

CODICI DI HUFFMAN

- CONSENTONO FATTORI DI COMPRESSIONE TRA IL 20% E IL 90%
- PROBLEMA: TROVARE UNA CODIFICA DI UN FILE DI CARATTERI IN MODO DA MINIMIZZARNE LA DIMENSIONE

ESEMPIO: FILE DI 100 CARATTERI

CAR.	FREQ.	COD1 (8 BIT)	COD2 (3 bit)	COD3	
a	45	00000000	000	0	45
b	13	00000001	001	101	39
c	12	00000010	010	100	36
d	16	00000011	011	111	48
e	9	00000100	100	1101	36
f	5	00000101	101	1100	20
100		800 bit	300 bit	224 bit	

$\xrightarrow{\text{Ripartito } 62.5\%}$
 LUNGHEZZA FISSA 72%

$\xrightarrow{\text{Ripartito } 25.3\%}$
 LUNGH. VARIABILE

ES. a b a c

COD2

002001000010

COD3

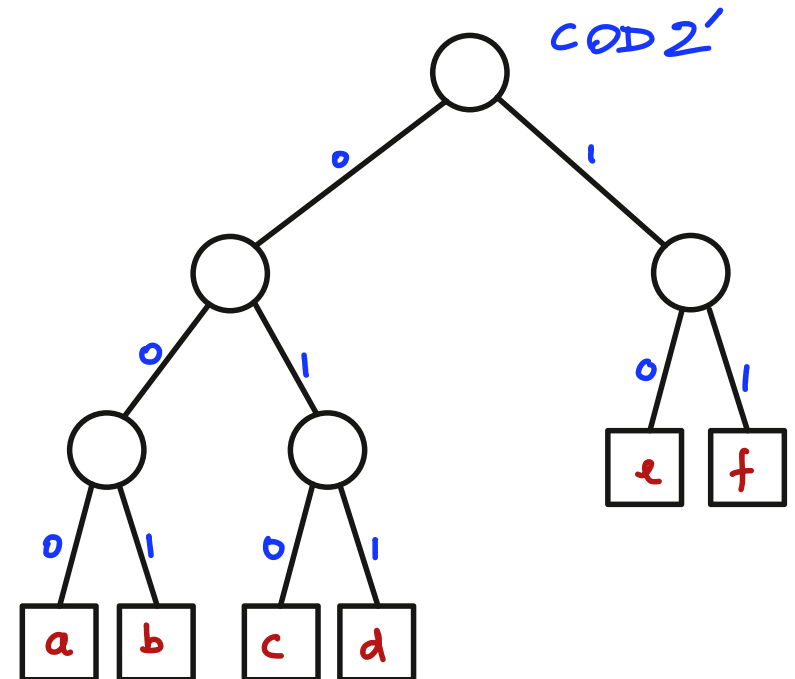
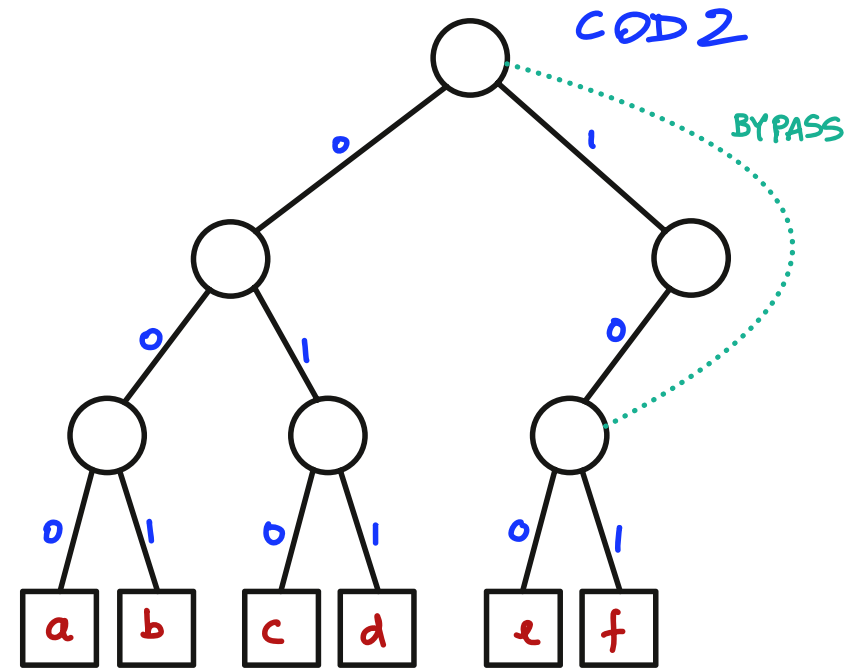
01010100

~25% IN MEMO

ALBERI DI DECODIFICA

CAR.	FREQ.	COD2	COD2'	
a	45	000	000	135
b	13	001	001	39
c	12	010	010	36
d	16	011	011	48
e	9	100	10	18
f	5	101	11	10
		300	286	286

RISPARMIO: $\frac{14}{300} \cdot 100 \approx 4.67\%$



ALBERI DI DECODIFICA

CAR.	FREQ.	COD2 (3 bit)	COD3
a	45	000	0
b	13	001	101
c	12	010	100
d	16	011	111
e	9	100	1101
f	5	101	1100

ES.

a b a c

COD 2

000 001 000 010

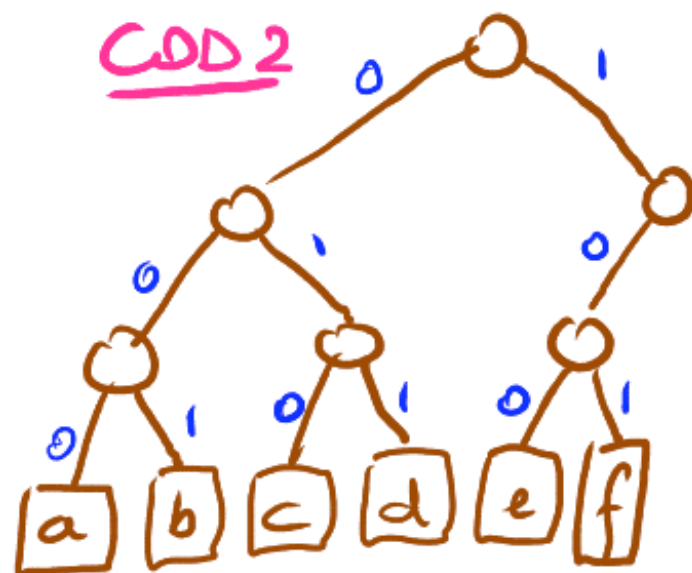
a b a c

COD 3

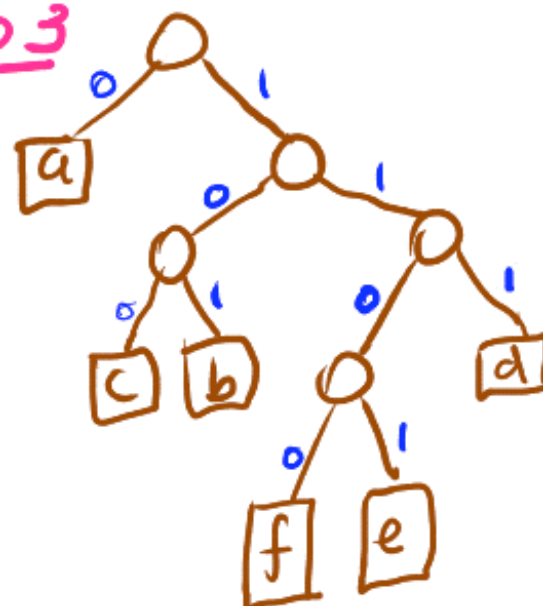
0 101 0100

a b a c

COD 2



COD 3

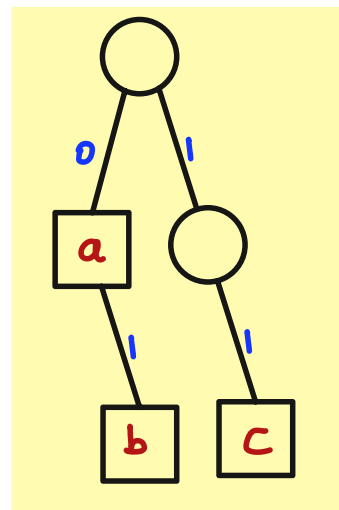


- CODICI PREFISSI: SONO CODICI IN CUI NESSUNA CODIFICA E' PREFISSO DI UN'ALTRA CODIFICA

ESEMPIO DI CODICE NON PREFISSO

a 0
b 01
c 11

0 1 1 1 1 1 1 1 1 1 1 1
b c c c a c



ESEMPIO DI CODICE NON PREFISSO **AMBIGUO**

a 0
b 1
c 01



?



ALBERI DI DECODIFICA

CAR.	FREQ.	COD2 (3 bit)	COD3
a	45	000	0
b	13	001	101
c	12	010	100
d	16	011	111
e	9	100	1101
f	5	101	1100

COSTO DELLA CODIFICA:

$$B(\text{cod}) = \sum_{c \in C} f(c) |\text{cod}(c)|$$

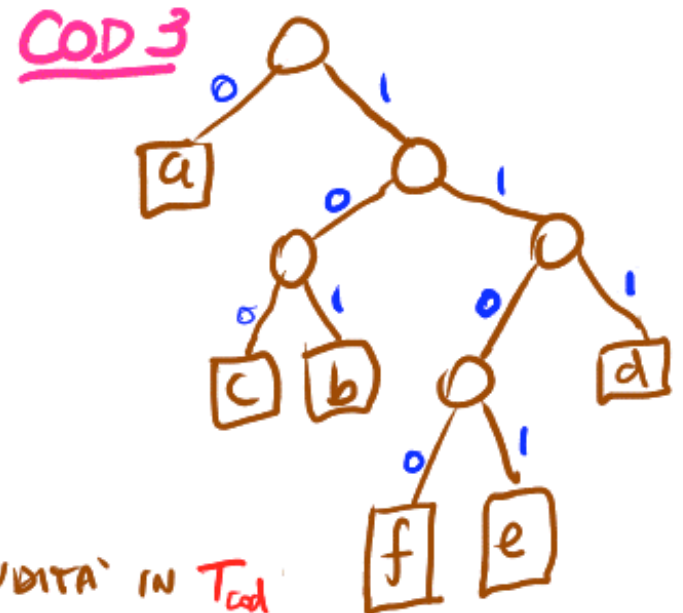
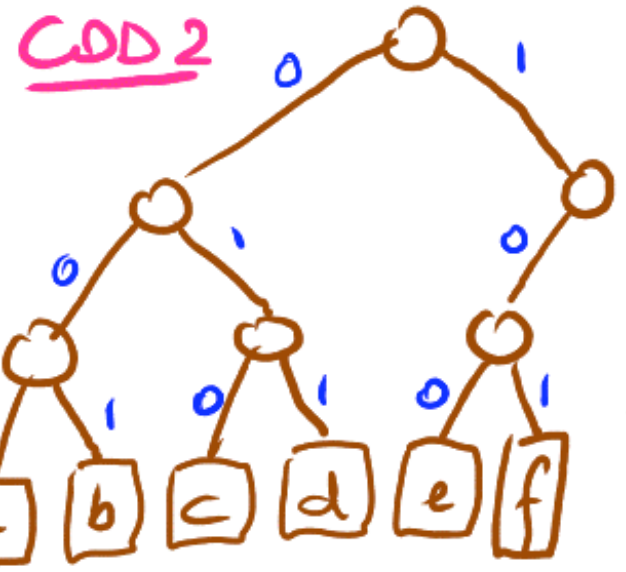
$$= \sum_{c \in C} f(c) d_{T_{\text{cod}}}(c) \stackrel{\text{def}}{=} B(T_{\text{cod}})$$

T_{cod} : ALBERO DI DECODIFICA

C : ALFABETO

$f: C \rightarrow \mathbb{N}$

$d_{T_{\text{cod}}}$: PROFONDITA' IN T_{cod}



$$B(\text{cod}_3) = 45 \cdot 1 + 13 \cdot 3 + 12 \cdot 3 + 16 \cdot 3 + 9 \cdot 4 + 5 \cdot 4 = 45 + 41 \cdot 3 + 14 \cdot 4 = 45 + 123 + 56 = 224$$

COSTO DEGLI ALBERI DI DECODIFICA

- C : ALFABETO
- $f : C \rightarrow \mathbb{N}$: FUNZIONE FREQUENZA
- $\text{cod} : C \rightarrow \{0,1\}^+$: CODIFICA DEI CARATTERI DI C
- T_{cod} : ALBERO DI DECODIFICA
- $B(\text{cod})$: COSTO DELLA CODIFICA DI UN TESTO
NELL'ALFABETO C CON FREQUENZA f
- $|\text{cod}(c)| \quad (c \in C)$: LUNGHEZZA DI $\text{cod}(c)$
- $d_{T_{\text{cod}}}(c) \quad (c \in C)$: PROFONDITA' DELLA FOGLIA DI T_{cod}
CONTENENTE IL CARATTERE c

SI HA:

$$B(\text{cod}) := \sum_{c \in C} f(c) \cdot |\text{cod}(c)| = \sum_{c \in C} f(c) \cdot d_{T_{\text{cod}}}(c)$$

POVIAMO

$$B(T_{\text{cod}}) := \sum_{c \in C} f(c) \cdot d_{T_{\text{cod}}}(c)$$

(COSTO DI T_{cod})

DUNQUE

$$B(T_{\text{cod}}) = B(\text{cod})$$

PROBLEMA: TRA TUTTI GLI ALBERI DI DECODIFICA RELATIVI AD UN SISTEMA (C, f) (DOVE $f: C \rightarrow \mathbb{N}$) DETERMINARE QUELLO DI COSTO MINIMO, CIOE' L'ALBERO BINARIO DI DECODIFICA T TALE CHE

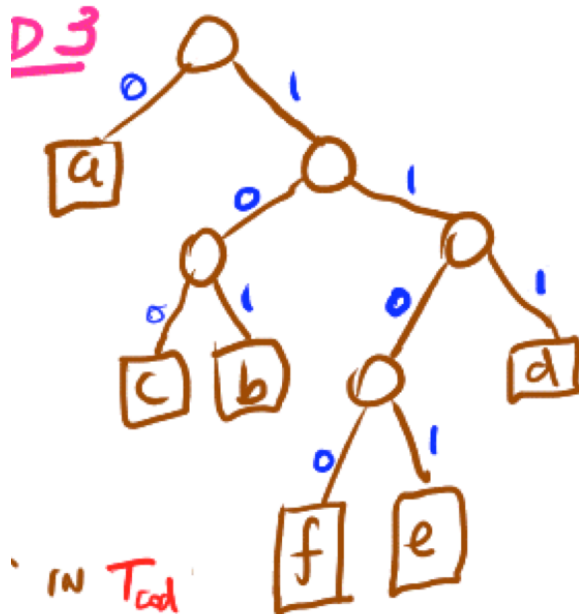
$$B(T) = \sum_{c \in C} f(c) d_T(c)$$

SIA MINIMO

NOTA: CIO' CORRISPONDE A CERCARE UNA CODIFICA DI COSTO MINIMO RELATIVA AL SISTEMA (C, f)

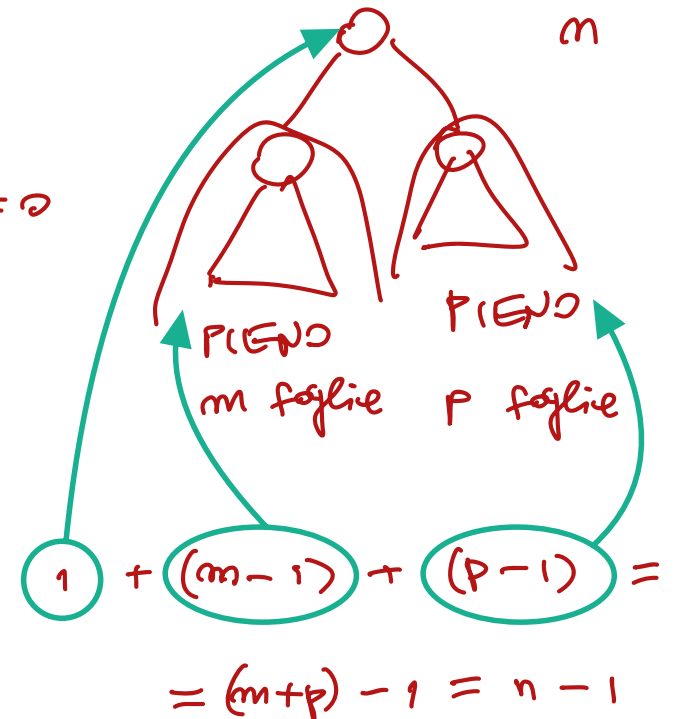
OSSERVAZIONE: POSSIAMO LIMITARE LA NOSTRA RICERCA
 AGLI ALBERI BINARI PIENI, QUELLI CIOE' PRIVI
 DI NODI INTERNI CON UN SOLO FIGLIO.

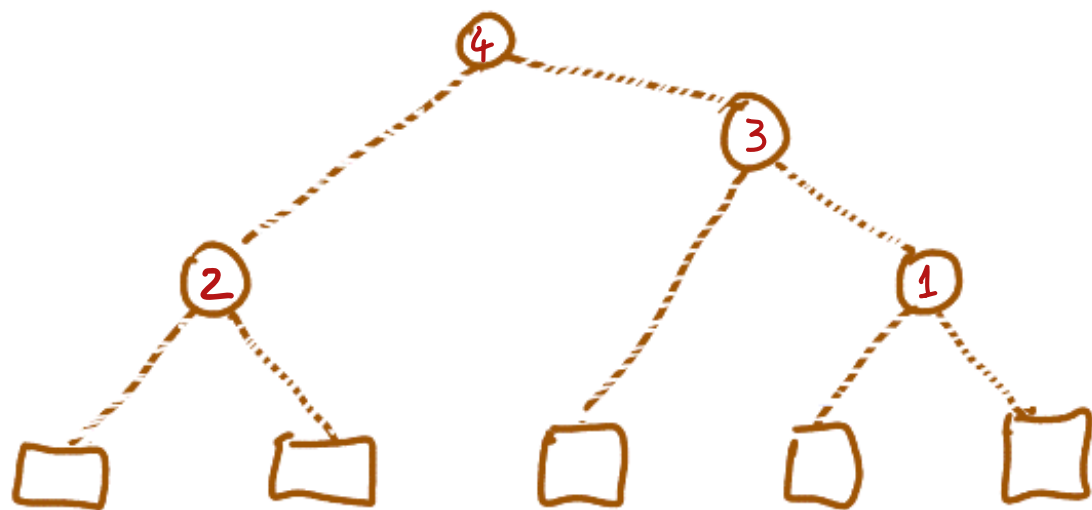
OSSERVAZIONE: IL NUMERO DI NODI INTERNI
 IN UN ALBERO BINARIO PIENO CON
 m FOGLIE E' $m-1$.



$$\square \quad n \approx 1$$

nodi interni = 0

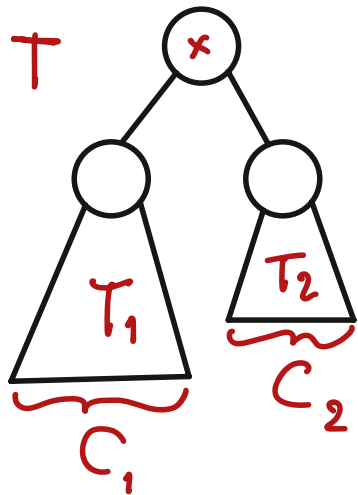




- PER COSTRUIRE UN ALBERO BINARIO PIENO CON m FOGLIE SI POSSONO EFFETTUARE $(m-1)$ OPERAZIONI DI MERGING (O FUSIONE)

COSTO DI UN'OPERAZIONE DI MERGING

DATO UN SISTEMA (C, f) , SIANO $C_1, C_2 \subseteq C$
TALI CHE $C_1 \cap C_2 = \emptyset$ E SIANO T_1 E T_2
ALBERI DI DECODIFICA DI C_1 E C_2 , RISPETT.



- L'ALBERO DI DECODIFICA T E' OTTENUTO DALLA FUSIONE DI T_1 E T_2 CON LA RADICE x
- INDICHIAMO L'OPERAZIONE DI **MERGING** DI T_1 E T_2 CON
 $\text{merging}(T_1, T_2)$ OPPURE $\text{merging}(x)$

COSTO DELL'OPERAZIONE DI MERGING DI T_1 E T_2

$$B(\text{merging}(T_1, T_2)) := \sum_{c \in C_1} f(c) + \sum_{c \in C_2} f(c)$$

SI HA:

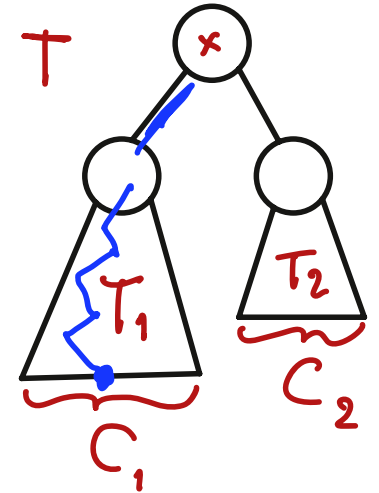
$$B(T) = \sum_{c \in C_1 \cup C_2} f(c) \cdot d_T(c)$$

$$= \sum_{c \in C_1} f(c) \cdot d_T(c) + \sum_{c \in C_2} f(c) \cdot d_T(c)$$

$$= \sum_{c \in C_1} f(c) \cdot (d_{T_1}(c) + 1) + \sum_{c \in C_2} f(c) \cdot (d_{T_2}(c) + 1)$$

$$= \sum_{c \in C_1} f(c) \cdot d_{T_1}(c) + \sum_{c \in C_2} f(c) \cdot d_{T_2}(c) + \sum_{c \in C_1} f(c) + \sum_{c \in C_2} f(c)$$

$$= B(T_1) + B(T_2) + B(\text{merging}(T_1, T_2))$$

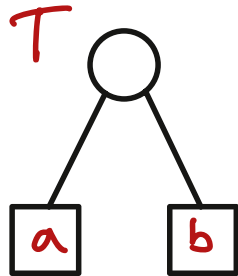


LEMMA IL COSTO BCT DI UN ALBERO DI DECODIFICA T
E' UGUALE ALLA SOMMA DEI COSTI DELLE OPERAZIONI
DI MERGING NECESSARIE A COSTRUIRE T .

DIM. PER INDUZIONE SUL NUMERO m DI NODI
INTERNI DI T .

PONIAMO: $\text{int}(T) :=$ INSIEME DEI NODI INTERNI DI T

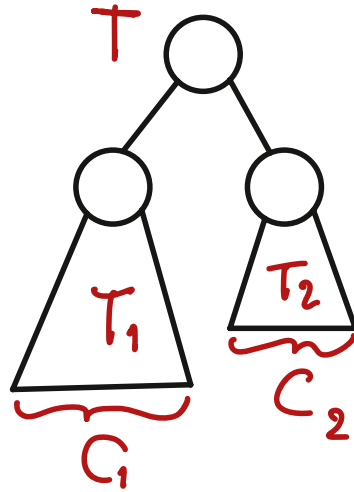
CASO BASE: $m = 1$



$$B(T) = f(a) + f(b) = B(\text{merging}(\text{root}(T)))$$

$$= \sum_{v \in \text{int}(T)} B(\text{merging}(v))$$

PASSO INDUTTIVO:



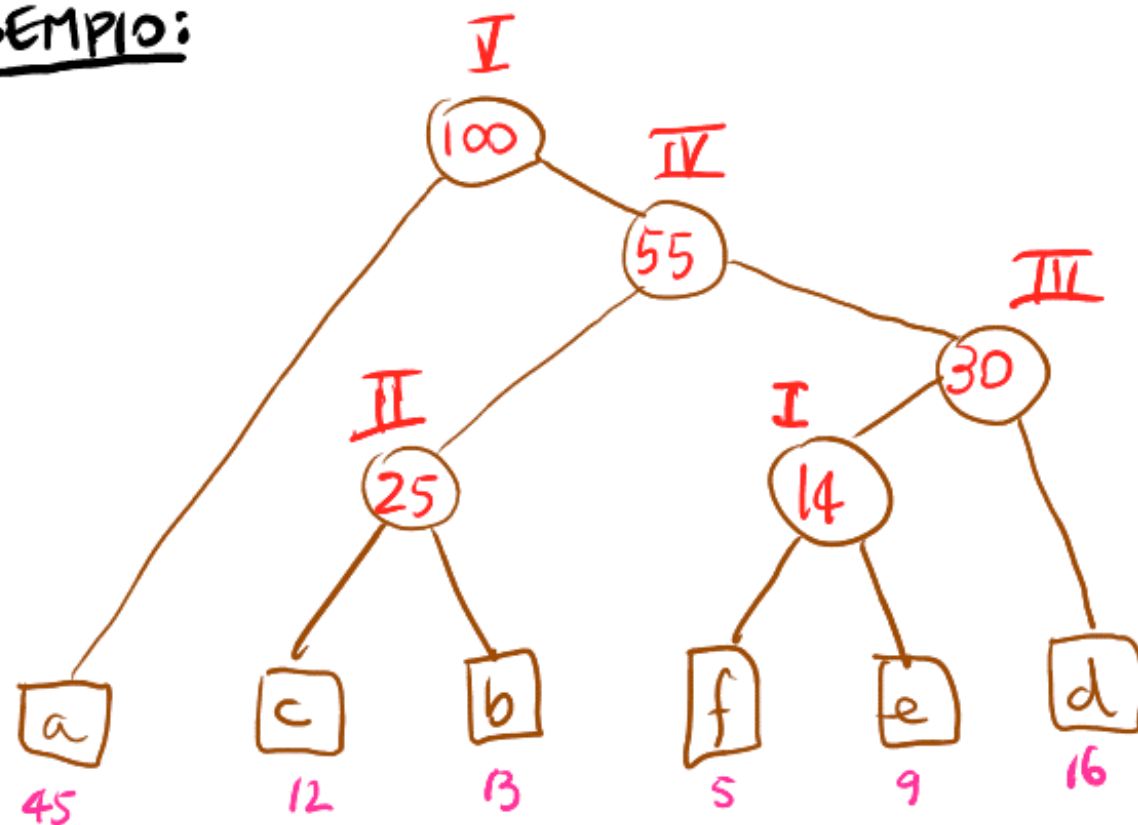
$$B(T) = B(T_1) + B(T_2) + B(\text{merging}(\text{root}(T)))$$

$$= \sum_{v \in \text{int}(T_1)} B(\text{merging}(v)) + \sum_{v \in \text{int}(T_2)} B(\text{merging}(v)) \\ + B(\text{merging}(\text{root}(T)))$$

$$= \sum_{v \in \text{int}(T)} B(\text{merging}(v))$$



ESEMPIO:



$$\begin{array}{r} 14 + \\ 30 + \\ 25 + \\ 55 + \\ \hline 100 \\ 224 \end{array}$$

- UNA POSSIBILE STRATEGIA "GREEDY" PER COSTRUIRE UN ALBERO DI COSTO MINIMO CONSISTE NELL'EFFETTUARE LE OPERAZIONI DI MERGING DI COSTO MINIMO

HUFFMAN(C, f)

$n := |C|$

$Q := \text{make_queue}(C, f)$

for $i := 1$ to $n-1$ do

- SI ALLOCA UN NUOVO NODO INTERNO z

$\text{left}[z] := x := \text{EXTRACT_MIN}(Q)$

$\text{right}[z] := y := \text{EXTRACT_MIN}(Q)$

$f[z] := f[x] + f[y]$

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

return $\text{EXTRACT_MIN}(Q)$

COMPLESSITA'

$(2n-1)$ EXTRACTMIN $O(n \log n)$

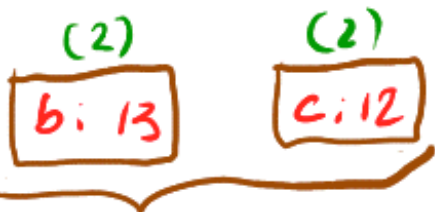
$(n-1)$ INSERT $O(n \log n)$

BUILDHEAP $O(n)$

$O(n \log n)$

ESEMPIO

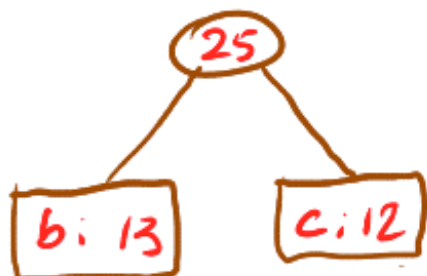
a:45



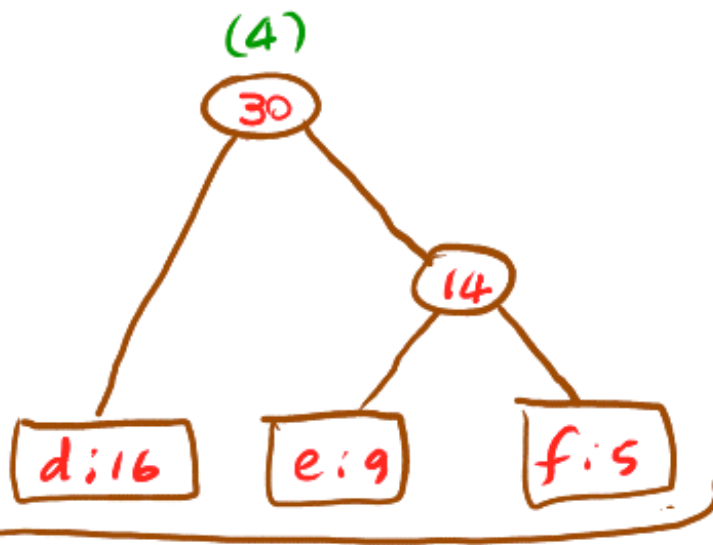
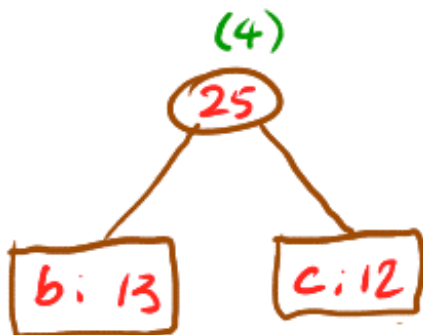
d:16

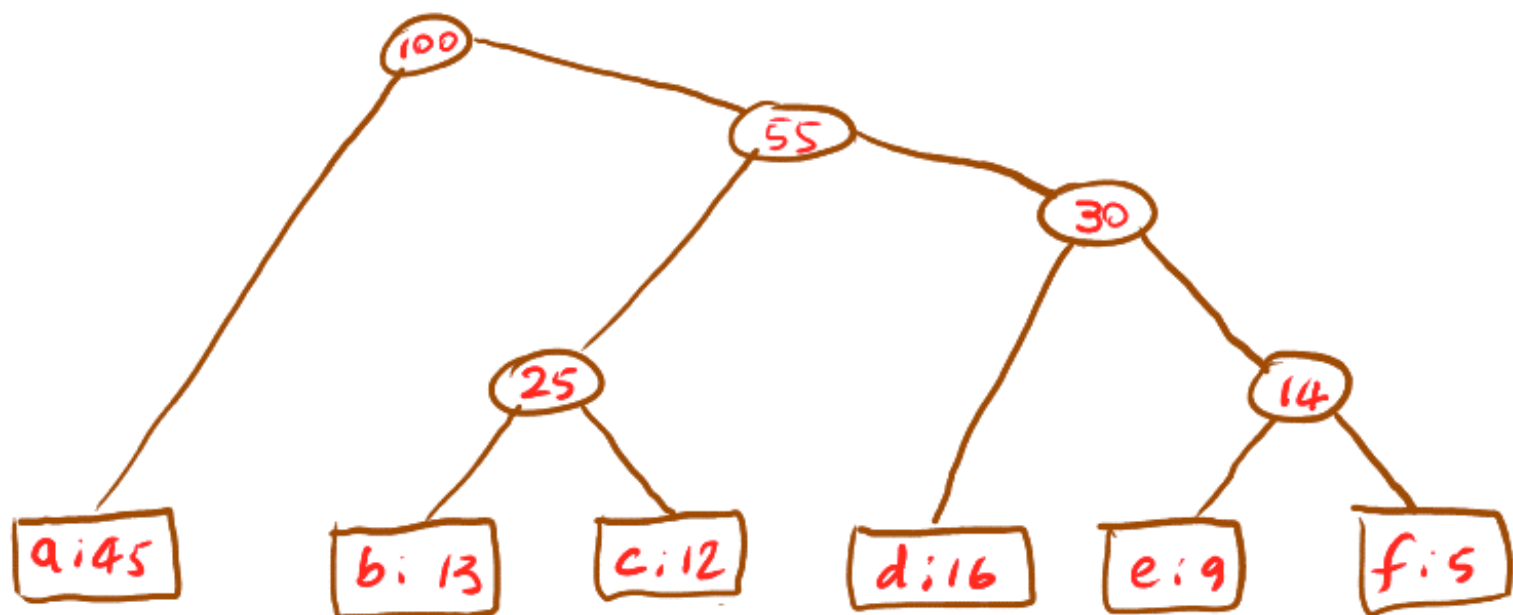
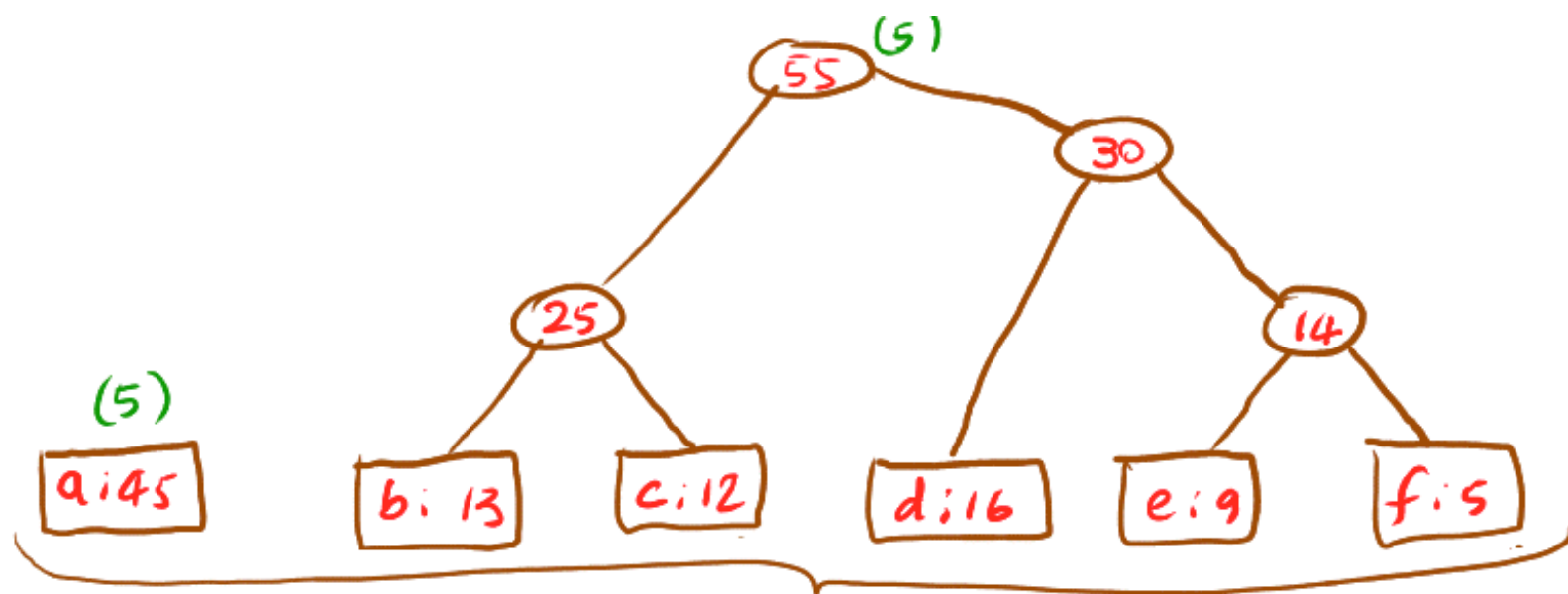


a:45



a:45





CORRETTEZZA DELL'ALGORITMO DI HUFFMAN

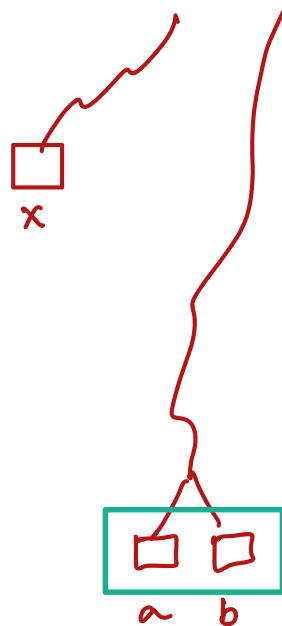
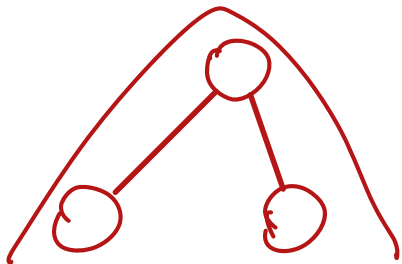
LEMMA

(PROPRIETA' DI SCELTA GREEDY)

SIA C UN ALFABETO ED $f: C \rightarrow \mathbb{N}$ UNA FUNZIONE FREQUENZA.

SIANO x ED y I DUE CARATTERI IN C DI FREQUENZA MINIMA.

ALLORA ESISTE UN CODICE OTTIMO PREFISSO PER C IN CUI LE CODIFICHE DI x ED y DIFFERISCONO SOLO PER L'ULTIMO BIT.



scambio
 $x \leftrightarrow a$
 $y \leftrightarrow b$

$$f(x) \leq f(y)$$

$$f(a) \leq f(b)$$

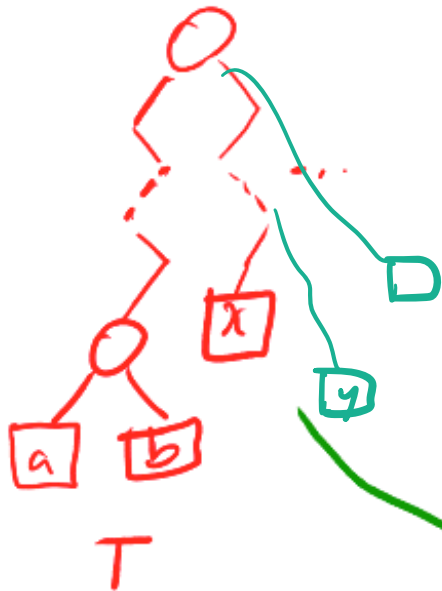
folge a profondita' massima

DIM. SIANO a E b DUE CARATTERI RESIDENTI SU FOGLIE
SORELLE DI PROFONDITA' MASSIMA IN UN ALBERO OTTIMO T .

SUPPONIAMO CHE $f(a) \leq f(b)$ E $f(x) \leq f(y)$.

ALLORA: $f(x) \leq f(a)$ E $f(y) \leq f(b)$.

SIA T' L'ALBERO OTTENUTO DA T SCAMBIANDO
I CARATTERI a ED x .



SI HA:

$$\begin{aligned} B(T) &= \sum_{c \in C} f(c) d_T(c) \\ &= \sum_{c \in C \setminus \{a, x\}} f(c) d_T(c) + f(a) d_T(a) + f(x) d_T(x) \\ &\quad + f(a) d_{T'}(a) + f(x) d_{T'}(x) \\ &\quad - f(a) d_{T'}(a) - f(x) d_{T'}(x) \quad \Big] = 0 \end{aligned}$$

$$\begin{aligned} &= \sum_{c \in C \setminus \{a, x\}} f(c) d_{T'}(c) + f(a) d_T(a) + f(x) d_T(x) \\ &\quad + f(a) d_{T'}(a) + f(x) d_{T'}(x) \\ &\quad - f(a) d_{T'}(a) - f(x) d_{T'}(x) \end{aligned}$$

$$\begin{aligned} &= \sum_{c \in C} f(c) d_{T'}(c) + f(a) d_T(a) + f(x) d_T(x) \\ &\quad - f(a) d_T(x) - f(x) d_T(a) \end{aligned}$$

$$\begin{aligned} &= B(T') + f(a) (d_T(a) - d_T(x)) \\ &\quad - f(x) (d_T(a) - d_T(x)) \end{aligned}$$

$$= B(T') + (f(a) - f(x)) \cdot (d_T(a) - d_T(x))$$

CIOE'

$$B(T) = B(T') + (f(a) - f(x)) \cdot (d_T(a) - d_T(x))$$

E PERTANTO:

$$B(T) - B(T') = \overset{\geq 0}{\boxed{f(a) - f(x)}} \cdot \overset{\geq 0}{\boxed{d_T(a) - d_T(x)}} \geq 0$$

(IN QUANTO $f(a) \geq f(x)$ E $d_T(a) \geq d_T(x)$)

DA CUI

$$B(T) \geq B(T') .$$

- SIA T'' L'ALBERO OTTENUTO DA T' SCAMBIANDO
I CARATTERI b ED y ,

- ANALOGAMENTE A QUANTO VISTO PRIMA, SI HA:

$$B(T') \geq B(T'') \quad (\text{IN QUANTO } f(b) \geq f(y) \text{ E } d_{T'}(b) \geq d_{T'}(y))$$

- PERTANTO: $B(T) \geq B(T'')$

- POICHE' T E' OTTIMO, $B(T'') \geq B(T)$, E QUINDI
 $B(T'')$ E' ANCH'ESSO OTTIMO

- INOLTRE IN T'' I CARATTERI x E y RISIEDONO SU
FOGLIE SORELLE E QUINDI I LORO CODICI DIFFERISCONO
SOLO PER L'ULTIMO BIT. ■

LEMMA (PROPRIETA' DELLA SOTTOSTRUTTURA OTTIMA)

- SIA C UN ALFABETO ED $f: C \rightarrow \mathbb{N}$ UNA FUNZIONE FREQUENZA.
 - SIANO x ED y I DUE CARATTERI IN C DI FREQUENZA MINIMA.
 - SIA $C' = (C \setminus \{x, y\}) \cup \{z\}$, CON $z \notin C$.
 - SIA $f': C' \rightarrow \mathbb{N}$ TALE CHE:
$$f'(c) = \begin{cases} f(c) & \text{SE } c \neq z \\ f(x) + f(y) & \text{SE } c = z \end{cases}$$
 - SIA T' UN ALBERO OTTIMO PER (C', f') .
 - SIA T L'ALBERO OTTENUTO DA T' SOSTITUENDO LA FOGLIA z CON UN NODO INTERNO AVENTE COME FIGLI DUE FOGLIE ETICHETTATE CON x ED y , RISPETTIVAMENTE.
- ALLORA T E' OTTIMO PER (C, f) .

DIMOSTRAZIONE

SI HA:

$$B(T) = \sum_{c \in C} f(c) \cdot d_T(c)$$

$$= \sum_{c \in C \setminus \{x, y\}} f(c) \cdot d_T(c) + f(x) \cdot d_T(x) + f(y) \cdot d_T(y)$$

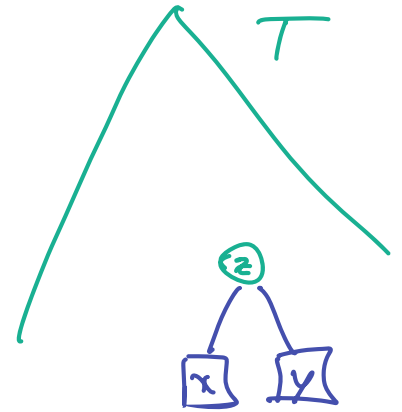
$$= \sum_{c \in C \setminus \{x, y\}} f'(c) \cdot d_{T'}(c) + f(x) \cdot (d_{T'}(z) + 1) + f(y) \cdot (d_{T'}(z) + 1)$$

$$= \sum_{c \in C \setminus \{x, y\}} f'(c) \cdot d_{T'}(c) + (f(x) + f(y)) \cdot d_{T'}(z) + (f(x) + f(y))$$

$$= \sum_{c \in C \setminus \{x, y\}} f'(c) \cdot d_{T'}(c) + f'(z) \cdot d_{T'}(z) + (f(x) + f(y))$$

$$= \sum_{c \in C'} f'(c) \cdot d_{T'}(c) + (f(x) + f(y))$$

$$= B(T') + (f(x) + f(y))$$



PERTANTO: $B(T') = B(T) - (f(x) + f(y))$

- SE T NON FOSSE OTTIMO PER (C, f) , ESISTEREBBE UN ALBERO T'' OTTIMO PER (C, f) TALE CHE:

$$B(T'') < B(T).$$

- GRAZIE AL LEMMA PRECEDENTE, POSSIAMO SUPPORRE CHE x E y SI TROVINO SU FOGLIE SORELLE IN T'' .
- SIA T''' OTTENUTO DA T'' , SOSTITUENDO IL PADRE DI x E y CON UNA FOGLIA z CON FREQUENZA $f(x) + f(y)$.
- ALLORA:
$$\begin{aligned} B(T''') &= B(T'') - f(x) - f(y) \\ &< B(T) - f(x) - f(y) \\ &= B(T') \end{aligned}$$

CONTRADDICENDO L'OTTIMALITA' DI T' PER (C', f') .

- PERTANTO T E' OTTIMO PER (C, f) . ■

"SI SENTIRANO CONTINUI CORI DI

PO PO PO PO PO PO PO PO PO PO PO PO PO PO PO PO

Handwritten list of items and their values:

- ✓ S → 2
- ✓ I → 6
- ✓ L → 5
- ✓ E → 1
- ✓ N → 4
- ✓ T → 2
- ✓ V → 1
- ✓ A → 1
- ✓ O → 17
- ✓ C → 2
- ✓ U → 1
- ✓ R → 1
- ✓ D → 1
- ✓ P → 14

Diagram illustrating groupings and sums:

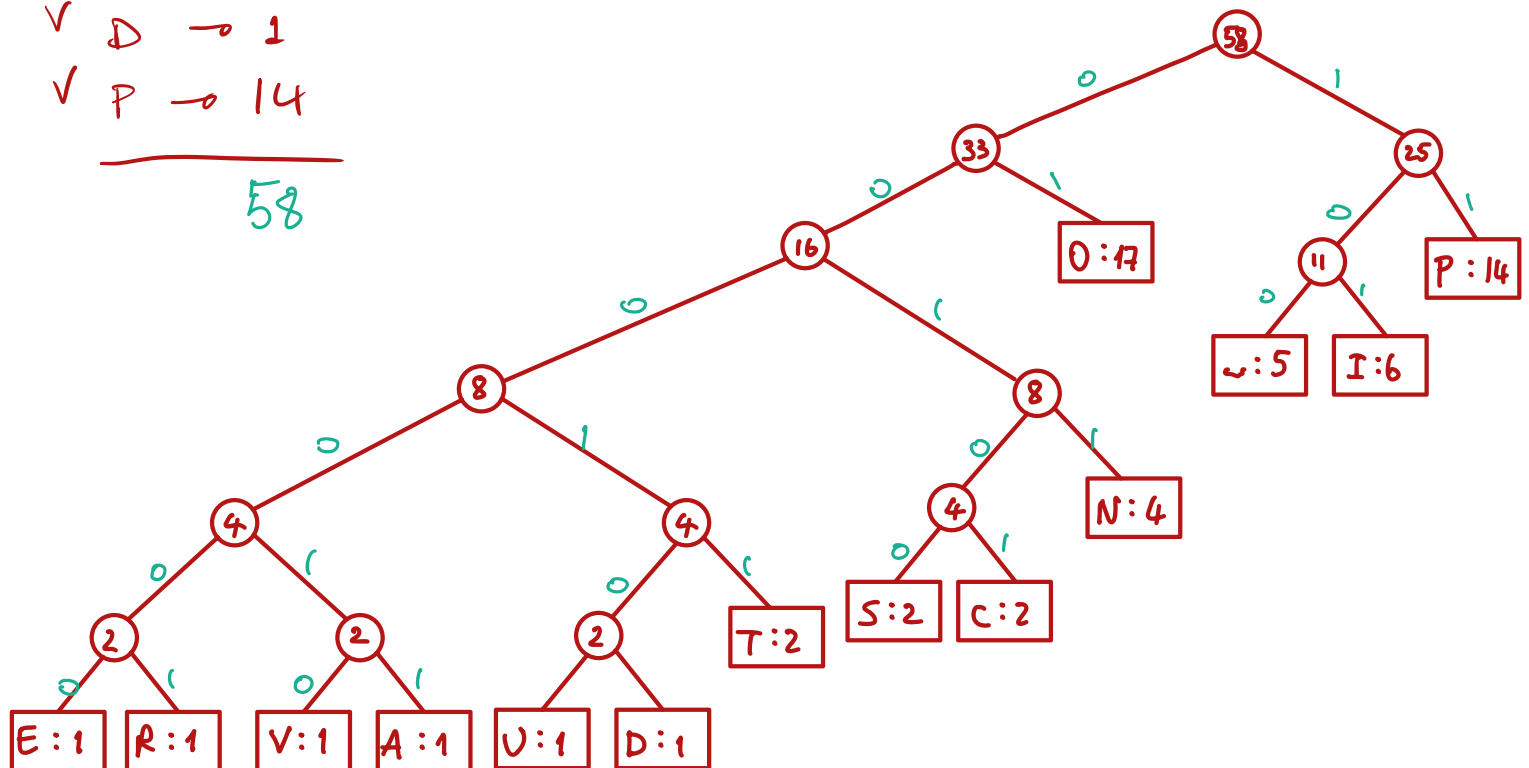
- Items S, I, and T are grouped together, with a sum of 10.
- Items L, E, N, and V are grouped together, with a sum of 10.
- Items A, O, and C are grouped together, with a sum of 20.
- Items U, R, and D are grouped together, with a sum of 18.

58

4 bit x carattere

$$4 \text{ bit} \times 58 = 232 \text{ bit}$$

in una codifica a length. first



$$6 \times 6 = 36$$

$$6 \times 5 = 30$$

$$4 \times 4 = 16$$

$$11 \times 3 = 33$$

$$31 \times 2 = 62$$

177 bit

$$\begin{array}{r} 232 - \\ 177 \\ \hline = 55 \end{array}$$

bit risparmiati su 232

$$\frac{55}{232} \cdot 100 \approx 23.71 \%$$

RIASSUMENDO:

LUNGHEZZA FISSA A 4 BIT	232 bit
CODICE PREFISSO OTTIMO	177 bit
RISPARMIO	55 bit
RISPARMIO PERCENTUALE	~ 23.71 %