



## Optimale Codes

Nachricht in einem Quellalphabet  $\Sigma_1 = \{a, b, c, A, L, !, \dots\}$   
 Codealphabet =  $\{0, 1\}$  feste Länge 8 Bits

ASCII

$a \rightarrow 01100001$

Gesucht ist ein Code  $C: \Sigma_1 \rightarrow \{0, 1\}^*$

- eindeutig entzifferbar:

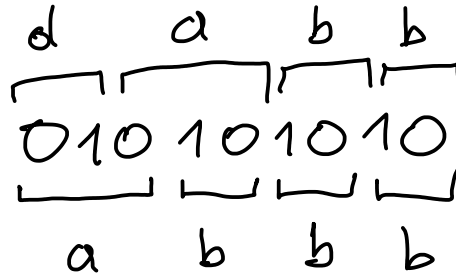
Jede codierte Nachricht  $\in \{0, 1\}^*$  darf sich auf höchstens eine Art in Codewörter zerlegen lassen.

$C_a = 010$

$C_b = 10$

$C_c = 11$

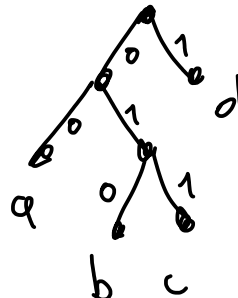
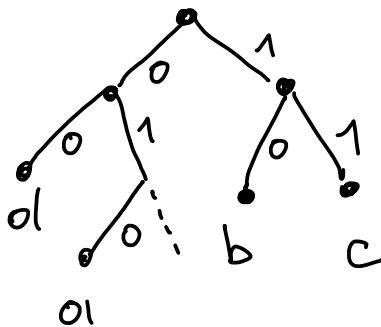
$C_d = 0100$



- stärkere Forderung: präfixfrei.

Kein Codewort ist Präfix eines anderen Codewortes

Codebaum



$l_a = 2 \quad P_a = 20$   
 $l_b = 3 \quad P_b = 17$   
 $l_c = 3 \quad P_c = 14$   
 $l_d = 1 \quad P_d = 77$

$\frac{1}{4} + \frac{1}{8} + \frac{1}{8} + \frac{1}{2} = 1$

Gegeben: Häufigkeiten  $p_1, p_2, \dots, p_n$  für die  $n$  Zeichen  $\in \Sigma$

Gesucht: Ein präfixfreier Code mit Codewortlängen  $l_1, l_2, \dots, l_n$

Gesamtlänge  $\sum_{i=1}^n p_i l_i \rightarrow \text{MIN}$

Äquivalent: Ein Codebaum mit  $n$  Blättern auf Tiefe  $l_1, l_2, \dots, l_n$ .

•  $p_1, \dots, p_n$  können auch relative Häufigkeiten / Wahrscheinlichkeiten sein.

Algorithmus von Huffman:

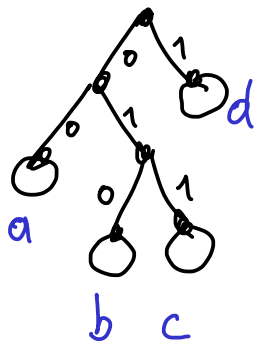
(1)  $p_i < p_j \Rightarrow l_j \leq l_i$

Bew. Ann.  $p_i < p_j \quad l_j > l_i$

$$\Sigma_i = p_i l_i + \dots + p_j l_j$$

$$\Sigma_i' = p_i l_j + \dots + p_j l_i$$

$$\Sigma_i - \Sigma_i' = \dots > 0$$



$l_1 = 1$   
 $l_2 = 2$   
 $l_3 = 3$   
 $l_4 = 3$

$p_a = 20$   
 $p_b = 17$   
 $p_c = 14$   
 $p_d = 77$

(1b) Ein optimaler Codebaum muss voll sein: jeder innerer Knoten hat 2 Kinder.

(2) In jedem vollen Binärbaum gibt es zwei Geschwisterknoten  $i, j$  mit  $l_i = l_j \geq l_k$  für alle anderen  $k$  (auf der tiefsten Ebene)

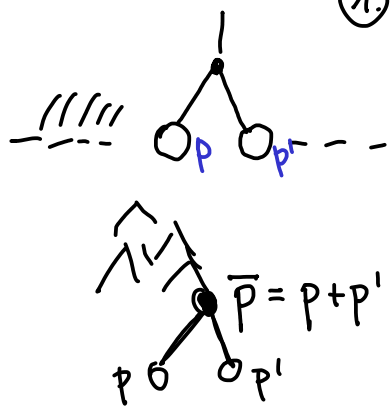
(1)+(2)

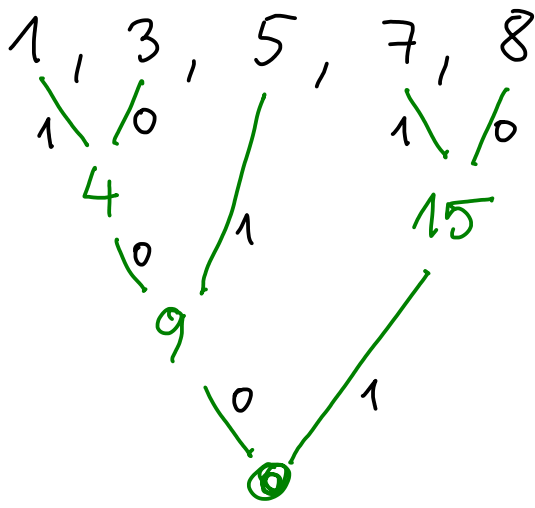
$p, p'$  seien die kleinsten Häufigkeiten.

Wir können annehmen, dass  $p$  und  $p'$  zwei Geschwisterknoten zugeordnet werden.

Ersetze  $p, p'$  durch  $\bar{p} = p + p'$ .

Konstruiere einen opt. Baum für die  $n-1$  entstehenden Häufigkeiten.





1:	001
3:	000
5:	01
7:	11
8:	10

Implementierung: Prioritätswarteschlange  $O(n \log n)$

Die Kraft-MacMillan'sche Ungleichung

$$(*) \quad \sum_{i=1}^n 2^{-l_i} \leq 1$$

ist notwendig und hinreichend für

- (a) die Existenz eines Binärbaums mit Blättern auf Tiefe  $l_1, \dots, l_n$   
(vgl. Aufgabe 44)
- (b) die Existenz eines präfixfreien Codes mit Wortlängen  $l_1, \dots, l_n$
- (c) die Existenz eines ~~präfixfreien~~ *eindeutig entzifferbaren* Codes mit Wortlängen  $l_1, \dots, l_n$

Satz von McMillan (1956)  $(c) \Rightarrow (*)$

$$C = \{00, 0101, 1111, 0111, 001100, 010010\} \quad \frac{1}{2^2} + \frac{1}{2^4} + \frac{1}{2^4} + \frac{1}{2^4} + \frac{1}{2^6} + \frac{1}{2^6} \leq 1$$

$l_1=2 \quad l_2=4 \quad l_3=4 \quad l_4=4 \quad l_5=6 \quad l_6=6$

$$L = C^* \quad A_k = |L \cap \{0,1\}^k| = \text{Anzahl der Nachrichten der Länge } k, \text{ die sich aus dem Codewörtern bilden lassen}$$

$$A_k \leq 2^k$$

$$A_0 = 1$$

$$A_k = 0, \quad k < 0$$

$$A_k = A_{k-2} + 3 \cdot A_{k-4} + 2 A_{k-6} \quad (k \geq 1)$$

$$A_k = \sum_{i=1}^n A_{k-l_i}$$

$$[A_k \sim y_0^k]$$

Annahme:  $1 < \sum_{i=1}^n 2^{-l_i}$

Finde  $y > 2$  mit:

$$1 < \sum_{i=1}^n y^{-l_i}$$

$$(l_1 \leq l_2 \leq \dots \leq l_n)$$

Behauptung:  $\bar{A}_k \geq c_0 \cdot y^k$  für  $k \geq 0$  und für eine Konstante  $c_0 > 0$

vollst. Induktion:

$$\bar{A}_k = \sum_{i=1}^n \bar{A}_{k-l_i} \quad (\text{für } k \geq l_n)$$

$$\stackrel{\text{I.V.}}{\geq} \sum_{i=1}^n c_0 y^{k-l_i} = c_0 \cdot y^k \underbrace{\sum_{i=1}^n y^{-l_i}}_{> 1} > c_0 y^k$$

Induktionsbasis: Wähle  $c_0 > 0$  klein genug sodass  $c_0 y^k \leq \bar{A}_k$  für  $k=0, 1, \dots, l_n-1$

$$\bar{L} = \{\epsilon, F, FF, \dots, F^{l_n-1}\} \cdot \mathbb{C}^*$$

$$\bar{A}_k = |\bar{L} \cap \{\epsilon, F\}^k| > 0$$

Bsp.:  $F \underbrace{00001100}_{\in \bar{L}}$   
 $\underbrace{\hspace{1.5cm}}_{\in \bar{L}}$

$$\bar{A}_k \leq 2^k + 2^{k-1} + 2^{k-2} + \dots + 2^{k-l_n+1} \leq 2^k \cdot 2$$

$$c_0 y^k \leq \bar{A}_k \leq 2 \cdot 2^k$$

$$\underbrace{\left(\frac{y}{2}\right)^k}_{> 1} \leq \frac{2}{c_0}$$

Widerspruch  $\square$