

# **High Anticipation**

Exploring Trends Between Public Perception and Player Value

> Alana Willis, Fiona Dunn, and Sahana Rayan October 25, 2020

# Guiding the Research

~ What is the best metric for defining player performance?

~ How do we quantify public perception?

~ Do any players or subset of players stand out?

~ Can we predict public perception using player performance and vise versa?

# Area of Focus

- Narrowed NBA players to 2018 first round draft class
  - All data had to be collected separately
  - Includes a randomized mix of players
    - International
    - Highly anticipated
    - Underrated



# Constructing the Data Set

#### Reddit Game Data Scraped Reddit user Scraped Box score data from comments NBA stats page during the 2018 using RedditExtractoR 2019 season 04 02 03 01 + **→** > Combined Web Hits

Scraped views from YouTube, Wikipedia, and Google (Images and News) Rows aligned by weeks a player participated in at least one game over rookie season.



# **Interesting Player Trends**



# Clustering



Scatter plot of Average Web hits and Sentiment Score

- Hierarchical Clustering was done using positive score, negative score, Wiki views, Google web hits, Google news hits, and Youtube hits
- Cluster 1: Not very popular, low sentiment
- Cluster 2: A bit more popular, high sentiment

### More EDA with the clusters



# XGBoost: Predicting Clusters

- These clusters were used as a response variable for a binary classification problem.
- XGboost model was trained using AUC as an evaluation metric and the test predictions has an AUC of 0.887



# Partial Dependence Plots



#### Partial Dependence plot with Average Points



# Final Thoughts

- Data wrangling and cleaning is HARD!
- There is a relationship between a player's performance, public sentiment, and overall popularity.
- One metric needs to be used to define public perception: Social Score and Clusters
- Average minutes played and average points scored are the most important in predicting popularity and sentiment clusters

# Continuing the Research

- Introduce more variables to the models.
  - Reddit: comment score and controversiality
- Expand to more draft classes and to more social media/chat rooms
  - Other subreddits, blogs, Facebook, ect
- Explore specific comments in the Reddit data
- Predicting public perception onto future draft classes based on college or international play
- Principal Component Analysis with sentiment and popularity data

### Contacts

### Alana Willis

- Winston-Salem State University
- awillis 117@rams.wssu.edu
- linkedin.com/in/alana-willi s-7a8014188/

#### Fiona Dunn

- Kenyon College
- dunn2@kenyon.edu
- linkedin.com/in/fiona-dunn-7 06035138/

#### Sahana Rayan

- Purdue University
- srayan@purdue.edu
- linkedin.com/in/sahana-raya n/







# References

Bresler, Alex. "NbastatR v0.1.1503." *NbastatR Package | R Documentation*, www.rdocumentation.org/packages/nbastatR/versions/0.1.1503.

Game Score, www.nba.com/resources/static/team/v2/thunder/statlab-gamescore-191201.pdf.

Massicotte, Philippe. "GtrendsR." Function | R Documentation, www.rdocumentation.org/packages/gtrendsR/versions/1.3.5/topics/gtrends.

Meissner, Peter. "Wikipediatrend v2.1.6." *Wikipediatrend Package | R Documentation,* www.rdocumentation.org/packages/wikipediatrend/versions/2.1.6.

Rivera, Ivan. Package 'RedditExtractoR,'1 May 2019, cran.r-project.org/web/packages/RedditExtractoR/RedditExtractoR.pdf.

Silge, Julia, and David Robinson. Text Mining with R, 7 Mar. 2020, www.tidytextmining.com/index.html.

# **Thank You!**



### - Carnegie Mellon University

- Rebecca Nugent
- Ron Yurko
- Beomjo Park
- Pratik Patil
- 2020 Cohort
- Atlanta Hawks
  - Maksim Horowitz

# Variables of Interest

- 1. Week: Mon-Sun of a week in the regular season of the NBA
- 2. **Game Score:** a rough measure of a player's productivity for a single game
- 3. **Popularity** = (wiki\_views + avg\_web\_hits + avg\_image\_hits + avg\_news\_hits + avg\_yt\_hits) / 5
- 4. **Social Score** = (sentiment\_week + 5 \* Popularity) / 6



#### **Rookie Season Sentiments**

# Interesting Player Trends





Trae Young's Average Game Score



Sentiment Score by 1-15 Picks

Weeks



Sentiment Score by 16-30 Picks

Luka Doncic Average Game Score vs Social Score 60 -. Social Score . 20 Anfernee Simons 0 30 20 0 10 Average Game Score

# Interesting Player Trends



Average Game Score by Player

Player Name

**Player Name** 

# **Clustering for Popularity Metrics**



•

3

• 4

2 .

Cluster • 1

Scatter plot of Average Web hits and Wiki views





# **Before Modeling**





Distribution of Average game score

# Decision Tree: Predicting Avg Game Score



# Decision Tree: Predicting Social Score



# In The Beginning

#### **Initial Research Questions**

How do players interact with their fan bases conditioned on the demographic of their home market?

Can we quantify "home town bias" among local media outlets for each NBA team?

Comparing public perception vs. actual player value (using player stats and social media data)

Packages

TidyCensus

NBAStatR

TwitteR

# Stage 1: Reddit Data

- Scraped all Reddit posts within subreddit r/nba using *RedditExtractoR* 
  - Filtered the posts and comments to match the dates of the rookie season (~October 2018 to April 2019)
- Used bing sentiment analysis on each of the comments to create a sentiment score
  - Resulted in over 2 million rows
  - Grouped comments about each player by week to reduce rows to 296

R
E
reddit

# Stage 4: Final Dataset

- The game data, Reddit data, and web data were combined to form a data set where each row's unique identifiers were the player name and the week
  - 206 observations by 25 variables
- This data set comprises of 3 types of data; Sentiment data (Reddit), popularity data (Wikipedia and Google trends), and basketball performance data (NBA stat)



# Stage 2 : Web Data

- Used *wikipediatrend* package to scrape daily Wikipedia page views for each player

   Filtered by the dates of the rookie season (~October 2018 to April 2019)
- Used *gtrends* package to scrape daily Google image hits, web hits, news hits and YouTube hits for each player
  - Filtered by the dates of the rookie season (~October 2018 to April 2019)
  - Google trends assigns scores from 0 to 100







# Stage 3: Game Data

- Scraped box score data from the NBA stats page
  - Only looking at 2018-2019 season for players drafted in top 30
- Each row in this data showed the player's performance for a game played in the season
- Game score: a measure of player performance in a given game

	BA	
S	A	S

# Random Forest: Predicting Game Score



- N\_trees = 40 80 with increments of 5
- M\_try = 1 to 6
- Min\_node\_size = 3, 5



# XGBoost: Predicting Avg Game Score



Best Tuning Parameters: nrounds = 120, eta = 0.025, max tree depth = 1



# **Clustering for Sentiment Metrics**



2

3 • 4

Cluster

#### Box plot of Average game score for the 4 clusters

