

Métodos estadísticos aplicados al baloncesto

Paula Moreno Blazquez

Enero 2022

Resumen

Hoy en día, el deporte es un hobby muy popular por todo el mundo. Desde pequeños, los niños practican algún tipo de deporte, especialmente aquellos que son de equipo. Eso nos lleva a querer saber más del deporte, más detalles, más información. Nos entra la curiosidad de "¿quién es el mejor jugador?", "¿Qué equipo es mejor?", o incluso intentar prever qué equipo ganará según sus resultados anteriores. Y gracias a los avances tecnológicos e informáticos, cada vez se nos facilita más poder seguir un deporte desde casa, ver la estadística de los deportistas e incluso hay plataformas o juegos que nos permiten ser, de manera virtual, managers de los clubs y, por lo tanto, nos facilitan mucha información que antes era más difícil de saber.

Eso hace que, de manera progresiva, también mejore el estudio y el análisis de cada deporte, y cada vez sea más específica para cada deporte, implementando nuevos recursos para mejorar los resultados. Pero, ¿son lo suficientemente eficaces los análisis que se realizan actualmente en Europa? ¿O dichos análisis están anticuados y requieren de una actualización?

Índice

1. Introducción	4
2. ¿Qué es GitHub?	4
2.1. Ventajas	4
2.2. ¿Qué es el control de versiones?	5
2.3. ¿Qué es Git?	5
2.3.1. Diferencias Git vs GitHub	5
3. El baloncesto	5
3.1. Historia y reglas básicas	5
3.2. Conceptos y definiciones básicas del baloncesto:	6
4. Posibles análisis realizables en el baloncesto	6
4.1. Boxscore	6
4.2. Play-by-play	7
4.3. Shot-charts	8
4.4. Graphic Stats	8
5. Estadístico Más/Menos (<i>Plus/Minus</i>)	9
5.1. Normal Plus Minus, PM	9
5.2. Adjusted Plus Minus, APM	9
5.3. Regularized Adjusted Plus/Minus, RAPM	9
5.4. Real Plus-Minus, RPM	9
6. Análisis y métodos	11
6.1. Descripción de la base de datos	11
Referencias	11
7. Anexo	12
7.1. Descripción de las variables	12
7.2. Code	13

1. Introducción

En este trabajo estudiaremos más a fondo el Baloncesto, el segundo deporte más popular de Europa (solo superado por el fútbol), y el cual tengo interés personal, ya que lo practico desde los 4 años.

Esta idea de estudio surgió del constante pensamiento de que los análisis actuales que se hacen en este deporte en Europa son bastante pobres a nivel informativo, puesto que se basan en conceptos muy básicos, y principalmente ofensivos (que vendría a ser el 50 % de un partido). Para que nos hagamos una idea, el estadístico por preferencia es el llamado *Valoración* y que se originó en 1991 (hace 30 años) y desde entonces nunca se ha modificado.

Es por eso que, considero que actualmente los análisis que se hacen de este deporte necesitan una actualización significativa para llegar a informar de todos aquellos datos que hoy en día sí se pueden recoger gracias a los avances tecnológicos, y de los cuales no se analizan por falta de dinero o porque se consideran poco relevantes.

El objetivo principal de este estudio es mejorar los análisis que se elaboran de cada partido, para poder encontrar una variable respuesta que nos diga que aportación al equipo tiene cada jugador personalmente, disminuyendo la diferencia de pesos que hay actualmente entre las aportaciones ofensivas y las aportaciones defensivas.

Este documento se estructura de la siguiente manera: a continuación, se realizará una breve explicación de los recursos informáticos que se han utilizado para realizar este estudio, seguidamente se explicará brevemente los conceptos de baloncesto que son necesarios para entender los tecnicismos del trabajo y se presentarán posibles análisis que se realizan. Finalmente, se describirá la base de datos con la que se ha trabajado y sus variables, y también se explicará en profundidad el análisis que se desarrollará en este trabajo, el Más/Menos Ajustado (*Adjusted Plus/Minus*, *APM*). En la sección de resultados presentaremos la resolución del análisis y finalmente discutiremos, en la sección de conclusiones, los resultados obtenidos.

2. ¿Qué es GitHub?

GitHub es una plataforma de alojamiento, propiedad de Microsoft, que ofrece a los desarrolladores la posibilidad de crear repositorios de código y guardarlos en la nube de forma segura, usando un sistema de control de versiones, llamado Git.

Como he comentado, facilita la organización de proyectos y permite la colaboración de varios desarrolladores en tiempo real. Es decir, nos permite centralizar el contenido del repositorio para poder colaborar con los otros miembros de nuestro grupo desde varios dispositivos.



GitHub está basada en el sistema de control de versiones distribuidas de Git, por lo que se puede contar con sus funciones y herramientas, aunque GitHub ofrece varias opciones adicionales y su interfaz es mucho más fácil de manejar, por lo que no es absolutamente necesario que las personas que lo utilizan tengan un gran conocimiento técnico.

2.1. Ventajas

Hay un gran número de razones por las que GitHub es una gran opción para el control y gestión de proyectos de código. Como por ejemplo:

- GitHub permite que alojemos proyectos en repositorios de forma gratuita
- Los repositorios son públicos por defecto. Sin embargo, GitHub te permite también alojar tus proyectos de manera privada
- Puedes crear y compartir páginas web estáticas con GitHub Pages
- Facilita compartir tus proyectos de una forma mucho más fácil y crear un portafolio
- Te permite colaborar para mejorar los proyectos de otros y a otros mejorar o aportar a los tuyos
- Ayuda reducir significativamente los errores humanos y escribir tu código más rápido con GitHub Copilot
- Te da control de versiones, una herramienta muy útil.

2.2. ¿Qué es el control de versiones?

Se le llama control de versiones a la administración de los cambios que se realizan sobre los elementos o la configuración de algún proyecto. En otras palabras, el control de versiones sirve para conocer y autorizar los cambios que hagan los colaboradores en tu proyecto, guardando información extra de qué están, incluyendo los cambios y cuándo se hicieron. Este control comienza con una versión básica del documento y luego va guardando los cambios que se hagan a lo largo del proyecto.

El control de versiones es una herramienta muy valiosa, pues con ella puedes tener acceso a las versiones anteriores de tu proyecto si es que en algún momento no llega a funcionar de forma correcta.

2.3. ¿Qué es Git?

Git es un software de control de versiones diseñado por Linus Torvalds, pensando en la eficiencia, la confiabilidad y compatibilidad del mantenimiento de versiones de aplicaciones cuando estas tienen un gran número de archivos de código fuente.

2.3.1. Diferencias Git vs GitHub

Entonces, ¿qué diferencia a Git de GitHub? La principal diferencia es que Git es un sistema que permite establecer un control de versiones, mientras que GitHub es una plataforma que ofrece un grupo de funciones que facilitan el uso de Git y la colaboración en tiempo real, así como el almacenamiento en la nube.

3. El baloncesto

3.1. Historia y reglas básicas

El baloncesto es un deporte de equipo que se originó en 1891, por James Naismith, profesor de educación física en la escuela YMCA de Springfield, Massachusetts, Estados Unidos.

James buscaba idear un deporte que sus alumnos pudieran practicar bajo techo, pues los duros inviernos en Massachusetts dificultaban la realización de ejercicio al aire libre, por lo que inventó el baloncesto utilizando unas cajas de melocotones y unos balones.

Con el paso de los años, este deporte, que empezó como actividad de colegio, ha ido evolucionando mucho, añadiendo más reglas, conceptos nuevos, límites de números de jugadores, se ha determinado tiempos de juego, las canastas tienen un valor distinto según la distancia, etc.

Actualmente, las normas más básicas de este deporte son:

- En las ligas superiores, hay un total de 4 cuartos de 10 minutos y pueden estar en pista 5 jugadores por equipo.
- No te puedes desplazar con la pelota en las manos, es obligatorio botar con una mano (si no será una infracción y conllevará la pérdida de pelota y saque de banda del equipo rival).
- Cada jugador puede realizar hasta un total de 5 faltas, que será penalizado con un saque de banda o con un tiro libre (dependerá de la situación). El jugador que realiza 5 faltas será expulsado del partido.
- El objetivo es encestar el máximo de puntos posibles, teniendo en cuenta que pueden sumar 1, 2 o 3 puntos, según la distancia.



3.2. Conceptos y definiciones básicas del baloncesto:

Para que podamos entender a que nos referimos en este trabajo, es necesario comprender unos conceptos básicos de vocabulario. Tendremos en cuenta los conceptos que se necesitan para realizar la valoración del jugador y/o del equipo que se utilizan en las estadísticas federadas.

- Puntos: Acumulación de tiros encestandos multiplicados por su valor, que cada jugador y/o equipo realiza durante el partido
- Minutos: Número de minutos que el jugador está en pista
- Falta: Acción en la que un defensor bloquea el avance de su rival sin tener control de balón o de manera no reglamentaria (empujar, agarrar...)
- Pérdidas de balón: cuando un equipo pierde el control del balón y pasa a ser del equipo rival.
- Rebotes: Recuperación de pelota después de que el tiro sea ejecutado, pero no haya encestado.
- Recuperación de balón: Cuando un equipo consigue robar el balón al equipo rival.
- Asistencia: Es un pase a un jugador que se encuentra en una posición de ventaja o que le ayuda a conseguir una canasta sin hacer ningún bote.
- Tapón: Bloqueo de un tiro en el aire.

4. Posibles análisis realizables en el baloncesto

Viendo la gran cantidad de datos que se pueden extraer de cada partido (y de cada equipo), se han ido creando análisis que recogen estos datos y los analizan para ayudarnos a identificar y desarrollar hipótesis sobre cada jugador y/o equipo.

4.1. Boxscore

El primer análisis que se hizo fue un *Box Score* (Caja de puntuación) donde se recopilaba únicamente los puntos de cada jugador según el valor de esta y las faltas realizadas. Posteriormente, se fue mejorando añadiendo conceptos como rebotes, tapones, pérdidas de balón, recuperaciones de balón... Y se añadió el estadístico (que acutalmente es por defecto) que se realiza a partir de todos estos datos: *Valoración* (en inglés PIR, *Performance Index Rating*) que engloba todo lo básico que pasa en el partido de manera individual y que, cuanto más positivo, mejor. Este estadístico se calcula utilizando la siguiente fórmula:

$$PIR = (Puntos + Rebotes + Asistencias + Robos + Tapones + FaltasRecibidas) - (TirosdeCampoFallados + TirosLibresFallados + TaponesRecibidos + Prdidias + FaltasRealizadas) \quad (1)$$

(en el Anexo 1 encontraréis la descripción de cada variable)

Posteriormente, se añadió la variable "Más/Menos" (P/M , *Plus/Minus*) que tiene que ver con la diferencia de puntos en el marcador durante el tiempo que el jugador esté en pista. Esta variable sirve para ver la contribución de los jugadores cuando están en pista. Todos los jugadores parten inicialmente con un 0, y según van entrando y saliendo de la pista, esta variable se va actualizando. Por ejemplo, los jugadores que son del quinteto inicial, empiezan con el marcador 0 - 0, y un $P/M = 0$. Si en el minuto 5, se substituye un jugador en cancha del equipo local (J1) por otro que está descansando (J2), y el marcador va 12 - 7, el P/M del J1 pasará a ser +5. Y si al cabo de 3 minutos, se sustituye el J2 por otro (J3) y el marcador ahora va 20 - 9, el P/M del J2 será +6 $((20 - 12) - (9 - 7) = 8 - 2 = +6)$.

4.3. Shot-charts

Este tipo de análisis es de los más visuales, ya que se realiza de una manera muy sencilla: se tiene como plantilla el dibujo de una pista de baloncesto de manera vectorial, y se va colocando cada tiro realizado en la posición del tiro, el jugador y si se encesta o no. De forma general se hace escribiendo el número del jugador en la posición desde donde se ejecuta el tiro, y si encesta, se hace un círculo alrededor del número.

El *Shot-Charts*, proporciona un output visual muy fácil de interpretar, ya que es se parece a un mapa de calor y, por lo tanto, podemos observar de una manera muy rápida desde que zonas de la pista es más efectivo el equipo y/o el jugador.

Una vez realizado, podemos obtener con facilidad el porcentaje de acierto del equipo y/o el jugador, o incluso, determinar el porcentaje de acierto por zonas.

De este análisis, es frecuente encontrar variantes: mapa de calor del equipo, mapa de porcentajes de aciertos por zonas de la pista...

4.4. Graphic Stats

5. Estadístico Más/Menos (*Plus/Minus*)

5.1. Normal Plus Minus, PM

5.2. Adjusted Plus Minus, APM

5.3. Regularized Adjusted Plus/Minus, RAPM

5.4. Real Plus-Minus, RPM

(explicar cada concepto y su manera de calcularlo o en base a que se calcula)

Tabla de diferencias ?

A tibble: 26,477 x 6

	season	game_code	play_number	play_type	away_player4	away_player5
	<dbl>	<dbl>	<dbl>	<chr>	<chr>	<chr>
1	2008	2	48	IN	BLU, DAVID	BLU, DAVID
2	2008	2	49	IN	BLU, DAVID	BLU, DAVID
3	2008	2	50	OUT	BLU, DAVID	BLU, DAVID
4	2008	2	51	OUT	BLU, DAVID	BLU, DAVID
5	2008	2	52	OUT	BLU, DAVID	BLU, DAVID
6	2008	2	53	CPF	BLU, DAVID	BLU, DAVID
7	2008	2	54	RPF	BLU, DAVID	BLU, DAVID
8	2008	2	55	FTM	BLU, DAVID	BLU, DAVID
9	2008	2	56	FTM	BLU, DAVID	BLU, DAVID
10	2008	2	57	TOV	BLU, DAVID	BLU, DAVID

... with 26,467 more rows

season	game_code	play_number	team_code
2008:92923	93 : 611	4 : 185	BAR : 5632
	100 : 606	295 : 185	OLY : 5555
	149 : 602	3 : 183	PAR : 5431
	96 : 599	2 : 182	PAN : 5381
	80 : 590	155 : 181	BAS : 5068
	165 : 590	203 : 181	(Other):65085
	(Other):89325	(Other):91826	NA's : 771

player_name	play_type	time_remaining	quarter
Length:92923	IN : 9401	Length:92923	Min. :1.000
Class :character	OUT : 9400	Class1:hms	1st Qu.:2.000
Mode :character	DRB : 8316	Class2:difftime	Median :3.000
	RPF : 7876	Mode :numeric	Mean :2.564
	CPF : 7389		3rd Qu.:4.000
	FTM : 5729		Max. :5.000
	(Other):44812		

points_home	points_away	play_info	seconds
Min. : 0.00	Min. : 0.00	Length:92923	Min. : 60
1st Qu.:11.00	1st Qu.:10.00	Class :character	1st Qu.: 660
Median :22.00	Median :19.00	Mode :character	Median :1260
Mean :22.89	Mean :21.07		Mean :1278
3rd Qu.:34.00	3rd Qu.:31.00		3rd Qu.:1860
Max. :75.00	Max. :69.00		Max. :2700

home_team	away_team	home	team_name
Length:92923	Length:92923	Mode :logical	Length:92923
Class :character	Class :character	FALSE:45714	Class :character
Mode :character	Mode :character	TRUE :46438	Mode :character
		NA's :771	

last_ft	and1
Mode :logical	Mode :logical
FALSE:89829	FALSE:92484
TRUE :3094	TRUE :439

[1] 2072

6. Análisis y métodos

6.1. Descripción de la base de datos

En este estudio vamos a analizar los datos de todas las jugadas registradas en partidos oficiales de la temporada 2018 de la Liga Nacional de ACB para buscar el Más/Menos Ajustado. Tenemos un total de $N = 92923$ jugadas, y 361 jugadores en 25 equipos.

Referencias

- [1] JOSEPH SILL, *Improved NBA Adjusted +/- Using Regularization and Out-of-Sample Testing*, PDF, 6 Marzo 2010.
- [2] HAPPY GIT, *Let's Git started*, url: <https://happygitwithr.com/index.html>, .

7. Anexo

7.1. Descripción de las variables

- MIN (*Minutes*): Minutos totales jugados
- PTS (*Points*): Puntos totales realizados
- 2FGA (*2-point Field Goals Attempted*): Número de tiros de 2 puntos intentadas
- 2FGM (*2-point Field Goals Made*): Número de tiros de 2 puntos anotadas
- 3FGA (*3-point Field Goals Attempted*): Número de tiros de 3 puntos ("triples") intentadas
- 3FGM (*3-point Field Goals Made*): Número de tiros de 3 puntos ("triples") anotadas
- FTA (*Free Throws Attempted*): Número de tiros libres intentados
- FTM (*Free Throws Made*): Número de tiros libres anotados

7.2. Code

Datos

```
> library(readr)
> library(dplyr)
> library(ggplot2)
> library(here)
> library(tidyverse)
> library(tidyr)
> library(dplyr)
> library(chron)           # Para CHR_to_seconds
> library(stringr)         # Para str_pad
> library(lubridate)
> library(reshape)
> library(tidyselect)
> library(corrplot)
> library(RColorBrewer)
> library("colorspace")
> library(graphics)
> library(rpart)
> library(ggplot2)
> library(car)             # VIF
> #ERROR SOLUCION
> pbp_2008 <- read_csv("pbp2008.csv")
> ## Check how many rows are affected by this
> bad_lineups <- pbp_2008 %>%
+   select(matches("_player[1-5]")) %>%
+   apply(1, function(x) max(table(x)) > 1)
> pbp_bad <- pbp_2008 %>%
+   filter(bad_lineups)
> pbp_bad %>%
+   select(season, game_code, play_number, play_type, away_player4, away_player5)

# A tibble: 26,477 x 6
   season game_code play_number play_type away_player4 away_player5
   <dbl>   <dbl>     <dbl> <chr>    <chr>         <chr>
1  2008     2         48 IN      BLU, DAVID BLU, DAVID
2  2008     2         49 IN      BLU, DAVID BLU, DAVID
3  2008     2         50 OUT     BLU, DAVID BLU, DAVID
4  2008     2         51 OUT     BLU, DAVID BLU, DAVID
5  2008     2         52 OUT     BLU, DAVID BLU, DAVID
6  2008     2         53 CPF     BLU, DAVID BLU, DAVID
7  2008     2         54 RPF     BLU, DAVID BLU, DAVID
8  2008     2         55 FTM     BLU, DAVID BLU, DAVID
9  2008     2         56 FTM     BLU, DAVID BLU, DAVID
10 2008     2         57 TOV     BLU, DAVID BLU, DAVID
# ... with 26,467 more rows

> ## Solucion
> source(here("R", "fix-lineups.R"))
> ## Function fix_lineups() only takes data from a single game,
> ## so I split the data and apply the function to each splitted data frame.
> pbp_2008_fixed <- split(pbp_2008, factor(pbp_2008$game_code)) %>%
+   map_df(fix_lineups)
> pbp_2009 <- read_csv("pbp2009.csv")
> pbp_2009_fixed <- split(pbp_2009, factor(pbp_2009$game_code)) %>%
+   map_df(fix_lineups)
> pbp_2010 <- read_csv("pbp2010.csv")
> pbp_2010_fixed <- split(pbp_2010, factor(pbp_2010$game_code)) %>%
```

```

+ map_df(fix_lineups)
> pbp_2011 <- read_csv("pbp2011.csv")
> pbp_2011_fixed <- split(pbp_2011, factor(pbp_2011$game_code)) %>%
+ map_df(fix_lineups)
> pbp_2012 <- read_csv("pbp2012.csv")
> pbp_2012_fixed <- split(pbp_2012, factor(pbp_2012$game_code)) %>%
+ map_df(fix_lineups)
> pbp_2013 <- read_csv("pbp2013.csv")
> pbp_2013_fixed <- split(pbp_2013, factor(pbp_2013$game_code)) %>%
+ map_df(fix_lineups)
> pbp_2014 <- read_csv("pbp2014.csv")
> pbp_2014_fixed <- split(pbp_2014, factor(pbp_2014$game_code)) %>%
+ map_df(fix_lineups)
> pbp_2015 <- read_csv("pbp2015.csv")
> pbp_2015_fixed <- split(pbp_2015, factor(pbp_2015$game_code)) %>%
+ map_df(fix_lineups)
> pbp_2017 <- read_csv("pbp2017.csv")
> pbp_2017_fixed <- split(pbp_2017, factor(pbp_2017$game_code)) %>%
+ map_df(fix_lineups)
> pbp_2018 <- read_csv("pbp2018.csv")
> pbp_2018_fixed <- split(pbp_2018, factor(pbp_2018$game_code)) %>%
+ map_df(fix_lineups)
> pbd_datos_fixed <- pbp_2008_fixed
> factor_names <- c('season', 'game_code', 'team_code', 'play_number', 'play_type')
> pbd_datos_fixed[,factor_names] <- lapply(pbd_datos_fixed[,factor_names] , factor)
> summary(pbd_datos_fixed[1:18])

```

season	game_code	play_number	team_code
2008:92923	93 : 611	4 : 185	BAR : 5632
	100 : 606	295 : 185	OLY : 5555
	149 : 602	3 : 183	PAR : 5431
	96 : 599	2 : 182	PAN : 5381
	80 : 590	155 : 181	BAS : 5068
	165 : 590	203 : 181	(Other):65085
	(Other):89325	(Other):91826	NA's : 771

player_name	play_type	time_remaining	quarter
Length:92923	IN : 9401	Length:92923	Min. :1.000
Class :character	OUT : 9400	Class1:hms	1st Qu.:2.000
Mode :character	DRB : 8316	Class2:difftime	Median :3.000
	RPF : 7876	Mode :numeric	Mean :2.564
	CPF : 7389		3rd Qu.:4.000
	FTM : 5729		Max. :5.000
	(Other):44812		

points_home	points_away	play_info	seconds
Min. : 0.00	Min. : 0.00	Length:92923	Min. : 60
1st Qu.:11.00	1st Qu.:10.00	Class :character	1st Qu.: 660
Median :22.00	Median :19.00	Mode :character	Median :1260
Mean :22.89	Mean :21.07		Mean :1278
3rd Qu.:34.00	3rd Qu.:31.00		3rd Qu.:1860
Max. :75.00	Max. :69.00		Max. :2700

home_team	away_team	home	team_name
Length:92923	Length:92923	Mode :logical	Length:92923
Class :character	Class :character	FALSE:45714	Class :character
Mode :character	Mode :character	TRUE :46438	Mode :character
		NA's :771	

last_ft	and1
Mode :logical	Mode :logical

```
FALSE:89829    FALSE:92484
TRUE :3094     TRUE :439
```

Categorización y edición de BBDD

```
> #CATEGORIZACION BBDD
> df <- pbd_datos_fixed
> N_df <- dim(df)[1] # Numero de jugadas
> M_df <- dim(df)[2] # Numero de variables
> df_fact <- df %>%
+   select(season, game_code, team_code, player_name, team_name) %>%
+   mutate(season = as.factor(season),
+          team_code = as.factor(team_code),
+          player_name = as.factor(player_name),
+          team_name = as.factor(team_name))
> num_teams <- length(levels(df_fact$team_name)) # Numero de equipos
> # Lista de jugadores:
> players <- df %>%
+   select(matches("player_name")) %>%
+   arrange(player_name) %>%
+   unique() %>%
+   drop_na()
> players <- players$player_name
> num_players <- length(players) # Numero de jugadores
> df_pbp <- df %>%
+   select(season, game_code, quarter, seconds, points_home, points_away, home_team, away_team, matches
```

textbfModificar variable tiempo

Para facilitar los calculos con el tiempo, se va a pasar los mm:ss a segundos.

- CHR_to_seconds: Para pasar la variable tiempo que nos llega como chracter a segundos.
- Print_MS: Que nos devolverá los segundos a formato mm:SS (se hará servir más adelante)

```
> CHR_to_seconds <- function(x){
+   a <- as.POSIXct(x, tz = '', format = "%H:%M:%S", usetz = FALSE)
+   tms <- secondss(format(a, "%H:%M:%S"))
+   s <- period_to_seconds(hms(tms))
+   return(s)
+ }
> Print_MS <- function(x){
+   t <- seconds_to_period(x)
+   sprintf('%02d:%02d:%02d', t@hour, minute(t), second(t))
+ }
```

Lista Jugadores-Equipo