

# PRA 1

---

## Contexto

Los datos han sido extraídos de la página oficial de la [FIFA](#) (Federación Internacional de Fútbol Asociado). Tratándose de información sobre el mundial de fútbol celebrado en Rusia en 2018 he considerado que la información suministrada por los propios organizadores del evento serían los más veraces.

## Título

"World Cup Russia 2018 stats and reviews"

## Descripción

El conjunto de datos contiene información sobre los partidos disputados en la copa del mundo de fútbol celebrada en Rusia en 2018. Cada partido es descrito mediante un conjunto de variables cuantitativas como pueden ser por ejemplo, número de goles, faltas, tiros a puerta, etc. (stats) y por un conjunto de variables que podrían considerarse como cualitativas y que consisten en las reseñas hechas por los periodistas encargados de cubrir el evento (reviews).

## Representación gráfica



## Contenido

Cada registro se corresponde con un partido. Y para cada uno de ellos se han recogido los siguientes datos:

- **id**: Identificador del partido
- **name\_home**: Nombre del equipo local
- **name\_away**: Nombre del equipo visitante
- **code\_home**: Código equipo local
- **code\_away**: Código equipo visitante
- **group**: Ronda clasificatoria
- **stadium**: Estadio
- **venue**: Localización
- **datetime**: Fecha y hora
- **headline**: Titular resumen
- **summary**: Reseña neutral
- **summary\_home**: Reseña del periodista local
- **summary\_away**: Reseña del periodista visitante
- **goals\_home**: Goles equipo local
- **goals\_away**: Goles equipo visitante
- **attempts\_home**: Lanzamientos equipo local
- **attempts\_away**: Lanzamientos equipo visitante
- **on-target\_home**: Lanzamientos entre los tres palos equipo local
- **on-target\_away**: Lanzamientos entre los tres palos equipo visitante
- **off-target\_home**: Lanzamientos fuera equipo local
- **off-target\_away**: Lanzamientos fuera equipo visitante
- **blocked\_home**: Paradas equipo local
- **blocked\_away**: Paradas equipo visitante
- **woodwork\_home**: Palos equipo local
- **woodwork\_away**: Palos equipo visitante

- **corners\_home**: Saques de esquina equipo local
- **corners\_away**: Saques de esquina equipo visitante
- **offsides\_home**: Fuera de juego equipo local
- **offsides\_away**: Fuera de juego equipo visitante
- **ball\_possession\_home**: Porcentaje de posesión equipo local
- **ball\_possession\_away**: Porcentaje de posesión equipo visitante
- **pass\_accuracy\_home**: Precisión pasa equipo local
- **pass\_accuracy\_away**: Precisión equipo visitante
- **passes\_home**: Pases equipo local
- **passes\_away**: Pases equipo visitante
- **passes\_completed\_home**: Pases completados equipo local
- **passes\_completed\_away**: Pases completado equipo visitante
- **distance\_covered\_home**: Distancia recorrida equipo local
- **distance\_covered\_away**: Distancia recorrida equipo visitante
- **balls\_recovered\_home**: Recuperaciones equipo local
- **balls\_recovered\_away**: Recuperaciones equipo visitante
- **tackles\_home**: Entradas equipo local
- **tackles\_away**: Entradas equipo visitante
- **blocks\_home**: Cortes equipo local
- **blocks\_away**: Cortes equipo visitante
- **clearances\_home**: Despejes equipo Local
- **clearances\_away**: Despejes equipo visitante
- **yellow\_cards\_home**: Tarjetas amarillas equipo local
- **yellow\_cards\_away**: Tarjetas amarillas equipo visitante
- **direct\_red\_cards\_home**: Tarjetas rojas directas equipo local
- **direct\_red\_cards\_away**: Tarjetas rojas directas equipo visitante
- **indirect\_red\_cards\_home**: Tarjetas rojas indirectas equipo local
- **indirect\_red\_cards\_away**: Tarjetas rojas indirectas equipo visitante
- **fouls\_committed\_home**: Faltas equipo local
- **fouls\_committed\_away**: Faltas equipo visitante
- **referee\_name**: Nombre del árbitro
- **referee\_contry**: País del árbitro
- **weather\_description**: Descripción del tiempo
- **weather\_temperature**: Temperatura (°C)
- **weather\_windspeed**: Velocidad del viento (km/h)
- **weather\_humidity**: Humedad (%)

## Agradecimientos

La propietaria de los datos es la FIFA. La Fédération Internationale de Football Association, más conocida por sus siglas **FIFA**, es la institución que gobierna las federaciones de fútbol en todo el planeta. Se fundó el 21 de mayo de 1904 y tiene su sede en Zúrich, Suiza. Forma parte del IFAB, organismo encargado de modificar las reglas del juego. Además, la FIFA organiza la Copa Mundial de Fútbol, los otros campeonatos del mundo en sus distintas categorías, ramas y variaciones de la disciplina, y los Torneos Olímpicos a la par del COI.

Me gustaría destacar la gran cantidad de información que la FIFA pone a disposición de los aficionados a través de su [web](#) y el orden y claridad con la que es presentada.

## Inspiración

El conjunto de datos podría ser utilizado en diferentes ámbitos. Algunos de ellos podrían ser:

- **Periodismo**: Los datos podrían ser usados para, por ejemplo, realizar un reportaje sobre el mundial que sacara a relucir los datos más curiosos/destacados.
- **Almacén de datos**: Los datos podrían ser cruzados con datos pertenecientes a otras fuentes para enriquecer un almacén de datos estadísticos.
- **Minería de datos**: En cuanto a su uso para proyectos de minería de datos podría ser interesante estudiar:
  - **Sesgos en información**: Podrían estudiarse la diferencias o similitudes entre las reseñas de cada uno de los periodistas (local/visitante) y la reseña neutral para extraer conclusiones sobre sesgos en el periodismo.
  - **Generador de reseñas**: Utilizando en conjunto las variables cuantitativas y las reseñas se podría intentar construir un sistema que generase reseñas de manera automática. Aunque seguramente se necesitarían más datos, el conjunto podría ser un punto de partida para estudiar la viabilidad del

proyecto.

- **Estilos de juego:** Podrían estudiarse los diferentes estilos de juego de las selecciones (con técnicas de clustering) y estudiar cuales son los más efectivos.

## Licencia

La licencia escogida ha sido **CC0: Public Domain License**. La razón por la que he escogido esta licencia, es que todo el trabajo ha sido realizado con el único objetivo de superar la asignatura "Tipología y ciclo de vida de los datos" por lo que si de alguna forma, alguien quiere usar los datos para cualquier finalidad, esta licencia le permitirá el uso de estos con el menor número de restricciones posibles.

## Código

El código fuente se puede encontrar dentro de la carpeta "src".

## Dataset

El dataset resultante se puede encontrar en la carpeta "csv".

## Recursos

- Subirats, L., Calvo, M. (2019). Web Scraping. Editorial UOC.
- Lawson, R. (2015). Web Scraping with Python. Packt Publishing Ltd. Chapter 2. Scraping the Data.