram on the right illustrates the architecture of the Graph Convolutional

Network, emphasizing the connections between nodes and their respective features. This approach leverages the relational information encoded in the graph to enhance the predictive capabilities of the model for node classification, specifically in the context of opponent team strength and recent performance in soccer matches.

CHAPTER 5

Result

In our comprehensive experimentation, we investigate various combinations of data augmentation techniques, feature selection methods, machine learning algorithms, and loss functions to evaluate their impact on predictive modeling. The primary metric for comparison is the accuracy score, and our objective is to discern optimal configurations for improved performance.

5.1 Features Selected

Building on the theoretical foundations outlined above, our experiment delves into the practical application of feature selection methodologies, specifically focusing on the empirical results obtained through exhaustive search, Lasso, and Mutual Information (MI). In line with our approach, we set a fixed subset size to maintain computational efficiency, leveraging the insights gained from the combinatorial challenge associated with the feature space. Our analysis of goal differences as the output variable reveals that, among numerous features, seven key features emerge as empirically selected. These features, identified through a thoughtful combination of theoretical underpinnings and practical considerations, represent a subset that effectively captures the essential predictors influencing goal differences.

The observed discrepancy in the features selected by Mutual Information (MI) compared to Lasso can be attributed to their distinct operational mechanisms. While MI operates in a sequential manner, sequentially adding features based on their incremental contribution to mutual information with the output variable, Lasso introduces a regularization term in the linear regression framework, promoting sparsity in the coefficient vector. The sequential

Feature id	Meaning
Att Pen	Touches in Attacking Penalty Area
Cmp%.4	Forward Pass Completion Rate
Clr	Clearance
DeadCmp.2	Dead Ball Completed
KP	Key Passes
CPA.1	Carries into Penalty Area
Opp.3	Attempted Crosses

Table 5.1: Glossary of Feature Selected by added MI

Feature id	Meaning
Att Pen	Touches in Attacking Penalty Area
Cmp%.4	Forward Pass Completion Rate
Clr	Clearance
DeadCmp.2	Dead Ball Completed
KP	Key Passes
Blocks.2	Number of times standing in the ball path
TklW.2	Tackles Won

Table 5.2: Glossary of Additional Features Selected by Lasso

nature of MI allows it to capture nuanced relationships and dependencies among features, revealing their collective impact on the output variable. In contrast, Lasso's regularization penalty tends to prioritize a subset of features that collectively optimize predictive performance. The choice between these approaches hinges on the data's underlying structure and the desired emphasis on individual versus collective feature importance. Consequently, the features selected by MI, through its sequential and cumulative approach, may exhibit superior performance in scenarios where the synergies among features are pivotal for predictive modeling precision.

The feature contributions to the model output reveal insightful patterns in the data. Among the attributes considered, "Att Pen_away" emerges as the most influential factor, contributing positively to the model's prediction. This suggests that the number of attempted penalties by the away team significantly impacts the outcome, potentially indicating a higher likelihood of scoring or affecting the game's dynamics.

Utilizing a force plot generated by SHAP (SHapley Additive exPlanations) offers a compelling visual representation of the underlying predictive process. This visualization elucidates the individual contributions of each feature towards the model's output, providing a comprehensive understanding of the factors shaping the prediction. The force plot shown in Fig. 5.1 succinctly displays the direction and magnitude of impact for each feature, facilitating the identification of key drivers influencing goal difference positively or negatively. Through this graphical depiction, it is easy to discern the most influential variables driving the model's predictions, thus enhancing the interpretability and validation of the model's performance.

However, with either method, the selected features align with our expectations based on intuition and basic understanding of the game, demonstrating the robustness of our feature selection process and its potential for enhancing predictive modeling precision. This empirical validation reinforces the significance of striking a balance between computational complexity and the attainment of a refined feature subset for optimal model performance.

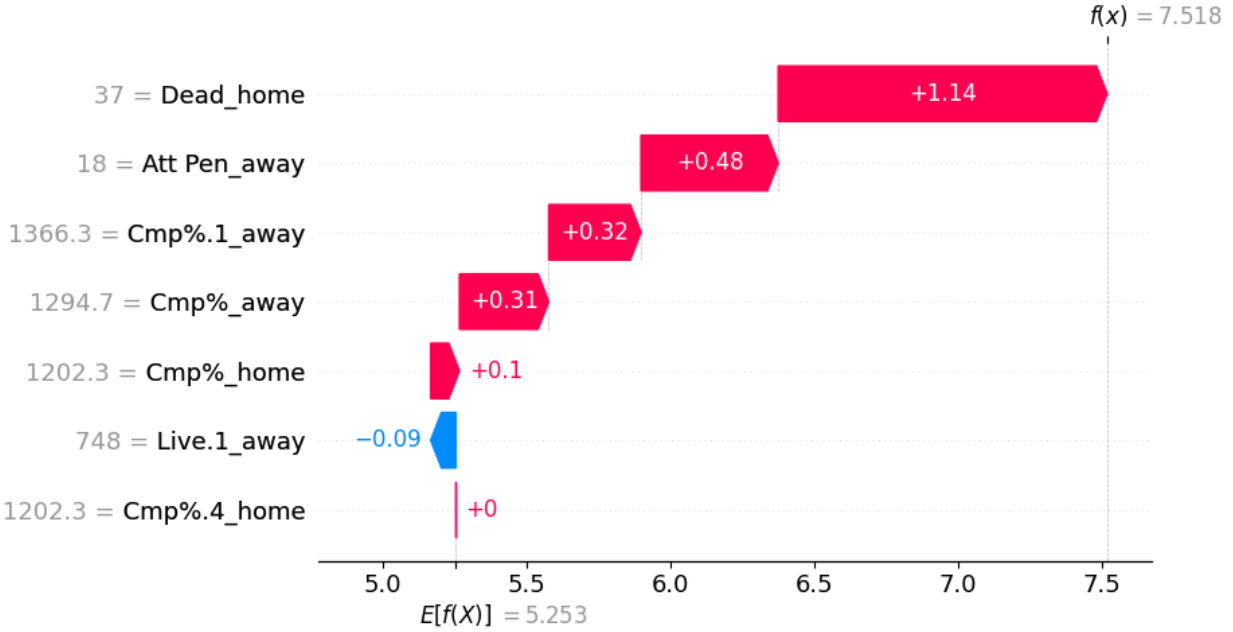


Figure 5.1: Force Plot of the Features' Shapley Value

5.2 Experiments

A crucial consideration in our experimental design is the order of applying data augmentation and feature selection. We recognize that if feature selection precedes data augmentation, there is a risk of introducing bias in the generated data. Augmenting the data before feature selection, however, allows us to capture more information, potentially enhancing the richness of the dataset and leading to more informed model selection. Table 5.3 encapsulates key performance metrics for each experimental configuration. The benchmark method we are comparing to is the betting agency's odd being converted to the accuracy metrics. Overall, betting agencies use a combination of data analysis, statistical modeling, expert input, market dynamics, and risk management techniques to determine prediction odds that reflect the probabilities of each outcomes of goal differences.

The combination of TabGAN, Mutual Information, and Graph Neural Network approaches yielded the highest accuracy in the experiment. However, it's essential to note that our benchmark, the betting agency, still outperforms all our combination models. Since

Step1	Step2	Method	Accuaracy
/	MI	Logistic Regression	0.2380
SMOTE	MI	Logistic Regression	0.2093
SMOTE	MI	KNN	0.1567
SMOTE	MI	XGB	0.2450
MI	SMOTE	XGB	0.2344
MI	TGAN	XGB	0.2621
TGAN	MI	XGB	0.2767
TGAN	MI	GNN	0.2890
		Betting Agency	0.2977

Table 5.3: Model Result with Combination of Methods

our results are not better, it’s challenging to identify a significant edge where our model distinctly surpasses the benchmark. Despite this, our experimentation indicates areas for potential enhancement, such as expanding the feature selection space or exploring ensemble model techniques, to achieve more competitive outcomes.

CHAPTER 6

Future Work

The study strategically fused cutting-edge artificial intelligence and machine learning techniques with the intricacies of soccer, focusing on the prediction of match outcomes. As a potential direction for future researchers, the incorporation of real-time data streams during matches is suggested, enabling dynamic model adjustments based on evolving game scenarios. With technological advancements, the increased computational power afforded by robust GPUs enables the exploration of a myriad of model architectures. This empowerment allows researchers to investigate not only the depth and complexity of CNNs but also the sequential and temporal dependencies captured by RNNs and the diverse capabilities of other advanced models like Transformer networks. Furthermore, the synergy between powerful GPUs and the surge in the volume and precision of data collection across multiple avenues is pivotal. The amalgamation of meticulously curated datasets, spanning player movements, team statistics, and contextual factors, establishes a foundation for model training characterized by richness and diversity. This wealth of high-quality data not only reinforces the resilience of the models but also facilitates the exploration of intricate relationships within the multifaceted realm of soccer. Advancements in data collection methodologies, including real-time player tracking systems and high-resolution game recordings, further elevate the potential for refining deep learning models. The precision in data acquisition empowers researchers to capture detailed aspects of player dynamics, fostering a more comprehensive and accurate representation of the intricate dynamics inherent in soccer matches. In the contemporary landscape of burgeoning data, there is an increasingly compelling need for more extensive Tracking Data Analysis, combining high-resolution videos with advanced computer vision techniques. The wealth of large-scale datasets available today presents a unique op-

portunity to delve deeper into the intricacies of player movements and object tracking. As sports analytics continues to evolve, the demand for richer, more nuanced data becomes paramount. The integration of high-resolution videos and sophisticated computer vision not only refines our understanding of player dynamics but also lays the foundation for more accurate and insightful analyses. This emphasis on leveraging large datasets and advanced tracking technologies underscores a pivotal direction for future research endeavors, promising enhanced precision and a deeper comprehension of the multifaceted aspects of sports events.

Another area ripe for improvement involves the exploration of ensemble models, combining predictions from diverse algorithms to enhance overall accuracy. To enhance interpretability, methodologies like SHAP values are recommended to illuminate the significance of different features in decision-making. The potential integration of human-in-the-loop approaches, involving expert insights in model training, presents an avenue for refinement. Transfer learning, the creation of user-friendly interfaces for stakeholders, and ethical considerations in predictive model applications are also critical areas for future exploration. By delving into these directions, researchers can contribute to the ongoing evolution of goal-oriented forecasting in soccer analytics.

BIBLIOGRAPHY

- [1] Assessing the performance of premier league goalscorers, 2012. <https://www.statsperform.com/resource/assessing-the-performance-of-premier-league-goalscorers>.
- [2] G. Anzer, P. Bauer, U. Brefeld, and D. Faßmeyer. Detection of tactical patterns using semi-supervised graph neural networks. 03 2022.
- [3] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [4] A. Jović, K. Brkić, and N. Bogunović. A review of feature selection methods with applications. In *2015 38th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*, pages 1200–1205, 2015.
- [5] L. F. Kozachenko and N. N. Leonenko. Sample estimate of the entropy of a random vector”. *Probl. Peredachi Inf*, 23(2):95–101, 1987.
- [6] A. Kraskov, H. Stögbauer, and P. Grassberger. Estimating mutual information. *Phys. Rev. E*, 69:066138, Jun 2004.
- [7] L. Maurath. Xg model - accuracy and goodness-of-fit. <https://www.thesignificantgame.com/portfolio/xg-model-accuracy-and-goodness-of-fit/>.
- [8] B. Ross. Mutual information between discrete and continuous data sets. *PloS one*, 9:e87357, 02 2014.
- [9] I. Tomek. An experiment with the edited nearest-neighbor rule. In *IEEE Transactions on Systems, Man, and Cybernetics*, 1976.
- [10] P. Xenopoulos and C. Silva. Graph neural networks to predict sports outcomes. In *2021 IEEE International Conference on Big Data (Big Data)*, pages 1757–1763, 2021.

- [11] L. Xu, M. Skoularidou, A. Cuesta-Infante, and K. Veeramachaneni. Modeling tabular data using conditional GAN. *CoRR*, abs/1907.00503, 2019.
- [12] L. Xu and K. Veeramachaneni. Synthesizing tabular data using generative adversarial networks, 2018.
- [13] H. Zou and T. Hastie. Regularization and Variable Selection Via the Elastic Net. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 67(2):301–320, 03 2005.