# Baskets and Ballers

## Using Machine Learning to Profile Pro Players

# The Team

Sneha

Ryan

Aaron

Shikha

Yun

★★★ 2023 ★★★

Using Machine Learning to Profile Pro Players

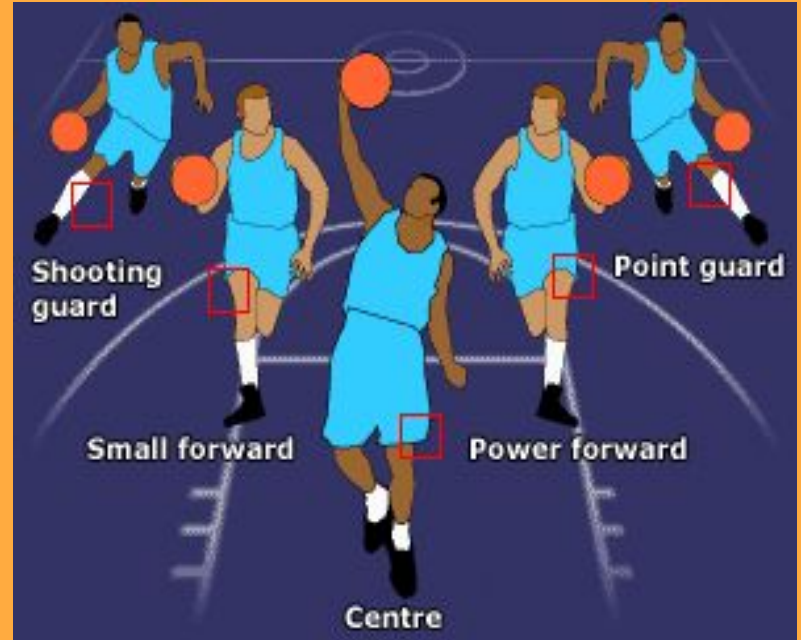**API - Pandas - Machine learning - Javascript - Flask API - Web App**

# Purpose

**Traditional groups for NBA players are based on their position in a game and the skills that position requires :**

1. **Point guard** - Assists, Steals, Three Point Shots
2. **Shooting guard** - Points, Three Point Shots, Free Throw%
3. **Small forward** - Points, Defense Rebound, Field Goal %
4. **Power forward** - Total Rebound, Field Goal %, Blocks
5. **Center** - Total Rebounds, Blocks, Field Goal %

In reality, no player has the perfect combination of skills for a specific position. So the game play is governed by how well the skills of each player compliments the others' on the team.

*Clustering allows us to group players based on observed unique set of skills represented through their performance statistics. These groups provide insights for talent evaluation, team composition, and strategic decision-making in the context of basketball management and coaching.*

# Player Stats

**19** metrics from **season 2021**, were used for each player who had at least **500 minutes** of playtime:

- **Minutes:** Total #minutes spent on court
- **PTS**: Total points scored through the season
- **FGM, FGA, FG%:** field goals made, attempted & %. All points made outside of free throws.
- **FTM, FTA, FTP%:** free throws made, attempted & %. Awarded after fouls on the player from the opposing team.
- TPM, TPA, TPM%: three-point field goals made, attempted & %.
- **totReb, offReb, defReb:** Possession after a missed shot by either offensive or defensive
- **Assists:** passes the ball to a teammate or defensive goaltending, that leads directly to a score
- **Steals:** Positive aggression that leads to turnover
- **Blocks:** legally deflects a field goal attempt from an offensive player to prevent a score
- **Turnovers:** ball stolen, out of bounds, pass intercepted, committing a violation/foul that leads to loss of team possession
- **pFouls:** rule breach that concerns illegal personal contact with an opponent

## Offensive Players

*Offensive Rebounds*

*Assists*

*Turnovers*

*Field Goals*

*Free  Throws*

## Defensive Players

*Defensive Rebounds*

*Blocks*

*Steals*

*pFouls*

# Data Sources and Cleaning

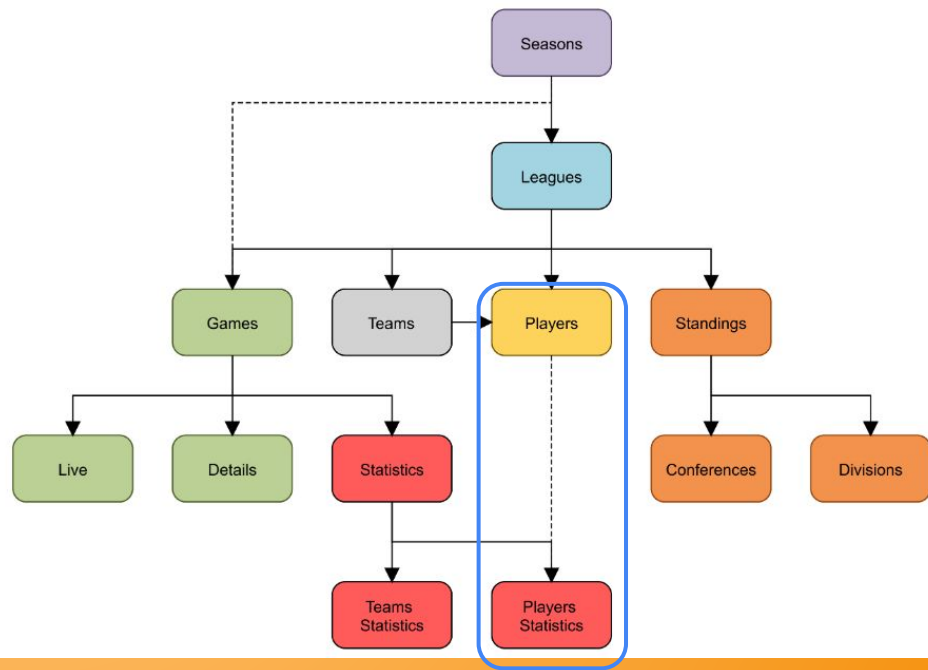**Source** - https://rapidapi.com/api-sports/api/api-nba

**Endpoints**:

- **Players** for demographic data
- **Player Statistics** for game metrics

**Steps:**

1. Player demographics from Players Endpoint
2. Player stats for above Player IDs
3. Missing value treatment and format changes
4. Aggregation at player level and recalculation of derived metrics
5. Filter players with at least 500 minutes of playtime

# Machine Learning



**Correlation Matrix**
- Field Goals , 3-Point goals , Free Throws were highly correlated - we hence retained only the derived % metrics FGP, FTP, TPP & points
- Offensive and Defensive Rebounds replaced the sum ie. Total Rebounds

**Data Normalization**

- Features were measured in different units.
- Used 'StandardScaler' to prevent any single feature from dominating the clustering process.
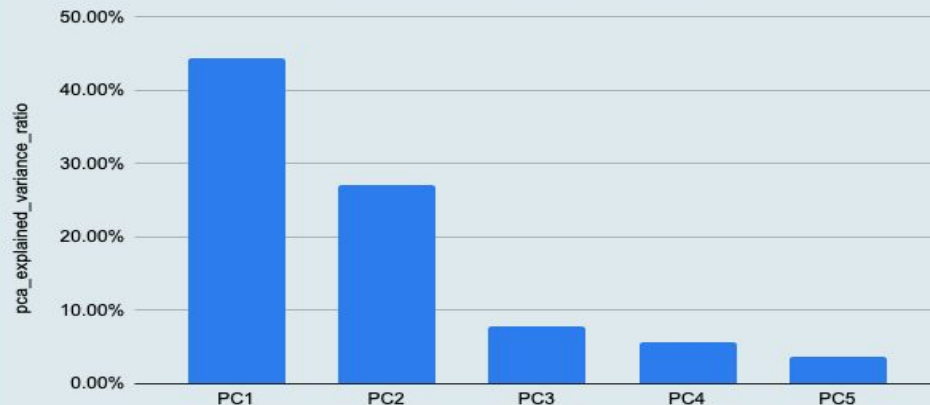
# Machine Learning



## Principal Component Analysis
- 11 metrics was the 'curse of dimensionality' in this dataset
- PCA was performed to reduce the dimensionality to 5 features hence improve the model speed and efficiency
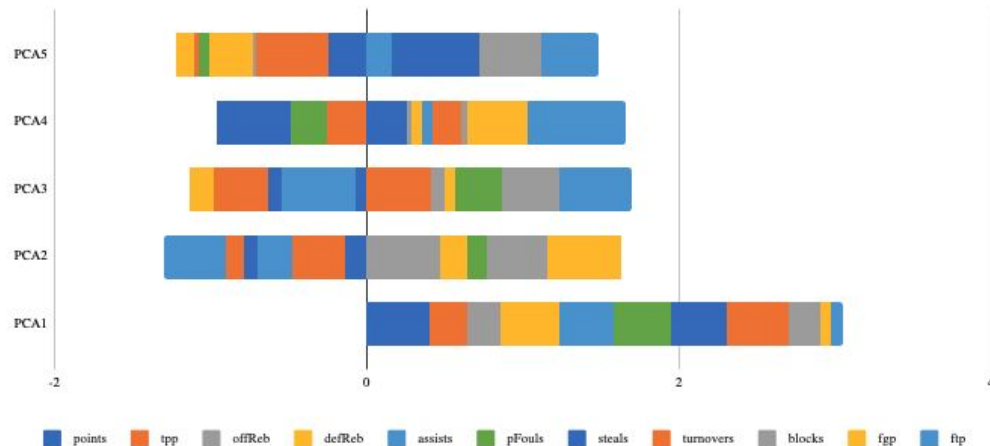
## Principal Components Results
- 5 principal components explained 88.6% of the variance
- PC1 and PC2 contributed to 71% of this variance
- PC1 and PC2 have very different distribution of weights across the 11 features as shown here



Weights assigned to each variable for each Principal Component

points  tpp  offReb  defReb  assists  pFouls  steals  turnovers  blocks  fgp  ftp

# Machine Learning
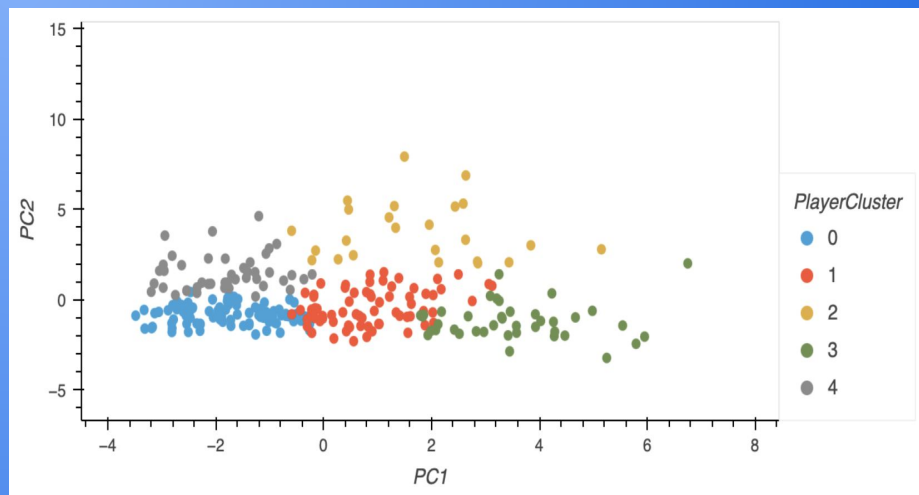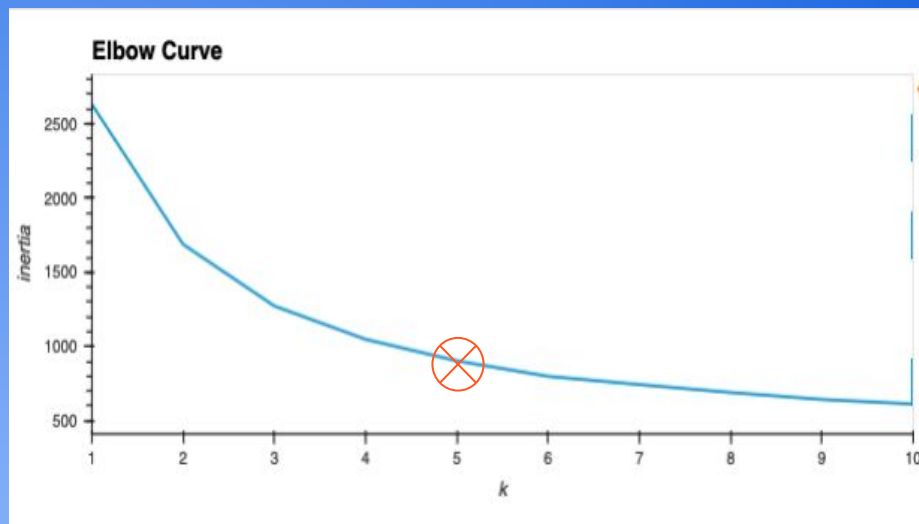
## K-Means Elbow Curve
- Using 5 PC features - the k means elbow curve determined 5 optimal clusters before drop in inertia was negligible

## K Means Clustering
- Using PC1 and PC2 that explained 71% of the variance - the distinction in clusters was clear from the scatter plot

## Data for App Visualizations:
- Data was aggregated at the Cluster, Position and Player level
- MinMaxScaler was used to ensure they could be rendered on a radar chart that requires only positive normalized values
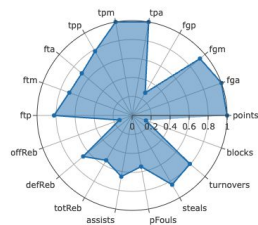


Elbow Curve

# Visualizations - Demo
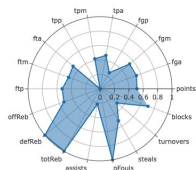
# Insights

- Skills that players presented did not necessarily match what is expected from their given position (e.g. Point Guard: there were higher free throw values than three point shots made)
- Clusters showed players had opposing metrics (e.g. Cluster 2  high offReb and blocks)
- PG, Cluster 1, & Cluster 3  groups with the highest assists also have  the highest turnover rates

# SASYR Team Neural Network Results:

✓ Open to ideas and discussion early in the process
✓ Success is not the only path to learning - as seen from the effort we put into creating spider charts in Tableau
✓ Positive attitude towards roadblocks
✓ Moral support is valuable support
✓ Recognize strengths and weaknesses are like sticks that band together to form a log

★ If it is convoluted- there's a better tool out there, use it!
★ Real world problem solving - require many iterations based on data understanding!

# Thank You