

Final Report - Team 103

Introduction/Background

Studies using previous performance data have shown that college experience may not improve performance in the NBA and that players with limited or no college experience perform better than those who attend college. Another argument is that a player who enters the NBA directly from high school lacks maturity, measured by one's ability to choose goals which highlights the need for important psychological needs, such as the need for autonomy, relatedness, and competence [1]. Fewer than 2 percent of NCAA student-athletes go on to become professional athletes even though on average NCAA student-athletes graduate at a higher rate than the general student body. Hence an analysis based on game statistics parameters like Win shares, VORP, defensive and offensive would provide better insights into a players development and how much college experience relates to the player's performance in the NBA.

Problem Statement

Perform a comprehensive analysis to predict the success of players drafted from NCAA, based on their relationship between college experience and their performance in NBA, using offensive and defensive statistics like point scores, assists, steals, blocks, field goals, Winshare and VORP.

Data Source

We sourced the NBA(1) and College Basketball(2) data for our project from 2 different websites.

1. The NBA data was downloaded from [Basketball Reference](#), and contains advanced statistics on each NBA player from 2011-2021 ([example link for 2021](#)). From this dataset we extracted the variables that can summarize overall success in the NBA. The variables we selected are: 'G' (games played in a season)', 'MP' (minutes played in a season), 'PER' (Player Efficiency Rating), 'OWS' (Offensive Win Shares), 'DWS' (Defensive Win Shares), 'WS' (Win Shares), 'WS.48' (Win Shares per 48 minutes), 'OBPM' (Offensive Box Plus Minus), 'DBPM' (Defensive Box Plus Minus), 'BPM' (Box Plus Minus), 'VORP' (Value Over Replacement Player), 'Player' (Player Name), 'Year' (year played) and 'Pos' (play position).

Example of nba advanced data from basketball-reference:

Data Modeling

Creating the Logistic Regression Model

We first create a logistic regression model using all variables except Position and School to predict a response variable or 0 or 1 for whether or not the player will be drafted by the NBA.

```
modell <- glm(Drafted ~ G + GS + MP + FG + FGA + FG_percent +  
  Twos + TwosA + Twos_percent + Threes + ThreesA +  
  Threes_percent + FT + FTA + FT_percent + ORB + DRB +  
  TRB + AST + STL + BLK + TOV + PF + PTS,  
  family=binomial(link='logit'), data=trainSet)
```

From the model's summary output (refer output in HTML file under folder "Code") we can identify the statistically significant variables. The statistically significant predictor variables are: G (Games), GS (Games Started), MP (Minutes Played), FG (Field Goals), FG_percent (Field Goal Percentage), FT_percent (Free Throw Percentage), DRB (Defensive Rebounds), AST (Assists), BLK (Blocks), and TOV (Turnovers).

Discussion on the Coefficients of our model

These statistically significant predictor variables do outline the attributes that a well-rounded NCAA basketball player should excel in to showcase their potential abilities as an NBA player. A player can be drafted to the NBA if they play and start in many games, and play for the majority of the game (with high Minutes Played). These statistics showcase the consistency, reliability, and health a player has throughout their collegiate career. Field Goals, Field Goal Percentage, Free Throws, and Free Throw percentage are good indicators for a player's ability to score points which, arguably, is one of the most important skills a player needs. Defensive rebounds, Assists, Blocks, and Turnovers are good indicators of the non-point based metrics that illustrate a strong defensive and team player.

One thing to note, which we will discuss later, is the line, "Coefficients: (4 not defined because of singularities)". The "Threes", "ThreesA", "TRB", and "PTS" have coefficients of NA. To understand this we will look into the correlations our predictor variables have with each other.

"Threes" (Three-point shots made) is almost perfectly correlated with "ThreesA" (0.99011326) which makes sense, the number of three points made has a strong linear relationship with how many three point shots are attempted.

"TRB" (Total Rebounds) is the sum of "ORB" (Offensive Rebounds) and "DRB" (Defensive Rebounds) which is why "TRB" is almost perfectly correlated with those predictor variables.

Train and Test the Model(s)

We will now use model1 to predict the Drafted response variable for the 30% test set we created earlier.

```
probabs1 <- predict(model1, testSet, type='response')
```

We then convert these predicted values to 0 or 1 using a threshold of 0.5. If the predicted value is less than 0.5 we say the player won't be drafted to the NBA, and if it is greater than 0.5 we say the player will be drafted to the NBA.

```
preds1 <- ifelse(probabs1 > 0.5, 1, 0)
```

We will construct a Confusion Matrix using the actual Draft values from the test set and the predicted values to assess the model's accuracy.

```
confusionMatrix(factor(preds1), factor(testSet$Drafted))
```

| | | |
|------------|------|-----|
| Reference | | |
| Prediction | 0 | 1 |
| 0 | 2903 | 111 |
| 1 | 7 | 4 |

| | |
|---------------------------------|--------------------------------|
| Accuracy : 0.961 | Sensitivity : 0.99759 |
| 95% CI : (0.9535, 0.9676) | Specificity : 0.03478 |
| No Information Rate : 0.962 | Pos Pred Value : 0.96317 |
| P-Value [Acc > NIR] : 0.6352 | Neg Pred Value : 0.36364 |
| | Prevalence : 0.96198 |
| | Detection Rate : 0.95967 |
| | Detection Prevalence : 0.99636 |
| | Balanced Accuracy : 0.51619 |
| Kappa : 0.0572 | |
| McNemar's Test P-Value : <2e-16 | 'Positive' Class : 0 |

The confusion matrix shows us that the model predicted the draft status for 96.1% of players correctly. However, there are 111 players that the model incorrectly did not draft, when they should have been drafted.

Our first model has an accuracy of 96.1% However, it did not draft 111 players that should have been drafted. Let's try lowering the threshold of our predictions to see if that make the model correctly draft more players.

We will lower the threshold to 0.4 and construct the Confusion Matrix.

```
preds2 <- ifelse(probabs1 > 0.4, 1, 0)
```

| Reference | | |
|------------|------|-----|
| Prediction | 0 | 1 |
| 0 | 2891 | 111 |
| 1 | 19 | 4 |

Accuracy : 0.957
 95% CI : (0.9492, 0.964)
 No Information Rate : 0.962
 P-Value [Acc > NIR] : 0.9276

Lowering the threshold did not improve accuracy, it lowered it 95.7%.

Let's create a second model using only the statistically significant predictor variables and see if it performs better than model1 (refer output in HTML file under folder "Code").

```
model2 <- glm(Drafted ~ G + GS + MP + FG + FG_percent + Twos + FT_percent +
  ORB + DRB + AST + STL + BLK + TOV,
  family=binomial(link='logit'), data=trainSet)
```

With model2 we will predict draft status using the testSet and using a threshold of 0.5 and construct the Confusion Matrix.

```
probabs2 <- predict(model2, testSet, type='response')
preds3 <- ifelse(probabs2 > 0.5, 1, 0)
```

| Reference | | |
|------------|------|-----|
| Prediction | 0 | 1 |
| 0 | 3396 | 119 |
| 1 | 11 | 3 |

Accuracy : 0.9632
 95% CI : (0.9564, 0.9691)
 No Information Rate : 0.9654
 P-Value [Acc > NIR] : 0.7849

 Kappa : 0.0373

 McNemar's Test P-Value : <2e-16

Sensitivity : 0.99677
 Specificity : 0.02459
 Pos Pred Value : 0.96615
 Neg Pred Value : 0.21429
 Prevalence : 0.96543
 Detection Rate : 0.96231
 Detection Prevalence : 0.99603
 Balanced Accuracy : 0.51068

 'Positive' Class : 0

According to the Confusion Matrix output our model's accuracy did improve to 96.32%. However, it still does not correctly draft 119 players.

Predicting the 2022 NBA Draft

Let's use model1 to predict which players from the 2022 NCAA season will be drafted to the NBA.

We used the same python script from earlier to get the relevant player information for the 2022 season only (ignoring draft status since that is what we want our model to predict)

After running our predictions with model1, we constructed a dataframe using player name's and the predicted value, 0 or 1. We sorted to list the 1's first.

| Name Drafted | |
|----------------------|---|
| Chet Holmgren | 1 |
| Trayce Jackson-Davis | 1 |
| Keegan Murray | 1 |
| JT Shumate | 1 |
| Norchad Omier | 1 |
| Jayveous McKinnis | 1 |
| Johni Broome | 1 |
| Max Abmas | 1 |
| Jalen Pickett | 1 |

Model1 predicted 9 players to be drafted by the NBA; Chet Holmgren, Trayce Jackson Davis, Keegan Murray, JT Shumate, Norchad Omier, Jayveous McKinnis, Johni Broome, Max Abmas, and Jalen Pickett.

Notable Draftees Prior to 2011

Let's see how model1 predicts 10 notable NCAA draft picks prior to 2011: (1) Stephen Curry, (2) Blake Griffin, (3) Derrick Rose, (4) James Harden, (5) Draymond Green, (6) Kevin Durant, (7) Chris Paul, (8) Brook Lopez, (9) Carmelo Anthony, and (10) Russell Westbrook.

Model1 did not predict Derrick Rose, James Harden, Draymond Green, Chris Paul, Brook Lopez, or Russell Westbrook from being drafted to the NBA.

How to make a more effective model

Our model only uses box score stats; there could be other factors that play a role in determining whether a player gets drafted or not, for example: the position that player plays, or maybe even their height.

The NBA teams draft order is decided by using the reverse order of their regular season record. The teams that did poorly will get to pick first, giving these teams a chance to get a high draft pick to improve their team's performance for the upcoming season.

However, what should also be considered is the team's current roster. Did a veteran player just retire? Are there players planning on retiring? Are there players with ending contracts? Are there any players with significant injuries that could affect them for future seasons? These scenarios could also influence what draft picks to choose that go beyond the scope of just box score stats.

Predicting NBA Success Through Advanced Stats

Another way we considered predicting NBA success from a college player's performance is by looking at if we can predict advanced NBA statistics that summarize a player's performance with college per game statistics. The NBA variables we chose to use are Win Shares (WS) and Value Over Replacement Player (VORP). Win Shares is a player statistic which attempts to divvy up credit for team success to the individuals on the team. This metric is designed to estimate the player's contribution in terms of wins. Value Over Replacement Player (VORP) converts the Box Plus/Minus (BPM) rate into an estimate of each player's overall contribution to the team, measured vs. what a theoretical "replacement player" would provide. Where a "replacement player" is defined as a player on minimum salary or not a normal member of a team's rotation. The two models we trained to predict success are a Decision Tree Regressor and a Linear Regression.

Decision Tree Regressor Models

To create our decision tree regressor we used the "rpart" library in R. We created a training and test dataset to build and test the model with. The training dataset contains 75% of the combined college & nba dataset, with the test dataset containing the other 25%. Initially we trained models that predicted WS and VORP using all of the college basketball statistics available (GS, MP, FGs, FGA, FG%, 2pt FGs, 2pt FGA, 2pt FG%, 3pt FGs, 3pt FGA, 3pt FG%, FT, FTA, FT%, ORB, DRB, TRB, AST, STL, BLK, Tov, PF, PTS). We then fine tuned each model by using Mallow's CP variable selection method. This improved each models' accuracy and gave us insight into what variables are important to NBA success.

Variable Selection:

For predicting Win Shares we found the following college basketball per-game statistics to be the most significant: FG (field goals), 2pt FG, 3pt FG, FTA (Free Throw Attempts), FT% (Free Throw Percentage), TRB (Total Rebounds), AST (Assists), STL (Steals), BLK (Blocks), and PTS (Points Scored). These stats give a good overview of impact in a basketball game both offensively and defensively. Interestingly, the variable selection method placed importance on counting stats for a player's shooting ability, favoring the total number of field goals a player makes vs statistics like FG% or 3pt FG%. It also placed importance on a player's free throw performance, where it did consider the percentage of shots being made. A player shoots a free throw as a penalty shot after being fouled, so a player gets to take standing still 15 feet from the basket with no defense being played. Perhaps the Mallow's CP algorithm found importance with efficiency in free throws because it better projects overall shot efficiency in the NBA than field

goals taken during the game. On the other hand, counting the number of made 2 point and 3 point field goals can show a college player's importance to their team's offense.

For predicting Value Over Replacement Player we found the following college basketball per-game statistics to be the most significant: MP (minutes played), FG, 2pt FG, 3pt FG, 3pt FGA (3pth Field Goal Attempts), 3pt FG%, FT (Free Throws made), FTA, DRB (Defensive Rebounds), TRB, AST, STL, BLK, TOV (Turnovers). Compared to the Win Shares model, the VORP variable selection placed more importance on 3 point field goals and also added in minutes played and turnovers. It makes sense to place importance on 3 point field goals. NBA offenses today are significantly more motion based, in an attempt to set up 3 point field attempts, compared to previous eras of the NBA. Turnovers can show how well a player can read offenses and defenses, which could align with the more motion based offenses as well.

Model Results

Win Share Model:

| Model | MAE | SSE |
|----------------------------|------|-----|
| All College Variables | 1.15 | 481 |
| Selected College Variables | 1.11 | 465 |

Value Over Replacement Player Model:

| Model | MAE | SSE |
|----------------------------|------|-----|
| All College Variables | 1.12 | 390 |
| Selected College Variables | 0.34 | 62 |

The Mallow's CP variable selection method successfully improved the accuracy of both of our models. There was great improvement for predicting VORP, going from an average absolute error of 1.12 points in Box Score Plus/Minus to 0.34. This means the decision tree model we trained was able to predict how many points a player adds or subtracts compared to a replacement-level player to within a fraction of a point. The WS model saw less of an improvement, going from an average absolute error of 1.15 games to 1.11. While the model didn't improve accuracy as much, it is still predicting Win Shares with only around 1 game of inaccuracy.

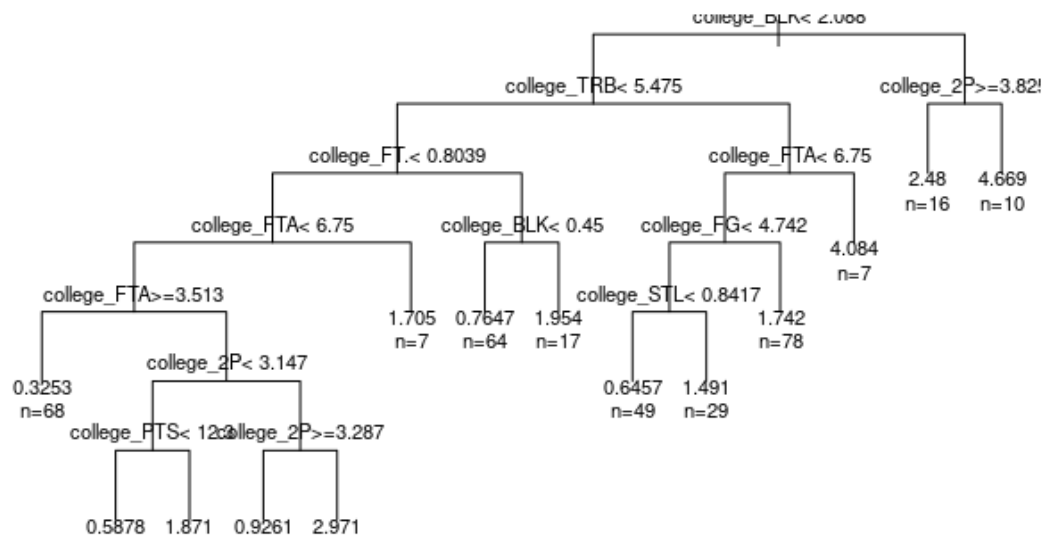
Summary based on Decision Tree Regressor Models

For both the Win Share and Value Over Replacement Player models favored college players with high counting stats, and only considered efficiency for shots like free throws. This makes sense because the best college players will have the ball in their hands often and accumulate a high volume of field goals made, assists, rebounds, etc. College players are still developing so even very talented players can still be inefficient, i.e. miss a significant amount of shots and turn the ball over. The VORP prediction model had better accuracy for a decision tree regressor, and

placed higher importance on 3 point performance which makes sense for the present day NBA, where 3s are more important than they are in college basketball or other eras of the NBA. The decision tree appears to split off defensive players by using variables like blocks and steals as split variables. Other important split variables are Free Throws Attempted and Minutes Played, showing how it's important to be given a large role in your team's rotation.

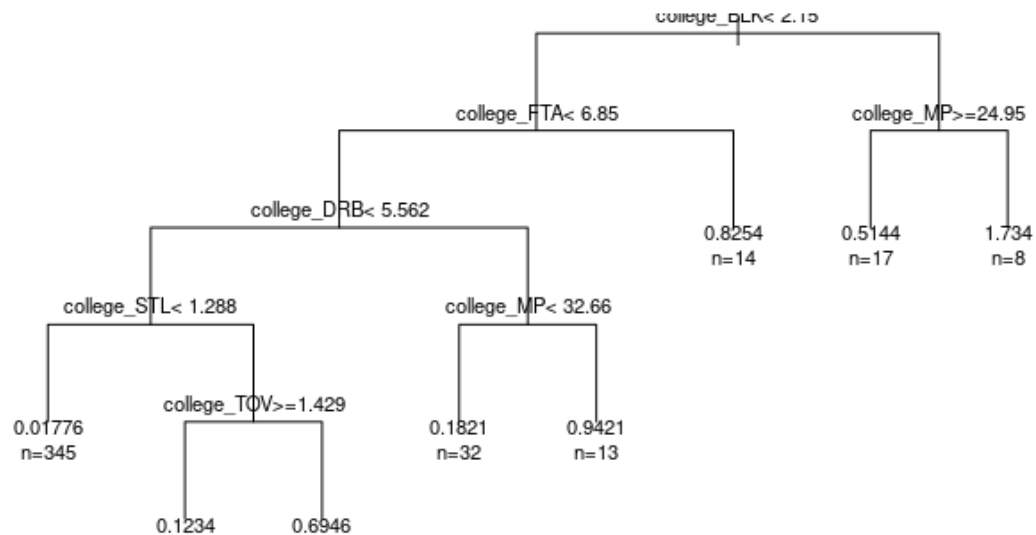
WS Decision Tree Regressor

NBA Win Share Decision Tree Regression



VORP Decision Tree Regressor

NBA VORP Decision Tree Regression



Linear Regression Models

The combined college and NBA dataset was used to create the training(75%) and test(25%) sets for linear regression models to predict the success of current NBA players using their college experience attributes, after excluding factor variables like "college_Player" and "college_Team". A bidirectional stepwise regression method (forward and backward selection) was run to identify college attributes which would be significant in predicting WinShare and VORP response variables. We chose the models with the lowest Akaike Information Criterion (AIC) value from the outputs generated by the stepwise regression method.

WinShare (WS) Regression Model

Fitting the WinShare regression model produced positive coefficients for college predictors - 2-point and 3-point attempts(college_2PA and college_3PA) total rebound (college_TRB), assists(college_AST), steals(college_STL), points (college_PTS) and blocks (college_BLK), while among the remaining predictors negative coefficients like turnovers (college_TOV) and field goals attempted (college_FGA) affect the overall WinShare measure of the player. However, assists, points, 2-point and 3-point attempts, and total rebounds are the statistically significant predictors which influence the WinShare measure of players in college and in their NBA careers.

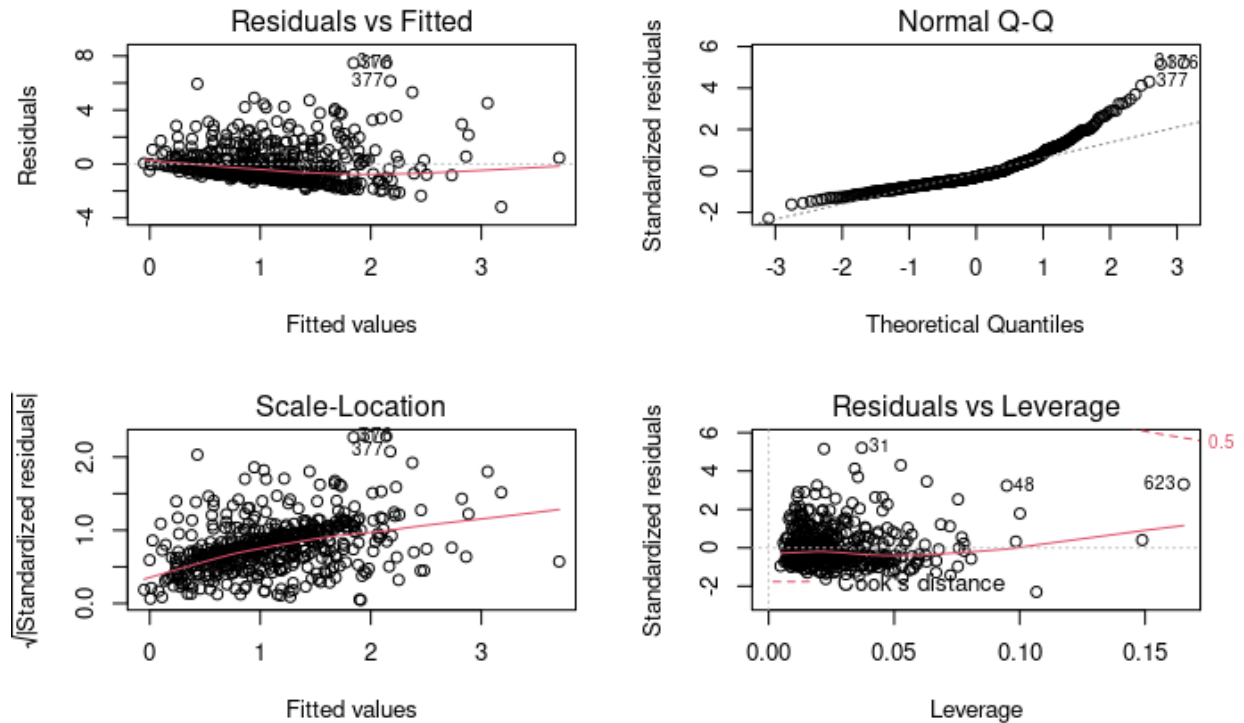


Fig. 1 : WinShare plot(s)

The residual plot shows a nonlinear relationship between the actual and predicted WinShare response variable and the Q-Q plot shows a right-skewed distribution which means that WinShare measure for most player's lies in a fixed range while there are several outlier players with a higher WinShare measure (Eg: Damian Lillard - 10.25). Outliers with greater than 2 standard deviations from the mean, as can be seen in the standardized residual vs. fitted values plot, points to the presence of an extremely unusual range of WinShare measure for some of the players.

Note: the residual std error indicates that on average a player's WinShare measure is 1.466 points away from the predicted values which clearly shows the model isn't an accurate one.

Call:

```
lm(formula = nba_WS ~ college_G + college_FGA + college_2PA + college_3P + college_3PA + college_DRB + college_TRB + college_AST + college_STL + college_BLK + college_TOV + college_PTS, data = as.data.frame(train_data))
```

Value-Over-Replacement (VORP) Regression Model

Fitting the VORP regression model produced statistically significant and positive coefficients like assists (college_AST), defensive rebounds (college_DRB), steals (college_STL), 2 & 3-PA (college_2PA & 3PA), which increases a player box-score estimate. But few significant negative coefficients like minutes played (college_MP), field goals attempted (college_FGA), 3-point % (college_3P) and turnovers (college_TOV) indicate an impact to their overall contribution to the team which could affect their NBA careers.

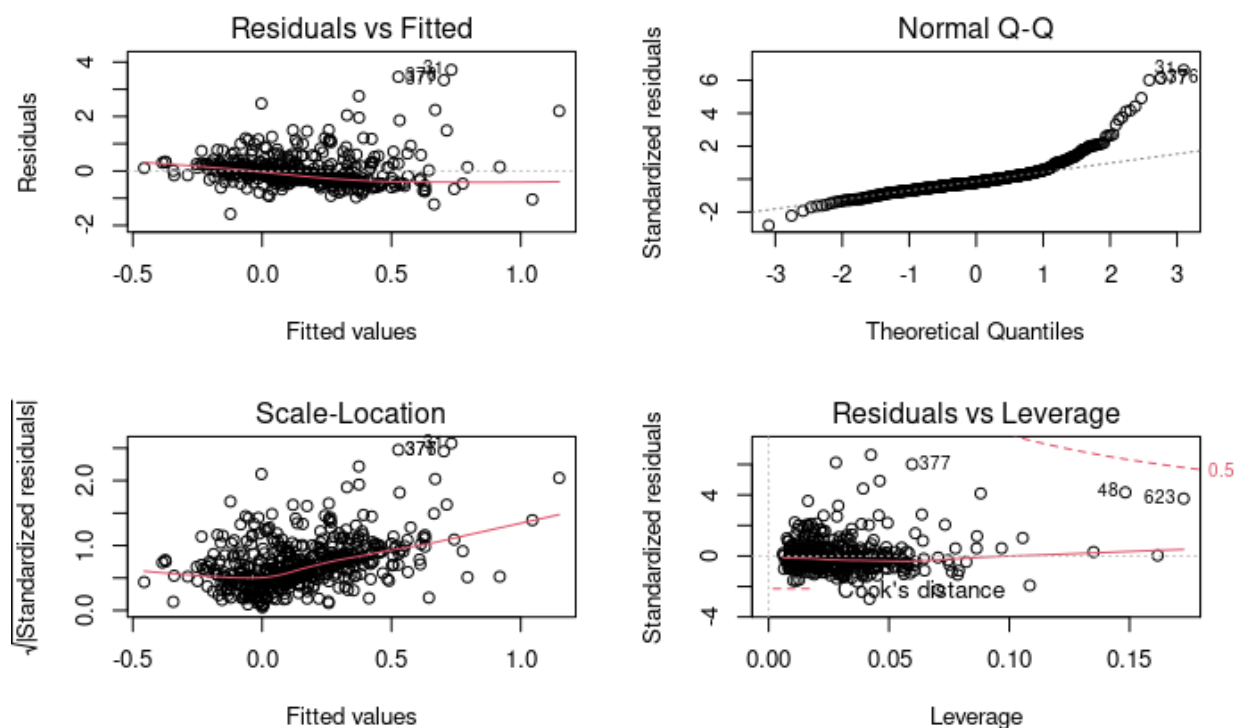


Fig. 2: VORP plot(s)

The residual plot shows a nonlinear relationship between actual and predicted VORP values and the Q-Q plot shows a right-skewed distribution which means that VORP score for most player's lies in a fixed range while there are several outliers who score more runs in a season (Eg: Damian Lillard - 4.76). Though a large number of datapoints appear concentrated near Cook's distance there are only a few influential observations that could significantly change regression the model, as can be seen in the residual vs. leverage plot. These influential points are players with very low VORP score from the mean (for e.g: Norris Cole [-0.5625], Jimmy Butler [3.125], Lonzo Ball [1.2])

Call:

```
lm(formula = nba_VORP ~ college_MP + college_FG + college_FGA + college_2PA + college_3PA + college_FT + college_FTA + college_DRB + college_AST + college_STL + college_BLK + college_TOV + college_3P., data = as.data.frame(train_data))
```

R^2 , Adjusted R^2 and p-value of the 2 linear models

| Models | Using all college predictors for each model | | | Using different selected college predictors for each model | | |
|----------|---|------------|----------|--|------------|-----------|
| | R^2 | Adj. R^2 | p-value | R^2 | Adj. R^2 | p-value |
| WinShare | 18.91% | 14.93% | 4.55e-12 | 13.52% | 11.45% | 7.108e-11 |
| VORP | 18.53% | 14.53% | 1.15e-11 | 13.42% | 11.17% | 2.413e-10 |

Both linear regression models have a low R-squared value and a low p-value (p-value ≤ 0.05). The low p-value for the regression models indicate that the model fits the data quite well even though their respective R-squared value is low.

Conclusion based on linear regression models

Output from the regression models show that the performance of players in college who are better in all offensive (assist, points, rebounds, field goal %, free throws) and defensive (steal, blocks, rebounds) areas may have an influence on their NBA careers but it also shows that players in college tend to commit more turnovers and fouls which affects their VORP and could have potential impact on their careers in NBA.

Interpretation

To know whether the results make sense, asking the tangential question of “how do basketball teams win games?” makes sense. Dean Oliver, the father of NBA basketball analytics, answered this question with what he called the "Four Factors of Basketball Success". These factors are efficient shooting, turnovers, rebounds, and free throw attempts.

When it comes to players being drafted, the logistical model found the most impactful independent variable to be high number of games played. The other variables in order of significance were blocks, assists, turnovers, minutes played, field goal percentage, field goals made, free throw percentage, defensive rebounds, and games started. As far as what is distinctly missing is anything regarding free throw attempts per field goal attempts.

Looking at the decision trees predicting future NBA VORP and WS, the variables that show up in both models are blocks, steals, and free throw attempts. Those last two were not in the drafted logistical model. Number of assists is missing from the decision trees despite being the third most significant factor in being drafted. Other factors that appear once in the decision tree models are two pointers made, free throw percentage, points, free throws made, defensive rebounds, total rebounds, and turnovers.

In the linear regression models, field goal attempts (2 point, 3 point, and total) appears in both VORP and WS models. Also, differentiating the linear regression models from the decision trees is that assists appears in both linear regression models while appearing in neither of the decision tree models. Again, steals and low turnovers makes an appearance in both the linear regression models like they did in the decision tree models. Other variables that appear in these models are blocks, points, 3 point percentage, defensive rebounds, and total rebounds.

So looking at these three different types of modeling there are some areas where the drafted players logistic model and the player performance, decision trees and linear regression, models differ. The main variables that differ are amount of games played, free throw attempts per field goal attempts, steals, and points. This makes sense as NBA teams have been known to look for

offensive stars over efficient scorers and defenders. The models were trained for accuracy, so looking for quality role players more than looking for highest potential.

As far as future work there are a few things we can add such as physical metrics, player positions, and on and off the court splits.

Reference

1. Zestcott, Colin & Dickens, Jessie & Bracamonte, Noah & Stone, Jeff & Harrison, C.. (2020). One and Done: Examining the Relationship Between Years of College Basketball Experience and Career Statistics in the National Basketball Association. *Journal of Sport and Social Issues*. 44. 019372352091981. 10.1177/0193723520919815.
2. Otten, Mark & Miller, Travis. (2015). A BALANCED TEAM WINS CHAMPIONSHIPS: 66 YEARS OF DATA FROM THE NATIONAL BASKETBALL ASSOCIATION AND THE NATIONAL FOOTBALL LEAGUE 1 ,2. *Perceptual and Motor Skills*. 121. 10.2466/30.26.PMS.121c25x4.