

Explicación código

Nicolas Andres Vera

506211068

```
import numpy as np
import pandas as pd
from pandas_profiling import ProfileReport
#El código importa las bibliotecas de NumPy y Pandas , también importa la biblioteca pandas_profiling y ProfileReport genera un informe de los datos utilizado por las bibliotecas

[ ] data = pd.read_csv("scopus.csv")
#El código carga el csv llamado scopus.csv

[ ] profile = ProfileReport(data, title="Profiling Report")
#genera el informe utilizando la librería pandas y le da un título

[ ] profile
#Se hace el llamada para ejecutar lo anterior
```

Summarize dataset: 100%  43/43 [00:06<00:00, 4.73it/s, Completed]

Generate report structure: 100%  1/1 [00:08<00:00, 8.62s/it]

Render HTML: 100%  1/1 [00:01<00:00, 1.56s/it]

Profiling Report			Overview	Variables	Interactions	Correlations	Missing values	Sample
	Number of observations	650	Numeric		2			
	Missing cells	3145	Unsupported		2			
	Missing cells (%)	20.2%						
	Duplicate rows	0						
	Duplicate rows (%)	0.0%						
	Total size in memory	122.0 KiB						
	Average record size in memory	192.2 B						

# Variables

```
[ ] import numpy as np
import pandas as pd

import requests
from bs4 import BeautifulSoup

from PIL import Image
from wordcloud import WordCloud, ImageColorGenerator
import matplotlib.pyplot as plt

import string
import nltk
from nltk.corpus import stopwords
#Basicamente importamos varias librerias que son:
#numpy y pandas para analizar los datos
#requests y BeautifulSoup para hacer solicitudes http y analizar un html
#PIL y WordCloud para crear una imagen de acuerdo a un texto
#matplotlib.pyplot para hacer graficos en python
#string y nltk procesar texto
#Todo esto se utilizara para generar la nube de palabras que veremos despues
```

```
[ ] nltk.download('stopwords')
#Se quitan palabras vacias o sin significado

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
True
```

```
[ ] punctuation=[]
for s in string.punctuation:
    punctuation.append(str(s))
sp_punctuation = ["¿", "¡", "¡¡", "!!!", "...", ":", "-", ">", "<"]

punctuation += sp_punctuation
#Se crea una lista de puntuaciones que pueden ser utiles para eliminarlos del texto porque no son relevantes
```





```
[ ] word_count={}
    for palabra in palabras:
        if palabra in word_count.keys():
            word_count[palabra][0]+=1
        else:
            word_count[palabra]=[1]
    #Se crea un diccionario wordl_count , con el bucle se verifica si la palabra esta en word_count, si esta suma y si no se agrega al diccionario
```

```
[ ]
```

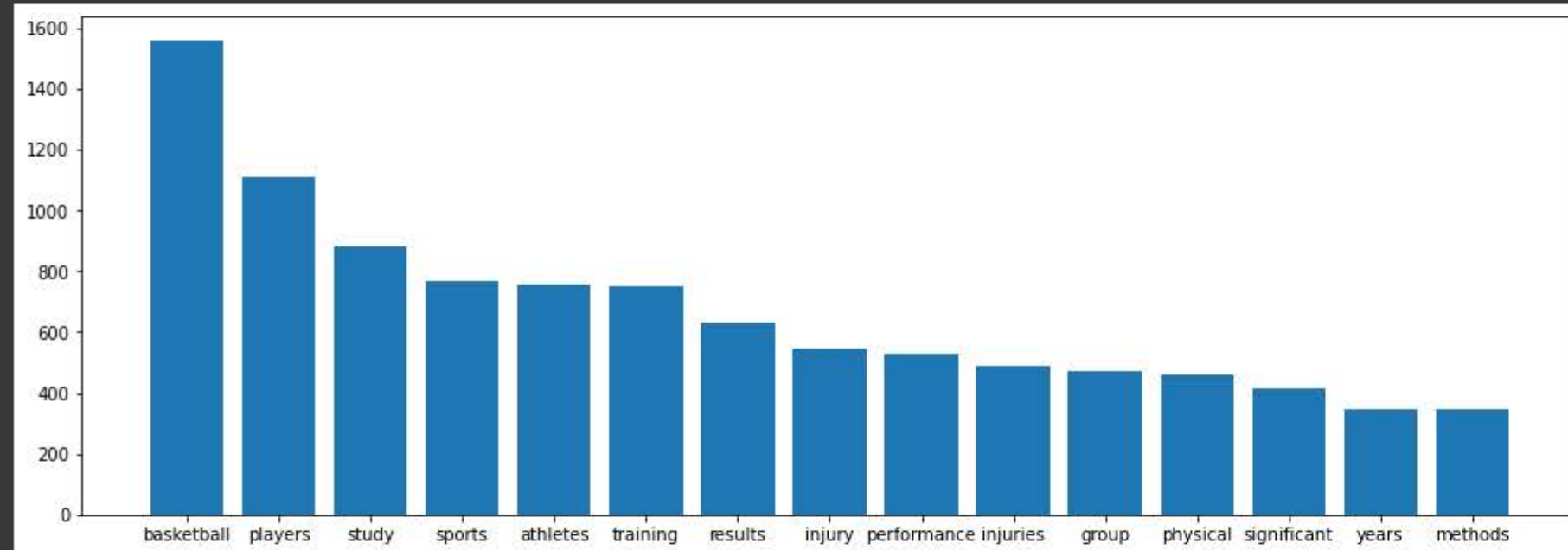
```
[ ] df = pd.DataFrame.from_dict(word_count).transpose()
df.columns=["freq"]
df.sort_values(["freq"], ascending=False, inplace=True)
df.head(10)
#, el código devuelve las primeras 10 filas del DataFrame panda ordenado, que corresponden a las palabras con mayor frecuencia
```

	freq
<b>basketball</b>	1562
<b>players</b>	1112
<b>study</b>	881
<b>sports</b>	768
<b>athletes</b>	754
<b>training</b>	748
<b>results</b>	631
<b>injury</b>	546
<b>performance</b>	530
<b>injuries</b>	486

```
[ ]
```

```
[ ] def plot_bar(data=df, top=5):
    fig = plt.figure()
    ax = fig.add_axes([0,0,2,1])
    ax.bar(x=df.iloc[:top,:].index, height = df.iloc[:top,0].values)
    plt.show()
    # Se crea una figura de Matplotlib y se agrega el eje de las barras, ademas se muestra el grafico
```

```
[ ] plot_bar(data=df, top=15)
#Llama a la funcion anterior y se hace un grafico de los 15 resultados que mas se repiten
```



```
[ ] word_cloud = WordCloud(height=800, width=800, background_color='white',max_words=150, min_font_size=5, collocation_threshold=10).generate(clean_texto)
word_cloud.to_file("n1.png")
plt.figure(figsize=(10,8))
plt.imshow(word_cloud)
plt.axis('off')
plt.tight_layout(pad=0)
plt.show()
#Se genera la nube de palabras teniendo en cuenta la biblioteca wordcloud , se dan parametros como la altura, se crea la imagen y se le dice donde guardarse , este es el resultado cuando ya limpiamos las palabras
#anteriormente
```

