

Codifica di Huffman / 23-04

$$H_0 = \log_2 |V_x|$$

dove V_x = # di simboli

quantità di informazione globale

è l'entropia max dei simboli x_i (equiprobabilità)

Detti $C(x) \forall x \in V_x$ una codifica

$$L_c(x) < H_0 \quad \forall x \in V_x$$

dati M bit, posso identificare file
 2^M file grandi: M bit

$$0 \quad 00 \quad 2^{2+1} - 2$$

$$1 \quad 01$$

$$10$$

$$11$$

$$1 \text{Gb} = 2^{30} \quad \text{file grandi: } 1 \text{Gb} = 2^{31} - 2$$

$$\text{n° di file} = V_x$$

$$H_0 = \log_2 (2^{31} - 2) \approx 31$$

$$0 \quad 00 \quad 000$$

$$1 \quad 01 \quad 001$$

$$10 \quad 010$$

$$11 \quad 100$$

$$101$$

$$110$$

$$111$$

tutti i file fino a 3 bit

con una compressione mi aspetto una riduzione
(a 2 bit)

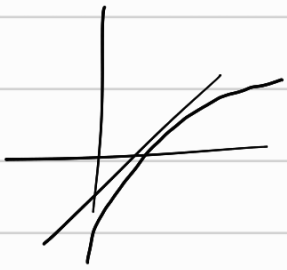
H_0 4 caselle con 2 bit, ma ho 7 caselle
di 3 bit, eppure devo codificare in
tutte caselle diverse (inevitabile)



unque può succedere che comprimendo
molte volte, la dimensione (di un file
già compresso) aumenti

La lunghezza attesa $L(C, V_x)$ non può essere minore di $H(x)$

$$\text{So } \log_2 e(t-1) \geq \log_2(t)$$



la retta sta
sopra il log

Le lunghezze della codifica
sono $L_1 \dots L_{|V_x|}$

Univocamente decifrabile
(file diversi in codifica diverse)

$$\text{Considero } Z = \sum 2^{-L_i}$$

$$q_i = \frac{2^{-L_i}}{Z} \quad \text{ciascuna termine della somma}$$

$Z \rightarrow$ sta normalizzando

$$\sum q_i = 1 \quad (\text{somma / somma})$$

$$\log_2 q_i = \log_2 2^{-L_i} - \log_2 Z$$

$$\log_2 q_i = -L_i - \log_2 Z$$

$$\text{per cui } L_i = -\log_2 q_i - \log_2 Z = \log \frac{1}{q_i} - \log Z = L_i$$

lunghezza
di ogni codifica

$$\Phi: \begin{matrix} q_1, \dots, q_{|V_x|} \\ p_1, \dots, p_{|V_x|} \end{matrix} \quad \text{def } t = \frac{q_i}{p_i}$$

$$\underbrace{\sum_i p_i \log \frac{1}{p_i}}_{H} - \sum_i p_i \log \frac{1}{q_i} = \sum_i p_i \left(\log \frac{q_i}{p_i} \right)$$
$$= \sum_i p_i \log_2 t_i$$

$$\sum p_i \log_2 t_i \leq \sum p_i (t_i - 1) \quad \text{per la concavità del log di cui sopra}$$

$$\quad \quad \quad \parallel \sum q_i - \sum 1$$

$$\quad \quad \quad 1 - 1 = 0$$

$$\downarrow$$

$$= 0 \quad \text{se } q_i = p_i$$

$$L(C, V_x) = \sum p_i L_i = \sum p_i (\log_2 \frac{1}{q_i} - \log_2 Z)$$

$$L_i = \log_2 \frac{1}{p_i} - \log_2 Z$$

$$\geq H(X)$$

VINCOLO INVARIABILE

se vado sotto l'entropia
qualcos'altro deve aumentare
(la lunghezza della sequenza)
x rispettare l'uguaglianza

Codifica di Huffman

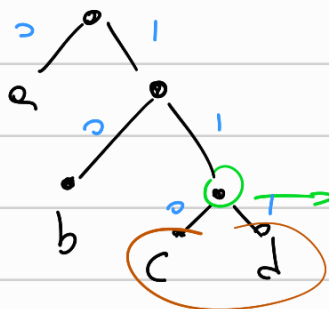
Costruisce un albero

STEP 1

STEP 2

STEP 3

X	p_i
a	$\frac{1}{2}$
b	$\frac{1}{4}$
c	$\frac{1}{8}$
d	$\frac{1}{8}$



a	$\frac{1}{2}$	a	$\frac{1}{2}$	(a, (b, (c, d)))	1
b	$\frac{1}{4}$	(b, (c, d))	$\frac{1}{2}$		
(c, d)	$\frac{1}{4}$				

$$p = \sum p_i \log_2 p_i$$

lung: meno
probabili

↓
possono
essere
lettere, numeri,
file

Ragionata attraverso una coda (estraggo man mano i nodi)

Se passo una **codifica** di 0 e 1 ottengo * dopo aver costruito l'albero
che è univoca e istantanea

Come input necessario di probabilità

V_x e p_i con $i = 1 \dots |V_x|$

Huffman è, per simbolo, la codifica ottimale



prende un simbolo alla volta

$\frac{1}{2}$	a	0
$\frac{1}{4}$	b	10
$\frac{1}{8}$	c	101
$\frac{1}{8}$	d	111

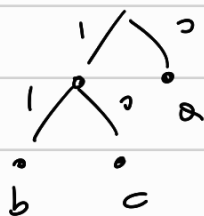
$$a b a c d = \underbrace{0}_{a} \underbrace{10}_{b} \underbrace{0}_{a} \underbrace{101}_{c} \underbrace{111}_{d}$$

analisi simbolo
x simbolo

$$\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{2}{8} \cdot 3 = 1,75$$

↓
lung. cod

$\frac{1}{2}$	a
$\frac{1}{4}$	b
$\frac{1}{4}$	c



a → 0
b → 11
c → 10

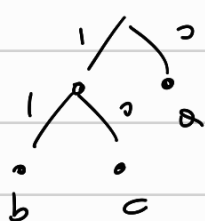
odi
simboli
↑

$$H = \log_2 V_x$$

$$H(x) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 2 = 1,5 \quad \text{lunghezza media}$$

↑
 $H = 1,58$, posso
comprimere ancora
(anche se di poco) *

$\frac{1}{3}$	a	0
$\frac{1}{3}$	b	11
$\frac{1}{3}$	c	10



$$H(x) = \frac{1}{3} + \frac{2}{3} + \frac{2}{3} = \frac{5}{3} = 1,6$$

$$H = 1,58 \quad \longleftrightarrow \quad \bar{L} \approx 1,67$$

ho un 10% di spreco,

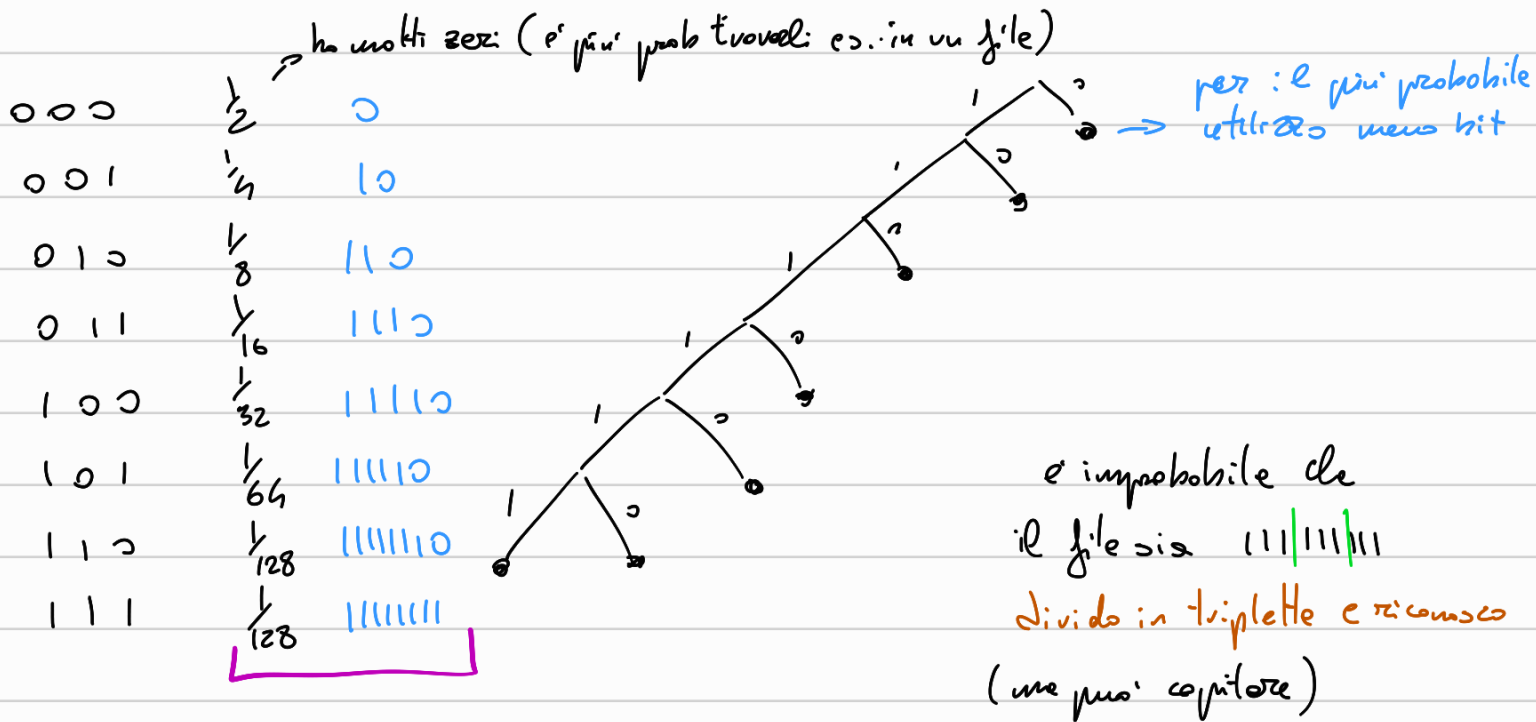
potrei comprimere ancora

ma non ho potenza di due più

Posso comprimere fino
all'entropia
(grezza) *

Quando ha molti simboli: (e quindi analizzare un simbolo per volta)

Huffman logica



8 simboli

3 bit

Quanti zeri?

$$\frac{3}{\text{zeri}} \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} + 2 \cdot \frac{1}{8} + \frac{1}{16} + \frac{2}{32} + \frac{1}{64} + \frac{1}{128} = 4$$

mi aspetto che questo numero sia $>$ n° di uni

perché e' più probabile che io abbia più zeri ($p = \frac{1}{2}$ per 000)

$4 \approx 2$ che sono più o meno il doppio degli uni

→ n° di zeri sulle codifiche: $\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \frac{1}{32} + \frac{1}{64} + \frac{1}{128}$

Debolezza di Huffman

- bisogna conoscere le probabilità a priori

- irrealistico il fatto che le probabilità siano

identicamente distribuite e che le var. casuali siano indipendenti

- Huffman si basa su scelte binarie; nel caso di sequenze di un solo simbolo, in cui l'entropia è vicina a zero (dove p di uno dei due risultati è circa 1) Huffman non comprime bene, determinando una L media lontana dall'entropia ottanta