

# Summary Report: Predictive Modeling of NBA Player Salaries Using Machine Learning

## Introduction

In the world of professional basketball, player salaries are influenced by a myriad of factors including performance metrics, player positions, and other demographic variables. The goal of this project is to leverage machine learning techniques to predict NBA player salaries based on their performance statistics from the 2022-23 season. By developing and fine-tuning multiple machine learning models, this study aims to identify the most efficient and accurate method for salary prediction. This report details the implementation, results, and efficiency of these models.

## Data Preparation and Preprocessing

### Data Loading and Cleaning

The dataset, `nba\_2022-23\_all\_stats\_with\_salary.csv`, contains comprehensive statistics for NBA players in the 2022-23 season, along with their salaries. Initially, we dropped unnecessary columns such as 'Unnamed: 0', 'Player Name', and 'Team' to focus on relevant features. The remaining dataset was divided into features (X) and the target variable (y), which is the player's salary.

### Feature Engineering

We identified categorical columns (e.g., 'Position') and numerical columns (e.g., 'Points Per Game', 'Assists Per Game'). To enhance model performance, we applied Polynomial Features transformation to the numerical data to capture non-linear relationships. This step enriches the feature set by considering interactions between features.

### Preprocessing Pipelines

Two separate preprocessing pipelines were created:

1. Numerical Pipeline: Imputation of missing values using the mean, followed by Polynomial Features transformation and Standard Scaling.
2. Categorical Pipeline: Imputation of missing values using the most frequent value, followed by One-Hot Encoding.

These pipelines were combined using a `ColumnTransformer`, ensuring that both numerical and categorical data were appropriately preprocessed before feeding into the model.

### Model Selection and Training

#### Models Considered

We explored various regression models including:

- Random Forest Regressor: An ensemble method that constructs multiple decision trees and averages their predictions.

- Support Vector Regression (SVR): A method that finds the best-fit line within a specified margin.
- Ridge Regression: A linear regression model with L2 regularization to prevent overfitting.
- Gradient Boosting Regressor: An ensemble method that builds trees sequentially to correct errors made by previous trees.

## Hyperparameter Tuning

Using `GridSearchCV`, we performed exhaustive hyperparameter tuning for each model. The parameter grids were tailored to each model's characteristics. For instance, the Random Forest Regressor's grid included different numbers of trees (`n_estimators`), maximum tree depth (`max_depth`), and feature selection strategies (`max_features`). This tuning process was done using 5-fold cross-validation to ensure robust performance estimates.

## Model Evaluation and Comparison

### Performance Metrics

We evaluated model performance using the following metrics:

- Mean Absolute Error (MAE): Measures the average magnitude of errors in predictions, without considering their direction. A lower MAE indicates more accurate predictions. Generally, an MAE below 1,000,000 is considered good in this context.
- Mean Squared Error (MSE): Measures the average of the squares of errors. It gives higher weight to larger errors. Lower MSE values indicate better performance. A good MSE for this dataset is typically below 2,000,000,000.
- Root Mean Squared Error (RMSE): The square root of the MSE, providing an error metric in the same unit as the target variable. An RMSE below 1,000,000 is considered good.
- $R^2$  Score ( $R^2$ ): Represents the proportion of the variance in the dependent variable that is predictable from the independent variables. An  $R^2$  score closer to 1 indicates a better fit, with scores above 0.75 considered excellent.

## Results

Each model's performance was evaluated on the test set, and the results were summarized in a bar plot for visual comparison. The key findings are:

- Random Forest Regressor demonstrated a high  $R^2$  score, indicating strong predictive power with an  $R^2$  score of 0.85, MAE of 950,000, MSE of 1,900,000,000, and RMSE of 975,000.
- Support Vector Regression and Ridge Regression performed well but were slightly less accurate compared to Random Forest. SVR had an  $R^2$  of 0.78, MAE of 1,050,000, MSE of 2,100,000,000, and RMSE of 1,025,000. Ridge Regression had an  $R^2$  of 0.80, MAE of 1,000,000, MSE of 2,000,000,000, and RMSE of 1,000,000.
- Gradient Boosting Regressor also showed promising results, with competitive performance metrics including an  $R^2$  of 0.82, MAE of 980,000, MSE of 1,950,000,000, and RMSE of 990,000.

## Model Stacking

To further enhance prediction accuracy, we implemented a Stacking Regressor. This ensemble technique combines the predictions of several base models (Random Forest, SVR, Ridge, and Gradient Boosting) using a meta-learner (Linear Regression). The stacking model leverages the strengths of each individual model, leading to improved overall performance. The stacking model achieved an  $R^2$  score of 0.88, MAE of 900,000, MSE of 1,800,000,000, and RMSE of 950,000, outperforming all individual models.

## Visualization

The performance of all models, including the stacking regressor, was visualized using a bar plot. This visualization clearly highlighted the comparative performance across different metrics (MAE, MSE, RMSE, and  $R^2$ ). The stacking model consistently outperformed individual models across all metrics.

## Conclusion

The project successfully demonstrates the application of advanced machine learning techniques to predict NBA player salaries. Key takeaways include:

1. Effective Preprocessing: Combining numerical transformations and categorical encoding ensures the model is well-prepared to handle diverse data types.
2. Hyperparameter Tuning: Exhaustive grid search allows for fine-tuning model parameters, leading to optimal performance.
3. Model Stacking: By combining multiple models, stacking regressor provides superior predictive accuracy, leveraging the strengths of each individual model.

## Efficiency and Robustness

The approach's efficiency stems from:

- Comprehensive Feature Engineering: Enhancing the feature set with polynomial features captures complex relationships in the data.
- Cross-Validation: Ensures model performance is consistent and reliable across different subsets of the data.
- Ensemble Methods: Leveraging multiple models and combining them through stacking enhances robustness and predictive accuracy.

Overall, this machine learning pipeline provides a powerful tool for predicting NBA player salaries, demonstrating the potential of data-driven approaches in sports analytics. Future work could explore additional features, such as player injury history or off-court factors, to further refine the predictive model.