

Compressão de dados

- Consiste em representar o texto (i.e. seq. de símbolos) original usando menos espaço

- Para isso, basta substituir os símbolos do texto por outros que possam ser representados usando um número menor de bits

	A	B	C	D	E	F
Frequência (x1000)	45	13	12	16	9	5
Código de comprimento fixo	000	001	010	011	100	101
Código de comprimento variável	0	101	100	111	1101	1100

↑ Frequência
↓ bits

- Código de prefixo: são univocamente decodificáveis e possuem codificação instantânea. Nunca um código vai ser prefixo do outro (não é dubl e não tem atraso)

Caractere	Prob.(%)	Cód.I	Cód.II	Cód.III	Cód.IV
A	50	00	0	0	0
B	25	01	1	10	01
C	15	10	00	110	011
D	10	11	11	111	0111
\bar{L}	-	2	1,25	1,75	1,875

- O Código III apresentou um comprimento médio \bar{L} menor que o do Código I, o que o torna mais eficiente.

$$\bar{L} = \sum_{c=1}^n l_c p_c \text{ bits/símbolo, onde}$$

l_c é o tamanho da palavra binária que representa o caractere c e p_c é a probabilidade de ocorrência do caractere c .

Código de Huffman

Algoritmo eficiente e amplamente usado para compressão de dados

Gera sempre código de prefixo

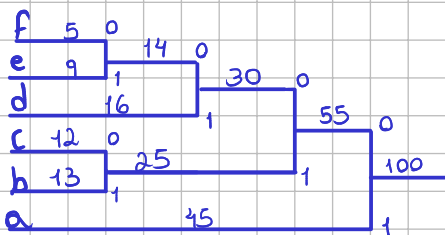
É gerada uma árvore binária de prefixo ótima ao final, ou seja, de comprimento médio mínimo: Árvore de Huffman

É um algoritmo guloso e consiste em:

- Manter uma floresta de árvores com peso igual à soma das frequências de suas folhas.
- Selecionar as árvores T_i e T_k de menores pesos e formar uma nova árvore com T_i e T_k .
- Repetir o passo anterior, até obter a árvore binária final.

Exemplo:

f:5 | e:9 | c:12 | b:13 | d:16 | a:45



a: 1
b: 011
c: 010
d: 001
e: 0001
f: 0000

Comprimento médio:

$$1.45 + 3.13 + 3.12 + 3.16 + 4.9 + 4.5 \\ 45 + 39 + 36 + 48 + 36 + 20 \\ 93 + 59 + 72 = 224 \text{ bits/símbolo}$$

- No algoritmo abaixo, assume-se que C é um conjunto de n caracteres e que cada caractere $c \in C$ é um objeto com um atributo $c.freq$ que dá a sua frequência.

HUFFMAN (C)

- $n = |C|$
- $Q = C$
- para $i = 1$ até $n - 1$
- Aloca um novo nó z
- $z.esquerda = x = \text{Extract-min}(Q)$
- $z.direita = y = \text{Extract-min}(Q)$
- $z.freq = x.freq + y.freq$
- $\text{Insert}(Q, z)$
- retorna $\text{Extract-min}(Q)$ // raiz da árvore

supondo ter um heap, temos: $\log n$

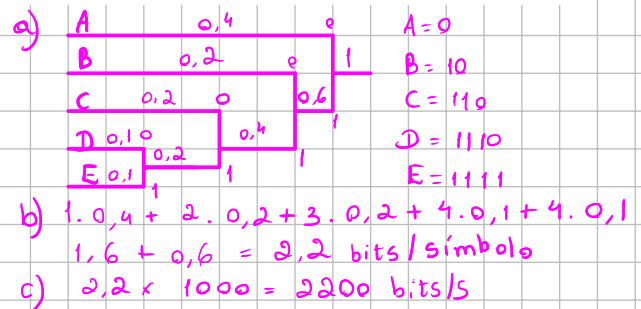
- A complexidade no tempo é $O(n \log(n))$ considerando a fila de prioridades Q implementada com heap binário.

- Com Huffman, podemos fazer pesquisas em um arquivo comprimido diretamente
- O tempo de codificação e decodificação ficam lineares
- Precisa do conhecimento, a priori, das probabilidades dos símbolos
- Tanto o encoder quanto o decoder precisam conhecer a árvore de codificação (ocupa espaço)

- Uma fonte de informação emite 1.000 símbolos/segundo. As probabilidades de ocorrência de cada símbolo do alfabeto são mostradas na tabela abaixo.

Símbolo da fonte	Probabilidade de ocorrência
A	0,4
B	0,2
C	0,2
D	0,1
E	0,1

- Construa uma tabela de codificação usando Huffman.
- Calcule o comprimento médio do código gerado.
- Calcule a taxa média de transmissão em bits/segundo após a codificação da fonte.

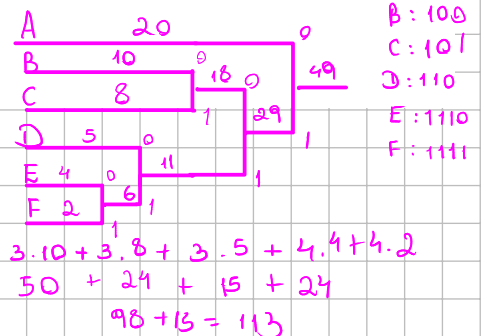


- Suponha que a tabela a seguir apresenta a frequência de cada letra de um alfabeto em uma string.

Quantos bits seriam necessários para representar essa string usando um código de Huffman?

Letra	A	B	C	D	E	F
Frequência	20	10	8	5	4	2

- 392.
- 147.
- 113.
- 108.
- 49.



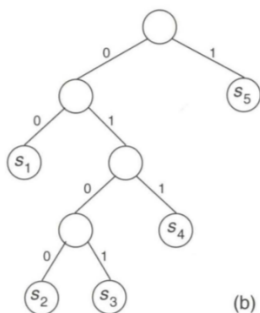
Suponha o texto S:

$S = [s_4 s_3 s_3 s_1 s_3 s_1 s_4 s_5 s_1 s_3 s_3 s_3 s_3 s_3 s_2 s_3 s_5 s_2 s_2 s_2 s_4]$

e os códigos dados pela árvore de prefixo T abaixo.

símbolo	código
s_1	00
s_2	0100
s_3	0101
s_4	011
s_5	1

(a)



(b)

a. Quantos dígitos binários são necessários para codificar o texto S usando a árvore de prefixo T ?

$$3 \cdot 2 + 4 \cdot 4 + 9 \cdot 1 + 3 \cdot 3 + 2 \cdot 1 = 69 \text{ bits}$$

b. O comprimento da sequência binária produzida no item (a) é mínimo considerando códigos de prefixo? Ou seja, é possível codificar o texto S com um número menor de dígitos usando uma árvore binária de prefixo? Explique.

Não, usando o alg. de Huffman, observamos que o comprimento pode ser reduzido a:

$$3 \cdot 3 + 3 \cdot 4 + 9 \cdot 1 + 3 \cdot 3 + 2 \cdot 2$$

$$27 + 12 + 6$$

$$27 + 12 = 45 \text{ bits}$$

Programação dinâmica

✓ é uma técnica de programação ascendente que armazena resultados parciais para acelerar algoritmos recursivos

✓ recursão com apoio de uma tabela

✓ resolve sequencialmente e armazena as soluções

✓ O truque é: quando a solução de um subproblema for requerida, ela já vai estar disponível na tabela

✓ O problema de otimização tem que ter isso para poder aplicar PD:

• Subestrutura ótima (ou otimalidade): e \rightarrow subproblema tem sol. ótima

\rightarrow é a estrutura usada para guardar essas soluções ótimas

• Superposição de subproblemas \rightarrow O alg. recursivo reexamina o mesmo subproblema várias vezes

Fibonacci

• Programação dinâmica: Armazena os valores das entradas menores e não os recalcula.

Fibonacci(n)

1. se ($n = 0$ ou $n = 1$)

2. então retorna (n)

3. senão

4. $F1 = 0$

5. $F2 = 1$

6. para $i = 1$ até $n - 1$ faça

7. $F = F1 + F2$

8. $F1 = F2$

9. $F2 = F$

10. retorna (F)

O Fibonacci recursivo tem complexidade

$$O(2^n)$$

mas com a programação dinâmica, temos:

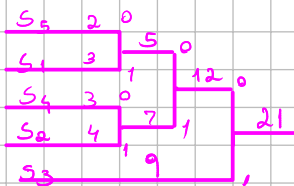
• O algoritmo mantém sempre em memória os dois últimos números da sequência e é executado em $O(n)$.

Pirâmide de números

• Objetivo: Calcular a rota, que começa no topo da pirâmide e acaba na base, com **maior soma**. Em cada passo, podemos ir diagonalmente para baixo e para a esquerda ou para baixo e para a direita.



Símb.	Freq.
s_4	3
s_3	9
s_1	3
s_5	2
s_2	4



$s_1: 001$

$s_2: 011$

$s_3: 1$

$s_4: 010$

$s_5: 000$

Então, o problema pode ser resolvido recursivamente.

Max(i, j)

% $P[i, j]$: j -ésimo número da i -ésima linha

% n : altura da pirâmide

1. se ($i = n$)

2. então soma = $P[i, j]$

3. senão

4. $esq = \text{Max}(i+1, j)$

5. $dir = \text{Max}(i+1, j+1)$

6. se ($esq > dir$)

7. então soma = $P[i, j] + esq$

8. senão soma = $P[i, j] + dir$

9. retorna (soma)

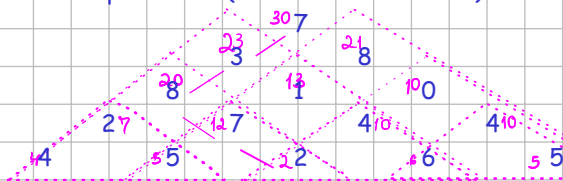
• Para resolver o problema basta executar $\text{Max}(1, 1)$.

• Mas dessa forma, temos o mesmo problema sendo executado várias vezes, dando em uma complexidade $O(2^n)$

• Para melhorar isso, podemos usar PD:

A ideia é guardar o valor calculado para cada subproblema em uma tabela para, quando for o caso, reaproveitá-lo.

• Além disso, começaremos a partir do fim! De baixo para cima (de forma ascendente)



• Tendo em conta a maneira como preenchemos a tabela, é possível aproveitar a matriz P .

Calcular(n)

% n : altura da pirâmide

1. Para $i = n - 1$ até 1

2. Para $j = 1$ até i

3. $esq = P[i+1, j]$

4. $dir = P[i+1, j+1]$

5. se ($esq > dir$)

6. então $P[i, j] = P[i, j] + esq$

7. senão $P[i, j] = P[i, j] + dir$

• Assim, a solução fica em $P[1, 1]$, sendo $n > 1$.

• E o tempo necessário para resolver o problema cresce de forma quadrática: $O(n^2)$.

Assim, a complexidade cai para $O(n^2)$

Algoritmos aproximados

✓ Problema do caixeiro viajante \rightarrow Complexidade $O(n!)$

✓ Como são problemas intratáveis, o algoritmo nem sempre dá a solução ótima

✓ A diferença é que, aqui, nos algoritmos aproximados, ele entrega uma solução dentro de um limite especificado (tem compromisso com a solução ótima)

✓ Dependendo do problema, o objetivo será minimizar ou maximizar o valor da solução ótima para um problema

Heurística

✓ Sem compromisso com a solução ótima

✓ Pode produzir um bom resultado, até mesmo a solução ótima, mas também pode não produzir solução alguma ou uma solução muito distante do ótimo.