

Tillegg til i går. Tallet 129 kan
lagres som

$129 = 128 + 1 = 10000001$,
altså ved å lagre de binære sifrene.

Men også ved å lagre de desimale
sifrene som tegn, '129', for eksempel
i UTF-8

Informations eksplorering

Kap 7.1
i komp.

- 1. Bøger. Typisk bøger har 300 ord pr. side, ca. 4 tegn pr. ord. Med 500 sider giver det ca. 600 000 tegn i hele bøger. Med ISO Latin eller UTF-8 kan vi regne 1 byte pr. tegn, altså tilsammen 600 000 bytes 600 KB, 0.6 MB. Med UTF-16 doubles dette til 1.2 MB.

2. Lyd. På en CD er lyden målt 44100
ganger pr. sekund. Hver måling lagres som
et heltall med 2 bytes. Siden vi har
stereo blir 4×44100 bytes pr. sekund;
tilsammen 176 KB pr. sekund. Det gir
10 MB pr. minutt og omtrent 40 MB per enlåt.
På en CD som inneholder en time musikk har
vi da 600 MB med data.

3. Film og bilder.

TV-bilde består av 576×720 punkter og i hvert punkt er det 24 bits (3 bytes) med farge info. Et bilde krever derfor 1,2 MB. For å få anstendig kvalitet 25 bilder pr. sekund, så ett sekund er da omtrent 31 MB, eller 1,9 GB pr. minutt og 112 GB pr. time.

Huffman koding

Grunnleggende ide. Bruk korte koder for tegn som forekommer ofte og lange koder for tegn som forekommer sjelden.

Noen begreper.

En tekst $X = \{x_1, x_2, x_3, \dots, x_n\}$ er en sekvens av symboler hentet fra et alfabet $\mathcal{A} = \{\alpha_1, \alpha_2, \dots, \alpha_m\}$

Antall ganger α_i forekommer i X angis med frekvensen $f(\alpha_i)$.

For kompressjon gir vi hvert tegn α_i en kode $C(\alpha_i)$. Og Z er den komprimerte teksten der hvert tegn i X er erstattet med sin kode

$$Z = \{C(x_1) C(x_2) C(x_3) \dots C(x_n)\}$$

Ex Anta at $X = DBACDBD$, $\mathcal{A} = \{A, B, C, D\}$

$$f(A) = 1, f(B) = 2, f(C) = 1, f(D) = 3.$$

$$C(D) = 0, C(B) = 1, C(C) = 01, C(A) = 10$$

Vi erstatter tegnene med koder og får

$$Z = 011001010 - 9 \text{ bits}$$

Ex 2. Samme tekst, men kodene

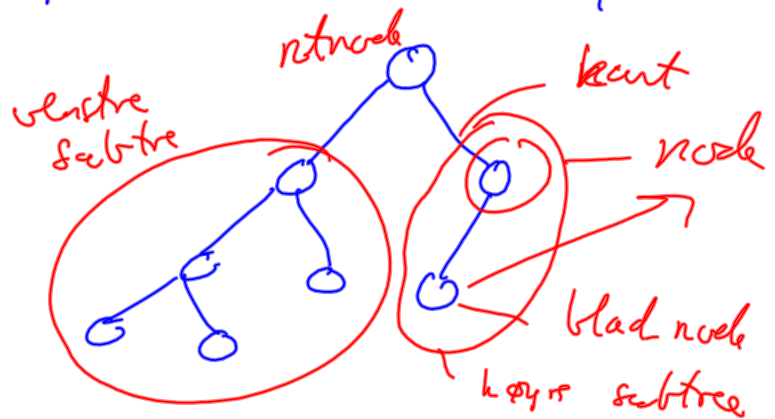
$$C(D) = 1, C(B) = 01, C(C) = 001, C(A) = 000$$

$$Z = 1010000011011 \quad 13 \text{ bits}$$

Binære træer.

Et træ består
for et binært
træ node.

en node og kanten,
går det to kanten fra



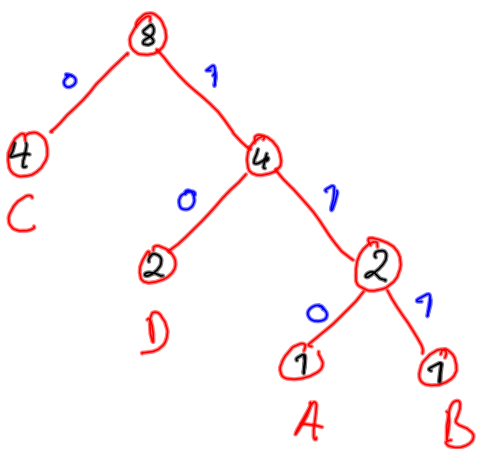
Definición av Huffman tre

Et Huffman tre er et binært tre som er assosiert med et alfabet $\{\alpha_i\}_{i=1}^n$ med frekvenser $f(\alpha_i)$

1. Hver bladnode er assosiert med nøyaktig en α_i
2. Hver node har en vekt
 - (a) Vekten av en bladnode er frekvensen til dens symbol
 - (b) Vekten til en annen node er summen av vektene til røttene i nodens subtrær
3. Alle noder som ikke er bladnoder har nøyaktig to barn.
4. Fra Huffman treet kan vi bygge koder.

Ex $X = C C D A C B D C$, $A = \{A, B, C, D\}$

$f(A) = 1$, $f(B) = 1$, $f(C) = 4$, $f(D) = 2$



$c(C) = 0$, $c(D) = 10$

$c(A) = 110$, $c(B) = 111$

Huffman algoritmen.

Gitt x basert på alfabet $\{\alpha_i\}_{i=1}^n$ med frekvenser $f(\alpha_i)$

1. Lag et en-node Huffman tre med hvert av de n symbolene α_i .

2. Gjenta inntil vi bare har $\alpha_1 \quad \alpha_2 \quad \alpha_3 \quad \dots \quad \alpha_n$ ett tre

(a) Velg to trær T_0 og T_1 med minimal vekt erstall dem med et nytt tre som har T_0 som venstre subtre og T_1 som høyre subtre.

3. Treet som er ~~ig~~ igjen til slutt er Huffman treet.

Huffman-algoritmen gir optimal koding blant binær tre algoritmer med prefix egenskap.