



Modeling the NBA Leap

Ryan Lewis



TABLE OF CONTENTS

01

**Business
Question**

02

Data Overview

03

**Exploratory
Data Analysis**

04

**Modeling
Techniques**

05

Results

06

Next Steps





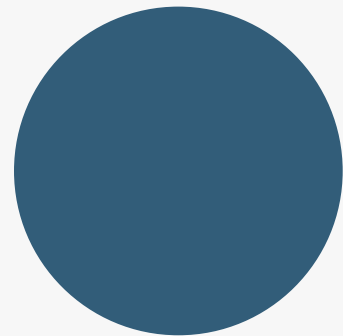
What is “The Leap”?

“A player’s identity typically begins to crystallize in his third or fourth NBA season. Young players have learned the ropes, and veterans have departed or aged, vacating heavy-duty roles that need filling. Everyone involved — players, agents, executives — looks to see what emerges as a player nears the expiration of his rookie contract.”
- Zach Lowe



Business Question

Based off an NBA players first three seasons, can you predict if they will make an All-NBA team in seasons 4 through 6?



Business Question Importance



Importance

Front office business planning: NBA players drafted in the first round of the draft can command contract extensions of up to 25% of a team's salary cap in 2020-21 season



Focus

In order to minimize false positives, our models will focus on the 'precision' metric -- this puts a higher emphasis on front offices correctly identifying actual true positives



Dataset

Stathead.com



Espn.com



Subset & Aggregate



Final Dataset

Stathead.com

Using Selenium, web scraped all seasonal and advanced player statistics dating back to 1947

ESPN.com

Merged in all player & team awards pulled from ESPN.com

Subset & Aggregate

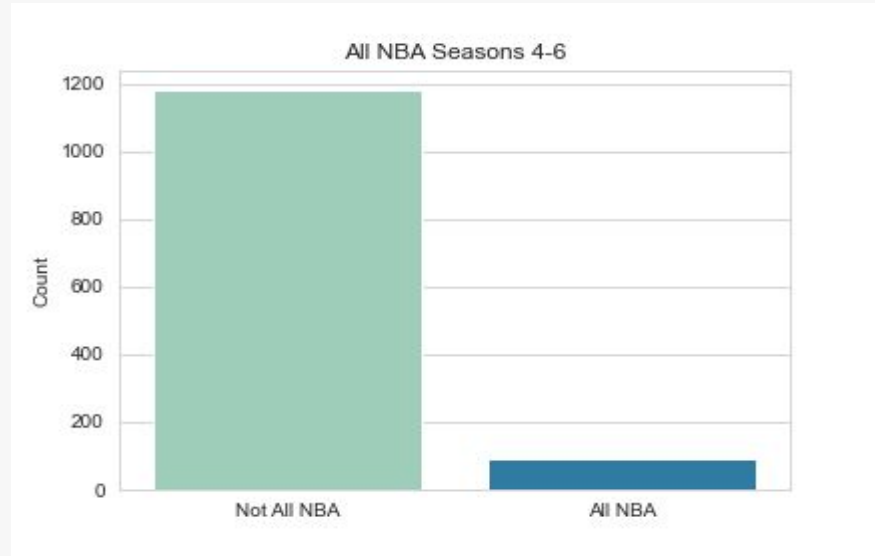
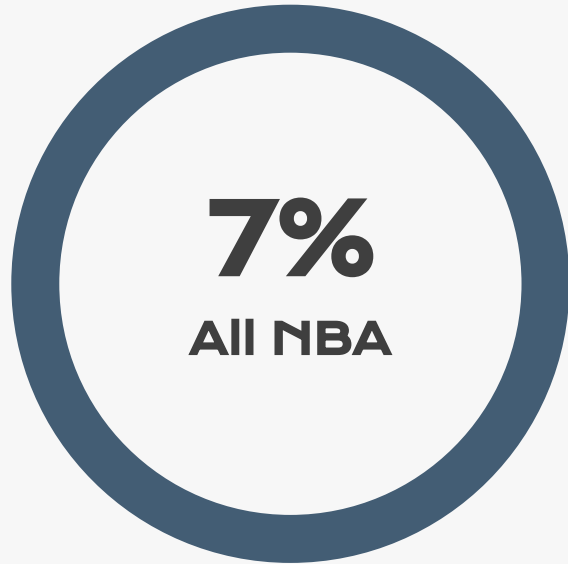
- Include only seasons after 1977
- Players who have played at least 6 seasons
- Aggregated so each row represented one player

Final Dataset

- 1273 qualified players
- 221 features; total & seasonal statistics



Exploratory Data Analysis



EDA - Basic Totals Pairplot



Statistic total from players first
3 seasons in the league

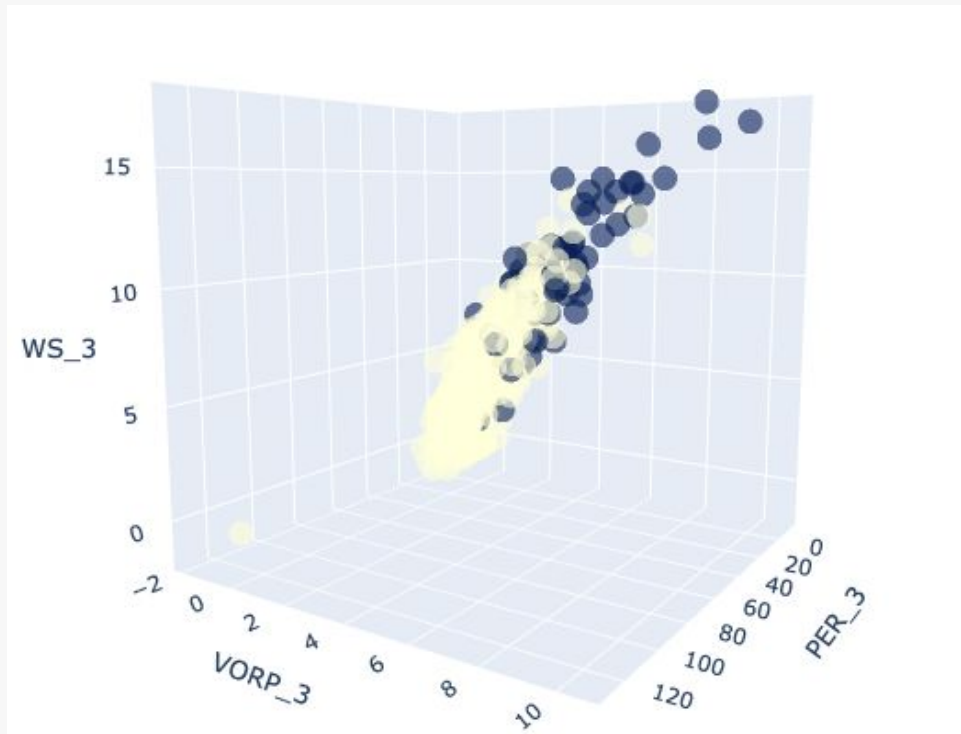


All-NBA players separate
themselves from their peers
statistically



EDA - Season 3 Advanced Statistics

- Season 3 is when players start to differentiate themselves
- Advanced statistics such as **'Win-Share', 'Player Efficiency Rating'** and **'Value Over Replacement Player'** are key identifiers of future stars



Modeling Techniques - RFE



After Feature Engineering

- 200+ columns of player data



Feature Selection

- Correlation Analysis
- Random Forest Feature Importance
- **Recursive Feature Elimination**



After RFE

- **VORP_3**
- **PER_3**
- TS%_3
- **PTS**
- TRB
- **PER_2**
- **PTS_3**
- MP_3
- PER_2
- SPG_3
- SPG_2
- PPG_2



Modeling

● Interpretability

Only focused on models that are easily interpretable

● Precision

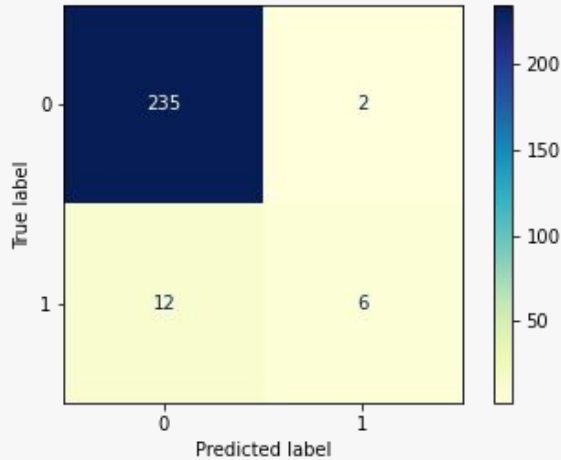
Most important metric was precision for our business question

	accuracy score	recall score	precision score	f1 score
name				
Logistic Regression	0.894118	0.833333	0.384615	0.526316
Logistic Regression w/ Resampling	0.886275	0.777778	0.358974	0.491228
Random Forest	0.945098	0.333333	0.750000	0.461538
Random Forest w/ Resampling	0.925490	0.555556	0.476190	0.400000
Decision Tree	0.894118	0.500000	0.333333	0.400000
Decision Tree w/ Resampling	0.909804	0.500000	0.391304	0.439024

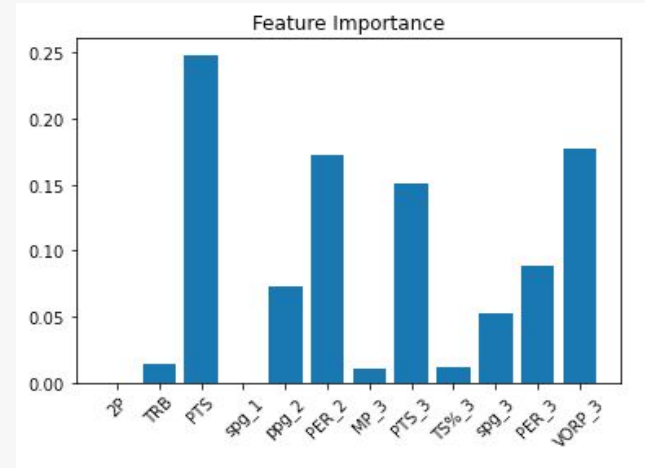


Modeling Results

Random Forest Classifier Model - TEST



Confusion Matrix



Feature Importance



Next Steps



Implement other resampling techniques to deal with class imbalances



Pull in additional categorical data such as draft pick position or team success



Increase the number of classes, break out All-NBA teams or include All Defensive teams



Look into generational trends, did All NBA players look statistically different in the 80's vs the 90's



Thank you!

