

MÁSTER EN BIG DATA APLICADO AL SCOUTING EN FÚTBOL

SPORTS DATA CAMPUS



INNOVATION CENTER
SEVILA FC



SPORTS DATA
CAMPUS



UCAM
UNIVERSIDAD
CATÓLICA DE MURCIA



UCAM
UNIVERSIDAD
CATÓLICA DE MURCIA



BIG DATA
International Campus

Índice

1.	Análisis del código	3
-----------	----------------------------------	----------

1. Análisis del código

```

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

```

En primer lugar, instalé las librerías necesarias e importé los paquetes requeridos para el análisis. Dejé las líneas de instalación comentadas con el fin de evitar su reinstalación ejecuciones posteriores.

```

18 # -----
19 # 2. Lectura y limpieza del dataset
20 # -----
21
22 getwd()
23
24 setwd("C:/Users/pepec/Documents/Master/Premaster/PM-Estadística/Modulo2")
25
26
27 # Leer el csv (ajustar el path si es necesario)
28 df <- read.csv("FBREF_players.csv", sep = ";")
29
30 # Filtro de defensas en La Liga con mínimo 684 minutos (20% de 38 partidos)
31 equipos_laliga <- c("Alavés", "Athletic Club", "Atlético Madrid", "Barcelona", "Betis",
32                   "Cádiz", "Celta Vigo", "Eibar", "Elche", "Getafe", "Granada",
33                   "Huesca", "Levante", "Osasuna", "Real Madrid", "Real Sociedad",
34                   "Sevilla", "Valencia", "Valladolid", "Villarreal")
35
36 df_defensas <- df %>%
37   filter(grepl("DF", Pos),
38          Squad %in% equipos_laliga,
39          Min >= 684)
40

```

A continuación, establecí el directorio de trabajo y procedí a leer los ficheros disponibles. Entre ellos, seleccioné el archivo .csv que contenía diversas métricas de rendimiento de jugadores pertenecientes a las cinco grandes ligas europeas.

Posteriormente, filtré los datos para quedarme únicamente con los equipos de LaLiga, y creé un nuevo *dataframe* denominado df_defensas, que contenía exclusivamente a los jugadores que actuaban en posiciones defensivas.

```

# -----
# 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()

# -----

```

Seleccioné las métricas de intercepciones, bloqueos, recuperaciones, duelos aéreos ganados, pases completados, pases clave y pases al área por 90 minutos, al considerarlas variables representativas del rendimiento defensivo.

```

1
2 # -----
3 # 4. Análisis exploratorio de los datos
4 # -----
5
6 # Estadísticas descriptivas
7 summary(df_metricas)
8
9 # Pairplot para ver relaciones
10 ggpairs(df_metricas)
11
12 # Matriz de correlación
13 cor_matrix <- cor(df_metricas)
14 corrplot(cor_matrix, method = "color", addCoef.col = "black")
15

```

Realicé un análisis exploratorio de datos (EDA) utilizando la función *summary()* para examinar la dispersión de cada variable. Complementé el análisis con la función *ggpairs()* para visualizar las relaciones entre métricas y *corrplot()* para representar la matriz de correlaciones y detectar posibles redundancias entre variables.

```

# -----
# 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
# -----

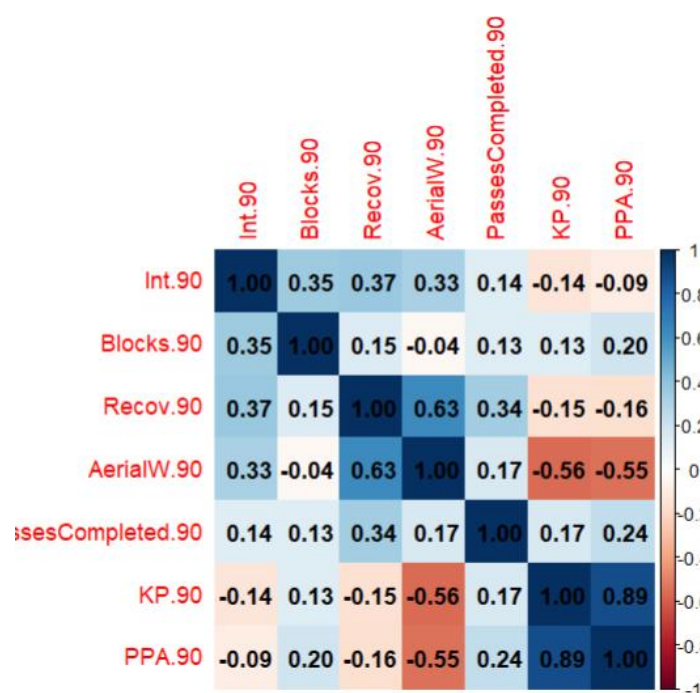
df_resultado <- df_defensas %>%
  filter(complete.cases(select(., all_of(mtricas)))) %>%
  mutate(PC1 = scores$PC1,
    PC2 = scores$PC2)

head(df_resultado[, c("Player", "Squad", "PC1", "PC2")])

```

Finalmente, normalicé todas las variables mediante la técnica Min-Max Scaling, con el objetivo de eliminar el efecto de las diferentes escalas y asegurar comparaciones consistentes entre métricas. Sobre este conjunto de datos estandarizados apliqué el Análisis de Componentes Principales (PCA) a través de la función *prcomp()*.

El resultado mostró que las dos primeras componentes principales explicaban aproximadamente el 65 % de la varianza total del dataset, lo que permitió reducir la dimensionalidad sin pérdida significativa de información. Por último, representé de forma gráfica tanto las variables como los jugadores en el nuevo espacio reducido, facilitando la interpretación visual de los patrones y relaciones existentes entre ellos.



Distribución de puntuaciones sobre las CPs

