# Final Project

## National Basketball Association

**Golden Gate University**
**MSBA 324: Data and Social Media Analytics**
**Professor Dr. Stephan Sorger**

**December 21, 2024**

**<u>Team Members</u>**
**Sylvester Ronith Reagan**
**Jamario Kelly**
**Rabiul Hasan**
**Suman Thapa**

# Agenda

| | |
|---|---|
| Situation | R. Hasan |
| Problem Statement | R. Hasan |
| Model Selection | R. Hasan |
| Solution Process | R. Reagan |
| Research | R. Reagan |
| Software | R. Reagan |
| Model Results | S. Thapa |
| Visualization | S. Thapa |
| Results Interpretation | S. Thapa |
| Situation Comparison | J. Kelly |
| Conclusion | J. Kelly |
| Recommendations | J. Kelly |

# Situation

- How the NBA leverages its top-scoring players to improve viewership and fan engagement rates.

- Intent is to increase TV ratings, live-streaming stats, and fan engagement through highlight reels and social media content featuring top scoring players.

- The top 10 NBA athletes by points per game (e.g., players like James Harden, Kobe Bryant, and Joel Embiid) generate significant attention through their consistent high-scoring performances.

- Material based on Sports Business Journal (Cannon, 2023)

# Problem Statement

Objectives of project:
Analyze how the NBA's leading scorers affect the league's overall viewership and fan engagement.
Examine the factors such as  media exposure (highlight reels, social media buzz) and player performance (points per game) affect the growth of NBA TV ratings and streaming figures.

Dependent Variable:
The dependent variable will be "average TV viewership per game" or "streaming numbers"  representing the number of people who watch NBA games involving the top 10 players, with a focus on those in which these players score at or above their season average.
Since NBA viewership is a significant source of income for the league and its broadcast partners, it is important to track how well elite scorers influence fan engagement and interest.

Numerical Threshold:
The project will be deemed successful if we can demonstrate a 5% increase in average TV viewership or streaming figures for games starring the top 10 NBA players over the league's baseline viewership.

# Model Selection

Model selection:
*Model*: Regression Analysis, Cluster Analysis and Time series Analysis
*Purpose*:
Cluster Analysis: Group NBA games by fan engagement (viewership/streaming) and player performance (points per game) to find game segments that draw the most viewers, particularly those with top-scoring players.
Regression Analysis: Assess whether player performance is a significant predictor of viewership changes, validated at a 95% confidence level, by looking at the relationship between the number of points earned by top players per game and the increase in viewership and streaming.
Time Series Analysis:

Reason for selecting the model:
By combining regression and cluster analysis, it is possible to find trends in viewing data and develop a targeted strategy that prioritizes games with top-scoring players, increasing TV ratings and fan engagement.

# Solution Process

| 1 | → | 2 | → | 3 | → | 4 | → | 5 | → | 6 |

Step 1. Collect and examine data on viewership and player performance for NBA matches, particularly focusing on the leading 10 players according to points scored per game. Collect information on television ratings, streaming statistics, and social media interaction for every game.

Step 2. Conduct regression analysis to evaluate the correlation between points per game for the leading players and the change in TV ratings or streaming figures. Compute the p-values for the variables to assess statistical significance and pinpoint the key factors influencing viewership.

Step 3. Perform cluster analysis to identify patterns in the data, categorizing games that showcase top-scoring players, and evaluate the impact on viewership. Document the key variables, including player performance (points scored per game), team achievements, and the location of the game.

Step 4. Perform Time Series Analysis after identifying the Trend and predict the future based on trend, marketing campaigns promoting top-producing players

Step 5. Research the impact of media exposure, including highlight videos and social media posts, in engaging  fans. Analyze trends and compare them with other sports leagues or media initiatives to identify successful approaches for enhancing audience engagement.

Step 6. Provide conclusions and suggestions for the NBA based on the analysis, including strategies for promoting games with top-scoring athletes to enhance television ratings and streaming figures.

# Research

Secondary Research:
- Project used sources from social media blogs or articles and company social media.
- Purpose was to provide additional data and insight for the case study.
- Sources are cited on slides where they are used and are listed in References section.

Primary Research:
- Data was obtained from Kaggle.com titled "NBA Database" by Wyatt Walsh.

# Research

Data Source
Kaggle.com " NBA Database"

Definitions (source: NBA website)
AST: Number of assists that lead directly to a made basket
REB: Number of recoveries made by team after missed shot
TS%:True shooting percentage
PTS: Total points scored.

# Research

## Snapshot of .csv file



| | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | player_na | team_abb | age | player_hei | player_we | college | country | draft_year | draft_rour | draft_num | gp | pts | reb | ast | net_rating | oreb_pct | dreb_pct | usg_pct | ts_pct | ast_pct | season |
| 2 | 0 Randy Livi | HOU | 22 | 193.04 | 94.80073 | Louisiana | USA | 1996 | 2 | 42 | 64 | 3.9 | 1.5 | 2.4 | 0.3 | 0.042 | 0.071 | 0.169 | 0.487 | 0.248 | 1996-97 |
| 3 | 1 Gaylon Ni | WAS | 28 | 190.5 | 86.18248 | Northwest | USA | 1994 | 2 | 34 | 4 | 3.8 | 1.3 | 0.3 | 8.9 | 0.03 | 0.111 | 0.174 | 0.497 | 0.043 | 1996-97 |
| 4 | 2 George Ly | VAN | 26 | 203.2 | 103.419 | North Car | USA | 1993 | 1 | 12 | 41 | 8.3 | 6.4 | 1.9 | -8.2 | 0.106 | 0.185 | 0.175 | 0.512 | 0.125 | 1996-97 |
| 5 | 3 George Mc | LAL | 30 | 203.2 | 102.0582 | Florida St | USA | 1989 | 1 | 7 | 64 | 10.2 | 2.8 | 1.7 | -2.7 | 0.027 | 0.111 | 0.206 | 0.527 | 0.125 | 1996-97 |
| 6 | 4 George Zic | DEN | 23 | 213.36 | 119.7483 | UCLA | USA | 1995 | 1 | 22 | 52 | 2.8 | 1.7 | 0.3 | -14.1 | 0.102 | 0.169 | 0.195 | 0.5 | 0.064 | 1996-97 |
| 7 | 5 Gerald Wi | ORL | 33 | 198.12 | 102.0582 | Tennessee | USA | 1985 | 2 | 47 | 80 | 10.6 | 2.2 | 2.2 | -5.8 | 0.031 | 0.064 | 0.203 | 0.503 | 0.143 | 1996-97 |
| 8 | 6 Gheorghe | WAS | 26 | 231.14 | 137.4384 | None | USA | 1993 | 2 | 30 | 73 | 10.6 | 6.6 | 0.4 | 6.9 | 0.098 | 0.217 | 0.185 | 0.618 | 0.024 | 1996-97 |
| 9 | 7 Glen Rice | CHH | 30 | 203.2 | 99.79024 | Michigan | USA | 1989 | 1 | 4 | 79 | 26.8 | 4 | 2 | 3.2 | 0.025 | 0.087 | 0.272 | 0.605 | 0.088 | 1996-97 |
| 10 | 8 Glenn Rob | MIL | 24 | 200.66 | 106.5941 | Purdue | USA | 1994 | 1 | 1 | 80 | 21.1 | 6.3 | 3.1 | -2.9 | 0.051 | 0.144 | 0.278 | 0.528 | 0.146 | 1996-97 |
| 11 | 9 Grant Hill | DET | 24 | 203.2 | 102.0582 | Duke | USA | 1994 | 1 | 3 | 80 | 21.4 | 9 | 7.3 | 6.9 | 0.049 | 0.232 | 0.283 | 0.556 | 0.356 | 1996-97 |
| 12 | 10 Gary Trent | POR | 22 | 203.2 | 113.398 | Ohio | USA | 1995 | 1 | 11 | 82 | 10.8 | 5.2 | 1.1 | 2.5 | 0.101 | 0.167 | 0.212 | 0.569 | 0.077 | 1996-97 |
| 13 | 11 Grant Lon | DET | 31 | 205.74 | 112.4908 | Eastern M | USA | 1988 | 2 | 33 | 65 | 5 | 3.4 | 0.6 | 4 | 0.096 | 0.15 | 0.154 | 0.523 | 0.058 | 1996-97 |
| 14 | 12 Greg Anth | VAN | 29 | 185.42 | 81.64656 | Nevada-La | USA | 1991 | 1 | 12 | 65 | 9.5 | 2.8 | 6.3 | -9.4 | 0.015 | 0.099 | 0.177 | 0.526 | 0.358 | 1996-97 |
| 15 | 13 Greg Dreil | DAL | 33 | 215.9 | 120.2019 | Kansas | USA | 1986 | 2 | 26 | 40 | 2 | 1.9 | 0.3 | -8 | 0.059 | 0.192 | 0.114 | 0.466 | 0.048 | 1996-97 |
| 16 | 14 Greg Foste | UTA | 28 | 210.82 | 113.398 | Texas-El P | USA | 1990 | 2 | 35 | 79 | 3.5 | 2.4 | 0.4 | -0.9 | 0.078 | 0.166 | 0.168 | 0.508 | 0.055 | 1996-97 |
| 17 | 15 Greg Grah | SEA | 26 | 193.04 | 82.55374 | Indiana | USA | 1993 | 1 | 17 | 28 | 3.3 | 0.5 | 0.4 | 3.6 | 0.013 | 0.063 | 0.245 | 0.476 | 0.1 | 1996-97 |

# Software

## Analysis
Regression Analysis in R
Cluster Analysis in R
Time Series Analysis in R

## Techniques
- linear model (lm)
- Clustering (kmeans)
- Time Series Forecasting (auto.arima)

```r
# Load necessary libraries
library(ggplot2)

library(dplyr)

# Load the dataset
nba_data <- read.csv("C:/Users/Ronit/OneDrive/Desktop/NBA dataset.csv")

# View the first few rows
head(nba_data)
```

|   | X<br><int> | player_name<br><chr> | team_abbreviation<br><chr> | a...<br><int> | player_height<br><dbl> | player_weight<br><dbl> | college<br><chr> | |
|---|---|---|---|---|---|---|---|---|
| 1 | 0 | Randy Livingston | HOU | 22 | 193.04 | 94.80073 | Louisiana State | |
| 2 | 1 | Gaylon Nickerson | WAS | 28 | 190.50 | 86.18248 | Northwestern Oklahoma | |
| 3 | 2 | George Lynch | VAN | 26 | 203.20 | 103.41898 | North Carolina | |
| 4 | 3 | George McCloud | LAL | 30 | 203.20 | 102.05820 | Florida State | |
| 5 | 4 | George Zidek | DEN | 23 | 213.36 | 119.74829 | UCLA | |
| 6 | 5 | Gerald Wilkins | ORL | 33 | 198.12 | 102.05820 | Tennessee-Chattanooga | |

6 rows | 1-8 of 22 columns

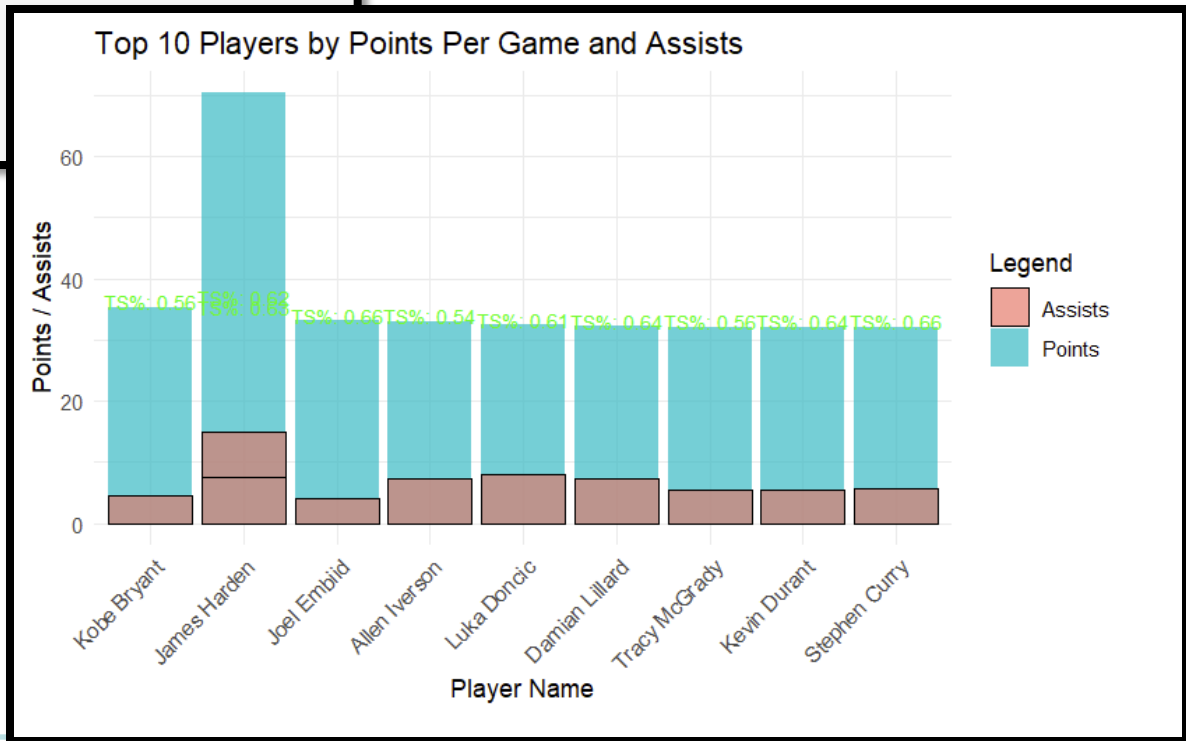# Model Results – Players Comparison

```
# Select the top 10 players based on points per game (pts)
top_players <- nba_data %>%
  arrange(desc(pts)) %>%
  slice(1:10) %>%
  select(player_name, pts, ast, ts_pct)

# View the top players
print(top_players)
```

| player_name<br><chr> | pts<br><dbl> | ast<br><dbl> | ts_pct<br><dbl> |
|---|---|---|---|
| James Harden | 36.1 | 7.5 | 0.616 |
| Kobe Bryant | 35.4 | 4.5 | 0.559 |
| James Harden | 34.3 | 7.5 | 0.626 |
| Joel Embiid | 33.1 | 4.2 | 0.655 |
| Allen Iverson | 33.0 | 7.4 | 0.543 |
| Luka Doncic | 32.4 | 8.0 | 0.609 |
| Damian Lillard | 32.2 | 7.3 | 0.645 |
| Tracy McGrady | 32.1 | 5.5 | 0.564 |
| Kevin Durant | 32.0 | 5.5 | 0.635 |
| Stephen Curry | 32.0 | 5.8 | 0.655 |

# Model Visualization - Players Comparison

```
# Create a bar chart
ggplot(top_players, aes(x = reorder(player_name, -pts))) +
  geom_bar(aes(y = pts, fill = "Points"), stat = "identity", alpha = 0.7) +
  geom_bar(aes(y = ast, fill = "Assists"), stat = "identity", alpha = 0.7, color = "black") +
  geom_text(aes(y = pts + 1, label = paste0("TS%: ", round(ts_pct, 2))), size = 3, vjust = 0.5, color = "green") +
  labs(
    title = "Top 10 Players by Points Per Game and Assists",
    x = "Player Name",
    y = "Points / Assists",
    fill = "Legend"
  ) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Top 10 Players by Points Per Game and Assists

# Model Results – Teams Comparison

```
# Group by team abbreviation and calculate average statistics for each team
team_comparison <- nba_data %>%
  group_by(team_abbreviation) %>%
  summarise(
    avg_pts = mean(pts, na.rm = TRUE),
    avg_ast = mean(ast, na.rm = TRUE),
    avg_reb = mean(reb, na.rm = TRUE),
    avg_ts_pct = mean(ts_pct, na.rm = TRUE)
  )

# View the summarized data
print(team_comparison)
```

| team_abbreviation | avg_pts | avg_ast | avg_reb | avg_ts_pct |
| --- | --- | --- | --- | --- |
| <chr> | <dbl> | <dbl> | <dbl> | <dbl> |
| ATL | 7.860137 | 1.700683 | 3.479271 | 0.5089362 |
| BKN | 8.564000 | 1.930000 | 3.536500 | 0.5239200 |
| BOS | 8.198824 | 1.835765 | 3.528000 | 0.5263953 |
| CHA | 8.123607 | 1.826230 | 3.521311 | 0.5099607 |
| CHH | 7.883146 | 1.865169 | 3.512360 | 0.4789101 |
| CHI | 8.141135 | 1.896217 | 3.610402 | 0.5018203 |
| CLE | 7.903111 | 1.760667 | 3.521778 | 0.5056067 |
| DAL | 8.120993 | 1.757111 | 3.445824 | 0.5145169 |
| DEN | 8.434813 | 1.906776 | 3.659579 | 0.5213107 |
| DET | 7.991408 | 1.792124 | 3.437709 | 0.5096539 |

# Model Visualization – Teams Comparison

```r
# Create a bar chart comparing teams by average points, assists, and rebounds
ggplot(team_comparison, aes(x = reorder(team_abbreviation, -avg_pts))) +
  geom_bar(aes(y = avg_pts, fill = "Average Points"), stat = "identity", alpha = 0.7) +
  geom_bar(aes(y = avg_ast, fill = "Average Assists"), stat = "identity", alpha = 0.7, color = "black") +
  geom_bar(aes(y = avg_reb, fill = "Average Rebounds"), stat = "identity", alpha = 0.7, color = "black") +
  labs(
    title = "Team Comparison by Average Points, Assists, and Rebounds",
    x = "Team Abbreviation",
    y = "Average Stats",
    fill = "Legend"
  ) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Team Comparison by Average Points, Assists, and Rebounds

# Model Results – Linear Regression

```
# Simple Linear Regression: points as a function of assists
simple_model <- lm(pts ~ ast, data = nba_data)

# Summary of the regression model
summary(simple_model)
```

```
Call:
lm(formula = pts ~ ast, data = nba_data)

Residuals:
    Min      1Q  Median      3Q     Max
-19.5648 -2.9505 -0.8724  2.3519 21.2496

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.16274    0.05649   73.69   <2e-16 ***
ast          2.21948    0.02204  100.72   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.497 on 12842 degrees of freedom
Multiple R-squared:  0.4413,    Adjusted R-squared:  0.4413
F-statistic: 1.014e+04 on 1 and 12842 DF,  p-value: < 2.2e-16
```

**Coefficients**

| Variables | coefficient | p –value |
|---|---|---|
| (Intercept) | 4.16274 | 2e-16 |
| AST | 2.21948 | 2.2e-16 |

# Model Visualization – Linear regression

```
# Plot the regression
ggplot(nba_data, aes(x = ast, y = pts)) +
  geom_point() +
  geom_smooth(method = "lm", col = "red") +
  labs(title = "Simple Linear Regression: Points vs Assists", x = "Assists", y = "Points")
```



Simple Linear Regression: Points vs Assists

# Model Results – Multiple Linear Regression

```
# Multiple Linear Regression: points as a function of assists, rebounds, and true shooting percentage
multiple_model <- lm(pts ~ ast + reb + ts_pct, data = nba_data)

# Show the regression model summary
summary(multiple_model)
```

**Coefficients**

| Variables | coefficient | p-value |
|-----------|-------------|---------|
| (Intercept) | -3.15179 | 2e-16 |
| AST | 1.76543 | 2e-16 |
| REB | 1.09159 | 2e-16 |
| TS% | 8.29919 | 2e-16 |

```
Call:
lm(formula = pts ~ ast + reb + ts_pct, data = nba_data)

Residuals:
     Min      1Q   Median      3Q     Max
-18.1709  -1.8225  -0.3797  1.4381  20.1827

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)  -3.15179    0.15360  -20.52   <2e-16 ***
ast           1.76543    0.01720  102.65   <2e-16 ***
reb           1.09159    0.01296   84.22   <2e-16 ***
ts_pct        8.29919    0.31062   26.72   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.382 on 12840 degrees of freedom
Multiple R-squared:  0.6842,    Adjusted R-squared:  0.6841
F-statistic:  9273 on 3 and 12840 DF,  p-value: < 2.2e-16
```

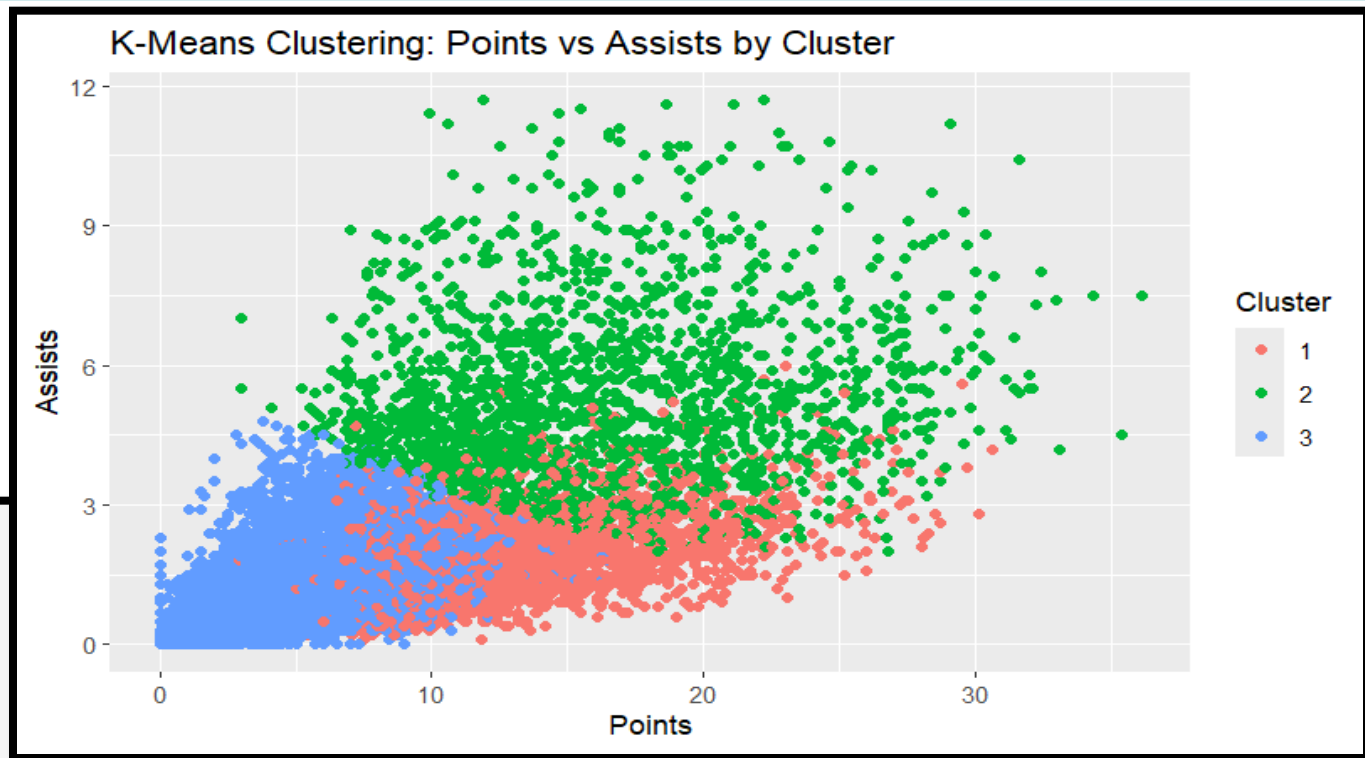# Model Visualization – K-Means Clustering



K-Means Clustering: Points vs Assists by Cluster

```r
# Select the columns for clustering
nba_clustering_data <- nba_data[, c("pts", "ast", "reb")]

# Normalize the data
nba_clustering_data <- scale(nba_clustering_data)

# Perform K-means clustering (for 3 clusters, for example)
set.seed(123)  # For reproducibility
kmeans_result <- kmeans(nba_clustering_data, centers = 3)

# Add the cluster assignment to the data
nba_data$cluster <- kmeans_result$cluster

# Visualize the clusters
ggplot(nba_data, aes(x = pts, y = ast, color = as.factor(cluster))) +
  geom_point() +
  labs(title = "K-Means Clustering: Points vs Assists by Cluster", x = "Points", y = "Assists", color = "Cluster")
```
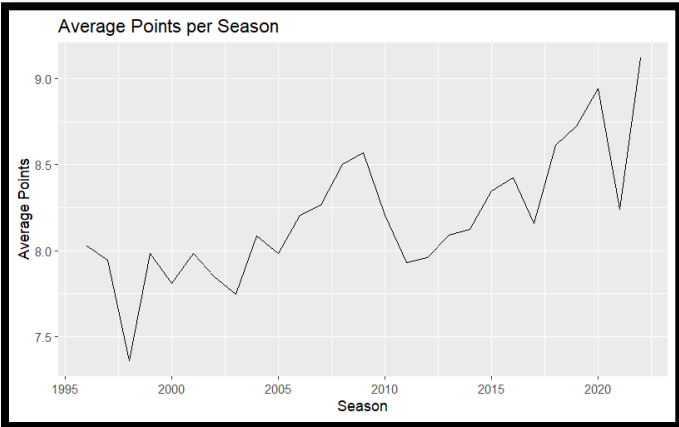
# Model Visualization – Avg points/season and Forecast

```r
# Aggregate average points per season
seasonal_data <- nba_data %>%
  group_by(season_numeric) %>%
  summarize(avg_pts = mean(pts, na.rm = TRUE))

# Create a time series object
ts_data <- ts(seasonal_data$avg_pts, start = min(seasonal_data$season_numeric), frequency = 1)

# Plot the time series
autoplot(ts_data) +
  ggtitle("Average Points per Season") +
  xlab("Season") +
  ylab("Average Points")
```
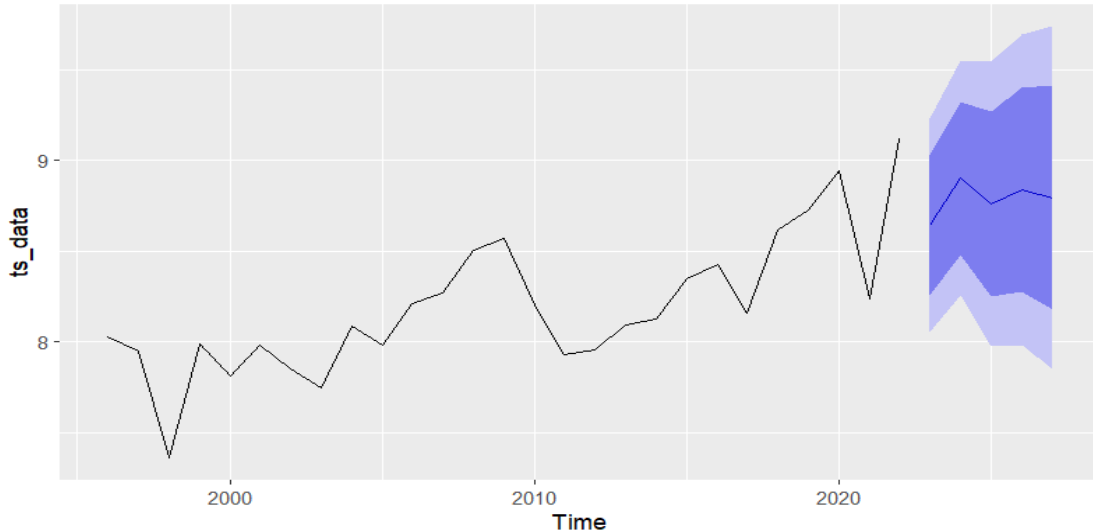
```r
model <- auto.arima(ts_data)
forecasted <- forecast(model, h = 5)
autoplot(forecasted)
```



Average Points per Season



Forecasts from ARIMA(1,1,0)

# Results Interpretation – Regression Analysis

- TS% is critical and highly significant. Player with high TS% scoring more points. These players are high contributor to team success.

- AST values positively relates to points scored. Players who can score and assist are important to team's success.

# Results Interpretation

Summarizing the output:

| Variables | coefficient | p-value |
|-----------|-------------|---------|
| (Intercept) | -3.15179 | 2e-16 |
| AST | 1.76543 | 2e-16 |
| REB | 1.09159 | 2e-16 |
| TS% | 8.29919 | 2e-16 |

```
Call:
lm(formula = pts ~ ast + reb + ts_pct, data = nba_data)

Residuals:
     Min      1Q  Median      3Q     Max
-18.1709 -1.8225 -0.3797  1.4381 20.1827

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -3.15179    0.15360  -20.52   <2e-16 ***
ast          1.76543    0.01720  102.65   <2e-16 ***
reb          1.09159    0.01296   84.22   <2e-16 ***
ts_pct       8.29919    0.31062   26.72   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.382 on 12840 degrees of freedom
Multiple R-squared:  0.6842,    Adjusted R-squared:  0.6841
F-statistic:  9273 on 3 and 12840 DF,  p-value: < 2.2e-16
```

Regression equation:

PTS = (intercept) + (coefficient of AST) * (AST) + (coefficient of REB) * (REB) + (coefficient of TS%) * (TS%)

# Results Interpretation

Regression equation:
PTS = (intercept) + (coefficient of AST) * (AST) + (coefficient of REB) * (REB) + (coefficient of TS%) * (TS%)

PTS = -3.15179 + 1.76543 * (AST) + 1.09159 * (REB) + 8.29919 * (TS%)

For average values:
PTS = -3.15179 + 1.76543 * (AST) + 1.09159 * (REB) + 8.29919 * (TS%)

# Results Interpretation – Cluster Analysis

- James Harden & Devin Booker are some players dominating scorecard but highly dependent on teammates for AST.

- Players like Chris Paul are pivotal playmakers for setting up other players for high impact plays.

- Rebounders are defensive assets (Ruby Gobert) and provides crucial second chance opportunities.

# Results Interpretation – Time Series Analysis

- League trend shows average point per game steadily increased. As per forecasted analysis, expectation is 5% growth in next five seasons.

- Increase in average points per game aligns with fans expectations for offensive gameplay.

# Results Interpretation

- If a player had zero assist, zero rebounds and zero TS%, he would end up with -3.15 points which is a nonrealistic scenario.
- Each additional assists per game would increase total points by 1.77 on average. (if Rebounds and TS % remain constant)
- Each addition Rebounds would increase the total points by 1.09. (if Assists and TS% remain constant).
- Each additional unit increase of TS% increase points by 8.3. TS% typically ranges from 0.45 to 0.70, if TS % increases by 0.01, the points increase by 0.083 meaning small improvement in TS% translate into significant rise of points.
- With R-squared value at 0.6842 meaning 68.42 % of variation in points are explained by assists, rebounds and TS %. And being p value at 2e-16 shows these predictors (assists, rebounds and TS% are significant).

# Situation Comparison

Comparison - National Football League (NFL)

- NFL uses top performing players like Patrick Mahomes and Tom Brady to draw massive viewers particularly during playoffs or Superbowl games. As expected, we can see significant increment in engagement around NFL games. (Ministryofsport, 2024).

- NFL achieved 12% TV rating increase after featuring teams like Kansas City Chiefs. Primetime games featuring popular teams and players captures large audiences and contributes to rising viewership. (Fisher, 2024).

- NFL partners with many streaming platforms to expand its digital presence and provides exclusive streaming options. Real-time highlights, behind-the-scene footages and player interviews streamed over these platforms engages audiences and increase viewership. (Ministryofsport, 2024).

# Conclusion

Problem Solved
- Top Scoring Players impact significantly on TV ratings and fan engagement.
- True Shooting % (TS%), Assists (AST) and Rebounds (REB) key drivers of player performance.
- 5% NBA viewership rise over next five seasons is forecasted. Media exposure and games schedule featuring top performing players are crucial factors to enhance audience numbers and viewership.

Lessons Learned
- Key and Top producing players are to be highlighted in promotional content to amplify fan's interest in games and encourage them to view games.
- Targeted marketing campaigns maximizes audience retention and increases viewership.
- Expanding digital presence plays significant role in increasing viewership.

# Recommendations

- Games should be scheduled and promoted featuring high-performing players or teams. It helps to optimize viewership.

- The marketing content should highlight players who has high TS%.

- Predictive insights such as forecasted analysis and insights should be used to guide promotional campaigns and develop strategies accordingly.

- Social media platforms are huge which can be leveraged to advertise the campaigns that highlights top-performing players. It helps to increase engagement.

# References

- Cannon, Julian. (2023). "Inside the NBA's social media and OOH strategies for the NBA Finals" . *DIGIDAY*. https://digiday.com/marketing/inside-the-nbas-social-media-and-ooh-strategies-for-the-nba-finals/
- Fisher, S. (January 7, 2024). NFL dominates what's left of live TV viewership. Axios. https://www.axios.com/2024/01/07/nfl-tv-ratings-live-events-viewership
- https://www.kaggle.com/datasets/wyattowalsh/basketball
- https://www.nba.com/stats/help/glossary#tspct
- https://ministryofsport.com/nfl-tv-ratings-surge-to-nine-year-high/