

1. DATA COMPRESSION & ENTROPY

- * LOSSLESS ~~CI~~ CI OCCUPIAMO DI QUESTA
- LOSSY (MULTIMEDIA)

OSSERVAZIONE

NON PUÒ ESISTERE UN
SUPER COMPRESSORE CHE
COMPIME QUALSIASI COSA
GU PASSO

ESEMPIO: PRENDIAMO TUTTI I FILE
DI LUNGHEZZA $\leq N$ BIT

NON È POSSIBILE CHE TUTTI I FILE
SIANO COMPRESSI A $< N$.

$$2 + 4 + 8 + \dots$$

$$2^N = 2^{N+1} - 2$$



TIPICI COMPRESSORI LOSSLESS

1. SYMBOL SUBSTITUTION
(HUFFMAN, MORSE)
2. PARSING (LZ77 LZ78
GZIP, XZ, ...)
3. TRASFORMAZIONI

ESEMPIO TRASFORMAZIONI USATE
PER LA COMPRESSIONE

1. MTF MOVE-TO-FRONT

STRINGA DA COMPRIMERE

BACCADACCADACCAAAAA

LISTA CARATTERI. SI PARTE CON ABCD →

BACD → ABCD → DABC → ADCB → CADB

OUTPUT:

114140.....

2. DELTA CODING

INVECE DI COMPRIMERE LA SEQUENZA A

1001 1005 1007 1010 ...

CONSIDERO IL PRIMO VALORE E LE DIFFERENZE:

1001 4 2 3

3. BWT

LA VEDREMO IN SEGUITO

SYMBOL SUBSTITUTION

CLASSIFICAZIONE DI CODICI BASATI

SU SYMBOL SUBSTITUTION

TABLE 5.1 Classes of Codes

| X | Singular | Nonsingular, But Not Uniquely Decodable | Uniquely Decodable, But Not Instantaneous | Instantaneous |
|---|----------|---|---|---------------|
| A | 0 | 0 | 10 | 0 |
| B | 0 | 010 | 00 | 10 |
| C | 0 | 01 | 11 | 110 |
| D | 0 | 10 | 110 | 111 |

Ci occuperemo principalmente di quelli istantanei (PREFIX-FREE)

Theorem 5.2.1 (Kraft inequality) For any instantaneous code (prefix code) over an alphabet of size D , the codeword lengths l_1, l_2, \dots, l_m must satisfy the inequality

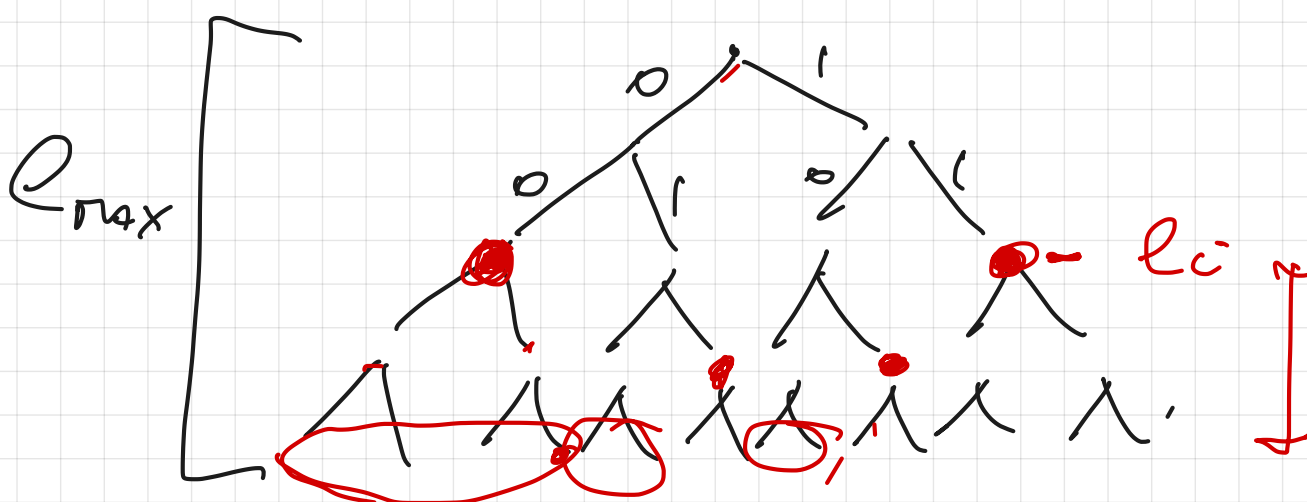
$$\sum_i D^{-l_i} \leq 1. \quad (5.6)$$

Conversely, given a set of codeword lengths that satisfy this inequality, there exists an instantaneous code with these word lengths.

DIMOSTRAZIONE KRAFT INEQUALITY

(IDEA) l_{\max} LUNGHEZZA

MASSIMA DEI CODICI



A OGM CODICE CORRISPONDE UN NODO DELL'ALBERO; ASSEGNAMO A QUEL CODICE TUTTE LE FOGLIE SOTTO QUEL NODO. OGM FOGLIA È ASSEGNATA AD AL + UN CODICE

SONO IL NUMERO DI PAGINE
ASSOCIATO AD OGNI CODICE

$$\sum_i 2^{\ell_{\max} - \ell_i} \leq \underline{2^{\ell_{\max}}}$$

$$\boxed{\sum_i 2^{-\ell_i} \leq 1} \leftarrow \text{KRAFT INEQUALITY}$$

PROBLEMA: STRINGA DI LUNGHEZZA N
I SINGOLI CARATTERI APPARTENONO

$$n_1 \quad n_2 \quad n_3 \quad \dots \quad n_m$$

VOLE

$$n_1 + n_2 + n_3 + \dots + n_m = N$$

CONSIDERO TUTTE LE POSSIBILI
CODIFICHE, UNA GENERICA CODIFICA
USA CODICI DI LUNGHEZZA

$$\ell_1 \quad \ell_2 \quad \ell_3 \quad \dots \quad \ell_m$$

LA DIMENSIONE DEL
FILE CODIFICATO SARA'

$$O = n_1 e_1 + n_2 e_2 + \dots + n_m e_m$$

HO UN COLO DELLA DIS. DI KRAFT

$$\sum 2^{-e_i} \leq 1$$

PROBLEMA DI OTTIMIZZAZIONE VINCOLATA

RISOLVIAMO CON MOLTIPLICATORI DI
LAGRANGE

$$F(e_1, e_2, \dots, e_m, \lambda)$$

$$= \sum n_i e_i + \lambda \left(\sum 2^{-e_i} - 1 \right)$$

LA TEORIA DICE CHE I PUNTI
DI MIN O MAX SONO DOVE $\partial F = 0$

OSSERVAZIONE

$$\frac{\partial F}{\partial \lambda} = \left(\sum_i 2^{-e_i} - 1 \right) = 0$$

QUESTO
È IL
VINCOLO

$$\frac{\partial F}{\partial e_i} = n_i + \lambda \frac{\partial}{\partial e_i} 2^{-e_i}$$

$$= n_i - \lambda 2^{-e_i} \cdot \ln 2 = 0$$

$$= \lambda 2^{-e_i} = \frac{n_i}{\ln 2}$$

PER OGNI i

$$2^{-e_i} = \frac{1}{\lambda} \frac{n_i}{\ln 2}$$

IMpongo IL VINCOLO:

$$\sum 2^{-e_i} = 1$$

OTTENGO:

$$\sum_i \frac{1}{\lambda} \frac{n_i}{\log_e 2} = 1$$

$$\Rightarrow \sum_i n_i = \lambda \log_e 2$$

$$\Rightarrow n = \lambda \log_e 2$$

CHE UNITO A

$$\lambda 2^{-e_i} = \frac{n_i}{\log_e 2}$$

MI DA LA FORMULA PER LA

LUNGHERZA OTTIMALE DEI CODICI

$$2^{-e_i} = \frac{n_i}{n}$$

EQUIVALENTE NIENTE

$$e_i = - \log_2 \frac{n_i}{n}$$

LA LUNGHERZA MINIMA

DELLA CODIFICA DI S

$$\sum_i -n_i \log \frac{n_i}{n} = n \cdot H_0(S)$$

ENTROPIA

DI ORDINE \emptyset

DI S

$$H_0(S) = \sum_i -\frac{n_i}{n} \log \frac{n_i}{n}$$

DEFINIAMO

$$\frac{n_i}{n} = P_i \quad \text{PROBABILITA'}$$

EMPIRICA DEL
SIMBOLO i

ESISTE UNA CATEGORIA DI CODICI

BASATI SULLE FREQUENZE $\frac{n_i}{n}$

CHE RIESCONO A COMPRIMERE

FINO A

$$n H_0(S) + \epsilon \in \mathbb{N}$$

PER FARE MEGLIO DI H_0

UNO DEI METODI È FARE LA
COMPRESSIONE BASATA SUL
CONTESTO.

LA COMPRESSIONE CHE POSSIAMO
OTTENERE CON CODICORD CHE
DIPENDONO DAL SIMBOLO PRECEDENTE

$$\sum_{i \in A} \sum_j -n_{ij} \log \frac{n_{ij}}{n_i}$$

n_{ij} = # DI j SUBITO DOPO LA i

$$= \sum_{i \in A} n_i \sum_{j \in A} - \frac{n_{ij}}{n_i} \log \frac{n_{ij}}{n_i}$$

$$= \sum_{i \in A} n_i H_0(S_i)$$

↓
INSIEME DEI SIMBOLI,
CHE SEGUONO IMMEDIATAMENTE
IL SIMBOLO i

DEFINIAMO

$$H_1(S) = \sum_{i \in A} \frac{n_i}{n} H_0(S_i)$$

L'OUTPUT DEL METODO SOPRA È
AL MEGLIO

$$n H_1(S)$$

$H_1(S)$ = ENTROPIA DI ORDINE 1
DI S

VALE

$$H_0(S) \leq H_1(S) \leq H_2(S) \leq \dots$$

METODI DI TIPO PPM

PERMETTONO DI USARE CODEWORD
CHE DIPENDONO DAI K SYMBOLI
PRECEDENTI (ES $K=4$)

SENZA DOVER TRASMETTERE
IN ANTICIPO LE CODEWORD.

MECCANISMO: COSTRUISCO LA
TABELLA DELLE FREQUENZE
DURANTE LA COMPRESSIONE

ESEMPIO. CONTESTO "COMP"

HO VISTO "COMP" 10 VOLTE

ED ERA SEGUITO DA

| | | | |
|---|---|-------|---------|
| R | 5 | VOLTE | 0 |
| E | 3 | VOLTE | 10 |
| O | 1 | VOLTA | 1110... |
| A | 1 | VOLTA | 1111... |

DEVE ESISTERE LA CODICEWORD

PER INDICARE "NUOVO"

COMPRESSORI BASATI SU PARSING
VENGONO CODIFICATE SOTTOSTRINGHE

I METODI DI FFRUSCONO PER
COME VIENE FATTO IL PARSING

E COME VENGONO CODIFICATE
LE SINGOLE FRASI DEL
PARSING.

LA FAMIGLIA DI PARSING
PIÙ NOTA È QUELLA LZ77
DI CUI ESISTONO NUMEROSE VARIANTI
CONSIDERIAMO LA SEGUENTE:

- (1) Find the longest prefix $S[i..j]$ of $S[i..n]$ that occurs in S starting before position i .
- (2) If $j \geq i$, that is, $S[i..j]$ is nonempty, then the next phrase is $S[i..j]$, and we set $i \leftarrow j + 1$.
- (3) Otherwise, the next phrase is the explicit symbol $S[i]$, which has not appeared before, and we set $i \leftarrow i + 1$.
- (4) If $i \leq n$, continue forming phrases.

ESEMPIO

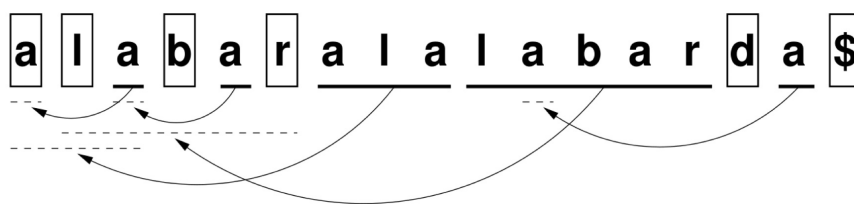


Fig. 2. Lempel-Ziv parse of $S = \text{alabaralabarda\$}$. Each phrase is either an underlined string, which appears before, or a boxed symbol. The arrows go from each underlined string to one of its occurrences to the left (which is underlined with a dashed line).