# Lempel–Ziv Compression (§6.4)

So far: Symbol codes achieve $H(X) \leq L(X,C) < H(X)+1$, but always $\geq 1 \frac{bit}{symbol}$.

By looking at large blocks, can achieve $H(P)$ for IID sources. ← both in the lossy and in the lossless scenario

Today: Lossless compression of "stream" of symbols that can emit $< 1 \frac{bit}{symbol}$

is asymptotically optimal for IID sources ($R \to H(X)$), and even is adaptive !

Variations are used in GIF, ZIP, PNG, ... (sometimes combined with Huffman)

---

**Lempel–Ziv Coding algo**

input: Stream that ends with special symbol $\bot$

* phrases ← [" "]   empty string

* While more to compress:

  – read symbols until we obtain "phrase" $\pi \notin$ phrases

    $\Rightarrow \pi = [\,\tau\,,\,x\,]$ where $\tau \in$ phrases, $x \in \mathcal{A}$

  – append $\pi$ to phrases

  – $k \leftarrow$ index of $\tau$ in phrases

  – write $(k, x)$ in bits

} Splits input into minimal distinct "phrases"

Use $\lceil \log(j) \rceil$ bits in $j$-th step ($j = 1,2,...$)   Can skip if $x = \bot$ (last step)

---

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

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

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

Q: Intuition how it works? Clear how to decompress?

Analysis? How well does it compress? Consider:

$$\ell = \#\text{bits of compression} \quad \& \quad R = \frac{\ell}{N} \quad \text{compression rate}$$

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

$$R \leq \log \#\mathcal{A} + O\left(\frac{1}{\log N}\right) \longrightarrow \log \#\mathcal{A} \qquad \boxed{\text{EX}}$$

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

Thus: LZ does no worse than not compressing at all! (for large N)

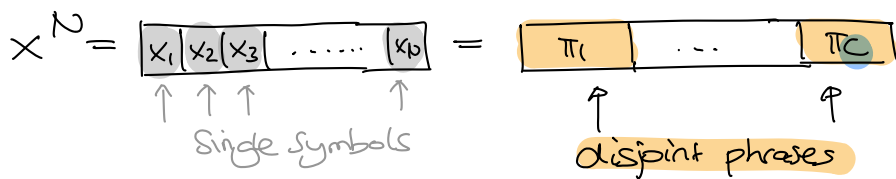* Average rate: Let $X^N = X_1 \cdots X_N \overset{\text{IID}}{\sim} P$.

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

Thus: For an IID source, LZ achieves entropy $H(P)$! (for large N)

This optimality holds even more generally for an "ergodic" source.

How to prove this?

optional reading ↓

Warmup: Fix source string $x^N$ and assume LZ compresses it into ○c phrases:

$$x^N = \boxed{x_1 | x_2 | x_3 | \; \cdots \cdots \; | x_D |} = \boxed{\pi_1 | \; \cdots \; | \pi_C |}$$

↑ ↑ ↑      ↑

Single symbols     ↑       ↑

disjoint phrases

$$\Rightarrow \boxed{\ell = \sum_{j=1}^{C} \left( \lceil \log(j) \rceil + \lceil \log \#\mathcal{A} \rceil \right)}$$

Ⓐ

dominant term

$$\leq C \cdot \log(C) + C\left(1 + \lceil \log \#\mathcal{A} \rceil\right)$$

Thus: Need to understand how number of phrases $C$ grows with $N$.

* Worst-case analysis? → Challenge exercise tomorrow. EX CLASS.

* We focus on average rate. key idea: Relate $C$ to $\log \dfrac{1}{P(X^N)}$ ∇₀

For simplicity: Assume all $P(x) \leq \frac{1}{2}$ ↩ — but arbitrary #$A$ ☺

① Classify phrases according to their probability:

$$\Pi_k = \left\{ \pi_i \ \bigg| \ 2^{-k-1} < P(\pi_i) \leq 2^{-k} \right\}$$
IID distribution

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

* any string $y^n$ has at most one prefix in any fixed $\Pi_k$

$$\left[ \text{if } y^n = \boxed{\pi_i | \cdots} = \boxed{\pi_j | \cdots} \text{ then } \pi_i = \boxed{\pi_j | \cdots} \text{ (or vice versa)} \right.$$
$$\left. \implies P(\pi_i) \leq P(\pi_j) \frac{1}{2} \ \zeta_b \right]$$

$$\implies 1 \geq Pr(X^N \text{ has prefix in } \Pi_k) = \sum_{\pi \in \Pi_k} Pr(X^N \text{ has prefix } \pi)$$
$$\geq \#\Pi_k \cdot 2^{-(k+1)}$$

$$\implies \boxed{\#\Pi_k \leq 2^{k+1}}$$

② How large can $P(X^N)$ be if we know it has $C$ phrases?

$$P(X^N) = \prod_i P(\pi_i) = \prod_k \prod_{\pi \in \Pi_k} P(\pi)$$
maximal if $\Pi_0, \Pi_1 \cdots$ as large as possible, ie. $\#\Pi_k = 2^{k+1}$

$$\leq \left(2^{-0}\right)^{2^{0+1}} \left(2^{-1}\right)^{2^{1+1}} \cdots \left(2^{-(L-1)}\right)^{2^L} \left(2^{-L}\right)^{C - \sum_{k=1}^{L} 2^k}$$

where $L$ is maximal with $\sum_{k=1}^{L} 2^k = 2^{L+1} - 2 \leq C$

$$\implies \boxed{L \approx \log(C)} \quad \text{More precisely: } \log(C) - 2 < L \leq \log(C+2) - 1 \leq \log(C).$$
if $C > 2$, ie. $N > 1$

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

Check by $\overset{\text{(=)}}{\text{induction}}$ $(L-2)2^{L+1} + 4 + L(c - 2^{L+1} + 2)$    $\overset{\circ}{\not{\text{Cancel}}}_\circ$

$$= -4 \cdot 2^{L} + 4 + L(c+2)$$

dominant term

$$\geq c \cdot \log c - 6c$$

$\Big)$ after some manipulations

③ Take expectation value and use $E\left[\log \frac{1}{P(x^N)}\right] = N \cdot H(P)$:

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

dominant term

Suppose we could only look at the "dominant" terms in Ⓐ, Ⓑ. Then:

$$E[R] = \frac{E[\ell]}{N} \overset{Ⓐ}{\underset{\sim}{\leq}} \frac{E[c \cdot \log c]}{N} \overset{Ⓑ}{\leq} H(P)$$

and we would be done!    (😊)

④ In reality, things are a bit more complicated:

$$E[R] = \frac{1}{N} E[\ell] \overset{Ⓐ}{\leq} \frac{1}{N} E[c \cdot \log c] + \frac{1}{N}(\lceil \log \#\omega \rceil + 1) E[c]$$

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

want that $\to 0 \ldots$

How to deal with $E[c]$?

$$E[c] \log E[c] \overset{Ⓑ}{\underset{\uparrow}{\leq}} E[c \cdot \log c] \leq (H(P)+6)N$$    since certainly $c \leq N$

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

...so $E[c]$ has to grow slower than linear! In fact:

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

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

Noon

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 $N \to \infty$. Then:

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

$$\geq (\gamma+1) N \left(1 - \underbrace{\frac{\log \log N}{\log N}}_{\to 0}\right)$$

$\zeta$ $\square$