T1.2: Nube de palabras

Sistemas de Big Data 26/11/23 – IES Fernando Wirtz David Fernández Reboredo

Índice

Tar	efa 1.2	3
	Exercicio Principal	
	Ampliación 1	
	Ampliación 2	

Tarefa 1.2

Exercicio Principal

Para comezar a realizar a tarefa deberemos realizar un pull sobre o repositorio proporcionado, para elo e tan fácil como facer click na icona da barra lateral esquerda como podemos observar na seguinte imaxe:



Posteriormente a realizar a seguinte acción deberemos ingresar o link que se nos proporcionó, na tarefa.

A partires de ahí xa poderemos traballar co repositorio proporcionado.

Unha vez realizado esto, deberemos a todas as librerías facer un conda install.

Deberemos facelo coas librerias seguintes:

conda install -c conda-forge wordcloud

```
- conda install numpy

- conda install PIL

- conda install wordclud

- conda install nltk

- conda install urllib3

- conda install bs4
```

Posteriormente a esto simplemente deberemos de correr o código, para unha mellor utilización agrupámolos en funcións

```
limparcodigo(texto):
    clean_texto=''
    punctuation=[]
    for s in string.punctuation:
        punctuation.append(str(s))
    sp_unctuation = ["¿", "|", """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, "", """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """, """", """, """, """", """, """, """, """", """, """, """, """, """,
```

O cal empregaremos para facer todo o filtrado de palabras

Por outro lado o xerador da nube de palabras

```
def generarNube(clean_texto):
    word_cloud = Wordcloud(height=800, width=800, background_color='white',max_words=150, min_font_size=5).generate(clean_texto)
    plt.figure(figsize=(10,8))
    plt.imshow(word_cloud)
    plt.axis('off')
    plt.tight_layout(pad=0)
    plt.show()
```

Usamos o noso ficheiro txt prefeito e realizaremos o seu parseo:

```
def parsear_txt():
    texto=''
    with open('defino.txt','r') as fichero:
        for linea in fichero:
            texto=texto+linea
            return texto+''
            fichero.close()

        vo.0s

texto=parsear_txt()
    texto=limparcodigo(texto)
    print(texto)
    generarNube(texto)
```

Finalmente como podemos observar na parte inferior correremos o código xerandonos o seguinte resultado:



Ampliación 1

Para a ampliación deberemos de facelo mediante unha URL

O seguinte código proporcionaranos a obtención do texto na URL ao obter nel os h4 da html:

```
def def urls(urlprop):
   ua = "Mozilla/5.0 (Linux; U; Android 2.2; en-us; Nexus One Build/FRF91) AppleWebKit/533.1 (KHTML, like Gecko) Version/4.0 Mobile Safari/
   h = {"User-Agent": ua}
   resultados=
   http = urllib3.PoolManager()
   medioDigital= urlprop
   r = http.request('GET', medioDigital, fields=None, headers=h)
   sopa = BeautifulSoup(r.data, "html.parser")
   web_solotexto = sopa.get_text()
    for linea in web_solotexto.split('\n'):
       aux=linea.strip()
        if aux and len(aux) > 50:
           salida += aux + '\n'
   print (salida)
    titulares = sopa.find_all('h4')
    for titular in titulares:
       resultados=titular.get_text().strip()
    return resultados + salida +
```

```
url=def_urls(url)
l reactor JET ha completado con éxito sus pruebas finales con deuterio y tritio. Es un hito crucial para la fusión nuclear
l reactor JET ha completado con éxito sus pruebas finales con deuterio y tritio. Es un hito crucial para la fusión nuclear
os técnicos de ITER iniciarán las pruebas de alta potencia con deuterio y tritio en 2035
as pruebas con plasma ionizado llevadas a cabo en JET son cruciales para que ITER y DEMO salgan bien
i todo sigue su curso como está previsto ITER (International Thermonuclear Experimental Reactor), el reactor experimental de fusión nuclear
o obstante, antes de iniciar las pruebas de alta potencia con el combustible final ITER deberá superar otros tests que también son cruciales
l reactor experimental JT-60SA reside en Naka, una pequeña ciudad no muy alejada de Tokio (Japón). Su propósito muy a grandes rasgos es llev
ET ha probado la fusión con el mismo combustible que utilizará ITER
l último gran hito de JET vio la luz pública el 9 de febrero de 2022. Ese día los científicos que lo operan anunciaron que habían logrado ba
a mejor arma de EEUU frente a China es una máquina europea: el equipo de litografía 'EUV-High NA' de ASML
os técnicos de JET han completado el programa DTE3, que es la tercera y última campaña de pruebas con plasma ionizado que contiene núcleos d
es que hace apenas unas horas los responsables de este reactor experimental de fusión nuclear han anunciado que sus científicos han complet.
lay muchas razones por las que esta campaña de pruebas con el mismo combustible que usará ITER tiene muchísima importancia. Una de las más re
demás han permitido refinar la administración del tritio, que es un isótopo radiactivo, y también han ayudado a los técnicos a comprender con
n Xataka: Estos son los plazos que maneja ITER actualmente para demostrar la viabilidad de la fusión nuclear
treaming Análisis Energía Espacio Móviles Xataka Movilidad Apple Samsung Inteligencia artificial China Empleo Windows 11 Ver más temas
atakaXataka MóvilXataka AndroidXataka Smart HomeApplesferaGenbetaMundo Xiaomi
ecotec Conga 12090, análisis: el mejor fregado que he probado en un robot aspirador merecía más
i tienes fotos viejas o borrosas, esta nueva herramienta IA de expertos españoles hace un upscaling milagroso (y hasta creativo)
```

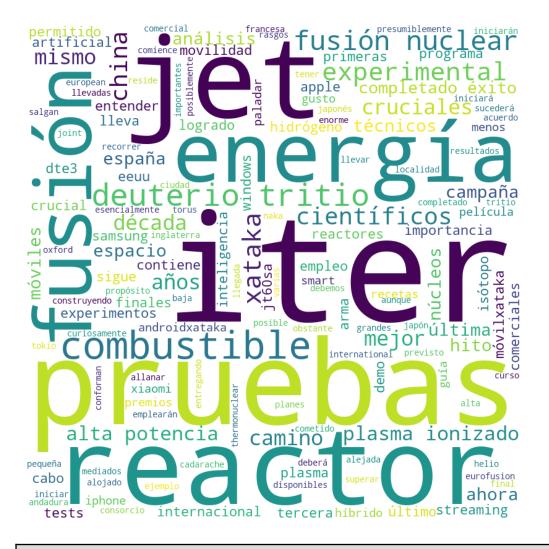
O resultado de este mostraranos unha serie de palabras separadas por un espacio

Posteriormente poderemos realizar as operacións de limpeza e creación da nube de palabras.

```
url=limparcodigo(url)
print(url)
```

generarNube(url)

Nos dará como resultado la siguiente nube de palabras



Ampliación 2

Para esta segunda parte deberemos crear un json coma o seguinte:

A continuación debemos de recorrer el json proporcionado sacar los h1 en los cuales se encuentra el titular del periódico (esto depende del periódico ya que algunos lo hacen en h2 otros en h1 otros en title...)

```
from bs4 import BeautifulSoup
import urllib3
import json

text = ''

ua = "Mozilla/5.0 (Linux; U; Android 2.2; en-us; Nexus One Build/FRF91) AppleWebKit/533.1 (KHTML, like Gecko) Version/4.0 Mobile Safari/533.1"
headers = {"User-Agent": ua}

with open("json.json", "r") as f:
    json_data = json.load(f)

urls = json_data["lista"]
for url_info in urls:
    url = url_info["link"]
    http = urllib3.PoolManager()
    request = http.request('GET', url, fields=None, headers=headers)
    soup = BeautifulSoup(request.data, 'html.parser')
    text += "".join([tag.text for tag in soup.find_all(['h1', 'h2'])])
```

Finalmente y al igual con la Ampliación 1 deberemos llamar desde el texto a la operación de limpiado y de generación de la nube de palabras. El resultado final es el siguiente:

