

Untercheidung Code feste Länge
 und solche mit variabler Länge

Feste Länge Codierung abgebildet
 $C: A \rightarrow B^n$

Var Länge Vorteil Häufig benutzte Zeichen
 können mit kürzerem Wert codiert. Beispiel

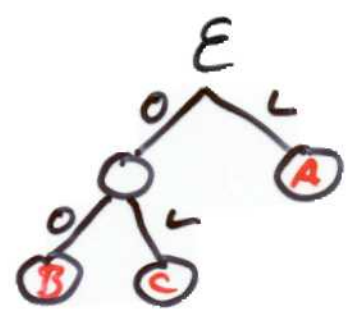
e im Deutschen
 \Rightarrow Gesamtlänge des zu codierten Textes
 minimal

Nachteil Bei Serienwort codierung ist u.U.
 die Trennung der Einzelwörter schwierig,
 wie besonders bei Störungen

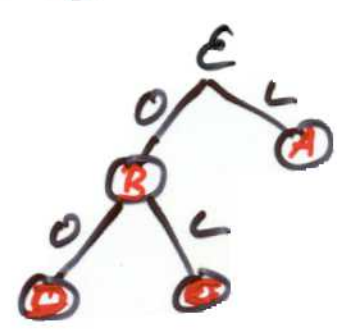
Serienwort codierung $C^*(\langle a_1 \dots a_n \rangle) =$
 $c(a_1) \circ c(a_2) \circ \dots \circ c(a_n)$

Beispiel
 a) $A \rightarrow L$
 $B \rightarrow 00$
 $C \rightarrow 0L$

Darstellung als 'Codebaum'



b) $A \rightarrow L$
 $B \rightarrow 0$
 $C \rightarrow 0L$
 $D \rightarrow 00$



Zu kodieren Zeichensequenz AA BCDA

LL O O L O O L; bei Kodierung
möglich: A A B B A B B A, aber auch
A A D A D A ..

Im Beispiel ~~zu~~ b) ist eine (hinreichend)
Voraussetzung zur Eindeutigkeit der Codierungs-
abbildung verletzt: Kein Codewort (CW)
darf identisch mit dem Anfang eines anderen
CW sein \Rightarrow Im Codebaum viel CW nur
an der Blätter ("Fano-Bedingung")

Erklärung unter Berücksichtigung der Auftritts-Wk

Ziel: gegebene Nachricht einer Quelle soll
durch Folge von CW codiert werden, ohne
Gesamtlänge kennt.

Plausibles Vorgehen: Berücksichtigung der
Häufigkeit, mit der ein bestimmtes Zeichen
von der Quelle emittiert wird.

Ob das Auftritts-Wk hin gegen nichts
bekannt, bleibt nur, allen Zeichen gleichlange

~~zu~~ CW zuzuordnen. Beispiel: ASCII-Code
mit 256 Zeichen, jedes Zeichen wird mit
ld 256 = 8 bit

Allgemein gilt: Um ein Zeichen aus einem
Zeichenvariablen A mit M Elementen zu
codieren, werden bei Binärcodierung
ld M bit benötigt; denn mit J bit lassen

sich \mathcal{L} auch darstellen

Begriffe

- a) Die Auftretens-WK eines Zeichens a ist die Anzahl n des Auftretens des Zeichens in einer Zeichenfolge der Länge N für $N \rightarrow \infty$, also $p_a = \lim_{N \rightarrow \infty} \frac{n}{N}$
- b) Der Entscheidungsgehalt eines Zeichens ist die Anzahl der Schritte, die nötig ist, ein Zeichen aus dem Zeichen vorat zu isolieren.
- c) Der Informationsgehalt eines Zeichens steht im umgekehrten Verhältnis zu seiner Auftretens-WK. Er wird definiert aus:
$$I_a := \log_2 \frac{1}{p_a} = -\log_2 p_a \quad \text{in bit}$$
- d) Der von einem Zeichen der Quelle zu erwartende mittlere Informationsgehalt wird als Entropie bezeichnet. ("Überschusswerthaltigkeit") Sie ist der ~~mittlere~~ Erwartungswert $H = E\{I_i\} = -\sum_{i=1}^N p_i \log_2 p_i$
- Je höher die Entropie der Quelle, desto größer der Informationsgehalt der von ihr ausgehenden Nachrichten. Entropie erreicht ihr

wahrscheinlich sind und ist dann $H_0 = I_i = \text{ld} M$
 die Entropie ist niedrig, wenn die Quelle nur
 das gleiche Zeichen ausgibt. In diesem Fall
 reichen wenige Bits, um Nachricht zu
 repräsentieren

Beispiel: • Quelle mit Zeichen '1' und '0'
 • Mit Auftritts - Wk p wird '1' gesendet,
 mit Wk $(1-p)$ wird '0' gesendet

Dann ist die Entropie

$$H = -p \text{ld} p - (1-p) \text{ld} (1-p)$$

$$\text{Für } p = 0,5 : H = 0,5 \cdot 1 + 0,5 \cdot 1 = 1$$

$$p = 0,25 : H \approx 0,8$$

Die Länge der CW wird klein, wenn
 jedes Bit den maximalen Informations-
 gehalt transportiert, wenn also für jedes Bit
 die Wk für '1' bei $p=0,5$ liegt. → Bit
 für Bit "maximale Überraschung".

Beim "Gang auf dem Codebaum" ergibt jedes
 Bit den Spielraum für noch mögliche Zeichen
 aus der Menge der CW muss weiter ein.

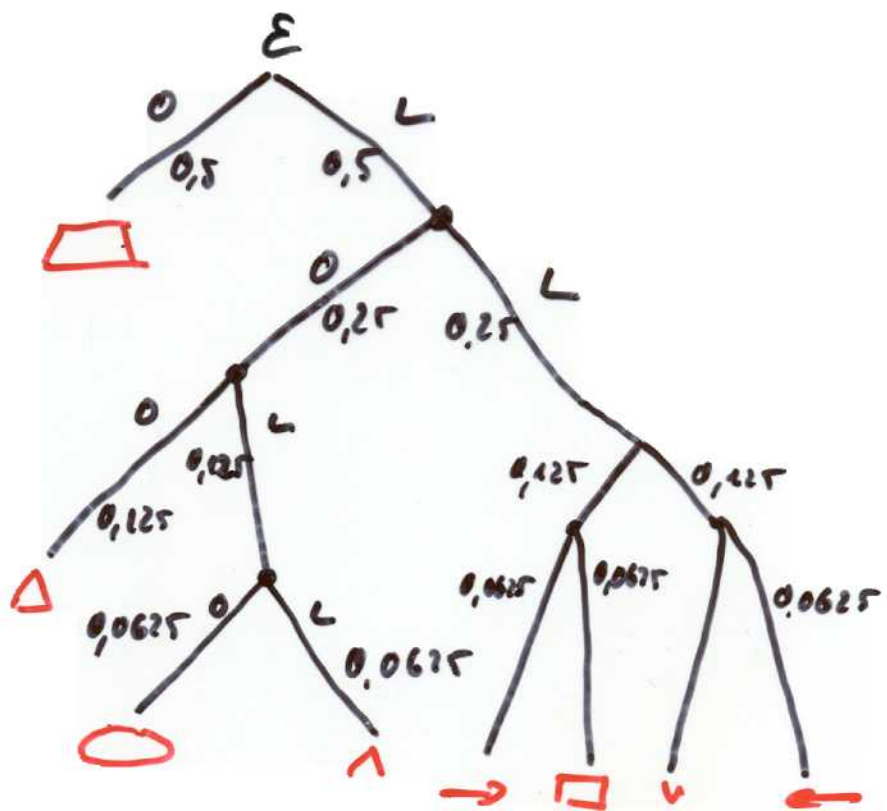
Bei einem (näherungsweise) optimalen Code
 muss daher jedes weitere Bit die noch ver-

gleich wahrscheinliche Teilmengen zerlegen

Beispiel 8 Zeichen seien zur Codierung bei bekannten Auftretens WK

Zeichen i	□	△	▣	○	∧	∨	→	←
P_i	0,5	0,1	0,05	0,03	0,09	0,08	0,02	0,08

Eine Möglichkeit zur Codierung Konstruktion eines Codebaums von der Wurzel her übergeht das sind die WK für die Restmengen der noch verbleibenden CW an jedem Knoten gleich auf die beiden Teilbäume verteilt Zuordnung der CW so dass ihre Auftretens WK möglichst nahe an der aufgeteilten "WK Masse" kommt.



Demm auf Code tabelle

Länge
des CW

Zichen	p_i	CW	$p_i \cdot L_i$
□	0,5	0	0,5
^	0,1	L00	0,3
▭	0,05	LL0L	0,2
○	0,03	L0LO	0,12
^	0,09	L0LL	0,36
v	0,08	LLLO	0,32
→	0,07	LL00	0,28
←	0,08	LLLL	0,32

$$\Sigma = 2,4 \text{ bit / Zeichen}$$

Un befriedigend: Abwang der geforderten WK bei der Feiserteilung im den Codebaum gegen uber dem tatsachlichen WK. Dies lat die letzte Sicherheit uber die tatsachliche Wahl des Optimums vermissen.

Deshalb bessere Moglichkeit: Konstruktieren des Baums "bottom-up" ausgehend von der Beobachtung, da die seltensten Zeichen am weitesten unten im Baum angeordnet sind. → siehe zur Kontrolle ubung

Literatur

- Original paper von C. Shannon von 1948 → Homepage
- Heise und Quattrocchi Springer, 4 Auflage

Finite Klasse von Kompressionsalgorithmen

Beobachtung: bestimmte Zeichenfolgen häufiger als andere.

Lexikon-basierte Codierung

Literatur:

Bell et. al. Modeling for Text Compression,
ACM Computing Surveys, Vol. 21, No. 4,
Dec. 89

Ersetzung von Zeichengruppen durch Ziffern auf "Codebuch"-Eintrag

Problem: a) Welche Zeichengruppen kommen in das Lexikon?

b) Wie wird der Text in Zeichengruppen aufgeteilt?