



## **Análisis exploratorio y de componentes principales**

# Índice

Análisis del código.....	3
1. Librerías.....	3
2. Lectura del fichero .....	3
3. Selección de métricas .....	4
4. Análisis Exploratorio de Datos (EDA).....	4
5. Análisis de Componentes Principales (ACP) .....	5

# Análisis del código

## 1. Librerías

```
1 # -----
2 # 1. Librerías
3 # -----
4 # install.packages("FactoMineR")
5 # install.packages("GGally")
6 # install.packages("corrplot")
7 # install.packages("factoextra")
8 # install.packages("caret")
9
10 library(dplyr)
11 library(ggplot2)
12 library(GGally)
13 library(corrplot)
14 library(FactoMineR)
15 library(factoextra)
16 library(caret)
17
```

Primeramente, instalé las librerías necesarias y tras ello las importo. Las he dejado comentadas para que no se vuelvan a instalar tras ejecutar el código.

## 2. Lectura del fichero

```
8 # -----
9 # 2. Lectura y limpieza del dataset
10 # -----
11
12 getwd()
13
14 setwd("C:/Users/pepec/Documents/Master/Premaster/PM-Estadística/Modulo2")
15
16
17 # Leer el csv (ajustar el path si es necesario)
18 df <- read.csv("FBREF_players.csv", sep = ";")
19
20 # Filtro de defensas en La Liga con mínimo 684 minutos (20% de 38 partidos)
21 equipos_laliga <- c("Alavés", "Athletic Club", "Atlético Madrid", "Barcelona", "Betis",
22                   "Cádiz", "Celta Vigo", "Eibar", "Elche", "Getafe", "Granada",
23                   "Huesca", "Levante", "Osasuna", "Real Madrid", "Real Sociedad",
24                   "Sevilla", "Valencia", "Valladolid", "Villarreal")
25
26 df_defensas <- df %>%
27   filter(grepl("DF", Pos),
28          Squad %in% equipos_laliga,
29          Min >= 684)
30
```

En primer lugar, seteamos ese path como nuestro directorio actual. Ahí, contenemos todos los archivos correspondientes al módulo 2 del Premáster Curso en Estadística y Matemáticas.

A continuación, leemos los archivos del directorio previamente establecido y leemos el archivo .csv que contiene diversas métricas de los jugadores de las 5 grandes ligas.

Selecciono los equipos de la liga española y después creo el dataframe df\_defensas que filtra por la posición de defensa.

### 3. Selección de métricas

```
# -----  
# 3. Selección de métricas  
# -----  
  
metricas <- c("Int.90", "Blocks.90", "Recov.90", "AerialW.90",  
              "PassesCompleted.90", "KP.90", "PPA.90")  
  
df_metricas <- df_defensas %>%  
  select(all_of(metricas)) %>%  
  na.omit()  
  
# -----
```

Elegimos las métricas de intercepciones, bloqueos, recuperaciones, duelos aéreos ganados, pases completados, pases clave y pases al área por 90 minutos.

### 4. Análisis Exploratorio de Datos (EDA)

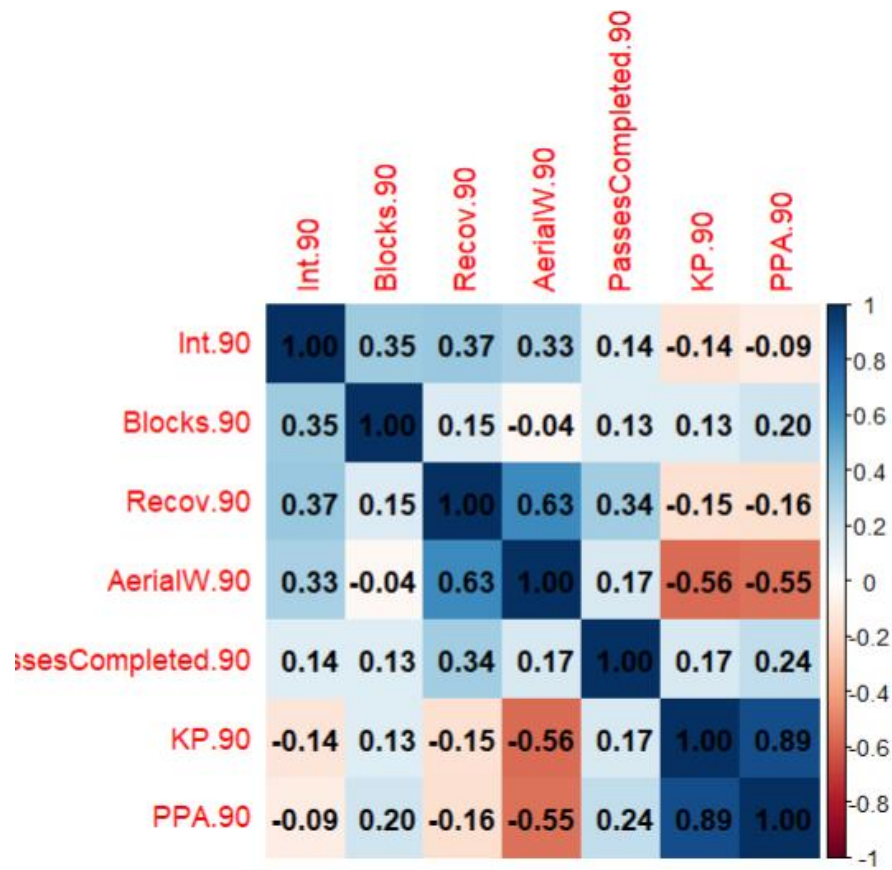
```
1  
2 # -----  
3 # 4. Análisis exploratorio de los datos  
4 # -----  
5  
6 # Estadísticas descriptivas  
7 summary(df_metricas)  
8  
9 # Pairplot para ver relaciones  
0 ggpairs(df_metricas)  
1  
2 # Matriz de correlación  
3 cor_matrix <- cor(df_metricas)  
4 corrplot(cor_matrix, method = "color", addCoef.col = "black")  
5
```

Realizamos un análisis básico con la función summary() para ver la dispersión de cada variable. A continuación, usamos ggpairs() para ver las relaciones entre métricas y corrplot() para mostrar la matriz de correlaciones.

## 5. Análisis de Componentes Principales (ACP)

```
# -----  
# 5. Análisis de componentes principales  
# -----  
  
# Escalar las variables (normalización Min-Max)  
preproc <- preProcess(df_mtricas, method = c("range"))  
df_normalizado <- predict(preproc, df_mtricas)  
  
# Aplicar PCA  
acp <- prcomp(df_normalizado, center = TRUE, scale. = TRUE)  
  
# Resumen de varianza explicada  
summary(acp)  
  
# Scree plot (varianza por componente)  
fviz_eig(acp, addlabels = TRUE, ylim = c(0, 60))  
  
# Correlación entre variables originales y componentes  
fviz_pca_var(acp, col.var = "contrib", repel = TRUE)  
  
# Puntuaciones de los jugadores sobre los dos primeros componentes  
fviz_pca_ind(acp,  
             geom.ind = "point",  
             pointshape = 21,  
             col.ind = "cos2",  
             palette = "viridis",  
             addEllipses = FALSE,  
             repel = TRUE)  
  
# Boxplot de las puntuaciones  
scores <- as.data.frame(acp$x)  
boxplot(scores, main = "Distribución de puntuaciones sobre las CPs")  
  
# -----  
# 6. Rating  
# -----  
  
df_resultado <- df_defensas %>%  
  filter(complete.cases(select(., all_of(mtricas)))) %>%  
  mutate(PC1 = scores$PC1,  
         PC2 = scores$PC2)  
  
head(df_resultado[, c("Player", "Squad", "PC1", "PC2")])
```

Normalizamos las variables utilizando Min-Max para evitar que las escalas distintas influyeran en los resultados. Después, aplicamos `prcomp()` para obtener las componentes principales. Observamos que los dos primeros explican aproximadamente el 65% de la varianza total del dataset. Representamos gráficamente las variables y los jugadores en el nuevo espacio reducido.



Distribución de puntuaciones sobre las CPs

