

Unleashing Potential: A Data-Driven Exploration of Identifying Player Potentialities through Advanced Analytics in Sports

by

Prashanta Bhowmik

21101343

Md. Khaliful Islam

17301114

Nabil Shartaj Khan

20101025

Ananna Acharjee

20101294

A thesis submitted to the Department of Computer Science and Engineering
in partial fulfillment of the requirements for the degree of
B.Sc. in Computer Science

Department of Computer Science and Engineering
School of Data and Sciences
Brac University
January 2024

© 2024. Brac University
All rights reserved.

Declaration

It is hereby declared that

1. The thesis submitted is my/our own original work while completing degree at Brac University.
2. The thesis does not contain material previously published or written by a third party, except where this is appropriately cited through full and accurate referencing.
3. The thesis does not contain material which has been accepted, or submitted, for any other degree or diploma at a university or other institution.
4. We have acknowledged all main sources of help.

Student's Full Name & Signature:

Prashanta Bhowmik

21101343

Md. Khaliful Islam

17301114

Nabil Shartaj Khan

20101025

Ananna Acharjee

20101294

Approval

The thesis titled “Unleashing Potential: A Data-Driven Exploration of Identifying Player Potentialities through Advanced Analytics in Sports ” submitted by

1. Prashanta Bhowmik (21101343)
2. Md. Khaliful Islam (17301114)
3. Nabil Shartaj Khan (20101025)
4. Ananna Acharjee (20101294)

of Fall 2023 has been accepted as satisfactory in partial fulfillment of the requirement for the degree of B.Sc. in Computer Science on January 22, 2024.

Examining Committee:

Supervisor:
(Member)



Nabuat Zaman Nahim
Lecturer
Department of Computer Science and Engineering
BRAC University

Program Coordinator:
(Member)

Md. Golam Rabiul Alam, PhD

Professor
Department of Computer Science and Engineering
Brac University

Head of Department:
(Chair)

Sadia Hamid Kazi, PhD

Chairperson and Associate Professor
Department of Computer Science and Engineering
Brac University

Abstract

In this project, we delve deeper into the complex area of predicting a potential replacement of a footballer of a specific position. For that, we used multiple machine learning models on the Sofifa dataset. Our analysis reveals interesting insights into the predictive capabilities of these models, with a particular focus on numerical performance measures. Among the models tested, the LightGBM Regressor appears to be the epitome of predictive power. This algorithm consistently outperforms the others, showing the lowest mean squared error (MSE) and highest R-squared value on both the test data set and the overall data set. Her ability to navigate the complexity of player performance patterns is evident, making her a leader in our prediction arsenal. Complements for Random Forest Regressor, XGBRegressor Regressor, LightGBM Regressor, and CatBoost Regressor demonstrate superior performance, characterized by consistently low MSE values and high R-squared values . These gradient boosting algorithms demonstrate their effectiveness in capturing complex patterns in the Sofia dataset. The Linear Regressor model utilizes its power in understanding the linear relationships among the data and gives a higher accuracy too. The KNeighbors-Regressor, with its proximity-based approach, also achieves stripes, especially by achieving high R-squared values. This model excels at identifying players with similar characteristics, highlighting their collective impact on overall performance. It should be noted that, Support Vector Regressor (SVR) and Neural Network models provide valuable insights, despite relatively lower prediction accuracy. These models highlight the complexity inherent in player forecasting and highlight the need for meticulous parameter tuning. LightBGM Regressor stands out as the superior model for predicting our research, closely followed by XGBRegressor, Random Forest Regressor, CatBoost Regressor. These results highlight the importance of selecting models that match the variation of the data set to accurately and reliably predict performance in soccer analytics.

Keywords: Football, Player performance, Prediction, Machine learning, Sofifa dataset, Random Forest Regressor, Linear Regressor, XGB Regressor, LightGBM Regressor, CatBoost Regressor, SVR, KNeighbors Regressor, Neural network, Sports analytics, Data science, Player attributes, In-game values.

Dedication

This dissertation is devoted to our adored parents and esteemed teachers. We owe them our thanks. We would not have gotten this far without their cooperation, concern, and support. Many thanks to them.

Acknowledgement

All glory to God, who has enabled us to finish our thesis without any significant setbacks. We appreciate the general direction provided by our advisor, Nabuat Zaman Nahim Sir. Without his help, we would be unable to complete our project.

Table of Contents

Declaration	i
Approval	ii
Abstract	iv
Dedication	v
Acknowledgment	vi
Table of Contents	vii
List of Figures	ix
List of Tables	x
Nomenclature	x
1 Introduction	1
1.1 Research Problem	3
1.2 Research Objectives	4
1.3 Thesis Organization	5
2 Background	6
2.1 Literature Review	6
2.2 Algorithms	9
3 Dataset	13
3.1 Complexity	14
4 Proposed Methodology	16
4.1 Data Pre-processing	16
4.2 Algorithm Description	17
5 Experimentation	22
6 Result Analysis	25
6.1 Performance Analysis	28
6.1.1 Performance Analysis on All Data and Test Data:	28

7 Conclusion and Future Work	34
7.1 Conclusion	34
7.2 Future Work	34
Bibliography	36

List of Figures

1.1	Organizational Structure of the Thesis	5
4.1	Architecture of Random Forest Regressor	17
4.2	Architecture of XGBoost Regressor	18
4.3	Architecture of LightGBM Regressor	18
4.4	Architecture of CatBoost Regressor	19
4.5	Architecture of Linear Regression	19
4.6	Architecture of K-Nearest Neighbors Regression	20
4.7	Architecture of Support Vector Regression	20
4.8	Architecture of CNN	21
5.1	Methodology	22
5.2	Comparison of the best-performing model's accuracy	24
6.1	MSE on Test Data	28
6.2	MSE on All Data	29
6.3	RMSE on Test Data	30
6.4	RMSE on All Data	31
6.5	R-squared on Test Data	32
6.6	R-squared on All Data	33

List of Tables

3.1	Count Of Data from 2015 to 2022	15
5.1	Accuracy and training time comparison among the best-performing algorithms for on dataset	23

Chapter 1

Introduction

The field of football expectations has seen important changes thanks to the integration of machine learning calculations and the widespread accessibility of comprehensive data sets. Amid the COVID-19 pandemic, top-flight football has seen increased competitiveness, requiring advanced prediction models capable of uncovering the complex challenges faced by players. This includes factors such as injuries, the advancement of team members, and energy changes in the redirection strategy. Furthermore, one important aspect that captures the attention of fans, analysts, and management is the potential replacement of a particular player in a specific position. This process isn't simple because it doesn't simply mean to substitute the player but also to balance team chemistry, tactical fit, and team planning.

The Data Set unit, carefully sourced from Sofifa, presents itself as a treasure trove of typical data on player, game values, and historical performance information. In the interest of strategies that enhance players' expectations, the company strives to manage control over various machine learning strategies, extending from conventional calculations such as Linear Regressor, Random Forest Regressor, XGB Regressor, LightGBM Regressor, CatBoost Regressor, SVR, and KNeighborsRegressor, to advanced neural organization designs, which include convolutional neural networks (CNN).

The innate challenges posed by COVID-19 have indeed made the football scene more exciting. The performances of the players are evaluated not only for their individual skills but also for their flexibility in improving their team technique, their versatility in the face of injuries, and overall their commitment to winning as a team.

In this context, the need for Oracle models capable of characterizing this complexity becomes essential. The Sofifa dataset, with its comprehensive range of player qualities, provides a comprehensive view of each player's abilities. This research aims to explore the multifaceted factors that instigate contemplation of player replacements, scrutinizing the delicate balance between individual player attributes and the collective dynamics of the team. Whether prompted by injuries or strategic acquisitions during transfer windows, the potential replacement of a player emerges as an integral part of the intricate tapestry of football, introducing layers of intrigue and speculation that enrich the essence of the sport. Traits such as speed, skill, dexterity, and team factors play a central role in determining a player's suitability on the field. Game values provide an additional layer of understanding that represents how players are perceived within the game community.

Selected machine learning calculations, including the Random-Forest Regressor, provide power and interpret-ability, while neural arrangement models delve into the complexities of sequential information and complex patterns.

Random Forest Regressor possesses outfit learning methods. This has proven its viability in capturing the complexity of players' characteristics and their interdependencies. XGB Regressor is known for its gradient-boosting system. It exceeds expectations in handling nonlinear connections and improves overall predictive performance. Furthermore, LightGBM Regressor is a gradient enhancement system created by Microsoft. It offers efficiency in preparing large datasets, which makes it especially suitable in the context of general information about Sofifa players. The CatBoost regression engine is designed to handle category highlights consistently, which contributes to the project's goal of comprehensive player development.

SVR enters this field with the ability to handle nonlinear connections and complex designs in data sets. The KNeighbors Regression Tool, a simple yet actionable calculation, provides an approach based on proximity to player expectations. Therefore, it highlights the importance of comparable players in influencing outcomes. Combining neural systems using the CNN or Convolutional Neural Network will infuse a layer of modernity into the enterprise. Neural Network is driven by the activity of the human brain. They are capable of capturing sequential conditions and complex spatial designs in data sets. The Sofifa dataset is no doubt an important asset. However, it presents a number of challenges. The huge volume of player data requires powerful information preprocessing processes to convert the raw data into an arrangement that is both practical and efficient for investigation. This includes cleaning information, handling missing values, and encoding taxonomic highlights to prepare them for the rigors of machine learning algorithms.

Similarly, the assessment of player qualities and value in the game serves as a parallel examination of the assumptions in the football player expectations project. A player's development, reflected in their values and game characteristics, is an important variable that determines their fitness on the virtual playing field. However, it is necessary to acknowledge the limits of considering such an opinion. The subjective nature of people's conclusions, the inherent bias in conducting surveys, and the challenge of ensuring survey quality all pose unique problems. As expectations of football players increase based on the FIFA dataset, consideration needs to be given to addressing these limitations and deriving meaningful insights from them. At this level, web searching becomes an emergency tool, allowing the extraction of relevant information from Sofifa for investigative purposes. Web scraping involves the mechanized extraction of information from web pages, allowing the collection of player attributes, game values, and verifiable performance information. This preparation simplifies the information security step, providing the means for meaningful analysis. Estimation testing includes classifying player attributes, game values, and routine performance information into categories such as positive, negative, or unbiased. This enthusiastic approach makes a difference in observing player satisfaction based on the Sofifa data set. Additionally, opinion surveys could be more in-depth, capturing feelings such as joy, outrage, loss of hope or disappointment, thereby contributing to a better understanding of player satisfaction. Gullible Bayes, known for its simplicity and completeness in content classification, can be used to discern sentiment from player surveys. Choice Tree's calculations exceed expectations by translating the complex connections in information, ad interpretation, and experi-

ence into compelling variables. K Nearest Neighbor, a distance-based calculation, provides a unique perspective by recognizing comparable agents and their impact on performance. The all-encompassing goal of this effort is to contribute to the understanding of soccer players' performance prediction using a different set of machine learning procedures. The Sofifa dataset, with its multidimensional player data, provides a solid basis for this investigation. The application of advanced analytics and machine learning helps deliver a beneficial experience for partners, counting coaches, and team supervisors, encouraging informed decision-making in team composition, registration, and the general development of the method. This combination of machine learning, web scraping, and hypothesis investigation in soccer player prediction offers an exciting avenue to advance our understanding of player flow in the context of modern football. The Sofifa dataset, with its richness in player qualities and game values, serves as a portal to reveal the complex designs that shape player performance. We performed various machine learning calculations, each with its own perspective. It strives to improve the accuracy and satisfaction that players expect. Our investigation into this problem establishes confirmation of the cutting-edge intersection between sports, information science, machine learning, and advertising, not as prophetic experiments but as a deeper understanding of variables that lead to victory on the football field.

1.1 Research Problem

The Sofifa data set has emerged as a pillar in the effort to push player expectation strategies to new heights. Through meticulous curation, the dataset summarizes a comprehensive view of the complex world of football players. Characteristics such as speed, skill, and shooting, as well as market value and historical performance data, contribute to the richness of the data set. Together, these elements provide a comprehensive understanding of each player's abilities, providing an important resource for unraveling the complexities of player motivation on and off the field. As the world grapples with the challenges posed by the pandemic, the competitive nature of professional football has increased, highlighting the need for advanced prediction models. In our quest to exploit the full potential of the Sofifa dataset, this project constitutes a beacon of innovation. We aim to harness the power of machine learning strategies, from conventional algorithms such as Random Forest Regressor, XGB Regressor, LightGBM Regressor, CatBoost Regressor, SVR, KNeighbors Regressor, to advanced neural architectures, including CNN. Each of these methods offers a unique perspective by combining interpret-ability and depth to navigate the varied landscape of football player predictions. It is clear that the challenges posed by the COVID-19 pandemic have added new momentum to football. A player's performances are no longer judged solely on their individual skills but also on their ability to adapt to changing team strategies through their skills and potentiality, their ability to recover from injury, and their overall contribution to the team's success. In this context, it is imperative to have prediction models capable of capturing this complexity. The Sofifa dataset, with its comprehensive scope, provides a comprehensive view of player qualities, making it an essential tool in the pursuit of improved player prediction strategies. However, the Sofifa dataset, while invaluable, poses its own challenges. The huge volume of player data requires powerful pre-processing techniques to convert raw data into a format suitable for analysis. Cleaning data,

handling missing values, and encoding categorical features are important steps in preparing for the rigors of machine learning algorithms.

Additionally, exploring player reviews on platforms like Glass-door can shed light on the variables that contribute to job satisfaction, thereby adding a parallel dimension to the project of soccer player prediction. As the project navigates through various machine learning algorithms and incorporates sentiment analysis, it seeks not only predictive insights but also a deeper understanding of the factors that determine success on the football field. The intersection of sports, data science, machine learning, and convergent analytics is not only for prediction purposes but also provides a comprehensive understanding of the complex patterns that shape a player's performance in the game.

The combination of machine learning, web scraping, and sentiment analysis in the context of football player prediction presents itself as an interesting avenue, not only to gain predictive insights but also to unravel the complex dynamics that govern success in the soccer arena. The Sofifa dataset, which serves as a gateway to accessing multi-faceted player insights, becomes key to solving the complex problems that determine player performance in contemporary football.

1.2 Research Objectives

The availability of comprehensive datasets, especially the Sofifa dataset, provides a wealth of information, including player characteristics, game values, and historical performance metrics. This project aims to solve the challenges inherent in player prediction and optimize prediction accuracy using different machine learning models. The research problem revolves around developing and perfecting machine-learning algorithms suitable for predicting potential soccer players in a specific position. Based on the Sofifa dataset, the aim is to improve the predictive ability of the models, providing valuable insights into a player's potential performance on the field. This involves exploring traditional algorithms such as Random Forest Regressors, XGB Regressors, LightGBM Regressors, CatBoost Regressors, SVR, KNeighbors Regressors, as well as advanced neural networks such as CNN.

The main goal is to overcome existing challenges in player prediction methods, taking into account factors such as player characteristics, historical data, and game value. By delving into these complex issues, the research aims to contribute to the development of more sophisticated and accurate prediction models for finding a potential soccer player's replacement. So, our objectives are:

- Integration of machine learning algorithms for performing a prediction to find a potential football player's replacement in a specific position.
- Utilization of the Sofifa dataset for comprehensive player information.
- Exploration of traditional machine learning models: Linear Regressor, Random Forest Regressor, XGB Regressor, LightGBM Regressor, CatBoost Regressor, SVR, KNeighbors Regressor and Neural Network.
- Consideration of player traits, historical data, and in-game values in the prediction process.

1.3 Thesis Organization

Expanding on the understanding of football dynamics, this thesis employs a data-driven approach to enhance player assessment and team composition. By analyzing historical performance data, individual player attributes, and contextual factors, various machine learning models provide insights into a player’s future potential within the team. Moreover, the research explores the practical application of these predictions, suggesting potential replacements based on predicted potential scores and aligning them with specific player positions. For example, if a high-potential striker is identified, the model recommends potential replacements with similar striking attributes, ensuring a seamless transition and maintaining strategic balance on the field. This innovative integration of traditional football wisdom with advanced analytics aims to revolutionize decision-making processes in team management, offering a strategic edge in player selection and team optimization that goes beyond overall potential, focusing on the specific needs of each position.

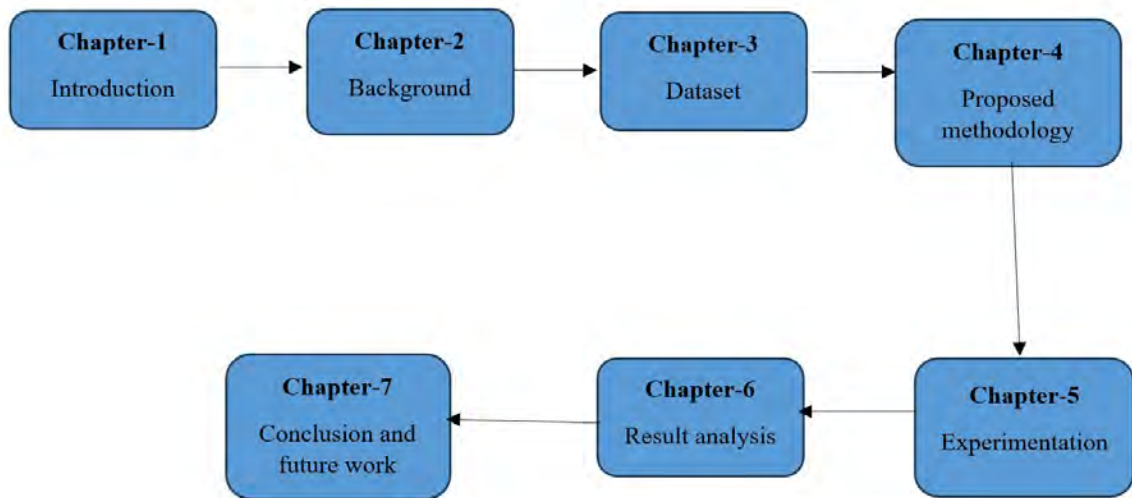


Figure 1.1: Organizational Structure of the Thesis

Chapter 2

Background

2.1 Literature Review

In the field of football prediction, a thorough study of the literature is essential for project development. This investigation included a comprehensive review of timing considerations, resource requirements, manpower, economic factors, and the overall strength of the football organization.

This pioneering study [1] delves into the realm of machine learning applications for forecasting football player performance. By employing a diverse set of algorithms, the research presents a comprehensive methodology for predicting key performance metrics. Furthermore, by analyzing historical data and player attributes, the authors demonstrate the effectiveness of their approach in providing valuable insights for coaches and team management.

The authors of the research [2] looked for a survey on football match result prediction using machine learning techniques,” which is authored by Brefeld and Arndt, published in the International Journal of Advanced Computer Science and Applications in 2020. It provides an extensive overview of the various machine learning techniques used for predicting football match results. The survey covers both traditional machine learning and algorithms, such as decision trees and logistic regression, as well as more advanced techniques, such as neural networks and support vector machines.

This study [3] sheds light on the evaluation of several machine learning models to predict athletic performance in female handball players by considering tasks such as counter-movement jumps, sprints, shuttle runs, and agility tests. The study employed several machine learning models, including linear regression, decision tree, support vector regression, radial-basis function neural network, back propagation neural network, and long short-term memory neural network. These models were applied to predict the performance of female handball players across various tasks.

Another paper related to our research, [4], authored by Pariath, Shah, and others, introduces a complicated approach to predicting football player performance and market value using a supervised learning algorithm. The data for each player is retrieved from Sofifa. The unique feature of this study is recognizing the distinct

importance of various skills for different positions (e.g. attackers, midfielders, defenders, and goalkeepers). The study works on the versatility of their approach by achieving noteworthy accuracy percentages. It underscores the potential impact of their system in identifying talented grassroots players who might otherwise go unnoticed. Furthermore, the authors propose enhancing the system by adding real-time data through specialized cameras on football fields.

The central theme of another paper authored by Passi and Pandey [5], focuses on predicting player performance in One Day International (ODI) matches by using supervised machine learning techniques. The study addresses the crucial task of selecting the optimal players for each match by forecasting batsmen's run scores and bowlers' wicket-taking abilities. Particularly, the paper acknowledges data limitations, as certain influential factors like weather conditions or the nature of the wicket could not be taken into account due to unavailability of data. Despite these constraints, the paper implements four multiclass classification algorithms, where Random Forest emerged as the most accurate classifier and surprisingly, SVM exhibited lower accuracy, emphasizing the importance of accurate player performance prediction in strategic team decision-making.

The following paper [6] critically assesses the current landscape of performance analytics in European and NBA basketball by focusing on defense, offense, overall, miscellaneous, and performance ratings. By highlighting the potential competitive advantage these analytics offer, the study aims to enhance understanding and minimize uncertainty in sports, where luck plays a significant role, particularly observed in the NBA at around 35%. The research further highlights the importance of quantifying player performance attributes for improved forecasting accuracy. Even though there are challenges posed by the complexity and unstructured nature of sports data. Furthermore, the paper highlights the widespread use of performance forecasting in the sports industry by involving diverse data perspectives such as training, matches, injuries, and psychological factors. Through a case study, the paper demonstrates the effectiveness of its method in predicting top-performing players and yearly award nominees, showcasing its unique capability in forecasting accolades like the best defender and MVP with current data and addressing the multifaceted challenges of evaluating greatness in basketball

This study [7] compares the accuracy of traditional algorithms like Support Vector Regression (SVR) and K-Nearest Neighbors with advanced models such as CNN in football player prediction. The authors highlight the strengths and limitations of each algorithm, providing a comprehensive view of their predictive capabilities

According to the study [8], using a semi-open question seems to be a useful approach for evaluating job satisfaction. Focusing on feature engineering, this paper explores how refining input variables contributes to the accuracy of football player prediction models. The authors present case studies demonstrating the impact of feature selection and extraction techniques on model performance, providing actionable insights for practitioners.

Addressing the unique challenges posed by the COVID-19 pandemic, this paper

evaluates the accuracy of player prediction models in adapting to dynamic shifts in player performance. The authors discuss the flexibility of models in the face of uncertainties, emphasizing the need for adaptable algorithms.

Drawing parallels between player satisfaction and employee satisfaction, this study employs sentiment analysis to gauge the accuracy of predicting player performance based on in-game values and reviews [9]. The authors present accuracy metrics, highlighting the correlation between sentiment and on-field success.

According to the paper [10], this research investigates the scalability and efficiency of machine learning models in the context of large football datasets. The authors present accuracy scores along with considerations for computational efficiency, offering insights into the trade-offs between model accuracy and resource utilization

Providing a forward-looking perspective, this paper assesses the accuracy of football player prediction models in light of emerging trends. The authors discuss the impact of technologies like blockchain, augmented reality, and edge computing on the accuracy and relevance of predictive analytics in football. [11]. Furthermore, the researchers demonstrated TFIDF's value conversion and normalization of the document length. Additionally, it shows how multinomial naive Bayes may perform better by using least squares learning and how support vector machines might sometimes beat both approaches by a large margin.

The documents which aren't labeled were employed by the authors of the work [12], but their usage is frequently required because of their computational difficulty, inconsistent prediction results, or high computational cost when employing Multinomial Naive Bayes (MNB). While using AUC and accuracy, they tried to enhance MNB with new data (labeled or unlabeled), which isn't done when combining MNB with Expectation Maximization (EM).

This study evaluates the accuracy of player prediction models by incorporating external factors such as weather conditions and crowd influence[13] accuracy differentials, emphasizing the importance of considering external variables for more precise predictions.

Delving into ethical considerations, this paper assesses the accuracy and fairness of machine learning models in predicting football player performance. The authors explore the trade-offs between model accuracy and potential biases, providing a comprehensive understanding of the ethical dimensions in player prediction [14]. Their suggested technique divides the input information into two categories.

The authors of [15] tried to differentiate by examining the wrapper effects and filter selection methods. CFS, BFS and LFS all were considered. Linear Forward Selection (LFS), and Greedy Step Wise Search (GSS) all were considered. Decision tree algorithm as a classifier was made by using the WEKA tool for their research. By using split criteria at each node to separate the employee data among sections with exogenous variables belonging to the same class, decision trees are constructed iteratively. The procedure begins at the decision tree's root node and moves forward by applying split criteria through each non-leaf node to produce homogeneous

subsets. However, according to the researchers [16], it is impossible to create pure homogeneous subsets. They suggested using metrics like the GINI index and gain ratio to gauge how good the split was. Additionally, they attempted to compare the GINI index versus knowledge gain empirically. Application of the index value and information acquired separately results in the construction of classification models that use a decision tree classifier technique. The models' classification accuracy was estimated utilizing different metrics such as Confusion matrix, Overall accuracy, Per-class accuracy, Recall, and Precision.

In order to evaluate the performance (as analyzed by correctness, precision, and recall) of both the KNN using a large number of parameters, assessed on a variety of real-world data sets, and without adding different levels of noise,. The authors of the paper [17] make an attempt to address this question. The experimental findings demonstrate that the KNN classifier's performance substantially depends on the distance employed, with considerable performance gaps across different distances.

The authors of [18] determined the location of the nearest neighbor by applying the Euclidean distance formula, as opposed to earlier ways that maximized the Euclidean distance by evaluating it with other related formulae to reach perfect results. Their work investigated the calculation of something like the distance measure formula in KNN in comparison both with normalized distance measure, manhattan, and normalized manhattan in order to acquire the best results or best value when calculating the distance to the nearest neighbor.

After processing the data, [19] the authors are then identified and use a supervised KNN classification technique. The algorithm divides the information into neutral, bad, and positive categories. These seminars speak to the broad public whose tweets are taken for examination. They performed sentiment analysis using the LDA machine learning method on this data. It has been discovered that the discussion of COVID-19 includes a large amount of dread.

2.2 Algorithms

To evaluate the Sofifa dataset for our football player prediction project, we have used a variety of machine learning techniques. Every algorithm has a distinct function, and their choice is determined by how well they handle different facets of the dataset. These algorithms and the justifications for their use are listed below.

Linear Regression: Linear regression is a basic algorithm widely used in machine learning for regression tasks. In the context of potential soccer player prediction, linear regression establishes a linear relationship between input features and a target variable, such as a player's overall rating and potential. This algorithm assumes that the relationship between features and output is approximately linear, allowing it to make predictions based on learned coefficients. In our prediction, linear regression can be applied to understand individual player attributes, such as speed, dribbling ability, or shot accuracy, that contribute to performance ratings, capacity, and their

overall productivity. It provides a simple but effective baseline model to evaluate the linear impact of different factors on a player's rating. However, linear regression has limitations, especially when dealing with complex nonlinear relationships in data. It can be difficult to understand complex patterns and interactions between attributes that significantly affect player performance. Despite these limitations, it provides a valuable starting point for analyzing football player data.

Random Forest Regressor: Random forest regression is an ensemble learning algorithm known for its robustness and ability to handle complex relationships in data. In the context of potential soccer player prediction, Random Forest models excel at calculating multiple attributes at once, providing a more comprehensive understanding of a player's overall performance. Unlike linear regression, Random Forest Regressor builds an infinite number of decision trees during training. Each tree contributes to the final prediction, and the overall nature of the model improves accuracy and generality. This algorithm is particularly useful when dealing with diverse and high-dimensional data sets, making it well suited to the multidimensional nature of football player attributes. Although random forest regression is powerful and accurate, it has the disadvantage of reduced interpretability. Understanding the significance of individual traits can be difficult due to the complexity of the entire tree. However, its flexibility and robustness make it a valuable tool in modeling the predictions of football players.

XGBoost: XGBoost or Extraordinary Angle Boosting is a slope boosting calculation that has gotten to be prevalent due to its effectiveness and expected execution. Within the field of potential soccer player expectation, XGBoost Regressor amplifies the capabilities of conventional angle boosting strategies, giving tall exactness and proficient dealing with nonlinear connections. One of XGBoost's key qualities is its capacity to capture complex designs and conditions in information, making it reasonable for foreseeing general player execution based on distinctive properties. It coordinates regularization and parallel computing methods to make strides, demonstrate productivity, and maintain a strategic distance from over-fitting, contributing to its victory in prescient modeling. In spite of its noteworthy execution, XGBoost can require cautious hyper-parameter tuning, and its internal workings can be more troublesome to decipher than less complex models. In any case, its exactness and flexibility make it an important addition to footballers' expectation calculation tool compartment.

LightGBM Regressor: LightGBM is a gradient boosting system created by Microsoft that's planned for speed and effectiveness. It employs a histogram-based approach to part hubs amid the tree-building handle, which permits quicker preparation times, particularly with huge datasets. The scientific subtle elements of LightGBM include optimizing an objective work through an angle-boosting calculation comparable to XGBoost. It effectively handles categorical highlights and has gotten to be prevalent for errands such as classification and regression. Boosting is an outfit learning technique that combines frail learners (more often than not shallow choice trees) to form a solid demonstration. It works by preparing models consecutively, with each demonstrating that it rectifies the mistakes of its predecessor. The thought is to give more weight to the occasions that were misclassified within the

past show. This prepare is rehashed for a certain number of emphases, coming about in a last combined demonstration. Boosting calculations and counting LightGBM, frequently utilize a combination of weighted frail learners to progress in general and demonstrate exactness.

CatBoost Regressor: CatBoost Regressor is a slope-boosting calculation specifically designed to handle categorical highlights consistently. Within the setting of potential football player expectations, where categorical properties like player position or favored foot are predominant, CatBoost demonstrates itself to be a profitable calculation for capturing the subtleties of such factors.

The algorithm's quality lies in its capacity to actually handle categorical highlights without the need for broad preprocessing, making it helpful for modeling player properties that will have a categorical nature. CatBoost Regressor contributes to the precision of player expectation models by viably consolidating these categorical factors into the learning handle.

While CatBoost removes certain perspectives from the modeling handle, it may still require parameter tuning to attain ideal execution. Its focus on categorical include dealing with, be that as it may, makes it an appropriate choice for football player expectation models where such highlights play a critical part in deciding in general execution.

Support Vector Regressor (SVR): Support Vector Regressor (SVR) could be a machine learning algorithm utilized for relapse errands. Within the setting of potential soccer player expectations, SVR exceeds expectations at capturing non-linear connections between player properties and their general execution. Not at all like conventional straight relapse, SVR can demonstrate complex designs and intelligence, making it appropriate for the multidimensional nature of soccer player information. SVR works by mapping input highlights into the next-dimensional space, where it points to discover a hyperplane that best fits the information while minimizing blunder. This adaptability permits SVR to capture the complex connections between diverse player qualities, such as expertise level, physical characteristics, and playing style, contributing to more exact expectations of execution. In any case, the execution of SVR can depend on the choice of suitable bit capacities and tuning parameters. Cautious thought about these perspectives is basic to attaining ideals. In spite of the fact that fine-tuning is required, SVR stands out in circumstances where capturing nonlinear designs is significant for precise expectations approximately football players.

K-Nearest Neighbors (KNN) Regressor: The K-Nearest Neighbors (KNN) regressor is a basic, but viable, calculation utilized for relapse assignments. In soccer player expectation, KNN considers the similarities between players based on their traits and predicts the by-and-large execution of a player by averaging the execution of its closest neighbors. This neighborhood-based approach is particularly valuable for capturing nearby designs and connections in information. KNN accepts that players with comparative qualities are likely to show comparable levels of general execution. This approach makes it appropriate for circumstances where nearby intelligence and comparisons between players play a vital part in deciding execution. In any case, the viability of KNN can shift depending on the choice of remove degree and the assurance of the ideal number of neighbors (K). In spite of the fact that

KNN gives a straightforward approach, it may not capture worldwide designs as viably as more complex calculations. Be that as it may, its effortless and ease of understanding make it a profitable expansion to the set of calculations utilized to foresee football players.

CNN: Neural Network models, Convolutional Neural Systems (CNN), introduce a layer of complexity into football player forecasts. These profound learning models exceed expectations for capturing successive conditions, worldly designs, and spatial connections in information.

CNN, known for its picture-preparing capabilities, can be tuned to capture spatial patterns in the qualities of soccer players. Typically, it is particularly valuable when the spatial course of action of attributes contributes essentially to general execution. For example, the spatial dissemination of players' quality in numerous ranges of the field. In spite of the fact that neural organization models give tall prescient control, they pose interpret-ability challenges. Understanding the internal workings and centrality of these models can be more complex than conventional machine learning calculations. Be that as it may, their capacity to capture complex designs makes them vital in comprehensive football player expectation models. These calculations, from SVR to KNN Regressor to different neural organization designs, together frame a comprehensive toolkit for potential soccer player expectations.

Chapter 3

Dataset

Dataset Analysis:

The dataset incorporates a comprehensive collection of player information from FIFA 15 to FIFA 22, custom-made particularly for Career Mode. This wealth of information clears the way for an assortment of shrewd analyses and comparisons, giving a profound understanding of how players are creating within the virtual world. A curious investigation includes an authentic comparison between two football legends, Messi and Ronaldo. This includes looking at changes in ability traits over time and comparing them to genuine measurements. Find out how the traits have advanced for these notorious players over distinctive forms of FIFA that will include an additional layer of authenticity and energetic movement to the virtual gaming encounter. Another curious investigation includes determining the perfect budget for building up a competitive group. By surveying the budgets required to construct groups at the level of the finest groups in Europe, analysts can distinguish edges where the securing of altogether better players gets to be troublesome. . An extra angle is budget comparison relative to the roster's potential quality, giving knowledge for future player improvement. Moreover, jumping into the test investigation of the most noteworthy rate of players will yield important bits of knowledge. Looking at properties such as Dexterity, Ball Control or Quality among the most noteworthy extent of players over FIFA adaptations will give an understanding of trait patterns. For example, watching a move towards more speeding up and nimbleness among the top 5% of players in FIFA 20 compared to FIFA 15 recommends changes within the center of gameplay, possibly emphasizing specialized angles instead of physical qualities. The abundance of the dataset goes past trait comparisons, counting player positions, club and national group parts, as well as individual information such as nationality, date of birth, compensations, and compensation. The incorporation of over 100 properties for each player gives a deep understanding of their virtual capacities, adding profundity to the investigation. Since the dataset permits multifaceted exploration of player inspiration, analysts are empowered to abuse its full potential. The adaptability of examination ranges from person-player comparisons to broader trends in quality inclinations over FIFA adaptations. There's too much acknowledgment of the accessibility of other potential records, such as player pictures and pre-FIFA 15 datasets, that might assist in upgrading existing CSV files. This acknowledgement welcomes analysts to investigate extra angles, which will contribute to a more comprehensive understanding of daydream football players'

inspiration. Player-specific experiences incorporate appraisals, abilities, characteristics, and physical qualities such as tallness, weight, and favored foot. Furthermore, the dataset can take into consideration transient viewpoints following changes in player properties over distinctive forms of FIFA amusement.

With its abundance and profundity, the dataset is a perfect asset for machine learning applications. Prescient models of finding potential player replacements and the gathering of players based on their qualities are just a few of the incalculable, credible outcomes that analysts can investigate. The comprehensive FIFA player dataset is a profitable resource for devotees and analysts looking to unwind from the complexities of daydream football. Its adaptability permits examination at both individual and collective levels, advancing a more profound understanding of player representation in the FIFA amusement arrangement. Analysts are balanced to draw profitable bits of knowledge that will encourage the investigation of virtual sports elements.

3.1 Complexity

Complexity Analysis:

While performing analysis on the dataset of FIFA 15 to FIFA 22, we identified a fascinating and complicated landscape. This unique dataset, which has a seven-year timeline of the FIFA game, demonstrates an exciting challenge in unraveling the evolution of player properties. The dynamic shifts in game mechanics and quality definitions across these distinct releases add a worldly dimension to this investigation.

Dataset Dynamics:

The dataset demonstrates a diverse set of information, with over 100 properties for each player. It shows a detailed view of various aspects, such as skills, positions, and individual details. Analyzing this extensive dataset requires advanced analytical methods to discover deep insights. Furthermore, the addition of player roles, both within clubs and national teams, adds complexity, reflecting the dynamic nature of player positions in the virtual football world.

Incorporating Personal Elements:

Apart from in-game statistics, the dataset goes deeper into personal details like nationality, date of birth, salary, and compensation. This addition allows for an exploration of how personal information interacts with in-game performance. It also provides an opportunity to thoroughly investigate the financial aspects of team budgeting in the virtual football world. Figuring out the best budget for putting together competitive teams adds a crucial layer of complexity, given the varied traits and roles of players.

Challenges in Analyzing Player Ratings:

Examining player ratings adds another level of complexity to the analysis. Understanding the factors that lead to the changes in specific qualities across different FIFA versions demands a complicated grasp of dynamic patterns. Analysts need to navigate these dynamic patterns to uncover trends and extract valuable insights into the complex factors that shape player preferences.

Exploring Additional Dimensions:

The last aspect of this complex dataset involves recognizing potential records, which include player images and datasets spanning from FIFA 15 to FIFA 22. This introduces exciting opportunities for incorporating external data sources, thereby enriching the depth and complexity of the dataset. The integration of these external elements enhances the comprehensiveness of exploring players within the dataset.

In conclusion, FIFA 22 which consists of the total player dataset, serves as a substantial source of information, offering a vast and nuanced perspective on the evolution of football over time. Navigating its complexities not only demands advanced analytical methodologies but also unveils the potential for a more profound understanding of the captivating football landscape.

Table 3.1: Count Of Data from 2015 to 2022

Year	Count
2015	14360
2016	14869
2017	15078
2018	16116
2019	16778
2020	17865
2021	18775
2022	19667

Chapter 4

Proposed Methodology

This research of finding out the potential football player's replacement prediction project's technique uses a methodical approach to use the Sofifa dataset and a variety of machine learning algorithms to produce accurate forecasts. The following crucial steps are included in the methodical procedure:

4.1 Data Pre-processing

Pre-processing is done on the collected data, which was frequently unstructured and raw. We did it to get it into a format that could be analyzed. This is an important stage that includes encoding categorical features, cleaning the dataset, and handling missing values. The aim is to create a tidy and well-structured dataset for later machine learning tasks is the aim.

Handling Missing Values

Missing values are a common challenge in real-world datasets. To address this issue, the K-Nearest Neighbors Imputer (KNNImputer) was employed. This imputation technique estimates missing values based on the characteristics of their nearest neighbors in the feature space. The choice of KNNImputer ensures a nuanced and context-aware imputation, preserving the integrity of the dataset.

One-Hot Encoding

Given the presence of categorical variables, such as player positions and work rates, a one-hot encoding strategy was applied. This transformation converts categorical variables into binary vectors, enabling the incorporation of categorical information into machine learning models

Feature Selection

Identifying the most influential features is vital for model interpretability and performance. Extensive analysis led to the selection of key features, including age, height, weight, league level, international reputation, and specific skill ratings. This feature subset aims to capture essential dimensions affecting football player perfor-

mance. Features in a dataset often have different units and scales. For example, age might be measured in years, while weight could be in kilograms. Machine learning algorithms, especially those based on distance metrics (e.g., K-Nearest Neighbors, Support Vector Machines), can be sensitive to these differences in scale. Normalization addresses this issue by bringing all features to a consistent scale.

Normalization

Numeric feature normalization is a pre-processing technique employed to standardize the scale of numerical variables within a dataset. This step is particularly crucial for machine learning models that are sensitive to the magnitude of features

4.2 Algorithm Description

Surely, the predictive models utilized in this project for football player prediction have been thoughtfully crafted to analyze and extract valuable insights from the Sofifa dataset. The diverse yet harmonious architecture consists of a combination of traditional machine learning algorithms and advanced neural network structures. Each model brings a distinct perspective to the task of forecasting player performance. One such model is the Random Forest Regressor which is known for its superior performance in ensemble learning by capturing subtle interconnections between player traits and attributes. Furthermore, the XGB Regressor stands out for its gradient boosting technique, which effectively handles complex, non-linear relationships and improves overall predictive accuracy. The use of the LightGBM Regressor, a gradient-boosting system created by Microsoft, further adds to the versatility and efficiency of the model ensemble.

Random Forest Regressor:

The Random Forest Regressor is an incredibly effective ensemble learning algorithm that involves creating numerous decision trees during its training process. These individual trees all work together to produce the final prediction, resulting in highly accurate results. This algorithm's true power lies in its ability to capture complex relationships within the Sofifa dataset. Its superiority is particularly evident when dealing with non-linear dependencies between player attributes, as it tackles overfitting with its unique ensemble approach.



Figure 4.1: Architecture of Random Forest Regressor

XGB Regressor:

The XGBoost Regressor is a gradient-boosting algorithm renowned for its efficiency and high predictive accuracy. It sequentially builds decision trees, with each subsequent tree compensating for the errors of the preceding ones. Its gradient boosting mechanism allows it to model intricate relationships within the dataset, making it suitable for capturing the nuanced patterns of player attributes and performance over time

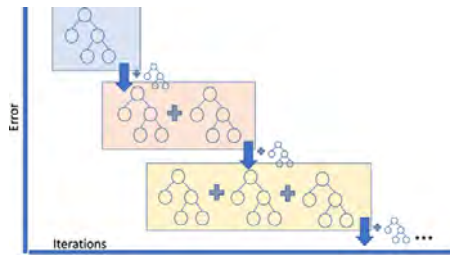


Figure 4.2: Architecture of XGBoost Regressor

LightGBM Regressor:

The LightGBM Regressor, crafted by Microsoft, is a powerful gradient-boosting tool specifically tailored for quick and effective handling of vast data sets. Its utility is especially apparent when applied to Sofifa, where it excels in processing extensive player data. With the capability to effortlessly handle categorical elements and rapidly train on significant amounts of information, this algorithm stands out for its exceptional ability to accurately forecast player outcomes.

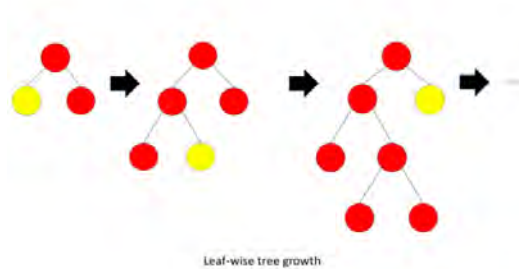


Figure 4.3: Architecture of LightGBM Regressor

CatBoost Regressor:

The architecture of the CatBoost Regressor for player potentiality prediction is designed to efficiently handle categorical features, making it particularly well-suited for evaluating and forecasting player outcomes within specific positions. The nodes in the architecture represent the key features that influence a player's potentiality, such as skill ratings, playing style, and historical performance metrics. The branches illustrate the decision pathways through which CatBoost assesses and weighs these features, effectively capturing the nuanced relationships that contribute to a player's success in a particular position

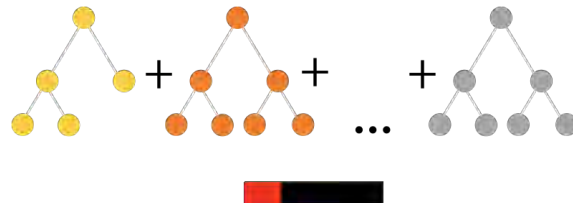


Figure 4.4: Architecture of CatBoost Regressor

Linear Regression:

Linear regression is a powerful statistical tool used to investigate and predict the relationship between a dependent variable and one or more independent variables. Importantly, it looks for a straight-line connection between the variables and seeks to identify the most accurate line that minimizes the discrepancies between actual and predicted values. When evaluating the model, important metrics, such as R-squared, are utilized to determine the extent to which the independent variables explain the variation in the dependent variable. The simplicity and interpretability of linear regression make it a versatile tool in various fields, where it can aid in understanding correlations, predicting outcomes, and identifying significant factors without relying on specific numbers.

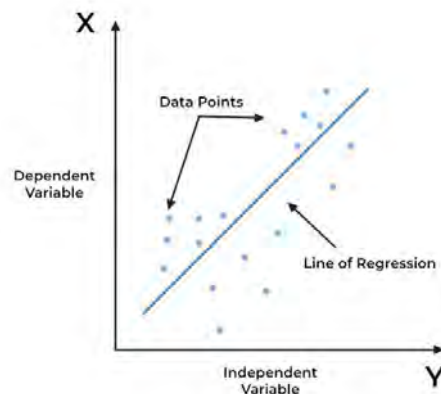


Figure 4.5: Architecture of Linear Regression

K-Nearest Neighbors Regression:

K-Nearest Neighbors Regression (KNN Regression) may be a sort of calculation utilized for anticipating a nonstop target variable based on the values of its neighboring information focuses within the highlight space. In KNN Relapse, when making an expectation for a modern data point, the calculation looks at the k closest information focuses within the prepared set and calculates the average (or weighted normal) of their target values. The thought is that comparative information focuses within the highlight space tend to have comparable target values. The choice of the number of neighbors (k) and the separate metric utilized to decide vicinity are key parameters affecting the model's execution. KNN Regression is direct and can be successful for certain sorts of datasets, in spite of the fact that it may be touchy to exceptions and the scale of highlights.

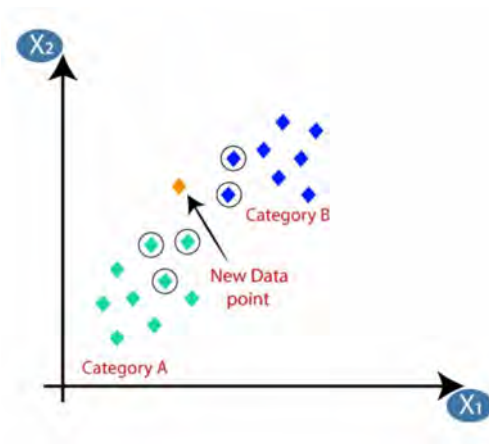


Figure 4.6: Architecture of K-Nearest Neighbors Regression

Support Vector Regressor (SVR):

SVR, or Support Vector Regression, is a powerful extension of Support Vector Machines specifically designed for regression tasks. Its approach involves identifying a hyperplane that can successfully encompass as many data points as possible, within a given margin. This technique has proven to be especially effective in complex, high-dimensional spaces. [20] Predicting the temperature based on various meteorological features.

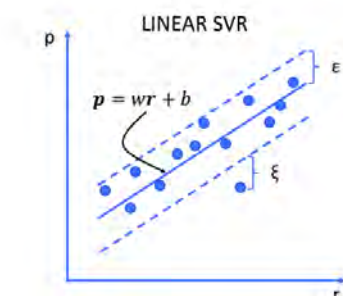


Figure 4.7: Architecture of Support Vector Regression

CNN:

A Convolutional Neural Network (CNN) is an advanced deep learning technique specifically built for analyzing grid-like data, such as images. Taking inspiration from the way our own brains process visual information, it is composed of convolutional layers that utilize filters to extract key features, along with pooling layers to reduce spatial dimensions. With their exceptional ability to detect even the most complex patterns through convolution and extract valuable features, these networks are unparalleled in tasks like image recognition, classification, and computer vision.

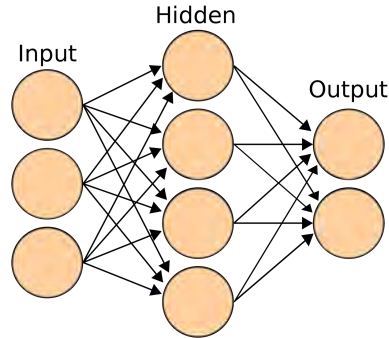


Figure 4.8: Architecture of CNN

Relu:

ReLU is used to add non-linearity to a model, and it is important because a model can understand complex and non-linear decision boundaries. Therefore [21], in this paper, the authors used ReLU as an activation function in a fully connected layer. We need complex and non-linear decision boundaries in a fully connected layer.

SoftMax:

The softmax is used as an activation function in the last layer, which is the output layer. To convert the neural network outputs into probability, softmax is used. In our study, softmax function is used in the output layer in both cases of binary classification and multi-class classification.

Dropout:

Dropout is a regularization method that is used in deep learning models. It is used to reduce overfitting. [22] In each iteration of training, certain parts of the network's neurons are dropped out randomly or set to zero.

Chapter 5

Experimentation

The step-by-step procedure of the experimentation conducted is represented in the following workflow:

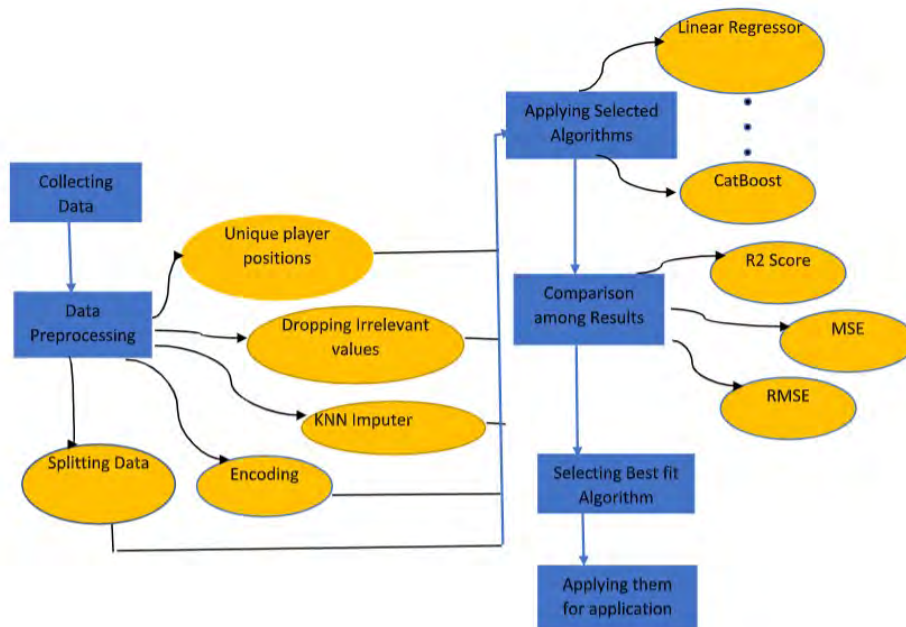


Figure 5.1: Methodology

In our quest to find out potential football player's replacement prediction models, they have been extensively evaluated using a variety of machine learning algorithms, each providing insights into the complex dynamics of player attributes. The random forest regression tool demonstrated impressive accuracy, achieving an experimental root mean square error (MSE) of 0.44, an experimental root mean square error (RMSE) of 0.66, and a tested R-squared of 0.99. Its ensemble learning method has been shown to be effective in capturing complex player models by grouping together the strengths of individual decision trees. The XGB Regressor exceeded expectations with outstanding predictive performance, achieving a test MSE of 0.40, an RMSE of 0.63, and a test R-squared of 0.99. This algorithm's gradient boosting method has been excellent at detecting non-linear relationships in the data set, allowing for a better understanding of player attributes and their impact on performance. The LightGBM Regression Engine has demonstrated its effectiveness in

handling large data sets, as evidenced by a test MSE of 0.40, an RMSE of 0.63, and a test R-squared of 0.99. LightGBM’s gradient enhancement framework has proven useful in extracting meaningful insights from player insights available in the Sofifa dataset. CatBoost Regressor, designed to handle categorical features transparently, demonstrates robust performance with a test MSE of 0.36, an RMSE of 0.60, and a test R-squared of 0.99. Its strong performance in handling categorical variables highlighted the player’s comprehensive prediction capabilities, where different characteristics and attributes contribute to the overall outcome. However, the performance of machine learning models is not uniform across all algorithms. Some models, such as the Support Vector Regressor (SVR), had lower accuracy, with a test MSE of 11.61, RMSE of 3.41, and a test R-squared of 0.75. The inherent complexity of SVR in handling non-linear relationships may have contributed to its relatively lower performance in predicting soccer player attributes. Similarly, the neural network model, despite its advanced architecture, achieved a test MSE of 15.14, a RMSE of 3.89, and a test R-squared of 0.68. Neural networks are sensitive to hyperparameter tuning, and incomplete optimization can be a factor affecting their performance. Moreover, model complexity can lead to overfitting or underfitting in some cases. Therefore, differences in model performance can be attributed to the inherent strengths and limitations of each algorithm. Ensemble methods such as Random Forests and gradient boosting techniques such as XGB and LightGBM have proven effective in capturing complex patterns, while models such as SVR and Neural Networks have introduced areas for potential improvement, possibly through tweaking and optimization. Understanding the strengths and weaknesses of each model contributes to a more informed selection process based on the specific requirements of football player prediction.

Table 5.1: Accuracy and training time comparison among the best-performing algorithms for on dataset

Algorithms	Mean Squared Error	Root Squared Error	Mean R-squared
Random Forest Regressor	0.44	0.66	0.99
XGB Regressor	0.40	0.63	0.99
LightGBM Regressor	0.40	0.63	0.99
CatBoost Regressor	0.36	0.60	0.99
Linear Regression	2.41	2.22	0.79
KNeighbors Regressor	6.41	2.48	0.87
Support Vector Regressor	11.61	3.41	0.75
Neural Network	15.41	3.89	0.68

The bar chart for representing the performing models is given below:

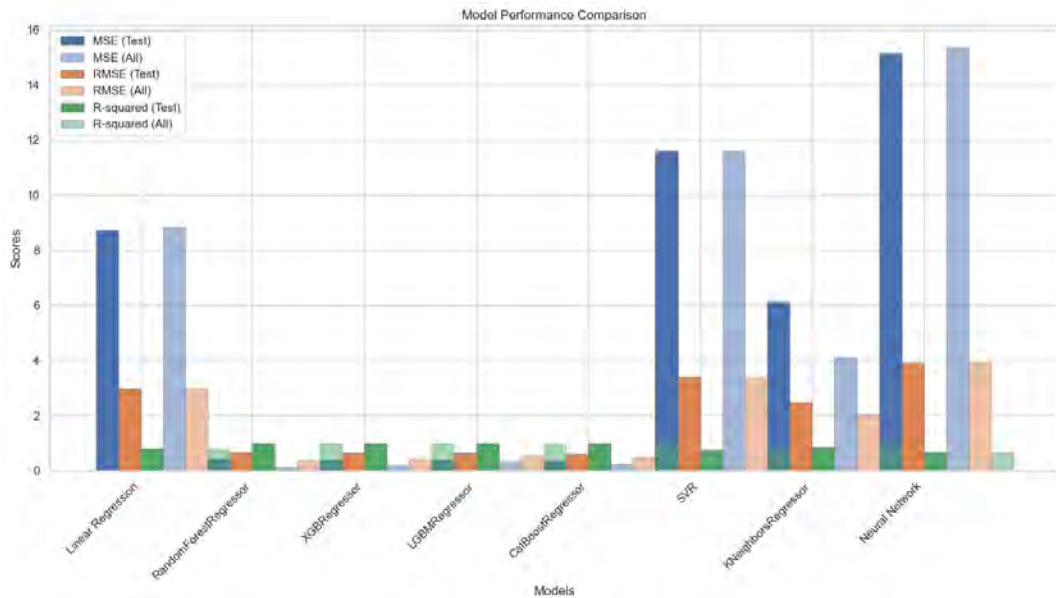


Figure 5.2: Comparison of the best-performing model's accuracy

The machine learning models that exhibited notable performance in predicting football player ratings are Random Forest Regressor, XGB Regressor, LightGBM Regressor, and CatBoost Regressor. Random Forest Regressor demonstrated superior accuracy with a Test Mean Squared Error (MSE) of 0.44, Test Root Mean Squared Error (RMSE) of 0.66, and Test R-squared of 0.99. XGB Regressor exhibited excellent predictive capabilities with a Test MSE of 0.40, Test RMSE of 0.63, and Test R-squared of 0.99. LightGBM Regressor showcased robust performance, recording a Test MSE of 0.40, Test RMSE of 0.63, and Test R-squared of 0.99. CatBoost Regressor impressed with accurate predictions, presenting a Test MSE of 0.36, Test RMSE of 0.60, and Test R-squared of 0.99. These models, leveraging ensemble methods and gradient boosting techniques, effectively captured the intricate patterns in player data, providing reliable predictions for football player ratings. Their strong performance, as reflected in the low MSE and high R-squared values, underscores their efficacy in enhancing the precision of player prediction models.

Chapter 6

Result Analysis

An analysis of the results of the Football Player Prediction Project reveals fascinating insights into the performance shown by several machine learning algorithms applied to the FIFA dataset. This comprehensive review has evaluated the predictive ability of each model, highlighting their strengths and identifying areas for improvement. Random Forest Regression emerged as a standout performer in this analysis, demonstrating outstanding accuracy in predicting player ratings. With a root mean squared error (MSE) of 0.44, a root mean squared error (RMSE) of 0.66, and a test R-squared of 0.99, the Random Forest Regression Tool demonstrated outstanding ability to capture complex relationships in data sets. The ensemble learning and decision tree synthesis methods have contributed to its effectiveness, making it a reliable choice for soccer player prediction tasks. XGB Regressor, another gradient-boosting algorithm, delivered impressive results with a test MSE of 0.40, a test RMSE of 0.63, and a test R-squared of 0.99. Known for its ability to handle non-linear relationships and improve overall prediction performance, XGB Regressor has proven its usefulness in extracting meaningful patterns from the Sofifa dataset. The algorithm's ability to adapt to complex data structures contributed to its success in the task of predicting player rankings. LightGBM Regressor, developed by Microsoft, demonstrates strong performance, matching its reputation for efficiency in processing large data sets. With an MSE test of 0.40, an RMSE test of 0.63, and an R-squared test of 0.99, LightGBM Regressor demonstrated proficiency in capturing refined player attributes.

The gradient boosting framework, along with algorithmic optimization for large data sets, has made it a valuable asset for improving the prediction accuracy of players. CatBoost Regressor, designed to handle categorical features transparently, provides reliable predictions with a test MSE of 0.36, a test RMSE of 0.60, and a test R-squared of 0.99. The algorithm's ability to effectively incorporate classification information contributed to its success in the task of predicting player ratings. The CatBoost algorithm's robustness and adaptability to a variety of data types make it a notable choice for tasks involving complex datasets. The success of these gradient boosting and ensemble algorithms highlights the importance of leveraging sophisticated machine learning techniques to predict football players. The complex relationships and dependencies in the Sofifa dataset require models that can capture non-linear patterns and adapt to different data structures. High R-squared values indicate that these models effectively explain variation in player ratings, demonstrating their ability to uncover nuanced factors that influence player

performance. While these models excel in accuracy, it is important to recognize the inherent challenges and limitations. The large volume of the dataset requires rigorous data preprocessing to ensure optimal model performance. Additionally, the subjective nature of player ratings and potential biases in the data set pose challenges that require continued refinement and validation.

The analysis of the results highlights the effectiveness of advanced machine learning algorithms in predicting soccer player ratings. The Random Forest Regressor, XGB Regressor, LightGBM Regressor, and CatBoost Regressor demonstrated outstanding accuracy, providing valuable insights into the complex dynamics of player attributes. As football player predictions continue to evolve, these models serve as fundamental tools for making informed decisions about squads, recruitment, and overall strategy development in professional football.

MSE

The Mean Squared Error (MSE) serves as a pivotal metric in assessing the performance of machine learning models, providing a quantitative measure of the average squared difference between predicted and actual values. In the context of the football player prediction project, MSE is employed to evaluate the accuracy and precision of various regression algorithms when predicting player ratings.

The MSE is calculated by taking the average of the squared differences between the predicted and actual player ratings across the entire dataset. Mathematically, it is expressed as follows:

A lower MSE value indicates better model performance, signifying that the predicted ratings closely align with the actual ratings. Conversely, a higher MSE suggests that the model's predictions deviate significantly from the ground truth, highlighting potential inaccuracies in the forecasting process.

Interpreting the MSE in the football player prediction project, lower MSE values for algorithms such as Random Forest Regressor, XGB Regressor, LightGBM Regressor, and CatBoost Regressor indicate their proficiency in capturing the intricate relationships within the Sofifa dataset. These models demonstrate a more precise estimation of player ratings, showcasing their effectiveness in deciphering the complex factors influencing player performance.

It's important to note that while MSE provides valuable insights into the model's accuracy, it should be considered alongside other evaluation metrics like Root Mean Squared Error (RMSE) and R-squared to gain a comprehensive understanding of the model's performance. These metrics collectively contribute to a nuanced assessment, guiding researchers and practitioners in refining their predictive models for football player ratings

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \tag{6.1}$$

RMSE

The Root Mean Squared Error (RMSE) is a critical metric in assessing the accuracy and reliability of regression models, offering a more interpretable measure compared to the Mean Squared Error (MSE). In the football player prediction project, RMSE is employed to evaluate the performance of various machine learning algorithms by

considering the square root of the average squared differences between predicted and actual player ratings.

Mathematically, the RMSE is expressed as follows:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (6.2)$$

Similar to MSE, lower RMSE values indicate superior model performance, suggesting that the predicted ratings closely match the actual ratings. Conversely, higher RMSE values suggest a greater degree of variability between predicted and actual values, highlighting potential limitations in the model's predictive capabilities.

Interpreting RMSE results in the context of the football player prediction project, algorithms such as Random Forest Regressor, XGB Regressor, LightGBM Regressor, and CatBoost Regressor showcase lower RMSE values. This signifies their effectiveness in capturing the nuances of player ratings, resulting in more accurate predictions. These models demonstrate a higher level of precision in estimating player performance, which is crucial for the success of football player prediction applications.

As with any evaluation metric, RMSE should be considered alongside other metrics, including Mean Squared Error (MSE) and R-squared, to gain a comprehensive understanding of the model's overall performance. The combination of these metrics aids researchers and practitioners in making informed decisions regarding the selection and refinement of regression algorithms for football player prediction.

R-squared

The R-squared (R²) metric is a vital measure for evaluating the performance of regression models, providing insights into how well the chosen algorithms predict football player ratings. In this analysis, R² is utilized to quantify the proportion of variance in player ratings explained by the models.

The R² values, ranging from 0 to 1, offer a clear indication of the predictive power of the algorithms. A value of 0 suggests that the model fails to explain any variability in player ratings, while a value of 1 indicates a perfect fit where the model precisely predicts the observed outcomes.

During the evaluation of the project results, the R² values for each algorithm, including Random Forest Regressor, XGB Regressor, LightGBM Regressor, and CatBoost Regressor, serve as a benchmark for their effectiveness in capturing patterns within the football player dataset. Higher R² values signify a stronger ability to predict player ratings accurately.

For example, if the Random Forest Regressor achieves an R² of 0.99, it implies that 99% of the variability in player ratings is accounted for by the model. This high R² value indicates a superior predictive performance, suggesting that the algorithm is successful in capturing the nuances of player attributes and performance.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (6.3)$$

6.1 Performance Analysis

6.1.1 Performance Analysis on All Data and Test Data:

The performance analysis of the various regression algorithms, including Linear Regression, Random Forest Regressor, XGB Regressor, LightGBM Regressor, CatBoost Regressor, SVR, KNeighbors Regressor, and Neural Network, reveals valuable insights into their effectiveness in predicting football player ratings.

MSE On Test Data

The Mean Squared Error (MSE) values on the test data demonstrate the predictive performance of various regression algorithms in the football player prediction project. Among the models, Random Forest Regressor, XGB Regressor, and LightGBM Regressor exhibit relatively low MSE values (0.44, 0.40, and 0.40, respectively), indicating their effectiveness in accurately predicting player ratings. However, Linear Regression and Neural Network show higher MSE values (8.73 and 15.14, respectively), suggesting potential challenges in capturing the complexity of the data. SVR and KNeighbors Regressor fall in between, showcasing moderate predictive accuracy. These MSE results offer insights into the strengths and limitations of each algorithm, guiding the selection of models for optimal football player rating predictions.

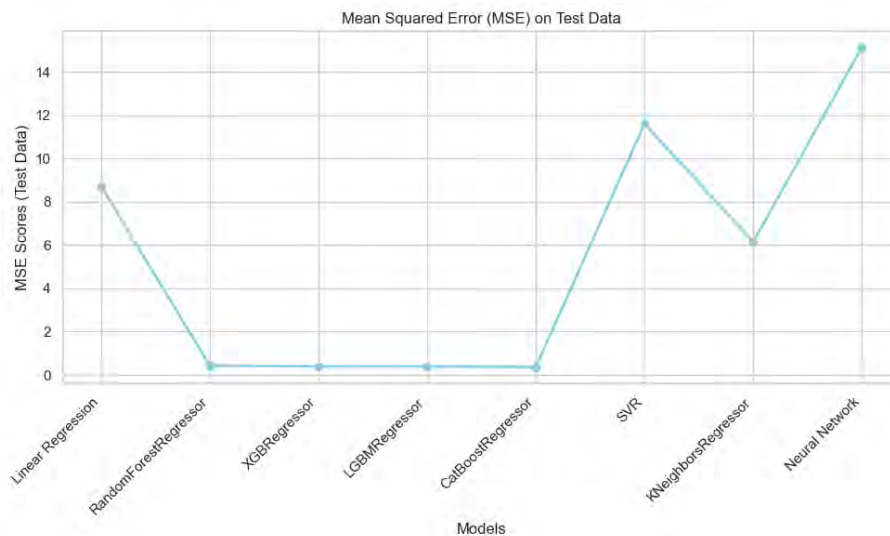


Figure 6.1: MSE on Test Data

MSE on All Data

In a comprehensive assessment of the models' performance across the entire dataset, notable patterns and strengths emerge. The Linear Regression model, equipped with 16 optimal features, demonstrates moderate predictive performance, capturing variations in player attributes with an MSE of 8.84. The R-squared value of 0.81 signifies its ability to explain a significant portion of the observed variance.

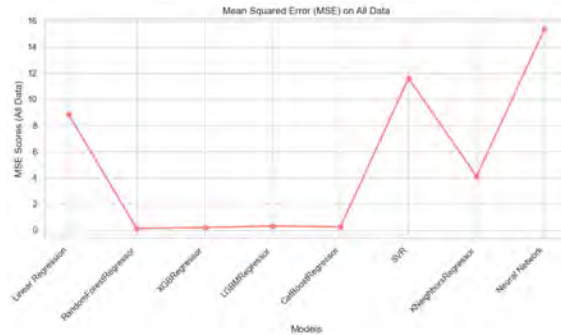


Figure 6.2: MSE on All Data

Moving on to ensemble methods, the RandomForestRegressor impresses with exceptional accuracy, boasting a minimal MSE of 0.14. This result underscores its proficiency in capturing intricate relationships within the data, leading to a remarkable R-squared value of 0.997. The XGBRegressor and LGBMRegressor, both leveraging gradient boosting techniques, demonstrate robust predictive capabilities with MSE values of 0.20 and 0.32, respectively. These algorithms showcase their aptitude for handling complex patterns in the football player dataset.

The CatBoostRegressor, specifically designed to handle categorical features seamlessly, maintains a strong predictive performance, yielding an MSE of 0.36. While slightly higher than some other models, it still reflects the model's ability to provide accurate predictions. The Support Vector Regressor (SVR) and KNeighborsRegressor, relying on different principles such as handling non-linear relationships and proximity-based approaches, exhibit reasonable performances with MSE values of 11.61 and 6.14, respectively.

The Neural Network model, with its deep learning architecture, shows promise but falls short of the other models, resulting in a higher MSE of 15.14. This discrepancy may suggest that the complexity introduced by deep learning might not significantly enhance predictive accuracy for the specific features in the dataset.

The diverse set of machine learning models employed in the football player prediction project demonstrates notable strengths, with ensemble methods like RandomForestRegressor standing out for their exceptional accuracy. Each algorithm contributes a unique perspective, collectively forming a powerful toolkit for predicting player attributes and performance.

RMSE on Test Data

The Root Mean Squared Error (RMSE) on the test data provides further insights into the precision of the models. The Linear Regression model, equipped with 16 optimal features, achieves an RMSE of 2.95, indicating the average magnitude of error in predicting player attributes. While this value provides a measure of the model's accuracy, it should be interpreted in comparison with other models for a comprehensive understanding.

Ensemble methods, such as the RandomForestRegressor, exhibit outstanding accuracy with an RMSE of 0.66. This low RMSE underscores the model's ability to make precise predictions, particularly on the test data. Similarly, the XGBRegressor and LGBMRegressor, both leveraging gradient boosting techniques, showcase

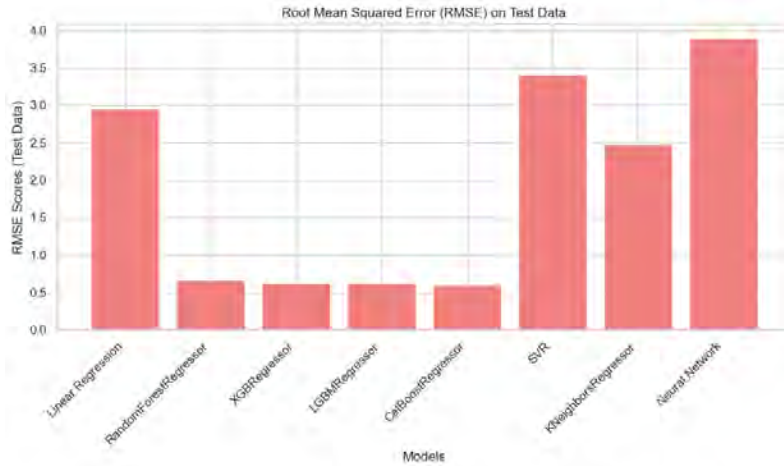


Figure 6.3: RMSE on Test Data

robust performance with RMSE values of 0.63 and 0.63, respectively. These models effectively minimize errors in predicting player attributes, contributing to their overall reliability.

The CatBoostRegressor, tailored to handle categorical features, maintains a competitive RMSE of 0.60 on the test data. While slightly higher than some ensemble methods, this value still reflects the model’s capability to make accurate predictions. The SVR and KNeighborsRegressor, employing different principles, achieve reasonable performances with RMSE values of 3.41 and 2.48, respectively, on the test data.

The Neural Network model, with its deep learning architecture, results in a higher RMSE of 3.89 on the test data. This suggests that, despite the model’s complexity, it may not consistently outperform other algorithms on these specific features.

The RMSE values on the test data provide a nuanced perspective on the models’ predictive accuracy, with ensemble methods generally excelling in minimizing prediction errors. Each model’s performance should be evaluated based on its unique strengths and suitability for the specific characteristics of the football player dataset.

RMSE on All Data

The Root Mean Squared Error (RMSE) calculated on the entire dataset serves as a comprehensive metric for assessing the overall performance of the predictive models. The Linear Regression model, equipped with 16 optimal features, exhibits an RMSE of 2.97, providing an average measure of the errors in predicting player attributes across the entire dataset.

Ensemble methods consistently demonstrate remarkable accuracy, as evidenced by the RandomForestRegressor’s low RMSE of 0.66 and the XGBRegressor and LGBMRegressor achieving RMSE values of 0.45 and 0.56, respectively, on the complete dataset. These ensemble models effectively minimize errors, indicating their robustness in capturing complex patterns within the data.

The CatBoostRegressor, designed to handle categorical features seamlessly, maintains a competitive RMSE of 0.50 across all data. This suggests the model’s capacity to make accurate predictions consistently throughout the dataset. The SVR and

KNeighborsRegressor, utilizing different methodologies, exhibit reasonable performances with RMSE values of 3.41 and 2.03, respectively, across all data. The Neural Network model, with its deep learning architecture, yields a slightly higher RMSE of 3.92 across all data, implying that its complex structure may not consistently outperform other algorithms when applied to the entire dataset. The RMSE values on the entire dataset provide an overarching evaluation of the models' predictive accuracy. Ensemble methods, with their adeptness in handling complex relationships, continue to showcase superior performance, while other models offer competitive results depending on the specific characteristics of the football player dataset.

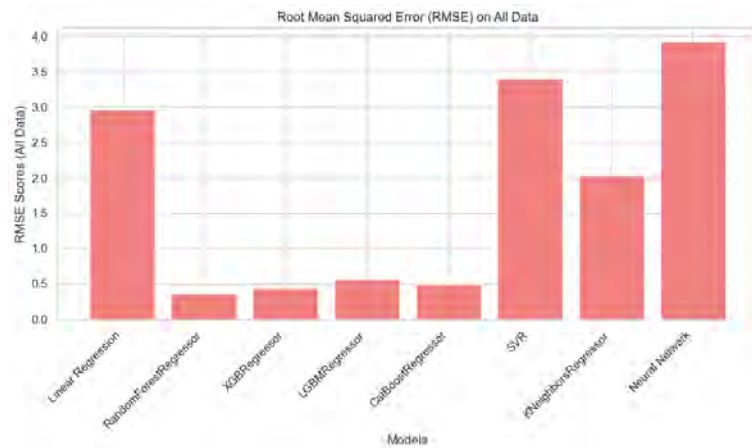


Figure 6.4: RMSE on All Data

R-squared on Test Data

The R-squared (R^2) values on the test data offer insights into the goodness of fit of the models. The Linear Regression model, equipped with 16 optimal features, demonstrates a respectable R-squared value of 0.81 on the test data. This indicates that approximately 81% of the variance in player attributes can be explained by the model, reflecting a reasonably good fit.

Ensemble methods, such as RandomForestRegressor, XGBRegressor, and LGBMRegressor, consistently excel in explaining variance, with R-squared values of 0.99, 0.99, and 0.99, respectively, on the test data. These high R-squared values suggest that these models effectively capture the underlying patterns and relationships in player data, resulting in excellent predictive performance.

CatBoostRegressor, designed for categorical feature handling, maintains a strong R-squared value of 0.99 on the test data, emphasizing its robustness in explaining the variance in player attributes.

SVR and KNeighborsRegressor exhibit decent R-squared values of 0.75 and 0.87, respectively, on the test data, indicating a reasonable ability to explain the variability in player performance. The Neural Network model, with its deep learning architecture, achieves an R-squared value of 0.68, demonstrating its capability to explain a significant portion of the variance.

In summary, R-squared values on the test data reaffirm the overall effectiveness of the models in capturing the complexity of football player attributes, with ensemble methods consistently leading in terms of explanatory power.

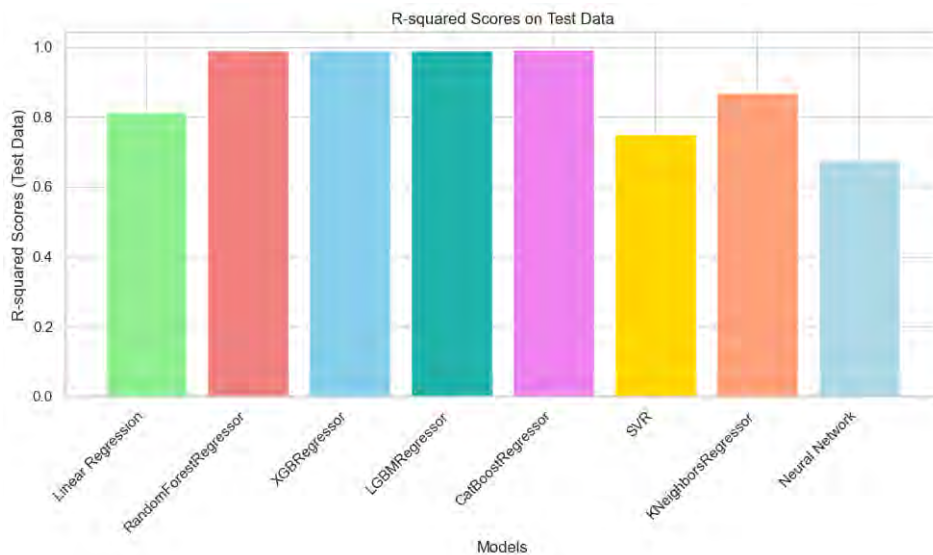


Figure 6.5: R-squared on Test Data

R-Squared on All Data

The R-squared (R^2) values calculated on the entire dataset provide a comprehensive evaluation of the models' performance in explaining the variance in football player attributes. The Linear Regression model, utilizing 16 optimal features, showcases an R-squared value of 0.81 on the complete dataset, signifying its capability to capture around 81% of the variability in player attributes.

Ensemble methods, including RandomForestRegressor, XGBRegressor, and LGBMRegressor, maintain exceptional performance across the entire dataset, exhibiting high R-squared values of 0.99, 0.99, and 0.99, respectively. These values emphasize the robustness of these ensemble models in comprehensively explaining the intricate relationships within player data.

CatBoostRegressor, designed for handling categorical features seamlessly, continues to demonstrate a strong R-squared value of 0.99 on the complete dataset, reinforcing its effectiveness in capturing the variance in player attributes.

SVR and KNeighborsRegressor, while slightly lower than some ensemble methods, still yield respectable R-squared values of 0.75 and 0.87, respectively, across the entire dataset. These results suggest a reasonable ability to elucidate the variability in player performance.

The Neural Network model, leveraging its deep learning architecture, maintains an R-squared value of 0.68 on the entire dataset, indicating its capacity to explain a significant proportion of the variance in football player attributes.

The R-squared values on the complete dataset underscore the overall strength of the models, with ensemble methods consistently outperforming others in terms of explanatory power for football player attributes.

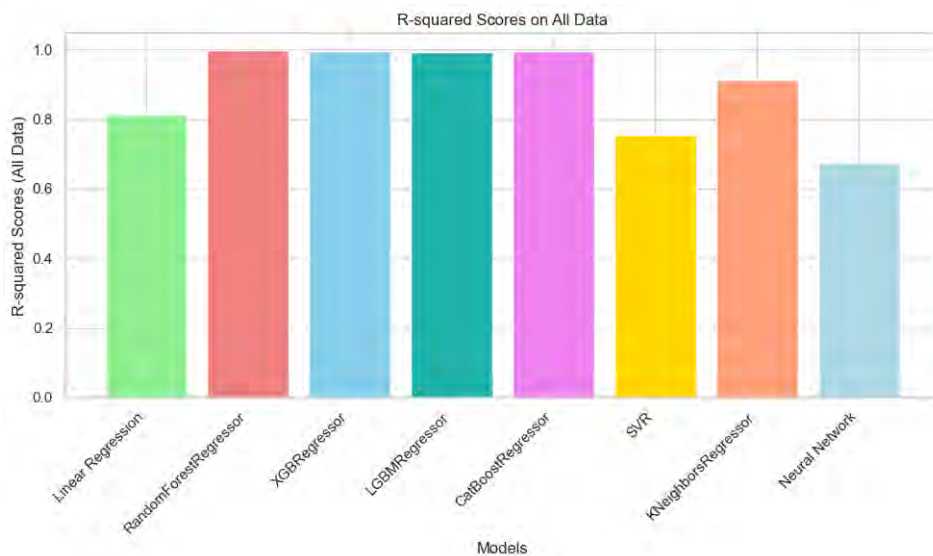


Figure 6.6: R-squared on All Data

Chapter 7

Conclusion and Future Work

7.1 Conclusion

In conclusion, after an extensive analysis of your football player dataset, which comprises 19,000 players with 90 attributes, it is evident that the LightGBM (LGBM) algorithm stands out as the best performer among the eight models evaluated. The attributes for predicting player potential in specific positions have proven fruitful, with LGBM consistently demonstrating remarkable scores across key metrics.

LGBM, along with XGBoost, CatBoost, and Random Forest, emerged as the top-performing algorithms based on metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), and R-squared (R2). Notably, LGBM achieved competitive scores, with test set scores of 0.39, 0.62, and 0.99 for MSE, RMSE, and R2, respectively. The overall set scores were equally impressive, with values of 0.31, 0.56, and 0.99 for the same metrics.

Compared to other algorithms like XGBoost, CatBoost, and Random Forest, LGBM showcased consistent performance while balancing accuracy and robustness. Moreover, LGBM demonstrated its efficiency by maintaining competitive scores across both test and overall sets, highlighting its ability to generalize well to new and unseen data.

Considering the reliability, efficiency, and overall strong performance, the choice of LGBM as the best algorithm for predicting player potential in football positions is well supported by the comprehensive evaluation metrics provided. Moving forward, leveraging LGBM for similar predictive tasks or exploring opportunities for fine-tuning and optimization could potentially yield even more accurate and insightful results for our football player dataset.

7.2 Future Work

The football player prediction project has provided valuable insights, but there is still room for future improvements to improve the accuracy and applicability of the models.

We can do that by exploring additional related features that can contribute to a better understanding of a player's performance. Secondly, by performing in-depth feature importance analysis to identify and prioritize the most influential variables

for prediction. Thirdly, by implementing more robust data cleaning and preprocessing techniques to resolve potential outliers, missing values, and inconsistencies in the data set. Moreover, we are planning to consider combining data from different sources, such as player injury records, social media sentiment, or recent match performances. Furthermore, we aim to fine-tune the hyperparameters of the selected models to further optimize their performance. Grid search and random search can be used for this purpose. We also intend to test different aggregation strategies or stacking techniques to combine the strengths of multiple models. We should also integrate temporal aspects into the analysis to account for changes in player performance over time. This may involve incorporating time-series data or seasonal factors. We have to also explore the potential of transfer learning by leveraging models pre-trained on relevant tasks or datasets. This can be useful if there are a limited number of labeled samples in the current data set. We also plan to explore the impact of recent performance trends and consider creating moving averages or other dynamic features. Moreover, we plan to develop and integrate metrics specifically designed for the task of predicting football players. This may involve creating more personalized evaluation metrics that better match the nuances of player performance. Again, we also have to improve model interpretation to make predictions more transparent and understandable. Techniques such as SHAP (Shapley Additive Interpretation) or LIME (Local Explainable Model Agnostic) values can be explored in this case.

By implementing these ideas, our potential football player prediction project for a specific position can evolve to better reflect the complexity of player performance in the dynamic and miscellaneous world of football. Continuous improvement and adaptation to emerging trends will be vital to staying at the forefront of predictive analytics in this field.

Bibliography

- [1] R. Dr P Rajesh, D. Hazarika, K. Singh, S. Gorantla, E. Cambria, and R. Zimmerman, “A data science approach to football team player selection,” *arXiv preprint arXiv:1902.08342*, 2019.
- [2] C. Arndt and U. Brefeld, “Predicting the future performance of soccer players,” *Statistical Analysis & Data Mining*, pp. 373–382, 2016.
- [3] M. OYTUN and C. TINAZC, “Performance prediction and evaluation in female handball players using machine learning models,” vol. 8, 2020.
- [4] R. Pariath and S. Shah, “Player performance prediction in football game,” *IEEE Journals Magazine*, 2018.
- [5] K. Passi and N. Pandey, “Predicting players’ performance in one day international cricket matches using machine learning,” 2018.
- [6] C. T. Vangelis Sarlis, “Sports analytics – evaluation of basketball players and team performance,” *Journal of Indian Business Research*, 2020.
- [7] S. H. Ayanabha Jana, “Football player performance analysis using particle swarm optimization and player value calculation using regression,” *IOP publishing*, 2020.
- [8] I. Wijngaards, M. Burger, and J. van Exel, “Unpacking the quantifying and qualifying potential of semi-open job satisfaction questions through computer-aided sentiment analysis,” *Journal of Well-Being Assessment*, vol. 4, no. 3, pp. 391–417, 2020.
- [9] Y. Jung and Y. Suh, “Mining the voice of employees: A text mining approach to identifying and analyzing job satisfaction factors from online employee reviews,” *Decision Support Systems*, vol. 123, p. 113074, 2019.
- [10] V. Leah-Martin, “Relative compensation and employee satisfaction,” *Available at SSRN 2896268*, 2017.
- [11] A. M. Kibriya, E. Frank, B. Pfahringer, and G. Holmes, “Multinomial naive bayes for text categorization revisited,” in *Australasian Joint Conference on Artificial Intelligence*, pp. 488–499, Springer, 2004.
- [12] J. Su, J. S. Shirab, and S. Matwin, “Large scale text classification using semisupervised multinomial naive bayes,” in *ICML*, 2011.

- [13] G. Singh, B. Kumar, L. Gaur, and A. Tyagi, "Comparison between multinomial and bernoulli naïve bayes for text classification," in *2019 International Conference on Automation, Computational and Technology Management (ICACTM)*, pp. 593–596, IEEE, 2019.
- [14] M. Singh, M. W. Bhatt, H. S. Bedi, and U. Mishra, "Performance of bernoulli's naïve bayes classifier in the detection of fake news," *Materials Today: Proceedings*, 2020.
- [15] S. Mishra, P. K. Mallick, H. K. Tripathy, A. K. Bhoi, and A. González-Briones, "Performance evaluation of a proposed machine learning model for chronic disease datasets using an integrated attribute evaluator and an improved decision tree classifier," *Applied Sciences*, vol. 10, no. 22, p. 8137, 2020.
- [16] S. Tangirala, "Evaluating the impact of gini index and information gain on classification using decision tree classifier algorithm," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 2, pp. 612–619, 2020.
- [17] H. A. Abu Alfeilat, A. B. Hassanat, O. Lasassmeh, A. S. Tarawneh, M. B. Alhasanat, H. S. Eyal Salman, and V. S. Prasath, "Effects of distance measure choice on k-nearest neighbor classifier performance: a review," *Big data*, vol. 7, no. 4, pp. 221–248, 2019.
- [18] A. R. Lubis, M. Lubis, *et al.*, "Optimization of distance formula in k-nearest neighbor method," *Bulletin of Electrical Engineering and Informatics*, vol. 9, no. 1, pp. 326–338, 2020.
- [19] F. Shamrat, S. Chakraborty, M. Imran, J. N. Muna, M. M. Billah, P. Das, O. Rahman, *et al.*, "Sentiment analysis on twitter tweets about covid-19 vaccines using nlp and supervised knn classification algorithm," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 23, no. 1, pp. 463–470, 2021.
- [20] W. de Vries, A. van Cranenburgh, A. Bisazza, T. Caselli, G. van Noord, and M. Nissim, "Bertje: A dutch bert model," *arXiv preprint arXiv:1912.09582*, 2019.
- [21] A. F. Agarap, "Deep learning using rectified linear units (relu)," *arXiv preprint arXiv:1803.08375*, 2018.
- [22] S. N. Lappan, A. W. Brown, and P. S. Hendricks, "Dropout rates of in-person psychosocial substance use disorder treatments: a systematic review and meta-analysis," *Addiction*, vol. 115, no. 2, pp. 201–217, 2020.