

Aluno: Vinícius Guimarães RA: 802431  
Vitor Enzo RA: 802123



1) O código de Huffman é um algoritmo guloso utilizado para realizar compressão de dados sem haver perda de informação. A ideia geral é que o algoritmo irá atribuir códigos binários para cada caractere de acordo com a probabilidade de ocorrência dos mesmos. Assim, um caractere que aparece várias vezes em um texto, terá código binário menor do que quando comparado a um caractere que aparece poucas vezes.

Por exemplo, considerando  $\Sigma = (a, b, c, d, e)$  como o alfabeto de símbolos a serem codificados e  $w_i$  a probabilidade de ocorrência de  $p_i$  nos dados, temos:

	a	b	c	d	e	f
w	0.37	0.12	0.09	0.25	0.10	0.07

→ Soma das probabilidades dá 1

Utilizando a forma tradicional, temos que utilizar 3 bits pois  $2^2 = 4$  e temos 6 caracteres, obtendo a seguinte codificação quando iniciada pelo símbolo a:

a = 000, b = 001, c = 010, d = 011, e = 100, f = 101

O que nos dá um tamanho médio do código de:

$$I(c) = \sum_{i=1}^n w_i \cdot \text{tamanho do código do caractere } p_i = 0,37 \cdot 3 + 0,12 \cdot 3 + 0,09 \cdot 3 + 0,25 \cdot 3 + 0,10 \cdot 3 + 0,07 \cdot 3 = \boxed{3}$$

Ou seja, um média é preciso 3 bits para codificar 1 caractere.

Porém, com a aplicação do algoritmo de Huffman, temos as seguintes codificações:

a = 11, b = 011, c = 001, d = 10, e = 010, f = 000

O que gera um tamanho médio de:

$$I(c) = 0,37 \cdot 2 + 0,12 \cdot 3 + 0,09 \cdot 3 + 0,25 \cdot 2 + 0,1 \cdot 3 + 0,07 \cdot 3 = \boxed{2,38}$$

Portanto, sendo o tamanho médio de cada caractere como sendo 2,38, temos que isso é aproximadamente 20,6% menor do que a codificação tradicional. Portanto, ao gerar códigos de acordo com a probabilidade de ocorrência em um texto, o algoritmo de Huffman consegue gerar códigos menores do que a forma tradicional.

2) Dado o código abaixo

Huffman( $\alpha$ ) {

$n = |\alpha|$

for each  $x_i \in \alpha$  → Estamos inserindo cada elemento do conjunto  $\alpha$  dentro de  $Q$  (Onde  $Q$  é uma estrutura de dados que serviria para pegarmos a árvore de menor probabilidade)

Insert( $Q, x_i$ )

for  $i = 1$  to  $n-1$  {

  Create a new node  $z$

$X = \text{ExtractMin}(Q)$

$Y = \text{ExtractMin}(Q)$

$z.\text{left} = X$

$z.\text{right} = Y$

$w(z) = w(X) + w(Y)$

  Insert( $Q, z$ )

}  
return  $\text{ExtractMin}(Q)$

→ Se temos 3 caracteres por exemplo, teremos que realizar duas junções de árvores

→ Começamos a criar um nó  $z$ , que será a raiz da árvore e adicionamos os dois menores elementos de  $Q$  dentro de  $z$ . Além disso, a probabilidade de  $z$  é atualizada para ser o somatório das probabilidades de  $X$  e  $Y$ .

→ Em seguida inserimos  $z$  dentro de  $Q$  para continuar o processo

→ Ao final da execução do for, teremos a árvore de Huffman, que será retornada

O código de Huffman é guloso pois a cada execução do for, são pegos da fila de prioridades os 2 nós que estão mais longe da raiz da árvore quanto comparado com aqueles que ainda estão na lista de prioridades, ou seja, aqueles com a maior prioridade (menor probabilidade)



③

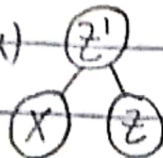
e	r	o	t	n	l	z	x
34	22	24	28	15	10	9	8

j=1

e	r	o	t	n	l	z	x
34	22	24	28	15	10	9	8

Escollidos: x, z

Frequência de  $z'$  =  $w(z) + w(x)$   
= 17



j=2

e	r	o	t	n	l	z'
34	22	24	28	15	10	17

Escollidos: l, n

Frequência de  $z''$  = 25

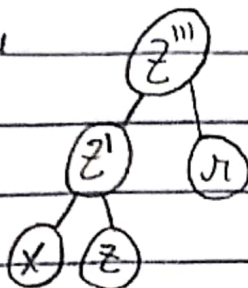


j=3

e	r	o	t	z'	z''
34	22	24	28	17	25

Escollidos: z', r

Frequência de  $z'''$  = 39

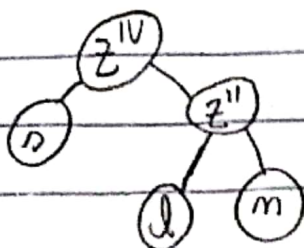


j=4

e	o	t	z''	z'''
34	24	28	25	39

Escollidos: o, z''

Frequência de  $z^{iv}$  = 49

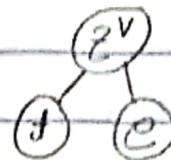


j=5

e	t	z'''	z^{iv}
34	28	39	49

Escollidos: t, e

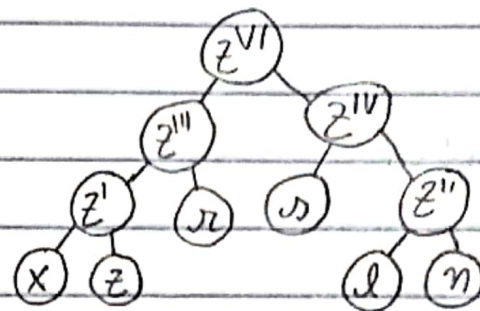
Frequência de  $z^v$  = 62



j=6

z'''	z^{iv}	z^v
39	49	62

Frequência de  $z^{vi}$  =  $39 + 49 = 88$

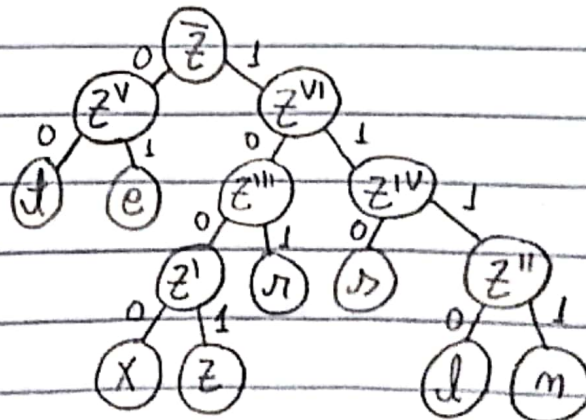


j=7

z^v	z^{vi}
62	88

Escollidos: z^v, z^{vi}

Frequência de  $\bar{z}$  =  $62 + 88 = 150$



Assim, tendo a árvore de Huffman calculada, obtemos os seguintes códigos binários para os símbolos:

$c = 01$ ,  $n = 101$ ,  $d = 110$ ,  $t = 00$ ,  $m = 1111$ ,  $l = 1110$ ,  $z = 1001$ ,  $x = 1000$

Se considerarmos a estratégia padrão, teríamos 3 bits para cada caractere.

Com Huffman temos:  $34.2 + 22.3 + 24.3 + 28.2 + 15.4 + 10.4 + 9.4 + 9.4 = 218.66$

Ou significa uma redução de cerca de 5%.<sup>150</sup>

# At2 - Algoritmos gulosos (Código de Huffman)

Questão 04)

$\alpha = (a, b, c, d, e, f, g, h)$

$w = (1, 1, 2, 3, 5, 8, 13, 21)$

Aplicando o algoritmo de Huffman:

$i=1: Q^{(1)} = \left[ \begin{array}{|c|c|c|c|c|c|c|c|} \hline a & b & c & d & e & f & g & h \\ \hline 1 & 1 & 2 & 3 & 5 & 8 & 13 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = a \\ y = b \end{array} \quad w(z^1) = 2$

$i=2: Q^{(2)} = \left[ \begin{array}{|c|c|c|c|c|c|c|} \hline z^1 & c & d & e & f & g & h \\ \hline 2 & 2 & 3 & 5 & 8 & 13 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = z^1 \\ y = c \end{array} \quad w(z^2) = 4$

$i=3: Q^{(3)} = \left[ \begin{array}{|c|c|c|c|c|c|} \hline z^2 & d & e & f & g & h \\ \hline 4 & 3 & 5 & 8 & 13 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = d \\ y = z^2 \end{array} \quad w(z^3) = 7$

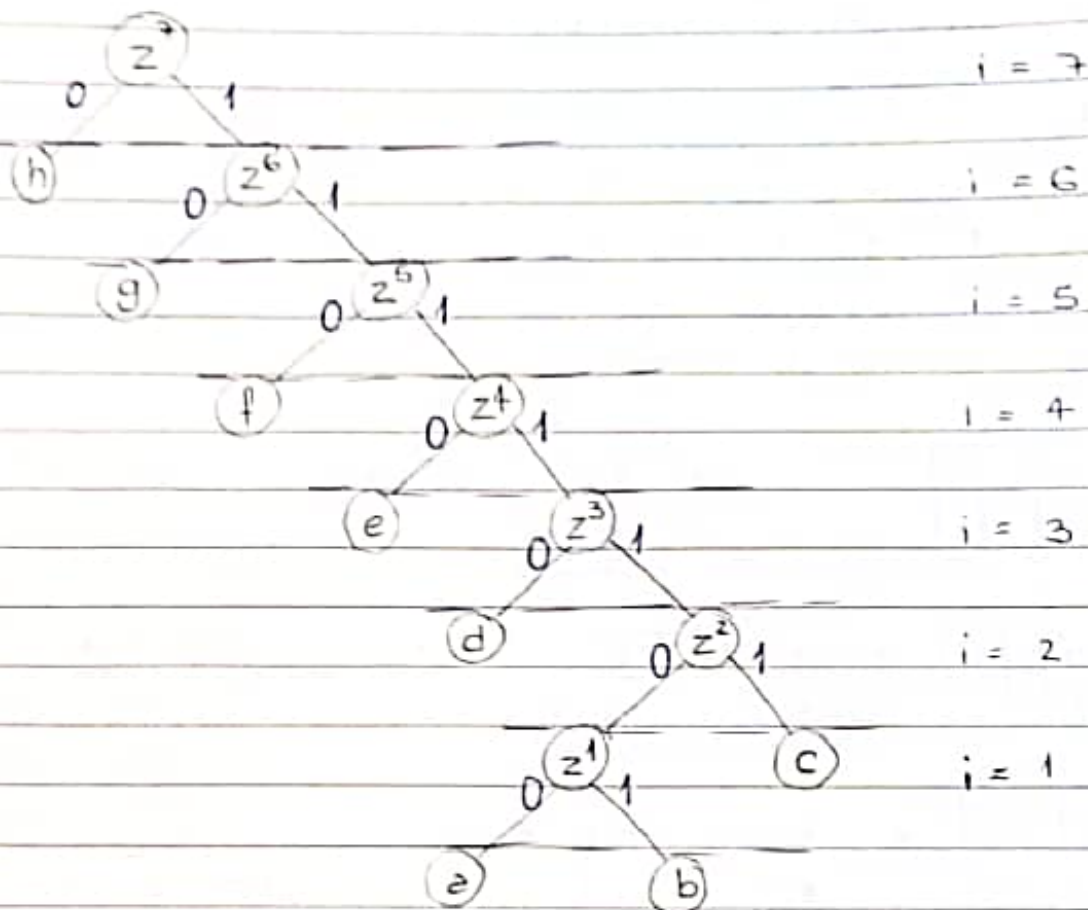
$i=4: Q^{(4)} = \left[ \begin{array}{|c|c|c|c|c|} \hline z^3 & e & f & g & h \\ \hline 7 & 5 & 8 & 13 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = e \\ y = z^3 \end{array} \quad w(z^4) = 12$

$i=5: Q^{(5)} = \left[ \begin{array}{|c|c|c|c|} \hline z^4 & f & g & h \\ \hline 12 & 8 & 13 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = f \\ y = z^4 \end{array} \quad w(z^5) = 20$

$i=6: Q^{(6)} = \left[ \begin{array}{|c|c|c|} \hline z^5 & g & h \\ \hline 20 & 13 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = g \\ y = z^5 \end{array} \quad w(z^6) = 33$

$i=7: Q^{(7)} = \left[ \begin{array}{|c|c|} \hline z^6 & h \\ \hline 33 & 21 \\ \hline \end{array} \right] \quad \begin{array}{l} x = h \\ y = z^6 \end{array} \quad w(z^7) = 54$





### Tabela de Símbolos

$h = 0$

$d = 11110$

$g = 10$

$c = 111111$

$f = 110$

$b = 1111101$

$e = 1110$

$a = 1111100$

Assim, a cadeia de bits de "fdheg" seria:

$fdheg = 11011110111010$

Generalizando, para  $n = 2$ , temos:

$$s_1 = 0$$

$$s_2 = 1$$

para  $n \geq 3$ :

$$s_1 = (n-3) \cdot '1' + '00'$$

$$s_i = (n-i) \cdot '1' + '0'$$

$$s_2 = (n-3) \cdot '1' + '01'$$

$$s_3 = (n-3) \cdot '1' + '1'$$

:

Questão 05) Sabe-se que, de acordo com Claude Shannon e sua teoria de informação, a quantidade de informação do símbolo  $\alpha_i$  é dada por

$$h(\alpha_i) = \log_2 \frac{1}{w_i}$$

ou seja, é inversamente proporcional à probabilidade

A entropia associada à distribuição de probabilidades do vetor  $\vec{w}$  é dada por:

$$H(\vec{w}) = \sum_i w_i \cdot h(\alpha_i)$$

$$= \sum_i w_i \cdot \log_2 \frac{1}{w_i}$$

$$= \sum_i w_i \cdot \log_2 w_i^{-1}$$

$$H(\vec{w}) = - \sum_i w_i \cdot \log_2 w_i$$

A prova, consiste em notar que ~~nao~~ o tamanho do código  $c$   $l(c)$  será igual a  $H(\vec{w})$  quando  $l(\alpha_i) = h(\alpha_i)$ .

Desse-se mostrar que  $H(\vec{w}) \leq l(c)$ . Iremos iniciar calculando  $H(\vec{w}) - l(c)$ ,

$$H(\vec{w}) - l(c) = \sum_j w_j \cdot \log_2 \frac{1}{w_j} - \sum_j w_j \cdot l(\alpha_j)$$

Note que

$$l(\alpha_j) = \log_2 2^{l(\alpha_j)}$$

Logo,

$$\begin{aligned} H(\vec{w}) - l(c) &= \sum_j w_j \cdot \log_2 \frac{1}{w_j} + \sum_j w_j \cdot \log_2 2^{-l(\alpha_j)} \\ &= \sum_j w_j \cdot \log_2 \frac{2^{-l(\alpha_j)}}{w_j} \end{aligned}$$

Realizando troca de base no logaritmo temos

$$\log_2 \left( \frac{2^{-l(\alpha_j)}}{w_j} \right) = \frac{\ln \left( \frac{2^{-l(\alpha_j)}}{w_j} \right)}{\ln 2} = \frac{1}{\ln 2} \cdot \ln \frac{2^{-l(\alpha_j)}}{w_j}$$

Portanto,

$$H(\vec{w}) - l(c) = \frac{1}{\ln 2} \sum_j w_j \ln \frac{2^{-l(\alpha_j)}}{w_j}$$

Tomaremos como verdade que  $\ln x \leq x - 1$ , segue

$$\begin{aligned} H(\vec{w}) - l(c) &\leq \frac{1}{\ln 2} \sum_j w_j \left( \frac{2^{-l(\alpha_j)}}{w_j} - 1 \right) \\ &\leq \frac{1}{\ln 2} \left( \sum_j 2^{-l(\alpha_j)} - \sum_j w_j \right) \\ &\leq \frac{1}{\ln 2} \left( \sum_j 2^{-l(\alpha_j)} - 1 \right) \end{aligned}$$

Seja  $l_{\max} = \max \{ l(\alpha_1), l(\alpha_2), \dots, l(\alpha_n) \}$ . Então,

$$\sum_{j=1}^{l_{\max}} \frac{l_{\max} - l_j}{2} = \frac{l_{\max} - 1}{2} + \frac{l_{\max} - 2}{2} + \dots + \frac{l_{\max} - l_{\max}}{2}$$

O número máximo de nós folhas é  $2^{l_{\max}}$ . Então,

$$\sum_{j=1}^{l_{\max}} \frac{l_{\max} - l_j}{2} < 2^{l_{\max}}$$

Dividindo ambos lados por  $2^{l_{\max}}$ , tem-se

$$\sum_{j=1}^{l_{\max}} \frac{l_{\max} - l_j}{2^{l_{\max}}} < 1$$

Portanto,

$$H(\vec{w}) - l(c) \leq 0$$

$$H(\vec{w}) \leq l(c)$$

no código de Huffman, quanto menor a probabilidade, maior a informação, o que resulta em  $l(c) = H(\vec{w})$ .