

$X$  r. v.

$$C: X \rightarrow \{0,1\}^*$$

$$x \in X \quad C(x)$$

$l(x) = |C(x)|$  ... length of encoding

Q:  $E[l(x)] \approx H(X)$  ... for good encoding

$y \in \{C(x), x \in X\} \dots$  code  
 $\rightarrow$  code word  $= C$

$C^*$  closure of code  $C$

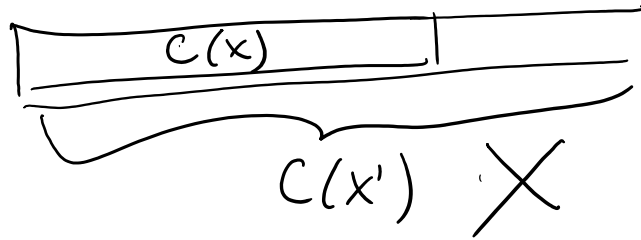
$$= \{y_1 y_2 \dots y_k, y_1, \dots, y_k \in C\}$$

$C$  unambiguous ... if  $\forall x_1, \dots, x_k \in X$   
 $x'_1, \dots, x'_k$

$$C(x_1)C(x_2)\dots C(x_k) \neq C(x'_1)C(x'_2)\dots C(x'_k)$$

$C$  prefix-free if  $\forall x, x' \in X$

$C(x)$  is not a prefix of  $C(x')$



- $C$  prefix free  $\Rightarrow$   $C$  unambiguous  
 ~~$\Rightarrow$~~   
~~2~~  
 (Exe)

$C$  suffix-free  $\Rightarrow$  ...

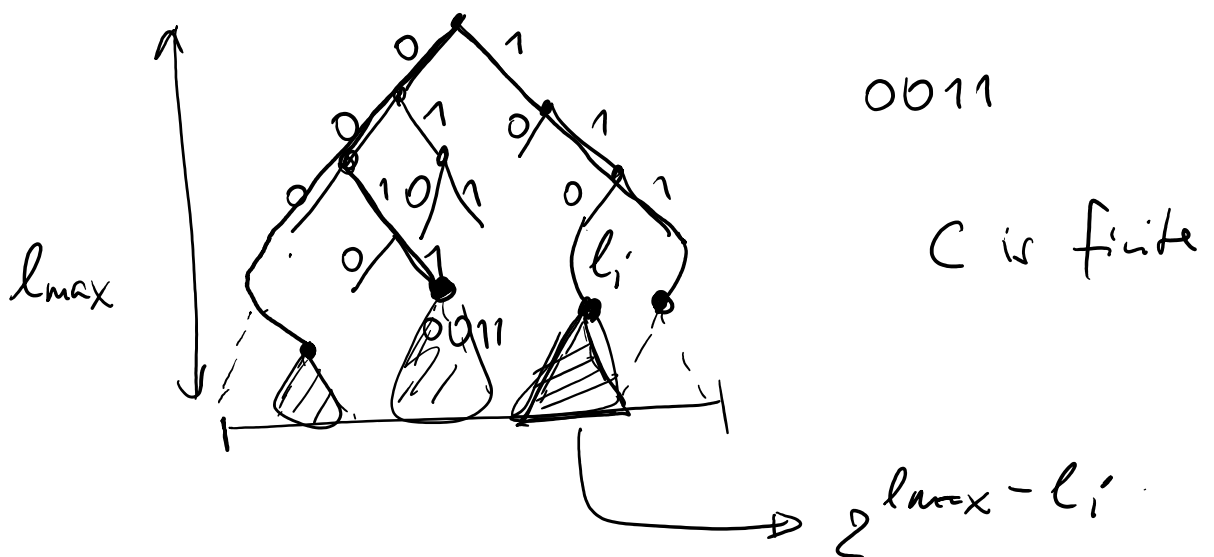
01  
0

Thm: (Kraft Inequality)

Let  $C$  be prefix-free code with code-lengths  $l_1, l_2, \dots$

$$\sum_i 2^{-l_i} \leq 1$$

Pf:



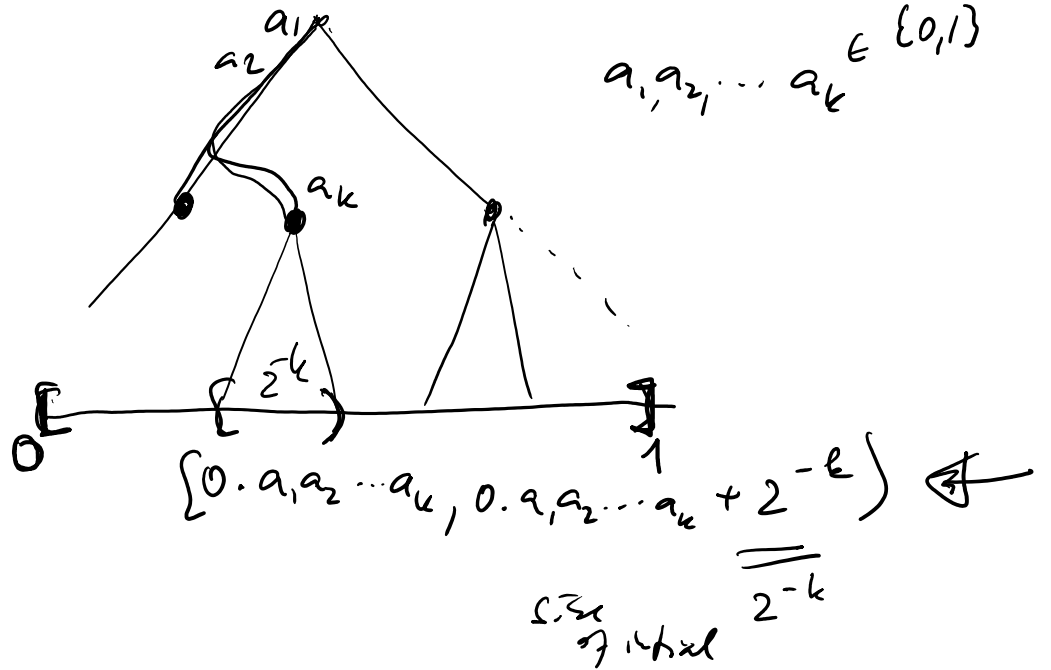
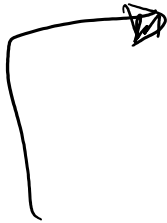
$2^{l_{max}}$  leaves

$$\sum_i 2^{l_{max} - l_i} \leq 2^{l_{max}}$$

$$2^{l_{max}} \sum_i 2^{-l_i} \leq 2^{l_{max}}$$

$$\sum_i 2^{-l_i} \leq 1$$

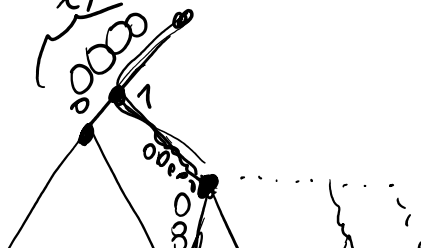
□

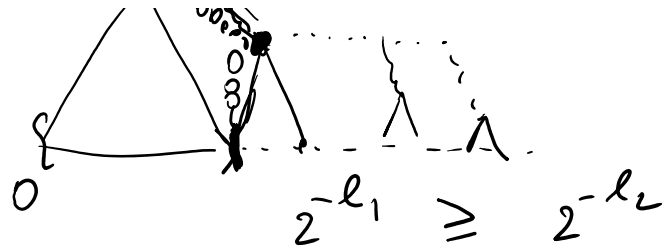


Claim If  $l_1, l_2, \dots$  satisfy  $\sum 2^{-l_i} \leq 1$

$\Rightarrow \exists$  code  $c$  prefix-free with lengths  $l_1, l_2, \dots$

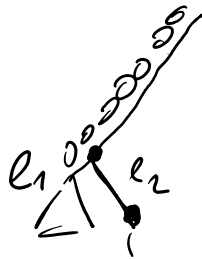
Pf:  $l_1 \leq l_2 \leq \dots$  sort



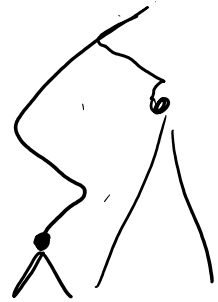


Starting place for interval  $l_i$   
is a multiple of  $l_i$ :

$\cong$  depth-first search



- sort codewords lexicographically



- optimal prefix-free code  $C$  for r.v.  $X$

$$C: X \rightarrow \{0,1\}^*$$

$$L = \mathbb{E}[|c(x)|]$$

$$L = \sum_{x \in X} p(x) \cdot l(x)$$

$$L - H(X) \geq 0$$

$$\rightarrow L - H(X) = \sum_{x \in X} p(x) \cdot l(x) - \sum_{x \in X} p(x) \cdot \log \frac{1}{p(x)}$$

$\downarrow$   $\uparrow$   
 $l(x)$   $\leftarrow$

$$= \sum_{x \in X} p(x) \lg \frac{2^{-l(x)} \cdot p(x)}{p(x)} = (4) \geq 0$$

want

$$c = \sum_x 2^{-l(x)} \leq 1$$

kraft

$$q(x) = \frac{2^{-l(x)}}{c}$$

$$\sum_x q(x) = 1$$

$$(4) = \sum_{x \in X} p(x) \lg p(x) \cdot 2^{+l(x)} \quad D(P \parallel q) \geq 0$$

$$= \sum_x p(x) \lg \frac{p(x)}{c \cdot q(x)} = \sum_x p(x) \cdot \lg \frac{p(x)}{q(x)} + \sum_x p(x) \lg \frac{1}{c} \geq 0$$

$\lg \frac{1}{c} \geq 0$   
 $c \leq 1$

$$L - H(X) \geq 0$$

$$\Rightarrow L \geq H(X) \quad \square$$

Shannon Code:

$$l(x) = \lceil \lg \frac{1}{p(x)} \rceil$$

Ex:

$$p(x) = 2^{-10}$$

$$\rightarrow l(x) = 10$$

$$\sum_x 2^{-l(x)} \leq 1$$

$$\sum_x 2^{-\lceil \lg \frac{1}{p(x)} \rceil} = \sum_x \frac{1}{2^{\lceil \lg \frac{1}{p(x)} \rceil}} \leq \sum_x \frac{1}{2^{\lg \frac{1}{p(x)}}} = \sum_x \frac{1}{\frac{1}{p(x)}}$$

$\Rightarrow \exists$  prefix-free code for  $l(x)$ 's.

$$\leftarrow \frac{1}{P(x)} = \sum P(x) = 1$$

$$\underbrace{\sum p(x) \log \frac{1}{p(x)}}_{H(x)} \leq \sum p(x) \cdot l(x) \leq \underbrace{\sum p(x) \left[ \log \frac{1}{p(x)} + 1 \right]}_{H(x) + 1}$$

$$H(x) \leq L \leq H(x) + 1 \quad \square$$

optimal upto one bit!

Ex:  $\rightarrow p(x_1) = 0.000001 \quad l(x_1) = 20 \quad \overline{1000000}$   
 $p(x_2) = 0.999999 \quad l(x_2) = 1 \quad \underline{0}$

2)  $p(x_1) = \frac{1}{3} \quad l(x_1) = l(x_2) = l(x_3) = 2$   
 $p(x_2) = \frac{1}{3}$   
 $p(x_3) = \frac{1}{3} \quad \underline{\text{opt: } 0, 10, 11}$

Huffman code:  $\leftarrow$  provably optimal  $\underline{\underline{L}}$

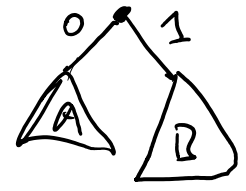
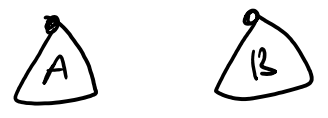
$$p_1 \geq p_2 \geq \dots$$

$$\boxed{p_{n-1} \geq p_n}$$

$$p_{n-1} + p_n$$

$$\underbrace{0 \wedge 1}$$

•  $p_{n-1} \geq p_{n-2}$



Fano code: preserves order of  $p_1, p_2, \dots, p_n$

$$\left[ \underbrace{p_1 + p_2 + \dots + p_k} \mid \underbrace{p_{k+1} + p_{k+2} + \dots + p_n} \right]$$

minimize difference

$$\min_k \left| \sum_{i=1}^k p_i - \sum_{i=k+1}^n p_i \right|$$



proceed recursively

Fact:  $L \leq H(X) + 2$

• 1 bit loss might be too large

pick k build a code for  $(X_1, X_2, \dots, X_k)$

$(x_1, x_2, \dots, x_k)$

$$H(x_1, \dots, x_k) = k H(x)$$

$$L_k = H(x_1, \dots, x_k) + 1$$

using  $\frac{1}{k}$  bits per r.v.

Thm: (McMillan): For any unambiguous  $C$

with code lengths  $l_1, l_2, \dots \leq l_{\max}$

$$\sum 2^{-l_i} \leq 1$$

Pf

$C^k$

... concatenations of  $k$  words from  $C$

$$\left( \sum_{x \in C} 2^{-l(x)} \right)^k = \sum_{x_1 \in C} \sum_{x_2 \in C} \dots \sum_{x_k \in C} 2^{-l(x_1) - l(x_2) - \dots - l(x_k)}$$

$$\sum_{\bar{x} \in C^k} 2^{-l(\bar{x})} = \sum_{m=1}^{k \cdot l_{\max}} w(m) \cdot 2^{-m} = (*)$$

$w(m)$  ... # of words of length  $m$  in  $C^k$

$$w(m) \leq 2^m$$

$$(*) \leq \sum_{m=1}^{k \cdot l_{\max}} 2^m \cdot 2^{-m} = k \cdot l_{\max}$$



$$\Rightarrow \left( \sum_{x \in C} 2^{-l(x)} \right)^k \leq k \cdot l_{\max} \quad k \geq 1$$

$$\sum_{x \in C} 2^{-l(x)} \leq \underbrace{\left( \underbrace{k \cdot l_{\max}}_{\text{fixed}} \right)^{\frac{1}{k}}}_{\xrightarrow{k \rightarrow \infty} 1}$$

$$\Rightarrow \sum_{x \in C} 2^{-l(x)} \leq 1 \quad \square$$

Q:

C code

(not necessarily prefix free)

X r.v.

$$\mathbb{E}[l_C(x)]$$

HW?