

perché la dimensione del file compresso è

$$B(T) \cdot \frac{|F|}{100} \rightarrow \text{n° caratteri file non compresso}$$

$\Rightarrow B(T)$, a meno di una costante (100), è ~~il fattore di~~ l'inverso del fattore di compressione: minore è $B(T)$ e migliore è il fattore di compressione.

(Osservazione: stiamo assumendo di dover associare una codeword a ogni singolo carattere, il che non è sempre la cosa migliore da fare:

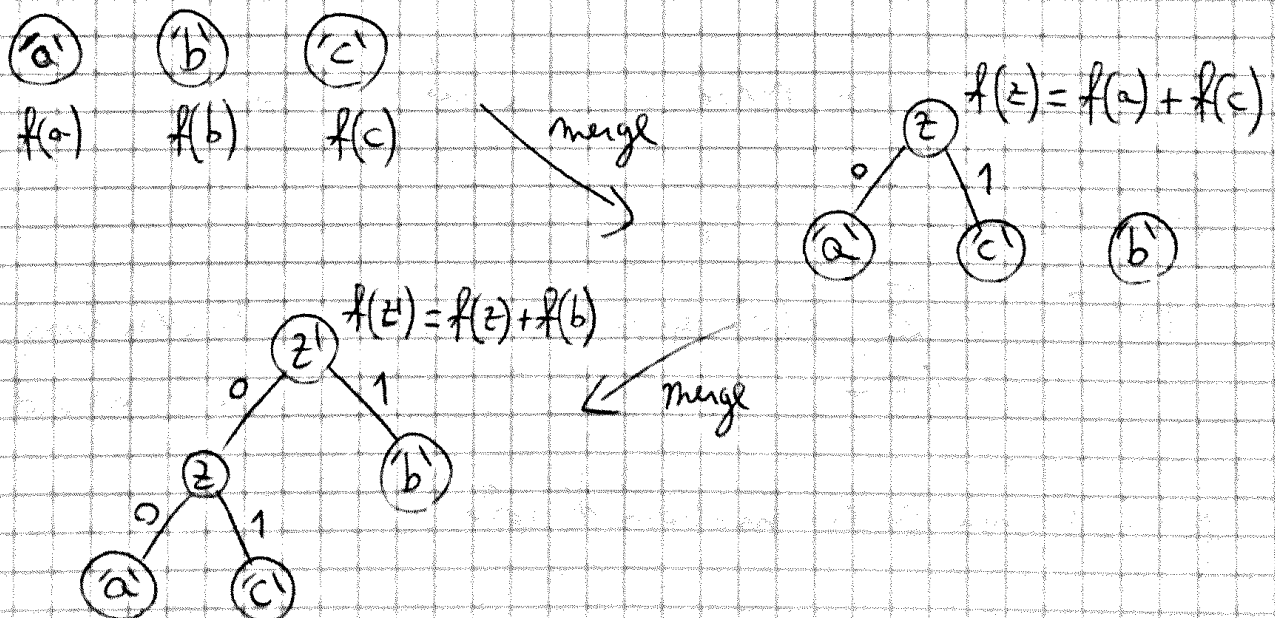
abbabbabbabb...

Ciò, vediamo il codice libero da prefissi ottimo tra tutti i codici che associano una codeword a ogni singolo carattere)

Osservazione: un albero ottimo è sempre pieno (cioè i nodi interni hanno 2 figli)

\Rightarrow spazio sottoproblemi: alberi pieni a n foglie

Goal: trovare un albero ottimo T , cioè che minimizza $B(T)$



sottoproblemi = alberi su meno caratteri

scelta greedy = merge tra nodi che hanno frequenza minore

Q = coda di priorità di nodi con chiave $f(z)$

attributi del nodo z : $left[z]$, $right[z]$, $f(z)$

HUFFMAN(C, f)

$n \leftarrow |C|$

$Q \leftarrow \emptyset$

for each $c \in C$ do

$z \leftarrow \text{NEW_NODE}()$

$f[z] \leftarrow f(c)$

$left[z] \leftarrow \text{null}$

$right[z] \leftarrow \text{null}$

$\text{INSERT}(Q, z)$

for $i \leftarrow 1$ to $n-1$ do

$x \leftarrow \text{EXTRACT_MIN}(Q)$

$y \leftarrow \text{EXTRACT_MIN}(Q)$

$z \leftarrow \text{NEW_NODE}()$

$left[z] \leftarrow x$

$right[z] \leftarrow y$

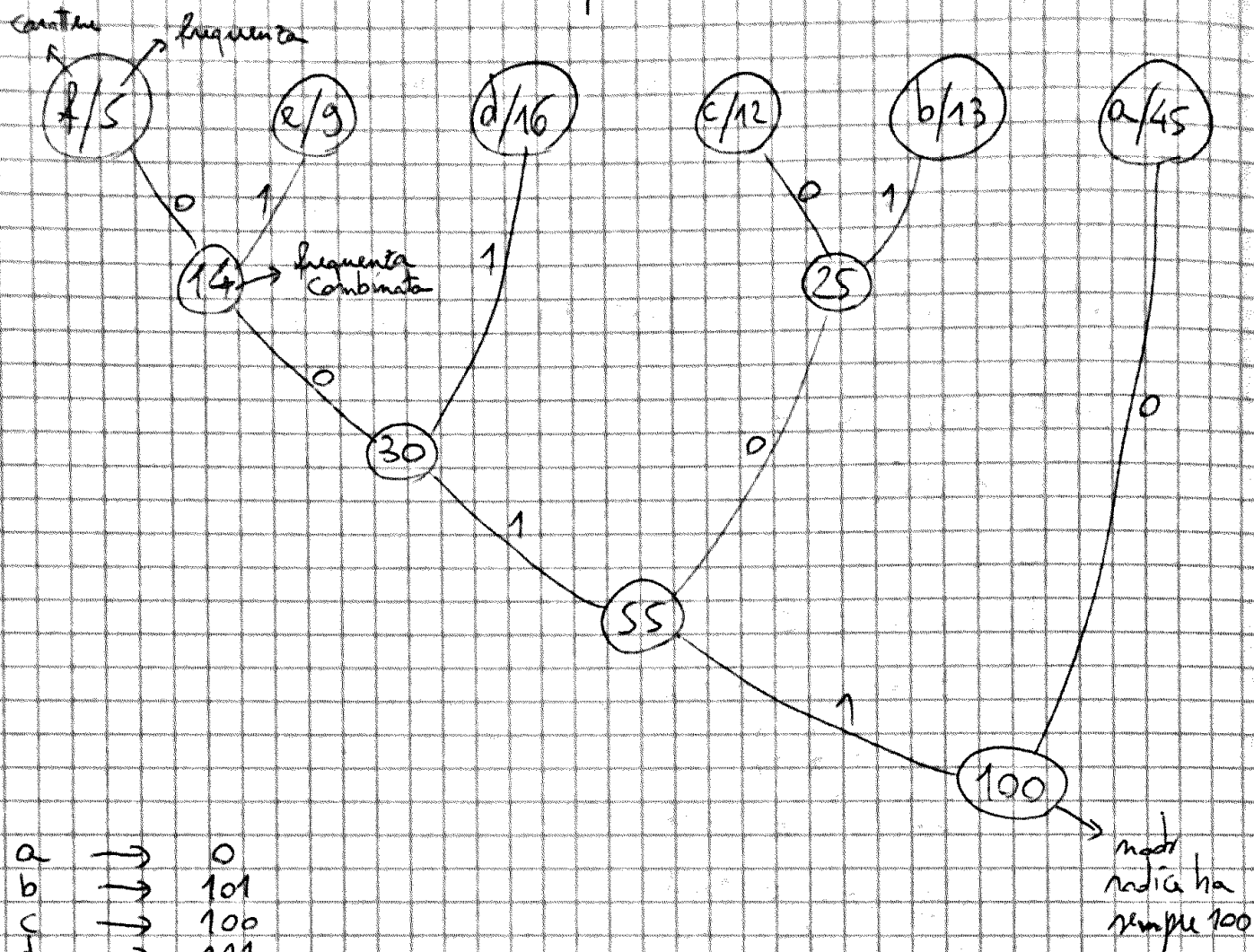
$f[z] \leftarrow f[x] + f[y]$

$\text{INSERT}(Q, z)$

return $\text{EXTRACT_MIN}(Q)$

Complessità: $\sum_{i=1}^{n-1} \Theta(\log(n-i)) = \Theta(n \log n)$

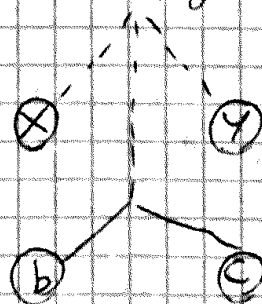
Funzionamento nel nostro esempio:



a	→	0
b	→	101
c	→	100
d	→	111
e	→	1101
f	→	1100

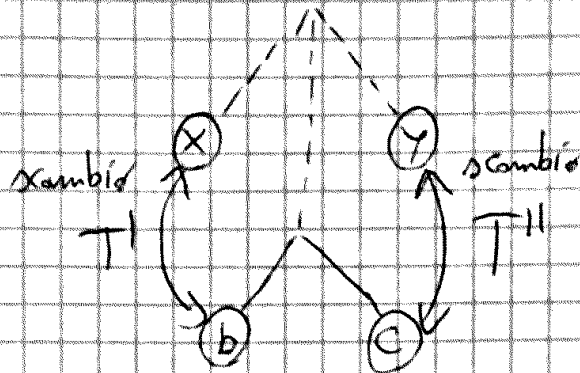
Proprietà di scelta greedy: sia C un alfabeto e siano x e y i caratteri in C di frequenza minore. Allora esiste un codice prefisso ottimo T in cui x e y sono foglie attaccate allo stesso padre.

Dimostrazione: sia T^* una soluzione ottima
siano b e c le foglie di profondità max in T^*



$$1) d_{T^*}(b) = d_{T^*}(c) \geq d_{T^*}(x), d_{T^*}(y)$$

$$2) f(x) \leq f(b) \\ f(y) \leq f(c)$$



$$T^* \rightarrow T' \rightarrow T''$$

deve mostrare che $B(T^*) \geq B(T') \geq B(T'')$

$$T^* \rightarrow T' : B(T^*) - B(T') = \sum_{c \in C} f(c) d_{T^*}(c) +$$

$$- \sum_{c \in C} f(c) d_{T'}(c) = f(b) d_{T^*}(b) +$$

$$+ f(x) d_{T^*}(x) - f(b) d_{T'}(b) - f(x) d_{T'}(x) =$$

$$\text{ma } d_{T'}(b) = d_{T^*}(x), d_{T'}(x) = d_{T^*}(b)$$

$$\Rightarrow = f(b) d_{T^*}(b) + f(x) d_{T^*}(x) - f(b) d_{T^*}(x) -$$

$$f(x) d_{T^*}(b) = f(b) (d_{T^*}(b) - d_{T^*}(x)) -$$

$$f(x) (d_{T^*}(b) - d_{T^*}(x)) =$$

$$= (d_{T^*}(b) - d_{T^*}(x)) (f(b) - f(x))$$

≥ 0

≥ 0

$$\Rightarrow B(T^*) = B(T')$$

✓
e poi lo stesso
per $T' \rightarrow T''$