# Exploring NBA Historical Data 1977-2020

Jerry Hong, Akshay Reddy, and Uziel Rios

12/9/2020

## Introducing our Data Set

This data set includes statistics of every NBA player from each team from the 1977-2020 (courtesy of FiveThirtyEight). Additionally, we included FiveThirtyEight's glossary to better understand what each category represents. You may notice that the five major stats are per 36 minutes instead of per game. One possible reason FiveThirtyEight measures this is to evaluate the player's performance on a per-minute basis. Plus, most starters average roughly 36-37 minutes per game, so rounding it to 36 minutes gives a decent representation of the player's production on the court.

```
nba_set <- read_csv("./nba-data-historical.csv")
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   player_id = col_character(),
##   name_common = col_character(),
##   pos = col_character(),
##   team_id = col_character(),
##   franch_id = col_character()
## )
```

```
## See spec(...) for full column specifications.
```

```
nba_set
```

```
## # A tibble: 20,059 x 40
##    player_id name_common year_id   age pos   team_id tmRtg franch_id     G   Min
##    <chr>     <chr>         <dbl> <dbl> <chr> <chr>   <dbl> <chr>     <dbl> <dbl>
##  1 youngtr01 Trae Young     2020    21 PG    ATL      -7.6 ATL          60  2120
##  2 huntede01 De'Andre H~    2020    22 SF    ATL      -7.6 ATL          63  2018
##  3 huertke01 Kevin Huer~    2020    21 SG    ATL      -7.6 ATL          56  1760
##  4 reddica01 Cam Reddish    2020    20 SF    ATL      -7.6 ATL          58  1551
##  5 collijo01 John Colli~    2020    22 PF    ATL      -7.6 ATL          41  1363
##  6 bembrde01 DeAndre' B~    2020    25 SG    ATL      -7.6 ATL          43   915
##  7 jonesda03 Damian Jon~    2020    24 C     ATL      -7.6 ATL          55   887
##  8 cartevi01 Vince Cart~    2020    43 SF    ATL      -7.6 ATL          60   876
##  9 parkeja01 Jabari Par~    2020    24 PF    ATL      -7.6 ATL          32   837
## 10 lenal01   Alex Len       2020    26 C     ATL      -7.6 ATL          40   745
## # ... with 20,049 more rows, and 30 more variables: `MP%` <dbl>, MPG <dbl>,
## #   `P/36` <dbl>, `TS%` <dbl>, `A/36` <dbl>, `R/36` <dbl>, `SB/36` <dbl>,
```

```
## #   `TO/36` <dbl>, `Raptor O` <dbl>, `Raptor D` <dbl>, `Raptor+/-` <dbl>,
## #   `Raptor WAR` <dbl>, `PIE%` <dbl>, `AWS%` <dbl>, `USG%` <dbl>, `AST%` <dbl>,
## #   `TOV%` <dbl>, `ORB%` <dbl>, `DRB%` <dbl>, `TRB%` <dbl>, `STL%` <dbl>,
## #   `BLK%` <dbl>, ORtg <dbl>, `%Pos` <dbl>, DRtg <dbl>, `2P%` <dbl>,
## #   `3P%` <dbl>, `FT%` <dbl>, 3PAr <dbl>, FTAr <dbl>
```

```
knitr::include_graphics("./Categories.png")
```

| Category | Description |
|----------|-------------|
| player_id | Basketball-Reference player ID |
| name_common | Player name |
| year_id | Season |
| age | Age as of Feb. 1 |
| pos | Position |
| team_id | Team ID |
| tmRtg | Team net rating |
| franch_id | Franchise ID |
| G | Games |
| Min | Minutes |
| MP% | % of available team minutes played |
| MPG | Minutes per game |
| P/36 | Points per 36 min (pace adj) |
| TS% | True Shooting % |
| A/36 | Assists per 36 min (pace adj) |
| R/36 | Rebounds per 36 min (pace adj) |
| SB/36 | Steals + Blocks per 36 min (pace adj) |
| TO/36 | Turnovers per 36 min (pace adj) |
| Raptor O | Offensive RAPTOR |
| Raptor D | Defensive RAPTOR |
| Raptor+/- | Overall RAPTOR |
| Raptor WAR | RAPTOR Wins Above Replacement |
| PIE% | Player Impact Estimate |
| AWS% | PIE%, but using Alternate Win Score as the metric |
| USG% | Usage Rate |
| AST% | Assist Rate |
| TOV% | Turnover Rate |
| ORB% | Offensive Rebounding Rate |
| DRB% | Defensive Rebounding Rate |
| TRB% | Total Rebounding Rate |
| STL% | Steal percentage |
| BLK% | Block rate |
| ORtg | Individual Pts created/100 possessions |
| %Pos | Share of team possessions used on court |
| DRtg | Individual Pts allowed/100 possessions |
| 2P% | 2-point FG% |
| 3P% | 3-point FG% |
| FT% | Free throw % |
| 3PAr | 3PA/FGA |
| FTAr | FTA/FGA |

```
str(nba_set)
```

```
## tibble [20,059 x 40] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
##  $ player_id  : chr [1:20059] "youngtr01" "huntede01" "huertke01" "reddica01" ...
##  $ name_common: chr [1:20059] "Trae Young" "De'Andre Hunter" "Kevin Huerter" "Cam Reddish" ...
##  $ year_id    : num [1:20059] 2020 2020 2020 2020 2020 2020 2020 2020 2020 2020 ...
##  $ age        : num [1:20059] 21 22 21 20 22 25 24 43 24 26 ...
##  $ pos        : chr [1:20059] "PG" "SF" "SG" "SF" ...
##  $ team_id    : chr [1:20059] "ATL" "ATL" "ATL" "ATL" ...
##  $ tmRtg      : num [1:20059] -7.6 -7.6 -7.6 -7.6 -7.6 -7.6 -7.6 -7.6 -7.6 -7.6 ...
##  $ franch_id  : chr [1:20059] "ATL" "ATL" "ATL" "ATL" ...
##  $ G          : num [1:20059] 60 63 56 58 41 43 55 60 32 40 ...
##  $ Min        : num [1:20059] 2120 2018 1760 1551 1363 ...
##  $ MP%        : num [1:20059] 65.1 62 54.1 47.6 41.9 28.1 27.2 26.9 25.7 22.9 ...
##  $ MPG        : num [1:20059] 35.3 32 31.4 26.7 33.2 21.3 16.1 14.6 26.2 18.6 ...
##  $ P/36       : num [1:20059] 29.3 13.5 13.6 13.7 22.7 9.5 12.1 11.9 20 16.4 ...
##  $ TS%        : num [1:20059] 59.5 52.1 53.6 50 65.9 50 71.2 47 56.5 59 ...
##  $ A/36       : num [1:20059] 9.2 1.9 4.2 2 1.5 3.1 1.4 1.9 2.4 2 ...
##  $ R/36       : num [1:20059] 4.2 5 4.5 4.9 10.7 5.8 8.1 4.9 8 10.8 ...
##  $ SB/36      : num [1:20059] 1.2 1.1 1.5 2 2.5 2.8 2.6 1.9 2.4 2.5 ...
##  $ TO/36      : num [1:20059] 4.8 1.8 1.7 2.2 1.9 2.3 1.1 1.3 2.5 1.8 ...
##  $ Raptor O   : num [1:20059] 7.1 -2.5 -0.4 -2.8 0 -4.2 -1.5 -2.7 0.9 -2.4 ...
##  $ Raptor D   : num [1:20059] -3.5 -1.3 -2.4 -0.1 -0.3 2.4 -2.6 -1.1 -1.7 1.6 ...
##  $ Raptor+/-  : num [1:20059] 3.6 -3.8 -2.8 -3 -0.3 -1.7 -4.2 -3.8 -0.8 -0.8 ...
##  $ Raptor WAR : num [1:20059] 7 -1.1 -0.1 -0.2 1.7 0.5 -0.7 -0.5 0.8 0.8 ...
##  $ PIE%       : num [1:20059] 17 5.9 8 5.9 15.6 6.4 8.3 4.8 11.2 11.2 ...
##  $ AWS%       : num [1:20059] 15.4 4.7 8.1 5 17.1 6.2 10.2 4.3 11.2 12.1 ...
##  $ USG%       : num [1:20059] 34.9 17.5 17.1 18.9 22.7 14 11.4 16.6 24 18.6 ...
##  $ AST%       : num [1:20059] 45.6 8 17.5 8 7.6 12.3 5.6 7.7 11.6 8.7 ...
##  $ TOV%       : num [1:20059] 16.2 12.1 12 13.6 10.1 19.4 11.5 9.4 12.4 11.7 ...
##  $ ORB%       : num [1:20059] 1.6 2.3 2.1 2.4 9 3.9 8.9 2.2 6.9 9.8 ...
##  $ DRB%       : num [1:20059] 11.5 13.1 12 12.7 24 14.1 16.2 13.2 17.9 23.6 ...
##  $ TRB%       : num [1:20059] 6.5 7.6 7 7.5 16.4 8.9 12.5 7.6 12.3 16.6 ...
##  $ STL%       : num [1:20059] 1.4 1 1.4 1.9 1.1 2.8 1.4 1.2 2.3 1.3 ...
##  $ BLK%       : num [1:20059] 0.3 0.7 1.3 1.5 4.1 1.7 3.8 2.4 1.5 3.7 ...
##  $ ORtg       : num [1:20059] 113.6 99.5 107.1 94.7 123.7 ...
##  $ %Pos       : num [1:20059] 36.1 16.9 17.2 18.3 21.6 14.6 11.7 15.9 23.3 18.7 ...
##  $ DRtg       : num [1:20059] 117 117 116 115 112 ...
##  $ 2P%        : num [1:20059] 50.1 45.4 45.3 42.8 64.2 54.6 70.4 45.1 60.1 62.7 ...
##  $ 3P%        : num [1:20059] 36.1 35.5 38 33.2 40.1 23.1 22.2 30.2 27 25 ...
##  $ FT%        : num [1:20059] 86 76.4 82.8 80.2 80 54.2 73.8 79.3 73.6 63 ...
##  $ 3PAr       : num [1:20059] 45.5 44.5 54.8 45.1 24.3 28.5 5.1 66.4 29.3 21.5 ...
##  $ FTAr       : num [1:20059] 44.8 21.1 10.5 22.7 24.8 21.1 47.2 9.5 18.3 31.2 ...
##  - attr(*, "spec")=
##   .. cols(
##   ..   player_id = col_character(),
##   ..   name_common = col_character(),
##   ..   year_id = col_double(),
##   ..   age = col_double(),
##   ..   pos = col_character(),
##   ..   team_id = col_character(),
##   ..   tmRtg = col_double(),
##   ..   franch_id = col_character(),
```

```
##   ..   G = col_double(),
##   ..   Min = col_double(),
##   ..   `MP%` = col_double(),
##   ..   MPG = col_double(),
##   ..   `P/36` = col_double(),
##   ..   `TS%` = col_double(),
##   ..   `A/36` = col_double(),
##   ..   `R/36` = col_double(),
##   ..   `SB/36` = col_double(),
##   ..   `TO/36` = col_double(),
##   ..   `Raptor O` = col_double(),
##   ..   `Raptor D` = col_double(),
##   ..   `Raptor+/-` = col_double(),
##   ..   `Raptor WAR` = col_double(),
##   ..   `PIE%` = col_double(),
##   ..   `AWS%` = col_double(),
##   ..   `USG%` = col_double(),
##   ..   `AST%` = col_double(),
##   ..   `TOV%` = col_double(),
##   ..   `ORB%` = col_double(),
##   ..   `DRB%` = col_double(),
##   ..   `TRB%` = col_double(),
##   ..   `STL%` = col_double(),
##   ..   `BLK%` = col_double(),
##   ..   ORtg = col_double(),
##   ..   `%Pos` = col_double(),
##   ..   DRtg = col_double(),
##   ..   `2P%` = col_double(),
##   ..   `3P%` = col_double(),
##   ..   `FT%` = col_double(),
##   ..   `3PAr` = col_double(),
##   ..   FTAr = col_double()
##   .. )
```

According to this structure, there are 20,059 rows and 40 columns. Each row represents a player of that
particular team's season. Each column represents one of the 40 variables in this data set.

```
summary(nba_set)
```

```
##   player_id          name_common          year_id          age
## Length:20059        Length:20059        Min.   :1977    Min.   :18.00
## Class :character    Class :character    1st Qu.:1991    1st Qu.:24.00
## Mode  :character    Mode  :character    Median :2002    Median :26.00
##                                         Mean   :2001    Mean   :26.69
##                                         3rd Qu.:2012    3rd Qu.:29.00
##                                         Max.   :2020    Max.   :44.00
##
##      pos              team_id              tmRtg           franch_id
## Length:20059        Length:20059        Min.   :-15.2000    Length:20059
## Class :character    Class :character    1st Qu.: -3.6000    Class :character
## Mode  :character    Mode  :character    Median : -0.1000    Mode  :character
##                                         Mean   : -0.2147
##                                         3rd Qu.:  3.4000
```

```
##                                         Max.    : 13.4000
##
##        G               Min              MP%             MPG
##  Min.   : 1.00   Min.   :    0   Min.   : 0.00   Min.   : 0.00
##  1st Qu.:25.00   1st Qu.:  304   1st Qu.: 7.80   1st Qu.:11.60
##  Median :54.00   Median :  991   Median :25.70   Median :19.50
##  Mean   :49.13   Mean   : 1162   Mean   :29.91   Mean   :20.13
##  3rd Qu.:75.00   3rd Qu.: 1901   3rd Qu.:49.40   3rd Qu.:28.70
##  Max.   :82.00   Max.   : 3638   Max.   :92.00   Max.   :44.50
##
##       P/36            TS%             A/36             R/36
##  Min.   :  0.00   Min.   :  0.00   Min.   : 0.000   Min.   : 0.000
##  1st Qu.: 11.00   1st Qu.: 47.50   1st Qu.: 1.600   1st Qu.: 4.000
##  Median : 14.00   Median : 51.90   Median : 2.600   Median : 6.100
##  Mean   : 14.23   Mean   : 50.65   Mean   : 3.223   Mean   : 6.684
##  3rd Qu.: 17.30   3rd Qu.: 55.50   3rd Qu.: 4.300   3rd Qu.: 9.000
##  Max.   :119.80   Max.   :150.00   Max.   :36.600   Max.   :79.600
##  NA's   :5        NA's   :95       NA's   :5        NA's   :5
##      SB/36           TO/36           Raptor O         Raptor D
##  Min.   : 0.000   Min.   : 0.000   Min.   :-82.200   Min.   :-64.8000
##  1st Qu.: 1.400   1st Qu.: 1.700   1st Qu.: -3.000   1st Qu.: -1.5000
##  Median : 1.800   Median : 2.300   Median : -1.200   Median : -0.5000
##  Mean   : 2.012   Mean   : 2.399   Mean   : -1.458   Mean   : -0.4234
##  3rd Qu.: 2.500   3rd Qu.: 2.900   3rd Qu.:  0.500   3rd Qu.:  0.7000
##  Max.   :40.800   Max.   :39.900   Max.   : 53.200   Max.   : 62.5000
##  NA's   :5        NA's   :5
##     Raptor+/-        Raptor WAR        PIE%             AWS%
##  Min.   :-103.100   Min.   :-7.400   Min.   :-68.600   Min.   :-150.900
##  1st Qu.:  -3.800   1st Qu.:-0.100   1st Qu.:  5.700   1st Qu.:   4.900
##  Median :  -1.500   Median : 0.400   Median :  8.200   Median :   8.200
##  Mean   :  -1.882   Mean   : 1.629   Mean   :  7.921   Mean   :   7.469
##  3rd Qu.:   0.600   3rd Qu.: 2.500   3rd Qu.: 10.600   3rd Qu.:  11.000
##  Max.   :  72.600   Max.   :24.400   Max.   : 82.800   Max.   : 134.900
##                                      NA's   :5        NA's   :5
##      USG%            AST%             TOV%             ORB%
##  Min.   :  0.00   Min.   :  0.0   Min.   :  0.00   Min.   :  0.000
##  1st Qu.: 15.30   1st Qu.:  6.4   1st Qu.: 11.20   1st Qu.:  2.500
##  Median : 18.60   Median : 10.4   Median : 14.00   Median :  5.100
##  Mean   : 18.92   Mean   : 13.1   Mean   : 14.84   Mean   :  6.086
##  3rd Qu.: 22.20   3rd Qu.: 17.7   3rd Qu.: 17.40   3rd Qu.:  8.900
##  Max.   :163.00   Max.   :182.3   Max.   :100.00   Max.   :130.500
##                                   NA's   :25
##      DRB%            TRB%             STL%             BLK%
##  Min.   :  0.00   Min.   :  0.000   Min.   : 0.00   Min.   : 0.000
##  1st Qu.:  8.90   1st Qu.:  6.000   1st Qu.: 1.10   1st Qu.: 0.400
##  Median : 12.90   Median :  9.200   Median : 1.50   Median : 0.900
##  Mean   : 13.87   Mean   :  9.975   Mean   : 1.64   Mean   : 1.463
##  3rd Qu.: 18.20   3rd Qu.: 13.400   3rd Qu.: 2.10   3rd Qu.: 2.000
##  Max.   :231.80   Max.   :114.900   Max.   :24.90   Max.   :77.800
##
##      ORtg            %Pos             DRtg             2P%
##  Min.   :  0.0   Min.   :  0.00   Min.   :  0.0   Min.   :  0.0
##  1st Qu.: 96.0   1st Qu.: 15.60   1st Qu.:103.7   1st Qu.: 42.6
##  Median :103.8   Median : 18.60   Median :107.0   Median : 47.0
```

```
##  Mean    :101.5    Mean    : 18.94   Mean    :106.7    Mean    : 46.1
##  3rd Qu.:110.3    3rd Qu.: 21.90   3rd Qu.:110.0    3rd Qu.: 50.7
##  Max.    :300.0    Max.    :187.80   Max.    :125.0    Max.    :100.0
##  NA's    :19       NA's    :14                         NA's    :165
##      3P%               FT%              3PAr              FTAr
##  Min.   :  0.00    Min.   :  0.00   Min.   :  0.00   Min.   :  0.0
##  1st Qu.: 12.50    1st Qu.: 66.50   1st Qu.:  0.40   1st Qu.: 19.4
##  Median : 30.00    Median : 75.00   Median :  6.70   Median : 28.2
##  Mean   : 25.52    Mean   : 72.45   Mean   : 16.89   Mean   : 31.5
##  3rd Qu.: 36.60    3rd Qu.: 81.50   3rd Qu.: 30.82   3rd Qu.: 38.9
##  Max.   :100.00    Max.   :100.00   Max.   :100.00   Max.   :600.0
##  NA's   :4342      NA's   :877      NA's   :107      NA's   :107
```

In this summary, there are a mix of basic and advanced statistics for each player. We found some statistics interesting to point out. One is that the defensive rating (DRtg) is relatively consistent across the 1st quarter through 3rd quarters. Another statistic to point is out that the mean for three-point field goal percentage (3P%) while accounting for 4342 missing values is 25.52%. It doesn't seem comparable to the NBA nowadays as we see more prolific three-point shooters.

## Determining which NBA Position on Average Scores the Most Points per Decade

Our first research question is to find out which position scores the most points on average from the decades 1980-2020. We did this by subsetting the data set where we only consider the five main positions (Point Guard (PG), Shooting Guard(SG), Small Forward(SF), Power Forward (PF), and Center(C)). This is to only consider the players who purely play on one of these positions as some have played multiple positions.
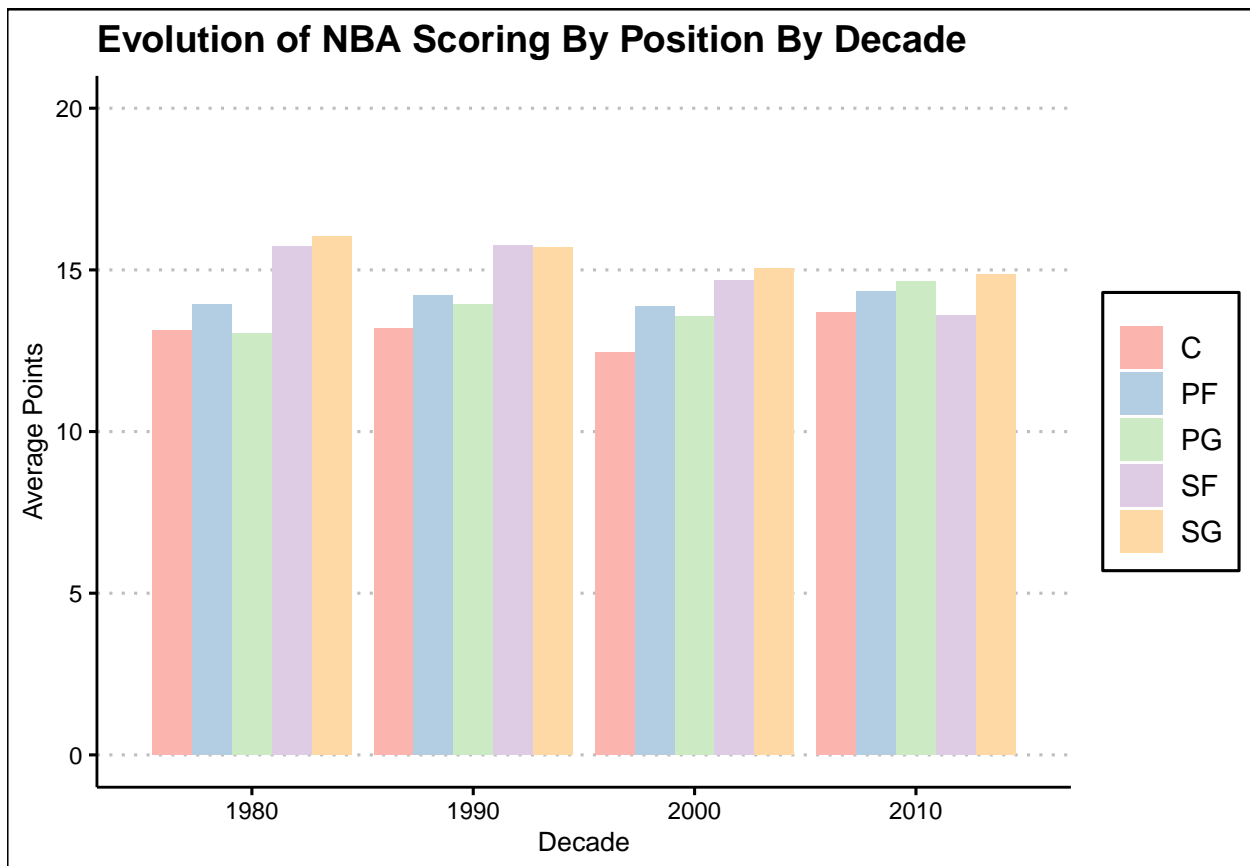
```r
#Akshay
nba_set %>%
  select(year_id, pos, 'P/36') %>%
  subset(pos %in% c("PG", "SG", "SF", "PF", "C")) %>%
  mutate(decade = floor(year_id/10) * 10) %>%
  group_by(decade, pos) %>%
  summarise(avg_pos_pts = mean('P/36', na.rm = TRUE)) %>%

  ggplot(aes(decade, avg_pos_pts)) +

  geom_bar(aes(fill = pos), position = "dodge", stat = "identity") +
  labs(title = "Evolution of NBA Scoring By Position By Decade", x = "Decade", y = "Average Points") +
    xlim(1975, 2015) +
    ylim(0,20) +
    theme_clean() +
    theme(legend.title = element_blank()) +
    scale_fill_brewer(palette = "Pastel1")
```

```
## 'summarise()' regrouping output by 'decade' (override with '.groups' argument)
```

```
## Warning: Removed 10 rows containing missing values (geom_bar).
```

## Evolution of NBA Scoring By Position By Decade



Based on this graph, we can see that there were no instances of drastic change throughout the course of this time period. The scoring average for all positions were always between the range of about 12.5 to 16 points. One trend that can be noticed is that shooting guards and small forwards scored less points on average as each decade passed by, while point guards scored more points. Another trend that is shown is that the power forward position seemed to be the most consistent throughout the four decades, as the scoring average for this position was around 14 points for all four decades. Furthermore, the 2010s is the decade in with the smallest scoring gaps among the five positions, in comparison to the other three decades.

To answer the question, there was not much evolution in regards to the scoring averages of each of the five positions. Although the scoring average for shooting guards slightly decreased in every decade, this position held the highest scoring average for every decade (with the exception of the 1990s, where it was edged out slightly by small forwards). On the other hand, centers usually scored the least amounts of points during each decade.

## Evolution of the Three-Point Shot

Our second question for this data set is how has the three-point shooting evolved over time. There are two ways we can go about answering this. One is to group the teams by division (Atlantic, Central, Southeast, Northwest, Pacific, and Southwest). Then, we can measure the average rate of three-point of attempts for each division. We can facet the plot to see how the three point usage has evolved for each division.

```
#Jerry
nba_set %>%
  mutate(Division = case_when(franch_id %in% c("MIL", "IND", "CHI", "CLE", "DET") ~ "CENTRAL",
                              franch_id %in% c("MIA", "ORL", "CHA", "WAS", "ATL") ~ "SOUTHEAST",
                              franch_id %in% c("DEN", "OKC", "UTA", "POR", "MIN") ~ "NORTHWEST",
```

```
                               franch_id %in% c("LAL", "LAC", "GSW", "SAC", "PHO") ~ "PACIFIC",
                               franch_id %in% c("HOU", "DAL", "SAS", "MEM", "NOH") ~ "SOUTHWEST",
                               franch_id %in% c("NYK", "TOR", "BOS", "NJN", "PHI") ~ "ATLANTIC")) %>%

  group_by(Division, year_id) %>%
  summarize(avg_3Usage = mean('3PAr', na.rm = TRUE)) %>%

  ggplot() +
  geom_line(aes(x = year_id, y = avg_3Usage, color = Division)) +
  labs(title = "Average Three-Point Usage by NBA Division",
       x = "Year",
       y = "Three-Point Usage (Percentage)") +
  facet_wrap(~ Division) +
  theme_clean() +
  theme(legend.position = "none")
```
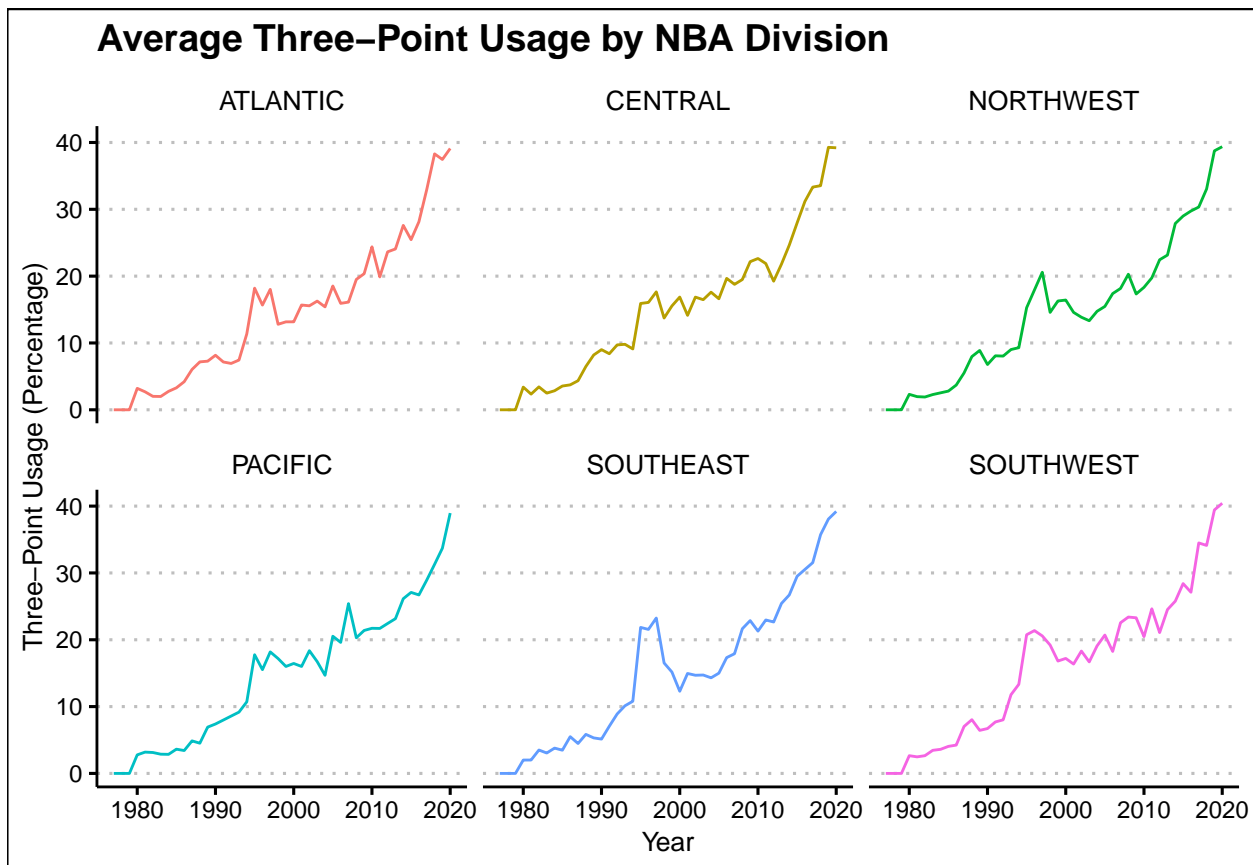
## `summarise()` regrouping output by 'Division' (override with `.groups` argument)



The graph above depicts the average three point usage for each of the six NBA divisions from 1977-2020. It is clear that for all divisions, there is a steep increase in three-point attempts since the introduction of the three-point line in 1979. We see a sizable increased rate in three-point attempts from 20% to up to 40% at around early to mid-2010's. This makes sense as the floor is more spaced and we have players like Stephen Curry, Klay Thompson, and James Harden, who have been some of the contributors of utilizing the three-point shot more often.

One thing to point out is that all divisions experienced dips in three point usage in the late 1990s to early 2000s. This is probably due to the more stifling defenses along the perimeter, which will discourage teams from attempting three-pointers. Additionally, centers and forwards were the focal point of offenses, whereas guards run the offense and feeding the ball to their centers and forwards. Since then, all divisions experienced increases in three-point usage, but the Central, Northwest, and Southeast divisions experienced much steeper increases in the mid 2010s, which was when three-point shooting really became mainstream.

Our second method is to create a heat map that shows the three-point usage for every franchise. Note that some franchises were founded later than others, which is why there are gaps within the heat map.

```
#Uziel
nba_set %>%
  group_by(franch_id, year_id) %>%
summarise('%Attempts' = mean('3PAr', na.rm = TRUE)) %>%
  ggplot(aes(year_id, franch_id)) +
  geom_tile(aes(fill = '%Attempts')) +
  labs(title = "Evolution of the 3 Point Shot by Franchise per Year",
       x = "Year",
       y = "Franchise",
       fill = "Attempt Percentage") +
  theme_clean()
```

## 'summarise()' regrouping output by 'franch_id' (override with '.groups' argument)
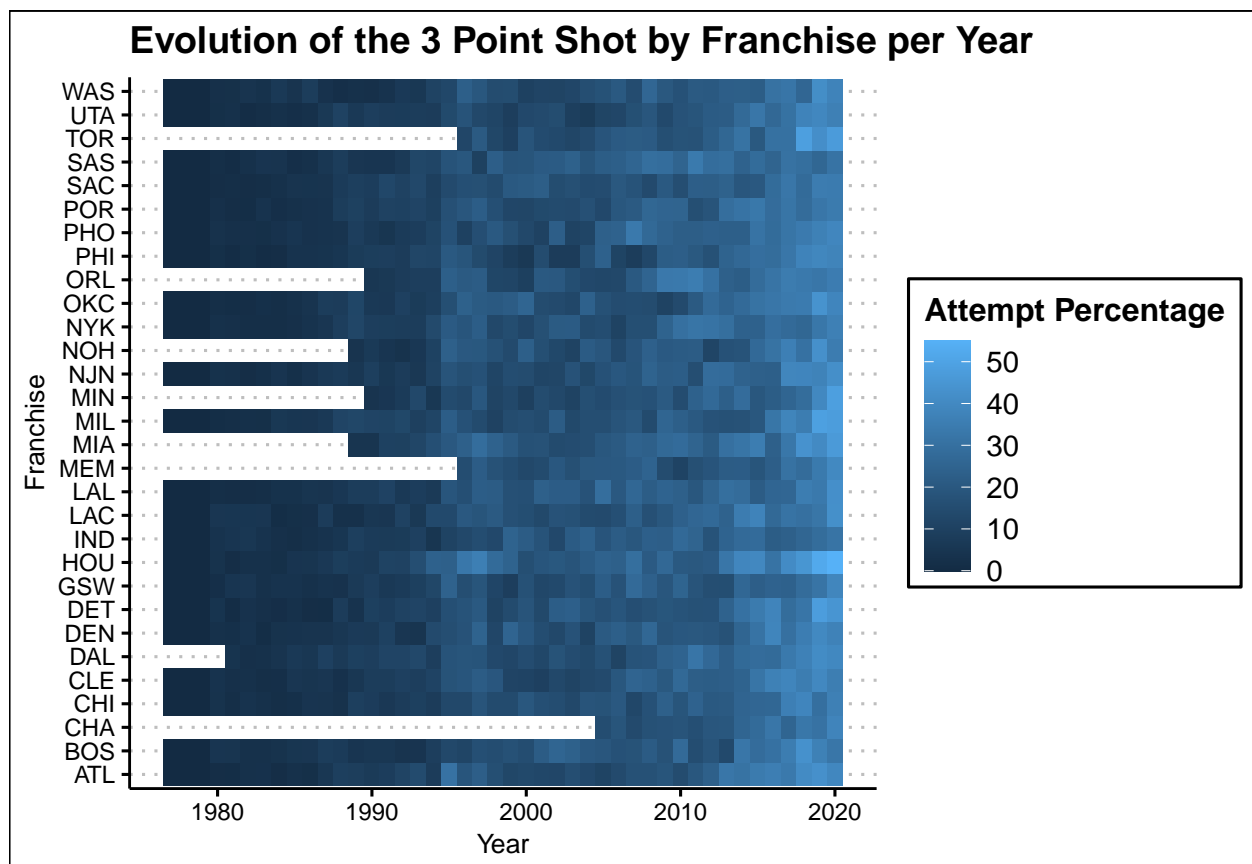


Figure Interpretation: From this figure, we can see that 3 point field goal attempts have increased drastically since the NBA introduced the three-point line. Notable franchises that rely heavily on these attempts today include HOU, MIN, and TOR.

**So, how has 3 point usage evolved over time?** In the beginning, the 3 point shot was virtually nonexistent with a field goal attempt percentage in the low single digits. Over time however, teams and coaches began to realize the efficiency of the shot and, as visualized in both graphs, 3 point shots increased significantly as a result across the board. Today, 3 point shots are a major part of basketball with some teams having nearly as much as 50% of their shots being 3 point attempts.