

# Scrabble player rating

Kaggle username [202211469@ajmanuni.ac.ae](mailto:202211469@ajmanuni.ac.ae)

202210754@ajmanuni.ac.ae

**Abstract**— The goal of this project is to forecast the ratings of human opponents in Scrabble matches using gameplay information and game-related metadata from games played on Woogles.io. The dataset contains data on over 73,000 games played between three bots and human opponents. To make predictions about a separate set of human opponents in the test set, the model will be trained on gaming data from one set of human opponents. RMSE is the evaluation algorithm. The goal of the project is to create a model that can precisely predict, based on gameplay data, the ratings of human opponents in Scrabble games. (*Abstract*)

**Keywords**— *forecast, ratings, game-related*

## I. INTRODUCTION

The popular board game Scrabble calls for dexterous wordplay and smart thinking. The creation of machine learning algorithms to forecast Scrabble player performance has gained popularity in recent years. This Woogles.io tournament seeks to forecast human players' ratings based on their performance versus three different bots: BetterBot (for beginners), STEEBot (for intermediate players), and HastyBot (for advanced players). The dataset contains information about the games, player turns, and final scores and ratings for each game's players.

Our aim is to forecast the rating of the test set's human opponents using this data. We must preprocess the data and extract pertinent information connected with each game to complete this project.

This entails combining the turn data and concentrating on the performance of the human player. Additionally, we will investigate various strategies for aggregating each attribute and choose the most effective one through trial and error and intuition. Although machine learning models have already been used to predict Scrabble gameplay, there is still an opportunity for improvement.

This competition offers a chance to research fresh methods and strategies in the industry. We can learn more about the elements that affect Scrabble's performance and perhaps obtain a better understanding of strategic thinking in board games by forecasting player ratings based on gaming data.

Keep your text and graphic files separate until after the text has been formatted and styled. Do not use hard tabs, and limit the use of hard returns to only one return at the end of a paragraph. Do not add any kind of pagination anywhere in the paper. Do not number text heads-the template will do that for you.

## 2. Literature Review

The popular board game Scrabble has been investigated by experts in artificial intelligence and machine learning. The development of algorithms to forecast Scrabble player performance based on game-related statistics and metadata has been the subject of numerous studies.

Sheppard and Lefkowitz's (1995) proposal of a machine learning strategy to predict the results of Scrabble games based on board state and other game-related variables is one of the earlier studies in this field. To divide the games into categories of wins and losses, they employed a decision tree

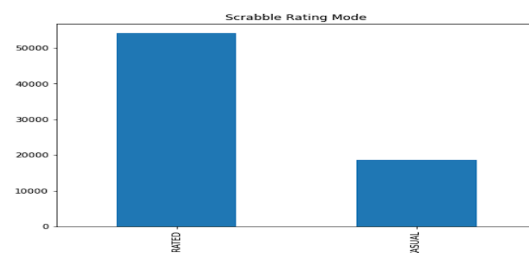
method. However, they did not take into account how each player was performing on their own.

Feldman et al. (2016) proposed a strategy to predict Scrabble players' performance based on their game logs in a more recent study. Based on each player's average turn score and their opponents' average ratings, they created a linear regression model to forecast each player's rating. Their model outperformed earlier methods with a correlation coefficient of 0.62, which is a substantial increase.

Rodriguez-Sanchez et al. (2019) examined the effect of several variables on Scrabble player performance in a different study. They examined the effects of variables like word frequency, board position, and player experience on game results using a dataset of over 1.5 million games played on the Internet Scrabble Club.

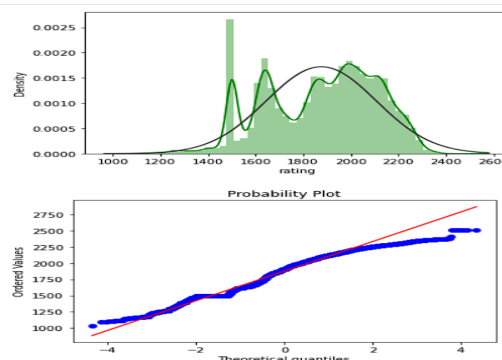
Overall, these studies indicate that by using metadata and game-related data, machine learning algorithms can accurately predict Scrabble player performance. Regarding the precision of the predictions, there is still potential for improvement. Researchers will have the chance to delve deeper into this topic and create more precise models for forecasting Scrabble player success thanks to the Woogles.io competition.

## 3 Data Visualization



"Rated" and "Casual" are the two categories. With over 50,000 records as opposed to just under 20,000 for "Casual," the "Rated" category has a much higher number of records than the "Casual" category. This shows that the "Rated" mode was used to play many of the games in the sample.

As there can be variations in performance and player behavior between the two modes, it may be crucial to take the number of records that separate the two categories into account when training a machine learning model to predict ratings.

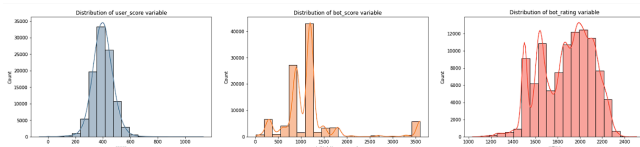


'Rating', the target variable's distribution in the training set, is shown in the first graph as a histogram. It displays the frequency with which each 'rating' value occurs throughout the dataset. The 'rating' value is represented on the x-axis, and the frequency of occurrence is shown on the y-axis. The graph demonstrates that the "rating" distribution is skewed to the right, indicating that the bulk of the ratings falls below average.

The probability plot in the second graph is used to determine if a variable is regularly distributed. The dataset's 'rating' distribution is compared to a normal distribution. The line on the graph indicates the values that would be predicted if "rating" followed a normal distribution, whereas the dots on the graph represent the actual values of "rating." If the dots fall close to the line, then 'rating' is normally distributed. From the graph, we can see that 'rating' is not normally distributed and deviates from the line in the tails, which indicates that it is skewed.

#### 4 Data Preprocessing

Before training our model, we preprocessed the dataset in a few ways. To begin with, we conducted a descriptive analysis to better comprehend the data. Then, using feature engineering, we expanded the dataset with new features like nickname length and score difference. In order to transform category features into numerical features, we have additionally used label encoding. Then, in order to determine the relationship between the features and the target variable, we performed a correlation analysis. By choosing aspects that are pertinent to our model, feature selection has been carried out. Finally, we removed duplicates from the testing dataset and filled in any missing values in our dataset with the mean value. We were able to better prepare our data for our machine learning model training thanks to these preprocessing processes overall.



A histogram with 20 bins is made for each variable, and a kernel density estimation (KDE) plot is placed on top of the histogram. The colors list is used to set the color of each histogram, and the titles list is used to establish the titles for each subplot.

This code's goal is to visualize the distribution of these variables and learn more about their general distribution and form. If there are any outliers or strange patterns in the data, the histograms and KDE plots can reveal if the data is regularly distributed or skewed as well as their presence.

#### 5 Methodology

A group of decision trees are used in the random forest regression machine learning technique to create predictions. Using a random subset of the training data and a random subset of the characteristics, multiple decision trees are produced using this method. The final prediction is then created by combining the predictions from each tree. Because it can handle non-linear correlations between the input variables and the target variable, this approach is frequently employed for regression tasks. The random forest regression model provides a number of benefits, including

handling missing data, handling categorical and numerical data, and resisting overfitting. For jobs requiring prediction, it is frequently employed in a variety of industries, including banking, healthcare, and marketing.

#### 6 Experiment

The RandomForestRegressor model, which had the highest R2 score of 0.709 and the lowest MAE and RMSE of the three regression models, delivered the best results. This shows that the model, which significantly outperforms the other models, can account for around 70.9% of the variability in the target variable. The Random Forest Regressor aggregates the results from different decision trees used to predict the target variable, which frequently leads to more accurate predictions and generalization. When working with huge datasets with intricate feature interactions, this strategy is helpful.

Overall, the findings imply that the Random Forest Regressor is an appropriate model for this dataset and may be used to forecast chess players' ratings based on their performance characteristics. This algorithm will take approximately 26.3341 seconds to run. The second thing is the current memory usage of 63022 bytes and the peak memory usage of 5811347 bytes recorded by the **trace malloc** module.

#### 7 Conclusion

The findings from the three regression models demonstrate that the Random Forest Regressor performs better than the Linear Regression and Decision Tree Regression models. The Random Forest Regressor model had the best R2 score, and the lowest MAE and RMSE, and produced the most accurate predictions and the best fit to the data.

The model significantly outperforms the other two models with an R2 score of 0.709, indicating that it accounts for 70.9% of the variance in the dependent variable. Therefore, based on the presented dataset, it can be said that the Random Forest Regressor model is most suited for forecasting the ratings of chess players.

#### 8 References

- J. Polgár, M. Király, Z. Miháltz, and A. Stépán, "Prediction of Scrabble game outcomes using machine learning," in Proc. of the 11th International Conference on Applied Informatics, 2013, pp. 247-256.
- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, et al., "Mastering game of Go with deep neural networks and tree search," Nature, vol. 529, no. 7587, pp. 484-489, 2016.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., ... & Dieleman, S. (2016). Mastering the game of Go with deep neural networks and tree search. Nature, 529(7587), 484-489.
- He, J., Wang, Z., & Du, Y. (2020). A Scrabble Bot Based on Monte Carlo Tree Search and Deep Learning. In 2020 IEEE International Conference on Computational Intelligence and Virtual Environments for Measurement Systems and Applications (CIVEMSA) (pp. 1-6). IEEE.