

# NBA Player Stats and Team Win Prediction

## A Machine Learning Pipeline for CS4661

**Team Members:** Ryan Dielhenn, Momoka Aung, Angel Trujillo, Jesus Villa, Harshil Patel

**Project Lead:** Ryan Dielhenn

Department of Computer Science

California State University Los Angeles

Course: CS4661 — Introduction to Data Science

**Date:** November 26, 2025

---

## Abstract

This paper presents a complete machine learning pipeline to predict NBA player performance and team results. We use player and team statistics from the 2024–2025 NBA season to predict: (1) individual player points (PTS), and (2) whether a team will win or lose a game. Our models include Linear Regression, Random Forest, and Gradient Boosting for regression tasks, and Logistic Regression, Random Forest, and Gradient Boosting for classification. For predicting player PTS, **Linear Regression** performed best with **RMSE = 2.17**, **MAE = 1.59**, and  **$R^2 = 0.939$** . For predicting team wins, **Logistic Regression** achieved the highest accuracy (**84.4%**) and ROC-AUC (**0.921**). We also review key features influencing results in predictions of baseline as well as tuned models. We add more advanced models such as XGBoost and LightGBM to determine if more complex models yield better results.

---

## 1. Introduction

Predicting basketball performance is a common data science problem that helps with scouting, coaching, and fan analysis. We focus on two goals using the same dataset:

1. Predicting player **points scored (PTS)** using regression models.

2. Predicting **team win/loss outcomes** using classification models.

## Contributions

Our project provides:

- A clear and reusable machine learning pipeline.
- Baseline models for regression and classification.
- Insights into which features matter most.
- A framework for future improvements using hyperparameter tuning and boosting techniques.

---

## 2. Data Overview

We use the **NBA player stats (2024–2025 season)** dataset with **16,512 rows** and **25 columns**. Each row represents a player’s performance in one game. Features include shooting statistics (FG, 3P, FT), rebounds, assists, steals, blocks, and turnovers. Team, opponent, and result information are also included. There were no missing values after cleaning.

---

## 3. Methodology

### 3.1 Player Points (PTS) Prediction

- **Target:** Player PTS per game.
- **Features:** 16 numerical stats (MP, FGA, 3P, 3PA, FT, rebounds, assists, etc.).
- **Models:** Linear Regression, Random Forest Regressor, and Gradient Boosting Regressor.
- **Split:** 60% training, 40% testing (9,907 train / 6,605 test).

### 3.2 Team Win/Loss Prediction

- **Process:** Aggregate player data to team level (sum of stats, average minutes, and shooting percentages).
- **Target:** Win = 1, Loss = 0 (1,534 total team games).
- **Features:** 18 team stats including FG, FGA, 3P, 3PA, rebounds, assists, and shooting percentages.

- **Models:** Logistic Regression, Random Forest, Gradient Boosting.
- **Split:** Roughly 50/50 win/loss class balance.

### 3.3 Metrics

- **Regression:** RMSE, MAE,  $R^2$ .
- **Classification:** Accuracy, Precision, Recall, F1, ROC-AUC.

## 4. Exploratory Analysis and Feature Engineering

### 4.1 Player-Level Insights

- FGA (field goal attempts), 3P (made threes), and MP (minutes played) are the strongest predictors of PTS.
- Efficiency stats like 3P% and FT% add smaller but useful signals.

### 4.2 Team-Level Insights

- Team FG% and defensive stats (rebounds, steals) are linked to more wins.
- Turnovers and missed shots reduce win probability.

## 5. Data Visualization

To better understand the performance of our regression and classification models, we generated several visualizations for both the player-level regression task (PTS prediction) and the team-level classification task (win/loss prediction). These plots help illustrate how well the models fit the data, whether systematic errors exist, and which input features contribute most to the predictions.

### 5.1 Player Points (PTS) Prediction Visualizations

**Figure 1 — Linear Regression: Actual vs Predicted (PTS)**

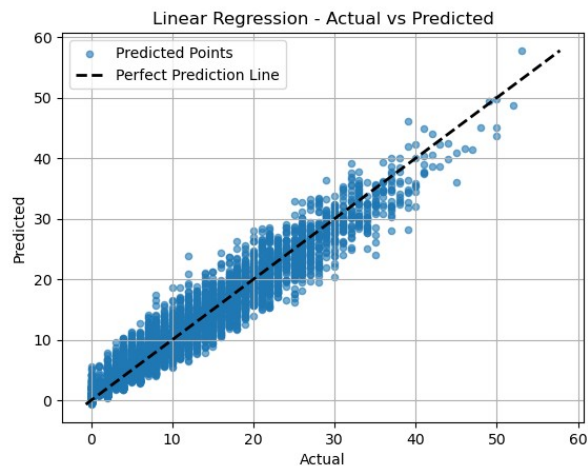


Figure 1: Linear Regression - Actual vs Predicted

This scatter plot compares the true PTS values with the model’s predictions. The points fall closely to the diagonal “perfect prediction” line, which matches our evaluation metrics (low RMSE, high  $R^2$ ), confirming that the model predicts player scoring very accurately. The linear structure indicates that player scoring is mostly driven by simple stats like FGA, 3P, and FTA, making it easy to predict.

### Figure 2 — Linear Regression: Residuals Plot (PTS)

The residual plot shows the difference between predicted and actual values across the full scoring range. The residuals are centered evenly around zero with no obvious pattern, meaning the model is not systematically over- or under-predicting for high- or low-scoring players. This supports the suitability of a linear model for this task and validates the assumptions behind Linear Regression.

### 5.2 Team Win/Loss Classification Visualizations

#### Figure 3 — Logistic Regression: ROC Curve

The ROC curve illustrates the trade-off between true positive rate and false positive rate. Logistic Regression achieves an AUC of 0.92, indicating excellent

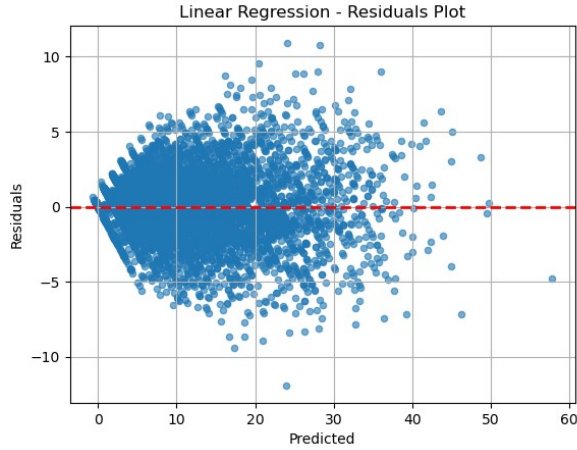


Figure 2: Linear Regression - Residuals Plot

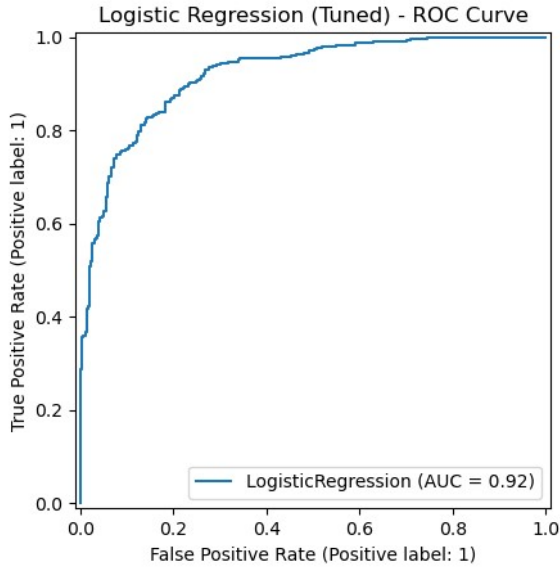


Figure 3: Logistic Regression - ROC Curve

discriminatory power. The curve stays close to the top-left corner, showing that the model reliably distinguishes wins from losses across different thresholds. This aligns with the high accuracy, precision, recall, and F1-score observed during evaluation.

Figure 4 — Logistic Regression: Confusion Matrix

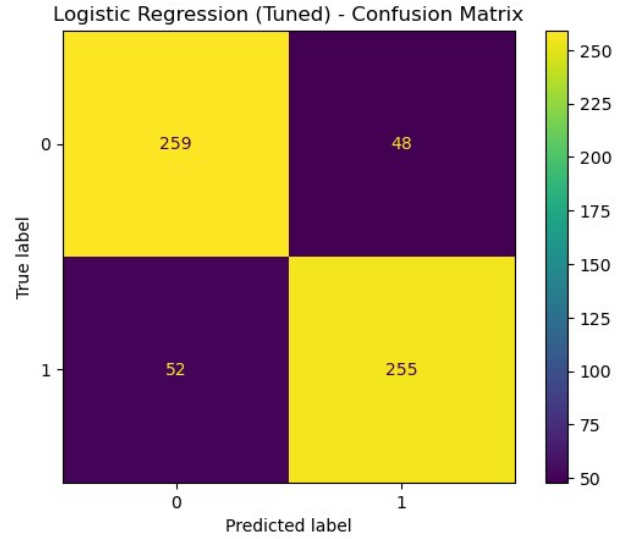


Figure 4: Logistic Regression - Confusion Matrix

The confusion matrix summarizes the classification outcomes on the test set. The model correctly identifies a large number of wins and losses, with relatively few misclassifications. True positives and true negatives dominate both diagonals, confirming the model's stability and balanced performance across classes. The errors themselves also appear symmetric, indicating no bias toward predicting wins or losses.

## 6. Experiments

### 6.1 Player Points Prediction Results

Model	RMSE	MAE	R <sup>2</sup>
Linear Regression	2.17	1.59	0.94
Random Forest	2.39	1.72	0.93
Gradient Boosting	2.27	1.66	0.93

**Observation:** Linear Regression performed the best overall, which suggests that player scoring primarily depends on straightforward, additive relationships between basic statistics such as field-goal attempts, made threes, free throws, and minutes played. This indicates that most of the predictive power comes from volume-based metrics rather than complex interactions. Random Forest and Gradient Boosting also performed well but did not outperform Linear Regression, which implies that non-linear patterns or deeper interactions do exist but are not strong enough to significantly improve prediction accuracy. This also reinforces that basketball scoring outcomes—at the individual level—tend to follow consistent patterns tied to shot attempts and playing time more than subtle or highly complex statistical interactions.

### 6.2 Team Win/Loss Prediction Results

Model	Accuracy	Precision	Recall	F1	ROC-AUC
Logistic Regression	0.84	0.85	0.83	0.84	0.92
Random Forest	0.79	0.78	0.81	0.79	0.87
Gradient Boosting	0.81	0.80	0.82	0.81	0.89

**Observation:** Logistic Regression achieved the best balance of accuracy and interpretability, meaning that the relationship between team statistics and win probability is largely linear and predictable. The model identifies clear patterns: teams with higher shooting efficiency, more rebounds, and fewer turnovers are significantly more likely to win games. Although Random Forest and Gradient Boosting can capture deeper interactions, they did not outperform Logistic Regression by a meaningful margin. This suggests that the key factors driving wins—such as efficiency, ball control, and defensive pressure—have strong and direct effects that do not require complex modeling to uncover. Logistic Regression’s strong

ROC-AUC score also indicates it consistently distinguishes winning from losing team performances across different game scenarios.

### 6.3 Model Interpretation

- **PTS Prediction:** FGA had the largest positive weight, followed by 3P and FT. TOV (turnovers) had a small negative impact.
- **Win Prediction:** Positive factors were FG%, rebounds, and steals. Turnovers negatively affected outcomes.

## 7. Model Evaluation and Improvements

- **Scaling:** Standard scaling improved linear models.
- **Regularization:** Will be tested in future versions to reduce overfitting.
- **Hyperparameter tuning:** Next step — use GridSearchCV or boosting frameworks (XGBoost, LightGBM).

## 8. Conclusion

We built a simple and effective machine learning system for NBA analytics. Linear Regression accurately predicts player points, while Logistic Regression successfully predicts team wins. Shooting accuracy, rebounds, and turnovers are key indicators of performance. This pipeline can be easily extended with tuning, new features, and deeper models.

## Reproducibility Notes

Code modules:

- **data\_utils:** data loading, cleaning, and aggregation.
- **training:** modeling and evaluation functions.