

# Technical Report: Final Project DS 5110: Introduction to Data Management and Processing

Team Members: Ashwin Satish, Joseph Russell  
Khoury College of Computer Sciences  
Data Science Program  
satish.ash@northeastern.edu, russell.jo@northeastern.edu

November 24, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Literature Review</b>	<b>2</b>
<b>3</b>	<b>Methodology</b>	<b>2</b>
3.1	Data Collection . . . . .	2
3.2	Data Preprocessing . . . . .	2
3.3	Analysis Techniques . . . . .	3
<b>4</b>	<b>Results</b>	<b>3</b>
<b>5</b>	<b>Discussion</b>	<b>3</b>
<b>6</b>	<b>Conclusion</b>	<b>3</b>
<b>7</b>	<b>References</b>	<b>4</b>
<b>A</b>	<b>Appendix A: Code</b>	<b>4</b>
<b>B</b>	<b>Appendix B: Additional Figures</b>	<b>4</b>

# 1 Introduction

As basketball fans, we thought it would be interesting to try to predict the outcomes of games using statistics. Using NBA data from the past 5 seasons, we gathered team and player metrics as well as historical game scores in order to determine how scores can be modeled from general statistics. Our process involved collecting multiple CSVs, compiling and preprocessing them in Python, performing exploratory data analysis and visualization, and ultimately building machine learning models to predict margins of victory for 6000 historical games.

# 2 Literature Review

Using basketball statistics to predict game outcomes involves a thorough understanding of how these statistics are recorded and their relevance to basketball analytics. While many of these are fairly intuitive, we researched advanced statistics that are intended to encapsulate overall player performance in different ways, such as Win Shares, PER (Player Efficiency Rating), and VORP (Value Over Replacement Player). Understanding the purpose of these metrics allows us to analyze players and teams more accurately and help build our model.

# 3 Methodology

The methods and techniques used in the project ranged from Data Collection, Data Preprocessing and Analysis for the same.

## 3.1 Data Collection

For collecting the appropriate team data, player data, match by match offensive and defensive data alongside the advanced player metrics which included some feature engineering, we referred to Basketball Reference (<https://basketball-reference.com>).

Our data collection method involved leveraging the site for all the information required and combining it into a single csv which had data for the past 5 seasons for the teams and players for each of the categories mentioned above.

## 3.2 Data Preprocessing

The steps we implemented for data preprocessing were

1. Handling missing values: We had many missing values present in our datasets, this is because the data given online was flawed and was not complete in nature. We had to replace numerical NaNs with appropriate means or medians where necessary.
2. Worked on correcting data inconsistencies: We had to standardize all values into appropriate formats for the model to understand. We implemented label encoding for the players, positions and teams. We also fixed data entry errors like incorrect player stats ( pts scored ; total shots attempted)
3. Duplicate records: There were several duplicate entries which caused errors in data collection/merging.

4. Encoding categorical data: We had to work on encoding categorical data like player names, positions and teams.
5. Scaling: Implemented robust scaling after trying and understanding how each scaling method affects our data accordingly.

### 3.3 Analysis Techniques

The analytical techniques and models used in the project are:

1. 3PT shoots v Win Margins: This visualization explains how 3 pointers affect win margins and win conditions in the game.
2. Points allowed v scored: This visualization explains how many points were scored by a team and how many were scored against the team.
3. Kmeans cluster: A cluster showing the two different categories or performances of players in the dataset. Implemented k using elbow method to find optimal clusters.

## 4 Results

The current results or numeric outputs from our models are as follows.

1. Silhouette Score: 0.621
2. Optimal k: 2
3. MS Error: 180.46
4. The cluster has a proper separation of the points, which shows that there is an appropriate basis for clustering the data
5. A proper correlation matrix for the data which allows us to understand the features and how they're correlated with each other.

## 5 Discussion

Our linear regression model provided a baseline mean squared error that we can work to improve. The metric will have a high floor as there is a lot of randomness in NBA games (good teams will often lose to bad teams unpredictably). We will try boosted trees next (XGBoost, LightGBM).

## 6 Conclusion

While NBA games can be very unpredictable, our model proves that team and player statistics can be used to understand to a degree how a game will turn out. Future models could try to incorporate time series data as game outcomes are often dependent on how teams and players perform in recent games more than older ones.

## 7 References

### References

## A Appendix A: Code

```
1 team_offense = pd.read_csv("drive/MyDrive/Colab Notebooks/SlamDunk/
    SlamDunkInsights_TeamData.csv")
2 team_defense = pd.read_csv("drive/MyDrive/Colab Notebooks/SlamDunk/
    TeamDefense.csv")
3 team_defense = team_defense.rename(columns={'season': 'Year'})
4 team_defense = team_defense.drop('Rk', axis=1)
5 team_defense['Team'] = team_defense['Team'].str.replace('*', '', regex=
    False)
6 team_offense['Team'] = team_offense['Team'].str.replace('*', '', regex=
    False)
7 team_defense = team_defense.drop(['G', 'MP'], axis=1)
8 team_offense = team_offense.drop(['G', 'MP'], axis=1)
9 team_df = pd.merge(team_offense, team_defense, on=['Team', 'Year'])
10 team_df.head()
```

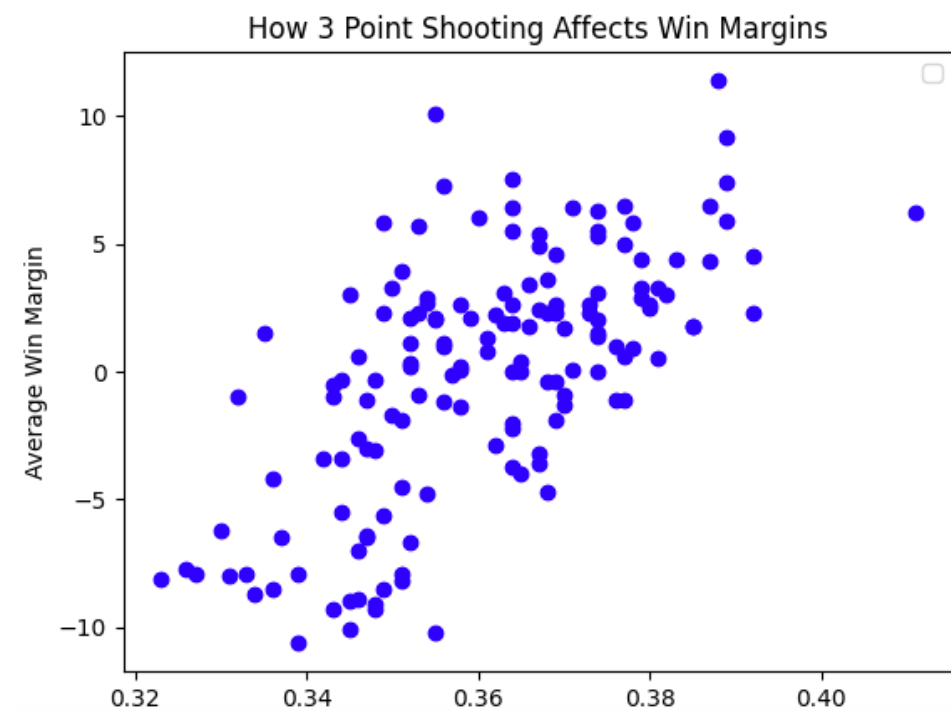
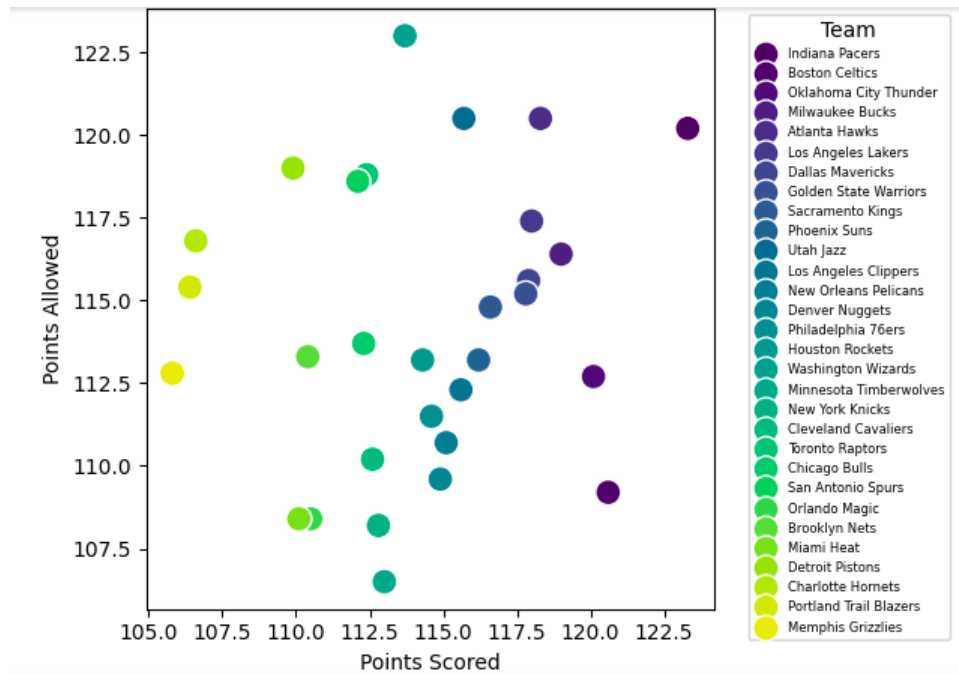
Listing 1: Preprocessing Team Data

```
1 from sklearn.decomposition import PCA
2
3 pca = PCA(n_components=3)
4 pca_data = pca.fit_transform(splayer_df[cts])
5
6 explained_variance_ratio = pca.explained_variance_ratio_
7 variance_df = pd.DataFrame(explained_variance_ratio, columns=['
    Explained Variance Ratio'], index=['PC1', 'PC2', 'PC3'])
8 print(variance_df)
9 cumulative_variance_ratio = pca.explained_variance_ratio_.cumsum()
10 print("Cumulative Explained Variance Ratio:")
11 print(cumulative_variance_ratio)
12 loadings = pd.DataFrame(pca.components_.T, columns=['PC1', 'PC2', 'PC3'
    ], index=cts)
13 print(loadings)
```

Listing 2: Dimensionality Reduction

## B Appendix B: Additional Figures

Include any additional figures or tables that support the analysis.



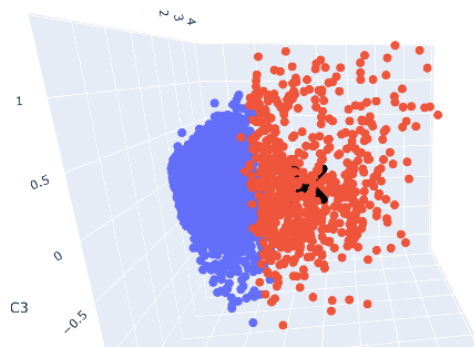


Figure 1: K-Means Cluster