# Information Theory & Coding
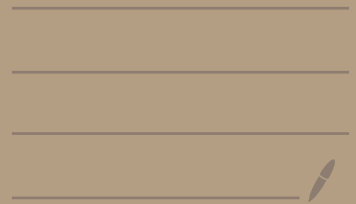
Oct 12th 2020

- Universal Compression

  - Lempel-Ziv method.

  - Finite-State Machine, I.L. compressors.

    o $\underline{FSM}$: with $m$ states.

    $g(\text{starting state}, s \text{ input letters } u_1 \cdots u_n) = t$ ← final state

    $f(\quad '' \quad , \quad '' \quad) = z_1 \cdots z_n$

    $\in \{0,1\}^{*}$

    $\rho(M, u_1 u_2 \cdots) = \limsup_{n \to \infty} \frac{1}{n} \overbrace{\text{output}}^{\#\text{ bits in the } M} (u_1 \cdots u_n)$

    o $\left( \underline{IL}: \text{ for any } s \in \mathcal{S}, \; u_1 \cdots u_n \neq u_1' \cdots u_\ell' \right.$

    $\qquad \text{either } g(s, u_1 \cdots u_n) \neq g(s, u_1' \cdots u_\ell') \longleftarrow$

    $\qquad \left. \text{or } f(s, u_1 \cdots u_n) \neq f(s, u_1' \cdots u_\ell') \longleftarrow \right)$

$\underline{\underline{Claim}}$: LZ compresses better than any IL, FSM.

$\underline{Def}$: given $u_1 \cdots u_n$ we say that $w_1 \cdots w_q$

$\left( \text{is a distinct parsing of } u_1 \cdots u_n \text{ if } \right.$

$\left. w_i \in \mathcal{U}^{*}, \; u_1 \cdots u_n = w_1 \cdots w_q \; \& \; w_i \neq w_j \atop \text{for } i \neq j \right)$

Example : $\mathcal{U} = \{a, b, c\}$

$$\mathcal{U}^* = \{ null, a, b, c, aa, ab, ac, ba, bb, bc, \cdots \}$$

$$\boxed{u_1 \cdots u_n = abaababc}$$

$$\omega_1 = a, \quad \omega_2 = b, \quad \omega_3 = aa, \quad \omega_4 = ba, \quad \omega_5 = bc$$

Note : LZ generates a distinct parsing

Recall (example LZ words do impson

$$a | b | a \; a | b \; a | b \; c | \; - \; -$$

$\omega_1 \quad \omega_2 \quad \omega_3$

$$\mathcal{J} = \{\cancel{a}, \cancel{b}, c\}$$

$$\begin{array}{lll} \cancel{aa} & \cancel{ba} & baa \\ ab & bb & bab \\ ac & \cancel{bc} & bac \end{array}$$

aaa
aab
aac

Lemma :

Suppose $u_1 u_2 \cdots \cdots$ is an infinite sequence, and let

$q(n)$ be the number of words in a distinct of $(u_1 \cdots u_n)$, Then

$$\lim_{n \to \infty} \frac{\boxed{q(n)}}{n} = 0. \quad \left( i.e. \; q(n) \text{ grows slower than } n. \right)$$

Pf : Let $u_1 \cdots u_n = \omega_1 \cdots \omega_q \qquad q = q(q)$

how many of these $\omega_i$'s can have length $\leq k-1$?

There are $1 + |u| + |u|^2 + \cdots + |u|^{k-1} = F(k)$

So $q - F(k)$ of the $w_i$'s have length $\geq k$

So $n \geq (q - F(k)) k$

$\Rightarrow q \leq \dfrac{n}{k} + F(k) \Rightarrow \dfrac{q(n)}{n} \leq \dfrac{1}{k} + \boxed{\dfrac{F(k)}{n}}$

$\displaystyle \lim_{n \to \infty} \dfrac{q(n)}{n} \leq \dfrac{1}{k} + 0$ as $k$ can be chosen arbitrarily large

we see $\displaystyle \lim_{n \to \infty} \dfrac{q(n)}{n} \begin{array}{c} \leq 0 \\ \geq 0 \end{array}$ //.

$\Bigg($ Also observe (as an aside) that we can make $q(n) \geq \sqrt{n}$

$\underbrace{u_1}_{w_1} \underbrace{u_2 u_3}_{w_2} \underbrace{u_4 u_5 u_6}_{w_3} \underbrace{u_7 \quad u_{10}}_{w_4} \cdots$ $\Bigg)$

Suppose now $u_1 \cdots u_n = w_1 \cdots w_q$ ($w$'s distinct)

is fed to a FSILM with $\leq m$ states.

$$s_i \in \mathcal{S}, \; |\mathcal{S}| \leq m$$

states: $s_0 \quad s_2 \quad s_3 \quad s_4 \quad \cdots \quad s_{q+1}$

input: $\quad w_1 \quad w_2 \quad w_3 \quad \cdots \quad w_q$

output: $\quad y_1 \quad y_2 \quad y_3 \quad \cdots \quad y_q$ — binary strings

Claim: in the collection $y_1, y_2 \cdots, y_q$

any binary string can occur at most $m^2$ times.
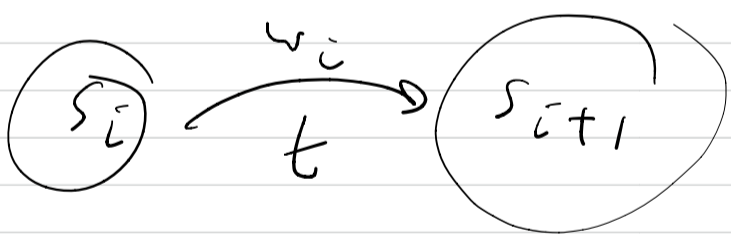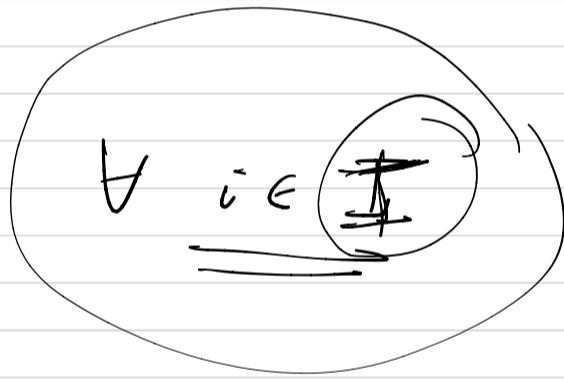
(ie. no binary string can occur $> m^2$ times).

why? : Suppose a binary string $t$ occurs $> m^2$ times

among $y_1 \cdots, y_q$, let $I = \{ i : y_i = t \}$

$|I| > m^2$. so:

$$f(s_i, \underline{w_i}) = t \qquad \forall \; i \in I$$

$$g(s_i, w_i) = s_{i+1}$$

$$s_i \xrightarrow[t]{w_i} s_{i+1} \qquad (s_i, s_{i+1}) \; \text{can take}$$
$$\text{on at most } m^2 \text{ values}$$

$$\exists \; i \triangle j \; \text{s.t} \qquad (s_i, s_{i+1}) = (s_j, s_{j+1}) = (\alpha, \beta)$$

$$f(\alpha, w_i) = t \qquad\qquad f(\alpha, w_j) = t$$

$$g(\alpha, w_i) = \beta \qquad\qquad g(\alpha, w_j) = \beta$$

Contradicts IL property

Summary : $u_1 \ldots u_q$ is a distinct pairs of

$u_1 \ldots u_n$

$\Rightarrow$ the output $y_1 \ldots y_q$ of the FSM has

the property that no binary string occurs $> m^2$

times in $y_1 \ldots y_q$.

Lemma : if $y_1 \ldots y_q$ is a collection of

binary string s.t no string $t$ can occur $> k$

times then :

- write $q = k \left[ \textcircled{1} + 2 + \ldots + 2^{j-1} \right] + \textcircled{r}$

  with $0 \leq r < \underline{k 2^j}$,

then $\sum_{i=1}^{q} \text{length}(y_i) \geq k \left[ 0 + 2 \cdot 1 + 4 \cdot 2 + \ldots 2^{j-1}(j-1) \right]$

$\phantom{then \sum_{i=1}^{q} \text{length}(y_i) \geq k} + rj$.

Example : Suppose 14 binary string s.t

no string occrs $> 3$ times then

$14 = 3 + 3 \cdot 2 + r$

$\phantom{14} = 3(1 + 2) + 5$

total length the 14 string is $\geq 3 \left[ 0 \cdot 1 + 1 \cdot 2 \right) + 2 \cdot 5$

$\phantom{total length the 14 string is} = 16$

<u>Pf</u>. set of binary strings: $\{0,1\}^* = \{$null, $0, 1, 00, 01,$
$(0, 11, 000, ..\}$

there:

1 string of length 0   (null string

2 " s " " 2 $(0, 1)$

4 " " " " 2 $(00, 01, 10, 11)$.

$q = (k) + k \cdot 2 + k \cdot 4 + \dots k \cdot 2^{j-1} + r$

$$0 \leq r < k 2^j.$$

total length of the "optimal" collection:

$0 \cdot k + 1 \cdot k \cdot 2 + 2 \cdot k 4$

$+ \dots (j-1) k 2^{j-1} + j r,$

<u>Corollary</u>: if $z_1 \dots z_q$ has the properties
in the previous lemma, then,

$$\sum_{i=1}^{q} \text{length}(z_i) \geq q \log_2 \frac{q}{8k}.$$

Pf : From the previous lemma

$$\boxed{q = k[2^j - 1] + r} \qquad 0 \le r < k2^j$$

$$\left( \sum_{i=0}^{j-1} x^i = \frac{x^j - 1}{x - 1} \right)$$

$$\text{total length} \ge k\left[ \underbrace{\sum_{i=0}^{j-1} i \, 2^i}_{2 + (j-2)2^j} \right] + rj$$

$$= (j-2)q + kj + 2r$$

$$\ge (j-2)q$$

$$q < k[2^{j+1} - 1] \le k \, 2^{j+1}$$

$$j+1 \ge \log_2 \frac{q}{k} \implies j - 2 \ge \log_2 \frac{q}{8k}$$

$$\implies \text{total length} \ge q \log_2 \frac{q}{8k} \quad //$$

Thm : Let $q^*\binom{u_1 \dots u_n}{M}$ be the largest number

of words in any distinct parsing of

$u_1 \dots u_n$. Then if $\underline{M}$ is a $\underline{\underline{\leq m}}$

state I.L.FS. M ,

length output$(\underline{M}, u_1 \dots u_n) \geq q^*\binom{u_1 \dots u_n}{M} \log_2 \dfrac{q^*\binom{u_1 \dots u_n}{M}}{\boxed{8 m^2}}$.

$\underline{Pf}$ : combine the lemma above

with $(\bcancel{**})$ above.   $/\!/$.

$\underline{Corollary}$ : For any I.L.FSM , $\underline{M}$,

$\boxed{\rho(\underline{M}, u_1 u_2 \dots)}$

$\geq \limsup\limits_{n \to \infty} \dfrac{1}{n} q^*\binom{u_1 \dots u_n}{M} \log_2 q^*\binom{u_1 \dots u_n}{M}$.

$\underline{Pf}$ : $\rho(\underline{M}, u_1 u_2 \dots)$

$= \limsup\limits_{n \to \infty} \dfrac{\text{length output}(\underline{M}, u_1 \dots u_n)}{n}$

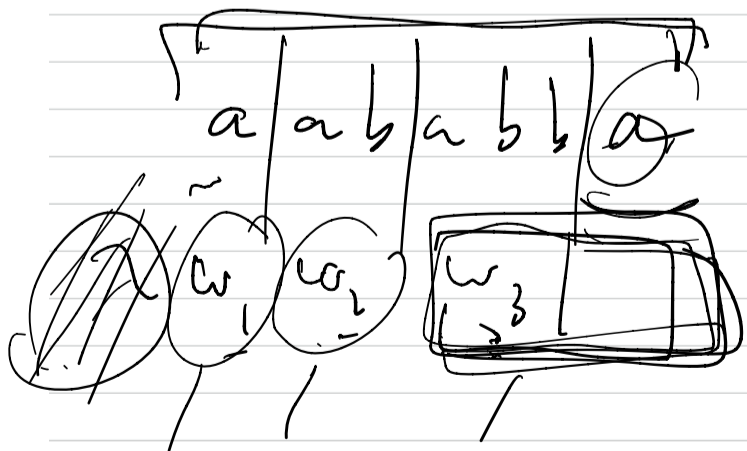$$\geq \limsup_{n \to \infty} \frac{1}{n} q^*(q) \log \frac{q^*(n)}{8m^2} \qquad \text{(lemma above)}$$

$$= \limsup_{n \to \infty} \frac{1}{n} q^*(n) \log q^*(n) - \underbrace{\frac{1}{n} q^*(n) \log(8m^2)}_{\to 0}$$

$$= \lim_{n \to \infty} \frac{1}{n} q^*(n) \log_2 q^*(n). \qquad \text{// .}$$

Corollary: $\boxed{\int_{FSM} \rho(u_1 u_2 \dots) \geq} \qquad \cdot \qquad \text{// .}$

To show:

$$P_{LZ}(u_1 u_2 \dots) \leq \limsup_{n \to \infty} \frac{1}{n} q^*(u_1 \dots u_n) \times \log_2 q^*(u_1 \dots u_n)$$

To prove this, remember how LZ operates:

$a|ab|ab|a$   $\tilde{w}_1$  $w_2$  $w_3$

# of words LZ has created so far

$$\leq q^*(u_1 \dots u_n)$$

How many bits has LZ produced?

$w_1$ is described by $\lceil \log_2 |u| \rceil$ bits.

$w_2$ is described by $\lceil \log_2 (|u| - 1 + (|u|)) \rceil$ bits

$\vdots$

$w_q$ is $\quad''\quad\quad''\quad$ $\lceil \log_2 (1 + q(|u|-1)) \rceil$

So the total output has length

$$\leq q \lceil \log_2 (1 + q(|u|-1)) \rceil$$

$$\leq q \log_2 2(1 + q(|u|-1))$$

$$\leq q \log_2 (2q|u|)$$

$$\leq q^*(u_1 \dots u_n)\left(\log_2 (2|u|) + \log_2 q^*(q_1 \dots q_n)\right)$$

So ?

$$\rho_{LZ}(u_1 u_2 \dots) \leq \limsup_{n \to \infty} \frac{1}{n} q^*(u_1 \dots u_n) \log_2 q^*(q_1 \dots q_n)$$
$$+ \frac{1}{n} q^*(q_1 \dots u_n) \log_2(2|u|)$$

$$= \left\{ \limsup_{n \to \infty} \frac{1}{n} q^*(u_1 \dots u_n) \log_2 q^*(u_1 \dots u_n) \right.$$

$Z_e$, Thm: for any $u_1 u_2 u_3 \ldots$

$$\rho_{LZ}(u_1 u_2 \ldots) \leq \rho_{FSM}(u_1 u_2 \ldots) \quad \forall$$

Corollary: if $u_1 u_2 u_3 \ldots$ is an

ergodic process, then

$$\rho_{LZ}(u_1 u_2 \ldots) \leq \text{entropy rate } \mathcal{H}$$
$$\text{of the process}$$
$$\text{with prob. } 1.$$