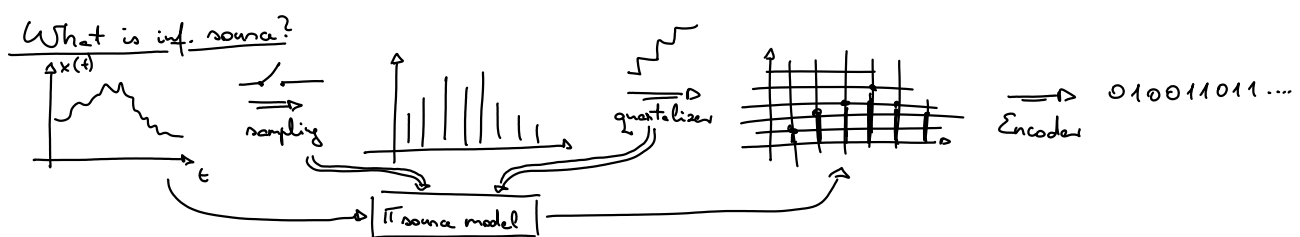
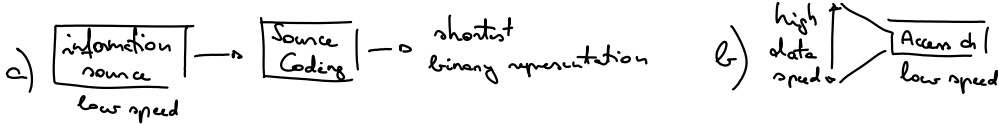


Data Compression - Source Coding

2016. november 9., szerda 16:20

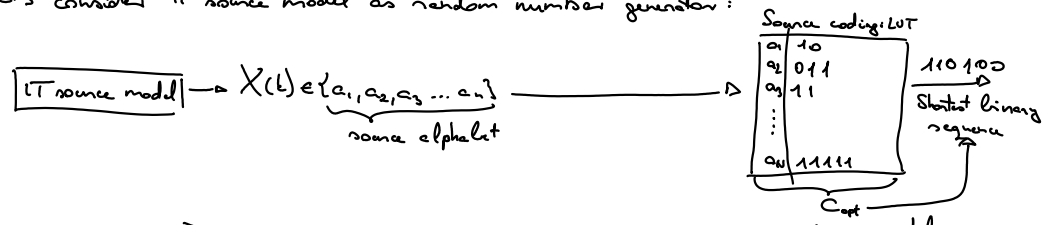


Data speed: $f_s \cdot n$ [bps]

- Speech: $8 \text{ kHz} \cdot 8 \text{ bit} \Rightarrow 64 \text{ kbps}$
- Music: $44.1 \text{ kHz} \cdot 16 \text{ bit} \Rightarrow 705.6 \text{ kbps}$
- Video: 84 Mbps

Compression \Rightarrow GSM: 13.2 kbps

Let's consider IT source model as random number generator:



$$P(X(L) = a_i | X(L-1) = a_1, \dots, X(0) = a_m) = P(X(L) = a_i) \Rightarrow \text{Memoryless model, Stationary}$$

Model:

$p_1, p_2, \dots, p_U \rightarrow p(x)$ ← probability
 $a_1, a_2, \dots, a_U \rightarrow x$ ← symbols
 $s_1, s_2, \dots, s_U \rightarrow c(x)$ ← code
 $l_1, l_2, \dots, l_U \rightarrow l(x)$ ← length of code
 expected value

Average code length: $L := E(l(x)) = \sum_x p(x) l(x) \Rightarrow$ Data speed: $f_s \cdot L$

→ Optimal source coding: $C_{opt} = \min_{\mathcal{C}} L$

Intuitive idea: → if $p(x)$ is large, $l(x)$ is small
 → if $p(x)$ is small, $l(x)$ is large

Questions:

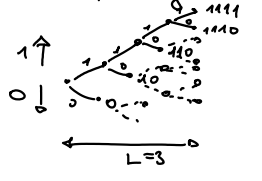
- 1) Unique decoding: how can I decide where the new codeword starts
- 2) Optimal $l(x)$?

⇒ Example: Inredacoffeenymuch ⇒ no breaks or spaces, but we can divide it to words

Why? → because it is a prefix-free code ⇒ This is not: $c_1 = 01$
 $c_2 = 0111$

→ How can we generate prefix-free codings?

⇒ with binary trees



number of missing nodes: $\sum_x 2^{L-l(x)}$
 all possible nodes 2^L


→ missing nodes should be lower than all nodes: $\sum_x 2^{L-l(x)} \leq 2^L$

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

Kraft-McMillan

information $I(x) = \log \frac{1}{p(x)} \Rightarrow H(x) = E I(x) = \sum_x p(x) \log \frac{1}{p(x)}$ ← entropy!!!
 random variable
 we multiply positive numbers, $l(x)$

$\rightarrow 0 \leq H(x) \leq \log N$ so $\pi(x) \geq 0$

Two distributions: 

I-divergence: $D(p||q) := \sum_x p(x) \log \frac{p(x)}{q(x)}$

$D(p,q) \geq 0$

$-D(p,q) = \sum_x p(x) \log \frac{q(x)}{p(x)} \leq \sum_x p(x) \left(\frac{q(x)}{p(x)} - 1 \right) = \sum_x q(x) - \sum_x p(x) = 1 - 1 = 0$

$D(p,q) = \sum_x p(x) \log \frac{p(x)}{1/N} - \sum_x p(x) \log (p(x) \cdot N) = \sum_x p(x) [\log p(x) + \log N] = \sum_x p(x) \log p(x) + \sum_x p(x) \log N = -H(x) + \log N \geq 0 \Rightarrow H(x) \leq \log N$

$\rightarrow p(x) = 1/N \Rightarrow H(x) = \log N$

Source code theorem: $L \geq H(x)$ (introduced by Shannon)

Preal: $p(x)$ source distribution

$q(x) = \frac{2^{-L(x)}}{\sum_y 2^{-L(y)}}$ $0 \leq q(x) \leq 1 \forall x$
 $\sum_x q(x) = 1$
 artificial

$D(p,q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = \sum_x p(x) \log \frac{p(x)}{\frac{2^{-L(x)}}{\sum_y 2^{-L(y)}}} = \sum_x p(x) \log \left\{ p(x) \cdot 2^{L(x)} \cdot \left(\sum_y 2^{-L(y)} \right) \right\} \leq \sum_x p(x) \log \left\{ p(x) \cdot 2^{L(x)} \right\} = -\sum_x p(x) \log p(x) + \sum_x p(x) \log 2^{L(x)} = H(x) + L \geq 0 \Rightarrow H(x) \leq L$

Objective is to approach $H(x)$ as much as it is possible:

Shannon-Fano coding:

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

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

$H(x) \leq L = \sum_x p(x) \lceil \log \frac{1}{p(x)} \rceil \leq \sum_x p(x) (\log \frac{1}{p(x)} + 1) = H(x) + 1$

Algorithm: given $p(x)$

- 1) $l(x) = \lceil \log \frac{1}{p(x)} \rceil, x \in \{a_1, \dots, a_n\}$
- 2) Binary tree $\rightarrow C$
- 3) LUT