

MÁSTER EN BIG DATA APLICADO AL SCOUTING EN FÚTBOL



SPORTS DATA CAMPUS



INNOVATION CENTER
SEVILLA FC



SPORTS DATA
CAMPUS



UCAM
UNIVERSIDAD
CATÓLICA DE MURCIA

Índice

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

1. Análisis del código

```

1
2 # -----
3 # 1. Librerías
4 #
5 library(dplyr)
6 library(ggplot2)
7 library(GGally)
8 library(caret)
9 library(factoextra)
0
1 #
2 # 2. Lectura y limpieza del dataset
3 #
4 setwd("C:/users/pepec/Documents/Master/Premaster/PM-Estadística/Modulo3")
5
6 # Leer el csv
7 df <- read.csv("FBREF_players.csv", sep = ";")
8
9 # Filtrar jugadores defensas de LaLiga con al menos 684 minutos jugados
0 equipos_laliga <- c("Alavés", "Athletic Club", "Atlético Madrid", "Barcelona", "Betis",
1 "Cádiz", "Celta Vigo", "Eibar", "Elche", "Getafe", "Granada",
2 "Huesca", "Levante", "Osasuna", "Real Madrid", "Real Sociedad",
3 "Sevilla", "Valencia", "Valladolid", "Villarreal")
4
5 df_defensas <- df %>%
6   filter(grepl("DF", Pos),
7         Squad %in% equipos_laliga,
8         Min >= 684)
0

```

Tras importar las librerías, establecí el directorio de trabajo mediante la función `setwd()`.

A continuación, leí los archivos del directorio previamente establecido y leemos el archivo .csv que contiene diversas métricas de los jugadores de las cinco grandes ligas que hayan disputado un mínimo de 684 minutos. Seleccioné los equipos de la liga española y después creo el dataframe `df_defensas` que filtra por la posición de defensa.

El resultado es un dataframe limpio y segmentado que servirá como base para el posterior análisis exploratorio y la aplicación de técnicas de clusterización.

```

# -----
# 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()

```

Tras filtrar los jugadores por liga y posición, seleccioné un conjunto representativo de KPIs defensivos y de construcción (intercepciones, bloqueos, recuperaciones, duelos aéreos ganados, pases completados, pases clave y pases al área por 90 minutos).

```

# -----
# 4. Análisis exploratorio
# ----

summary(df_metricas)
ggpairs(df_metricas)
boxplot(scale(df_metricas), main = "Boxplot normalizado de métricas seleccionadas")

```

En esta fase realicé un análisis exploratorio de los datos con el objetivo de comprender la distribución y relación entre las métricas seleccionadas. Primero utilice la función `summary()` para obtener una visión general de la dispersión y tendencia central de cada variable, identificando posibles valores atípicos o diferencias de escala entre indicadores.

Después, apliqué `ggpairs()` para visualizar las relaciones bivariadas entre las métricas seleccionadas, lo que permite detectar patrones de correlación o redundancia entre variables.

Por último, empleé `boxplot()` sobre los datos normalizados (`scale(df_metricas)`) con el fin de comparar gráficamente la variabilidad de las distintas métricas y evaluar su comportamiento antes de aplicar técnicas de clusterización.

```

# -----
# 5. Normalización de los datos
# ----

preproc <- preprocess(df_metricas, method = c("range"))
df_normalizado <- predict(preproc, df_metricas)
# -----

```

Para evitar que las variables con diferentes escalas afectaran al análisis, apliqué una normalización Min-Max a todas las métricas mediante la función `preProcess()`. Esto permitió homogeneizar los valores en un rango común y garantizar una comparación equitativa entre variables.

```

# -----
# 6. Número de clusters
# ----

fviz_nbclust(df_normalizado, kmeans, method = "wss") +
  geom_vline(xintercept = 3, linetype = 2)

```

Para determinar el número óptimo de grupos, implementé el método del codo (Elbow Method), calculando la suma de errores cuadráticos dentro de los grupos (WSS). El punto de inflexión identificado en el gráfico correspondió a K = 3, por lo que se estableció la segmentación en tres clusters.

```

# -----
# 7. K-Means
# -----
set.seed(123)
km <- kmeans(df_normalizado, centers = 3, nstart = 25)

# -----
# 8. visualización de clusters
# -----
fviz_cluster(km, data = df_normalizado)

```

Finalmente, apliqué el algoritmo K-Means con 25 inicializaciones distintas para asegurar la convergencia del modelo y evitar resultados dependientes de la inicialización aleatoria. El proceso permitió segmentar a los jugadores en tres grupos funcionales, diferenciados según su perfil de rendimiento defensivo y capacidad de construcción de juego.

