

# III. Data compression

## III.1. Symbol codes

Motivation: in natural languages letters come with unequal probabilities (e.g. Hawaiian  $p("a") \approx 26\%$  whereas in Swedish any letter  $\leq 10.5\%$ .)

→ Texts can be compressed by encoding them so that more likely letters are represented by shorter binary strings ...

- Def. 1
- $A^+ := \bigcup_{n \in \mathbb{N}} A^{*n}$  = set of finite strings/words  $a_1 \dots a_n$  with  $a_i \in A, n \in \mathbb{N}$
  - $A^* = A^+ \cup \{\emptyset\}$  ( $\emptyset$  = empty word)
  - Let  $X$  be a random variable with range  $\mathcal{X}$ . A "symbol code" for  $X$  is a map  $C: \mathcal{X} \rightarrow A^+$ .  $C(x)$  is the set of "codewords",  $l_x \in \mathbb{N}$  the respective length and  $L(C) := \sum_{x \in \mathcal{X}} p(x) l_x$  the average length.
  - A symbol code  $C$  is called

(i) "non-singular" iff  $C$  is injective,



(ii) "uniquely decodable" iff the map  $\mathcal{X}^+ \rightarrow A^+$  defined by



$x_1 \dots x_n \mapsto C(x_1) \dots C(x_n)$  is injective,

(iii) "prefix-free" (or sometimes a "prefix code") iff

there is no pair  $x \neq x'$  s.t.  $C(x) = C(x')a$  for some  $a \in A^*$

## Examples:

- non-singular but not uniquely decodable:

$$x_1 \mapsto 0, x_2 \mapsto 1, x_3 \mapsto 01$$

- uniquely decodable but not prefix-free:

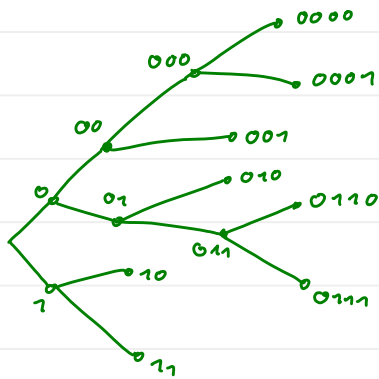
$$x_1 \mapsto 0, x_2 \mapsto 01$$

- prefix-free:  $x_1 \mapsto 0, x_2 \mapsto 10, x_3 \mapsto 110$

- ASCII is prefix-free since all codewords have equal length
- Morse code is not prefix-free (overcome by time gaps between letters)

## Code trees:

A "code tree" for a prefix-free code is an  $|A|$ -ary tree of depth  $\max\{L_x\}_{x \in X}$ ,



so that  $\{\text{leaves}\} = \{\text{codewords}\}$

(for non-prefix free codes intermediate vertices could also correspond to codewords)

## Thm.: (Kraft's inequality)

- For any prefix-free code  $C: X \rightarrow A^+$  with codeword lengths  $L_1, \dots, L_{|X|}$

$$\boxed{\sum_{x \in X} |A|^{-L_x} \leq 1} \quad \text{holds.}$$

- Conversely, given  $L \in \mathbb{N}^{|X|}$  satisfying this inequality there exists a prefix-free code with those codeword lengths.

proof:  $L_m := \max\{L_x\}_{x \in X}$ . Consider the

$|A|$ -ary tree of depth  $L_m$  with  $|A|^{L_m}$  leaves.

If a codeword has length  $L_x$ , this reduces the number of leaves by at least  $(|A|^{L_m - L_x} - 1)$ .

$$\rightarrow |X| \leq |A|^{L_m} - \left( \sum_{x \in X} |A|^{L_m - L_x} - 1 \right)$$

↑  
total # of leaves

$$\Leftrightarrow \sum_{x \in X} |A|^{-L_x} \leq 1$$

Converse: shown by explicit construction via code tree ...

□

Generalization due to McMillan:

Thm.: Let  $\mathcal{X}$  be countable infinite. Every uniquely decodable code satisfies Kraft's inequality.

Conversely, if  $\sum_{x \in \mathcal{X}} |A|^{-l_x} \leq 1$  holds for some  $l_x \in \mathbb{N}$ , then there exists a prefix-free code with these codeword lengths  $(l_x)_{x \in \mathcal{X}}$ .

Def.: A symbol code  $C: \mathcal{X} \rightarrow A^+$  is called "optimal" w.r.t. a random variable  $X$  iff the average codeword length  $L(C)$  is minimal among all symbol codes in  $(A^+)^{\mathcal{X}}$ .

Thm.: (entropy bound for optimal codes)

Let  $C: \mathcal{X} \rightarrow A^+$  be a uniquely decodable code for a random variable  $X$ .

Then  $L(C) \geq \frac{H(X)}{\log |A|}$  with equality iff  $l_x = -\log_{|A|} p(x) \forall x \in \mathcal{X}$ .

Conversely, there exists a prefix-free code for which:

$$L(C) < \frac{H(X)}{\log |A|} + 1$$

proof: lower bound:

$$L(C) - \frac{H(X)}{\log |A|} = \sum_{x \in \mathcal{X}} p(x) \left( l_x - \log_{|A|} \frac{1}{p(x)} \right)$$

$$= \sum_{x \in \mathcal{X}} p(x) \left( \log_{|A|} p(x) - \log_{|A|} (|A|^{-l_x}) \right)$$

$$= \underbrace{\left( \sum_{x \in \mathcal{X}} p(x) \log_{|A|} \frac{p(x)}{q(x)} \right)}_{= D(p||q) \geq 0} - \underbrace{\log_{|A|} \sum_{y \in \mathcal{X}} |A|^{-l_y}}_{\text{Kraft: } \leq 1}$$

$q(x) := \frac{|A|^{-l_x}}{\sum_{y \in \mathcal{X}} |A|^{-l_y}}$

◦ equality requires

$$(i) D(p||q) = 0 \Leftrightarrow \forall x: p(x) = q(x)$$

$$\text{and (ii)} \sum_{x \in \mathcal{X}} |A|^{-L_x} = 1$$

which is equivalent to  $p(x) = q(x) = |A|^{-L_x} \forall x \in \mathcal{X}$

$$\text{and thus } L_x = \log_{|A|} \frac{1}{p(x)} \forall x \in \mathcal{X}$$

◦ upper bound: define a set of integers  $L_x := \lceil \log_{|A|} \frac{1}{p(x)} \rceil$ .

Then  $L \in \mathcal{N}^{|\mathcal{X}|}$  satisfies Kraft's inequality since

$$\sum_{x \in \mathcal{X}} |A|^{-L_x} \leq \sum_{x \in \mathcal{X}} |A|^{-\log_{|A|} p(x)} = 1$$

→ there exists a prefix-free code with codeword lengths  $L$ .

$$L(C) = \sum_{x \in \mathcal{X}} p(x) \lceil \log_{|A|} \frac{1}{p(x)} \rceil < \sum_{x \in \mathcal{X}} p(x) \left( 1 + \log_{|A|} \frac{1}{p(x)} \right)$$

"  $1 + \frac{H(X)}{\log |A|}$  □

remark: the resulting code is called "Shannon code"

Thm. 1 (Entropy rate bounds)

Let  $\{X_i\}_{i \in \mathcal{N}}$  be a stationary stochastic process. Then

(i) for all  $N \in \mathcal{N}$  and all uniquely decodable codes  $C: \mathcal{X}^N \rightarrow \{0,1\}^+$ :

$$H(\{X_i\}) \leq \frac{1}{N} L(C)$$

(ii)  $\forall \varepsilon > 0 \exists N \in \mathcal{N}$  and a prefix-free code  $C: \mathcal{X}^N \rightarrow \{0,1\}^+$  such that

$$\frac{1}{N} L(C) < H(\{X_i\}) + \varepsilon$$

proof: (i)  $H(\{X_i\}) \leq \frac{1}{n} H(X_1, \dots, X_n) \leq \frac{1}{n} L(c)$

$\uparrow$  non-increasing sequence due to stationarity       $\uparrow$  entropy bound

(ii)  $\forall \delta > 0 \exists N \in \mathcal{N} : \frac{1}{n} H(X_1, \dots, X_n) < H(\{X_i\}) + \delta$

$\uparrow$  convergence towards entropy rate due to stationarity

$\exists C: X^n \rightsquigarrow \{0,1\}^+$  such that

$$\frac{1}{n} L(c) < \frac{1}{n} H(X_1, \dots, X_n) + \frac{1}{n} < H(\{X_i\}) + \underbrace{\delta + \frac{1}{n}}_{\varepsilon}$$

$\uparrow$  entropy bound

□