

Contributed Factors in Predicting Market Values of Loaned Out Players of English Premier League Clubs

Muhammad Daffa Arviano Putra¹, Deshinta Arrova Dewi²,
Wahyuningdiah Trisari Putri³, Retno Hendrowati⁴, Tri Basuki Kurniawan⁵

Department of Informatics-Faculty of Engineering Science, Paramadina University, Jakarta, Indonesia^{1,3,4}
Faculty of Data Science and Information Technology-INTI International University, Nilai, Malaysia^{1,2}
Postgraduate Program of Information Technology-Bina Darma University, Palembang, Indonesia⁵

Abstract—The top tier of the English football league division is occupied by the English Premier League (EPL). It has become a global phenomenon with exhilarating skills and has been one of the most-watched professional football leagues on the planet. The possibility of a player temporarily playing for a club other than the one to whom they are now contracted is known as a "loan player" in the English Premier League (EPL) hence, each player has a market value. Market value is an estimate of how much a player costs when a club wants to buy his contract from another club. The purpose of this study is to determine the factors that influence a player's market value at the conclusion of a loan period. With the Transfermarkt player transfer record dataset for the years 2004 through 2020, we use linear regression analysis. Our study found that a football player's market worth at the end of a loan period is influenced by several aspects, including market value at the beginning, goals, appearances, and total loan.

Keywords—Data analytics; predicting market value; English Premier League; loaned out players; consumption; resource use

I. INTRODUCTION

Football is one of the most popular team sports worldwide [1]. According to research conducted by Nielsen Sports, more than 40% of people aged 16 or older in countries with high populations and large markets said they are "interested" or "very interested" in following football [2]. The most prominent and well-known football league in the world is the English Premier League (EPL). The EPL has a lot of fans all over the world. According to the official website of the English Premier League, the cumulative global audience of EPL for season 2018/2019 was over 3 billion [3]. Due to the huge fans and excitement of the EPL, it has succeeded in attracting the interest of investors. Matchday revenue, broadcast deals, and commercial activity are some of the ways a club can generate a lot of profit [4].

To maximize revenue, the club must be popular among the viewers and have winning matches. Therefore, every owner strives to strengthen their club squad to be competitive and have a successful season. One way to strengthen a club is to buy good and talented players. Some of the examples we saw recently when Manchester City bought Jack Grealish from Aston Villa with a figure of €117.50m, Chelsea bought Romelu Lukaku from Inter Milan for €115.00m, and Manchester United bought Jadon Sancho for €85.00m from the German club, Borussia Dortmund [5].

Each player has a market value. Market value is an estimate of how much a player costs when a club wants to buy his contract from another club [1]. The price does not apply if a club only loans a player. If a club wants to loan a player, they are most likely only required to pay the player's salary for the duration of the loan. However, the player will still carry a market value that can go up or down when playing at the loan club. During the loan period, the player might perform extraordinarily and make many appearances and therefore, his market value will increase and vice versa. In every transfer window, every club in the EPL is likely to loan out their players to make room for their squad or give players a chance to build a reputation at another club. For this study, we use the market value published by the Transfermarkt website. The market values provided by the website are economically relevant and are viewed as having a fine reputation in the sports industry [6]. A study by Peeters revealed that Transfermarkt crowd valuation is referenced privately by club officials during player contract negotiations because it is more accurate than other valuations such as FIFA ranking and the ELO rating [7].

The academic community has found the football transfer market to be an engaging subject [8]. Additionally, we believe that additional investigation into the market value of football players would be an intriguing topic to pursue. Fortunately for us, the information required to do so is readily accessible through websites devoted to the sport of football. The general contributing aspects to a player's market value are not often covered in academic publications, nevertheless. Since the market value at the conclusion of the loaned time in EPL is greater than average, we are looking for contributing causes for that higher market value in this study.

A. Data Availability Statement

The data collected for the model is from Transfermarkt, a German Website that provides football information and data, such as scores, league tables, club squads, and many more using web scraping. Football-related research has used this website as its source data. The website incorporates crowdsourcing to estimate a player's market value in several professional football leagues. This means that every person can join the community and discuss the market value of any football player.

II. PREVIOUS STUDIES

There were studies on the subject of prediction of the market value of a football player that aligned with this research. Such as one from Singh and Lamba that suggested consistency, popularity, crowd estimation, and performance parameters enhance the prediction accuracy of the market value of football players [9]. While Felipe et. al. stated in their paper that the playing position (attacking midfielders) and age of the player (born in the first quarter of the year) are the most economically valued in terms of current value and maximal value [10]. Müller, Simons, and Weinmann stated that there are three categories of indicators of the market value of a player. There are the player characteristics (age, height, position, footedness, nationality), player performance (playing time, goals, assists, passing, dribbling, dueling, fouls, and cards), and player popularity (news, internet links) [1]. Further study on the popularity component, Frenger et. al. suggested that social media activities significantly influenced a football player's market value on the site [11].

A study on the player performance is also done by Richau et. al. which emphasized the actual performance of a football player to determine their market value, the paper stated that the actual performance is measured through individual player's age, minutes played, offense, defense, and team and analyzed using boosted regression trees [12]. They found that individual player performance indicator does not have the highest influence on the market value; instead, they found that team dimension average rank influences market value. Along this line, Metelski tried to find factors affecting the value of football players in the transfer market for the Polish Football League using descriptive statistics and several statistical tests. The writer uses the position on the pitch, age at transfer, year of transfer, the destination country, and selling club as the indicators. They found that the age of the player is a significant factor in the football players' values [13]. Moreover, Behravan and Razavi used the FIFA 20 dataset of various performance ratings of 18,278 players. Their novelty is the use of an automatic clustering algorithm in the first phase which they called APSO-clustering, and further training of a hybrid regression method called PSO-SVR for each cluster [14] that can estimate the players' value with an accuracy of 74%.

III. METHOD

Fig. 1 illustrates the methods we use for this study. These steps will be more specifically described afterward.

B. Data Collection

The data collected for the model is from Transfermarkt, a German Website that provides many footballs information and data, such as scores, league tables, club squads, and many more using web scraping. Football-related research has used this website as their source data, such as in [14] [15][16][6][17]. The website incorporates crowdsourcing to estimate a player's market value in several professional football leagues. This means that every person can join the community and discuss the market value of any football player. Everyone can suggest a market value for a player with good arguments and reasons to justify the player's estimation [6]. However, not everyone's opinion has the same value. The Transfermarkt website has

data on past loan players' performance in the EPL from the season 2004/2005 to 2021/2022. However, as the EPL 2021/2022 season is still ongoing at the time of the writing of this paper, we will limit the data to season 2020/2021.

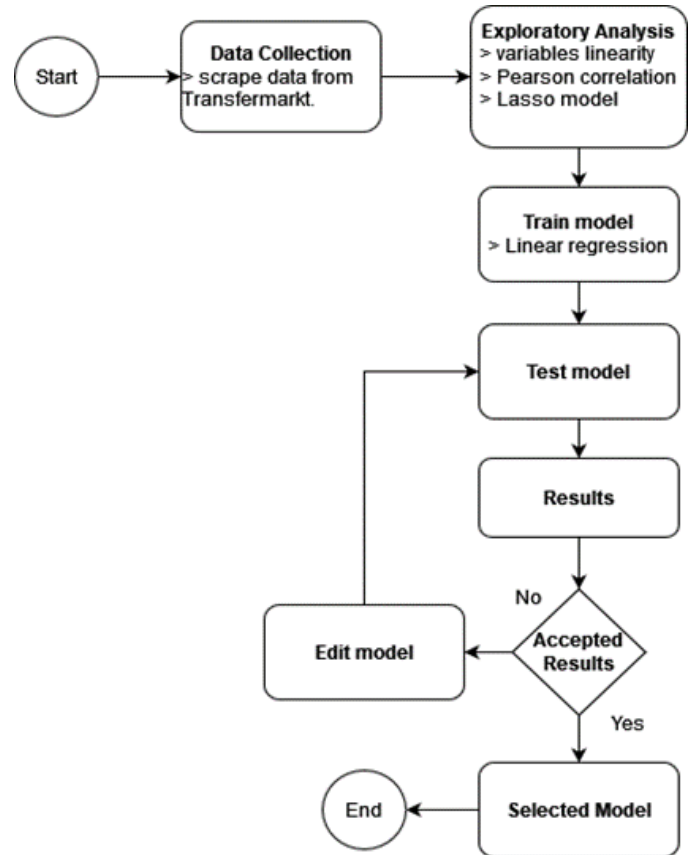


Fig. 1. Flowcharts of an applied method for the research.

There are 10 variables in the data as captured in Table I. In this study, we want to predict the last variable, which is the total market value of the loaned-out players from a club in a season. Therefore, the first nine variables are labeled as the explanatory variables and the last variable is our target or dependent variable.

C. Exploratory Analysis

In this section, we explore the data that we have. Our data has 339 rows and 10 columns, which means we have a total record of 339 data of all clubs who loaned out their players from season 2004/2005 to 2020/2021. Every season consists of 20 clubs and because the EPL uses a promotion and relegation system, there can be more than 20 unique clubs in the data. The first thing we have done is to list and count every unique club in the data. The outcome is that we have 40 unique clubs that will be trained in the Linear Regression model. We start exploring our data with the perspective of the relationship between variables, the data distribution, the correlation among variables, Least Absolute Shrinkage and Selection Operator (LASSO) in identifying features that may be the contributing factors to predicting market values of the loaned players.

TABLE I. SUMMARY OF THE AFM INFORMATION OF CDS QDS

Variable Name	Description
name	The club's name
year	The year of loaned out players data of a club (ranging from 2004 to 2020)
total_loan	The total loaned-out players from the club in the respective year
average_loan (in years)	The average number of loan periods of all loaned out players from the club in the respective year
appearances	The total number of appearances from all loaned-out players from the club in the respective year
starting_formation	The total number of appearances in the starting formation from all loaned-out players from the club in the respective year
goals	The total number of goals from all loaned-out players from the club in the respective year
average_minutes_played	The average minutes played by all loaned-out players from the club in the respective year (the maximum is 90 as a football match is played for 90 minutes)
market_value_at_start (in M €)	The total market values at the start of the loan period from all loaned-out players from the club in the respective year
market_value_at_end (in M €)	The total market values at the end of the loan period from all loaned-out players from the club in the respective year

D. Checking the Relationship between Variables using Scatter Plots

We start by checking the bivariate relationships between eight variables of our data. The total loan personnel, average loan time in years, number of appearances made by the players, starting formation of the players, goals made, average minutes played, market value at the start, and the last variable is the market value at the end. As we can see from the scatter plot in Fig. 2, there are three types of relationships shown, discrete relationships, random, and linear.

The discrete plot was obtained from the total loaned variables as there were only two values in these variables as the players were loaned only for 1 year or 2 years. The random relationship is shown by the total loaned players' variable against average minutes played, market value at the start, and market value at the end. The random relationship we see with average minutes played against all other variables. While the rest of the bivariate combinations showed some kind of linear relationship.

In this paper, we focused on the relationship between the market value at the end with the rest of the variables, so we found that the market value at the end has a strong linear relationship with the market value at the start. We explore the relationship more in the sections below.

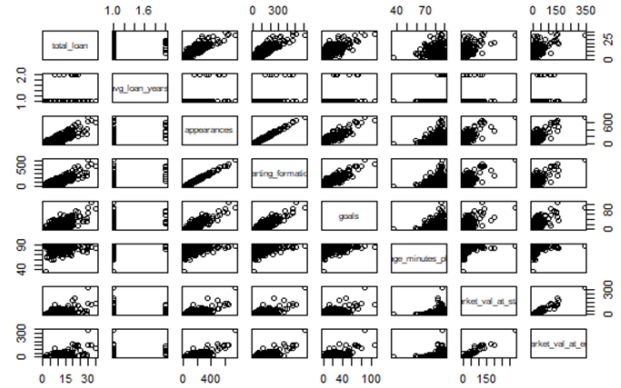


Fig. 2. Bivariate relationship of variables data.

E. Checking the Univariate Distribution using Histogram

To see if the observed data represent a random sample from the population; we use the histogram to check the distribution. Fig. 3 below shows the distribution of seven variables that we consider from the data; we did not include the name and year variables because they are categorical. The initial histogram showed left and right skewed data, so we do a log transformation on the variables to stabilize the variance, such as seen in [18] and [19]. Ensuring the log transformation, the average minutes played variable is still left skewed, while the average loan years is discrete.

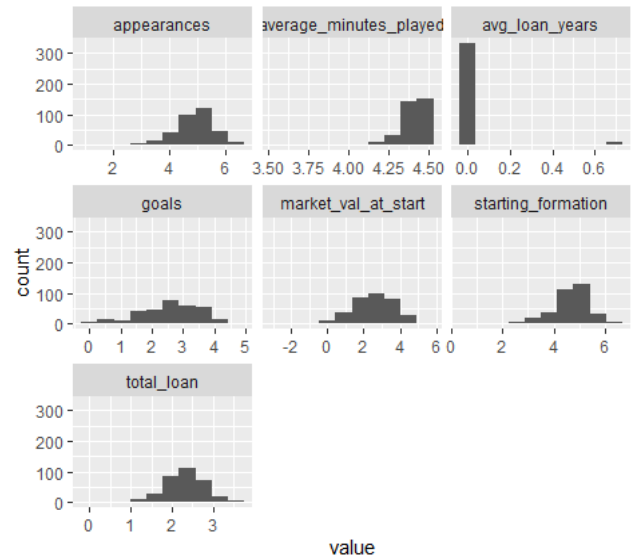


Fig. 3. Log transformed data distribution of market value predictors.

F. Pearson Correlation

The next thing we do in our exploratory step is to see the Pearson correlation to determine the precise extent or degree of any connection between any two variables, indicating its presence or absence.

that model 1 has the least R2 values, at an average of 88.9%, 87.9%, and 87.2% for respective 3, 5, and 10 folds. The percentage showed how well the dependent variable, market value at the end can be accounted for by the nine predictors. Model 2 came second with the average R2 values of 89.6%, 88.7%, and 88.3% accounted for by seven predictors, and model 3 has the highest average R2 values of 90.5%, 90.3%, and 90.1% accounted for by five predictors.

In terms of accuracy, the mean RMSE for model 1 for 3, 5, and 10 folds, respectively are 8.69525, 7.91934, and 7.24353. For Model 2, the mean RMSE are 8.10107, 7.51161, and 6.84796, respectively. For Model 3, the mean RMSE are 0.33339, 0.32948, 0.31620, respectively. Model 3 accuracy is higher than the two previous models. Therefore, we can conclude that removing the two variables average minutes played, and average loan years, is the right decision. With model 3 we come up with this linear equation:

$$\begin{aligned} \text{market_val_at_end} = & (-2.282 - 0.04653 * \text{total_loan}) + (0.5237 * \\ & \text{appearances}) - (0.00007325 * \text{starting_formation}) + (0.1056 * \\ & \text{goals}) + (0.8716 \text{ market_val_at_start}) \end{aligned} \quad (1)$$

With the use of the aforementioned equation, we can observe that while every variable affects the market value at the end, the market value at the beginning, objectives, appearances and total loan all has positive correlations and significant contributions. This differs slightly from our initial hypotheses based on the Pearson correlation, which was initial market value, early appearances, and initial formation.

The majority of the contributing factors to the loan player in the ELP have generally been identified by our investigation. With this investigation and its outcome, we have discovered a previously unknown association. We have identified which variables are connected to or have the strongest relationships with, and we may be able to identify patterns within the dataset as a result of this understanding.

V. CONCLUSION

In this study, we identified the elements that contributed to a football player's market worth in the English Premier League at the end of the loan period. Market value at launch, goals, appearances and total loan is among them. The market worth of a football player after a loan has been made is something we can forecast using exploratory research and a linear regression model. To discover the best predictor, we tested three distinct models. Our study revealed that the predictors differed slightly from what we had initially thought. Although linear regression is simple to comprehend and explain, we believe the model is adequate for use in this investigation. In future experiments, we hope to incorporate more data into our model, as the current data is very limited. We can also implement other models to better understand the contributing factors of a football player's market value.

REFERENCES

- [1] O. Müller, A. Simons and M. Weinmann, "Beyond crowd judgments: Data-driven estimation of market value in association football," *European Journal of Operational Research*, vol. 263, no. 2, pp. 611-624, 2017.
- [2] "World Football Report 2018," Nielsen Sports, 2018. [Online]. Available at: <https://www.nielsen.com/wp-content/uploads/sites/3/2019/04/world-football-report-2018.pdf>.
- [3] "Premiere League Global Audience on The Rise," Premier League, 2019. [Online]. Available: <https://www.premierleague.com/news/1280062>.
- [4] T. Dima, "The Business Model of European Football Club Competitions," *Procedia Economics and Finance*, vol. 23, no. Oct. 2014, pp. 1245-1252, 2015.
- [5] "Premiere League - Transfer records," Transfermarkt, 2021. [Online]. Available: https://www.transfermarkt.com/premier-league/transferrekorde/wettbewerb/GB1/plus/galerie/0?saison_id=2021&land_id=alle&ausrichtung=&spielerposition_id=alle&altersklasse=&leih=&w_s=s&zuab=zu.
- [6] S. Herm, H.-M. Callsen-Bracker and H. Kreis, "When the crowd evaluates soccer players' market values: Accuracy and evaluation attributes of an online community," *Sport Management Review*, vol. 17, no. 4, pp. 484-492, 2013.
- [7] T. Peeters, "Testing the Wisdom of Crowds in the field: Transfermarkt valuations and international soccer result," *International Journal Forecasting*, vol. 34, no. 1, pp. 17-29, 2018.
- [8] D. Matesanz, F. Holzmayer, B. Torgler, S. L. Schmidt and G. J. Ortega, "Transfer market activities and sportive performance in European first football leagues: A dynamic network approach," *PLOS ONE*, vol. 13, no. 12, pp. 1-16, 2018.
- [9] P. Singh and P. S. Lamba, "Influence of crowdsourcing, popularity and previous year statistics in market value estimation of football players," *Journal of Discrete Mathematical Sciences and Cryptography*, vol. 22, no. 2, pp. 113-126, 2019.
- [10] J. L. Felipe, A. Fernandez-Luna, P. Burillo, L. E. de la Riva, J. Sanchez-Sanchez and J. Garcia-Unanue, "Money Talks: Team Variables and player Positions that Most Influence the Market Value of Professional Male Footballers in Europe," *Sports Policy and Finance*, vol. 12, no. 9, pp. 10-17, 2020.
- [11] M. Frenger, F. Follert, L. Richau and E. Emrich, "Follow me... on the relationship between social media activities and market values in the German Bundesliga," *Saarbrücken*, 2019. [Online]. Available: <http://www.soziooekonomie.org>.
- [12] L. Richau, F. Follert, M. Frener and E. Emrich, "Performance indicators in football: The importance of actual performance for the market value of football players," *SCIAMUS Sport and Manag*, vol. 4, pp. 41-61, 2019.
- [13] A. Metelski, "Factors affecting the value of football players in the transfer market," *Journal of Physical Education and Sport*, vol. 21, no. 2, pp. 1150-1155, 2021.
- [14] I. Behravan and S. M. Razavi, "A novel machine learning method for estimating football players' value in the transfer market," *Soft Computing*, vol. 25, no. 3, pp. 2499-2511, 2021.
- [15] H. Adiwiyana, H. I. Adiwiyana and Harywaman, "Factors that Determine the Market Value of Professional Football Players in Indonesia," *J. Din. Akunt*, vol. 13, no. 1, pp. 51-61, 2021.
- [16] R. Stanojevic and L. Gyarmati, "Towards Data-Driven Football Player Assessment," *IEE Int. Conf. Data Min. Work. ICDMW*, vol. 0, pp. 167-172, 2016.
- [17] M. He, R. Cachucho and A. Knobbe, "Football Player's Performance and Market Value," in *Proc 2nd Work. Sport. Anal. Eur. Conf. Mach. Learn. Princ. Pract. Knowl. Discov. Databases (ECML PKDD)*, 2015.
- [18] D. Curran-Everett, "Explorations in statistics: the log transformation," *Adv. Physiol. Educ.*, vol. 42, no. 2, pp. 343-347, 2018.
- [19] H. Son, C. Hyun, D. Phan and H. J. Hwang, "Data analytic approach for bankruptcy prediction," *Expert Systems with Applications*, vol. 138, 2019.
- [20] Z. Yan and Y. Yao, "Variable selection method for fault isolation using least absolute shrinkage and selection operator (LASSO)," *Chemom. Intell. Lab. Syst.*, vol. 146, pp. 136-146, 2015.
- [21] S. Tian, Y. Yu and H. Guo, "Variable selection and corporate bankruptcy forecasts," *Journal of Banking & Finance*, vol. 52, no. December, pp. 89-100, 2015.

- [22] P. Ghosh, S. Azam, M. Jonkman and A. Karim, "Efficient Prediction of Cardiovascular Disease Using Machine Learning Algorithms with Relief and LASSO Feature Selection Techniques," *IEEE Access*, vol. 9, pp. 19304-19326, 2021.
- [23] M. R and R. R, "LASSO: A Feature Selection Technique in Predictive Modeling for Machine Learning," 2016 IEEE International Conference on Advances in Computer Application (ICACA), pp. 18-20, 2016.