DS 4400: Machine Learning and Data Mining I Spring 2022

Team Members: Cameron Gleichauf (@CamGleichauf), Justin Chen (@jstinchen) Video

# **NCAA Basketball Prospect Grader**



### **Problem Description**

Every year NBA teams spend countless resources and time trying to determine the best college and international prospects that have entered the NBA draft. Draft selections can shape the course of a franchise for years to come and make for some of the most important decisions a franchise can make. Take the 2018 draft for example, where 4 of the top 5 picks have ranged from above average starters to bonafide franchise players. Then look at the #2 overall pick in the draft, Marvin Bagley, who struggled to find his footing in Sacramento and was eventually traded to Detroit for rotation pieces Trey Lyles and Josh Jackson. The importance of drafting well is showcased here, as each of the four other teams drafting in the top 5 made the playoffs this year, while Sacramento missed the playoffs for the 16th consecutive year.

Even after putting large amounts of time, energy and effort into scouting, many franchises still fail to choose the 'best' players possible or the 'best' players for their team needs. These decisions affect the success of the franchise for years to come and should be made with the most possible context information possible. Our machine learning algorithm aims to resolve this issue by collecting a large set of relevant feature variables that contribute to the success of draft prospects in the NBA. Our goal is to predict the quality of a potential NCAA D1 prospect's NBA career on a continuous scale based on evidence from recent prior drafts and patterns in the data.

The problem we are aiming to resolve is a regression problem that predicts a continuous output variable representing the estimated NBA win shares for each NCAA prospect. We will use regularized linear regression, decision tree regression, random forest regression, adaboost ensemble regression and a multi-layer perceptron to attempt to predict each training row's value on a continuous scale.

The resources, money and time that are spent on scouting is enormous and to solely use 'the eye test' in evaluating players will almost certainly lead to suboptimal draft selections. This is why the use of machine learning as an aid for decision making is incredibly important for NBA front offices. Machine learning is by no means a perfect system either, but its use as a supplemental factor when drafting players can provide huge benefits in selecting optimal players that may be valued less if only the eye test is used. Furthermore, the application of this machine learning algorithm can be used to revolutionize the way draft prospects are valued in the NBA by allowing us to make unbiased evaluations based on player talent that are repeatable over a large period of time. The use of machine learning algorithms to quantify the quality of a player can provide benefits in the range of millions and millions of dollars to teams if used correctly and in supplementation of classic scouting techniques.

# **Dataset and Exploratory Data Analysis**

### **Dataset statistics**

- Around 23,000 records of collegiate players and their NBA contributions (win shares)
- Our preliminary dataset before implementing recursive feature elimination and regularization includes about 46 different features
- Another set of 73 prospects and their statistics for the 2022 NBA Draft were collected. This list was based off all the collegiate prospects on <u>ESPN's Top 100</u> Draft Board
- Statistics and biographical information for the two datasets were scraped from <u>Sports Reference</u> using BeautifulSoup4 and wrangled with Pandas

### Feature definition and some insights

- Features are split into three main categories: counting statistics, advanced statistics and metadata regarding the player's physical attributes and position
- Positional labels are encoded
- We analyzed the covariance between various feature variables and attempted to eliminate unnecessary variables in an attempt to make the model easier to train and interpret.

# **Exploratory data analysis**

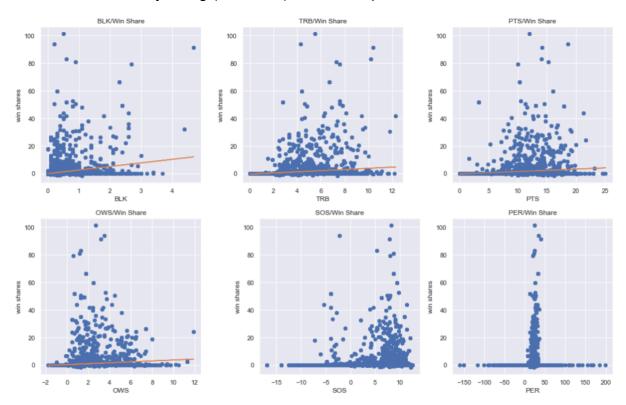
In order to apply recursive feature elimination and to determine the covariance between various feature variables, we plotted the Pearson correlation to the features dataframe. Here is the resulting graph depicting a subset of the features:



- Since we have 46 feature variables we are starting with, the graph is quite cluttered and hard to make sense of. To simplify the model, we applied recursive feature elimination to reduce the total number of feature variables and get the most relevant features. Not many of the variables were eliminated after applying

recursive feature selection. One possible reason for this is that, based on our feature correlation map, most of the feature variables were positively correlated to the target variable and thus feature elimination was not highly suitable for this particular model.

- Another technique we utilized was one-hot encoding. This allowed us to turn categorical features such as the player's position into numerical features that could be used as part of our model.
- We also ranked all of the features by individual correlation to win shares to get a rudimentary idea of each of them's importance
  - Blocks (0.23818), rebounds (0.191122), points (0.182362), strength of schedule (0.142153), offensive win shares (0.117306), and player efficiency rating (0.118170) were a couple features that stood out



- With these key, high-correlation statistics, we also looked at how different positions averaged in these categories. It can be seen that big-men triumph guards in blocks and player efficiency, while margins for rebounds, offensive win-share, and points aren't too major.

position	TRB	BLK	ows	PTS	PER
Guard	2.152826	0.133528	1.08446	6.118033	10.809212
Forward/Center	3.424525	0.498334	1.024874	5.393287	12.979949

- All the features were positively correlated with win shares except for 3-point attempt rate (3PAr) and turnover percentage (TOV%).
  - 3PAr's negative correlation can be attributed to the idea of taking efficient shots and not settling for threes. Being able to score around the basket and from distance also shows a player's well-roundedness, so if a player is only taking 3-pointers that is not ideal.
  - TOV%'s negative correlation can be explained easily you don't want a player giving the ball to the other team

# **Approach and Methodology**

### Machine learning models/Metrics:

We trained and tested 5 separate models after randomly splitting the training and testing set with sklearn.model\_selection import train\_test\_split function. These were the metrics we obtained for each model, in decreasing order of quality on the testing set:

- 1. Random Forest Regression (Best Model)
  - a. Training r2 score: 0.930
  - b. Testing r2 score: 0.601
  - c. Testing MSE: 7.659
- 2. Decision Tree Regression
  - a. Training r2 score: 0.558
  - b. Testing r2 score: 0.449
  - c. Testing MSE: 10.574
- 3. AdaBoost Ensemble Regression
  - a. Training r2 score: 0.467
  - b. Testing r2 score: 0.320
  - c. Testing MSE: 13.057
- 4. Multi-Layer Perceptron Regression
  - a. Training r2 score: 0.467
  - b. Testing r2 score: 0.307
  - c. Testing MSE: 13.309
- 5. Ridge Regression (Worst Model)
  - a. Training r2 score: 0.158
  - b. Testing r2 score: 0.145
  - c. Testing MSE: 16.41

### **Discussion and Result Interpretation**

We found random forest regression to be the most effective model on the testing set, with an r^2 value of just over 0.6. Initially, the r^2 value for random forest regression was closer to 0.5 but after applying cross-validation to find the optimal max depth hyperparameter (which was 15), we were able to increase the testing r2 value by almost 0.1. Random forest regression models are composed of multiple decision trees and average the results of these trees together. In this way, we increase the previously small sample size of the decision trees that tend to catch much of the noise from the training data and the model is able to generalize much better. While decision trees tend to overfit to the training data, random forest regression is able to combat this by taking the average of a larger number of decision trees.

The second most effective model was the decision tree regression model. Initially, the model overfitted the data by a significant margin, with the training data having an r^2 value of 1 and the testing data having an r^2 value of 0.13. However, after applying cross validation with the hyperparameter of maximum tree depth, we found the optimal value to be 6. Setting 6 as the max depth value, we achieved a significantly higher testing r^2 value of 0.449 but also a significantly lower training r^2 value of 0.558. This is acceptable because the goal is to have the highest testing r2 score, we do not care about optimizing the training r2 score.

The next best performing model was AdaBoost Ensemble Regression with a training r2 score of 0.467 and a testing r2 score of 0.320. We applied cross validation on this model as well to find the optimal number of estimators to be 10.

The second worst performing model was the Multi-Layer Perceptron Regression model, with a testing r2 score of 0.307 and training r2 score of 0.467.

The worst performing model by a wide margin was the Ridge Regression model. It had a training r2 score of 0.158 and a testing r2 score of 0.145. One possible reason why it performed so poorly is because it had an extremely high bias, assuming that the model was linear. This is likely an incorrect assumption given the complex nature of the problem and thus resulted in low r^2 scores.

# Model's Draft Board (Top 20 Prospects)

Model Rank	name	position	school	ESPN ranking	NBA Win Share Prediction
1	Walker Kessler	Forward	UNC/Auburn	19	31.38
2	Mark Williams	Center	Duke	15	29.12
3	Chet Holmgren	Center	Gonzaga	1	14.43
4	Trevion Williams	Forward	Purdue	36	12.31
5	Jabari Smith	Forward	Auburn	3	12.31
6	Jalen Duren	Center	Memphis	6	11.97
7	Tari Eason	Forward	Cincinnati/LSU	13	10.82
8	Bryce McGowens	Guard	Nebraska	22	10.37
9	Malaki Branham	Guard	OSU	15	8.87
10	Kennedy Chandler	Guard	Tennessee	16	8.13
11	Scotty Pippen Jr	Guard	Vanderbilt	69	7.83
12	Blake Wesley	Guard	Notre Dame	17	7.66
13	Paolo Banchero	Forward	Duke	2	7.65
14	TyTy Washington	Guard	Kentucky	12	6.68
15	Ochai Agbaji	Guard	Kansas	10	6.44
16	Keegan Murray	Forward	Iowa	5	5.99
17	AJ Griffin	Forward	Duke	7	5.93
18	EJ Liddell	Forward	OSU	18	5.93
19	Justin Lewis	Forward	Marquette	29	5.87
20	Jabari Walker	Forward	Colorado	45	5.31

It is clear that the model favored forwards and centers, looking at both the pre-implementation data exploration and the results. Blocks and rebounds were highly correlated with win-shares and bigger players collect those statistics a lot easier than guards.

For the most part, all the prospects who entered our top-20 list were all ranked in the same range as ESPN. It was interesting to see Walker Kessler and Mark Williams, both mid-first round picks, ranked so highly. Another immediate observation was how the top 7 prospects were all big men – that's a testament to their ability to grab rebounds, block shots, and make close-range shots at a high percentage (and therefore efficiency).

There were two surprising outliers in our top-20 predictions: Trevion Williams and Scotty Pippen Jr. Both benefited from above average strength of schedules (with Williams' coming in at a whopping 10). Pippen excelled in efficiency, had the second most steals, and the most points in our dataset – all features that the model weighed heavily. Williams was a good rebounder, ranked 11th in offensive win shares, and 8th in player efficiency rating.

That follows the trend of the NBA, valuing their forwards being versatile – scoring both inside and out – as "stretch bigs." A lot of veteran forwards have also improved their three-point shooting as the league and analytics shifted to put a bigger emphasis on the skill. Some players that fit this description include 2021-22 Season's MVP finalists – Giannis Antetokounmpo, Joel Embiid, and defending MVP Nikola Jokic – with rising stars and prospects like Evan Mobley, Jaren Jackson Jr, and DeAndre Ayton.

Since almost all of the features were positively correlated, the model rewards players who collect higher statistics playing harder competition (strength of schedule). This is reflected in all of the top prospects coming out of traditional basketball powerhouses (Duke, Kentucky, Kansas, Gonzaga, etc) and being high-usage, star players on their teams during their collegiate career.

#### References

- Using Machine Learning to Predict Careers of 2019 NBA Draft Picks An article detailing a similar project that attempts to classify possible NCAA prospects into one of several categories: MVP winners, All-NBA First Team players, All-NBA Second or Third Team players, and All-Stars, good starters, good role players, role players, fringe bench players, and players who flamed out of the league. This was one major difference between our model and the model outlined by this article, as our model used regression to predict total win shares in the NBA while their model attempted to classify players into one of a few groups. The article reviews the different types of models they attempted to use, such as decision trees or random forests. They ran into similar issues as us with decision trees, as the models tended to perform very well on the training sets but not as well on the testing sets, something known as overfitting. These issues were addressed by the random forest models because they are composed of multiple decision trees and average the results of these trees together. In this way, we reduce the small sample size of the decision trees that tend to catch much of the noise from the training data and the model is able to generalize much better.
- Advanced NBA Stats for Dummies: How to Understand the New Hoops Math An article that gives an overview and summary of many advanced NBA metrics that we will use as feature variables in our model training. The article discusses the tendency of advanced metrics to overvalue big men as they tend to have higher shooting percentages, blocks and rebounds as compared to guards. This was a relevant reference for us because we were attempting to find the 'best' target variable to assign to NBA players to quantitatively measure their success in the league.
- NBA Win Shares A deep dive into the statistic known as win shares. A system developed by mathematician Bill James where three win shares represent one win. Win shares are calculated via aggregations performed on various counting statistics as outlined in the article above. This was a relevant reference for us because we were attempting to decide on which advanced metric to use as the target variable to most effectively measure an NBA player's success in the league.
- Sloan Sports Analytics Conference Discussion Regarding Use of Advanced Statistics and Metrics in NBA for Drafting A video of a panel from the Sloan Sports Analytics Conference discussing the use of advanced metrics and analytics as guides for drafting players in the NBA draft. It discusses the current use of analytics at the time (2014) and its possible uses in the future as it pertains to NBA front offices making decisions.

### Conclusion

We set out to create a Machine Learning model that could predict the best players in a given NBA draft class coming from the NCAA. We tried various models throughout the process and found random forest regression to be the most effective with an overall r^2 value of just over 0.6. While this may not seem incredibly high in comparison to some other types of problems that are commonly addressed by machine learning algorithms, drafting in the NBA is an inherently uncertain enigma. Many things cannot be entirely quantified by statistics, such as a player's mental makeup or the circumstances around them while playing in college. Additionally, teams must be able to develop and play their players effectively in order to turn them into quality NBA players. Some teams have better track records of doing so than others, and this may in turn affect what the players ultimately turn into. Furthermore, basketball is a team sport and players can contribute meaningfully, like setting screens, being a good defender, and creating scoring opportunities for teammates, without it showing up in the statistics.

There are a few ideas we came up with to further improve the quality of the model from its current state. The idea that comes to mind first is to emphasize the importance of different statistics based on a player's position. For example, rebounding and blocks should be more important for a center than for a guard. Vice versa, three point shooting and assist rate should be valued more highly for point guards than centers. Another idea that comes to mind to improve the algorithm is to include a player's age as part of the algorithm. Older players tend to have a lower ceiling than younger players in college when it comes to NBA potential and this is something we failed to include in our model. Another feature variable that could be of use is to include a player's ranking coming out of high school, as this is yet another indicator of their overall talent level. Some features that would be nice to have but are difficult to quantify are mental makeup and defensive prowess, although metrics for defensive prowess are improving as of late. Another thing to consider is potentially using a different target variable instead of career Win Shares. One downside to career Win Shares is that they tend to favor big men over guards. Additionally, for recent players such as Anthony Edwards or Ja Morant who have been excellent but for a short period of time, their career win shares will be low and in turn not necessarily indicative of the quality of player they truly are. One remedy for this would be to use WS/40, although this has its downsides as well.

Overall, the model performs fairly well on testing data and projects two of ESPN's top three players in the 2022 NBA draft (Chet Holmgren and Jabari Smith) to be in the top five players of the draft, with the third player (Paolo Banchero) finishing 13th. The model does tend to favor big men, with projected mid-first round picks Walker Kessler and Mark Williams projected as the top two players in the draft. Two prospects that our model indicates as underrated are Trevion Williams of Purdue and Scotty Pippen Jr. of Vanderbilt.