

EXPLAINABLE ARTIFICIAL INTELLIGENCE MODEL FOR IDENTIFYING MARKET VALUE IN PROFESSIONAL SOCCER PLAYERS

A PREPRINT

 **Chunyang Huang***

Zhili College
Tsinghua University
Beijing 100084, China
cy-huang20@mails.tsinghua.edu.cn

 **Shaoliang Zhang**

Research Centre for Athletic Performance and Data Science & Division of Sports Science and Physical Education
Tsinghua University
Beijing 100084, China
zslief@mail.tinghua.edu.cn

November 27, 2023

ABSTRACT

This study introduces an advanced machine learning method for predicting soccer players' market values, combining ensemble models and the Shapley Additive Explanations (SHAP) for interpretability. Utilizing data from about 12,000 players from Sofa, the Boruta algorithm streamlined feature selection. The Gradient Boosting Decision Tree (GBDT) model excelled in predictive accuracy, with an R-squared of 0.901 and a Root Mean Squared Error (RMSE) of 3,221,632.175. Player attributes in skills, fitness, and cognitive areas significantly influenced market value. These insights aid sports industry stakeholders in player valuation. However, the study has limitations, like underestimating superstar players' values and needing larger datasets. Future research directions include enhancing the model's applicability and exploring value prediction in various contexts.

Keywords Market Value, Feature Selection, Explainable Machine Learning Models

1 Introduction

Soccer, often referred to as 'the beautiful game', holds preeminence as a global sport, captivating a diverse audience across various cultural and geographical landscapes. Its allure transcends the physical boundaries of the playing field, significantly bolstering a multi-billion-dollar economic industry. This industry is characterized by varied revenue streams, including broadcasting rights, ticket sales, merchandising, sponsorships, and notably, the transfer market. The latter is subject to extensive financial scrutiny, underscoring its economic impact [Dobson, 2001].

The process of estimating a player's market value is a critical aspect of the economic operations within football, shaping the sport's financial landscape. These valuations are crucial during transfer negotiations, reflecting a player's potential impact on a team's performance and economic success. High-profile transfers, involving substantial financial investments, not only affect a club's prestige but also its appeal, making accurate player valuation essential for maximizing financial returns [Wolfers and Zitzewitz, 2004]. Furthermore, a player's market value influences wage

*Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

policies within clubs, with top-tier players commanding higher salaries. This, in turn, affects the financial stability and budgetary planning of the clubs. Additionally, the collective market value of players plays a pivotal role in a club's overall valuation, particularly during ownership transitions or financial negotiations, highlighting the significance of precise player valuation in the strategic fiscal management of football clubs.

In the contemporary era, dominated by 'big data', analytical methodologies have become indispensable in sports, particularly in the context of player market valuation. These methodologies inform decision-making processes related to player recruitment and selection. The emergence of sophisticated online platforms such as Sofifa, WhoScored, and Transfermarkt, offering comprehensive player data, has transformed this area. The integration of machine learning algorithms with these data sources aims to enhance the precision of player market valuations [Baboota and Kaur, 2019]. For instance, Mustafa A. AL-ASADI and Sakir Tasdemir [2022] analyzed FIFA 20 video game data from Sofifa.com, employing four regression models—linear regression, multiple linear regression, decision trees, and random forests—to evaluate players' market values. This study demonstrated that the random forest model outperformed traditional statistical models in predicting players' market values, with players' potential, international reputation, age, and height being key determinants of individual market value. Similarly, McHale and Holmes [2023] analyzed transfer details and crowd-sourced player ratings from transfermarkt.com and sofifa.com over nine seasons. The addition of advanced player rating systems from sofifa.com, beyond player performance profiles, significantly improved the predictive accuracy of the models, with XGBoost showing a substantial increase in predictive accuracy compared to linear regression models. Moreover, Yang et al. [2022] compiled a longitudinal transfer dataset from the top five European leagues over 14 seasons from transfermarkt.de. The study employed generalized, quantile additive models, and random forests to predict players' market values, concluding that machine learning models, particularly random forests, are superior in evaluating market values. The study also identified non-linear predictors of transfer fees, such as buying-club expenditure and selling-club income, as having significant impacts on transfer fees, overshadowing players' anthropometric characteristics and technical performance. Despite the advancements in machine learning techniques that improve predictive accuracy, the challenge of attaining clear interpretability among various models at both the global and local scales persists.

In light of these developments, this study adopts ensemble machine learning models, concentrating on the most important factors influencing athletic performance. Utilizing the SHapley Additive exPlanations (SHAP) method, it achieves transparent and comprehensive interpretive visualizations from both local and global perspectives. These insights identify the primary features affecting athlete performance, offering a detailed quantitative analysis of individual player metrics that contribute to their market value.

2 Methods

2.1 Sample

Building upon the findings of previous research, this study conducts an in-depth analysis of the dataset available on Sofifa.com, a widely recognized resource among FIFA football manager enthusiasts. The website offers extensive data, including but not limited to, nuanced player ratings, team compositions, and a variety of other statistics pertinent to the game. Further, it provides detailed information on aspects such as player positions and their preferred playing foot.

In the course of our research, data was systematically extracted relating to approximately 12,000 players from Sofifa.com, as recorded on January 5, 2023. This dataset is comprehensive, encompassing a multitude of attributes including players' names, market values, wages, overall ratings, potential, and an additional 34 features. Among these, 29 features are applicable to outfield players, while 5 are uniquely relevant to goalkeepers. The table 2 in Appendix A details these features, offering a thorough overview essential for our analysis

In the initial phase of our methodology, the dataset underwent a rigorous cleansing process to rectify any instances of missing values and to categorize the data into two primary classifications: outfield players and goalkeepers. The dataset exhibited a wide range of player values, stretching from €15,000 to €190 million. Figure 1 a graphically represents the distribution of these player values.

Upon examining the distribution in Figure 1 a, a notable aggregation near the lower end of the value spectrum was observed. This pattern suggests that a limited number of high-value players disproportionately influence the extension of the horizontal axis in the graph. This distributional characteristic aligns with the concept of player popularity and the 'superstar' effect, as discussed in seminal works by Adler [1985], Franck and Nüesch [2012] and Rosen [1981]. However, these elements are extraneous to the performance-centric focus of our analysis, as delineated by Müller et al. [2017]. Consequently, data points representing player values exceeding €25 million, constituting approximately 3% of the dataset, were excluded, as illustrated in Figure 1b.

The revised distribution, as depicted in Figure 1b, revealed a significant skewness towards the lower value range. To rectify this skewness and foster a more symmetrical distribution conducive to effective statistical modeling, we employed

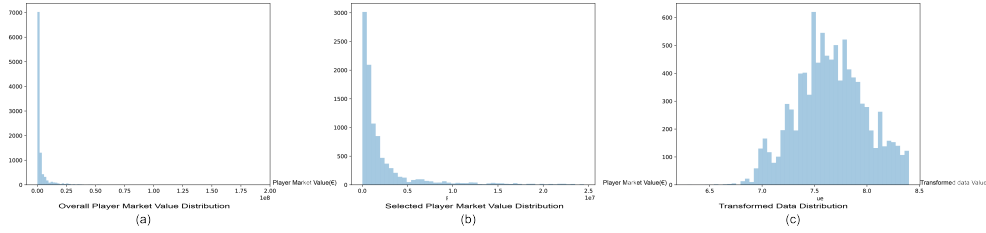


Figure 1: The distribution of origin value(1a), deducted value(1b) and transformed data(1c)

the Box-Cox transformation, a technique well-established in statistical literature [Box and Cox, 1964, Osborne, 2019]. This transformation was implemented utilizing the Scipy library within the Python programming environment [Virtanen et al., 2020]. The outcome of this transformation, which significantly improved the symmetry of the data distribution, is displayed in Figure 1c.

2.2 Feature selection

The initial step in our analytical process involves a critical task: feature selection. Our dataset includes up to 29 attributes related to football players' performance. Within this set, certain features may contribute minimally to the predictive accuracy of our model, presenting several challenges. Firstly, an excess of features can lead to increased computational time and escalate the computational power requirements. Secondly, the elimination of non-essential features, as suggested by Kohavi and John [1997], is likely to enhance the precision of our machine learning model.

To address these concerns, this study employs the Boruta algorithm for feature selection. This algorithm offers a comprehensive solution to issues extending beyond minimal-optimal feature sets [Daoud, 2017, Kursu and Rudnicki, 2010, Nilsson et al., 2007]. The Boruta algorithm, a wrapper built around the random forest classification algorithm, aims to identify all potentially significant features in the dataset, including those that might otherwise be overlooked. It operates through an iterative process, comparing the importance of original features with shadow features - which are randomized permutations of the original set - to ascertain their relative significance.

2.3 Model selection

In the pursuit of developing an optimal predictive model for assessing the market value of football players, this study rigorously evaluates a variety of learning algorithms. The algorithms chosen for this analysis include Adaboost [Freund and Schapire, 1997], LightGBM [Ke et al., 2017], [Ho, 1995], Gradient Boosting Decision Tree (GBDT) [Friedman, 2001], CatBoost (Dorogush et al., 2018) [Dorogush et al., 2018], and XGBoost [Chen and Guestrin, 2016]. These were selected for their relevance and potential efficacy in forming the foundational structure of the models.

The methodological approach centers on ensemble learning algorithms, a class of meta-algorithms that synthesize the methodologies of individual models to construct a comprehensive predictive framework [Webb and Zheng, 2004]. The rationale behind selecting ensemble modeling lies in its multiple benefits, which include diminishing variance and bias while concurrently enhancing the overall predictive accuracy of the model [Sagi and Rokach, 2018]. Ensemble learning, by amalgamating predictions from multiple models, inherently demonstrates a capacity for superior performance in comparison to singular model approaches.

2.4 Model development and tuning

In the construction of the predictive model, the dataset underwent a randomized splitting process: 80% was allocated for training and model validation, while the remaining 20% was designated for testing purposes. To maintain the analytical integrity and avert any potential data leakage, a rigorous methodology was employed. This method entailed conducting both imputation and feature selection exclusively on the training set prior to undertaking any evaluative measures. Such a procedural approach was integral in ensuring that the test dataset, reserved solely for the ultimate evaluation of classifier efficacy, remained devoid of any influences that could introduce bias into the model's performance assessment.

To optimize the hyperparameters of each ensemble learning model, we employed a 5-fold cross-validation technique coupled with Grid Search. This rigorous process allowed us to effectively tune the models and obtain reliable and unbiased performance estimates [Yarkoni and Westfall, 2017]. The training dataset was initially partitioned into ten distinct subsets of equal size through random assignment. Out of these subsets, nine were used for training the

model, while the remaining subset served the purpose of validation. This iterative procedure was repeated for all ten feasible combinations. Ultimately, a prognostic model was constructed by leveraging discerning features and refining hyperparameters for optimal performance [Raschka, 2018].

2.5 Model evaluation

The evaluation of various machine learning algorithms constitutes a critical component of this study. In this context, a comprehensive assessment was conducted utilizing multiple metrics to gauge the accuracy of predictions pertaining to player market value. Key among these metrics are the R-squared value, which quantifies the proportion of variance in the dependent variable that is predictable from the independent variables, and the Root Mean Squared Error (RMSE), which provides a measure of the average magnitude of the prediction errors. By employing both R-squared and RMSE, we aim to offer a comprehensive and multifaceted evaluation of our model's performance.

2.6 Model interpretation

In the realm of machine learning, models are often characterized as 'black boxes,' presenting challenges in interpreting the underlying mechanisms driving their predictive accuracy, particularly in the context of evaluating player market values [Linardatos et al., 2021]. To surmount these interpretability hurdles, Lundberg and Lee introduced the Shapley Additive exPlanations (SHAP) approach. This method employs Shapley values, a well-established metric for assessing feature importance, to elucidate the outputs of any machine learning model. SHAP facilitates both global and local interpretability, enabling an analysis of how individual input characteristics contribute positively or negatively to the predictive outcomes.

For global interpretability, this study utilizes the SHAP beeswarm plot and feature importance measures. The beeswarm plot ranks features based on their influence within the predictive model, with the y-axis representing the features and the x-axis indicating the respective SHAP values. Each feature is depicted by rows of colored dots, where red indicates high values and blue denotes low values, thus providing a clear visual representation of the impact each feature has on the model's output [Rommers et al., 2020].

In the pursuit of local interpretability, the SHAP force plot was implemented to specifically elucidate the predicted market value of individual players. This analytical tool clarifies the contribution of each feature to the prediction for a given sample, graphically illustrating the trajectory from the base value (average predicted value) to the final predicted outcome [Hiabu et al., 2023]. Features that positively influence the prediction, leading to higher values, are represented in red, while those that negatively impact the prediction, resulting in lower values, are shown in blue. This visualization offers a granular perspective of the model's predictions, underscoring the significance of each feature in ascertaining the market value of the players under consideration.

Additionally, to assess the SHAP value matrix for each feature, the Partial Dependence Plot (PDP) was generated. The PDP elucidates the marginal effect of a specific feature on the predicted outcome by averaging out the effects of all other features. It essentially portrays the average prediction for varying values of the targeted feature, thereby isolating its influence from other contributing factors.

3 Result

The methodological framework of this study is depicted in Figure 2, illustrating the comprehensive study design inclusive of data collection, feature screening, model development and validation, and model evaluation and interpretation.

3.1 Feature selection

In the initial phase of feature selection, we commenced with a set of 29 features. Utilization of the Boruta algorithm, implemented through the BorutaShap package in Python, allowed for a reduction in the feature set to 22, as indicated by the green bars in Figure 3. Acceleration, Heading accuracy, Defensive awareness, Vision, Volleys, Sprint speed, Long passing, Positioning, Standing tackle, Dribbling, FK accuracy, Short passing, Interceptions, Penalties, Finishing, Reactions, Ball control, Stamina, Crossing, Strength, Shot power, Sliding tackle.

3.2 Model development and evaluation

The outcomes of the cross-validation analysis and evaluation on test set are detailed in Table 1. Within the ensemble of six learning algorithms evaluated, the Gradient Boosting Decision Tree (GBDT) model demonstrated superior performance, achieving the highest R-Squared value of 0.889. The CatBoost model was a close second with an

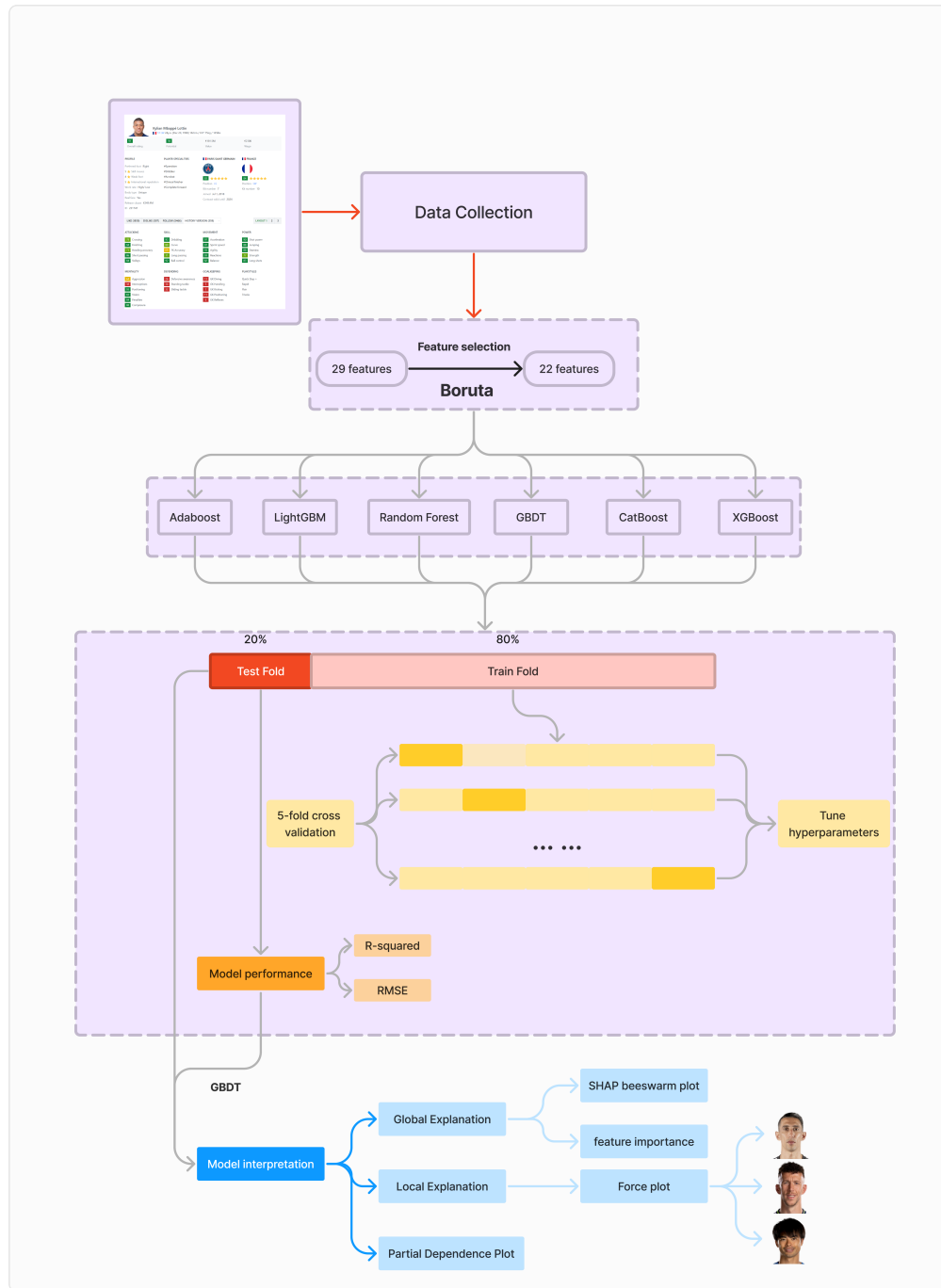


Figure 2: Flow chart of the overall study design with data collection, feature screening, model development and validation, and model evaluation and interpretation

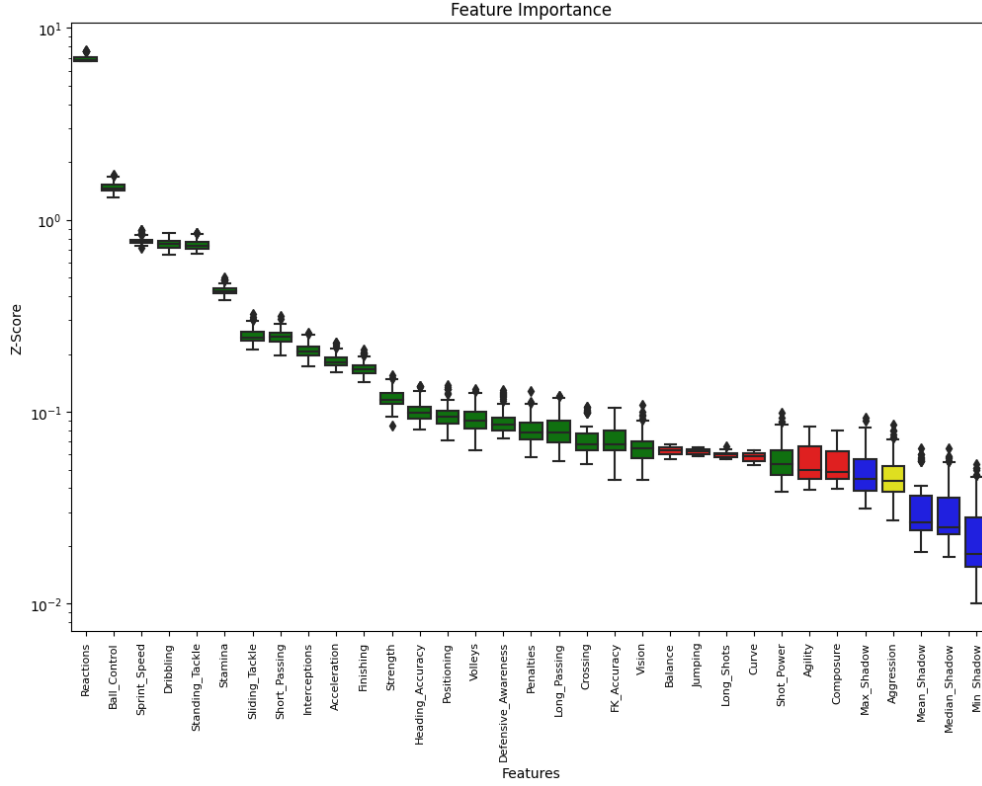


Figure 3: Feature selection using the Boruta algorithm for key features

R-Squared of 0.887, followed by LightGBM with 0.885. The Random Forest and XGBoost models obtained R-Squared values of 0.877 and 0.861, respectively, while the AdaBoost model had the lowest R-Squared at 0.773. Regarding the Root Mean Squared Error (RMSE), the GBDT model surpassed all counterparts, recording an RMSE of 3,060,228.569. The subsequent models, CatBoost (4,715,039.662), LightGBM (3,249,280.179), Random Forest (3505068.8371), XGBoost (3,320,149.832), and AdaBoost (4,442,839.041), followed in ascending order of RMSE values. Notably, within the test set, the GBDT model sustained its advantage, achieving the highest R-squared value of 0.901 and the lowest RMSE of 3,221,632.175, indicative of its robustness in predicting player market values.

Table 1: The cross-validation analysis and evaluation on test set

Model Name	Cross Validation		Test Set	
	Mean R ²	Mean RMSE	Mean R ²	Mean RMSE
AdaBoost	0.764	4442839.041	0.752	4657341.125
GBDT	0.878	3221632.175	0.901	3221632.175
LightGBM	0.877	3249280.179	0.886	3157342.577
Random Forest	0.856	3505068.837	0.856	3547133.366
XGBoost	0.871	3320149.832	0.888	3127608.099
CatBoost	0.742	4715039.662	0.742	4715039.662

3.3 Global and local interpretation of the ML Model

In the study, the SHAP (Shapley Additive Explanations) beeswarm plot and feature importance for the Gradient Boosted Decision Trees (GBDT) model, as illustrated in Figure 4, were employed to identify the features with the most significant influence on the prediction model. The analysis revealed that nine variables—Ball control, Reactions, Short passing, Sprint speed, Finishing, Interceptions, Dribbling, Sliding Tackle, and Acceleration—held the most substantial predictive power, significantly impacting a player’s market value.

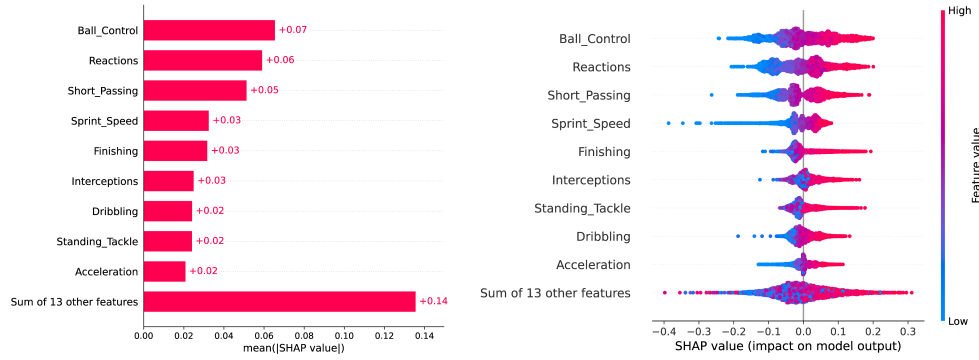


Figure 4: The GBDT model’s interpretation. The importance ranking of the different variables according to the mean (|SHAP value|) (a); The importance ranking of different risk factors with stability and interpretation using the optimal model (b). The higher SHAP value of a feature is given, the higher market value of players would have. The red part in feature value represents higher value.

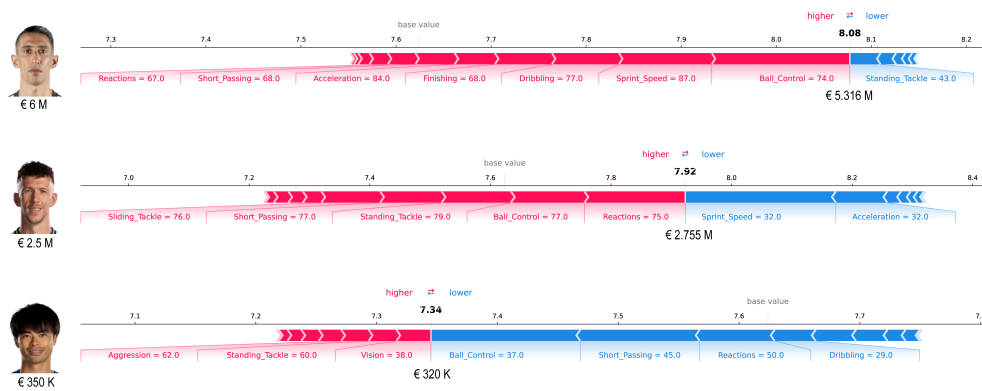


Figure 5: SHAP force plot for predicting players’ market value. (A) SHAP forces plot to correctly predict market value of Ángel Fabián. (B) SHAP forces plot to correctly predict market value of Ivan Perišić (C) SHAP force plot correctly predict market value of Teruki Miyamoto. Red features mean pushing the prediction higher market value and blue features pushing the prediction lower market value.

This detailed investigation into the predicted market values for players such as Ángel Fabián, Ivan Perišić, and Teruki Miyamoto, as illustrated in Figure 5, highlights the practical applicability and accuracy of the Gradient Boosted Decision Trees model in real-world scenarios. For the first player, Ángel Fabián, the predicted market value, post Box-Cox transformation, was approximately €6 Million. This prediction, derived from a transformed value of 8.08, closely aligns with his actual market value of €5.31 million. In the case of Ivan Perišić, the predicted market value following the Box-Cox transformation was approximately €2.5 Million, based on a transformed prediction of 7.92. This estimate is near the real market value of €2.75 million. Finally, for Teruki Miyamoto, the model estimated a market value of approximately €3.20 Million after the transformation, calculated from a transformed prediction of 7.34, which is in close proximity to the actual market value of €3.50 million.

3.4 Partial Dependence Plot

The Partial Dependence Plots (PDP) for these features are presented in Figure 6, offering a detailed examination of their marginal effects on the predicted outcomes. It was observed that the SHAP values of Ball Control, Reactions, and Sprint Speed generally exhibited an increasing trend in correspondence with the escalation of the respective market values. This trend underscores a direct relationship between these specific player attributes and their market valuation, thereby highlighting their critical importance in the predictive model.

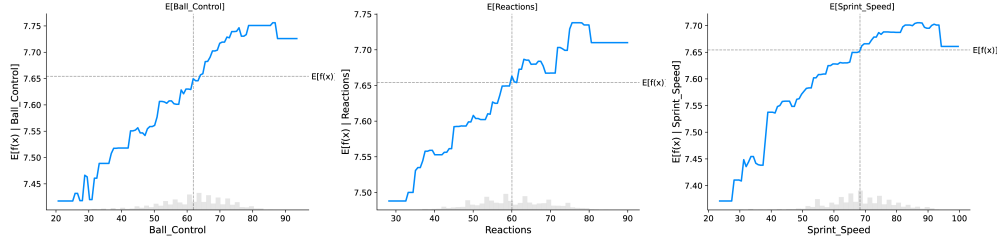


Figure 6: partial dependence plots of top features in skill, cognitive, and fitness dimensions

4 Discussion

The aim of this study is to ensemble machine learning models, concentrating on the most important factors influencing athletic performance. Utilizing the SHapley Additive exPlanations (SHAP) method, it achieves transparent and comprehensive interpretive visualizations from both local and global perspectives. In previous research, methodologies such as static statistical methods, linear regression, and conventional machine learning approaches have demonstrated limited effectiveness in achieving outstanding predictive performance. Moreover, these methods often fall short in providing in-depth insights into each feature’s contribution, thereby limiting the explainability of the prediction model. In contrast, the present study, utilizing data scrupulously gathered from the Sofifa website, employs an ensemble machine learning model. This model not only achieves exceptional accuracy in predicting market values but also leverages Shapley Additive Explanations (SHAP) to gain both global and local insights into all features initially selected by the Boruta method. The study successfully identifies key features that are crucial based on their global explanation and their response tendencies to varying feature values.

The evaluation of feature influences is of paramount importance in the estimation of players’ market and transfer values. Traditionally, club managers have considered three key dimensions in assessing players and their market worth: skills, fitness, and cognition. Our study reveals distinct insights within these dimensions. Specifically, in the skills dimension, Ball Control, Short Passing, Finishing, Interceptions, Dribbling, and Tackling emerge as the most influential factors. In the fitness dimension, Sprint Speed and Acceleration are identified as having the greatest impact. Notably, within the cognition dimension, Reactions are found to hold the most significant position. These findings regarding the impact of various features can guide managers to focus on the most influential aspects during the player valuation process. By understanding the relative importance of these attributes, club managers can make more informed decisions, leading to more accurate and effective player assessments and strategic planning in the realm of player acquisitions and transfers.

The exemplary performance of our model in evaluation phases suggests that the predicted market values closely align with actual figures, thereby validating the effectiveness of our prediction approach. However, it is crucial to highlight that this enhanced prediction performance is achieved at the cost of some interpretability of the model output. In our methodology, we employed the Box-Cox transformation to convert the original exponential distribution into a single-peak distribution. This transformation involved converting the actual market values into transformed data, which facilitated an increase in prediction accuracy. Consequently, the direct output of our model is in the form of predicted transformed data. This necessitates an inverse Box-Cox transformation to revert these predictions back to the estimated market values. Furthermore, the SHAP values for all features and their summary, in relation to a baseline of predicted transformed data, are uniformly expressed in the same unit as the transformed data. This uniformity is particularly evident in local explanations where the baseline, alongside the feature responses in the force plot of an example player, aligns with the unit of the transformed data. This consistency in unit representation across different model outputs ensures a coherent understanding and interpretation of the prediction results, albeit with some complexity in the transformation process.

Our study, while providing valuable insights, is subject to several limitations. Firstly, our current model is unable to accurately estimate the market value of superstar players, as this requires consideration of a broader range of societal and social factors that extend beyond the scope of our study and the data utilized. Consequently, future research should aim to include data pertaining to these additional factors and employ cross-disciplinary analytical methods to enhance the accuracy of superstar valuation.

Secondly, our data, sourced from public websites, encompasses a variety of features across different dimensions. However, this data set pales in comparison to the more comprehensive data captured by commercial providers, which can include over two hundred metrics per player per game. This disparity undoubtedly places limitations on the depth and breadth of insights our study can provide. In light of these limitations, future research endeavors should aim to apply our analytical pipeline to more robust datasets or to different problems at varying levels of analysis. These levels

may include individual player assessments, team and league evaluations, and explorations in other fields where the prediction of value and the influence of various features are pertinent. Such expansions would not only enhance the applicability of our current model but also provide a more comprehensive understanding of value prediction in diverse contexts.

5 Conclusion

In conclusion, this study adeptly integrates advanced machine learning techniques and feature importance analysis, providing a nuanced understanding of player market value prediction within the realm of sports analytics. Our findings indicate that the Gradient Boosting Decision Tree (GBDT) model demonstrates a distinct advantage over other artificial intelligence algorithms in this context. The study identified key features that significantly influence player valuation, categorized into three dimensions: skill (encompassing Ball Control, Short Passing, Finishing, Interceptions, Dribbling, and Standing Tackle), fitness (including Sprint Speed and Acceleration), and cognitive (specifically, Reactions). The quantification of the impact of these features offers valuable insights for player evaluation, enhancing the precision and efficacy of market value assessments.

Looking ahead, there is considerable scope for future research to broaden this framework. By encompassing a more diverse range of features and applying the model in varied contexts, subsequent studies could achieve greater applicability and yield deeper insights. This expansion would not only reinforce the findings of the current study but also contribute to the ongoing advancement of machine learning applications in sports analytics and player valuation.

References

- Stephen Dobson. *The economics of football*, volume 10. 2001.
- Justin Wolfers and Eric Zitzewitz. Prediction Markets. *Journal of Economic Perspectives*, 18(2):107–126, June 2004. ISSN 0895-3309. doi:10.1257/0895330041371321. URL <https://www.aeaweb.org/articles?id=10.1257/0895330041371321>.
- Rahul Baboota and Harleen Kaur. Predictive analysis and modelling football results using machine learning approach for English Premier League. *International Journal of Forecasting*, 35(2):741–755, April 2019. ISSN 0169-2070. doi:10.1016/j.ijforecast.2018.01.003. URL <https://www.sciencedirect.com/science/article/pii/S0169207018300116>.
- Mustafa A. AL-ASADI and Sakir Tasdemir. Predict the Value of Football Players Using FIFA video game data and Machine Learning Techniques. *IEEE Access*, pages 1–1, January 2022. doi:10.1109/access.2022.3154767. MAG ID: 4214589397 S2ID: b85f0efabc3fdce4bd997449f22eefe56e50b319.
- Ian G. McHale and Benjamin Holmes. Estimating transfer fees of professional footballers using advanced performance metrics and machine learning. *European Journal of Operational Research*, 306(1):389–399, April 2023. ISSN 03772217. doi:10.1016/j.ejor.2022.06.033. URL <https://linkinghub.elsevier.com/retrieve/pii/S0377221722005082>.
- Yanxiang Yang, Joerg Koenigstorfer, and Tim Pawlowski. Predicting transfer fees in professional European football before and during COVID-19 using machine learning. *European Sport Management Quarterly*, pages 1–21, December 2022. ISSN 1618-4742, 1746-031X. doi:10.1080/16184742.2022.2153898. URL <https://www.tandfonline.com/doi/full/10.1080/16184742.2022.2153898>.
- Moshe Adler. Stardom and Talent. *The American Economic Review*, 75(1):208–212, 1985. ISSN 0002-8282. URL <https://www.jstor.org/stable/1812714>. Publisher: American Economic Association.
- Egon Franck and Stephan Nüesch. Talent and/or Popularity: What Does It Take to Be a Superstar? *Economic Inquiry*, 50(1):202–216, 2012. ISSN 1465-7295. doi:10.1111/j.1465-7295.2010.00360.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1465-7295.2010.00360.x>. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1465-7295.2010.00360.x>.
- Sherwin Rosen. The Economics of Superstars. *The American Economic Review*, 71(5):845–858, 1981. ISSN 0002-8282. URL <https://www.jstor.org/stable/1803469>. Publisher: American Economic Association.
- Oliver Müller, Alexander Simons, and Markus Weinmann. Beyond crowd judgments: Data-driven estimation of market value in association football. *European Journal of Operational Research*, 263(2):611–624, December 2017. ISSN 03772217. doi:10.1016/j.ejor.2017.05.005. URL <https://linkinghub.elsevier.com/retrieve/pii/S0377221717304332>.

- G. E. P. Box and D. R. Cox. An Analysis of Transformations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 26(2):211–252, 1964. ISSN 0035-9246. URL <https://www.jstor.org/stable/2984418>. Publisher: [Royal Statistical Society, Wiley].
- Jason Osborne. Improving your data transformations: Applying the Box-Cox transformation. *Practical Assessment, Research, and Evaluation*, 15(1), November 2019. ISSN 1531-7714. doi:<https://doi.org/10.7275/qbpc-gk17>. URL <https://scholarworks.umass.edu/pare/vol15/iss1/12>.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, Aditya Vijaykumar, Alessandro Pietro Bardelli, Alex Rothberg, Andreas Hilboll, Andreas Kloeckner, Anthony Scopatz, Antony Lee, Ariel Rokem, C. Nathan Woods, Chad Fulton, Charles Masson, Christian Häggström, Clark Fitzgerald, David A. Nicholson, David R. Hagen, Dmitrii V. Pasechnik, Emanuele Olivetti, Eric Martin, Eric Wieser, Fabrice Silva, Felix Lenders, Florian Wilhelm, G. Young, Gavin A. Price, Gert-Ludwig Ingold, Gregory E. Allen, Gregory R. Lee, Hervé Audren, Irvin Probst, Jörg P. Dietrich, Jacob Silterra, James T. Webber, Janko Slavič, Joel Nothman, Johannes Buchner, Johannes Kulick, Johannes L. Schönberger, José Vinícius de Miranda Cardoso, Joscha Reimer, Joseph Harrington, Juan Luis Cano Rodríguez, Juan Nunez-Iglesias, Justin Kuczynski, Kevin Tritz, Martin Thoma, Matthew Newville, Matthias Kümmerer, Maximilian Bolingbroke, Michael Tartre, Mikhail Pak, Nathaniel J. Smith, Nikolai Nowaczyk, Nikolay Shebanov, Oleksandr Pavlyk, Per A. Brodtkorb, Perry Lee, Robert T. McGibbon, Roman Feldbauer, Sam Lewis, Sam Tygier, Scott Sievert, Sebastiano Vigna, Stefan Peterson, Surhud More, Tadeusz Pudlik, Takuya Oshima, Thomas J. Pingel, Thomas P. Robitaille, Thomas Spura, Thouis R. Jones, Tim Cera, Tim Leslie, Tiziano Zito, Tom Krauss, Utkarsh Upadhyay, Yaroslav O. Halchenko, Yoshiki Vázquez-Baeza, and SciPy 1.0 Contributors. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods*, 17(3):261–272, March 2020. ISSN 1548-7105. doi:10.1038/s41592-019-0686-2. URL <https://doi.org/10.1038/s41592-019-0686-2>.
- Ron Kohavi and George H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97(1):273–324, December 1997. ISSN 0004-3702. doi:10.1016/S0004-3702(97)00043-X. URL <https://www.sciencedirect.com/science/article/pii/S000437029700043X>.
- Jamal I. Daoud. Multicollinearity and Regression Analysis. *Journal of Physics: Conference Series*, 949(1):012009, December 2017. ISSN 1742-6596. doi:10.1088/1742-6596/949/1/012009. URL <https://dx.doi.org/10.1088/1742-6596/949/1/012009>. Publisher: IOP Publishing.
- Miron B. Kurşa and Witold R. Rudnicki. Feature Selection with the Boruta Package. *Journal of Statistical Software*, 36:1–13, September 2010. ISSN 1548-7660. doi:10.18637/jss.v036.i11. URL <https://doi.org/10.18637/jss.v036.i11>.
- Roland Nilsson, José M. Peña, Johan Björkegren, and Jesper Tegnér. Consistent Feature Selection for Pattern Recognition in Polynomial Time. *Journal of Machine Learning Research*, 8(21):589–612, 2007. ISSN 1533-7928. URL <http://jmlr.org/papers/v8/nilsson07a.html>.
- Yoav Freund and Robert E Schapire. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. *Journal of Computer and System Sciences*, 55(1):119–139, August 1997. ISSN 0022-0000. doi:10.1006/jcss.1997.1504. URL <https://www.sciencedirect.com/science/article/pii/S002200009791504X>.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/6449f44a102fde848669bdd9eb6b76fa-Abstract.html>.
- Tin Kam Ho. Random decision forests. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, volume 1, pages 278–282 vol.1, August 1995. doi:10.1109/ICDAR.1995.598994. URL <https://ieeexplore.ieee.org/document/598994>.
- Jerome H. Friedman. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics*, 29(5): 1189–1232, 2001. ISSN 0090-5364. URL <https://www.jstor.org/stable/2699986>. Publisher: Institute of Mathematical Statistics.
- Anna Veronika Dorogush, Vasily Ershov, and Andrey Gulin. CatBoost: gradient boosting with categorical features support, October 2018. URL <http://arxiv.org/abs/1810.11363>. arXiv:1810.11363 [cs, stat].

- Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, pages 785–794, New York, NY, USA, August 2016. Association for Computing Machinery. ISBN 978-1-4503-4232-2. doi:10.1145/2939672.2939785. URL <https://dl.acm.org/doi/10.1145/2939672.2939785>.
- G.I. Webb and Z. Zheng. Multistrategy ensemble learning: reducing error by combining ensemble learning techniques. *IEEE Transactions on Knowledge and Data Engineering*, 16(8):980–991, August 2004. ISSN 1558-2191. doi:10.1109/TKDE.2004.29. Conference Name: IEEE Transactions on Knowledge and Data Engineering.
- Omer Sagi and Lior Rokach. Ensemble learning: A survey. *WIREs Data Mining and Knowledge Discovery*, 8(4): e1249, 2018. ISSN 1942-4795. doi:10.1002/widm.1249. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/widm.1249>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/widm.1249>.
- Tal Yarkoni and Jacob Westfall. Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspectives on Psychological Science*, 12(6):1100–1122, 2017. Publisher: Sage Publications Sage CA: Los Angeles, CA.
- Sebastian Raschka. Model evaluation, model selection, and algorithm selection in machine learning. *arXiv preprint arXiv:1811.12808*, 2018.
- Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable AI: A Review of Machine Learning Interpretability Methods. *Entropy*, 23(1):18, January 2021. ISSN 1099-4300. doi:10.3390/e23010018. URL <https://www.mdpi.com/1099-4300/23/1/18>. Number: 1 Publisher: Multidisciplinary Digital Publishing Institute.
- Nikki Rommers, Roland Rössler, Evert Verhagen, Florian Vandecasteele, Steven Verstockt, Roel Vaeyens, Matthieu Lenoir, Eva D’hondt, and Erik Witvrouw. A Machine Learning Approach to Assess Injury Risk in Elite Youth Football Players. *Medicine & Science in Sports & Exercise*, 52(8):1745, August 2020. ISSN 0195-9131. doi:10.1249/MSS.0000000000002305. URL https://journals.lww.com/acsm-msse/fulltext/2020/08000/a_machine_learning_approach_to_assess_injury_risk.12.aspx.
- Munir Hiabu, Joseph T. Meyer, and Marvin N. Wright. Unifying local and global model explanations by functional decomposition of low dimensional structures, February 2023. URL <http://arxiv.org/abs/2208.06151>. arXiv:2208.06151 [cs, math, stat] version: 2.

Appendices

A Feature description

Table 2: Feature description

feature name	Description	data type
name		string
overall rating		integer
potential		integer
value	The estimated market value of a player	integer
wage		integer
Crossing	the ability to deliver accurate crosses	integer
Finishing	the ability to score goals from various positions and situations.	integer
Heading_Accuracy	the ability to accurately direct headers towards the goal or teammates	integer
Short_Passing	the ability to accurately perform short passes to teammates	integer
Volleys	the skill in striking the ball out of the air, often when attempting to score	integer
Dribbling	the ability to control the ball and move past opponents.	integer
Curve	the ability to apply spin and curve to passes and shots.	integer
FK_Accuracy	the skill in taking accurate free kicks.	integer
Long_Passing	the ability to accurately perform long passes to teammates.	integer
Ball_Control	the skill in controlling the ball when receiving it or when dribbling.	integer
Acceleration	the ability to quickly reach their top speed.	integer
Sprint_Speed	the top running speed.	integer
Agility	the ability to change direction quickly while maintaining control of the ball.	integer
Reactions	the ability to quickly respond to in-game situations	integer
Balance	the ability to maintain their balance when under physical pressure or when making quick movements.	integer
Shot_Power	the ability to generate powerful shots on goal.	integer
Jumping	the ability to jump high, which can be crucial in aerial duels or headers.	integer
Stamina	the ability to maintain their performance level throughout the match without becoming fatigued.	integer
Strength	the physical strength, which can help in challenges or holding off opponents.	integer
Long_Shots	the ability to accurately shoot from long distances.	integer
Aggression	the tendency to engage in physical challenges and apply pressure on opponents.	integer
Interceptions	the ability to read the game and intercept passes from the opposing team.	integer
Positioning	the ability to be in the right place at the right time, both offensively and defensively.	integer
Vision	the ability to see and make creative passes or plays during a match.	integer
Penalties	the skill in taking penalty kicks.	integer
Composure	the ability to remain calm and focused under pressure, which can impact their decision-making and performance.	integer
Defensive_Awareness	the ability to read the game defensively and position themselves effectively to make interceptions or tackles	integer
Standing_Tackle	the ability to effectively perform standing tackles to dispossess opponents without committing fouls.	integer
Sliding_Tackle	the skill in executing sliding tackles to win the ball from opponents while minimizing the risk of fouls or injuries.	integer

B Optimal hyperparameters for compared algorithms

Table 3: Optimal hyperparameters for compared algorithms

Model Name	Hyperparameters
AdaBoost	learning_rate=0.1, loss='linear', n_estimators=150
GBDT	learning_rate=0.1, max_depth=3, n_estimators=900
LightGBM	max_depth=3, n_estimators=900
Random Forest	max_depth=15, min_samples_leaf=3, min_samples_split=3, n_estimators=700
XGBoost	learning_rate=0.3, max_depth=3
CatBoost	depth=16, subsample=0.9