

Huffman

Problema della compressione dei dati

"Ridurre le dimensioni del file"

Supponiamo di avere un testo T

Σ : $(T) = 100$

a: 20

b: 7

c: 3

d: 10

e: 15

f: 1

g: 30

h: 4

i: 5

l: 5

Carattere: numero di
occorrenze
nel testo

ogni carattere consume 10 bit

Per codificare tutto abbiamo bisogno: $100 \cdot 10 = 1000$ bit

"10 bit sono troppi"

cod

a: 20 0000

b: 7 0001

c: 3 0010

d: 10 0011

e: 15 0100

f: 1 0101

g: 30 0110

h: 4 0111

i: 5 1000

l: 5 1001

Codifichiamo "Dieci"

D i e c i
0011 - 1000 - 0100 - 0010 - 1000

e l'intero computerizza e
memorizza queste codifiche

Se le facilità di codificare e decodificare non è importante possiamo usare un dimensionamento dinamico per i nostri caratteri usando meno bit per i caratteri più frequenti e di più per quelli meno frequenti

	COD	COD N
e : 20	0000	
l : 4	0001	
- c : 3	0010	01001
- d : 10	0011	01
- e : 15	0100	1
f : 1	0101	
g : 30	0110	
h : 4	0111	
- i : 5	1000	1100
l : 5	1001	

la nuova codifica di "Dieci" è:

01110010100111000

↑
decifrare questa cosa è difficile finché non si ha le lunghezze fine

Dobbiamo stabilire dei codici fidi detti "prefissi" per rendere la decodifica fattibile, questi codici hanno queste caratteristiche:

• una codifica non deve avere un prefisso di un'altra codifica

	COD	COD (aggiungendo 1)
e : 20	0000	10
l : 4	0001	1100
c : 3	0010	1101
d : 10	0011	000001
e : 15	0100	001
f : 1	0101	000000
g : 30	0110	01
h : 4	0111	111
i : 5	1000	0001
l : 5	1001	00001

400

298

abbiamo perso un sacco di tempo per creare le nuove colonne

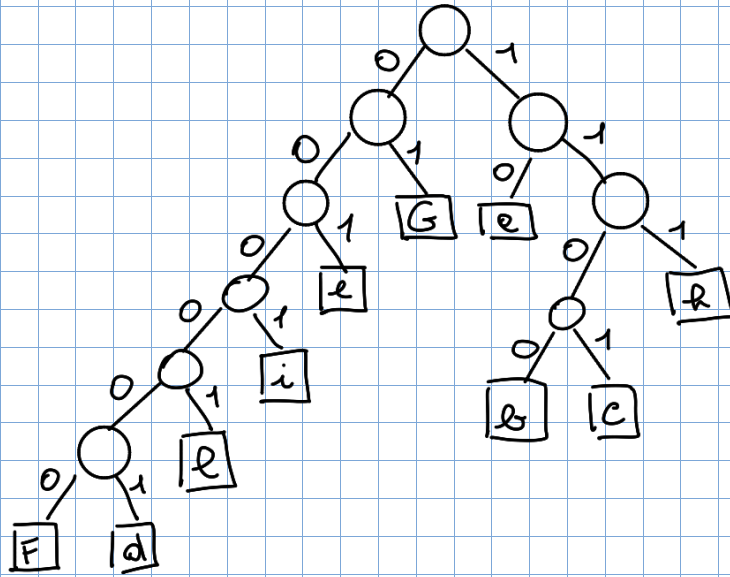
la codifica di "Dieci" diventa:

000001 0001 001 1101 0001

anche se ha dimensioni variabili riusciamo a decodificare

Huffman ci dà un modo facile creare le tabelle in modo facile e le fa tornare in albero.

albero delle nostre codifiche:



L'obiettivo di creare l'albero che codifica meglio i nostri caratteri:

alfabeto

$$\forall c \in \Sigma : f(c) = \text{freq}$$
$$d_+(c) = \text{profondità}$$

$$B(+) = \sum_{c \in \Sigma} f(c) \cdot d_+(c)$$

↓

Costo di T

} usiamo queste funzione per valutare il nostro albero

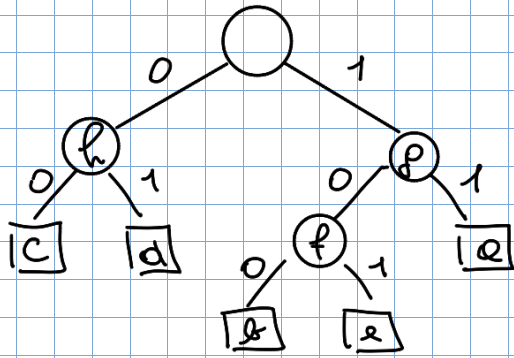
↳ quanti bit usa l'albero T

Dobbiamo trovare una strategia che minimizza $B(+)$ lo faremo ricorsivamente, greedy.

Ma come riduciamo in sottoproblemi? Cerchiamo di definire un caso base:

Prendiamo due lettere del nostro alfabeto e le sostituisco con una

e
 b
 c
 d
 e
 f
 g
 h



de questo albero
 nascono i codici prefissi

Scegliamo 2 caratteri e lo sostituiamo con 1 solo
 Come trasformiamo i 2 caratteri migliori da scegliere?

Dimostrazione sottostrutture ottime
 Σ soluzione T

$$B(T) = B(T') + f(e) + f(b)$$

$$\Sigma' = \Sigma - \{e, b\} \cup \{z\} \quad \text{soluzione } T'$$

$$B(T) \sim B(T')$$

$$f(z) = f(e) + f(b)$$

$$d_+(e) = d_+(z) + 1$$

$$d_+(b) = d_+(z) + 1$$

$$\begin{aligned}
 B(T) &= \sum_{c \in \Sigma} f(c) \cdot d_+(c) = \left[\sum_{c \in \Sigma'} f(c) d_+(c) \right] - \underbrace{f(z) \cdot d_+(z)}_{f(e) + f(b)} + \\
 &+ f(e) d_+(e) + f(b) d_+(b) \\
 &\quad \nwarrow \quad \nearrow \\
 &\quad d_+(z) + 1 \quad \quad \quad = f(e) + f(b)
 \end{aligned}$$

$$\begin{aligned}
 & - (f(e) + f(b)) d_+(z) + f(e)(d_+(z) + 1) - f(b)(d_+(z) + 1) \\
 & = f(e) + f(b)
 \end{aligned}$$

suppongo che le sottostrutture ottime non esiste

$$\exists T'' \quad B(T'') < B(T')$$

$$B(T''') = B(T'') + f(e) + f(b)$$

Come individuiamo le scelte migliori durante la costruzione dell'albero di Hoffmann

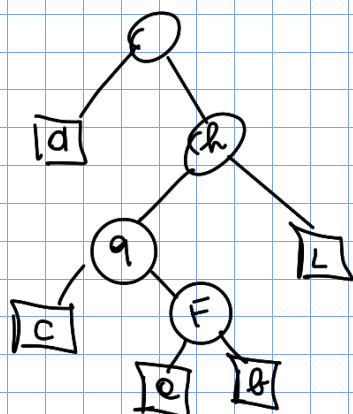
$e : 5$

$b : 5$

$c : 10$

$d : 40$

$e : 40$

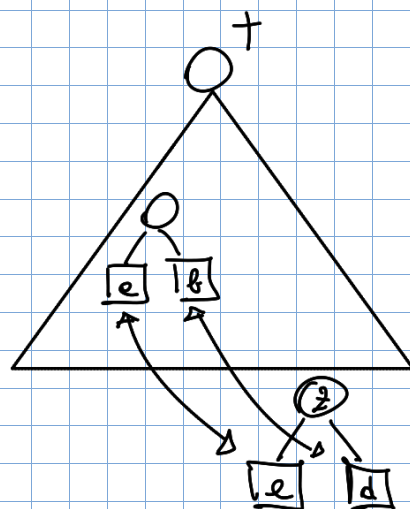


Σ

$e, b \in \Sigma$

$f(e) \leq f(c) \quad \forall c \in \Sigma$

$f(b) \leq f(c) \quad \forall c \in \Sigma - \{e\}$



$$b(+') = \sum_{c \in \Sigma} f(c) \quad d_{+'}(c) = \sum_{c \in \Sigma - \{e, b, e, d\}}$$

$$f(c) d_{+'}(c) + \begin{cases} f(e) d_{+'}(e) \\ f(b) d_{+'}(b) \\ f(e) d_{+'}(e) \\ f(d) d_{+'}(b) \end{cases}$$

$$f(e) d_{+'}(e) + f(b) d_{+'}(b) + f(e) d_{+'}(e) + f(d) d_{+'}(b)$$

$$d_{+'}(e) (F(e) + F(b)) + d_{+'}(e) (F(e) + F(d))$$

$b(+') \geq b(+)$ per ottenere una contraddizione

$$b(+') - b(+)$$

$$\begin{aligned} b(+') - b(+)&= f(e) d_{+'}(e) + f(b) d_{+'}(b) + f(e) d_{+'}(e) + f(d) d_{+'}(d) \\ &= f(e) d_{+'}(e) - f(b) d_{+'}(b) - f(e) d_{+'}(e) - F(d) d_{+'}(d) \end{aligned}$$

$$f(e)(d_{+1}(e) - d_+(e))$$

$$f(e)(d_+(e) - d_{+1}(e))$$

$$f(e) \cdot (d_{+1}(e) - d_+(e) - f(e)(d_{+1}(e) - d_+(e)))$$

$$(d_{+1}(e) - d_+(e))(f(e) - f(e))$$

HUFFMAN(Σ, f)

$Q \leftarrow$ new Priority Queue

foreach $c \in \Sigma$ do:

$x \leftarrow$ new Node(c)

insert x in Q

} $m \log m$

for $i \leftarrow 1$ to $|\Sigma| - 1$ do

$x \leftarrow$ extractMin(Q)

$y \leftarrow$ extractMin(Q)

$z \leftarrow$ new Node

$f(z) = f(x) + f(y)$

left(z) = x

right(z) = y

insert z in Q

} $\log m$

} $m \log m$

$$m = |\Sigma|$$

~~m: 4~~

~~n: 8~~

~~o: 10~~

~~p: 15~~

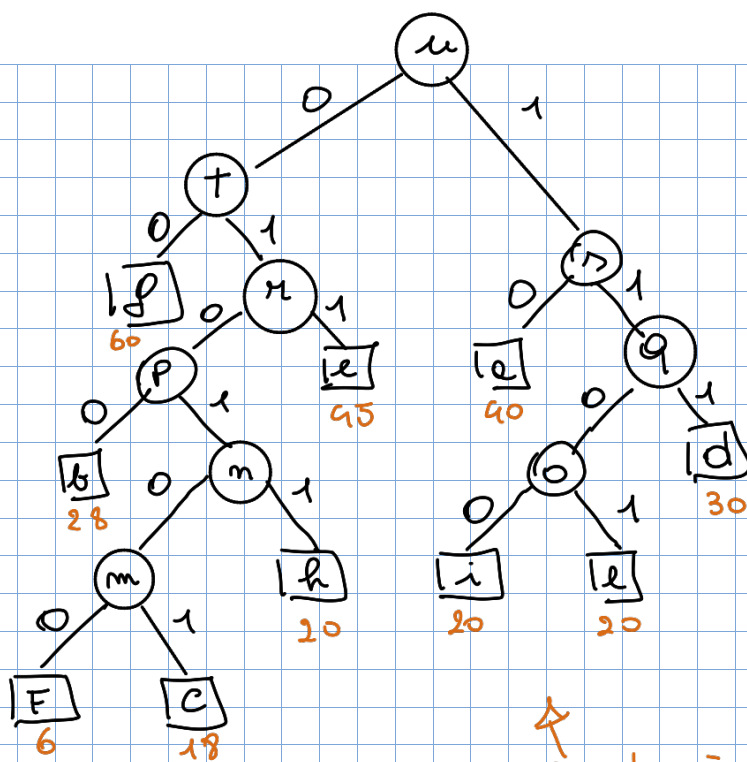
~~q: 20~~

~~r: 30~~

~~s: 40~~

~~t: 60~~

u: 100



↑
questo è l'albero
ottimo

	cod	cod
e: 20	0000	10
b: 4	0001	1100
c: 3	0010	1101
d: 10	0011	000001
e: 15	0100	001
f: 1	0101	000000
g: 30	0110	01
h: 4	0111	111
i: 5	1000	0001
l: 5	1001	00001
	400	298