



Instituto Politécnico Nacional
Escuela Superior de Cómputo



Práctica 4

Códigos Huffman

Acevedo Juárez Sebastián

Introducción:

Para comprimir la información se busca representar un conjunto de símbolos (o cadena de símbolos) obtenidos a partir de un cierto alfabeto, usando el menor número de bits posible, pero preservando en todo momento la capacidad de descomprimir o decodificar la información. En general, el sistema que realiza el proceso directo lo llamamos compresor o codificador, mientras que el que reconstruye los datos originales (o una aproximación a ellos si realizamos compresión con pérdidas) lo llamamos descompresor o decodificador. El algoritmo de codificación/compresión Huffman se propuso en 1952 como una forma sencilla y óptima de mapear cada símbolo de un alfabeto con un código (codeword) de longitud óptima. De esta forma, para comprimir cada símbolo de la cadena, simplemente debemos usar el código que se ha calculado mediante Huffman. Para conseguir esta asignación óptima, los símbolos se representan con códigos cuya longitud es inversamente proporcional a la probabilidad del símbolo. De esta forma, los símbolos cuya frecuencia es menor en la cadena se representan con códigos más largos, mientras que los símbolos con mayor frecuencia de apariciones representan con códigos más cortos.

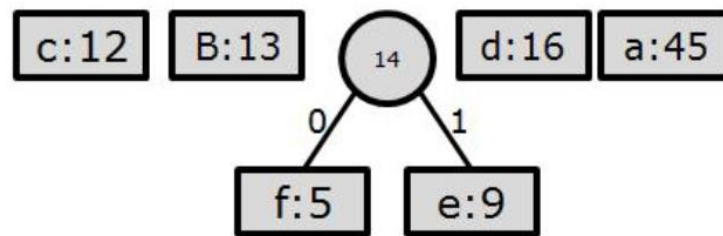
Marco Teórico

Consiste en la creación de un árbol binario en el que se etiquetan los nodos hoja con los caracteres, junto a sus frecuencias, y de forma consecutiva se van uniendo cada pareja de nodos que menos frecuencia sumen, pasando a crear un nuevo nodo intermedio etiquetado con dicha suma. Se procede a realizar esta acción hasta que no quedan nodos hoja por unir a ningún nodo superior, y se ha formado el árbol binario. Posteriormente se etiquetan las aristas que unen cada uno de los nodos con ceros y unos (hijo derecho e izquierdo, respectivamente, por ejemplo. El código resultante para cada carácter es la lectura, siguiendo la rama, desde la raíz hacia cada carácter (o viceversa) de cada una de las etiquetas de las aristas.

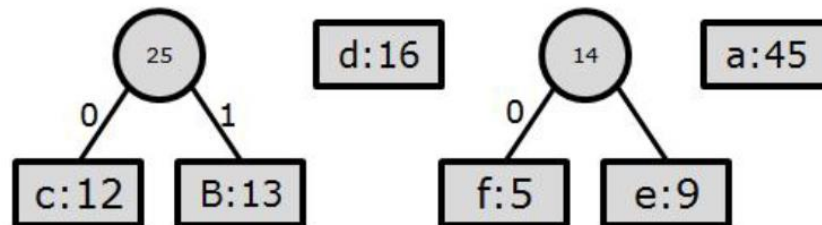
f:5	e:9	c:12	B:13	d:16	a:45
-----	-----	------	------	------	------

Representación de lo antes mencionado:

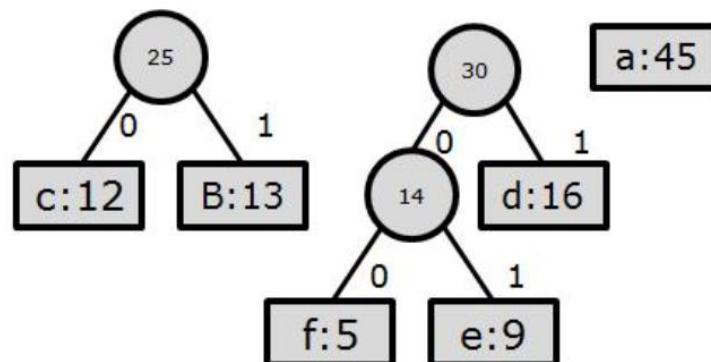
Seleccionamos los dos símbolos con menor frecuencias y los agregamos a un nodo con la suma de ambos, recalcando que el de menor frecuencia ira a la izquierda con el valor 0 en su rama.



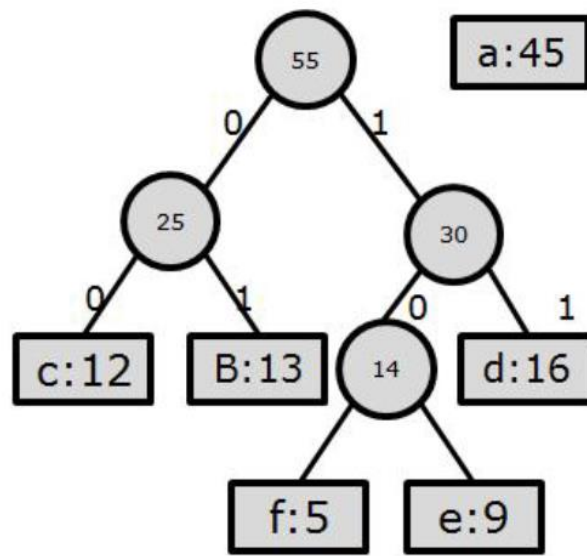
Repetimos el proceso de seleccionar los dos menores sumando sus cantidades.



Observemos que, ahora los dos menores incluyen a un nodo que hemos creado así que usamos este junto con el segundo menor y se crea una nueva rama y un nuevo nivel.



Seguimos con este proceso hasta terminar con los nodos restantes y obtener el árbol completo



Implementación:

Para poder implementar este problema, realizamos los siguientes pasos:

1. Desplegar menú con los archivos disponibles a codificar.
2. Ejecutar función madre Huffman() del archivo seleccionado
 - a. Obtener frecuencia de los caracteres.
 - b. Crear un árbol binario en base a las frecuencias.
 - c. Crear un código en base al árbol previamente creado.
 - d. Codificar el archivo seleccionado.
 - e. Decodificar la compresión y crear un nuevo archivo con el mensaje decodificado.

Al usar una imagen, los pasos tienen algún cambio, debido a que primero se tiene que convertir la imagen a un mapeado, para así poder escribirlo en un txt y realizar los pasos mencionados.

La imagen es representada mediante una matriz que representan los pixeles de una imagen. Al trabajar con una foto en escala de grises, cada celda de la matriz representa la intensidad de color. Este mapa se imprime en un nuevo archivo, al cual se le aplica Huffman.

Al ejecutar la función Huffman(), el archivo de entrada es un archivo ya existente. Sin embargo, el archivo de salida no existe si no se ha ejecutado previamente el programa. En este archivo, se encuentra el mensaje decodificado.

Para “ADN2.txt”

```
Seleccione una opción
1 .- ADN1.txt
2 .- ADN2.txt
3 .- ADN3.txt
4 .- Texto.txt
2

===== Codificación =====

0111010110010010111101000010001111111110101110011110010110010110010110011001111010101011101111010111101111101100101111100
10111111110010100110101000101011111110111110101111001100111011011100110110100010010110110010110110101101100010110010101
01101010111100111100100111000000111001110110000101000101000000100110001010011110011111011010111101100010111101010111111001101
00011010011001101001011011001110101001101000101011101101011101011011010010110101110110111111111110110110101001111011011110
10100111011001010010111110010110110010111000011000011001010010110011110000110011110011110010110101010001010111100110111110100
100101001011001101111011010110000000101000001110110011110010010101101011001001011101000010001111111111010111100111100101100101100
10110010110011001111101010101110111101011110111110110010111111001011111100101001101010001010111111101111101011111001100111011011
1001101101000100101101101001011011010110100101000101101001010110110111100111100100111000001111001110110000101000100000010
011000101001111100111110110101111011000101111101010111111100110100011010011001101001011011001110101001101000101011101101011101101
101001011010111011011111111111110110110101001111011101111101101001110110010100101111100101101101001011111001011011010010110
011111000011001111001111001011010100010101011110011011111101001001010011011011101101100000001010000011101100111110010010110
011100000000000000000000000000001111000000101001010110100001110110010010111101000010001111011001011001011001100111110101011110111
10101111011111011001011111001011111110010100110101000101011111101111101011010001001011011010010110110101101101011011101011001
001011110100001000110011110010110010110010110011001111101010101110111101011110111110110010111111001010011010100010
101111111011111010111100110011011011100110110100010010110110100101011011011100111100101110110000101000101000000
10011000101001111001111101101011101100010111101010111111001101000110100110011011001110110011010001010111101101011101001
10110101001111101110111110110110011011001010011111100101101101001011111000011001111001111001011010110001010111110011011111101
00101011010110101000110011110101011100000001110010100001100000001000100110000110110100001000011110011001111111000000101000011101
1001111100100101
```

```
===== Arbol Binario =====

Nota:
' ' = espacio
\n = salto de línea
{} = Nodo
(a,23) = (letra,frecuencia)
[0] = Camino izquierda
[1] = Camino derecha

[TOTAL] { 1214 }
  [0] { 474 }
    [0] ( T , 223 )
    [1] ( A , 251 )
  [1] { 740 }
    [0] ( G , 343 )
    [1] ( C , 397 )

===== CASO ADN2.txt =====

Archivo de entrada = ADN2.txt || Archivo de salida = newADN2.txt

Bits totales sin comprimir: 9712
Bits totales después de comprimir: 2428
Porcentaje total de compresión = 25.0 %
Tiempo de compresión y decompresión: 2.0505 milisegundos



| Símbolo | Frecuencia | Código |
|---------|------------|--------|
| A       | 251        | 01     |
| C       | 397        | 11     |
| G       | 343        | 10     |
| T       | 223        | 00     |


```

Para “ADN3.txt”

[illegible]

```

===== Arbol Binario =====
Nota:
' ' = espacio
\n = salto de linea
{} = Nodo
(a,23) = (letra,frecuencia)
[0] = Camino izquierda
[1] = Camino derecha

[TOTAL] { 773 }
    [0] { 306 }
        [0] ( T , 141 )
        [1] ( A , 165 )
    [1] { 467 }
        [0] ( G , 215 )
        [1] ( C , 252 )

===== CASO ADN3.txt =====

Archivo de entrada = ADN3.txt || Archivo de salida = newADN3.txt

Bits totales sin comprimir: 6184
Bits totales después de comprimir: 1546
Porcentaje total de compresión = 25.0 %
Tiempo de compresión y decompresión: 1.4894 milisegundos

```

Símbolo	Frecuencia	Código
A	165	01
C	252	11
G	215	10
T	141	00

Para “Texto.txt”

Seleccione una opción

- ```
1 .- ADN1.txt
2 .- ADN2.txt
3 .- ADN3.txt
4 .- Texto.txt
4
```

```
===== Codificacion =====
```

[illegible]



[illegible]

===== CASO Texto.txt =====

Archivo de entrada = Texto.txt || Archivo de salida = newTexto.txt

Bits totales sin comprimir: 15096

Bits totales después de comprimir: 8162

Porcentaje total de compresión = 54.07 %

Tiempo de compresión y decompresión: 3.9375 milisegundos

| Símbolo | Frecuencia | Código      |
|---------|------------|-------------|
| ' '     | 293        | 111         |
| (       | 3          | 000011011   |
| )       | 3          | 000011100   |
| ,       | 19         | 1011111     |
| -       | 3          | 000011101   |
| .       | 13         | 0100011     |
| 3       | 1          | 00001111010 |
| A       | 6          | 01000100    |
| B       | 1          | 00001111011 |
| D       | 8          | 10111100    |
| E       | 1          | 01101011100 |
| H       | 2          | 0110101111  |
| I       | 1          | 01101011101 |
| N       | 6          | 01000101    |
| T       | 3          | 000011111   |
| V       | 1          | 0000110000  |
| W       | 2          | 000011001   |
| \n      | 4          | 011010110   |
| a       | 124        | 1000        |
| b       | 24         | 010000      |
| c       | 69         | 11010       |

|   |     |            |
|---|-----|------------|
| d | 66  | 10110      |
| e | 188 | 001        |
| f | 44  | 00011      |
| g | 42  | 00010      |
| h | 51  | 01001      |
| i | 129 | 1001       |
| j | 1   | 0000110001 |
| k | 3   | 011010100  |
| l | 54  | 01100      |
| m | 34  | 101110     |
| n | 129 | 1010       |
| o | 139 | 1100       |
| p | 39  | 00000      |
| q | 2   | 000011010  |
| r | 106 | 0101       |
| s | 77  | 11011      |
| t | 120 | 0111       |
| u | 32  | 011011     |
| v | 13  | 0110100    |
| w | 8   | 10111101   |
| x | 3   | 011010101  |
| y | 19  | 000010     |
| z | 1   | 0000111100 |

## Para “Imagen.PNG”

```

Seleccione una opción
1 .- ADN1.txt
2 .- ADN2.txt
3 .- ADN3.txt
4 .- Texto.txt
5 .- Imagen.png
5
C:\Users\user\AppData\Local\Programs\Python\Python39\lib\site-packages\PIL\Image.py:975: UserWarning: Palette images with
Transparency expressed in bytes should be converted to RGBA images
 warnings.warn(
[[255 255 255 ... 255 255 255]
 [255 255 255 ... 255 255 255]
 [255 255 255 ... 255 255 255]
 ...
 [255 255 255 ... 255 255 255]
 [255 255 255 ... 255 255 255]
 [255 255 255 ... 255 255 255]]

```

[illegible]

===== Arbol Binario =====

Nota:

' ' = espacio

\n = salto de linea

{ } = Nodo

(a,23) = (letra,frecuencia)

[0] = Camino izquierda

[1] = Camino derecha

```
[Inicio] { 619705 }
 [0] (' ', 302253)
 [1] { 317452 }
 [0] { 132700 }
 [0] { 26084 }
 [0] { 11093 }
 [0] { 4753 }
 [0] { 2107 }
 [0] (4 , 1020)
 [1] { 1087 }
 [0] (\n , 389)
 [1] (1 , 698)
 [1] (8 , 2646)
 [1] { 6340 }
 [0] (6 , 2898)
 [1] { 3442 }
 [0] { 1647 }
 [0] (0 , 712)
 [1] (3 , 935)
 [1] (7 , 1795)
 [1] (2 , 14991)
 [1] (9 , 106616)
 [1] (5 , 184752)
```

===== CASO p1cadena.txt =====

Archivo de entrada = p1cadena.txt || Archivo de salida = newp1cadena.txt

Bits totales sin comprimir: 4957640

Bits totales después de comprimir: 1126410

Porcentaje total de compresión = 22.72 %

Tiempo de compresión y decompresión: 682.2523 milisegundos

| Símbolo | Frecuencia | Código   |
|---------|------------|----------|
| ' '     | 302253     | 0        |
| 0       | 712        | 10001100 |
| 1       | 698        | 10000011 |
| 2       | 14991      | 1001     |
| 3       | 935        | 10001101 |
| 4       | 1020       | 1000000  |
| 5       | 184752     | 11       |
| 6       | 2898       | 100010   |
| 7       | 1795       | 1000111  |
| 8       | 2646       | 100001   |
| 9       | 106616     | 101      |
| \n      | 389        | 10000010 |

Tablas de Resultados

Texto.txt

| Simbolo | Frecuencia | Codigo      |
|---------|------------|-------------|
| ' '     | 293        | 111         |
| (       | 3          | 000011011   |
| )       | 3          | 000011100   |
| ,       | 19         | 1011111     |
| -       | 3          | 000011101   |
| .       | 13         | 0100011     |
| 3       | 1          | 00001111010 |
| A       | 6          | 01000100    |
| B       | 1          | 00001111011 |
| D       | 8          | 10111100    |
| E       | 1          | 01101011100 |
| H       | 2          | 0110101111  |
| I       | 1          | 01101011101 |
| N       | 6          | 01000101    |
| T       | 3          | 000011111   |
| V       | 1          | 0000110000  |
| W       | 2          | 000011001   |
| \n      | 4          | 011010110   |
| a       | 124        | 1000        |
| b       | 24         | 010000      |
| c       | 69         | 11010       |
| d       | 66         | 10110       |
| e       | 188        | 001         |
| f       | 44         | 00011       |
| g       | 42         | 00010       |
| h       | 51         | 01001       |
| i       | 129        | 1001        |
| j       | 1          | 0000110001  |
| k       | 3          | 011010100   |
| l       | 54         | 01100       |
| m       | 34         | 101110      |
| n       | 129        | 1010        |
| o       | 139        | 1100        |
| p       | 39         | 00000       |
| q       | 2          | 000011010   |
| r       | 106        | 0101        |
| s       | 77         | 11011       |
| t       | 120        | 0111        |
| u       | 32         | 011011      |
| v       | 13         | 0110100     |
| w       | 8          | 10111101    |
| x       | 3          | 011010101   |
| y       | 19         | 000010      |
| z       | 1          | 0000111100  |

ADN1.txt

| Simbolo | Frecuencia | Codigo |
|---------|------------|--------|
| A       | 349        | 00     |
| C       | 408        | 01     |
| G       | 469        | 11     |
| T       | 461        | 10     |

ADN2.txt

| Simbolo | Frecuencia | Codigo |
|---------|------------|--------|
| A       | 251        | 01     |
| C       | 397        | 11     |
| G       | 343        | 10     |
| T       | 223        | 00     |

ADN3.txt

| Simbolo | Frecuencia | Codigo |
|---------|------------|--------|
| A       | 165        | 01     |
| C       | 252        | 11     |
| G       | 215        | 10     |
| T       | 141        | 00     |

Imagen.PNG

| Simbolo | Frecuencia | Codigo   |
|---------|------------|----------|
| ' '     | 302253     | 0        |
| 0       | 712        | 10001100 |
| 1       | 698        | 10000011 |
| 2       | 14991      | 1001     |
| 3       | 935        | 10001101 |
| 4       | 1020       | 1000000  |
| 5       | 184752     | 11       |
| 6       | 2898       | 100010   |
| 7       | 1795       | 1000111  |
| 8       | 2646       | 100001   |
| 9       | 106616     | 101      |
| \n      | 389        | 10000010 |

## Conclusión

En la práctica se llevó a cabo un método de codificación (compresión y descompresión) para el cual se formuló una serie de algoritmos recursivos y por demás para poder tener una compresión satisfactoria, nos percatamos que es una codificación sin pérdida de datos. Teniendo un porcentaje de compresión promedio del 30%. La asignación de claves para cada carácter depende de la cantidad de veces que se repita cada carácter en el texto mientras se repita más veces el carácter, tendrá una clave más pequeña mientras menos se repita el carácter tendrá una clave más grande. Independientemente del archivo a comprimir, lo que determina que tan compacto sea el archivo resultante, será el número de veces que se repita los caracteres y símbolos que estén en el mismo. Con esto es evidente que una compresión sumamente eficiente es cuando la gran mayoría de los caracteres se repiten. De forma particular tomando en cuenta el lenguaje que se utilizó (Python) nos apoyamos en algunas librerías que facilitaron el análisis de datos, así como su acomodo e impresión de estos. Observamos que de igual forma se pudieron comprimir imágenes PNG con el mismo principio a coste de operaciones mayor.