

1.- Creación del diccionario (Combinaciones Binarias).

Felipe Sánchez Martínez
Escuela Superior de Cómputo, IPN

Septiembre de 2023

1. Planteamiento del problema.	1
2. Marco Teórico.	3
2.1. Los Conceptos Centrales de la Teoría de Autómatas	3
2.1.1. Alfabetos	3
2.1.2. Cadenas	3
2.1.3. Lenguajes	4
3. Desarrollo del problema.	5
3.0.1. Approach.	5
3.0.2. Complejidad	6
3.0.3. Código fuente.	6
3.0.4. Casos de prueba.	11
4. Conclusión.	19
5. Referencias.	21

Planteamiento del problema.

La representación binaria es una parte fundamental de las ciencias de la computación, ya que forma la base de la comunicación y el procesamiento de datos en sistemas digitales. En el mundo de la informática, cada número, carácter o símbolo se puede expresar mediante una combinación de dígitos binarios, es decir, unos y ceros.

La siguiente práctica muestra el planteamiento de la creación de un diccionario que aborda la representación de todas las secuencias binarias de longitud 'k', así como sus antecesoras, para poder formar su diccionario. El diccionario consistirá en una lista exhaustiva de todas las combinaciones posibles de dígitos binarios de cierta longitud 'k', desde las secuencias más simples hasta las más complejas. Por ejemplo, comenzando con la cadena vacía ϵ , hasta llegar a tener todas las combinaciones posibles dadas por el valor de 2^k y sus antecesoras $2^{k-1} + 2^{k-2} + \dots + 2^1$.

Así se podrá formar el lenguaje binario, ampliamente utilizado en la programación y la electrónica, que se basa en estas combinaciones de dígitos binarios. Cada secuencia en el diccionario representa una palabra en este lenguaje, y su comprensión es esencial para entender cómo funcionan las operaciones binarias, como la adición, la multiplicación y la representación de datos.

En esta práctica, exploraremos la generación de este diccionario y analizaremos cómo se relaciona con el lenguaje binario, mediante su representación en una gráfica de los números de 1's y cómo se comporta al aplicarle una función logarítmica.

2.1. Los Conceptos Centrales de la Teoría de Autómatas

2.1.1. Alfabetos

Un alfabeto es un conjunto finito y no vacío de símbolos. Convencionalmente, utilizamos el símbolo Σ para representar un alfabeto. Algunos ejemplos comunes de alfabetos son:

1. $\Sigma = \{0, 1\}$, el alfabeto binario.
2. $\Sigma = \{a, b, \dots, z\}$, el conjunto de todas las letras minúsculas.
3. El conjunto de todos los caracteres ASCII o el conjunto de todos los caracteres ASCII imprimibles.

2.1.2. Cadenas

Una cadena (o a veces palabra) es una secuencia finita de símbolos elegidos de algún alfabeto. Por ejemplo, 01101 es una cadena del alfabeto binario $\Sigma = \{0, 1\}$. La cadena 111 es otra cadena elegida de este alfabeto.

La Cadena Vacía

La cadena vacía es la cadena que no contiene ningún símbolo. Esta cadena, denotada como ϵ , es una cadena que puede ser elegida de cualquier alfabeto.

Longitud de una Cadena

A menudo es útil clasificar las cadenas por su longitud, es decir, el número de posiciones para los símbolos en la cadena. Por ejemplo, 01101 tiene una longitud de 5. Es común decir que la longitud de una cadena es 'el número de símbolos' en la cadena; esta afirmación es aceptada coloquialmente pero no es estrictamente correcta. Así, en la cadena 01101, solo hay dos símbolos, 0 y 1, pero hay cinco posiciones para símbolos, y su longitud es 5. Sin embargo, generalmente se puede usar el número de símbolos cuando se quiere decir "número de posiciones". La notación estándar para la longitud de una cadena w es $|w|$. Por ejemplo, $|011| = 3$ y $|\epsilon| = 0$.

2.1.3. Lenguajes

Un conjunto de cadenas, todas elegidas de algún Σ^* , donde Σ es un alfabeto particular, se llama un lenguaje. Si Σ es un alfabeto y $L \subseteq \Sigma^*$, entonces L es un lenguaje sobre Σ . Nótese que un lenguaje sobre Σ no necesita incluir cadenas con todos los símbolos de Σ , por lo que una vez que hemos establecido que L es un lenguaje sobre Σ , también sabemos que es un lenguaje sobre cualquier alfabeto que sea un superconjunto de Σ .

La elección del término "lenguaje" puede parecer extraña. Sin embargo, los lenguajes comunes pueden verse como conjuntos de cadenas. Un ejemplo es el inglés, donde la colección de palabras en inglés legales es un conjunto de cadenas sobre el alfabeto que consiste en todas las letras. Otro ejemplo es C, o cualquier otro lenguaje de programación, donde los programas legales son un subconjunto de las posibles cadenas que se pueden formar a partir del alfabeto del lenguaje. Este alfabeto es un subconjunto de los caracteres ASCII. El alfabeto exacto puede diferir ligeramente entre diferentes lenguajes de programación, pero generalmente incluye letras mayúsculas y minúsculas, dígitos, signos de puntuación y símbolos matemáticos.

Sin embargo, también existen muchos otros lenguajes que aparecen cuando estudiamos autómatas. Algunos son ejemplos abstractos, como:

Desarrollo del problema.

3.0.1. Approach.

Se decidió implementar el algoritmo utilizando dos ciclos 'for'. Por ejemplo, para un valor de $k = 3$, donde k representa la longitud de la cadena, la salida deseada debería tener el siguiente formato:

{ ϵ , 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, 101, 110, 111}

En este caso, el primer ciclo abarcaría el número máximo que se puede formar con k bits (2^k), mientras que el segundo ciclo completaría el conjunto hasta ese punto.

```

1 number_zeros_left = 0
2
3 for i in range(1, k_length + 1):
4     number_zeros_left += 1
5     bits = 2**i
6     for j in range(bits):
7         if i == k_length and j == bits - 1:
8             output.write(bin(j)[2:].zfill(number_zeros_left))
9         else:
10            output.write(bin(j)[2:].zfill(number_zeros_left)+ ', ')

```

Una vez que se haya completado el código fuente, se desea graficar la cantidad de unos por cadena. Para esto, se utilizará la función de peso de Hamming para reducir el tiempo de ejecución en el cálculo del número de unos por cadena. Los resultados se almacenarán en un archivo '.csv' para su posterior manipulación con la biblioteca Pandas.

```

1 def hammingWeight(n):
2     """
3     :type n: int
4     :rtype: int

```

```

5      """
6      mask = 1
7      count = 0
8      while n != 0:
9          if (mask & n == 1):
10             count = count + 1
11             n = n>>1
12
13     return count

```

3.0.2. Complejidad

Complejidad Temporal

La complejidad del algoritmo se divide en dos componentes principales: la generación de las cadenas y la función de Hamming.

Para la generación de las cadenas, la complejidad temporal es:

$$T : O(2^k + 2^{k-1} + 2^{k-2} + \dots + 2^1 + 2^0) = O(2^k)$$

Por otro lado, al considerar la función de Hamming, la complejidad temporal es:

$$T : O(k \cdot 2^k)$$

En conjunto, la complejidad temporal total del algoritmo es:

$$O(2^k) + O(k \cdot 2^k) = O(2^k)$$

Complejidad Espacial

Dentro del algoritmo, la generación del output no se considera en la complejidad espacial, lo que resulta en una complejidad espacial de:

$$S : O(1)$$

Sin embargo, al momento de graficar, se utiliza memoria RAM, lo que eleva la complejidad espacial a:

$$S : O(2^k)$$

El uso de un archivo .csv en combinación con las bibliotecas Pandas y Matplotlib permite una gestión más eficiente de la memoria en RAM al evitar cargar grandes conjuntos de datos completamente en memoria. Sin embargo, es importante tener en cuenta que la lectura y escritura de archivos .csv desde y hacia el disco pueden tener un impacto en el rendimiento en términos de velocidad de acceso al disco."

3.0.3. Código fuente.

```

1 import random
2 import time
3 # graph libs
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 # args
7 import argparse
8 import math
9
10 def main(args):
11
12     while True:
13
14         k_length = input("Pls give me the 'k' for the alphabet or just press
15             'enter': ")
16         if k_length == '':
17             k_length = random.randint(0, 1000)
18         else:
19             k_length = int(k_length)
20
21         print(f"Your k is: {k_length}")
22
23         # OUTPUT FILES
24         output = open('Outputf_BinaryStrings.txt', 'w', encoding='utf-8')
25
26
27
28         # BASIC INFO
29         output.write("{0, 1}")
30         output.write(f"^{k_length} = ")
31         output.write("{\u03B5, ")
32
33
34         # LOGIC
35         number_zeros_left = 0
36
37         for i in range(1, k_length + 1):
38             number_zeros_left += 1
39             bits = 2**i
40             for j in range(bits):
41                 if i == k_length and j == bits - 1:
42                     output.write(bin(j)[2:].zfill(number_zeros_left))
43                 else:
44                     output.write(bin(j)[2:].zfill(number_zeros_left)+ ', ')
45
46         output.write('}')

```

```

47
48 bits = 2**k_length
49 # CLOSE FILES
50 output.close()
51
52 # ""
53 # k =3
54
55 # output : {e, 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, 101, 110,
56             111}
57
58 # T:  $0(2^k + 2^{k-1} + 2^{k-2} + \dots + 2^1 + 2^0) = 0(2^k)$ 
59 # O: 0(1) only output.txt
60
61 # ""
62
63 # # PLOT
64 # ''
65 # Data frame:
66 # # k = 3
67 # # bits = 2**3
68 # bits = 8
69 # chain      number_of_1s
70 # 1          0          -> epsilon
71 # 2          0          -> '000'
72 # 3          1          -> '001'
73 # 4          1          -> '010'
74 # 5          2          -> '011'
75 # 6          1          -> '100'
76 # 7          2          -> '101'
77 # 8          2          -> '110'
78 # 9          3          -> '111'
79
80 # ''
81 # OUTPUT FILES
82 df = open('Binary_Strings.txt', 'w')
83 blog = open('Binary_Strings_Log.txt', 'w')
84
85 # BASIC INFO
86 df.write('chain' + ',' + 'number_of_1s\n')
87 df.write('1' + ',' + '0 \n')
88
89 blog.write('chain' + ',' + 'number_of_1s_base10\n')
90
91
92

```

```

93     for i in range(2, bits + 2):
94         df.write(str(i) + ',' + str(hammingWeight(i - 2)) + '\n')
95         if hammingWeight(i - 2) != 0:
96             blog.write(str(i) + ',' + str(round(math.log10(hammingWeight(i -
97                 2)), 4)) + '\n')
98
99     # CLOSE FILES
100    df.close()
101    blog.close()
102
103    # SAVE DATA FRAME
104    df = pd.read_csv('Binary_Strings.txt')
105
106    # Plotting
107    plt.plot(df['chain'], df['number_of_1s'], 'bo')
108    plt.title('Binary Strings')
109    plt.xlabel('Chain')
110    plt.ylabel('Number of 1s')
111    plt.show()
112
113    # SAVE DATA FRAME
114    blog = pd.read_csv('Binary_Strings_Log.txt')
115
116    # Plotting
117    plt.plot(blog['chain'], blog['number_of_1s_base_log10'], 'bo')
118    plt.title('Binary Strings base Log10')
119    plt.xlabel('Chain')
120    plt.ylabel('Number of 1s')
121    plt.show()
122
123    if input("Do you want to continue? (y/n): ") == 'n':
124        break
125    print("-"*60)
126
127    def hammingWeight(n):
128        """
129        :type n: int
130        :rtype: int
131        """
132        mask = 1
133        count = 0
134        while n != 0:
135            if (mask & n == 1):
136                count = count + 1
137            n = n>>1
138

```

```

139     return count
140
141
142
143
144 def parse_args():
145     # setup arg parser
146     parser = argparse.ArgumentParser()
147
148     # random default value
149     default_k_length = random.randint(0, 1000)
150
151
152     # add arguments
153     parser.add_argument("k_length",
154                         type=int, help="number of primes numbers to check",
155                         default=default_k_length, nargs='?')
156
157     # parse args
158     args = parser.parse_args()
159
160     # return args
161     return args
162
163 # run script
164 if __name__ == "__main__":
165     # add space in logs
166     print("\n\n")
167     print("*" * 60)
168     start = time.time()
169
170     # parse args
171     args = parse_args()
172
173     # run main function
174     main(args)
175
176     end = time.time()
177     print("\n\n")
178     print("Total time taken: {}s (Wall time)".format(end - start))
179     # print("max number of K length: {}".format(args.k_length))
180     # add space in logs
181     print("*" * 60)
182     print("\n\n")

```

3.0.4. Casos de prueba.

Sin nignun Input.

- Output en archivo .txt:

$\{0, 1\}^{11} = \{\epsilon, 0, 1, 00, 01, 10, \dots, 11111111011, 11111111100, 11111111101, 11111111110, 11111111111\}$

- Output en archivo .csv:

```
chain,number_of_1s
1,0
2,0
3,1
4,1
5,2
6,1
7,2
8,2
9,3
10,1
11,2
12,2
...
2043,9
2044,9
2045,10
2046,9
2047,10
2048,10
2049,11
```

- Output en archivo .csv en base log10:

```
chain,number_of_1s_base10
3,0.0
4,0.0
5,0.3010299956639812
6,0.0
7,0.3010299956639812
8,0.3010299956639812
9,0.47712125471966244
10,0.0
...
2043,0.9542425094393249
2044,0.9542425094393249
2045,1.0
2046,0.9542425094393249
```

```
2047,1.0
2048,1.0
2049,1.0413926851582251
```

- **Mensaje en la terminal:**

```
C:\Users\CristoRey\Escom\SEM_5>py Binary_Strings.py
*****
Pls give me the 'k' for the alphabet or just press 'enter':
Your k is: 11
Do you want to continue? (y/n): n

Total time taken: 16.247236251831055s (Wall time)
*****
```

- **Gráfica:**

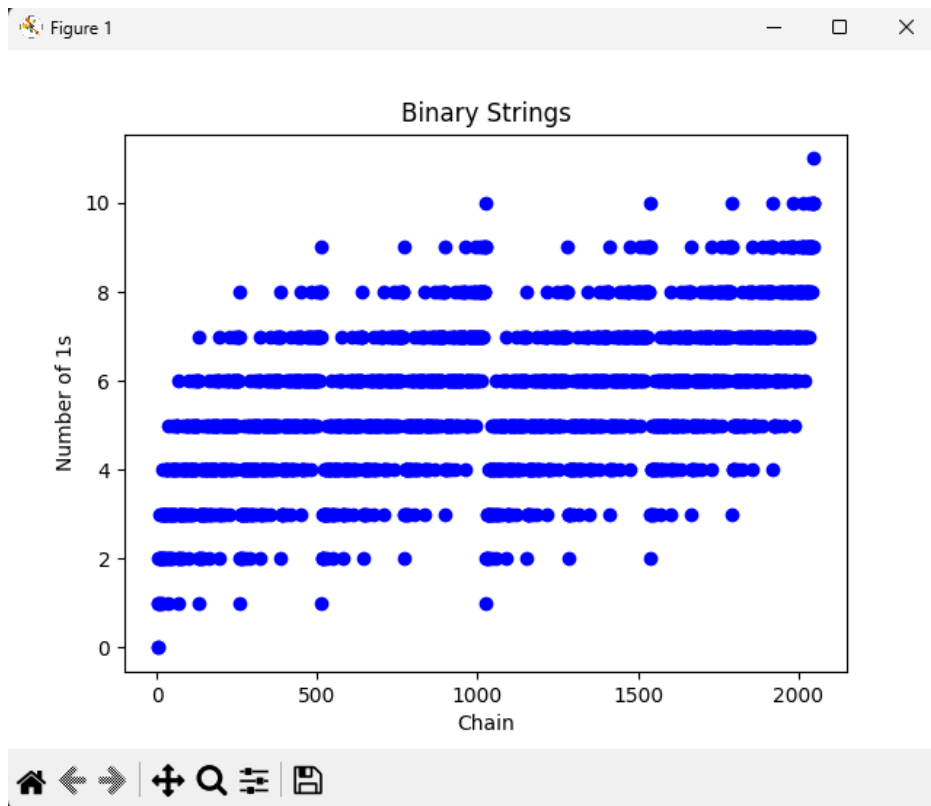


Figura 3.1: Gráfica de puntos con $K = \text{random}(11)$.

- Gráfica base log10:

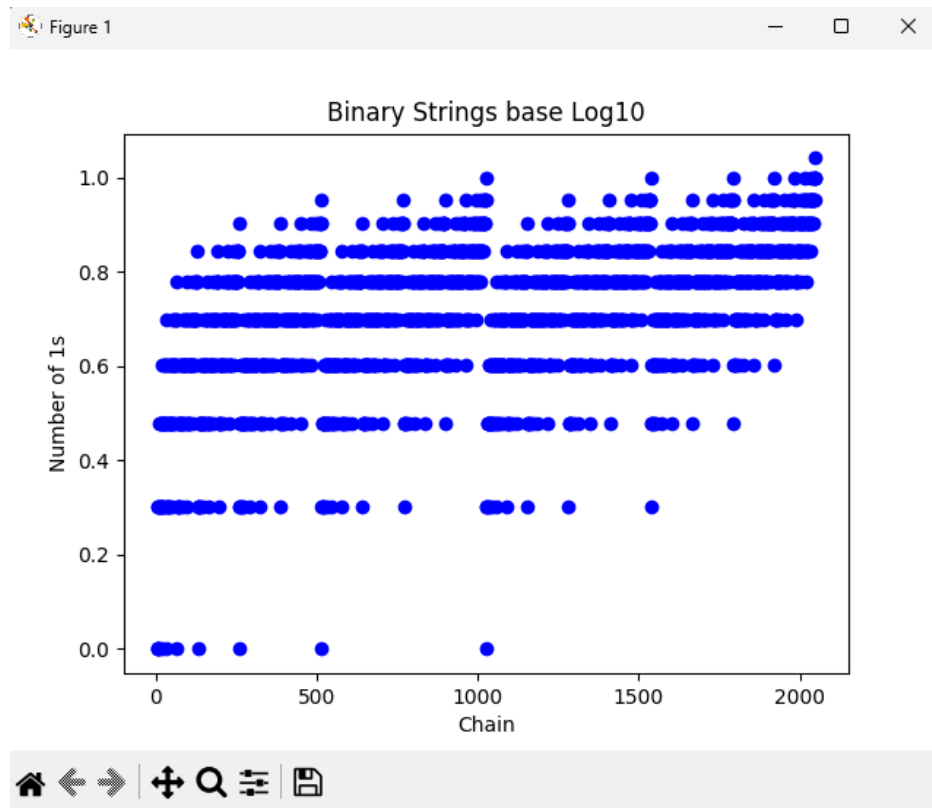


Figura 3.2: Gráfica de puntos con $K = 11$ y log10.

Tamaño de la cadena igual a 3

- Output en archivo .txt:

$$\{0, 1\}^3 = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, 101, 110, 111\}$$

- Output en archivo .csv:

```
chain,number_of_1s
1,0
2,0
3,1
4,1
5,2
6,1
7,2
8,2
9,3
```

- **Output en archivo .csv en base log10:**

```
chain,number_of_1s
chain,number_of_1s_base10
3,0.0
4,0.0
5,0.3010299956639812
6,0.0
7,0.3010299956639812
8,0.3010299956639812
9,0.47712125471966244
```

- **Mensaje en la terminal:**

```
C:\Users\CristoRey\Escom\SEM_5>py Binary_Strings.py
*****
Pls give me the 'k' for the alphabet or just press 'enter': 3
Your k is: 3
Do you want to continue? (y/n): n
```

```
Total time taken: 9.055204629898071s (Wall time)
*****
```

- **Gráfica:**

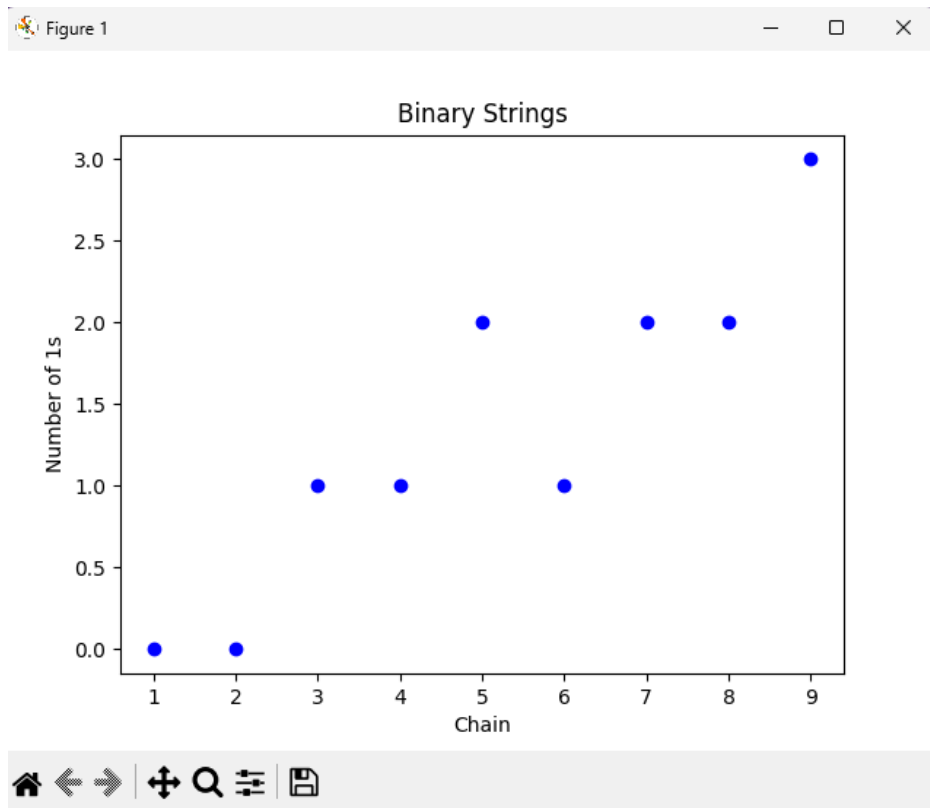


Figura 3.3: Gráfica de puntos con $K = 3$.

- **Gráfica base log10:**

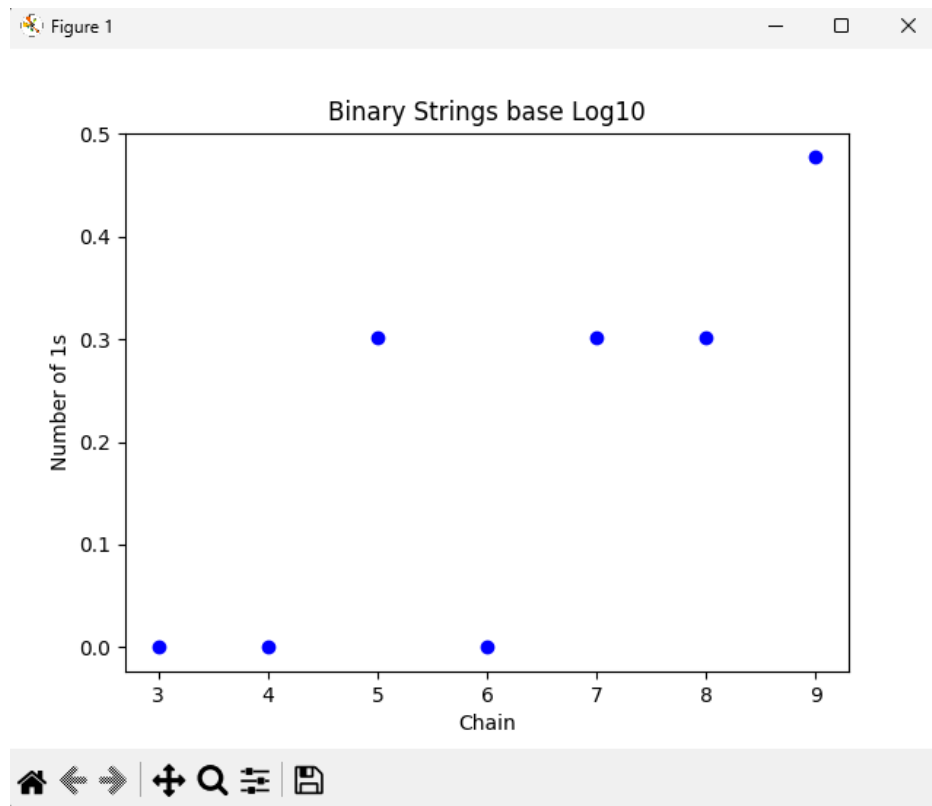


Figura 3.4: Gráfica de puntos con $K = 3$ y \log_{10} .

Tamaño de la cadena igual a 28

- **Output en archivo .txt:**

Usando un valor de K igual a 28, se crea un archivo .txt con un tamaño superior a los 15 GB, lo cual no es posible abrir en un editor de código convencional.

- **Mensaje en la terminal:**

```
C:\Users\CristoRey\Escom\SEM_5>py Binary_Strings.py
*****
Pls give me the 'k' for the alphabet or just press 'enter': 28
Your k is: 28
Do you want to continue? (y/n): n

Total time taken: 2203.679901599884s (Wall time)
*****
```

- Gráfica:

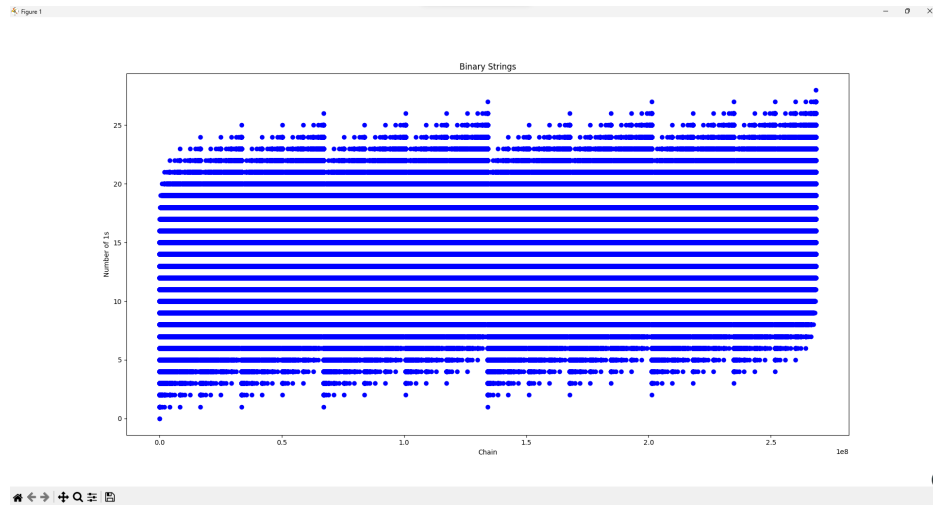


Figura 3.5: Gráfica de puntos con $K = 28$.

- Gráfica base log10:

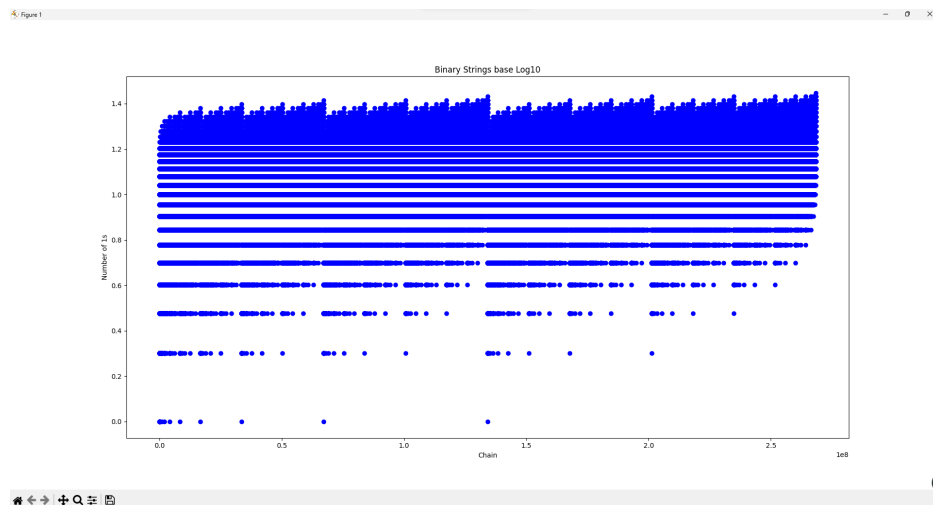


Figura 3.6: Gráfica de puntos con $K = 28$ y log10.

Conclusión.

La creación de un diccionario que comprende todas las posibles combinaciones de longitud 'k' en el sistema binario nos ha proporcionado una visión clara y exhaustiva de las secuencias binarias, desde las más simples hasta las más complejas, y nos ha permitido comprender cómo se generan y representan las palabras en el lenguaje binario.

Dentro de la programación del problema y su complejidad, nos dimos cuenta de que para 'k' igual a 28, el número de combinaciones totales crece de manera gigantesca en cuanto al almacenamiento en el disco. Aunque utilizamos herramientas para manejar una gran cantidad de datos y algoritmos para reducir la complejidad temporal, podemos decir que, aun así, es bastante tardado visualizar una gráfica. Además, observando la gráfica, es demasiado complejo poder visualizar una función con el número de 1's de las cadenas y de igual forma con el algoritmo. Del mismo modo, si usamos una 'k' más baja, nos resulta igual de complicado poder visualizar una función dentro de la gráfica. Posiblemente, utilizando algoritmos enfocados en encontrar el comportamiento de la gráfica, podremos descubrir su naturaleza.

CAPÍTULO 5

Referencias.

- Hopcroft, J. E., Motwani, R., and Ullman, J. D. (2006). Introduction to Automata Theory, Languages, and Computation (3rd ed.). Pearson.