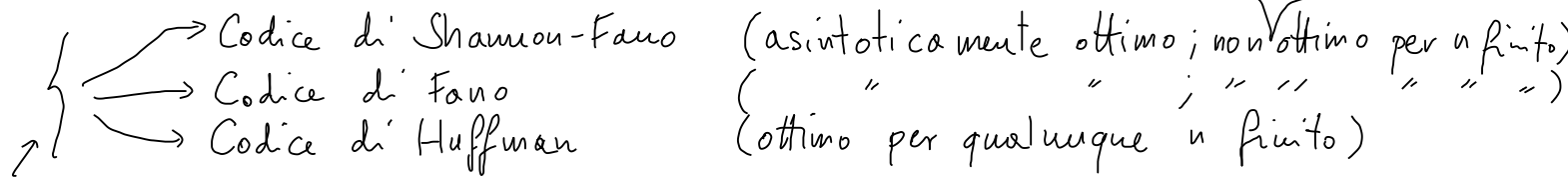
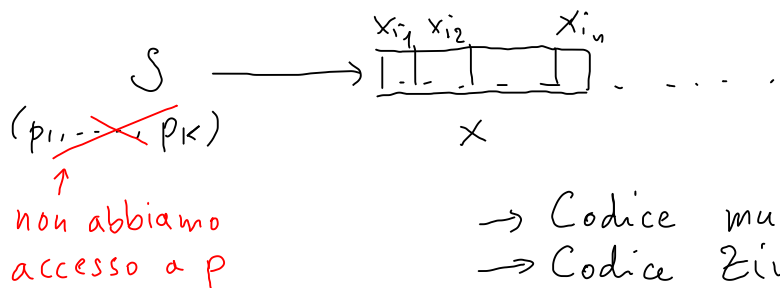


CODIFICA DI SORGENTE

Sorgente S
(d.d.p. (p_1, p_2, \dots, p_k))



↑ Codici basati sulla descrizione statistica della sorgente (la d.d.p. (p_1, \dots, p_k))



Codifica universale: asintoticamente ottima ma indipendente dalla d.d.p. della sorgente

- Codice multinomiale
- Codice Ziv-Leempel (LZ77, LZ78) (alla base del formato di compressione zip)

Avere codici efficienti senza conoscere la d.d.p. della sorgente S

=

$$A = \{x_1, \dots, x_k\}$$

↓ ↓

n_1 n_k

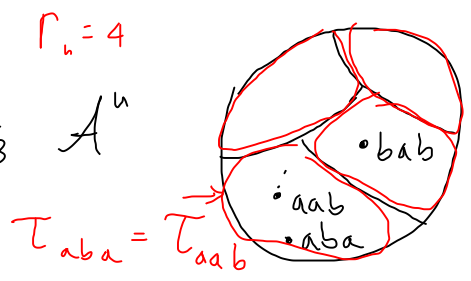
Una stringa $x = (x_{i_1}, x_{i_2}, \dots, x_{i_n})$ ha un tipo esatto

$T_x = (n_1, n_2, \dots, n_k)$ dove n_i è il numero di occorrenze del simbolo x_i nella stringa x

Es. $A = \{a, b\}$, $x = aba \rightarrow T_x = (2, 1)$

$x = aab \rightarrow T_x = (2, 1)$

$(0,3) \rightarrow 0$
 $(3,0) \rightarrow 1$
 $(2,1) \rightarrow 2$
 $(1,2) \rightarrow 3$
 $n=3$



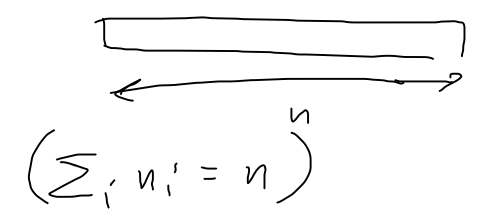
A^n è partizionato in classi di equivalenza date dai tipi delle sequenze.
 - Il numero di classi di equivalenza è indicato con Γ_n

Osservazione : $\Gamma_n \leq (n+1)^K$ numero di K -ple di interi con valore tra 0 e n

Inoltre sia T_x l'insieme delle sequenze in A^n che hanno lo stesso tipo di x

$$\begin{aligned}
 & (\cdot, \cdot, \cdot, \cdot) \\
 T_x &= (n_1, n_2, \dots, n_K) \\
 & \downarrow \\
 & 0 \leq n_1 \leq n \quad n+1 \text{ m.o.} \\
 & 0 \leq n_i \leq n
 \end{aligned}$$

$$|T_x| = \binom{n}{n_1, n_2, \dots, n_K} = \frac{n!}{n_1! n_2! \dots n_K!} \quad (\text{Coefficiente multinomiale})$$



Codifica multinomiale : per codificare $x \in A^n$, usa la parola a lunghezza variabile

$$\varphi(x) = \varphi_p(T_x) * \varphi_s(x)$$

\uparrow
 prefisso che dipende solo dal tipo di x

\uparrow
 suffisso (dipende da x)

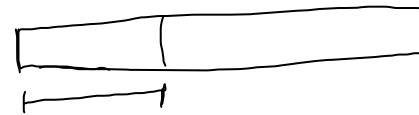
- $\varphi_p(T_x)$ è il numero d'ordine di T_x in un elenco lessicografico di tutti i tipi (identifica il tipo).
- $\varphi_s(x)$ è il numero d'ordine di x

in un elenco lessicografico dell'insieme T_x di tutte le sequenze con lo stesso tipo di x (identifica x se già conosco T_x)

- 1 (0, 3)
- 2 (1, 2)
- 3 (2, 1)
- 4 (3, 0)

Parola di codice:

$$w = \underbrace{\varphi_p(T_x)}_{\text{lunghezza costante}} * \underbrace{\varphi_s(x)}_{\text{lungh. variabile}}$$



$$A^n \xrightarrow{\varphi} B^+$$

Sia $l_p = |\varphi_p(T_x)|$ la lungh. del prefisso (costante)
 $l_s(x) = |\varphi_s(x)|$ la lungh. del suffisso (variabile)

$$A^n \xrightarrow{\varphi_p} B^{l_p}$$

$|A| = K$

$|B| = D$

$|B^{l_p}| = D^{l_p}$

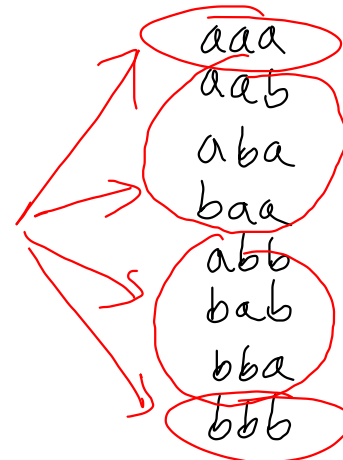
Devo avere $D^{l_p} \geq \Gamma_n \rightarrow l_p = \lceil \log_D \Gamma_n \rceil \leq \lceil \log_D (u+1)^k \rceil \leq \lceil K \log_D (u+1) \rceil$

$$D^{l_s(x)} \geq |T_x| \rightarrow l_s(x) = \lceil \log_D |T_x| \rceil$$

Esempio. $A = \{a, b\}$. $B = \{0, 1\}$ $D = 2$ $K = 2$ $n = 3$

$$\Gamma_n = 4 : (0, 3), (1, 2), (2, 1), (3, 0) ; l_p = \lceil \log_2 4 \rceil = 2$$

Tipo	n-pla	$\varphi_p(T_x)$	$l_s(x)$	$\varphi_s(x)$	$\varphi(x) = \varphi_p(T_x) * \varphi_s(x)$
0	(0,3) } bbb	00	$\lceil \log_2 1 \rceil = 0$	—	00
1	} <u>abb</u> ⁰ <u>bab</u> ¹ <u>bba</u> ²	01	$\lceil \log_2 3 \rceil = 2$	00	0100
		01	" 2	01	0101
		01	" 2	10	0110
2	} baa aba aab	10	$\lceil \log_2 3 \rceil = 2$	00	1000
		10	"	01	1001
		10	"	10	1010
3	(3,0) } aaa	11	$\lceil \log_2 1 \rceil = 0$	—	11



Cosa succede per $n \rightarrow \infty$?

$$\mathbb{E}[L^{(n)}] = \sum_{x \in A^n} p(x) |q(x)| = \sum_{x \in A^n} p(x) (l_p + l_s(x)) =$$

lunghezza
parole di
codice

$$= l_p \sum_{x \in A^n} p(x) + \sum_{x \in A^n} p(x) l_s(x)$$

$$= \lceil \log_D \Gamma_n \rceil + \sum_{j=1}^{\Gamma_n} \underbrace{\sum_{x \in \mathcal{Z}_j} p(x)}_{\text{non dipende da } x} \lceil \log_D |\mathcal{Z}_j| \rceil$$

$$= \lceil \log_D \Gamma_n \rceil + \sum_{j=1}^{\Gamma_n} \lceil \log_D |\mathcal{Z}_j| \rceil \Pr(\mathcal{Z}_j)$$

Prob. complessiva
delle sequenze
di tipo \mathcal{Z}_j

$$(\Pr[X \in \mathcal{Z}_j])$$

!!

$$\Pr(\mathcal{Z}_j)$$

$$(\lceil x \rceil < 1+x) \rightarrow < 1 + \log_D \Gamma_n + \sum_{j=1}^{\Gamma_n} (1 + \log_D |\mathcal{Z}_j|) \Pr(\mathcal{Z}_j)$$

$$= 1 + \log_D \Gamma_n + 1 + \boxed{\sum_{j=1}^{\Gamma_n} \Pr(\mathcal{Z}_j) \log |\mathcal{Z}_j|}$$

Consideriamo la v.a. Z (intera) che rappresenta il numero d'ordine del tipo \mathcal{Z}_j nella lista dei tipi

Consideriamo la seguente entropia condizionata:

$$H(X|Z) = \sum_{j=1}^{\Gamma_n} \Pr(Z=j) H(X|Z=j) = \sum_{j=1}^{\Gamma_n} \underbrace{\Pr(\tau_j)}_{= \Pr(Z=j)} H(X|Z=j)$$

Tutte le sequenze di tipo τ_j sono equiprobabili, quindi:

$$H(X|Z=j) = \log_2 |\tau_j| \quad \Rightarrow \quad H(X|Z) = \boxed{\sum_{j=1}^{\Gamma_n} \Pr(\tau_j) \log_2 |\tau_j|}$$

$$\begin{aligned} \text{Riprendendo i conti, } \mathbb{E}[L^{(n)}] &< 2 + \log_2 \Gamma_n + H(X|Z) \\ &\leq 2 + \log_2 \Gamma_n + H(X) \\ &\leq 2 + K \log_2(n+1) + H(X). \end{aligned}$$

Quindi il tasso della codifica è:

$$R = \frac{\mathbb{E}[L^{(n)}]}{n} \leq \frac{2}{n} + K \frac{\log_2(n+1)}{n} + \frac{H(X)}{n} = \underbrace{\frac{2}{n}}_{\downarrow 0} + K \underbrace{\frac{\log_2(n+1)}{n}}_{\downarrow 0} + \frac{H(X)}{n} = \frac{2}{n} + K \frac{\log_2(n+1)}{n} + \frac{H(X)}{n}$$

Per $n \rightarrow \infty$

\Rightarrow La codifica è asintoticamente ottima.

Codifica Ziv-Lempel:

codifica "strutturale"; cerca di identificare ripetizioni

