

# Final Report for MA678 Project

Daniel(Chen) Xu BUID:U49903384

12/10/2021

## **Abstract:**

The National Basketball Association (NBA) is a professional basketball league in North America. The league is composed of 30 teams (29 in the United States and 1 in Canada) and is one of the four major professional sports leagues in the United States and Canada. It is the premier men's professional basketball league in the world. The league was founded in New York City on June 6, 1946, as the Basketball Association of America (BAA). It changed its name to the National Basketball Association on August 3, 1949, after merging with the competing National Basketball League (NBL). The NBA's regular season runs from October to April, with each team playing 82 games. As the premier professional basketball league in the world, the change in style of basketball in the NBA led to the change in style of modern basketball. Here rises problems: Which statistical data of basketball best captures the changing style of modern basketball and how these types of data determine the injury situation for players of each season in today's style of basketball? To address those problems. I use Exploratory Data Analysis to find the factors related to the change of play style and use some statistics to build a multilevel model. However, the model shows that the variables all have a slight impact on the injury situation for each player and is slightly different between seasons as the style of basketball changes as time passes. This report consisted of 6 main parts: Abstract, Introduction, Method, Result and Discussion.

## **Introduction:**

The game of basketball has worldwide appeal. It requires speed, athleticism, skill and the ability to stay calm in the most hectic moments of the game. Basketball has undergone many changes since James Naismith invented the game to give his students something to do when the cold weather prevented them from playing sports outside. As the representative basketball league in the world, in order to gradually improve the ornamental and competitive basketball. The NBA has gone through several landmark rule changes. These rule changes have directly or indirectly affected the development and style of basketball. For example, the establishment of the 3-point line, the establishment of the defensive 3-second violation, the expansion of the 3-second zone, etc. From my point of view, two changes contributed to the style of basketball now, the first one is setting the hand check as a blocking foul outside the 3-point line and the second is the establishment of the 3-point line. The hand check is a defensive maneuver that is made on the ball-handler that a defender would extend an arm and use his hand to initiate contact with the ball-handler. Without defenders using their hands to initiate contact outside the 3-point line, players nowadays have more space to take a relatively more efficient offensive choice: shoot 3-point. I will present how the average number of 3 point attempts per game change from the Season 2010 - 2011 to the Season 2019 - 2020 in my EDA part. When it comes to professional sports, injuries are a perennial problem. For an NBA player, the degree of injury in one season directly determines his performance in this season or even the whole career. In this report, I will analyze that in the current style of basketball, what factors affect the average injuries level of NBA professional players each season.

Therefore I use a multilevel model to see what and how factors may influence the average injuries level of a player in a season. Before that, I clean the data and combine some information collected from Kaggle. I will summarize my process for data cleaning in the following part.

## Method

### Data Cleaning and Processing:

The main data set is published on Kaggle: Basketball Players Stats per Season - 49 Leagues includes 11,000 players details and stats per Season from the 1999-2000 season through the 2019-2020 season. And in order to analyze the injury situation for players, I also found a data set on Kaggle: NBA Injuries from 2010-2020 includes detail on every injury in the NBA from the beginning of the 2010-2011 season through the end of the 2019-2020 season.

To clean the stats dataset: Firstly, I subset the original dataset by only selecting players who played for NBA from the season 2010-2011 to the season 2019-2020; Secondly, I created a new variable age to store every player's age at that time because I think age is one of important factors that have effect on the average injuries level for players and it can be a factor to distinguish same player from different seasons; Thirdly, I calculate average stats for every players for each season and store them in a new dataset because from my perspective, using average data as factors can more intuitively show how the players' court performance will change with the change of basketball style. Using aggregate data is too general, and if these averages are a significant factor in the season's injury levels, then more detailed adjustments can be made to players' performance to prevent potential injuries; Finally, I separate my dataset into Regular Season and Playoffs and paste season and players name to a new variable join\_c to prepare for joining datasets. I will only focus on the regular season.

To clean the injuries dataset: Firstly, I classify each injury by body part; Secondly, I standardize dates and summarize all the dates into 10 seasons; Thirdly, paste season and players name columns to a new variable join\_c to prepare for the joining of datasets.

To wrangle datasets for further analysis: After I cleaned two datasets, I left join them together by the variable: join\_c and select the variables that I need for further analysis. And I also assigned levels for injuries according to the standard estimated recovery time and calculated average injury level for each player by generate the mean levels for injuries each season. Below is the table for injury level corresponding to each injury body part:

| Injury body part | Injury level | Injury body part | Injury level |
|------------------|--------------|------------------|--------------|
| None             | 1            | Torso            | 11           |
| Rest             | 2            | Back             | 12           |
| Illness          | 3            | Leg              | 13           |
| Finger           | 4            | Foot             | 14           |
| Hand             | 5            | Toe              | 15           |
| Face             | 6            | Ankle            | 16           |
| Arm              | 7            | Groin            | 17           |
| Shoulder         | 8            | Hip              | 18           |
| Neck             | 9            | Knee             | 19           |
| Head             | 10           | Achilles         | 20           |

Then, the final version of my cleaned dataset has 2399 observations and 30 variables. I will only choose the variables I listed above for my further analysis.

### Exploratory Data Analysis:

For the first step of my exploratory data analysis, I present the change of number of 3 point attempt per game for each players from the season 2010-2011 to the season 2019-2020. In order to show the trend of change more intuitively, I use violin and boxplot to show the distribution per season and use an arrow to connect each mean point to show the changing trend. We can conclude from the boxplot in Figure. 1, the number of three-point attempts is steadily increasing and there is a clear trend of growth start from the season 2014-2015. And there's a steady stream of players averaging more than 10 three-point attempts per

game in a season. For the violin plot, we can conclude that more and more players began to attempt the 3-point shot in the game and gradually made the 3-point shot a regular scoring option. Now we know that shooting more three-pointers is a game changing trend for today's basketball game. However, how this factor affect the average injuries level for a player? Let's explore more plot to select factors to build my multilevel model.

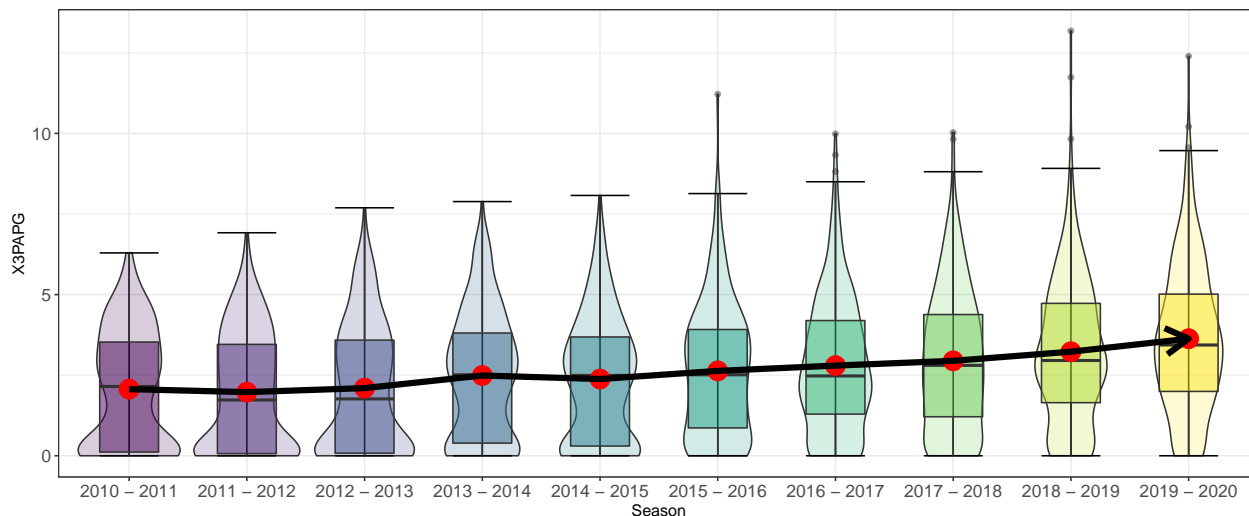


Figure 1: Number of 3 Points attempt per game distribution from season 2010-2011 to season 2019-2020

Because of the length of the report, I only picked four factors: `MPG`, `GP`, `X3PAPG` and `PFPG` to display their correlation with the average injuries level. I will use table to show the reason I choose those four factors:

| Variable Names      | Explanation                                   | Reason   |
|---------------------|---|--|
| <code>MPG</code>    | Average minutes a player play per game        | Shows a player's fatigue per game  |
| <code>GP</code>     | The number of games a player play in a season | Shows a player's fatigue per season  |
| <code>X3PAPG</code> | Average number of 3 points attempt per game   | The on-court statistics that best represent the evolution of modern basketball |
| <code>PFPG</code>   | Average number of personal fouls per game     | The contact during play time that is most likely to injure a player            |

From Figure.2, although the relationships are not linear, we can conclude that as a player's fatigue index increases, his likelihood of serious injury increases.

From Figure.3, although the relationships are not linear, in a certain range, we can conclude that the average injury index has a positive impact with the number of 3-point attempts and the number of fouls. In other words, the more three-point attempts and personal fouls a player makes within a certain range, the more likely he is to suffer injuries.

**Model Fitting:**

**Result:**

**Model build:**

Here are explanations of variables I selected for model building:

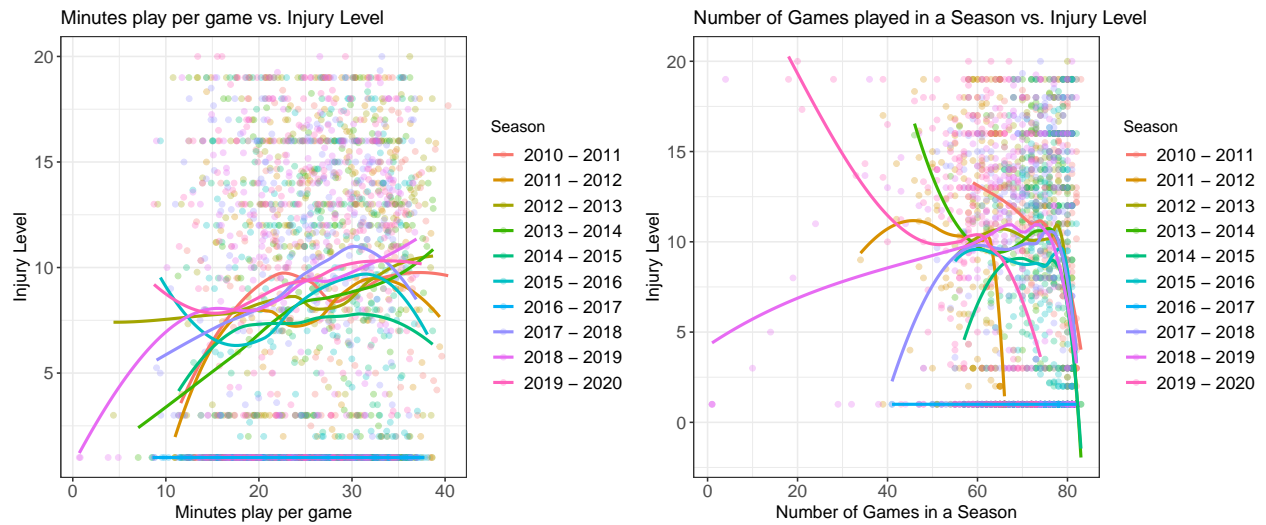


Figure 2: Correlation between the playing time per season and the average injury level

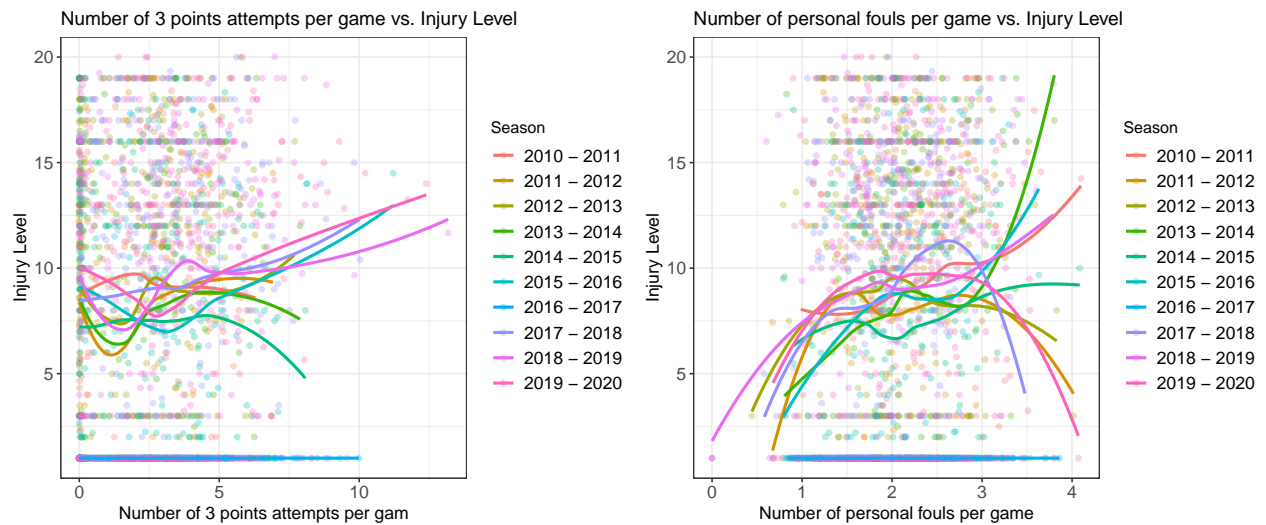


Figure 3: Correlation between some of the average stats per season and the average injury level

| Variable Names | Explanation  |
|----------------|--|
| Season         | NBA regular season year range  |
| Age            | Player's Age   |
| PER            | Player Efficiency Rating   |
| GP             | The number of games a player play in a season                                    |
| MPG            | Average minutes a player play per game   |
| PFPG           | Average number of personal fouls per game  |
| X3PAPG         | Average number of 3 points attempt per game                                      |
| height_cm      | Player's height in centimeter  |
| weight_kg      | Player's weight in kilogram  |
| injury_level   | A level summarizes the average injuries situation for each player of each season |

From the summary of my model, we can conclude that the formula:

$$InjuryLevel = 9.93 + 0.014 \cdot height_{cm} + 0.013 \cdot weight_{kg} - 0.0016 \cdot Age + 0.14 \cdot MPG + 0.041 \cdot X3PAPG - 0.16 \cdot GP + 0.38 \cdot PFPG + 0.0$$

#### Model Validation:

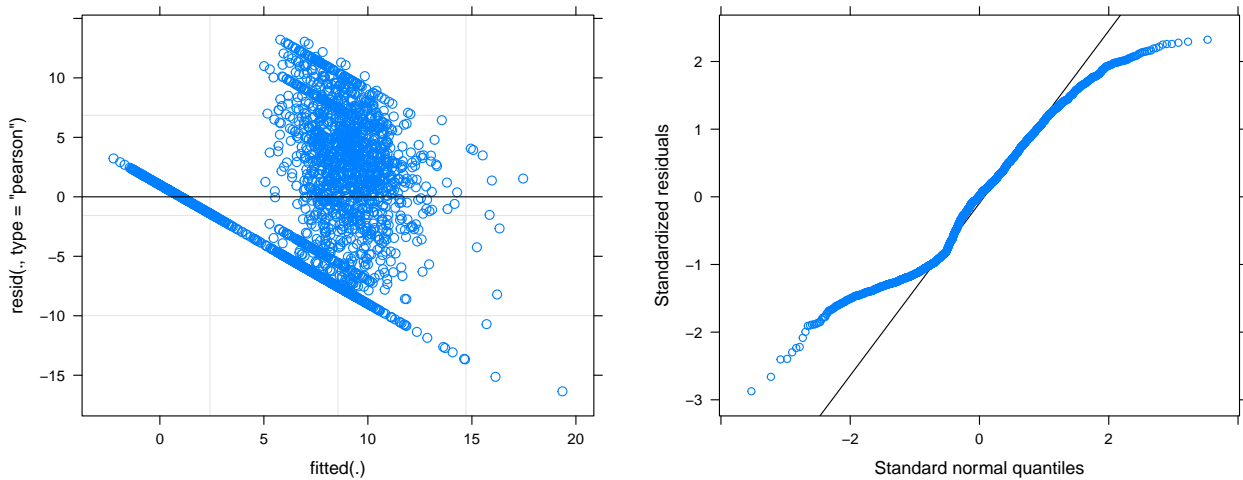


Figure 4: Residual plot and Q-Q plot.

#### Discussion:

#### Citation:

1. <https://www.sportsrec.com/358083-the-history-of-basketball-for-kids.html>
2. <https://dunkorthree.com/hand-checking-rule-nba/>
3. <https://www.kaggle.com/jacobbaruch/basketball-players-stats-per-season-49-leagues>
4. <https://www.kaggle.com/ghopkins/nba-injuries-2010-2018>

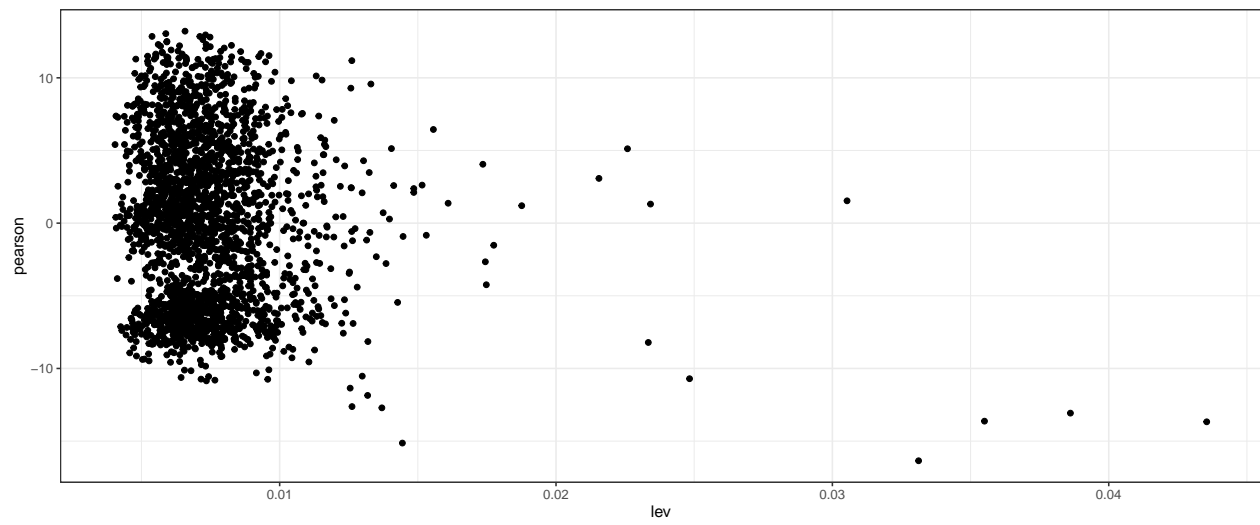


Figure 5: Residuals vs Leverage.

## Appendix

### More EDA