

# Lempel-Ziv Compression (§6)

Last time: Symbol codes ( $C: \mathcal{A} \rightarrow \{0,1\}^*$ ), Kraft's inequality ( $\sum_x 2^{-l_x} \leq 1$ ),  $H(X) \leq L(X, C) < H(X) + 1$  for UD codes, achievable by  $l_x = \lceil \log \frac{1}{P(x)} \rceil$  optimal code via Huffman algo

Today: Compression algos that operate on "stream" of symbols, can emit  $< 1$  bit/symbol, are asymptotically optimal for IID sources ( $\mathbb{R} \rightarrow H(X)$ ), but are also adaptive!

Used in GIF; uugy  
 Similar to LZ78 which is used in ZIP, PNG, ...  
 (together with Huffman)

## Lempel-Ziv's coding algo (LZ78):

- assume: stream ends with special symbol  $\perp$
- \* phrases  $\leftarrow [\emptyset]$
  - \* While more to compress:
    - read symbols until we obtain "phrase"  $\pi \notin \text{phrases}$
    - $\Rightarrow \pi = \boxed{\tau} \boxed{x}$  where  $\tau \in \text{phrases}, x \in \mathcal{A}$
    - append  $\pi$  to phrases
    - $k \leftarrow$  index of  $\tau$  in phrases
    - write  $(k, x)$  in bits

$\rightarrow$  [HW]

} splits input into minimal distinct "phrases"

use  $\lceil \log(j) \rceil$  bits in  $j$ -th step ( $j=1,2,\dots$ )      can skip if  $x = \perp$  (last step)

Example: Let's compress A|B|B A|B A A|B A A B|A B|A  $\perp$ :

Step	0	1	2	3	4	5	6	7
phrases	$\emptyset$	A	B	BA	BAA	BAA B	AB	A $\perp$
$(k, x)$	-	(0, A)	(0, B)	(2, A)	(3, A)	(4, B)	(1, B)	(1, $\perp$ )
Compression #bits for k		-, 0	0, 1	10, 0	11, 0	100, 1	001, 1	001, -
		0	1	2	2	3	3	3

$\Rightarrow$  😞 14 bits compressed into 20 bits... but the principle is sound 😊

**Q:** Intuition? Clear how to decompress?

**Analysis?** Let  $R = \frac{\ell}{n}$  the compression rate, where  $\ell = \# \text{bits of compression}$ .

\* **Average case:** For IID source, where  $X^N = X_1 \dots X_N \stackrel{i.i.d.}{\sim} P$ :

$$E[R] \leq H(P) + O\left(\frac{1}{\log N}\right) \rightarrow H(P)$$

\* **Worst case:** For any string  $x^N = x_1 \dots x_N$ ,

$$R \leq \log \#A + O\left(\frac{1}{\log N}\right) \rightarrow \log \#A$$

$f = O(g)$  means  $\exists c > 0: f(N) \leq c \cdot g(N) \forall N$

**Warmup:** Fix  $x^N$  and assume LZ decomposes it into  $c$  phrases:

$$x^N = x_1 \dots x_N = \pi_1 \dots \pi_c$$

↑                    ↑  
disjoint phrases

$$\Rightarrow \ell = \sum_{j=1}^c (\lceil \log c_j \rceil + \lceil \log \#A \rceil)$$
$$\leq c \cdot \log(c) + c(1 + \lceil \log \#A \rceil)$$

need to understand

To make progress, we need to upper-bound  $c$ . For worst-case analysis, we would try to bound  $c$  in terms of  $N$  → **EX CLASS**.

We focus on the average case, so want to relate  $c$  to  $P(x^N)$  ∇

For simplicity: **Assume all  $P(x) \leq \frac{1}{2}$**  ← but arbitrary  $\#A$  ∩

For our fixed string  $x^N$ , consider:

$$\Pi_k = \{ \pi_i \mid 2^{-k-1} < P(\pi_i) \leq 2^{-k} \}$$

← classify phrases according to probability

\* for any phrase:  $P(\pi) = \Pr(Y^N \text{ has prefix } \pi)$

\* any  $y^N$  has at most one prefix in  $\Pi_k$  (if both  $\pi_i$  &  $\pi_j$  are prefix then  $\pi_i = \pi_j * \dots * (\text{or vice versa}) \Rightarrow P(\pi_i) \leq P(\pi_j) \cdot \frac{1}{2}$  (↙))

$$\Rightarrow 1 \geq \Pr(Y^N \text{ has prefix in } \Pi_k) \geq \sum_{\pi \in \Pi_k} P(\pi) \geq \#\Pi_k \cdot 2^{-k-1}$$

Thus:  $\#\Pi_k \leq 2^{k+1}$

How large can  $P(x^N)$  be if we know it has  $c$  phrases?

$$P(x^N) = \prod_k \prod_{\pi \in \Pi_k} P(\pi) \quad \left. \begin{array}{l} \text{maximal if } 2^{k+1} \text{ phrases in } \Pi_k \text{ (} \forall k \end{array} \right\}$$

$$\leq (2^{-0})^{2^{0+1}} (2^{-1})^{2^{1+1}} \dots (2^{-(L-1)})^{2^L} (2^{-L})^{c - (2^{L+1} - 2)}$$

where  $L$  is maximal with  $\sum_{k=0}^L 2^k \equiv 2^{L+1} - 2 \leq c$ . Note:

①  $c \geq 2^{L+1} - 2 \Rightarrow L \leq \log(c+2) - 1 \leq \log(c)$  if  $c \geq 2$ , i.e.  $N > 1$

②  $c < 2^{L+2} - 2 < 2^{L+2} \Rightarrow L \geq \log(c) - 2$

Thus:

$$\log \frac{1}{P(x^N)} \geq \sum_{k=1}^L (k-1) 2^k + L(c - 2^{L+1} + 2)$$

check by inclusion  $\Rightarrow (L-2)2^{L+1} + 4 + L(c - 2^{L+1} + 2)$

$$= -4 \cdot 2^L + 4 + cL + 2L$$

①  $\geq -4c + \cancel{4} + c(\log c - 2) + 2(\log c - \cancel{2})$

②  $\geq -4c + \cancel{4} + c(\log c - 2) + 2(\log c - \cancel{2})$

$$= c \cdot \log c - 6c$$

Take expectation values and use  $H(x^N) = N \cdot H(P)$ :

$$\boxed{N \cdot H(P) \geq E[c \cdot \log c] - 6E[c]} \quad (*)$$

$\Rightarrow E[R] = \frac{1}{N} E[R] \stackrel{\text{warpup}}{\leq} \frac{1}{N} E[c \cdot \log c] + \frac{1}{N} (\lceil \log \#A \rceil + 1) \cdot E[c]$

$$\leq H(P) + O\left(\frac{1}{N} E[c]\right)$$

How to deal with  $E[c]$ ?

want that  $\rightarrow 0$

$$E[c] \log E[c] \leq E[c \cdot \log c] \leq (HCP + G) N \quad \text{since certainly } c \leq N$$

$\uparrow$  Jensen:  $f(x) = x \cdot \log x$  is convex

...so  $E[c]$  has to grow slower than linear  $\triangleright$  (in fact:

$$\Rightarrow E[c] = O\left(\frac{N}{\log N}\right) \text{ and so we arrive at}$$

$$\Rightarrow E[R] \leq HCP + O\left(\frac{1}{\log N}\right)$$

QED

Why is this true? Assume that  $f(N) \cdot \log f(N) \leq \gamma \cdot N$  for large  $N$ .

We claim that  $f(N) < (\gamma + 1) \frac{N}{\log N}$ . Indeed, otherwise we have

$f(N) \geq (\gamma + 1) \frac{N}{\log N}$  for a subsequence of  $N \rightarrow \infty$ . Then:

$$f(N) \cdot \log f(N) \geq (\gamma + 1) \frac{N}{\log N} \log \left( (\gamma + 1) \frac{N}{\log N} \right)$$

$$\geq (\gamma + 1) N \left( 1 - \frac{\log \log N}{\log N} \right)$$

$\rightarrow 0$

$\downarrow$