# Lempel-Ziv Compression (§6.4)

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

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 $\perp$

* phrases ← [∅]
* While more to compress:
  - read symbols until we obtain "phrase" $\pi \notin$ phrases
    $\implies \pi = [\tau, x]$ where $\tau \in$ phrases, $x \in \omega$      } Splits input into minimal distinct "phrases"
  - append $\pi$ to phrases
  - $k$ ← index of $\tau$ in phrases
  - write $(k,x)$ in bits

Use $\lceil \log_2 j \rceil$ bits in $j$-th step ($j = 1, 2, ...$)      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 | ∅ | A | B | BA | BAA | BAAB | AB | A$\perp$ |
| $(k,x)$ | — | (0,A) | (0,B) | (2,A) | (3,A) | (4,B) | (1,B) | (1,$\perp$) |
| Compression | —,0 | 0,1 | 10,0 | 11,0 | 100,1 | 001,1 | 001,— | |
| #bits for $k$ | 0 | 1 | 2 | 2 | 3 | 3 | 3 | |

$\implies$  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:

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

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

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

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

EX

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

\* Average rate: Let $X^N = X_1 \cdots X_N \overset{IID}{\sim} P$.
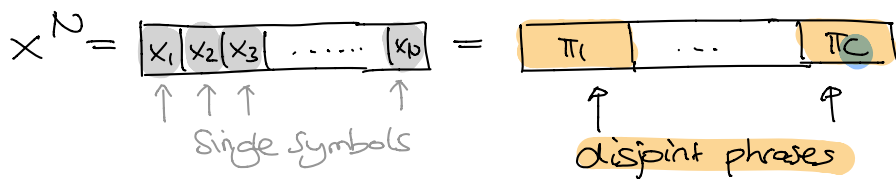
$$\mathbb{E}[R] \leq H(P) + O\left(\frac{1}{\log N}\right) \longrightarrow H(P)$$

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?

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

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

↑ ↑ ↑        ↑
Single symbols          disjoint phrases

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

Ⓐ

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. <u>key idea</u>: Relate $c$ to $\log \frac{1}{P(x^N)}$ ▽

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

① Classify phrases according to their probability:

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

⟵ 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 $\Pi_k$

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

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

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

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. $\#\Pi_k = 2^{k+1}$

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

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

$$\Longrightarrow \boxed{L \approx \log(c)}$$

More precisely: $\log(c) - 2 < L \leq \log(c+2) - 1 \leq \log(c)$.

if $c > 2$, ie. $N > 1$

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

Check by induction $\stackrel{(=)}{} \underbrace{(L-2) 2^{L+1}}_{} + 4 + L\left(c - \underbrace{2^{L+1}}_{} + 2\right)$   $\triangleright$ Cancel $\stackrel{\circ}{\circ}$

$= -4 \cdot 2^L + 4 + L(c+2)$   $\Big)$ after some manipulations

dominant term

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

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:

$E[R] = \frac{E[\ell]}{N} \stackrel{Ⓐ}{\lesssim} \frac{E[c \cdot \log c]}{N} \stackrel{Ⓑ}{\lesssim} H(P)$

and we would be done!   (◕‿◕)

④ In reality, things are a bit more complicated:

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

$\leq H(P) + \underbrace{O\left(\frac{1}{N}\right) \cdot E[c]}_{\text{want that} \rightarrow 0 \dots}$

How to deal with $E[c]$?

$\underbrace{E[c] \log E[c] \stackrel{Ⓐ}{\leq} E[c \cdot \log c] \stackrel{Ⓑ}{\leq} (H(P) + 6) N}_{}$   Since certainly $c \leq N$

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

... so $E[c]$ has to grow slower than linear $\stackrel{\circ}{\circ}$ In fact:

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

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

Noo1

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:

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

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

$\to 0$

$\notin$   $\square$