

Universidad de San Buenaventura

Facultad ingeniería de sistemas



**UNIVERSIDAD DE
SAN BUENAVENTURA
CALI**

Taller 5 Corte 3

Análisis de algoritmos

Presenta:

Juan Felipe Hurtado Villani

Samuel Martínez

Codificación Huffman:

Prefacio (se puede omitir):

La codificación de Huffman puede comprimir datos de forma muy eficaz: normalmente puede ahorrar entre un 20% y un 90% de espacio, y la tasa de compresión específica depende de las características de los datos. Consideramos que los datos se comprimen como una secuencia de caracteres. De acuerdo con la frecuencia de cada personaje, el algoritmo codicioso de Huffman construye la representación binaria óptima del personaje.

Suponga que queremos comprimir un archivo de datos de 10 caracteres. La siguiente tabla muestra los caracteres que aparecen en el archivo y su frecuencia. En otras palabras, solo aparecieron 6 caracteres diferentes en el archivo, de los cuales el personaje a apareció 45.000 veces.

	a	b	c	d	e	f
Frecuencia (miles)	45	13	12	16	9	5
Codificación de longitud fija	000	001	010	011	100	101
Codificación de longitud variable	0	101	100	111	1101	1100

- Un problema de codificación de caracteres. Un archivo de 100.000 caracteres contiene solo 6 caracteres diferentes de la A la F. La frecuencia de aparición se muestra en la tabla anterior. Si especificamos una palabra de código de 3 dígitos para cada carácter, podemos codificar el archivo con una longitud de 300 000 bits. Pero usando la codificación de longitud variable en la tabla anterior, podemos codificar el archivo con solo 224,000 bits. Entonces vemos que esta es la codificación de caracteres óptima para el archivo. Así que conocemos la codificación de Huffman.

La codificación de Huffman es un tipo de codificación de longitud de palabra variable. Huffman diseñó un algoritmo codicioso para construir el código de prefijo óptimo, que se llama código de Huffman. Árbol de Huffman: El árbol de Huffman también se denomina árbol binario óptimo. Es el árbol binario con la menor longitud de ruta ponderada WPL entre todos los árboles binarios compuestos por n nodos de hojas ponderados.

A continuación, llegamos a comprender el proceso específico del algoritmo de Huffman

Nota: Cada vez que se construye el árbol binario, se ha ordenado en orden creciente de frecuencia de letras.

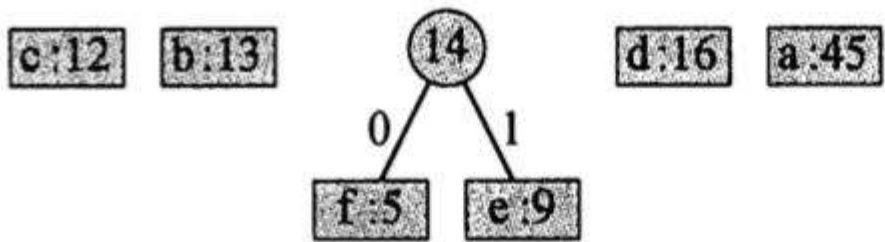
Paso: cada vez que se construye un árbol binario, los dos árboles de baja frecuencia se fusionan

1. Tomamos prestado el ejemplo anterior para obtener el nodo hoja de letra inicial:



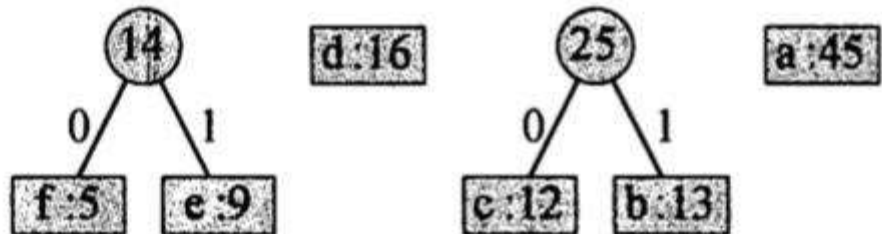
(a)

2. Combina f y e:



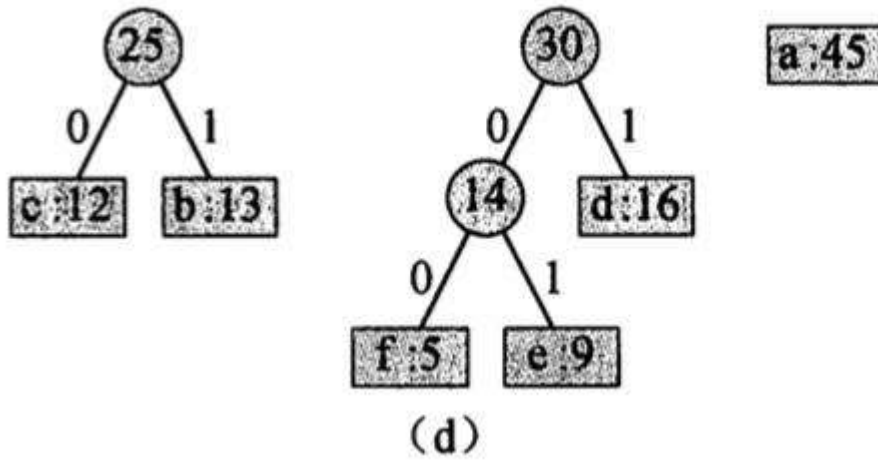
(b)

3. Combina cyb:

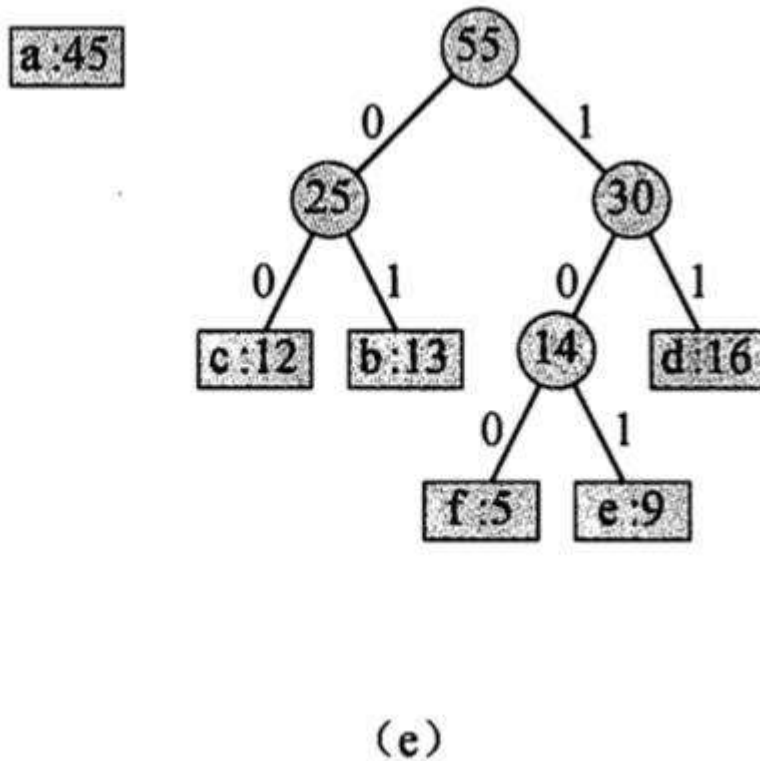


(c)

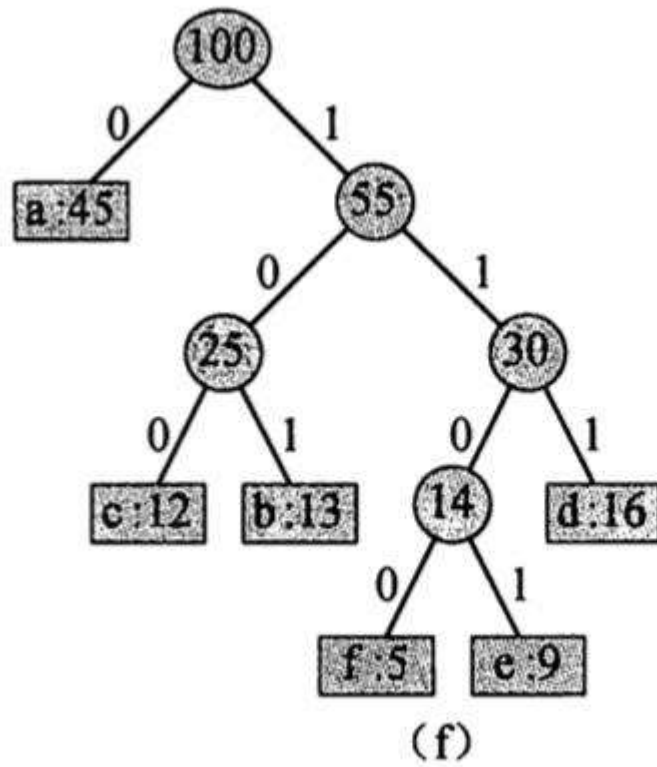
4. Combina 14 y d:



5. Combinar 25 y 30:



6. Combine 45 and 55, and finally obtain a complete Huffman tree:



Entonces, el código de Huffman es:

El código de a es: 0

El código de b es: 101

El código de c es: 100

El código de d es: 111

El código de e es: 1101

El código de f es: 1100

Código para implementar la codificación Huffman

1. Prepara la tabla de probabilidad para cada personaje.

```
1 def findTheCharFrequency(text):
2     result = dict()
3     with open(text, 'r') as f:
4         for line in f.readlines():
5             line = line.lower()
6             for i in line:
7                 if i.isalpha():      #Determina si el personaje es una letra inglesa
8                     if i in result:
9                         result[i] += 1
10                    else:
11                        result.update({i:1})
12     return result
```

2. Crea una clase de nodo

```
1 class Node(object):
2     def __init__(self, name=None, value=None):
3         self.name = name
4         self.value = value
5         self.lchild = None
6         self.rchild = None
```

3. Revertir el establecimiento del árbol de Huffman con nodos

```
1 #Crear árbol Huffman
2 class HuffmanTree(object):
3     # Según la idea del árbol de Huffman: basado en el nodo, construyo el árbol de Huffman al revés
4     def __init__(self, char_Weights):
5         self.Leaf = [Node(k,v) for k, v in char_Weights.items()]
6         while len(self.Leaf) != 1:
7             self.Leaf.sort(key=lambda node: node.value, reverse=True)
8             n = Node(value=(self.Leaf[-1].value + self.Leaf[-2].value))
9             n.lchild = self.Leaf.pop(-1)
10            n.rchild = self.Leaf.pop(-1)
11            self.Leaf.append(n)
12        self.root = self.Leaf[0]
13        self.Buffer = list(range(10))
```

Genera códigos con pensamiento recursivo

```
1 | def Hu_generate(self, tree, length):
2 |     node = tree
3 |     if (not node):
4 |         return
5 |     elif node.name:
6 |         print(node.name + La codificación de Huffman de 'es:', end='')
7 |         for i in range(length):
8 |             print(self.Buffer[i], end='')
9 |         print('\n')
10 |     return
11 |     self.Buffer[length] = 0
12 |     self.Hu_generate(node.lchild, length + 1)
13 |     self.Buffer[length] = 1
14 |     self.Hu_generate(node.rchild, length + 1)
```

Salida de código Huffman

```
1 | def get_code(self):
2 |     self.Hu_generate(self.root, 0)
```

Salida

```
1 | if __name__ == '__main__':
2 |     text = r'123.txt'
3 |     result = findTheCharFrequency(text)
4 |     print(result)
5 |     tree = HuffmanTree(result)
6 |     tree.get_code()
```

Salida de muestra:

```
1 La codificación Huffman de e es: 000
2 La codificación de Huffman de f es: 001
3 La codificación de Huffman de r es: 0100
4 La codificación de Huffman de c es: 0101
5 El código de s de Huffman es: 0110
6 El código de Huffman de i es: 0111
7 La codificación de Huffman de m es: 1000
8 La codificación de Huffman de t es: 1001
9 El código de Huffman de o es: 1010
10 La codificación de Huffman de n es: 1011
11 La codificación de Huffman de w es: 1100000
12 El código Huffman de l es: 1100001
13 La codificación Huffman de x es: 1100010
14 La codificación de Huffman de b es: 11000110
15 La codificación de Huffman de y es: 11000111
16 La codificación de Huffman de d es: 11001
17 La codificación de Huffman de un es: 1101
18 La codificación de Huffman de h es: 1110
19 La codificación de Huffman de p es: 111100
20 La codificación Huffman de g es: 111101
21 La codificación de Huffman de u es: 11111
```


Decodificación Huffman:

Pues resulta que decodificar es mucho más fácil que codificar. Sólo es necesario seguir los siguientes pasos:

1. Cargar el diccionario a la memoria
2. Empezar a leer bits del archivo binario (saltarse el primer 1).
3. Leer bits de uno en uno, y checar si la cadena de bits está en el diccionario
4. Si encontramos la cadena de bits en el diccionario, escribir al archivo de salida
5. Si es el final, terminar. Si no, ir al paso 3.

Como ven, es bastante simple: Leer bits, checar en el diccionario y si es el final, terminar.

Un ejemplo

Antes de presentar código, vamos a hacer un ejemplo rápido. Imaginemos que tenemos un archivo binario que contiene "10000101010111111110", y un diccionario que dice {'a':0, 'b':10, 'c':111, Fin:110}. Empezamos:

1. Cargamos el diccionario: dic = {'a':0, 'b':10, 'c':111, Fin:110}
2. Nos saltamos el primer bit. Entrada: "0000101010111111110"
3. Leemos un bit, y checamos en el diccionario. Bits leídos: "0".
Entrada: "000101010111111110"
4. ¿Está 'bits leídos' disponible en el diccionario? Sí: 'a':0. Salida: a.
5. Los siguientes tres bits son "000", o tres a. Salida: aaaa.
Entrada: "101010111111110"
6. Leemos un bit. Bits leídos: "1". Entrada: "010101111111110". No está en el diccionario.
7. Leemos otro bit. Bits leídos "10". Entrada: "10101111111110". Diccionario: 'b':10.
Salida: aaaab.
8. Los siguientes cuatro bits son "1010", o dos b. Salida: aaaabbb.
Entrada: "111111110"

Y de continuar así, es posible ver que los bits que quedan en la entrada "111111110" corresponden a cc Fin.


```

1  from node import node
2  """
3  CREAR UN ARCHIVO Y ENVIAR UNA PALABRA EN EL
4  """
5
6
7
8
9  # función de utilidad para imprimir huffman
10 # códigos para todos los símbolos en el nuevo
11 # árbol de Huffman creado
12
13
14 def printNodes(node, val=''):
15     # código huffman para el nodo actual
16     newVal = val + str(node.huff)
17
18     # si el nodo no es un nodo de borde
19     # luego atraviesa su interior
20     if(node.left):
21         printNodes(node.left, newVal)
22     if(node.right):
23         printNodes(node.right, newVal)
24     # si el nodo es un nodo de borde, entonces
25     # mostrar su código huffman
26     if(not node.left and not node.right):
27         file.write(f"{node.symbol} ----- {newVal}\n")
28
29
30 # personajes para huffman tree
31 chars = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n',
2  'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z',
3  '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', '.', ',', '']
32
33 freq = []
34
35 filename=str(input("Ingrese el nombre del archivo (nombre.extension): "))
36 f=open(filename,"r")
37 frase=f.read().rstrip()
38 f.close
39 for x in range(len(chars)):
40     counter=0
41     for i in range(len(frase)):
42         if frase[i]==chars[x]:
43             counter+=1
44     freq.append(counter)

```

```

45
46 nodes = []
47
48 # conversión de caracteres y frecuencias
49 # en los nodos del árbol de huffman
50 for x in range(len(chars)):
51     if freq[x]!=0:
52         nodes.append(node(freq[x], chars[x]))
53 numNodos=len(nodes)
54 while len(nodes) > 1:
55     # ordena todos los nodos en orden ascendente
56     # según su frecuencia
57     nodes = sorted(nodes, key=lambda x: x.freq)
58
59     # elige 2 nodos más pequeños
60     left = nodes[0]
61     right = nodes[1]
62
63     # asignar valor direccional a estos nodos
64     left.huff = 0
65     right.huff = 1
66
67     # combina los 2 nodos más pequeños para crear
68     # nuevo nodo como padre
69     newNode = node(left.freq+right.freq, left.symbol+right.symbol, left, right)
70
71     # elimine los 2 nodos y agregue sus
72     # padre como nuevo nodo entre otros
73     nodes.remove(left)
74     nodes.remove(right)
75     nodes.append(newNode)
76
77     # ¡Huffman Tree está listo!
78     file = open("decodificacion.txt","w")
79     file.write("Numero de nodos: "+str(numNodos)+"\n")
80     file.write("Profundidad maxima del arbol generado: "+str(numNodos)+"\n")
81     file.write("Simbolo -Codigo\n")
82     printNodes(nodes[0])
83     file.close()
84
85
86
87

```

Este es el archivo que usaremos para decodificar un archivo con extensión “.txt”, el cual leerá y calculará el numero de nodos, profundidad del árbol y cada letra/símbolo con su respectivo código.

```
node.py x
1 class node:
2     def __init__(self, freq, symbol, left=None, right=None):
3         # frequency of symbol
4         self.freq = freq
5
6         # symbol name (character)
7         self.symbol = symbol
8
9         # node left of current node
10        self.left = left
11
12        # node right of current node
13        self.right = right
14
15        # tree direction (0/1)
16        self.huff = ''
```

Este archivo es el que usaremos para mover dentro del archivo los nodos del árbol.

```
file.txt x
1 jefgylkguyfthdyftdyug , jljol-grgug, hgju jhvcahennuacpp98126128ttdf. .
Console - Shell
Ingreso el nombre del archivo (nombre.extension): file.txt
```

Este es el archivo que usamos de ejemplo que será el que se abrirá, y usaremos para la decodificación Huffman.

```
decodificacion.txt x
1 |Numero de nodos: 30
2 |Profundidad maxima del arbol generado: 30
3 |Simbolo -Codigo
4 |t ----- 0000
5 |a ----- 00010
6 |e ----- 00011
7 |p ----- 00100
8 |r ----- 00101
9 |v ----- 00110
10 |1 ----- 00111
11 |8 ----- 01000
12 |----- 01001
13 |, ----- 01010
14 |b ----- 010110
15 |c ----- 010111
16 |i ----- 011000
17 |k ----- 011001
18 |w ----- 011010
19 |x ----- 011011
20 |0 ----- 011100
21 |3 ----- 011101
22 |6 ----- 011110
23 |9 ----- 011111
24 |f ----- 1000
25 |h ----- 1001
26 |j ----- 1010
27 |u ----- 1011
28 |y ----- 1100
29 |d ----- 11010
30 |o ----- 11011
31 |2 ----- 11100
32 |. ----- 11101
33 |g ----- 1111
34
```

Este es el archivo file.txt decodificado con Huffman.