

# Análisis Data

Braian Felipe Ramírez Ortiz  
Fundación Universitaria Konrad Lorenz

## I. INTRODUCCIÓN

Dado el taller 3 y 4 de estructura de datos, se propone un texto para su análisis. Se pide desarrollar dos algoritmos: uno que recorra el documento reconociendo cada una de las palabras del mismo e identificando cuántas veces se repiten las palabras en él, imprimiendo como resultado una lista con cada palabra del documento mostrando la cantidad de veces que esta se encuentra dentro de este.

El segunda parte del algoritmo analiza el documento reconociendo la posición de cada una de las palabras dentro de este y retornando como resultado una nueva lista con las palabras del texto y su respectivos indices de ubicación.

## II. RESULTADOS

Para analizar el comportamiento del algoritmo se emplearon las bibliotecas time y psutil para estudiar la complejidad espacial y temporal, se ejecutó el algoritmo tres veces obteniendo varios datos de tiempo y memoria para asi generar un contraste de datos, se tienen encuestas las especificaciones tecnicas de la maquina en la cual se desarrollo el código.

Maquina	Sistema Operativo	Procesador	Memoria ram	Disco SSD
Lenovo V14 G2 ALC	Windows 11	Ryzen 7	40 GB	512GB

Figura 1.

```

C:\Users\User\Python\proyecto\datos\analisis\datos\python.exe C:\Users\User\Python\proyecto\datos\analisis\datos.py
estas son todas las palabras que contiene el documento con el numero de veces que aparecen en el mismo:
{'la': 20, 'programación': 9, 'en': 24, 'Python': 2, 'en': 20, 'clave': 3, 'para': 23, 'el': 22, 'trabajo': 2, 'con': 2, 'datos': 11, 'los': 9, 'programa':
estas son las cinco palabras mas usadas en el documento:
{'la': 20, 'la': 20, 'en': 24, 'en': 24, 'para': 23}
lista de las palabras con los indices de cada palabra:
{'la': [0, 2, 5, 7, 10, 11, 15, 16, 17, 18, 20, 27, 28, 30, 31, 35, 37, 39, 41, 45, 47, 48, 50, 51, 54, 55, 45, 61, 63, 65], 'programación': [8, 6, 6,
el tiempo de ejecución del algoritmo es de 0.014004819872718 segundos
el espacio en memoria consumido durante la ejecución del algoritmo es de 4896 bytes

Process finished with exit code 0

```

Figura 2. Ejecución 1

```

C:\Users\User\Python\proyecto\datos\analisis\datos\python.exe C:\Users\User\Python\proyecto\datos\analisis\datos.py
estas son todas las palabras que contiene el documento con el numero de veces que aparecen en el mismo:
{'la': 20, 'programación': 9, 'en': 24, 'Python': 2, 'en': 20, 'clave': 3, 'para': 23, 'el': 22, 'trabajo': 2, 'con': 2, 'datos': 11, 'los': 9, 'programa':
estas son las cinco palabras mas usadas en el documento:
{'la': 20, 'la': 20, 'en': 24, 'en': 24, 'para': 23}
lista de las palabras con los indices de cada palabra:
{'la': [0, 2, 5, 7, 10, 11, 15, 16, 17, 18, 20, 27, 28, 30, 31, 35, 37, 39, 41, 45, 47, 48, 50, 51, 54, 55, 45, 61, 63, 65], 'programación': [8, 6, 6,
el tiempo de ejecución del algoritmo es de 0.014004819872718 segundos
el espacio en memoria consumido durante la ejecución del algoritmo es de 12600 bytes

Process finished with exit code 0

```

Figura 3. Ejecución 2

```

C:\Users\User\Python\proyecto\datos\analisis\datos\python.exe C:\Users\User\Python\proyecto\datos\analisis\datos.py
estas son todas las palabras que contiene el documento con el numero de veces que aparecen en el mismo:
{'la': 20, 'programación': 9, 'en': 24, 'Python': 2, 'en': 20, 'clave': 3, 'para': 23, 'el': 22, 'trabajo': 2, 'con': 2, 'datos': 11, 'los': 9, 'programa':
estas son las cinco palabras mas usadas en el documento:
{'la': 20, 'la': 20, 'en': 24, 'en': 24, 'para': 23}
lista de las palabras con los indices de cada palabra:
{'la': [0, 2, 5, 7, 10, 11, 15, 16, 17, 18, 20, 27, 28, 30, 31, 35, 37, 39, 41, 45, 47, 48, 50, 51, 54, 55, 45, 61, 63, 65], 'programación': [8, 6, 6,
el tiempo de ejecución del algoritmo es de 0.014004819872718 segundos
el espacio en memoria consumido durante la ejecución del algoritmo es de 12600 bytes

Process finished with exit code 0

```

Figura 4. Ejecución 3

## III. CÓDIGO EN PYTHON

```

1 import time
2 import psutil
3
4 def contar_palabras(lista):
5     cache = []
6     contador_palabras = {}
7
8     for linea in lista:
9         for palabra in linea.split():
10             if palabra in cache:
11                 contador_palabras[palabra] += 1
12             else:
13                 cache.append(palabra)
14                 contador_palabras[palabra] = 1
15
16     return contador_palabras
17
18 def palabras_mas_usadas(documento):
19     cinco_mayores =
20         dict(sorted(documento.items(), key =
21             lambda item: item[1], reverse = True)
22             [:5])
23     return cinco_mayores
24
25 def buscar_palabra(documento):
26     cache = []
27     indice = {}
28     for contador_linea, linea in
29         enumerate(documento):
30             for palabra in linea.split():
31                 if palabra in cache:
32                     indice[palabra].append(contador_linea)
33                 else:
34                     indice[palabra] =
35                         [contador_linea]
36                     cache.append(palabra)
37
38     return indice

```

En el análisis temporal se obtuvo como resultado un tiempo promedio de 0,014 segundos en la ejecución del algoritmo, mientras que en el análisis espacial se obtuvo 0,62 mega bytes en promedio de la memoria empleada por el algoritmo en su ejecución.

## IV. CONCLUSIÓN

Con la implementación de este tipo de algoritmos, destaca lo fácil que se vuelve el analizar textos en búsqueda de información relevante en un breve período de tiempo. Teniendo esto como punto de partida, podemos emplear estos algoritmos en el estudio de grandes bancos de datos en los cuales se guarda una gran cantidad de información. Si su análisis se realizara de forma manual, emplearía una gran cantidad de tiempo.

Con el análisis temporal y espacial, podemos encontrar formas para mejorar el rendimiento de nuestros algoritmos con

el fin de hacerlos más eficaces en su ejecución, consumiendo la menor cantidad de recursos posibles.