

# Data Compression

Def •  $\mathcal{X}$ : raw <sup>data?</sup> information?  
 • source code or code  $C: \mathcal{X} \rightarrow \{0,1\}^*$  → only consider binary code.  
 •  $C(x)$ : codeword of  $x$ ,  $l(x)$ : length of  $C(x)$ .

Example • ~~Assume~~  $\mathcal{X} = \{A, B, C, D\}$

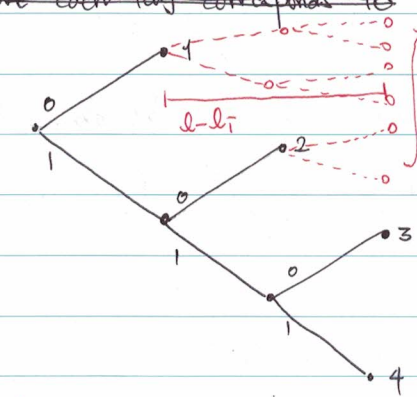
## Classes of codes

- ①  $C(T) = 0$   $C(L) = 0$   $C(E) = 0$   $C(Z) = 0$  → Singular
- ②  $C(T) = 0$   $C(L) = 010$   $C(E) = 01$   $C(Z) = 10$  → Nonsingular but not uniquely decodable. Consider  $010 \rightsquigarrow 'L'$  OR  $'TE'$  OR  $'E'$ .
- ③  $C(T) = 10$   $C(L) = 00$   $C(E) = 11$   $C(Z) = 110$   
 → Nonsingular, uniquely decodable, not instantaneous.  
 Consider  $110$  and  $1100$  → Need to read the whole string.
- ④  $C(T) = 0$   $C(L) = 10$   $C(E) = 110$   $C(Z) = 111$  → Instantaneous.

WLOG,  $\mathcal{X} = \{1, 2, \dots, m\}$  and write  $l_i := l(i) \forall i \in \mathcal{X}$ .

Thm 5.2.1 (Kraft inequality)  $C$  is instantaneous  $\iff \sum_{i \in \mathcal{X}} 2^{-l_i} \leq 1$ .

pf) ~~Assume~~ ( $\Rightarrow$ ) As  $C$  is instantaneous,  $\exists$  corresponding complete binary tree where each leaf corresponds to s.t.



draw imaginary full binary tree.

→ (# leaves) =  $2^l$

(# leaves in  $i$ 's subtree) =  $2^{l-l_i}$  mutually disjoint

→  $\sum_{i \in \mathcal{X}} 2^{l-l_i} \leq 2^l$

⇒  $\sum_{i \in \mathcal{X}} 2^{-l_i} \leq 1$ .

$l$ : max length depth (starting from 0)

( $\Leftarrow$ )  $l := \max l_i$ , so we then have

$$\sum_{i \in \mathcal{X}} 2^{-l_i} \leq 2^{-l}$$

Use the same technique.  $\square$

Problem • Given a random variable  $X$  on  $\mathcal{X}$ , find an <sup>Instantaneous</sup> source code  $C$  that minimizes the expected code length,

$$\mathbb{E}[l(X)] = \sum_{x \in \mathcal{X}} p(x) l(x).$$

Optimal code Formulate the following opt problem:

$$\min_{\mathcal{L}} \sum_{i \in \mathcal{X}} p_i l_i$$

$$\text{s.t. } \sum_{i \in \mathcal{X}} 2^{-l_i} \leq 1 \quad (\text{by Kraft inequality})$$

$$l_i \in \mathbb{Z}_+, \quad \forall i \in \mathcal{X}.$$

• Relaxation:

$$l_i \geq 0, \quad \forall i \in \mathcal{X}.$$

• Further simplification: if  $l_i < 0$ ,  $2^{-l_i} > 1$   
 $\Rightarrow$  hence, we can drop <sup>the</sup> sign constraints

$$\min_{\mathcal{L}} \sum_{i \in \mathcal{X}} p_i l_i \quad \text{s.t. } \sum_{i \in \mathcal{X}} 2^{-l_i} \leq 1. \quad \lambda$$

• Want to find a "lower bound" of the above prob.

$$\min_{\mathcal{L}} \underbrace{\sum_{i \in \mathcal{X}} p_i l_i}_{\geq 0} + \underbrace{\lambda \left( \sum_{i \in \mathcal{X}} 2^{-l_i} - 1 \right)}_{\leq 0 \text{ for feas } \mathcal{L}} \quad \text{s.t. } \lambda \geq 0$$

• ~~Given~~ Fix  $\lambda \geq 0$ , the obj func is <sup>differentiable</sup> "convex" while  $\mathcal{L} \in \mathbb{R}^m$ . The min is obtained when the gradient = 0.



$$\forall i \in \mathcal{X} \quad \frac{\partial L}{\partial l_i} = p_i - \lambda \cdot 2^{-l_i} \ln 2 = 0 \Rightarrow 2^{-l_i} = \frac{p_i}{\lambda \ln 2}$$

As  $\sum_{i \in \mathcal{X}} 2^{-l_i} \leq 1$ , (and for some reason),  $= 1$ . ← plug in.

$$\Rightarrow \sum_{i \in \mathcal{X}} \frac{p_i}{\lambda \ln 2} = 1 \Rightarrow \lambda = \frac{1}{\ln 2}. \quad \text{plug back}$$

We thus have  $\forall i \in \mathcal{X}$ ,  $l_i = \log \frac{1}{p_i}$ . ~~optimal code~~

~~If  $\log \frac{1}{p_i}$  are all integers, optimal code? What if not?~~

~~Shannon code~~  ~~$l_i = \lceil \log \frac{1}{p_i} \rceil$~~

~~Relaxation~~

Thm For any instantaneous code  $C$  for  $X$ ,  $\mathbb{E}[l(X)] \geq H(X)$ .

pf)  $l_i \leftarrow \log \frac{1}{p_i}$  is an opt soln to the relaxation. whose obj val is  $\sum_{i \in \mathcal{X}} p_i \log \frac{1}{p_i} = H(X)$ .  $\square$

If  $\{\log \frac{1}{p_i}\}$  are all integers, definitely opt! What if not?

Shannon code  $l_i \leftarrow \lceil \log \frac{1}{p_i} \rceil$ .

Thm This code is an instantaneous code &  $\mathbb{E}[l(X)] \leq H(X) + 1$ .

pf)  $l_i$  indeed satisfies Kraft inequality since

$$\sum_{i \in \mathcal{X}} 2^{-l_i} = \sum_{i \in \mathcal{X}} 2^{-\lceil \log \frac{1}{p_i} \rceil} \leq \sum_{i \in \mathcal{X}} 2^{-\log \frac{1}{p_i}} = \sum_{i \in \mathcal{X}} p_i = 1.$$

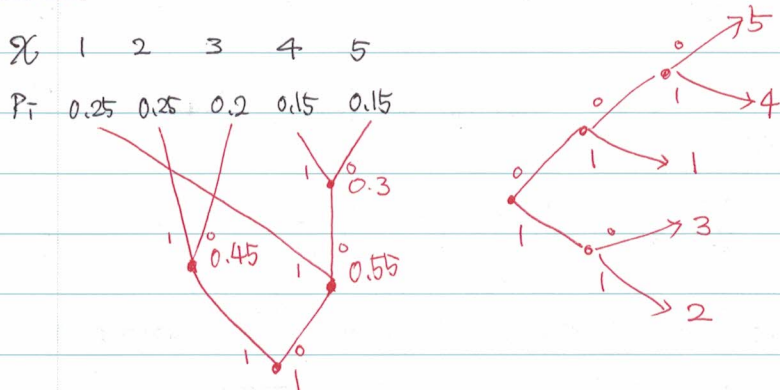
As  $l_i < \log \frac{1}{p_i} + 1$ ,

$$\sum_{i \in \mathcal{X}} p_i l_i < \sum_{i \in \mathcal{X}} p_i (\log \frac{1}{p_i} + 1) = H(X) + 1. \quad \square$$

Suboptimal Consider  $\mathcal{X} = \{1, 2\}$  w/  $p_1 = 1 - \epsilon$  &  $p_2 = \epsilon$ .

Then  $l_1 \leftarrow 1$  &  $l_2 \leftarrow$  Very big. while  $C(1) = 0$   $C(2) = 1$  is sufficient.

# Huffman code



This code is an integral opt code. Greedy alg.  $\mathbb{P}$  omitted.

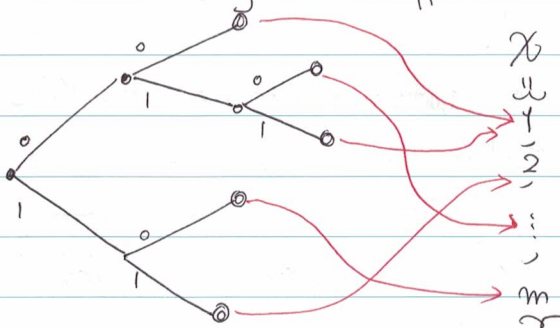
## Discrete rand var gen

So far,  $X \sim p \rightarrow \sqrt{\text{find } C \text{ s.t. } \min \mathbb{E}[l(X)]}$ .

Q. #fair coin tosses to generate  $X \sim p$ ?

Thm For any alg to gen  $X \sim p$ ,  $\mathbb{E}(\# \text{tosses}) \geq H(X)$ .

pf) The exec of the alg can be mapped to a complete binary tree.



Let  $\mathcal{Y} := (\text{leaves of the tree})$  & consider mapping  $\sigma: \mathcal{Y} \rightarrow \mathcal{X}$ .

~~Observe~~  $\forall y \in \mathcal{Y}$ ,  $P_y$  be the prob of  $y$  in tree.

Obsv we have

$$\forall x \in \mathcal{X}, \sum_{y \in \sigma^{-1}(x)} P_y = P_x.$$

$$\mathbb{E}[\# \text{ tosses}] = \sum_{y \in \mathcal{Y}} p'_y \cdot \lg y = \sum_{y \in \mathcal{Y}} p'_y \cdot \lg \frac{y}{p'_y}$$

$$= \sum_{x \in \mathcal{X}} \left( \sum_{y \in \mathcal{T}(x)} p'_y \lg \frac{y}{p'_y} \right) \quad \text{since } p'_y \leq p_x$$

$$\geq \sum_{x \in \mathcal{X}} \lg \frac{1}{p_x} \left( \sum_{y \in \mathcal{T}(x)} p'_y \right) = H(X). \quad \square$$

• If  $\{p_x\}$  are all power of two, easy to find an opt alg. What if not?

Split into powers of two! e.g.,

$$\frac{7}{8} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8}, \quad \frac{1}{3} = \frac{1}{4} + \frac{1}{12} = \frac{1}{4} + \frac{1}{12} + \frac{1}{12} + \dots$$

& treat each as distinct outcomes of  $\mathcal{Y}$ .

Thm For this alg,  $\mathbb{E}[\# \text{ tosses}] \leq H(X) + 4$

pf) 
$$\mathbb{E}[\# \text{ tosses}] = \sum_{y \in \mathcal{Y}} p'_y \cdot \lg \frac{1}{p'_y}$$

$$= \sum_{x \in \mathcal{X}} \left( \sum_{y \in \mathcal{T}(x)} p'_y \cdot \lg \frac{1/p'_y}{p_x} \right)$$

$$= \underbrace{\sum_{x \in \mathcal{X}} \left( \sum_{y \in \mathcal{T}(x)} p'_y \right) \cdot \lg \frac{1}{p_x}}_{= H(X)} + \underbrace{\sum_{x \in \mathcal{X}} \left( \sum_{y \in \mathcal{T}(x)} p'_y \cdot \lg \frac{1}{p'_y} \right)}_{\leq 4?}$$

Suffices to show that  $\forall x \in \mathcal{X}$

$$\sum_{y \in \mathcal{T}(x)} p'_y \cdot \lg \frac{1}{p'_y} \leq 4 p_x$$

Spse  $\frac{1}{2^{2k}} \leq p_x < \frac{1}{2^{2k-1}}$  for some  $k$ . Then, a subset of  $\frac{1}{2^{2k-1}}, \frac{1}{2^{2k}}, \dots$  constitute  $p'_y$ 's; hence,

$$(LHS) \leq \sum_{k=1}^{\infty} \frac{1}{2^{2k}} \cdot \lg \frac{1}{2^{2k}} < \sum_{k=1}^{\infty} \frac{k}{2^{2k}}$$



As  $\sum_{k=1}^{\infty} \frac{k}{2^k} = 2$ , we have

$$(\text{LHS}) < \frac{2}{2^e} \leq 4P_d \quad \text{since } P_d \geq \frac{1}{2^{e+1}} \quad \square$$

In textbook,  $\mathbb{E}[\# \text{ tests}] \leq H(X) + 2$ . given.