

30.08.05.

Last Time;

① Kraft's inequality:

② Lemma

a)  $p_i > p_k, \quad l_k \leq l_i$

b) Two least probable symbols have same codeword length

c) Two longest codewords  $\equiv p_m, p_{m-1}$  & they differ only in the last bit

Huffman Codes:

$p_1 \geq p_2 \geq \dots \geq p_m$

Code  $C_m$  which satisfies

a) to c)

$\{p_1, p_2, \dots, p_{m-2}, p_{m-1} + p_m\}$

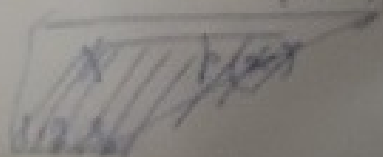
$C_{m-1}$ : Remove last bit in codewords for  $m-1$  &  $m$  in  $C_m$  & assign the common prefix to the common prefix to the new symbol  $\equiv p_{m-1} + p_m$ .

$$L(C_m) = \sum_i p_i l_i$$

$$= L(C_{m-1}) + p_{m-1} + p_m$$

→ It does not depend on the code

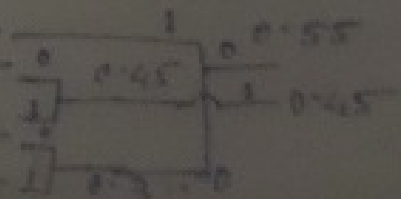
Reduce till we have to find optimal code for 2 symbols which is  $\{0, 1\}$



e.g.

	X
01	1
10	2
11	3
000	4
001	5

	$p(x)$
01	0.25
10	0.25
11	0.2
000	0.15
001	0.15



Cons: a)  $O(|X|)$  - ~~dependent~~  $\Rightarrow$  high complexity.

b) We need a probabilistic source model.

### Arithmetic Coding

\* Lempel-Ziv Algorithm:

(Penalty is that: it approaches entropy slowly, than the case when we know the statistics (Huffman Coding))

$X = \{0, 1\}$ .

Algo: Given a source string, we parse it into shortest phrases that have not occurred before.

1 0 11 01 01 00010

Dictionary: 1, 0, 11, 01, 010, 00, 10.

(shortest phrases we have to take)

### Dictionary

<u>Index</u>	<u>Code</u>	<u>Index</u>	<u>Phrases</u>	<u>Code</u>
000 $\leftarrow$		0	NULL	
001		1	1 $\rightarrow$	000 1
010		2	0 $\rightarrow$	000 0
011		3	11 $\rightarrow$	001 0
100		4	01 $\rightarrow$	010 0
101		5	010 $\rightarrow$	100 0
110		6	00	
111		7	10	

$\rightarrow$  address

1, 0

[Number of phrases logarithmic]

Decoded similarly  
decoder will also have a dictionary

No error (channel) or adversary is assumed in communication

$$\underline{2^{1/2}} \approx 1.414 \text{ limit}$$

$n$

$$C(n) = \# \text{ phrases}$$

$$L(x_1^n) = C(n) (\log(C(n)) + 1)$$

\* defn: A distinct parsing is a parsing of  $x_1^n$  in which no two phrases are equal.

\* lemma: For a distinct parsing

[This lemma does not depend upon ~~the~~ statistics]

$$C(n) \leq \frac{n}{(1 - \epsilon_n) \log(n)}$$

$$\text{where } \epsilon_n = \min \left\{ 1, \frac{\log(\log(n)) + 4}{\log(n)} \right\}$$

$$\rightarrow 0 \text{ as } n \rightarrow \infty$$

Proof:

$n_k$  = sum of the lengths of all phrases with length  $\leq k$

$$= \sum_{j=1}^k j 2^j = (k-1) 2^{k+1} + 2$$

[Any phrase of longer length will decrease number of sequences.]

$$n = n_k$$

$$C(n_k) < 2^{k+1} < \frac{n_k}{(k-1)}$$

~~We are trying to prove~~

$$\frac{2^{k+1}}{(k-1)} < \frac{n_k}{(k-1)}$$

$$n = n_k + \Delta$$

$$0 \leq \Delta < n_{k+1} - n_k = k 2^{k+2} - (k-1) 2^{k+1} \\ = 2^{k+1} \{ 2k - (k-1) \} = (k+1) 2^{k+1}$$

$$C(n) \leq C(n_k) + \frac{\Delta}{k+1} \leq \frac{n_k}{k-1} + \frac{\Delta}{k+1} \\ \leq \frac{n_k}{k-1} + \frac{\Delta}{k-1} = \frac{n_k + \Delta}{k-1} = \frac{n}{k-1}$$

$$\boxed{C(n) \leq \frac{n}{(k-1)}, \quad n_k \leq n < n_{k+1}}$$

~~$$k-1 \leq n < n_{k+1}$$~~

$$n \geq n_k$$

$$n_k = (k-1) 2^{k+1} + 2 \geq 2^k$$

$$k \leq \log(n)$$

$$n < n_{k+1}$$

$$n_{k+1} = k 2^{k+2} + 2 \leq (k+2) 2^{k+2}$$

$$2^{k+2} \geq \frac{n}{(k+2)} \geq \frac{n}{(\log(n) + 2)}$$

$$k+2 \geq \log(n) - \log(\log(n) + 2)$$

$$\Rightarrow k-1 \geq \log(n) \left[ 1 - \frac{\log(\log(n) + 2) + 3}{\log(n)} \right], \quad n \geq 4$$

$$n \geq 4 \Rightarrow \left[ \log(n) \left[ 1 - \frac{\log(\log(n) + 2) + 3}{\log(n)} \right] \right] \left\{ \because \log(n) + 2 \geq 2 \log(n) \right\}$$

$$= (1 - \epsilon_n) \log(n), \quad n \geq 4$$



$$\therefore C(n) \leq \frac{n}{(k-1)} \leq \frac{n}{(1-\epsilon_n) \log(n)} \quad (*) \quad \left[ \begin{array}{l} \text{Combinatorics} \\ \text{lemma} \end{array} \right]$$

### \* Markov Process:

A process  $\{X_n\}_{n=-\infty}^{+\infty}$  is said to be a stationary Markov process of order  $k$  if,

a) It is stationary

$$b) P(X_1^n = x_1^n | X_{-\infty}^{n-1}) = P(X_1^n = x_1^n | X_{n-k}^{n-1} = x_{n-k}^{n-1})$$

$$= P(X_1^n = x_1^n | X_{n-k}^{n-1} = x_{n-k}^{n-1})$$

(Markov Property)

[Only last  $k$  samples will suffice]

$$H(\{X_n\}) = \lim_{n \rightarrow \infty} H(X_n | X_1^{n-1})$$

$$= H(X_0 | X_{-k}^{-1})$$

$$\rightarrow H(X_0 | X_k^{-1})$$

[Using the fact that the process is stationary]

[ $k$ -order Markov process]

$$* \mathcal{I}_k(x_1^n | x_{-k}^{-1}) = P(X_1^n = x_1^n | X_{-k}^{-1} = x_{-k}^{-1})$$

$$= P(X_n = x_n | X_{n-k}^{n-1} = x_{n-k}^{n-1})$$

$$P(X_1^{n-1} = x_1^{n-1} | x_{-k}^{-1})$$

↳ [Since Markov Process, it depends on the last  $k$ ]

$$= \prod_{j=1}^n P(X_j = x_j | X_{j-k}^{j-1} = x_{j-k}^{j-1})$$

It's possible to extend a Markov process of order  $k$  to ergodic process

$$= \prod_{j=1}^n p(x_j | x_{j-k}^{j-1})$$

$$\log \left( \mathcal{P}_k \left( X_1^n | X_{-k}^{-1} \right) \right) = -\frac{1}{n} \sum_{j=1}^n \log \left( p(x_j | x_{j-k}^{j-1}) \right)$$

$$\xrightarrow{\text{Ergodicity}} -E \left[ \log \left( p(X_0 | X_{-k}^{-1}) \right) \right]$$

$$= H(X_0 | X_{-k}^{-1}) \quad \text{a.s.}$$

This is A.E.P. [Not all population is important]

\* Suppose  $x_1^n$  is parsed by LZ into  $y_1, y_2, \dots, y_{C(n)}$

$\nu_i$  = time index at which phrase  $y_i$  starts

$$= x_{\nu_i}^{\nu_{i+1}-1}$$

$$S_i = x_{\nu_i-k}^{\nu_i-1}$$

$$S_1 = x_{-(k-1)}^0$$

$$S_i \in \mathcal{X}^k$$

$C_{l,s}$  = # phrases  $y_i$  with length  $l$  and state  $S_i$

$$\boxed{x_1 \ x_2 \ x_3 \ x_4} \quad C_{2x_1} = 1 \quad \begin{matrix} l=1,2,\dots \\ s \in \mathcal{X}^k \end{matrix}$$

$y_2$

$$C(n) = \sum_{l=1}^{\infty} \sum_{s \in \mathcal{X}^k} C_{l,s}$$

$$\sum_{l,s} C$$

$$n = \sum_{l,s} l C_{l,s}$$

Homework due Tuesday