ITECS : CS 229r

LECTURE 5

TODAY : - MARKOVIAN SOURCES

- LZ Compression Theorem

— ✂ —

Source : My notes from Spring 2006 → | Gallager's Notes 1994

— ✗ —

RECALL MARKOVIAN SOURCES (= Hidden Markov Model) :



$X_i \ldots X_n \ldots$ is a hidden Markov Model

if $Z \in M$ ($k \times n$ state mc) &

- Hidden Markov Model given by matrix $M \in \mathbb{R}^{k \times k}$

& $P_x^{(1)} \ldots P_x^{(k)}$ dist on $\Omega$

- ~~Set.~~ Output of HMM $= X_1 \ldots X_n \ldots$ ~~where~~ generated as follows

~~where $X_i$ gene~~ $Z_1 \sim \pi(m)$

↑
stationary

$Z_{i-1} \xrightarrow{m} Z_i$

$Z_1, Z_2 \ldots Z_n \ldots$ generated according to ↑

$X_i \sim P_x^{(Z_i)}$ .

[all mc's in lecture "irreducible" + "aperiodic"]

# Entropy of Chain:

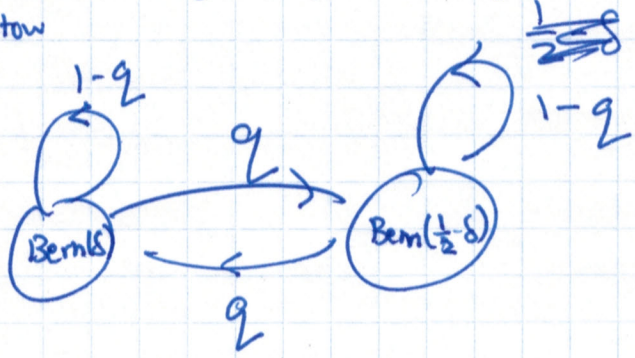$$\lim_{n \to \infty} H(X_n \mid X_1 \ldots X_{n-1}) \overset{\triangle}{=:} H(m)$$

ⓔ limit exists since

$$H(X_n \mid X_1 \ldots X_{n-1}) \leq H(X_n \mid X_2 \ldots X_{n-1})$$

$$= H(X_{n-1} \mid X_1 \ldots X_{n-2}) \quad \Big\} \text{ time invariant}$$

Note limit exists but not known to be computable.

Eg. $H(m) = ?$ (as function of $q, \delta$).
below



Theorem: $\forall$ markov chain $M$, $\forall \epsilon$, $\exists n_0 \ \forall n \geq n_0$

$$\Pr_{X_1 \ldots X_n} \left[ \ |LZ(X_1 \ldots X_n)| \ > \ (H(m) + \epsilon) \cdot n \right] \leq \epsilon.$$

[ Very "qualitative". Finite length analysis = "great project" ]

## Key Ingredients in Proof:

① Universal-Huffman-Compressor achieves "optimal compression"

② ⟹ Finite State Compressor Compresses.

③ LZ no worse than finite state compressor.

—————————————————— x ——————————————————

"Universal -Huffman- Compressor" $C_{Huff, \ell} (X_1 \ldots X_n)$; $X_i \in \Omega$

① ~~Comp for every~~

divide $X = B_1 \ldots B_{n/\ell}$        $B_i \in \Omega^\ell$

② ~~Com~~ Compute frequencies ~~for~~ $\{f_w\}_{w \in \Omega^\ell}$

among $B_1 \ldots B_{n/\ell}$

③ ~~Huffman-Code~~ $(B_1 \ldots B_{n/\ell})$;

$Z_i = Huffman (B_i ; \{f_w\})$

④ Output $\overset{\{f_w\}}{Z_1} \ldots Z_{n/\ell}$

Theorem'': $\forall\, m.c\ M,\ \forall\epsilon,\ \exists n_0\ \forall n \geq n_0$ ④

$$\Pr_{X_1 \ldots X_n}\left[\,\big|\,\mathrm{Univ\cdot Huff}(X_1 \ldots X_n)\,\big| > \big(H(m)+\epsilon\big)\cdot n\right] < \epsilon$$

———— ✕ ————

key insight to proof of Theorem': AEP holds

for Markov Sources

———— ✕ ————

### AEP Theorem. $\forall\, m.c.\ M,\ \forall\epsilon\ \exists\ \ell\ \&\ set$

$$S \subseteq \Sigma^\ell \quad \text{s.t. the following hold}$$

① $|S| \leq 2^{(H(m)+\epsilon)\cdot n}$

$\frac{1}{2^{(H(m)+\epsilon)n}} \underset{(X_1 \ldots X_\ell)\in S}{\Pr} \left[ (X_1 \ldots X_\ell) \right.$

② $\forall\, (\alpha_1 \ldots \alpha_\ell) \in S$

$$\frac{1}{2^{(H(m)+\epsilon)n}} \leq \Pr\left[ (X_1 \ldots X_\ell) = (\alpha_1 \ldots \alpha_\ell)\right] \leq \frac{1}{2^{(H(m)-\epsilon)n}}$$

③ $\underset{(X_1 \ldots X_\ell)}{\Pr}\left[ (X_1 \ldots X_\ell) \notin S\right] \leq \epsilon$ .

Proof ~~#~~ of AEP Theorem :  Omitted

Ideas :  ~~① Consider States $Z_{\ell}, Z_{2\ell}, Z_{3\ell}, \ldots$~~

① Prove for Markov Chains ~~by~~ decoupling walk
into sequences that end at state ① & do not
contain ① in between. Sequences now i.i.d.

$$Z_1 \ldots Z_a, Z_{a+1} \ldots Z_b, \ldots Z_c \ldots$$

$$\underset{①}{\overset{"}{Z_a}} \qquad \underset{①}{\overset{"}{Z_b}} \qquad \underset{①}{\overset{"}{Z_c}}$$

↑
lengths i.i.d.                    ⟶  both    } ⟶ ratio
Expected entropy iid  ⟶     converges    converges

② To get to Homm's , insert $Z_t, Z_{2t} \ldots$ at periodic
intervals ;  result in slightly higher entropy but
not much higher ;  now we have Markov
chain !

AEP Theorem $\implies$ Theorem'

Essentially implies w.p. $1-\epsilon$ we get strings$_n$ of frequency $B_i$

$$\approx 2^{-H(m)\cdot\log\ell} \implies \text{Shannon coding will associate strings}$$

$$\text{of length} \leq (H(m)+\epsilon)\cdot\log\ell$$

w.p. $\epsilon$   use string of length $\approx \log\cdot\ell$

$$\implies \text{Expected compression length}$$

$$= \frac{n}{\ell}\left(\epsilon + H(m)+\epsilon\right)\ell$$

$$= (H(m)+2\epsilon)\, n. \qquad \boxtimes$$

[Full proof will involve union bound over $S$...]

———×———

## Part II : Finite State Compressors:

$C =$ finite state compressor if

$$\boxed{X_1 - - - - \qquad\qquad\qquad X_n}$$

↑ ← Left to right read access

(F.S.m)

↓ ← left to right write access.

( Compression $(X_1...X_n)$

**Proposition**: Univ.-Huffman$^\ell$ ~~is~~ <ins>can be converted to</ins> $\underset{\wedge}{a}$ $\Omega^{2^\ell}$-state Finite State Compressor. [ with preprocessing of $M$ ].

~~**Proof**: Figure out $S$ ~~before~~ expected (proof obvious).~~

~~**Proof**:~~

**Theorem'** $\Rightarrow$ $M$ can be compressed by $\Omega^{2^\ell}$-state compressor.

**Theorem''**: if $M$ compressed by $S$-state compressor

then $M$ compressed by Lempel-Ziv.

Key Ingredient in this proof

$$C(X_1 \dots X_n) \hat{=} \max \left\{ t \mid \exists Y_1 \dots Y_t \text{ distinct in } ~~\text{set}~~ \Omega^* \right.$$
$$\left. \text{s.t.} \quad X_1 \dots X_n = Y_1 \circ Y_2 \circ \dots \circ Y_t \right\}$$

Intuition:

~~Rough idea~~ ① $C(X_1 \dots X_n) \approx \underset{\text{high}}{\frac{n}{\log n}} \Rightarrow X_1 \dots X_n$ not compressible by any thing

② $C(X_1 \dots X_n) =$ Small $\Rightarrow$ $LZ$ compresses well.

# Formal Claims

Let $t = C_{\emptyset}(x_1 .. x_n)$  then

① ~~$\exists$~~ $n \geq t \log \dfrac{t}{4}$      (intermediate)

         [$t$ is growing]

② $\forall$ $s$ state compressors $C^s$

$$|C^s(x_1.. x_n)| \geq t \cdot \log\left(\frac{t}{s^2}\right)$$

$$= t \log t - t \log s^2$$

$$\geq (1 - o(1))\, t \log t \qquad \left[\begin{array}{l} \text{if } s = O(1) \\ \& \; t = \omega(1) \end{array}\right]$$

③ $|C^{LZ}(x_1 \cdots x_n)| \leq (1 + o(1))\, t \log t$

$\Rightarrow$ $|C^{LZ}(x_1.. x_n)| \leq (1 + o(1)) \cdot |C^s(x_1.. x_n)|$

$$\leq (1 + o(1)) \cdot H(m) \cdot n \quad \left[\begin{array}{l} \text{using} \\ \text{Theorem'} \end{array}\right]$$

# Proofs of Claims

① Follow from # strings of length $\ell \leq \Sigma^\ell$;

$Y_1 \ldots Y_t$ distinct

$$\Rightarrow \quad \sum |Y_i| = \sum_R R \cdot \{\# \, i \text{ s.t. } |Y_i| = k\}$$

$$\geq \sum_{R=0}^{\lfloor \log_\Sigma \ell \rfloor} R \cdot \Sigma^k + (\ell - \Sigma \Sigma^{\lfloor \log_\Sigma \ell \rfloor})(\log \ell + 1)$$

$$\geq \quad \ell \cdot \log_\Sigma \ell \cdot (1 - o(1))$$

③
③̶ $\quad |C^{LZ}(X_1 \ldots X_n)| \leq t(\lceil \log t \rceil + 1)$

$$\leq \quad t \cdot (1 + o(1)) \cdot \log t$$

② Let $X_1 \ldots X_n = Y_1 \circ Y_2 \circ \cdots Y_t$ , $Y_i$ distinct

— Let $Y_1' \ldots Y_t'$ be outputs of FSM while traversing $Y_1 \ldots Y_t$.

— Note $Y_1' \ldots Y_t'$ need not be distinct; However (**key**) if $Y_i$ & $Y_j$ both start in state $a$ & end in $b$, then $Y_i' \neq Y_j'$

$-\quad T_{ab} \triangleq \{ j \mid \text{FSM starts at } a \text{ } \& \text{ ends at } b \text{ on } Y_j \}$

$-\quad \mathcal{L}_{ab} \triangleq \sum\limits_{j \in T_{ab}} |Y_j'|)$

$-\quad |\mathcal{L}_{ab}) \geq |T_{ab}| \log \dfrac{|T_{ab}|}{4}$ $\qquad$ (by ① applied to binary strings).

$-\quad |C^s(X_1 .. X_n)| = \sum\limits_{j=1}^{n} |Y_j'|)$

$\qquad = \sum\limits_{a,b} \cancel{\sum\sum} \mathcal{L}_{ab}$

$\qquad \geq \sum\limits_{a,b} |T_{ab}) \log \dfrac{|T_{ab}|}{4}$

$\qquad \geq \dfrac{t}{|S|^2} \cdot \log \dfrac{t}{S^2 \cdot 4} \cdot |S|^2$

$\qquad = t \cdot \log \dfrac{t}{S^2 \cdot 4} \qquad \boxtimes$