

# Codifica aritmetica / 05-05

Una debolezza di Huffman è sapere la distribuzione di probabilità a priori. Tale codifica si adatte invece sulla frequenza dei simboli presenti, applicandosi su stream di byte. Chiede però pochi valori si analizzano

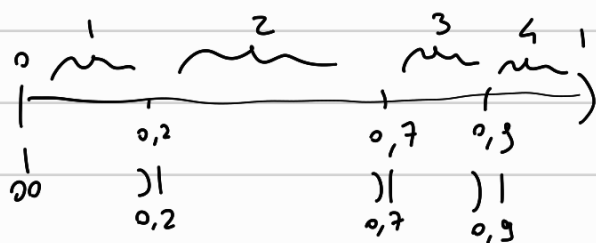
Dati quattro valori  $x_1=1$   $x_2=2$   $x_3=3$   $x_4=4$  e 321124  
 $p_1=0,2$   $p_2=0,5$   $p_3=0,2$   $p_4=0,1$

$$p(x|6) = p(3)p(2)p(1)p(1)p(2)p(4)$$

↓  
n° simboli

↓ se ho tanti simboli improbabili, anche  $p(x|6)$  sarebbe piccolo

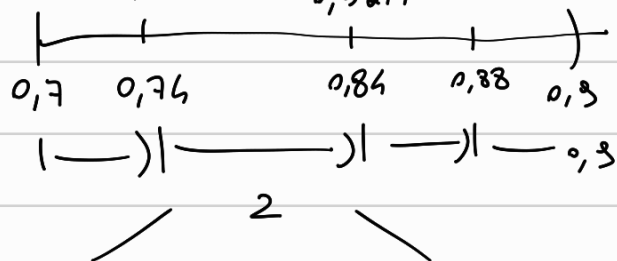
L'intervallo sarà compreso fra  $[0, 1)$



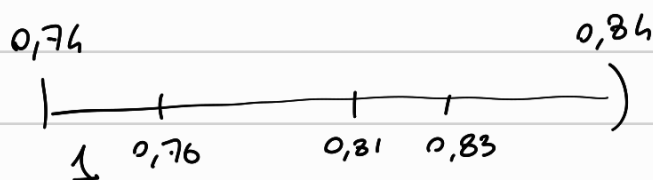
distribuendo le probabilità nell'intervallo

Analizzo la sequenza, ho un 3 quindi ci si aspetta che cada fra 0,7 e 0,9

→ di  $0,9 - 0,7 = 0,2$  nuovo intervallo  
 $0,5$  di  $0,2$



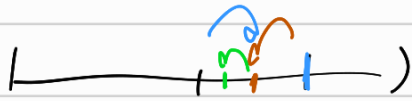
distribuisco le prob. precedenti ridimensionando su questo intervallo



e continuo: finire non finisco non conosco gli estremi

[ 0,742607428 )

l'idea è l'approccio binario (maggiore o minore dello zero?)



0.101111000100

0,7426758125

risparmiare gli zeri dopo la virgola 12 simboli  
codifica aritmetica

$$H = \sum p_i \log \frac{1}{p_i} = 1,76 \quad \text{per sei caratteri} = 10,56 \quad \text{vicino a 12}$$

avere lunghezza media di caratteri  
che mi aspetto

→ migliorare la precisione con seq più lunghe

$\Phi_k$  è l'intervallo di iterazione  $k$

$\Phi_0 = [0, 1)$  iniziale

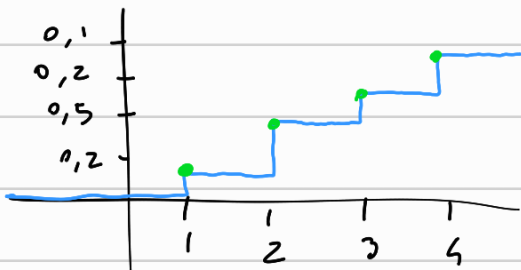
$\Phi_N$  la meta dell'esplorazione

$\Delta_k = \beta_k - \alpha_k$   
empirico intervallo

$\Phi_k = [\alpha_k, \beta_k)$

cdf

• il valore è assunto in quel  $x$



## Algoritmo

(stringhe di numeri naturali, la cdf è su dei numeri)

for  $k = 1 \dots N$

$$\mathbb{D}_k [\alpha_k, \beta_k) = [\alpha_{k-1} + d_{k-1} \text{cdf}(x_{k-1}), \alpha_{k-1} + d_k \text{cdf}(x_k))$$

Iniziali  $\mathbb{D}_0 = [0, 1)$   $d_0 = 1$

$$\mathbb{D}_1 = [0 + 1 \cdot 0,7, 0 + 1 \cdot 0,3)$$

↳ il 1° simbolo della seq. è 3

$$\dots [0,7426, 0,7428)$$

## Correttezza

$$\log_2 \frac{1}{p(x|6)} = \log_2 \frac{1}{p(3)} + \log_2 \frac{1}{p(2)} + \dots$$

$$\log_2 \frac{1}{p(x|N)} = (\cancel{x} 3) \log \frac{1}{p(3)} + (\cancel{x} 4) \log \frac{1}{p(4)} + (\cancel{x} 1) \log \frac{1}{p(1)} + (\cancel{x} 2) \log \frac{1}{p(2)}$$

$$\frac{\log_2 \frac{1}{p(x|N)}}{N} \approx p(3) \log \frac{1}{p(3)} + p(1) \log \frac{1}{p(1)} + \dots = H(X)$$

quindi codice ottimale

Decodifica → devo sapere sempre quanti valori ho

$\mathcal{D} = \mathcal{D}_0$  rappresentazione delle codifica

for  $k = 1 \dots N$

$$x_k = \{x \in V_x \text{ cdf}(x_{k-1}) < \mathcal{D}_{k-1} < \text{cdf}(x)\}$$

$$\mathcal{D}_k = \frac{\mathcal{D}_{k-1} - \text{cdf}(x_{k-1})}{p(x_k)}$$

va a vedere 0,74267578125 dare cade : fra 2 e 3

$X_k = 3$  al primo giro

Non è necessario sapere la probabilità a priori

0 1  
01 | 1100100001....  
b<sub>1</sub> b<sub>2</sub> b<sub>3</sub>...

Codifica binaria  $p_0(0) = p_0(1) = \frac{1}{2}$

$$p_k(0) = p_{k-1}(0) + \frac{1}{k+2} ((1-b_k) - p_{k-1}(0))$$

$$p_k(1) = p_{k-1}(1) + \frac{1}{k+2} (b_k - p_{k-1}(1))$$

(si basa sullo stesso concetto di media  $m_{11} = m_{10} + \frac{1}{11}(24 - m_{10})$   
 $\rightarrow \frac{27 \cdot 10 + 24}{11}$ )

Tale processo viene eseguito a runtime, senza attesa