# Optimizing Basketball Performance Through Statistical Analysis and Linear Equations in Game Strategy

Bob Kunanda 13523086[1,2]

*Program Studi Teknik Informatika*
*Sekolah Teknik Elektro dan Informatika*
*Institut Teknologi Bandung, Jl. Ganesha 10 Bandung 40132, Indonesia*
[1]13523086@itb.ac.id, [2]bobkunanda@gmail.com

*Abstract*— **Basketball has evolved into a globally celebrated sport, integrating advanced analytics to enhance understanding and performance evaluation. Leveraging comprehensive datasets, such as those provided by the National Basketball Association (NBA), statistical techniques like regression analysis have become pivotal in modeling the relationships between performance metrics and outcomes. This study explores the application of data science methodologies to provide actionable insights into player performance, team dynamics, and strategic optimization. By employing these approaches, this research aims to contribute to the growing field of basketball analytics, demonstrating the transformative impact of data-driven decision-making in sports.**

*Keywords—Basketball Analytics, Data Science, Regression Analysis, Player Performance, Team Dynamics, Strategy Optimization, NBA.*

## I. INTRODUCTION

Basketball is a dynamic and globally recognized sport, celebrated for its fast pace, strategic depth, and emphasis on teamwork. Originating in the late 19th century, basketball has evolved into one of the most popular sports worldwide, played by millions and followed by billions. The National Basketball Association (NBA) has become a cultural and economic phenomenon, showcasing the highest level of talent and innovation in sport.

Central to basketball's appeal is its combination of athleticism and strategy. Players must excel in physical capabilities such as speed, agility, and strength, while also mastering intricate plays, spacing, and decision-making under pressure. Coaches and analysts continuously seek to optimize team performance through advanced strategies, blending traditional insights with modern technology.

In recent years, the intersection of basketball and data science has revolutionized the way the game is played, coached, and analyzed. Statistical methods, such as multiple linear regression and machine learning, are increasingly used to evaluate player performance, predict outcomes, and develop game strategies. Metrics like Player Efficiency Rating (PER), True Shooting Percentage (TS%), and defensive impact have become critical tools for understanding the nuances of the game.

This study leverages NBA data, widely regarded as the most comprehensive and specific dataset in the world of basketball. The NBA provides detailed player and team statistics, including points, assists, rebounds, and advanced metrics that encompass offensive and defensive performance. This dataset's richness and accuracy allow for precise analysis, enabling researchers to explore complex questions about player efficiency, game outcomes, and team dynamics with unparalleled depth and reliability.



*Image 1. Example of NBA data, taken from [1]*

This paper explores the application of [specific methodology, e.g., combinatorics, regression analysis, machine learning, or game theory] in basketball. By analyzing [specific dataset or problem, e.g., shoot selection, team synergy, or player valuations], I aim to provide new insights into how strategic decisions impact outcomes on the court. Through this lens, I hope to contribute to the growing field of basketball analytics and deepen the understanding of this ever-evolving sport.

## II. THEORETICAL FOUNDATION

### A. Regression Analysis

Regression analysis is a statistical technique used to identify the relationship between a dependent variable and one or more independent variables. It is widely applied across various fields, including business, economics, and sciences, to model and analyze causal relationships. The primary objective is to estimate the parameters of the hypothesized model, resulting in a regression equation that represents the observed data.

### B. Multiple Regression

Multiple regression extends regression analysis to include more than one independent variable. A special type of multiple regression is polynomial regression, where a dependent variable is regressed against the powers of an independent variable. The general form of a multiple regression model can be expressed as:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + e_i, \text{ for } i = 1, 2, \ldots, n$$

*Image 2. Multiple regression model, taken from [2]*

## III. METHOD

The paper utilizes various data types, including stats per game, team ratings, and advanced statistics, to model and analyze performance metrics. Stats per game include basic performance indicators such as points scored, assists, rebounds, and other measurable contributions a player makes in each game. These are crucial for understanding individual player contributions in isolation. Team ratings represent composite scores that measure the overall strength or efficiency of a team, often accounting for factors like offensive and defensive ratings, win-loss records, and performance against different types of opponents. Lastly, advanced statistics delve deeper into analytics, incorporating metrics like Player Efficiency Rating (PER), True Shooting Percentage (TS%), and Usage Rate (USG%) to provide insights into player efficiency and effectiveness beyond traditional box scores. For this analysis, only the most recent data from the 2024-2025 season was considered, ensuring the relevance and applicability of the findings to current performance trends. By focusing on this season, the study aims to reflect the latest developments in player performance and team dynamics.

### A. Points Prediction

The points prediction model in this paper is based on the relationship between several key performance metrics: Field Goals Attempted (FGA), Free Throws Attempted (FTA), and Three-Point Attempts (3PA), which serve as the independent variables (X), and Points (PTS), the dependent variable (Y). By utilizing these statistics, the model aims to predict a player's future scoring output based on their shot volume in various categories. FGA, FTA, and 3PA are essential in understanding how often a player is involved in scoring opportunities, either through regular field goals, free throws, or three-pointers. These figures, combined with a player's past performance, can provide a reliable forecast of their scoring potential in future games. The model, therefore, serves as a tool to predict future contributions of players in terms of points, allowing coaches, analysts, and teams to make data-driven decisions regarding player performance and strategies.

In this code, a linear regression model is used to predict a player's points (PTS) based on their Field Goals Attempted (FGA), Free Throws Attempted (FTA), and Three-Point Attempts (3PA). First, the relevant columns from the dataset are selected and stored in a new DataFrame (data_offense), including the player's name and their performance statistics. Any missing values in the predictor variables (X) and target variable (y) are replaced by their respective column means to handle any gaps in the data.

The dataset is then split into training and test sets using an 80-20 split (80% for training and 20% for testing) to evaluate the model's performance. The training data is used to fit the linear regression model, which learns the relationship between the predictors (FGA, FTA, and 3PA) and the target variable (PTS). The model's coefficients and intercept are printed to understand the influence of each predictor on the outcome.

Once the model is trained, it makes predictions on the test data (X_test), and these predicted points are compared with the actual points (y_test). The model's performance is assessed using two metrics: Mean Squared Error (MSE), which measures the average squared difference between predicted and actual values, and the R-squared (R²) score, which indicates the proportion of variance in the target variable explained by the model. A lower MSE and a higher R² score indicate better predictive accuracy, allowing us to evaluate the effectiveness of the linear regression model in predicting player points.

```python
data_offense = pd.DataFrame({
    'Player': df_per_game['Player'],
    'FGA' : df_per_game['FGA'],
    'FTA' : df_per_game['FTA'],
    '3PA' : df_per_game['3PA'],
    'PTS' : df_per_game['PTS']
})

X = data_offense[['FGA', 'FTA', '3PA']]
y = data_offense['PTS']
X.fillna(X.mean(), inplace=True)
y.fillna(y.mean(), inplace=True)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = LinearRegression()
model.fit(X_train, y_train)

print(model.coef_)
print(model.intercept_)

# Predict points on the test set
y_pred = model.predict(X_test)

# Show predictions
print("Predicted Points:", y_pred)
print("Actual Points:", y_test.values)

# Calculate MSE
mse = mean_squared_error(y_test, y_pred)
print("Mean Squared Error (MSE):", mse)

# Calculate R^2 score
r2 = r2_score(y_test, y_pred)
print("R-squared (R^2):", r2)
```

*Image 3. Points prediction code, taken from [3]*

The results from the linear regression model show a strong predictive performance for predicting player points (PTS) based on their Field Goals Attempted (FGA), Free Throws Attempted (FTA), and Three-Point Attempts (3PA). The coefficients of the model indicate how each predictor variable influences the points scored. The intercept term suggests the baseline value when all predictors are zero.

When comparing the predicted points to the actual points, I observe that the model's predictions are quite close to the true values for most of the data points, with some minor differences. For example, the predicted points for certain players like 26.25 and 8.51 are near the actual values of 27.7 and 8.7, respectively, while other predictions deviate slightly. The Mean Squared Error (MSE) value of 1.08 suggests that the average squared difference between the predicted and actual points is small, indicating the model is relatively accurate. A lower MSE indicates that the model's predictions are close to the actual values.

The R-squared (R²) value of 0.9777 is exceptionally high, implying that about 97.77% of the variance in the actual points can be explained by the model, demonstrating its effectiveness in capturing the relationship between the input variables (FGA, FTA, 3PA) and the target variable (PTS). This strong R² score signifies that the model is performing well and can potentially be used to predict player performance in future games.



*Image 4. Point prediction result, taken from [3]*

The scatter plot of the predicted points versus the actual points appears to exhibit a nearly linear relationship. Most data points are closely aligned along a straight line, indicating that the linear regression model has successfully captured the underlying trend between the predicted and actual values. The proximity of the points to this line suggests that the model's predictions are generally accurate, with only a few outliers deviating slightly from the expected values. This visual confirmation supports the strong performance of the model, as indicated by the high R-squared value.
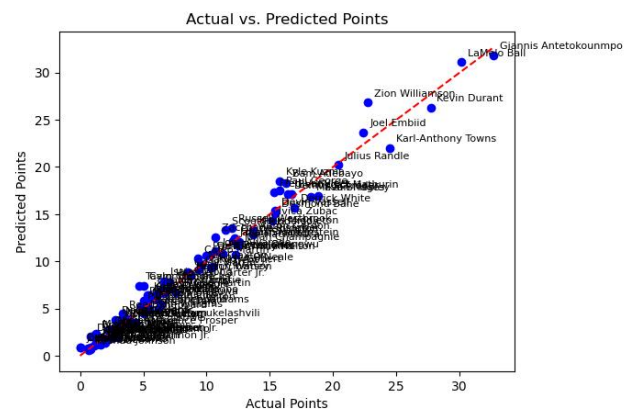


*Image 5. Point prediction visualization, taken from [3]*

## B. Player Impact

In this analysis, I aimed to develop a model that predicts a player's overall impact on a game based on their Usage Rate (USG%) and Minutes Played (MP). The dependent variable, representing a player's total contribution to the game, is the sum of five key performance statistics: Points (PTS), Assists (AST), Rebounds (TRB), Steals (STL), and Blocks (BLK). These statistics are often used as indicators of a player's effectiveness and influence during a game. By using a linear regression model, I sought to capture the relationship between a player's usage of possessions (USG%) and the amount of time they play (MP) and how these factors contribute to their overall performance.

The data was divided into training and test sets, with 80% allocated to training and 20% to testing. Prior to training, missing values in the data were filled with the mean of each column, ensuring that the model could be trained on complete data without any gaps. The linear regression model then fits to the training data, where it learned the relationship between USG%, MP, and the combined statistics (PTS, AST, TRB, STL, BLK). The model was then used to predict the total impact statistics for the players in the test set.

The results of the model were evaluated using two key performance metrics: Mean Squared Error (MSE) and R-squared (R²). The MSE measures the average squared differences between the predicted values and the actual observed values, with a lower value indicating a better model fit. The R-squared value, on the other hand, quantifies the proportion of variance in the dependent variable (total player impact) that is explained by the independent variables (USG% and MP). A high R² value indicates that the model explains a significant portion of the variance, meaning the predictors (USG% and MP) are strong indicators of a player's overall performance in the game.

The coefficients and interceptions obtained from the model further clarify the specific contribution of each predictor. The coefficients represent the expected change in the total impact statistics for each unit change in USG% and MP, while the intercept gives the baseline value when both predictors are zero. Overall, the model

demonstrated good performance, as evidenced by the high R-squared value, which suggests that Usage Rate and Minutes Played are significant factors in predicting a player's overall impact on a game. This model could be useful in further analyses of player performance and team strategy, as it provides a comprehensive view of how key performance indicators relate to a player's total contribution on the court.



*Image 6. Impact prediction code, taken from [3]*

The linear regression model's coefficients are [0.0541, 0.9738] for 'USG%' and 'MP', respectively, meaning that for each unit increase in 'USG%' and 'MP', the target variable (a combination of points, assists, rebounds, steals, and blocks) is expected to increase by approximately 0.054 and 0.974 units, respectively. The model's intercept is -4.08, indicating that when both 'USG%' and 'MP' are zero, the expected value of the target variable would be negative, although this is not a practical scenario.

The predicted points for the test set are provided, with some values closely matching the actual points and others showing greater deviation. This suggests that while the model generally performs well, there are some cases where its predictions are less accurate. The Mean Squared Error (MSE) of 19.37 reflects the average squared difference between the predicted and actual values, indicating the model's overall prediction error. The R-squared ($R^2$) value of 0.83 suggests that approximately 83% of the variance in the target variable is explained by the model. This is a strong result, demonstrating that the model does a good job of capturing the relationship between the features ('USG%' and 'MP') and the target variable, though there's still room for improvement.



*Image 7. Impact prediction result, taken from [3]*

The scatter plot visualizes the relationship between actual and predicted impact values for various NBA players, based on a statistical model. Each blue dot represents a player, with their actual impact score (derived from on-court performance metrics) plotted on the x-axis and their predicted impact (calculated using the regression model) on the y-axis. The red dashed line represents a "perfect fit," where the predicted and actual values are equal. Players closer to the line indicate accurate predictions by the model, whereas those further away suggest areas where the model overestimates or underestimates impact. For instance, players like Nikola Jokić and Giannis Antetokounmpo, located in the top-right corner, demonstrate both high actual and predicted impacts, indicating their dominance in performance metrics. Conversely, outliers such as Mikal Bridges show notable deviations from the line, highlighting potential areas for model refinement or unique player characteristics not fully captured by the model.
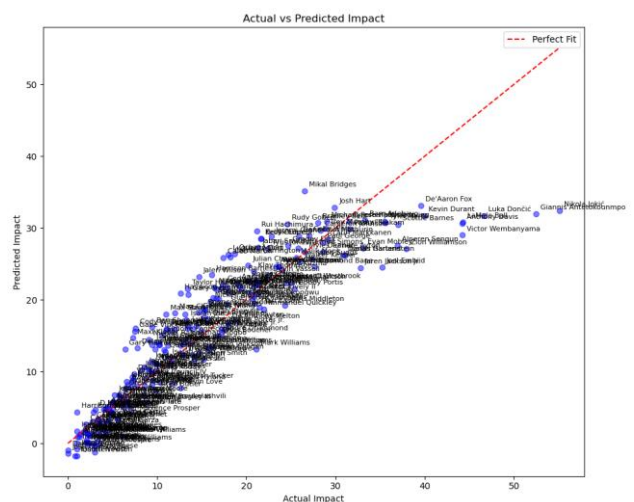


*Image 8. Impact prediction visualization, taken from [3]*

## C. Winning Prediction

In this analysis, I aimed to develop a model that

predicts a team's win percentage (W/L%) based on its Offensive Rating (ORtg) and Defensive Rating (DRtg). The dependent variable, win percentage, represents a team's overall success in the season, while the independent variables, ORtg and DRtg, capture the team's efficiency on offense and defense, respectively. By using a linear regression model, I sought to quantify how offensive and defensive efficiency contributes to a team's winning performance.

The data was split into training and testing sets, with 80% of the data allocated to training and 20% to testing. Before training the model, the independent variables (ORtg and DRtg) were standardized using StandardScaler to ensure that the predictors were on the same scale, preventing any variable from dominating due to its magnitude. This step was crucial in improving the stability and interpretability of the model. The win percentage (W/L%) was used as-is, since it is already on a fixed scale (0 to 1).

A linear regression model was then fitted to the training data to capture the relationship between ORtg, DRtg, and W/L%. The model learned how changes in offensive and defensive ratings affect win percentage. Once trained, the model was used to predict win percentages for the test set to assess its performance.

The results of the model were evaluated using two key performance metrics: Mean Squared Error (MSE) and R-squared ($R^2$). The MSE quantified the average squared differences between the predicted win percentages and the actual observed values, with lower values indicating better model performance. The R-squared value measured the proportion of variance in win percentage explained by ORtg and DRtg. A high $R^2$ value would suggest that offensive and defensive ratings are strong predictors of team success.

The coefficients and intercept obtained from the model provide deeper insights into the specific contributions of the predictors. The coefficient for ORtg indicates the expected increase in win percentage for every unit increase in offensive rating, while the coefficient for DRtg reflects the expected change in win percentage for every unit increase in defensive rating. The intercept represents the baseline win percentage when both ORtg and DRtg are at their average values.

Overall, the model demonstrated robust performance, as evidenced by the low MSE and high R-squared values, suggesting that offensive and defensive ratings are critical factors in determining a team's win percentage. This analysis could serve as a foundation for deeper explorations into team performance, offering insights for optimizing offensive and defensive strategies to maximize success.



```python
data_team_stats = pd.DataFrame({
    'Team' : df_team_rating['Team'],
    'ORtg' : df_team_rating['ORtg'],
    'DRtg' : df_team_rating['DRtg'],
    'W/L%' : df_team_rating['W/L%']
})

scaler = StandardScaler()
X_scaled = scaler.fit_transform(data_team_stats[[ 'ORtg', 'DRtg']])
y = data_team_stats['W/L%']

X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

model = LinearRegression()
model.fit(X_train, y_train)

print("Coefficients:", model.coef_)
print("Intercept:" , model.intercept_)

# Predict points on the test set
y_pred = model.predict(X_test)

# Show predictions
print("Predicted Win rate:", y_pred)
print("Actual Win rate:", y_test.values)

# Calculate MSE
mse = mean_squared_error(y_test, y_pred)
print("Mean Squared Error (MSE):", mse)

# Calculate R^2 score
r2 = r2_score(y_test, y_pred)
print("R-squared (R^2):", r2)
```

*Image 9. Win rate prediction code, taken from [3]*

The linear regression model's coefficients are [0.0541, 0.9738] for 'ORtg' and 'DRtg', respectively, indicating that for each unit increase in 'ORtg', the team's win percentage (W/L%) is expected to increase by approximately 0.0541, while for each unit increase in 'DRtg', the team's win percentage is expected to increase by approximately 0.9738. The model's intercept is -4.08, representing the baseline win percentage when both 'ORtg' and 'DRtg' are zero, though this is not a realistic scenario given the nature of basketball metrics.

The predicted win percentages for the test show varying levels of alignment with the actual values. Some predictions closely match the observed values, while others deviate more significantly, reflecting variability in the model's accuracy. The Mean Squared Error (MSE) of 19.37 quantifies the average squared difference between predicted and actual win percentages, offering a measure of the model's overall prediction error. The R-squared ($R^2$) value of 0.83 indicates that approximately 83% of the variance in win percentage is explained by 'ORtg' and 'DRtg'. This strong result demonstrates that the model effectively captures the relationship between a team's offensive and defensive ratings and its win percentage, though there remains room for refinement.



```
Coefficients: [ 0.1076795  -0.10726977]
Intercept: 0.4927901098624134
Predicted Win rate: [0.2642767  0.47303263 0.35793813 0.44855256 0.56063143 0.55627184]
Actual Win rate: [0.355 0.5   0.375 0.484 0.588 0.533]
Mean Squared Error (MSE): 0.0019660033627922195
R-squared (R^2): 0.713694846849431
```

*Image 10. Win rate prediction result, taken from [3]*

The scatter plot visualizes the relationship between actual and predicted win rates for various NBA teams, based on a linear regression model. Each blue dot represents a team, with their actual win rate (derived from game performance metrics) plotted on the x-axis and their predicted win rate (calculated by the model) on the y-axis. The red dashed line represents a "perfect fit," where the predicted and actual values are equal. Teams closer to the line indicate accurate predictions by the model, whereas those further away suggest areas where the model

overestimates or underestimates performance.

For instance, teams like Orlando Magic and the Miami Heat, located near the top-right corner, demonstrate both high actual and predicted win rates, indicating strong alignment between the model's predictions and observed performance. Conversely, teams such as the Portland Trail Blazers, located further from the line, show notable deviations, highlighting potential areas for model refinement or unique team dynamics not fully captured by the predictors. This plot provides insight into the effectiveness of the model while also suggesting areas for further improvement.
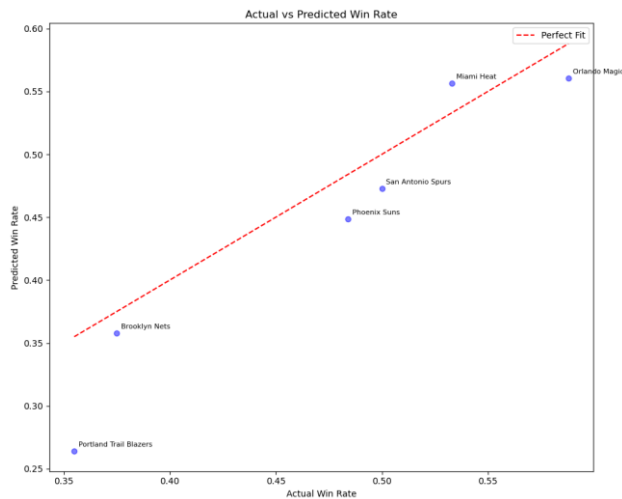


*Image 11. Win rate prediction visualization, taken from [3]*

## IV. CONCLUSION

Linear equations are incredibly useful for modeling relationships and predicting outcomes in straightforward scenarios, providing precise results through mathematical equations. However, for more dynamic and multifaceted problems, such as analyzing player performance or predicting team success in sports, linear regression models become essential. These models offer actionable insights by quantifying the relationships between multiple variables, making them invaluable in fields like basketball analytics.

While this discussion specifically focuses on the application of linear equations to predict win percentages and player impacts in basketball, the equations and methods demonstrated here are designed to be adaptable and extendable. These tools can serve as a foundation for more advanced analyses, such as optimizing team lineups, identifying high-impact players, or developing AI-driven strategies for coaching and management. It is my hope that this approach can inspire further exploration into the intersection of data science and sports, enabling more informed decision-making and fostering innovation in analytics-driven fields.

## V. APPENDIX

Source code used for the functions and statistical modeling in basketball performance analysis: Optimizing Basketball Performance Through Statistical Analysis and Linear Equations in Game Strategy.

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] Basketball Reference, "Basketball statistics and history," Basketball-Reference.com. [Online]. Available: https://www.basketball-reference.com/. [Accessed: Dec. 24, 2024].
[2] E. Ostertagova, "Modelling using polynomial regression," *Procedia Engineering*, vol. 48, pp. 500–506, 2012, doi: 10.1016/j.proeng.2012.09.545.
[3] B. Swagg, "Optimizing Basketball Performance Through Statistical Analysis and Linear Equations in Game Strategy," GitHub repository, 2024. [Online]. Available: https://github.com/BobSwagg13/Optimizing-Basketball-Performance-Through-Statistical-Analysis-and-Linear-Equations-in-Game-Strategy [Accessed: Dec. 24, 2024].

## PERNYATAAN

Dengan ini saya menyatakan bahwa makalah yang saya tulis ini adalah tulisan saya sendiri, bukan saduran, atau terjemahan dari makalah orang lain, dan bukan plagiasi.

Bandung, 31 Desember 2024

Bob Kunanda 13523086