

Exempel på entropi.

Teksten 'Then the hen began to eat'  
ble med Huffman koding kodet med 75 bits.  
Teksten består av 25 tegn, altså 3 bits/tegn.  
Hva er entropien for denne teksten?  
(egentlig for alfabetet).

Hvis  $P_i = P(\alpha_i)$  - sannsynligheten til  $\alpha_i$   
er entropien definert som

$$H(P_1, \dots, P_n) = - \sum_{i=1}^n P(\alpha_i) \log_2 P(\alpha_i), \quad 2^{\log_2 x} = x$$

For vår tekst får vi

$$\log_2 x = \frac{\ln x}{\ln 2}$$

$$H(P_1, \dots, P_n) = 2,93$$

Det betyr at vi i beste fall kunne  
klare oss 2,93 · 25 bits, 74 bits

## Aritmetisk koding

Ex. Anta at  $\mathcal{A} = \{0, 1\}$ ,  $P(0) = 0.9$ ,  $P(1) = 0.1$

Huffmankoding gir  $C(0) = 0$ ,  $C(1) = 1$  så vi trenger 1 bit pr. tegn.

$$\begin{aligned} \text{Entropi: } H(P_1, P_2) &= -(P(0) \log_2 P(0) + P(1) \log_2 P(1)) \\ &= 0.47 \end{aligned}$$

Ideen bak aritmetisk koding.

Aritmetisk koding assosierer ulike tekster med ulike delintervall av  $[0, 1]$ . Bredden til delintervallet er lik sannsynligheten til teksten og den aritmetiske koden er ett tall i intervallet, representert i 2-tallsystemet.

Eksempel på aritmetisk koding.

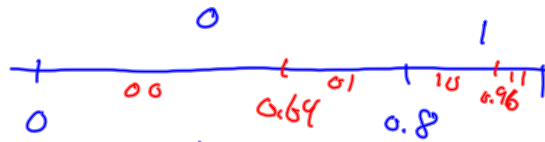
La  $X = 00100$ , da er  $P(0) = 0.8$ ,  $P(1) = 0.2$

Vi legger med  $a$  dele  $[0,1]$  i forhold til sannsynlighetene

Siden første tegn

er 0 ligger koden

i  $[0, 0.8)$



For  $a$  komme videre må vi

dele dette intervallet på samme måte,

i  $[0, 0.64)$  og  $[0.64, 0.8)$ . Siden tegn nr. 2

er 0 harner vi i  $[0, 0.64)$ .

Dette må så deles på samme måte:

i  $[0, 0.512)$  og  $[0.512, 0.64)$ .

Siden tegn 3 er 1

harner vi i

$[0.512, 0.64)$ .



Dette må deles i  $[0.512, 0.6144)$ ,  $[0.6144, 0.64)$

Tegn nr. 4 er 0 så vi harner i

$[0.512, 0.6144)$ .

Dette deles i  $[0.512, 0.59392)$

og  $[0.59392, 0.6144)$

Tegn 5 er 0 så vi ender til slutt

i  $[0.512, 0.59392)$ .

Vi ønsker et tall i 2-tall systemet i

dette intervallet. Tallet i 2-tallssystemet

med minst nerner som ligger

i dette intervallet er  $9/16$

$$\frac{9}{16} = 0.5625 = 0.1001_2 \quad \left(\frac{1}{2} + \frac{1}{16}\right)$$

Aritmetisk kode: 1001

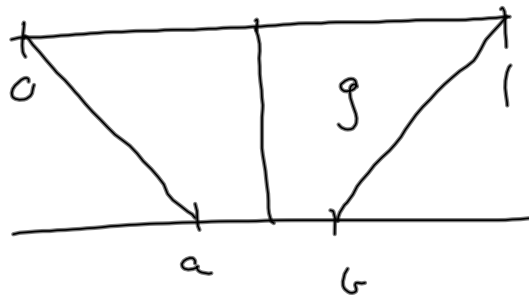
Lineære funktioner:

La  $[a, b]$  være givet med  $a < b$ .

Funktionen  $g(x) = a + x(b-a)$  afbilder

$[0, 1]$  på  $[a, b]$ . Særligt

$$g(0) = a, \quad g\left(\frac{1}{2}\right) = \frac{a+b}{2}, \quad g(1) = b.$$



Kumulative sandsynlighedsfordelinger:

La  $A = \{\alpha_1, \dots, \alpha_n\}$  være et alfabet med sandsynligheder  $p_i = P(\alpha_i)$ .

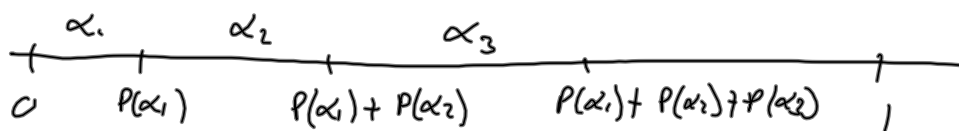
Da defineres ved

$$F(\alpha_j) = \sum_{i=1}^j P(\alpha_i), \quad j=1, 2, \dots, n.$$

$$L(\alpha_j) = F(\alpha_j) - P(\alpha_j) = F(\alpha_{j-1}), \quad j=2, 3, \dots, n.$$

Med  $n$  symboler er første opdeling af  $[0, 1]$

$$[0, F(\alpha_1)], [F(\alpha_1), F(\alpha_2)], [F(\alpha_2), F(\alpha_3)] \dots$$



$\alpha_i$  har intervallet  $[L(\alpha_i), F(\alpha_i)]$

Generell algoritme:

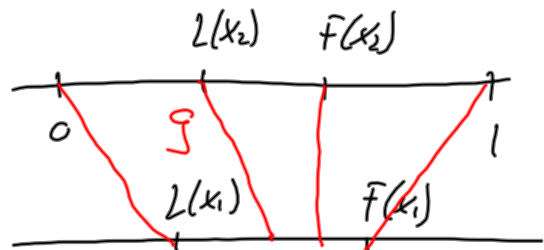
Anta at teksten er  $X = \{x_1, x_2, \dots, x_m\}$

Første tegn er  $x_1$ . Dette svarer til intervallet  $[L(x_1), F(x_1)]$ .

Neste symbol er  $x_2$ . Når skal  $[L(x_1), F(x_1)]$  deles som  $[0, 1]$ .

Så intervallet  $[L(x_2), F(x_2)]$  må plasseres riktig i  $[L(x_1), F(x_1)]$ .

Hå løse  $g$  som sender  $[0, 1]$  på  $[L(x_1), F(x_1)]$ . Da vil intervallet  $[L(x_2), F(x_2)]$  automatisk bli sendt til riktig delintervall av  $[L(x_1), F(x_1)]$  av  $g$



Algoritme:

1. Sett  $[a_0, b_0] = [0, 1]$ .

2. For  $k=1, \dots, m$

(a) Definer  $g_k = a_{k-1} + z(b_{k-1} - a_{k-1})$

(b) Sett  $[a_k, b_k] = [g_k(L(x_k)), g_k(F(x_k))]$

Den aritmetiske koden er midtpunktet  
til  $[a_m, b_m]$  trunkert til

$$\left\lceil -\log_2 [P(x_1) P(x_2) \dots P(x_m)] \right\rceil + 1$$

binære siffer.  $\lceil \cdot \rceil$  - tak funksjonen