



Don't DIY PII.

Isolate, protect, and govern sensitive data.

Try

DATA SCIENCE

Predicting the outcome of games with Machine Learn

How we used (and you can too) machine learning to k the role statistics play in sports.

Josh Weiner

Jan 5, 2021 19 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

When deciding on a final project for our Big Data class, my partners Jack Rosener, Jackson Joffe and I had an interest in sports with the principles learned this semester. After a few days of discussion, we set out to determine what would aim to predict the outcome of NBA games. In implementing our goal, we found it helpful to divide into the following steps with the following questions:

[Submit an Article](#)

- 1. Scraping Relevant Data** – where do we gather team and player statistics over several seasons?
- 2. Cleaning and Processing the Data** – how can we combine our scraped data so that it is both usable?

3. **Feature Engineering** – what additional metrics can we append to our datasets that would help any user or ML model to better understand and predict, respectively,



Keep PII out of breaches

Treat PII differently: Isolate. Protect. Govern.

Re:

relations within the data that may better inform our predictions?

5. **Predictions** – which models and features were useful for us in developing an accurate prediction focus on team or aggregated individuals' statistics

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Before we get into the nitty-gritty of our workflow, it's a good moment to review and acknowledge the other work done on this exact topic. First, in 2013 [Renato Torres from the University of Wisconsin-Madison](#) set out to accomplish a similar goal: to predict specific season outcomes of NBA games using different machine learning models. He used several of the same features that are included in our project, primarily feature reduction techniques to remove multicollinearity from the available data, and also explored different models to explore those with the highest accuracy. Like our project, his selected features included team statistics, but unlike our project had a particular focus on home and away winning percentage at home and away. (We will explore this type of analysis later.) A lot of other fantastic work on this topic has been done before and can be read here:

[Cheng, Ge & Zhang, Zhenyu & Kyebambe, Moshi & Kimbugwe. \(2016\). Predicting the Outcome of NBA Games Based on the Maximum Entropy Principle. Journal of Quantitative Analysis in Sports, 12\(2\), 1–12.](#)

[Predicting the Outcomes of NBA Games. North Carolina State University.](#)

University. Fayad, Alexander. Building My First Machine Learning Model | NBA Prediction Algorithm. Towards Data Science. The Complete History of the NBA. FiveThirtyEight.

Through our research, we found that the best published model had a prediction accuracy of 74.1% (for playoff outcomes), with most others achieving an upper bound between 66–72% accuracy. Most published research, too, focused on predicting playoff scores – which may lend towards biased teams are more consistent in a number of stats regular season, and playoff game expected outcome experience less variance as a result. Critically, the upset rate across the entire season in the NBA is higher than the playoffs, the upset rate – as defined by team regular season record winning— drops to 22% (which means most NBA-playoff prediction models usually get it wrong). Because our project looked to predict the outcome of a game and is playoff agnostic, we were looking to build a model that could reach and hopefully beat a 67.9% accuracy by doing so predict some upsets.

Feel free to follow along here or view our files on GitHub.

Scraping our Data

We scraped our data from available information — which has extremely detailed team and player statistics for every game played since the 2008–2009 season. Scraping the data took several days due to rate-limiting (as we had to query the NBA API).

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

game over 12 seasons), and was initially compiled into a JSON format and finally saved to csv files.

Author's Note: We accessed data from Synergy Sports with funding from the University of Pennsylvania. As an unfortunate consequence of this, we are not at liberty to provide public access to our dataset. However, we have compiled a list of alternative datasets through which one can replicate and improve upon our data: <https://www.basketball-reference.com/leagues>

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Cleaning the Data

We now had both player stats and team stats for each season saved as separate csv files. Our next step was to merge all this data and combine it into two large dataframes containing the player stats for the past 12 seasons and one containing the team stats. Once created, we would use dataframes to remove invalid statistics (negative values, columns that served little purpose to us (charge taken/committed, for example)).

We then saved these new dataframes to the folder which allowed us (and you) to skip the laborious steps of both scraping and cleaning the data within a notebook's runtime:

Search this file...

	Date	Game_ID	Season	H_Team	A_Team	H_Team_Record	
1	0	2020-08-14	732811	2019 - 2020	Houston Rockets	Philadelphia 76ers	44 - 27
2	1	2020-08-14	732814	2019 - 2020	Toronto Raptors	Denver Nuggets	52 - 19
3	2	2020-08-14	732812	2019 - 2020	Indiana Pacers	Miami Heat	44 - 28
4	3	2020-08-13	732808	2019 - 2020	Orlando Magic	New Orleans Pelicans	32 - 40
5	4	2020-08-13	732809	2019 - 2020	Phoenix Suns	Dallas Mavericks	33 - 39
6	5	2020-08-13	732806	2019 - 2020	Los Angeles Lakers	Sacramento Kings	52 - 18
7	6	2020-08-13	732805	2019 - 2020	Brooklyn Nets	Portland Trail Blazers	35 - 36
8	7	2020-08-13	732807	2019 - 2020	Memphis Grizzlies	Milwaukee Bucks	33 - 39
9	8	2020-08-13	732810	2019 - 2020	Utah Jazz	San Antonio Spurs	43 - 28
10	9	2020-08-13	732804	2019 - 2020	Boston Celtics	Washington	

NBA_Combined_Team_Stats.csv hosted with ❤ by GitHub

LATEST

EDITOR'S PICKS

Search this file...

	PlayerName	PlayerTeam	OpposingTeam	Date
1	PJ Tucker	Houston Rockets	Philadelphia 76ers	202
2	Hamidou Diallo	Oklahoma City Thunder	Los Angeles Clippers	202
3	Danilo Gallinari	Oklahoma City Thunder	Los Angeles Clippers	202
4	Nerlens Noel	Oklahoma City Thunder	Los Angeles Clippers	202
5	Steven Adams	Oklahoma City Thunder	Los Angeles Clippers	202
6	Devon Hall	Oklahoma City Thunder	Los Angeles Clippers	202
7	Kevin Hervey	Oklahoma City Thunder	Los Angeles Clippers	202
8	Dennis Schroder	Oklahoma City Thunder	Los Angeles Clippers	202
9	Andre Roberson	Oklahoma City Thunder	Los Angeles Clippers	202
10	Terrance Ferguson	Oklahoma City Thunder	Los Angeles Clippers	202

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

NBA_Combined_Player_Stats.csv hosted with ❤ by GitHub

Feature Engineering

This is where the fun stuff begins. Our primary goal is to make all of the available data understandable: game-level statistics for an entire team don't help us much unless we can extract data in a higher-level analysis that leads us to our goal of predicting wins and losses. To that end, we sought to find different features which we would use in understanding how teams progressed and regressed throughout each season.

1. Elo Ratings

This is perhaps the best existing method to relativize NBA team strength and performance over many seasons. The way Elo Ratings are calculated is simple: all teams start at a median score of 1500 and are either given or subtracted points based on the final score of each game and where it was played with weights being given to point difference, upsets, and location. In essence, it's a more sophisticated win-loss record. LATEST prediction models don't look at Elo Ratings but amalgamate a simple win-loss record with several factors. We wanted to use Elo to appropriately weight quality losses), while also recognizing that not all teams are equal. The exact formula is as follows: If R_i is the initial rating of a team, its Elo rating after its played its games is defined as follows:

$$R_{i+1} = k * (S_{team} - E_{team} + R_i)$$

We calculate ELO for each team and each game for every season of

Here, S_{team} is a state variable: 1 if the team wins, 0 if it loses. E_{team} represents the expected win probability for the team, which is represented as:

$$E_{team} = \frac{1}{1 + 10^{\frac{opp_elo - team_elo}{400}}}$$

k is a moving constant, dependent on both the win probability and difference in Elo ratings:

[LATEST](#)
[EDITOR'S PICKS](#)
[DEEP DIVES](#)
[NEWSLETTER](#)
[WRITE FOR TDS](#)
[Sign in](#)
[Submit an Article](#)

$$k = 20 \frac{(MOV_{winner} + 3)^{0.8}}{7.5 + 0.006(elo_difference_{winner})}$$

It's also important to note that Elo ratings carry over from season-to-season (as all teams are not created equal, good teams tend to stay good or at least gradually decline – very rarely do teams drop onto or off of the map). If R represents the final Elo of a team in one season, it's Elo Rating at the beginning of the next season is approximately:

$$(R * 0.75) + (0.25 * 1505)$$

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



Here we can actually see that Elo Ratings track quite well with teams' performances in years in which Golden State and Cleveland appeared and dueled in the NBA Finals. We can also see what was widely confirmed by most people at the time: that the Western Conference was much tougher than the East – as exhibited in the wins on Elo for the Warriors versus the Cavaliers. We can also see how far these championship seasons and as they all suffered from roster losses and inj

2. Recent Team Performance (Avg. stats over 10 most recent games)

These are pretty self-explanatory, we are simply looking to average the stats for each team over their last 10 games. To do this we wrote a simple function that would calculate a sliding average for a given team's stats and a window of n games:

```

1  #given a team and a date, this method will return that team's LATEST
2
3  def get_avg_stats_last_n_games(team, game_date, season_team):
4      prev_game_df = season_team_stats[season_team_stats['Date'] == game_date]
5
6      h_df = prev_game_df.iloc[:, range(3, 43, 2)]
7      h_df.columns = [x[2:] for x in h_df.columns]
8      a_df = prev_game_df.iloc[:, range(4, 44, 2)]
9      a_df.columns = [x[2:] for x in a_df.columns]
10
11     df = pd.concat([h_df, a_df])
12     df = df[df['Team'] == team]
13     df.drop(columns = ['Team'], inplace=True)
14
15     return df.mean()
16
17 recent_performance_df = pd.DataFrame()
18
19 for season in team_stats['Season'].unique():
20     l = ['Date', 'Game_ID', 'Season', 'H_Team', 'A_Team']
21     other = list(team_stats.columns[9:47])
22     cols = l + other
23
24     season_team_stats = team_stats[team_stats['Season'] == season]
25
26     season_recent_performance_df = pd.DataFrame()
27
28     for index, row in season_team_stats.iterrows():
29         game_id = row['Game_ID']
30         game_date = row['Date']
31         h_team = row['H_Team']
32         a_team = row['A_Team']
33

```

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

```

34     h_team_recent_performance = get_avg_stats_last_n_games(h_team, game_date, season)
35     h_team_recent_performance.index = ['H_Last_10_Avg_' + x for x in h_team_recent_perfor
36
37     a_team_recent_performance = get_avg_stats_last_n_games(a_team, game_date, season)
38     a_team_recent_performance.index = ['A_Last_10_Avg_' + x for x in a_team_recent_perfor
39
40     new_row = pd.concat([h_team_recent_performance, a_team_recent_performance], sort=False)
41     new_row['Game_ID'] = game_id
42
43     season_recent_performance_df = season_recent_performance_df.append(new_row, ignore_index=True)
44     season_recent_performance_df = season_recent_performance_df[new_row.index]
45
46
47     recent_performance_df = pd.concat([recent_performance_df, new_row])

```

LATEST

EDITOR'S PICKS

DEEP DIVES

After saving this data into a new dataframe, we separate each game (which contains stats for both home and away teams) into its own rows by team, which allows us to group-by and aggregate team stats much more easily. This simplifies existing features. Finally, we added a 'Wins' column to include the most critical measure of a team's performance: wins and losses.

NEWSLETTER

WRITE FOR TDS

Sign in

[Submit an Article](#)

3. Recent Player Performance (Avg. stats over games)

We create our player_recent_performance dataframe using similar methods to the above section, this time for individual players as opposed to teams. This creates a dataset showing each player's performance over the past 10 games.

4. Player Season Performance

We also sought to include the average player statistics for the season: unlike teams, players themselves get injured.

and out of the rotation and it's perhaps more critical for us to understand how player performances in individual games track with their averages. We will use this later in our models to see if it will allow for accurate predictions on the team-level.

5. Player Efficiency Ratings (PER)

Critically, just as we had done with teams via Elo Ratings, we wanted to be able to relativize player performance that combines seemingly unrelated statistics. One could use [Hollinger's Player Efficiency Rating](#) to predict team performance by the aggregated PER of its players. In the NBA, it is easy for players to experience inflated or deflated statistics (such as points per game) by virtue of the amount of playing time they get, minutes played, or versus starters, number of games played, or outlier performances. We did not want to rely solely on averages simply because of their ability to skew the problem by weighting certain in-game statistics such as number of minutes played, which creates a metric that does not reflect player performance relative to the number of minutes played. Thus for each player, we added a column for PER according to the following formula:

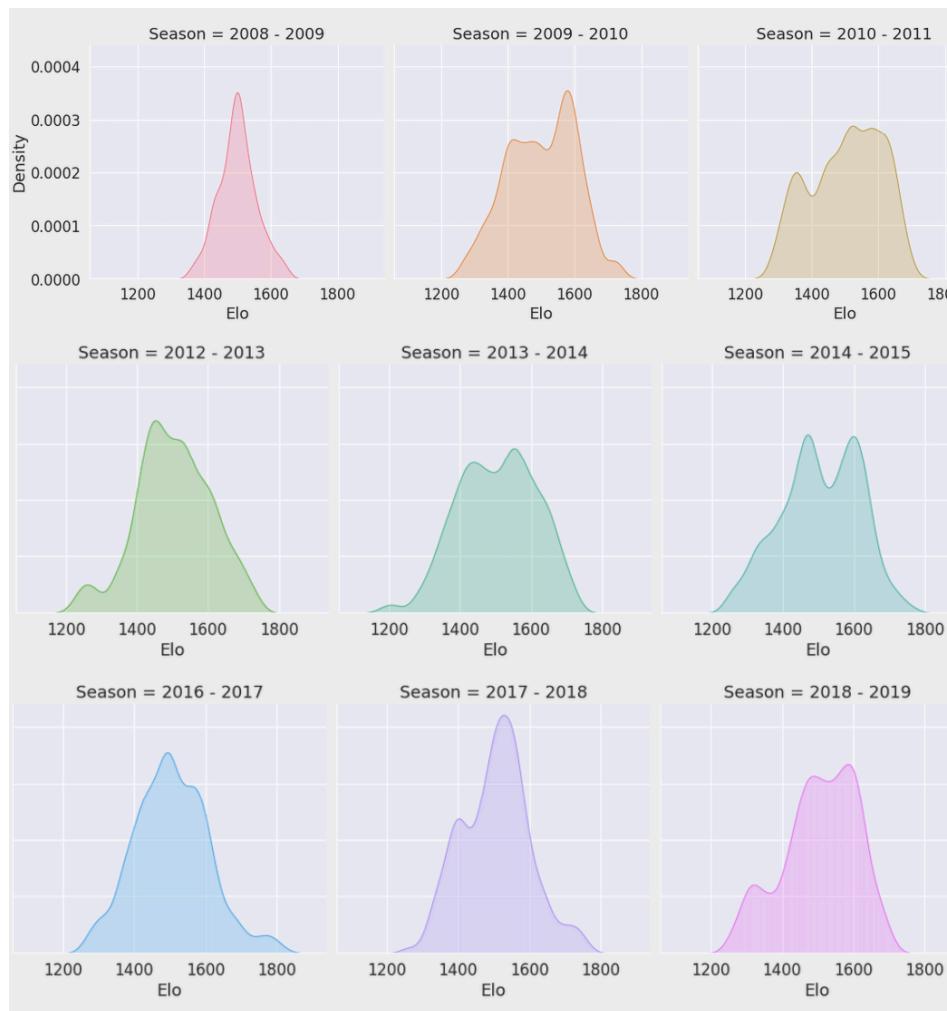
$$\begin{aligned} \text{PER} = & (\text{FGM} \times 85.910 + \text{Steals} \times 53.897 + 3\text{PTM} \\ & 46.845 + \text{Blocks} \times 39.190 + \text{Offensive_Reb} \times 39.190 \\ & 34.677 + \text{Defensive_Reb} \times 14.707 - \text{Foul} \times 17.174 \\ & 20.091 - \text{FG_Miss} \times 39.190 - \text{TO} \times 53.897) \times (1, \end{aligned}$$

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[NEWSLETTER](#)[WRITE FOR TDS](#)[Sign in](#)[Submit an Article](#)

Data Analysis

Our data analysis was centered around the use of Elo Ratings as our test metric. Essentially, could we be confident that Elo correlates with and correctly aggregates other statistics? Furthermore, would it be more appropriate for us to predict game outcomes using team stats (Elo Ratings) or averaged player stats (PER Ratings)?

First, let's explore the density of Elo Ratings across the NBA on a per-season basis. This tells us a little about the distribution of Elo Ratings across the league: if we can see Elo Ratings app well-matched. Otherwise, we see large disparities in development of "super-teams".

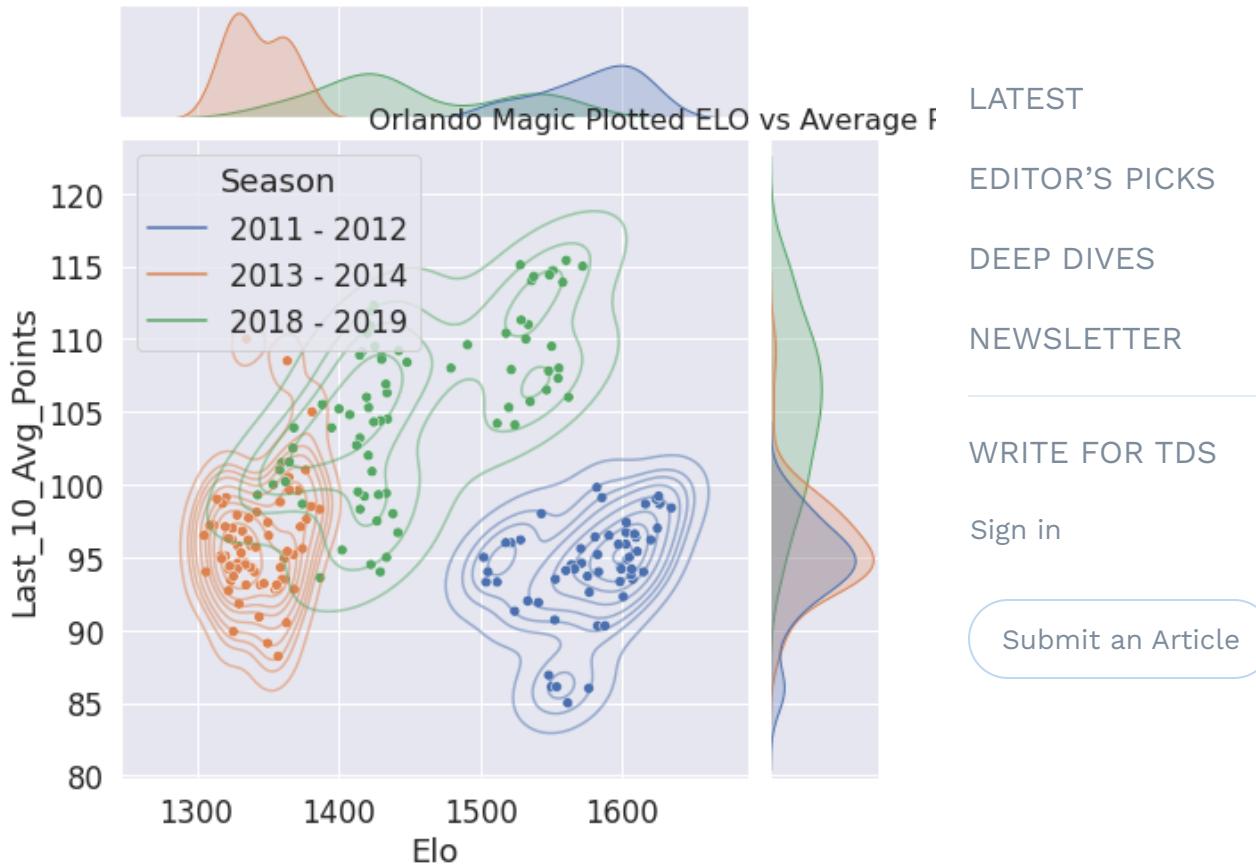


Twelve seasons of league Elo densities. (Image by Aut

[LATEST](#)
[EDITOR'S PICKS](#)
[DEEP DIVES](#)
[NEWSLETTER](#)
[WRITE FOR TDS](#)
[Sign in](#)
[Submit an Article](#)

Moving away from an understanding of Elo Ratings from a league-perspective, we endeavored to see how Elo Ratings tracked against an individual team's performance in other statistics.

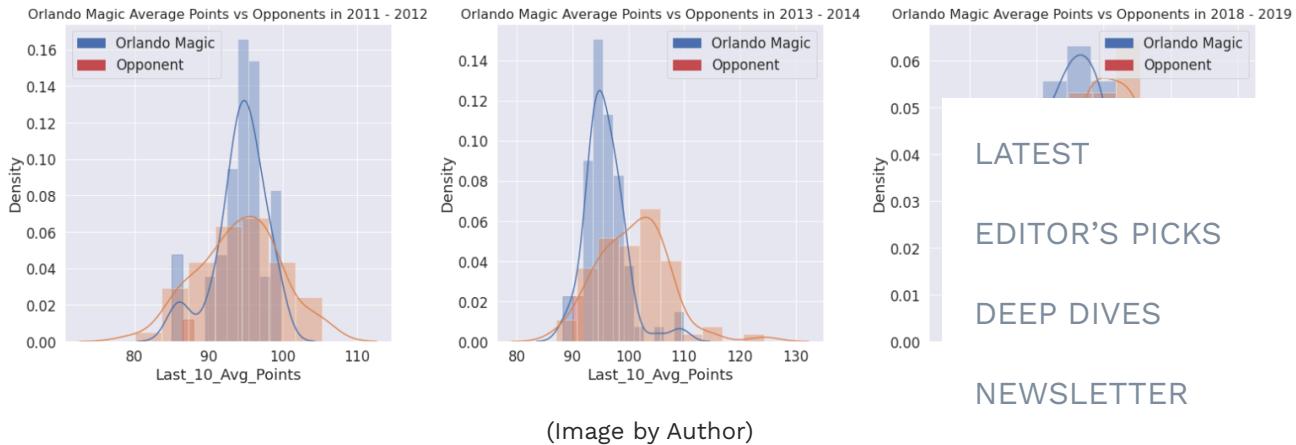
First, we looked to plot the distribution of Elo for a random team against the average number of points scored in recent games:



(Image by Author)

We can actually see from this that there is some correlation between the average number of points a team scores and its Elo Rating – the higher the average points scored over a window of games the higher the Elo Rating seen. However, we can also see that the Elo may also have variance across similar scoring figures. So, to better understand how Elo Ratings track with points scored, we examine how the average points scored compares to season averages.

league – from there we can determine if points scored improves Elo, provided that high scoring is relative to the rest of the league. To do this, let's look at that same team for the same seasons and plot the distribution of points scored against its opponents.



This confirms our suspicions, as we can see that the distribution of average points is greater than the opponents, or is more concentrated at an equal level when the Elo is higher for those seasons. When the Elo is lower for those seasons, the average points scored are even or lesser value, those seasons' Elo rating for the team are lower. Therefore, average points scored is a solid determinant of predicting game outcomes, ***relativized***. This demonstrated for us that Elo was a better determinant in predicting wins for us than by design a relative statistic.

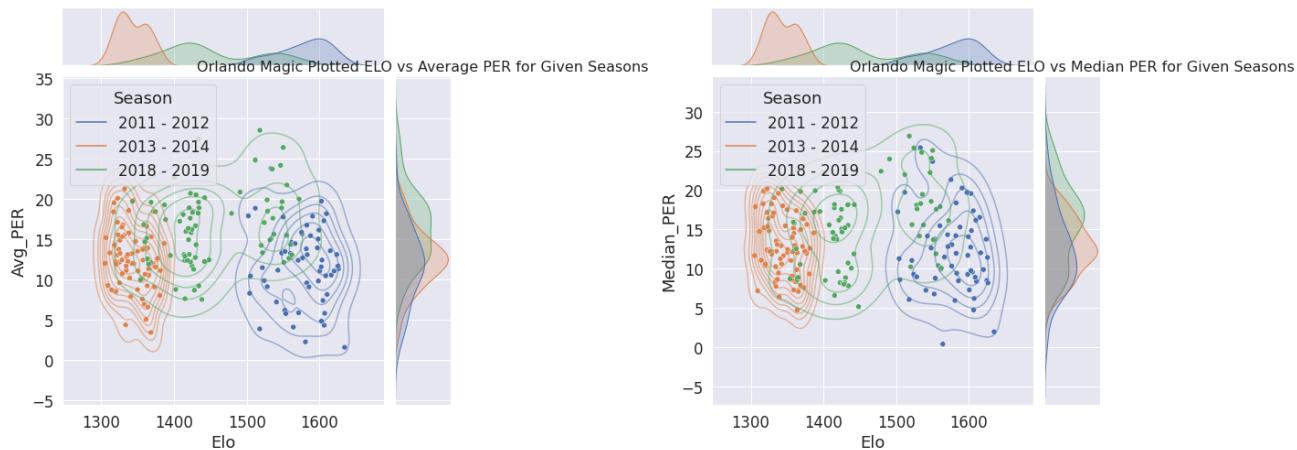
Shifting away from team statistics, we sought to see if Elo tracks better with player performance than team performance. To do this, we took a similar approach and plotted Elo Ratings with average points scored for a random team, this time with PER.

LATEST
EDITOR'S PICKS
DEEP DIVES
NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

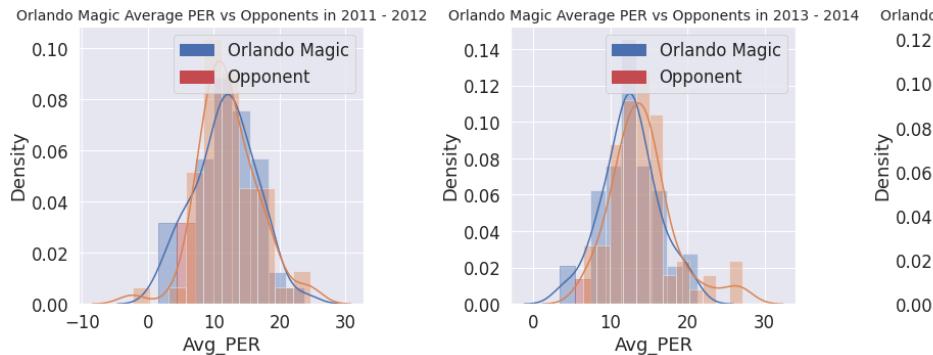


(Image by Author)

LATEST

From the plotted data, we can see that aggregate PER compared to opponents doesn't show much of a correlation with the strength of a team as determined by Elo. Points scored translates better – which makes sense as a player's efficiency isn't necessarily tied to scoring points – and points scored against opponents is a factor of winning a game and therefore impacting Elo.

We can see this further by mapping the Orlando Magic and median PER ratings against its opponents for three seasons, and find that there is almost no relationship between PER averages or medians and team strength.



From these and the above plotted distribution, we see that Average PER – while Elo Ratings throughout the season – generally show us little about how individual team strength when tracked against opponents. (Image by

EDITOR'S PICKS

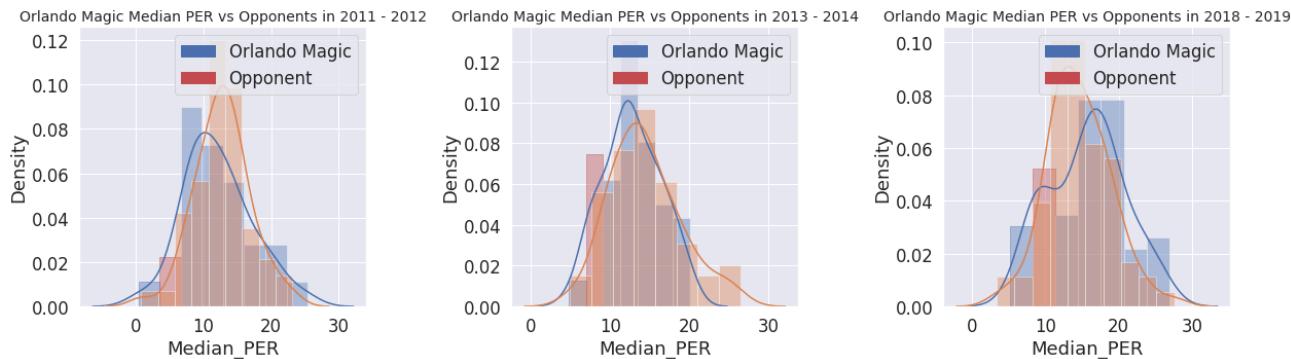
DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article



Median PER ratings show even less of a correlation to Elo Ratings throughout a given season. Here we can observe that in winning seasons (2011–12), the Orlando Magic had a lower median PER than its opponents for most games, yet they had their highest Elo Ratings and best record in recent history.

[LATEST](#)
[EDITOR'S PICKS](#)
[DEEP DIVES](#)
[NEWSLETTER](#)
[WRITE FOR TDS](#)
[Sign in](#)
[Submit an Article](#)

Predicting the Outcome of Games Based on Team Statistics and Elo Ratings

Our first step here was to split our data into feature and target columns. Reading from our dataset, once split, we used sklearn to randomly split our data into train and test sets in an 80:20 ratio.

```

1  from sklearn.model_selection import train_test_split
2
3  final_team_stats = pd.read_csv('/content/drive/MyDrive/CIS5'
4
5  cols = [0,1,2,3,4,7,8, 18, 37]
6  final_team_stats.drop(final_team_stats.columns[cols],axis=1
7
8  features = final_team_stats.drop(columns = 'Label')
9  label = final_team_stats['Label']
10
11 x_train, x_test, y_train, y_test = train_test_split(features, label)

```

model_prep.py hosted with ❤ by GitHub

The first model we aimed to use to predict the outcome of an NBA game was a Logistic Regression model. Unlike a Linear Regression model which predicts outcomes on a range of values between (and sometimes outside) 0 and 1, Logistic Regression models aim to group predictions into binary outcomes. Since we are predicting wins and losses, this type of classification suits us perfectly.

To begin, we used a simple non-parameterized I LATEST
team stats and Elo Ratings as parameters using

After playing around with some hyperparameter that using max_iter=131 and verbose=2 slightly increased initial testing accuracy to 66.95%. Definitely not parameterized model and very close to our desired accuracy. However, we sought to see if we could tune hyperparameters to improve our overall accuracy. We would try out many combinations of possible hyperparameters by our data to give us the absolute best weights for our model.

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn import metrics
3
4 # create a simple, non-parameterized Logistic Regression model
5 model = LogisticRegression(random_state=42)
6 model.fit(x_train, y_train)
7
8 y_pred = model.predict(x_test)
9 print(metrics.accuracy_score(y_test, y_pred))
10
11 from pprint import pprint
12 # Look at parameters used by our current forest
13 print('Parameters currently in use:\n')
14 pprint(model.get_params())
15
16 # create complex Logistic Regression with max_iter=131
```

LATEST
EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

```

17 log_model = LogisticRegression(max_iter=131, verbose=2, random_state=42)
18 log_model.fit(x_train, y_train)
19 y_pred_log = log_model.predict(x_test)
20 print(metrics.accuracy_score(y_test, y_pred_log))

```

team_logistic_regression.py hosted with ❤ by GitHub

[view raw](#)

We accomplished this using cross-validation: because we only have a vague idea of the parameters we might want to use, our best approach is to narrow our search is and evaluate a wide range of values for each hyperparameter.

LATEST

Using RandomizedSearchCV, we searched among **= 59,400** possible settings – and so the most efficient way to do this would be to take a random sample of the values.

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

Running our model with the best parameter values, however, samples actually decreased the accuracy of our model which showed us that while random sampling helps us narrow down our hyper parameter tuning within a distribution, we still have to explicitly check all combinations with GridSearchCV.

WRITE FOR TDS

Sign in

[Submit an Article](#)

```

1 from sklearn.model_selection import GridSearchCV
2
3 # Create the parameter grid based on the results of grid search
4 # Penalty type
5 penalty = ['l1', 'l2', 'elasticnet', 'none']
6 # Solver type
7 solver = ['lbfgs', 'liblinear']
8 # Maximum number of iterations
9 max_iter = [int(x) for x in np.linspace(start = 80, stop =
10 # Multi class
11 multi_class = ['auto', 'ovr']
12 # Verbosity
13 verbose = [0, 1, 2]
14 # l1 ratio
15 l1_ratio = [0, 0.8, 0.9, 1]
16 # C
17 C = [0.5, 0.75, 1.0, 1.25, 1.5]

```

```

18
19 # Create the param grid
20 param_grid = {'penalty': penalty, 'solver': solver, 'max_iter':max_iter,
21     'multi_class':multi_class, 'verbose':verbose, 'l1_ratio':l1_ratio,
22     'C':C
23 }
24 pprint(param_grid)
25
26 # Instantiate the grid search model with 2-fold cross-validation
27 log_grid_search = GridSearchCV(estimator = LogisticRegression(random_state=42), para
28
29 # Fit the grid search to the data
30 log_grid_search.fit(x_train, y_train)                                LATEST
31 best_log_grid = log_grid_search.best_estimator_                      EDITOR'S PICKS
32 best_log_grid.fit(x_train, y_train)
33 y_pred_best_log = best_log_grid.predict(x_test)                     DEEP DIVES
34 print(metrics.accuracy_score(y_test, y_pred_best_log))

```

grid_search.py hosted with ❤ by GitHub

In this case, implementing GridSearch only marginally improves our accuracy with our LR model.

The second model we looked to implement was [RandomForestClassifier](#), which can be efficiently used for both regression and classification. In this case, we will see how the RandomForestClassifier can build a proper decision tree to determine the given team-stats.

[Submit an Article](#)

```

1 from sklearn.ensemble import RandomForestClassifier
2
3 rf = RandomForestClassifier(random_state=42)
4 rf.fit(x_train, y_train)
5 y_pred_rf = rf.predict(x_test)
6 print(metrics.accuracy_score(y_test, y_pred_rf))
7
8 from pprint import pprint
9 # Look at parameters used by our current forest
10 print('Parameters currently in use:\n')
11 pprint(rf.get_params())

```

team_random_forest.py hosted with ❤ by GitHub

Immediately, we get that the RandomForestClassifier reaches an initial accuracy of 66.95%, which again is pretty good. Like with the LR model, we attempted to tune the hyperparameters to give us more accurate results – first using RandomizedSearchCV.

```

1  from sklearn.model_selection import RandomizedSearchCV
2  import numpy as np
3
4  # Number of trees in random forest
5  n_estimators = [int(x) for x in np.linspace(start = 0, stop = LATEST
6  # Number of features to consider at every split
7  max_features = ['auto', 'sqrt']                                     EDITOR'S PICKS
8  # Maximum number of levels in tree
9  max_depth = [int(x) for x in np.linspace(0, 100, num = 6)]        DEEP DIVES
10 max_depth.append(None)
11 # Minimum number of samples required to split a node
12 min_samples_split = [2, 5, 10]
13 # Minimum number of samples required at each leaf node
14 min_samples_leaf = [1, 2, 4]                                         NEWSLETTER
15 # Method of selecting samples for training each tree
16 bootstrap = [True, False]
17 # Create the random grid
18 random_grid = {'n_estimators': n_estimators,
19                 'max_features': max_features,
20                 'max_depth': max_depth,
21                 'min_samples_split': min_samples_split,
22                 'min_samples_leaf': min_samples_leaf,
23                 'bootstrap': bootstrap}
24 pprint(random_grid)
25
26 # Use the random grid to search for best hyperparameters
27 # First create the base model to tune
28 rfc = RandomForestClassifier(random_state=42)
29 # Random search of parameters, using 2-fold cross validation
30 # search across 100 different combinations, and use all ava
31 rfc_random = RandomizedSearchCV(estimator = rfc, param_distr
32 # Fit the random search model
33 rfc_random.fit(x_train, y_train)
34 y_pred_rfr_random = rfc_random.predict(x_test)
35 print(metrics.accuracy_score(y_test, y_pred_rfr_random))

```

randomized_search.py hosted with ❤ by GitHub

[Submit an Article](#)

[WRITE FOR TDS](#)

[Sign in](#)

Unlike with the LR model, we find that RandomizedSearch improves our hyperparameter tuning, giving us a better accuracy of 67.15%.

Running GridSearchCV in a similar manner to what we did above, we also sought to explicitly test $2 \times 6 \times 2 \times 3 \times 3 \times 5 = 1080$ combinations of settings instead of randomly sampling a distribution of settings. GridSearch also gave us an improvement from the base RandomForestClassifier, with an accuracy of 67.1

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

Predicting the Outcome of Games Based on Individual Player Statistics and Score

We then took a different approach to predicting game to see if we can achieve any better performance. We will use a larger dataset of individual player statistics that we will train a model to predict how many points a player will score in a given game. We will predict this based on their average season stats up until the game we are trying to predict, as well as their average performance over the previous 10 games. We have already created this data in the feature engineering section above. We will also make use of Elo ratings in order to account for home advantage, as presumably the higher rating of the opposing team, the fewer points a player will score. Once we have this data, we can use it to train a machine learning model to predict the outcome of each game.

predict how many points a team will score in a game by summing the predicted number of points of each individual player will score. With this information we will be able to predict which team will score more points and thus win the game.

Before we run our models, we need to clean our data slightly. For some games in this dataset, we have the statistics for one teams' players, but not for the other team – generally ~~only for the first~~ game that other team plays in the season. Thus, ~~LATEST~~ all these games from the dataset.

```
1 game_id_value_counts = final_player_stats['GameID'].value_counts()
2 valid_game_ids = [x for x in final_player_stats['GameID'] if
3
4 final_player_stats = final_player_stats[final_player_stats['
```

cleaning_player_stats.py hosted with ❤ by GitHub

[EDITOR'S PICKS](#)

[DEEP DIVES](#)

[NEWSLETTER](#)

[WRITE FOR TDS](#)

[Sign in](#)

[Submit an Article](#)

Unlike with the above games, we can't randomly split up our data into train and test sets. We are looking to use individual player statistics to predict the final score of a team, therefore all players playing in the same game together. To do this, we will split up our train and test sets by game so players from the same game stay together. About 80% of the games will be in the train set and 20% will be in the test set:

```
1 import random
2 games = set(final_player_stats['GameID'].unique())
3
4 train_set_game_ids = random.sample(games, 9700)
5
6 x_train = features[features['GameID'].isin(train_set_game_ids)]
7 y_train = label[label['GameID'].isin(train_set_game_ids)].drop(['GameID'], axis=1)
8
9 x_test = features[features['GameID'].isin(train_set_game_ids)]
10 y_test = label[label['GameID'].isin(train_set_game_ids)] ==
```

[train_test_split.py](#) hosted with ❤ by GitHub[view raw](#)

Instead of using a Logistic Regression model, for player scoring we will use a Linear Regression model as we are looking to predict a range of possible values (points scored) instead of simply predicting a win or a loss. Our RMSE (Root Mean Squared Error) for all players was 5.56, or the equivalent of each player making or missing around 2–3 baskets game arc averages.

On the test set, we grouped each team's predict each game and compared it with their actual score. Computing the numbers of games won versus the on predicted scoring gave us a ratio of 1483/252 of **58.66%**. Clearly, and as we realized earlier with distributions of teams versus their opponents, a performance is too variable of a determinant to predict the outcome of games – especially when team performance which tends to be more consistent across games.

[...](#)[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[NEWSLETTER](#)[WRITE FOR TDS](#)[Sign in](#)[Submit an Article](#)

Conclusions and Future Considerations

As avid NBA fans, we felt that creating a model to predict the outcome of NBA games would be an interesting project. It taught us a lot about building classifiers for predicting discrete outcomes. We were able to utilize many of the skills learned in our Big Data Analytics class for this project, including data scraping, data cleaning, feature analysis, building and tuning models, and hyperparameter tuning – and want to thank Professor Hui Li for his excellent work in teaching throughout the semester.

Our Random Forest Regression model, with parameters optimized through RandomSearchCV, gave us the highest testing accuracy of 67.15%. It is slightly higher than the Logistic Regression model, and it is much higher than the Linear Regression model based on individual player statistics. Optimizing parameters using GridSearchCV and RandomizedSearchCV was time consuming and computationally costly, and it resulted in only marginal changes in testing accuracy. If we had more time, we could spend less time optimizing parameters and more time on the model.

The best NBA game prediction models only accurately predict the winner about 70% of the time, so our logistic regression and random forest classifier are both very close to the upper bound of predictions that currently exist. If we had more time, we could explore other models and see just how much higher the accuracy we could get. Some of those candidate models include Naïve Bayes Classifier, linear discriminant analysis, convolutional neural networks, and support vector machines.

Hopefully, you enjoyed reading about our work and enjoyed making it – and learned something from it!

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

WRITTEN BY

Josh Weiner

See all from Josh Weiner

Data Science

Machine Learning

Python

Share This Article

Towards Data Science is a community publication. Submit your insights to reach our global audience and earn through the TDS Author Payment Program!

LATEST

Write for TDS

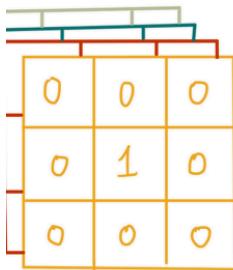
EDITOR'S PICKS

DEEP DIVES

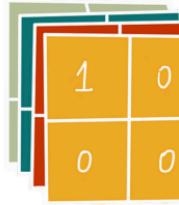
NEWSLETTER

WRITE FOR TDS

Sign in

Submit an Article

max
pooling



ARTIFICIAL INTELLIGENCE

ARTIFICIAL INTELLIGENCE

Implementing Convolutional Neural Networks in TensorFlow

Step-by-step code guide to building a Convolutional Neural Network

Shreya Rao

August 20, 2024 6 min read

How to Forecast Time Series

A beginner's guide to time series reconciliation

Dr. Robert Kübler

August 20, 2024 1



DATA SCIENCE

Hands-on Time Series Anomaly Detection using Autoencoders, with Python

Here's how to use Autoencoders to detect signals with anomalies in a few lines of...

Piero Paialunga

August 21, 2024 12 min read

DATA SCIENCE

Solving a Constrained Project Scheduling Problem with Quantum Annealing

Solving the resource constrained project scheduling problem (RCPSP) with D-Wave's hybrid constrained quadratic model (CQM)

Luis Fernando PÉREZ ARMAS, Ph.D.

August 20, 2024 29 min read



MACHINE LEARNING

3 AI Use Cases Chatbot)

Feature engineering unstructured data

Shaw Talebi

August 21, 2024 7

ARTICLES

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS



DATA SCIENCE

Back To Basic Regression ar

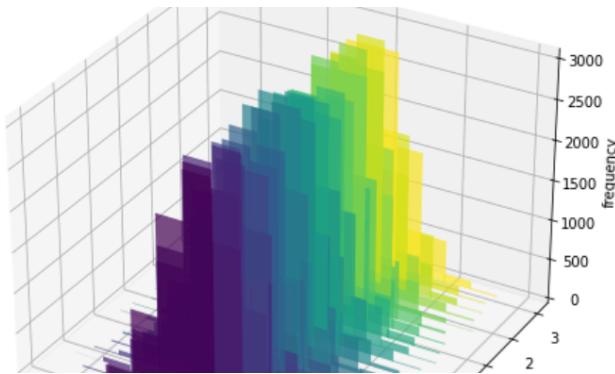
An illustrated guide to learning concepts:

Shreya Rao

February 3, 2023

Sign in

Submit an Article



DATA SCIENCE

Must-Know in Statistics: The Bivariate Normal Projection Explained

Derivation and practical examples of this powerful concept

Luigi Battistoni

August 14, 2024 7 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

NEWSLETTER

WRITE FOR TDS

Sign in



Submit an Article

Your home for data science and AI. The world's leading publication in analytics, data engineering, machine learning, and artificial intelligence.

© Insight Media Group, LLC 2025

Subscribe to Our Newsletter

WRITE FOR TDS · ABOUT · ADVERTISE · PRIVACY POLICY

DO NOT SELL OR SHARE MY PERSONAL INFORMATION