

# Predicting the Next Defensive Player of the Year in the NBA

Joshua Samuel  
[jsamuel41@tamu.edu](mailto:jsamuel41@tamu.edu)

Alex LeGresley  
[alexlegresley@tamu.edu](mailto:alexlegresley@tamu.edu)

Texas A&M University

# Contents

<b>1</b>	<b>Importance of Study</b>	<b>1</b>
<b>2</b>	<b>Exploratory Data Analysis</b>	<b>1</b>
2.1	Description of Dataset	1
2.2	Data Cleaning	1
2.3	Summary Statistics	1
<b>3</b>	<b>Methodology</b>	<b>1</b>
3.1	Stage 1	2
3.2	Stage 2	2
<b>4</b>	<b>Stage 1 Model</b>	<b>2</b>
4.1	Initial Results	2
4.2	Tuning	2
4.3	Post-Tuned Results	3
4.4	Logistic Model Comparison	4
<b>5</b>	<b>Stage 2 Model</b>	<b>4</b>
5.1	Tuning	4
5.2	Post-Tuned Results	4
5.3	Linear Model Comparison	5
<b>6</b>	<b>Current Season Results</b>	<b>5</b>
<b>7</b>	<b>Final Conclusions &amp; Future Modifications</b>	<b>6</b>
<b>8</b>	<b>References</b>	<b>7</b>
<b>9</b>	<b>Code &amp; Data</b>	<b>7</b>
<b>10</b>	<b>Code Appendix</b>	<b>7</b>

# 1 Importance of Study

We seek to predict the winner of this year’s NBA’s Defensive Player of the Year Award, in addition to identifying other players who will be in contention for the award and how they will stack up. The Defensive Player of the Year (DPOY) award is one of the most coveted awards given for the performance of NBA players during the regular season, second only to the overall Most Valuable Player of the Year Award. With betting continuing to become more and more prevalent in sports, the amount of money won and lost regarding regular season awards continues to increase dramatically. Furthermore, player contracts often contain lucrative incentives for winning regular season awards, nevertheless the intangible boost to one’s legacy and status in the game. Thus, statistical prediction of the DPOY award is of potential importance to a variety of parties.

## 2 Exploratory Data Analysis

### 2.1 Description of Dataset

The dataset used in this analysis was aggregated from various collections provided by Basketball Reference. The data was obtained in CSV and Excel format. It comprises of four primary datasets spanning the NBA seasons from 1982-1983 to 2022-2023. The first dataset encompasses standard statistics (such as blocks, steals, etc.) of all active players who participated in at least one game during this period. The second dataset consists of advanced metrics (including box plus/minus, win shares, etc.) of the same players. The third dataset, covering the same time frame, lists every player who received at least one vote for the Defensive Player of the Year award in each respective season. The fourth and final dataset adds the respective team standings position for every player in the dataset.

### 2.2 Data Cleaning

Duplicates were removed from each dataset, specifically instances where a player switched teams within a season. Instead of retaining separate observations for each team, the statistics for such players were aggregated for the entire season. The team rank associated with players who played for multiple teams during a season was taken to be the average of the teams ranks they played for weighted by the number of minutes played for each team. Ultimately, each row of the dataset was a unique combination of player and season. A few columns provided in the Basketball-Reference dataset did not have meaningful associated names and thus were dropped since they could not be confidently regarded as a known metric. This led to a total of 49 meaningful covariance and one response variable.

While the details of the statistical models will be covered in detail later, data setup was necessary to be congruent with their limitations. All but one of the covariates was quantitative, with the player position being multi-class categorical. The XGBoost implementation provides built-in support for handling categorical variables, however the Scikit-Learn implementation of linear models does not. Dummy variables for player position were created for use with the linear models. Furthermore, many of the quantitative covariates had missing values. The XGBoost implementation also provides built-in support for handling missing values, however the Scikit-Learn linear models implementation does not. Median imputation was used to fill the missing values for use with the linear models. We chose median imputation due to the median being more robust to skewness, as we anticipated some of the covariates would be skewed.

### 2.3 Summary Statistics

The distribution of the share of Defensive Player of the Year Award votes is highly right-skewed, with a vast majority of players having a vote share of zero. There is nearly a 30:1 ratio of players who did not receive a DPOY award vote to those who received at least one vote. This important characteristic of the data is accounted for in our methodology.

## 3 Methodology

To address the significant imbalanced nature of the voting for the DPOY award due to very few players receiving any votes, we will employ a two stage process.

### 3.1 Stage 1

The first stage model aims to classify whether a player is likely to receive any votes for the Defensive Player of the Year award. We will build the stage 1 model first with XGBoost, then compare the performance with a simpler logistic regression model. For this model, we will transform the response variable, which is the DPOY vote share for each player in each season, to a binary outcome: 1 if they received any share, that is their vote share is non-zero, and zero otherwise. The classes will be very unbalanced, with many zeros and a few ones. We will balance the classes in the statistical model training process by adjusting the class weights such that the observations of players who received a vote will be weighted higher. Specifically, the class weights will correspond to a 30:1 ratio as found during the exploratory data analysis process.

Model performance will be evaluated using the metric area under the precision recall curve (AUCPR). This metric is advantageous in scenarios with significant class imbalance since it will penalize false negatives. This is necessary because a model that only predicts zero would have a high accuracy but ultimately not be useful. We will use cross-validation when evaluating potential models, as it provides a statistically more reliable estimation of the performance than a single train/test split. Specifically, we will use stratified k-fold cross-validation. Stratification is necessary due to the class imbalance, as it helps ensure an approximately equal representation of classes in each group compared to standard k-fold cross-validation. Furthermore, the model will output probabilities of belonging to the positive class, that is receiving any vote share, and it is up to us to determine whether the probability is high enough to expect the player to receive any vote share.

We will tune this probability threshold using the F1 score metric, which is also a robust metric in cases of class imbalance. Ultimately, the stage 1 model will act as a filter, removing those who are unlikely to receive any votes. We can then take those who we expect to receive at least one vote and inspect on a more granular level to determine the specific amount of vote share.

### 3.2 Stage 2

The second stage model predicts the amount of share in the Defensive Player of the Year Award for players classified in the previous model as projected to receive a share. We will build the stage 2 model with XGBoost, then compare the performance with a simpler linear regression model. The player predicted to have the highest share in a given season is expected to win the award.

We will work with the response variable as a continuous quantitative variable and while utilizing the same covariates as the first model. Model performance will be evaluated using the metrics root mean squared error (RMSE) and  $R^2$ . RMSE is a robust standardized metric that indicates how well the model is able to predict the response variable, while  $R^2$  indicates how well the covariates explain the variation in the response. Furthermore, we will use cross-validation when evaluating potential models. Specifically, k-fold cross-validation since the response is quantitative.

## 4 Stage 1 Model

XGBoost classifier models were fit. The dataset consists of 17900 observations (all players appearing in at least one game from the 1982-1983 to 2022-2023 seasons) with 49 covariates (48 quantitative, 1 categorical) and 1 binary response. The class weights were adjusted to account for the imbalance and stratified cross-validation employed with evaluation by AUCPR. The built-in implementation handling of categorical predictors and missing observations was also used.

### 4.1 Initial Results

We initially fit the model with default hyperparameters. This yielded an average AUCPR from cross-validation of 0.478 and an average F1 score from cross-validation of 0.469. These values suggest that the model is similar to a coin flip in terms of determining whether a player is likely to have received at least one DPOY vote, and thus there is much room for improvement in the predictive performance of the model.

### 4.2 Tuning

We turned to tuning the hyperparameters of the model as a source of potential improvement. We identified the XGBoost hyperparameters learning rate, maximum depth of trees, and number of estimators as being

most likely to result in significantly enhanced predictive performance. We chose a randomized search technique due to the balance between computational efficiency and ability to search a wide space of potential values for the hyperparameters. Specifically, randomized search helps find optimal hyperparameter values by randomly sampling with replacement from a given statistical distribution of hyperparameter values. This randomized search process was ran with 75 iterations of stratified 5-fold cross-validation. The hyperparameters leading to the model with the maximum AUCPR were considered optimal.

### 4.3 Post-Tuned Results

The average AUCPR from cross-validation for our best model was 0.609 with the best parameters being 0.153 for the learning rate, 11 for the max tree depth, and 162 for the number of estimators. We next tuned the probability decision threshold for considering a player to receive at least one DPOY vote. Stratified 5-fold cross-validation was used to test thresholds in increments of 0.05 with evaluation by F1 score. The average best threshold from cross-validation was 0.35 with a corresponding average F1 score from cross-validation of 0.581.

The feature importance chart [Figure 1](#) shows the relative importance of each feature in the model. Features with higher importance have more influence on the model's predictions, while those with lower importance have less. It is not surprising to see defensive metrics such as Personal Fouls (PF), Defensive Box Plus Minus (DBPM), and Defensive Win Share (DWS) near the top in terms of influence. However, it is surprising to see Steal % (STL%) with the least influence, as this is a commonly cited defensive metric.

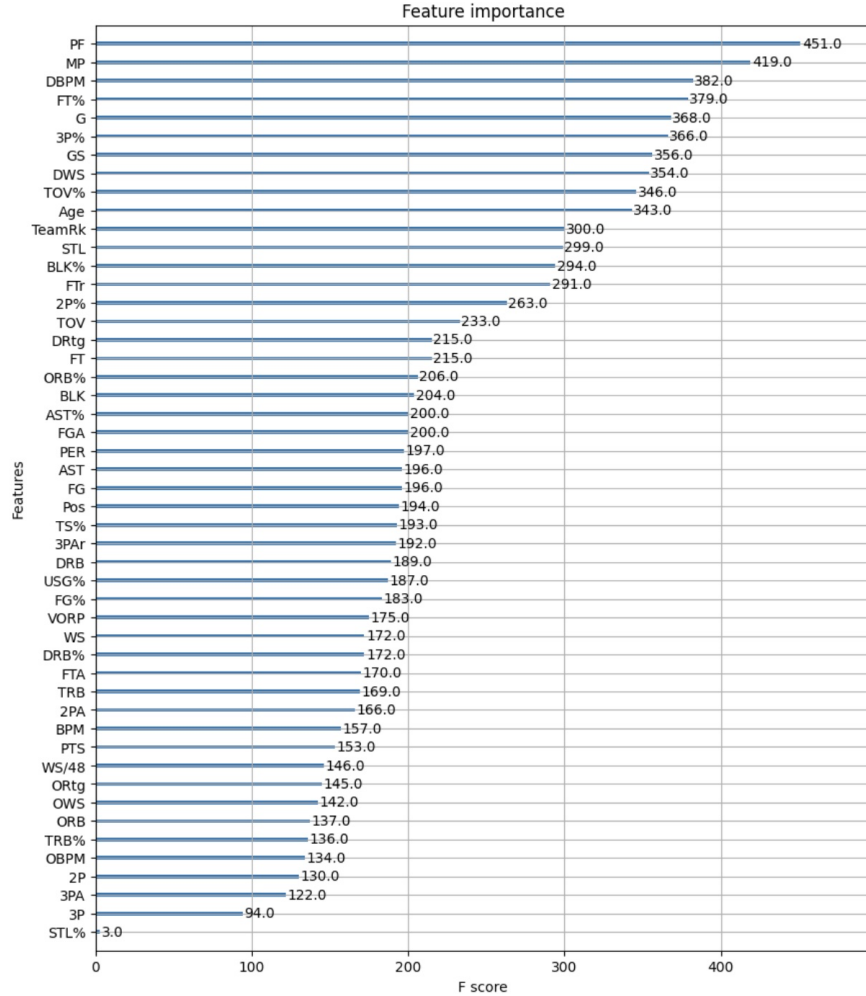


Figure 1: Stage 1 feature importance for tuned XGBoost classifier model

## 4.4 Logistic Model Comparison

To investigate whether a complex model such as XGBoost is necessary, we compared the performance to a simpler logistic regression model. Class weights were adjusted to account for the imbalance. Dummy variables were created for the categorical Position covariate and missing data for quantitative covariates was imputed using the median. Random search with 75 iterations of stratified 5-fold cross-validation was used to tune the C hyperparameter, which controls the aggressiveness of the regularization process. The optimal model was that which maximized AUCPR.

The best logistic regression model had an average AUCPR from cross-validation of 0.496. The probability decision threshold was optimized by employing stratified 5-fold cross-validation to test thresholds in increments of 0.05 with evaluation by F1 score. The maximum average F1 score from cross-validation was 0.524 achieved at a threshold of 0.90.

The AUCPR and F1 score for the logistic regression model was slightly lower than those of the XGBoost model. Thus, we conclude that the XGBoost model has slightly better predictive performance.

## 5 Stage 2 Model

XGBoost regressor models were fit. The dataset consists of 575 observations (players receiving any vote share from the 1982-1983 to 2022-2023 seasons) with the same 49 covariates (48 quantitative, 1 categorical) as the stage 1 model and 1 quantitative response. K-fold cross-validation and evaluation by RMSE was used to select the best model. The built-in implementation handling of categorical predictors and missing observations was also used.

### 5.1 Tuning

We tuned the same XGBoost hyperparameters as in stage 1, those being learning rate, max tree depth, and number of estimators. Randomized search was used with 75 iterations of 5-fold cross-validation with evaluation by -RMSE. The use of -RMSE rather than RMSE is an implementation detail consequence of the Scikit-Learn Unified Scoring API. The interpretation is the same, that is a score closer to zero represents less prediction error and is desirable. The model that minimized -RMSE was considered optimal.

### 5.2 Post-Tuned Results

The average -RMSE from cross-validation for our best model was -0.153 with the best parameters being 0.097 for the learning rate, 3 for the max tree depth, and 94 for the number of estimators. Furthermore, the average  $R^2$  from cross-validation was 0.292. This means that our model explains only about 29.2% of the variability in DPOY vote share, and thus there is definitely room for improvement.

The feature importance chart [Figure 2](#) shows the relative importance of each feature in the XGBoost model. Advanced defensive metrics such as Defensive Win Shares (DWS) and Defensive Box Plus Minus (DBPM) were considered the two most important features which is not surprising. However, Blocks (BLK) and Defensive Rebounding Percentage (DRB%) being low in terms of importance is surprising given that they are considered by many to be important defensive statistics. Field Goal Percentage (FG%) being ranked 7th in importance illustrates how the efficiency of a player on the offensive end still plays a factor in whether they win the award.

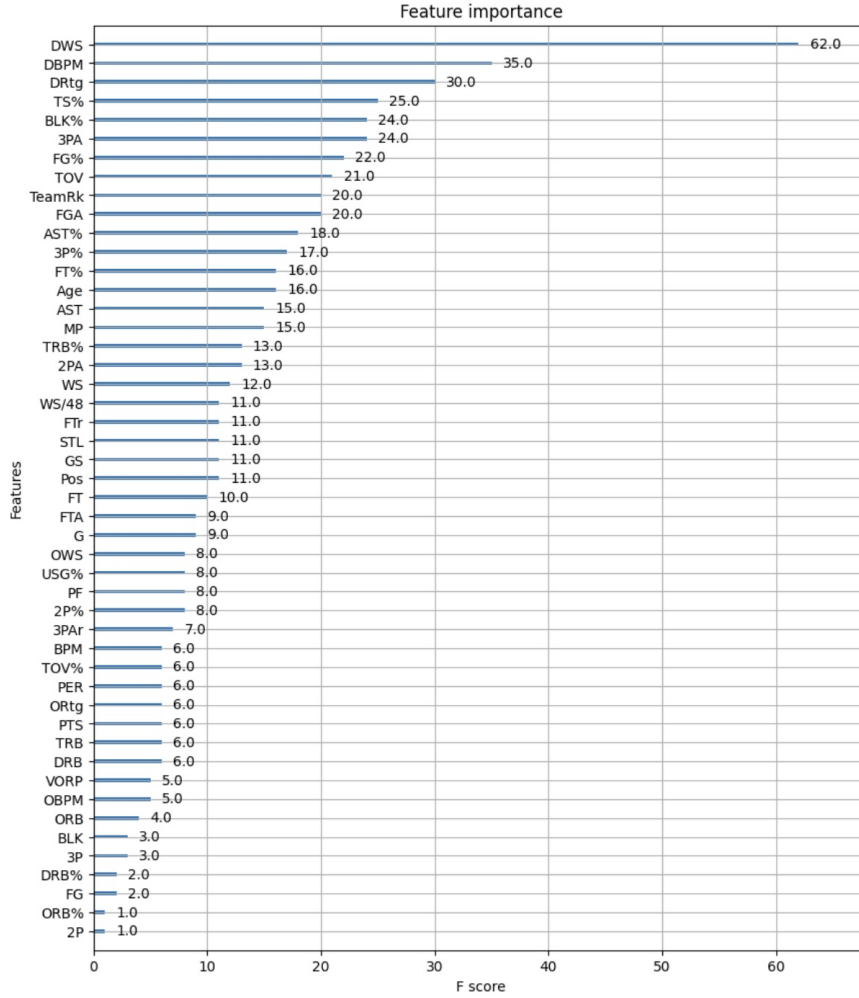


Figure 2: Stage 2 feature importance for XGBoost regressor model

### 5.3 Linear Model Comparison

To investigate whether a complex model such as XGBoost is again necessary, we compared the performance to a simpler linear regression model. Dummy variables were created for the categorical Position covariate and missing data for quantitative covariates was imputed using the median. Reverse Feature Elimination (RFE) was used to identify the optimal number of predictors to include in the model. This was run with 75 iterations of 5-fold cross-validation and the best model identified based on minimizing -RMSE.

The best linear regression model had an average -RMSE from cross-validation of -0.167 and an average  $R^2$  from cross-validation of 0.200. Thus, the linear regression model has slightly higher prediction error and explains slightly less of the variability in vote share compared to the XGBoost regressor model.

## 6 Current Season Results

In line with our goal, we tested the stage 1 and stage 2 XGBoost models and the stage 1 and stage 2 linear models with the current NBA season data to predict this year's Defensive Player of the Year. The current season has concluded so all data is available, however the award has not yet been given out. The predictions using the stage 1 and stage 2 XGBoost models are shown in [Table 1](#), the stage 1 and stage 2 linear models in [Table 2](#), and the current Vegas betting odds from the DraftKings Sportsbook as of 4/17/2024 in [Table 3](#).

Our XGBoost process performed quite well as the top 4 DPOY candidates according to Vegas are accounted for in the top 8 for our preferred process. Our XGBoost's top 2 candidates are the same as Vegas however

the order is different. Our linear model process had a few of the same projected candidates as our XGBoost process. However, there are two very significant outliers at the top of the linear process’s ranking, as background knowledge of the NBA season reinforces that Markquis Nowell and Hamidou Diallo are in reality nowhere close to being in consideration for the award. The two players currently play in the development league (G League) of their respective teams and will occasionally be called up to play for the actual NBA team.

Table 1: XGBoost predicted DPOY vote share

<b>Place</b>	<b>Player</b>	<b>Vote Share</b>
1	Victor Wembanyama	0.345
2	Rudy Gobert	0.283
3	Chet Holmgren	0.185
4	Kristaps Porzingis	0.182
5	Anthony Davis	0.144
6	Giannis Antetokounmpo	0.137
7	Shai Gilgeous-Alexander	0.135
8	Bam Adebayo	0.062
9	Nic Claxton	0.008

Table 2: Linear model predicted DPOY vote share

<b>Place</b>	<b>Player</b>	<b>Vote Share</b>
1	Markquis Nowell	0.587
2	Hamidou Diallo	0.351
3	Rudy Gobert	0.321
4	Chet Holmgren	0.152
5	Victor Wembanyama	0.147
6	Anthony Davis	0.143
7	Shai Gilgeous-Alexander	0.142
8	Isaiah Hartenstein	0.131
9	Josh Hart	0.130
10	Nikola Jokic	0.124
11	Derrick White	0.116
12	Herbert Jones	0.110
13	Shaquille Harrison	0.108
14	Giannis Antetokounmpo	0.089
15	Brook Lopez	0.080
16	Kentavious Caldwell-Pope	0.073
17	Jrue Holliday	0.069

Table 3: DraftKings Sportsbook DPOY odds as of 4/17/2024

<b>Player</b>	<b>Odds</b>
Rudy Gobert	-3,000
Victor Wembanyama	+1,000
Bam Adebayo	+20,000
Anthony Davis	+20,000
Jarret Allen	+20,000
Derrick White	+20,000

## 7 Final Conclusions & Future Modifications

Our two-stage XGBoost model process predicts with reasonable accuracy which players will be near the top of the DPOY award voting based off data from a given regular season. We have also been able to identify which



metrics contribute most (and least) significantly in predicting how much of a share in DPOY votes a player will receive. We have seen how our preferred tuned XGBoost models outperform our more elementary linear models. In the future, we would consider taking into account the years of NBA experience a player has, since typically the DPOY award is not given to rookies. While we did have an "age" predictor, an additional "years of experience" predictor would be helpful given players enter the league at various ages.

A limitation of our model is that the final predictions are not restricted such that the vote shares sum to one. This is an issue if the specific vote share itself is highly desired to be accurate. In our case, we cared more about the relative rankings of the players in terms of the order they will finish in the voting. In the future, we would like to refine the model to address this limitation and thus produce more realistic vote shares.

## 8 References

- [1] Sports Reference LLC. Basketball-Reference.com - Basketball Statistics and History. <https://www.basketball-reference.com/>.
- [2] Salah, AJ. "2024 NBA Defensive Player of the Year Odds - Gobert Huge Favorite." Covers, 17 Apr. 2024. <https://www.covers.com/nba/nba-defensive-player-of-the-year-odds>.

## 9 Code & Data

The data files and a Python notebook containing the code used to produce all of the analyses with explanatory code comments can be obtained from <https://github.com/alexlegresley/STAT482-Capstone>. The code is also listed below in the [Code Appendix](#).

## 10 Code Appendix

```
1  # -----
2  # Imports
3
4  import pathlib
5  import os
6  import pandas as pd
7  import numpy as np
8  import functools
9  import operator
10
11 import scipy.stats as stats
12 from matplotlib import pyplot as plt
13
14 import xgboost as xgb
15 from sklearn.model_selection import RandomizedSearchCV, StratifiedKFold, KFold,
16     ↪ cross_val_score
17 from sklearn.metrics import f1_score
18 from sklearn.linear_model import LogisticRegressionCV, LinearRegression
19 from sklearn.feature_selection import RFECV
20
21 # -----
22 # Utility functions
23
24 def computeWeightedAverage(row):
25     combos = basicStats.loc[(basicStats["Season"] == row["Season"]) & (basicStats["Tm"] !=
26     ↪ "TOT") & (basicStats["Player"] == row["Player"])]
```

```

26     rkKeys = zip(combos["Season"], combos["Tm"])
27     weights = combos["MP"] / combos["MP"].sum()
28     return np.average([rkMapping[k] for k in rkKeys], weights = weights)
29
30
31 def optimizeThreshold(model, X, y):
32     best = []
33     for trainIndex, testIndex in StratifiedKFold(n_splits = 5).split(X, y):
34         mod = model.fit(X.iloc[trainIndex], y.iloc[trainIndex])
35         res = []
36         for th in np.arange(0.05, 1, 0.05):
37             yPred = mod.predict_proba(X.iloc[testIndex])[:, 1] > th
38             res.append((th, f1_score(y.iloc[testIndex], yPred)))
39             best.append(max(res, key = operator.itemgetter(1)))
40
41     return np.average([th for th, _ in best]).round(2), np.average([f1 for _, f1 in best])
42
43
44 def withDummyPos(X):
45     posFeatures = pd.get_dummies(X.Pos, drop_first = True, prefix = "pos")
46     return pd.concat([X, posFeatures], axis = 1).drop("Pos", axis = 1).apply(lambda x:
47     ↪ x.fillna(x.median()) if pd.api.types.is_numeric_dtype(x) else x)
48
49 def createModelData(basic_stats, adv_stats, standings):
50     basic = pd.read_csv(basic_stats)
51     adv = pd.read_csv(adv_stats)
52     stand = pd.read_csv(standings, header = 1)
53
54     basic.loc[:, "Player"] = basic["Player"].map(lambda x: x.rstrip("*"))
55     adv.loc[:, "Player"] = adv["Player"].map(lambda x: x.rstrip("*"))
56
57     tmMapping = {
58         "Atlanta Hawks": "ATL",
59         "Boston Celtics": "BOS",
60         "Brooklyn Nets": "BRK",
61         "Charlotte Hornets": "CHO",
62         "Chicago Bulls": "CHI",
63         "Cleveland Cavaliers": "CLE",
64         "Dallas Mavericks": "DAL",
65         "Denver Nuggets": "DEN",
66         "Detroit Pistons": "DET",
67         "Golden State Warriors": "GSW",
68         "Houston Rockets": "HOU",
69         "Indiana Pacers": "IND",
70         "Los Angeles Clippers": "LAC",
71         "Los Angeles Lakers": "LAL",
72         "Memphis Grizzlies": "MEM",
73         "Miami Heat": "MIA",
74         "Milwaukee Bucks": "MIL",
75         "Minnesota Timberwolves": "MIN",
76         "New Orleans Pelicans": "NOP",
77         "New York Knicks": "NYK",
78         "Oklahoma City Thunder": "OKC",

```

```

79         "Orlando Magic": "ORL",
80         "Philadelphia 76ers": "PHI",
81         "Phoenix Suns": "PHO",
82         "Portland Trail Blazers": "POR",
83         "Sacramento Kings": "SAC",
84         "San Antonio Spurs": "SAS",
85         "Toronto Raptors": "TOR",
86         "Utah Jazz": "UTA",
87         "Washington Wizards": "WAS"
88     }
89     stand["Tm"] = stand["Team"].map(tmMapping)
90     rkMapping = dict(zip(stand["Tm"], stand["Rk"]))
91
92     def computeWeightedAverage(row):
93         combos = basic.loc[(basic["Player"] == row["Player"]) & (basic["Tm"] != "TOT")]
94         weights = combos["MP"] / combos["MP"].sum()
95         return np.average([rkMapping[k] for k in combos["Tm"]], weights = weights)
96
97     basic["TeamRk"] = [rkMapping[rowKey] if (rowKey := row["Tm"]) in rkMapping.keys() else
98         ↪ computeWeightedAverage(row) for _, row in basic.iterrows()]
99
100     basic.drop_duplicates(subset = "Player", keep = "first", inplace = True)
101     adv.drop_duplicates(subset = "Player", keep = "first", inplace = True)
102
103     basic.drop(["Rk", "Tm", "Unnamed: 29", "Player-additional"], axis = 1, inplace = True)
104     adv.drop(["Rk", "Pos", "Age", "Tm", "G", "MP", "Unnamed: 19", "Unnamed: 24",
105         ↪ "Player-additional"], axis = 1, inplace = True)
106
107     combined = pd.merge(basic, adv, how = "left", on = "Player")
108     combined["Pos"] = combined["Pos"].astype("category")
109
110     return combined
111
112 def predictVoteShare(data, stage1, threshold, stage2):
113     probSomeShare = stage1.predict_proba(data.drop("Player", axis = 1))[:, 1]
114
115     noShare = data.loc[~(probSomeShare > threshold)].copy()
116     noShare["Vote_Share"] = 0
117
118     someShare = data.loc[probSomeShare > threshold].copy()
119     someShare["Vote_Share"] = stage2.predict(someShare.drop("Player", axis = 1))
120
121     return pd.concat([someShare[["Player", "Vote_Share"]], noShare[["Player",
122         ↪ "Vote_Share"]]], ignore_index = True).sort_values("Vote_Share", ascending =
123         ↪ False).reset_index(drop = True)
124
125 def matchDummiesPos(X, XTrain):
126     presentDummies = pd.get_dummies(X, prefix = "pos")
127     res = pd.DataFrame()
128     for d in [d for d in pd.get_dummies(XTrain, prefix = "pos").columns if
129         ↪ d.startswith("pos_")]:
130         if d in presentDummies.columns:

```

```

128         res[d] = presentDummies[d]
129     else:
130         res[d] = False
131     res.drop(columns = res.columns[0], axis = 1, inplace = True)
132     return pd.concat([X, res], axis = 1).drop("Pos", axis = 1).apply(lambda x:
133         ↪ x.fillna(x.median()) if pd.api.types.is_numeric_dtype(x) else x)
134
135 # -----
136 # Data cleaning
137
138 basicStats = pd.concat(
139     (pd.read_csv(f).assign(Season = os.path.basename(f).removesuffix(".csv")) for f in
140     ↪ (pathlib.Path()/ "Data"/ "Basic_Stats").glob("*.csv")),
141     ignore_index = True
142 )
143 advStats = pd.concat(
144     (pd.read_csv(f).assign(Season = os.path.basename(f).removesuffix(".csv")) for f in
145     ↪ (pathlib.Path()/ "Data"/ "Advanced_Stats").glob("*.csv")),
146     ignore_index = True
147 )
148 dpoyVoting = pd.concat(
149     (pd.read_excel(f, header = 1).assign(Season =
150     ↪ os.path.basename(f).removesuffix(".xlsx")) for f in
151     ↪ (pathlib.Path()/ "Data"/ "DPOY_Voting").glob("*.xlsx")),
152     ignore_index = True
153 )
154
155 expStand = pd.concat(
156     (pd.read_csv(f, header = 1).assign(Season = os.path.basename(f).removesuffix(".csv"))
157     ↪ for f in (pathlib.Path()/ "Data"/ "Standings").glob("*.csv")),
158     ignore_index = True
159 )
160
161 teamMappings = {
162     "Atlanta Hawks": "ATL",
163     "Boston Celtics": "BOS",
164     "Brooklyn Nets": "BRK",
165     "Charlotte Bobcats": "CHA",
166     "Charlotte Hornets": "CHO",
167     "Chicago Bulls": "CHI",
168     "Cleveland Cavaliers": "CLE",
169     "Dallas Mavericks": "DAL",
170     "Denver Nuggets": "DEN",
171     "Detroit Pistons": "DET",
172     "Golden State Warriors": "GSW",
173     "Houston Rockets": "HOU",
174     "Indiana Pacers": "IND",
175     "Kansas City Kings": "KCK",
176     "Los Angeles Clippers": "LAC",
177     "Los Angeles Lakers": "LAL",
178     "Memphis Grizzlies": "MEM",
179     "Miami Heat": "MIA",
180     "Milwaukee Bucks": "MIL",
181     "Minnesota Timberwolves": "MIN",

```

```

176     "New Jersey Nets": "NJN",
177     "New Orleans Hornets" : "NOH",
178     "New Orleans Pelicans": "NOP",
179     "New Orleans/Oklahoma City Hornets": "NOK",
180     "New York Knicks": "NYK",
181     "Oklahoma City Thunder": "OKC",
182     "Orlando Magic": "ORL",
183     "Philadelphia 76ers": "PHI",
184     "Phoenix Suns": "PHO",
185     "Portland Trail Blazers": "POR",
186     "Sacramento Kings": "SAC",
187     "San Antonio Spurs": "SAS",
188     "San Diego Clippers": "SDC",
189     "Seattle SuperSonics": "SEA",
190     "Toronto Raptors": "TOR",
191     "Utah Jazz": "UTA",
192     "Vancouver Grizzlies": "VAN",
193     "Washington Bullets": "WSB",
194     "Washington Wizards": "WAS"
195 }
196 expStand["Tm"] = expStand["Team"].map(teamMappings)
197 rkMapping = dict(zip(zip(expStand["Season"], expStand["Tm"]), expStand["Rk"]))
198
199 basicStats.loc[:, "Tm"] = basicStats["Tm"].apply(lambda x: "CHO" if x == "CHH" else x)
200 advStats.loc[:, "Tm"] = advStats["Tm"].apply(lambda x: "CHO" if x == "CHH" else x)
201
202 basicStats["TeamRk"] = [rkMapping[rowKey] if (rowKey := (row["Season"], row["Tm"])) in
203     ↪ rkMapping.keys() else computeWeightedAverage(row) for _, row in basicStats.iterrows()]
204
205 basicStats.drop_duplicates(subset = ["Season", "Player"], keep = "first", inplace = True)
206 advStats.drop_duplicates(subset = ["Season", "Player"], keep = "first", inplace = True)
207
208 basicStats.loc[:, "Player"] = basicStats["Player"].map(lambda x: x.rstrip("*"))
209 advStats.loc[:, "Player"] = advStats["Player"].map(lambda x: x.rstrip("*"))
210
211 basicStats = basicStats.drop(["Rk", "Tm", "Unnamed: 29", "Player-additional"], axis = 1)
212 advStats = advStats.drop(["Rk", "Pos", "Age", "Tm", "G", "MP", "Unnamed: 19", "Unnamed:
213     ↪ 24", "Player-additional"], axis = 1)
214 dpoyVoting = dpoyVoting[["Season", "Player", "Share"]]
215
216 combinedData = functools.reduce(
217     lambda left, right: pd.merge(left, right, how = "left", on = ["Season", "Player"]),
218     [basicStats, advStats, dpoyVoting]
219 )
220 combinedData["Share"] = combinedData["Share"].fillna(0)
221
222 # -----
223 # Stage 1
224
225 model1Data = combinedData.drop(columns = ["Season", "Player"])
226 model1Data["Pos"] = model1Data["Pos"].astype("category")
227 model1Data["ShareBinary"] = [1 if x != 0 else 0 for x in model1Data["Share"]]

```

```

228 X = model1Data.drop(["Share", "ShareBinary"], axis = 1)
229 y = model1Data["ShareBinary"]
230
231 classWeightRatio = y.value_counts()[0] / y.value_counts()[1]
232
233
234 # XGBoost classifier
235
236 mod = xgb.XGBClassifier(
237     scale_pos_weight = classWeightRatio,
238     eval_metric = "aucpr",
239     enable_categorical = True
240 )
241 paramsDist = {
242     "learning_rate": stats.uniform(0.01, 0.20),
243     "max_depth": stats.randint(3, 12),
244     "n_estimators": stats.randint(50, 350)
245 }
246 rsMod = RandomizedSearchCV(
247     estimator = mod,
248     param_distributions = paramsDist,
249     scoring = "average_precision",
250     n_iter = 75,
251     cv = StratifiedKFold(n_splits = 5),
252     refit = True
253 )
254 rsMod.fit(X, y)
255 print(f"Best parameters: {rsMod.best_params_}")
256 print(f"Average cross validation score (AUCPR) from best model: {rsMod.best_score_}")
257 tunedMod = rsMod.best_estimator_
258
259 _, ax = plt.subplots(figsize = (10, 12))
260 xgb.plot_importance(tunedMod, ax = ax)
261
262 bestThXGB, bestF1XGB = optimizeThreshold(tunedMod, X, y)
263 print(f"Average best threshold: {bestThXGB}")
264 print(f"Average best F1 score: {bestF1XGB}")
265
266
267 # Logistic regression
268
269 modLR = LogisticRegressionCV(
270     class_weight = "balanced",
271     scoring = "average_precision",
272     solver = "newton-cholesky",
273     cv = StratifiedKFold(n_splits = 5),
274     refit = True
275 )
276 modLR.fit(withDummyPos(X), y)
277 print(f"Average cross validation score (AUCPR): {np.average(modLR.scores_[1])}")
278
279 bestThLR, bestF1LR = optimizeThreshold(modLR, withDummyPos(X), y)
280 print(f"Average best threshold: {bestThLR}")
281 print(f"Average best F1 score: {bestF1LR}")

```

```

282
283
284 # -----
285 # Stage 2
286
287 model2Data = model1Data.drop(["ShareBinary"], axis = 1)
288 model2Data = model2Data.loc[model2Data["Share"].gt(0)]
289 model2Data.reset_index(drop = True, inplace = True)
290
291 X = model2Data.drop(["Share"], axis = 1)
292 y = model2Data["Share"]
293
294
295 # XGBoost regressor
296
297 mod2 = xgb.XGBRegressor(
298     eval_metric = "rmse",
299     enable_categorical = True
300 )
301 paramsDist = {
302     "learning_rate": stats.uniform(0.01, 0.20),
303     "max_depth": stats.randint(3, 12),
304     "n_estimators": stats.randint(25, 300)
305 }
306 rsMod2 = RandomizedSearchCV(
307     estimator = mod2,
308     param_distributions = paramsDist,
309     scoring = "neg_root_mean_squared_error",
310     n_iter = 75,
311     cv = KFold(n_splits = 5),
312     refit = True
313 )
314 rsMod2.fit(X, y)
315 print(f"Best parameters: {rsMod2.best_params_}")
316 print(f"Best cross validation score (neg RMSE): {rsMod2.best_score_}")
317 tunedMod2 = rsMod2.best_estimator_
318 print(
319     f"Average cross validation score (r2): "
320     f"{np.average(cross_val_score(tunedMod2, X, y, scoring = "r2", cv = KFold(n_splits =
321     ↪ 5)))}"
322 )
323
324 _, ax = plt.subplots(figsize = (10, 12))
325 xgb.plot_importance(tunedMod2, ax = ax)
326
327 # Linear regression
328
329 modLR2 = RFECV(
330     estimator = LinearRegression(),
331     min_features_to_select = 1,
332     step = 1,
333     scoring = "neg_root_mean_squared_error",
334     cv = KFold(n_splits = 5)

```

```

335 )
336 modLR2.fit(withDummyPos(X), y)
337 print(
338     f"Average cross validation score (negative RMSE): "
339     f"{np.average(modLR2.cv_results_["mean_test_score"])}"
340 )
341 print(
342     f"Average cross validation score (r2): "
343     f"{np.average(cross_val_score(modLR2, withDummyPos(X), y, scoring = "r2", cv =
    ↪ KFold(n_splits = 5)))}"
344 )
345
346
347 # -----
348 # Current year prediction
349
350 currentYear = createModelData(
351     basic_stats =
    ↪ (pathlib.Path()/"Data"/"Current_Year"/"2023-2024_Basic_Stats.csv").resolve(),
352     adv_stats = (pathlib.Path()/"Data"/"Current_Year"/"2023-2024_Adv_Stats.csv").resolve(),
353     standings = (pathlib.Path()/"Data"/"Current_Year"/"2023-2024_Standings.csv").resolve()
354 )
355
356
357 # XGBoost models
358
359 resultsXGB = predictVoteShare(data = currentYear, stage1 = tunedMod, threshold = bestThXGB,
    ↪ stage2 = tunedMod2)
360 resultsXGB.loc[resultsXGB["Vote_Share"].gt(0)]
361
362
363 # Linear models
364
365 resultsLR = predictVoteShare(
366     data = matchDummiesPos(currentYear, model1Data),
367     stage1 = modLR,
368     threshold = bestThLR,
369     stage2 = modLR2
370 )
371 resultsLR.loc[resultsLR["Vote_Share"].gt(0)]

```