

ADS: Algorithmen und Datenstrukturen 2

Teil 1

Prof. Dr. Gerhard Heyer

Institut für Informatik
Abteilung Automatische Sprachverarbeitung
Universität Leipzig

3. April 2019

[Letzte Aktualisierung: 03/04/2019, 08:03]

Einführung in allgemeine *Entwurfsprinzipien* für Algorithmen, die sich für viele Problemstellungen anwenden lassen.

- Greedy-Algorithmen
- Dynamische Programmierung
- Probabilistische und heuristische Algorithmen (für NP-schwierige Probleme)
- Maschinelles Lernen

Anwendungsschwerpunkte

- Datenkompression
- Graphen und Netzwerke
- Rucksackproben, Travelling Salesman, ...

Zwei Ziele bei der Formulierung von Algorithmen

- Effizienz bezüglich *Zeit* ODER
- Effizienz bezüglich *Speicherplatz*

Im folgenden:

- Datenkompression mit Hilfe effizienter *Kodierung*
- für eine *Speicher*-effiziente Darstellung

Zwei grundsätzlich unterschiedliche Typen der Kompression

- **verlustfrei** (lossless)
Codierung eindeutig umkehrbar.
- **verlustbehaftet** (lossy)
Originaldaten können i.A. nicht eindeutig aus der komprimierten Codierung zurückgewonnen werden.
Beispiele: JPEG, MPEG (mp3)

- **Laufängenkodierung:** Aufeinanderfolgende identische Zeichen werden zusammengefasst.
- **LZW, LZ77, gzip:** Mehrfach auftretende Zeichenfolgen werden indiziert (*Wörterbuch/Codetabelle*), um dann durch ein einzelnes Zeichen dargestellt zu werden.
- **Huffman-Kodierung:** Häufiger auftretende Zeichen werden mit weniger Bits codiert.
- **Burrows-Wheeler:** Zeichenketten werden in besonders gut komprimierbare Zeichenketten vorverarbeitet.

Allgemeines Prinzip:

Ausnutzen von **Redundanz**

(Abweichungen von zufälliger Zeichenfolge, wiederkehrende Muster)

- Einfachster Typ von Redundanz:
Läufe (runs) = lange Folgen sich wiederholender Zeichen.
- Beispiel:
AAAABBBBAABBBBBCCCCCCCCDABCBAABBBBCCCD.

- Einfachster Typ von Redundanz:
Läufe (runs) = lange Folgen sich wiederholender Zeichen.
- Beispiel:
AAAABBBBAABBBBBCCCCCCCCDABCBAABBBBBCCCD.
- Ersetze jeden Run durch seine Länge und das wiederholte Zeichen:
4A3B2A5B8CDABC3A4B3CD
- Lohnt sich nur für Läufe mit Länge > 2 .

**Wie können wir Text-Zeichen von Längenangaben unterscheiden?
Insbesondere, wenn alle Zeichen (auch Ziffern) im zu codierenden
Text vorkommen dürfen.**

- Wähle als **Escape-Zeichen** ein Zeichen Q, das in der Eingabe wahrscheinlich nur selten auftritt.
- Jedes Auftreten dieses Zeichens besagt, dass die folgenden beiden Buchstaben ein Paar (Zähler, Zeichen) bilden.
- Ein solches Tripel wird als *Escape-Sequenz* bezeichnet.
- Zählerwert i wird durch i -tes Zeichen des Alphabets dargestellt.
- Codierung erst für Läufe ab Länge 4 sinnvoll.
- Auftreten des Escape-Zeichens Q selbst muss speziell codiert werden, z.B. als Run der Länge 1 oder $Q_{_}$ ($_ =$ Leerzeichen), etc.
- Beispiel:

Text **AAAABBBBAABBBBBCCCCCCCCDABCQC.**

Laufängencodierung: Codierung mit Escape-Zeichen

- Wähle als **Escape-Zeichen** ein Zeichen Q, das in der Eingabe wahrscheinlich nur selten auftritt.
- Jedes Auftreten dieses Zeichens besagt, dass die folgenden beiden Buchstaben ein Paar (Zähler, Zeichen) bilden.
- Ein solches Tripel wird als *Escape-Sequenz* bezeichnet.
- Zählerwert i wird durch i -tes Zeichen des Alphabets dargestellt.
- Codierung erst für Läufe ab Länge 4 sinnvoll.
- Auftreten des Escape-Zeichens Q selbst muss speziell codiert werden, z.B. als Run der Länge 1 oder $Q_{_}$ ($_ =$ Leerzeichen), etc.
- Beispiel:
Text AAAABBBBAABBBBBBCCCCCCCCDABCQC.
Codierung QDABBBAAQE BQH CDABCQ_ $_C$

Für Sequenzen von Bitwerten $\in \{0, 1\}$:

- Beobachtung: Läufe von 0 und 1 wechseln sich ab.
- Nutze das aus: gebe nur noch Laufängen an.
- Beispiel:

Sequenz	0001110111101100011111
Kodierung	3 3 1 4 1 2 3 5

Was ist, wenn die Sequenz mit 1 beginnt?

- Ansatzpunkt für Kompression: In Sprache kommen Zeichen mit unterschiedlicher Häufigkeit vor, z.B. **E** häufiger als **Y**.
- Beispiel: **ABRACADABRA**
- Standardcodierung (unkomprimiert) mit 5 Bits pro Zeichen:
00001 00010 10010 00001 00011 00001 00100 00001 00010 10010 00001
- Häufigkeit der Buchstaben:

A	B	C	D	R
5	2	1	1	2
- **Idee:** verwende für häufige Zeichen kurze Bitsequenzen, für seltene dafür längere. Minimiere so die Gesamtzahl der benötigten Bits.

Wie kann man, trotz unterschiedlich langer Bitfolgen, erkennen welche Bits welches Zeichen codieren?

Gegeben sei folgende Codierung mit 0 und 1 (bzw. kurz oder lang wie im Morsealphabet):

$$E = 0$$

$$T = 1$$

$$A = 01$$

Was bedeutet 0101?

Buchstabenkombination ETA AA ETET AET

Ohne Sonderzeichen keine eindeutige Zerlegung!

Präfixcode (Präfix-freier Code)

- Kein Codewort ist Präfix eines anderen Codeworts.
- Durch Präfix-Freiheit: Codierung eindeutig umkehrbar.
- Im Beispiel (ABRACADABRA):

Buchstabe	A	B	C	D	R
Häufigkeit	5	2	1	1	2
Code	0	100	1010	1011	11

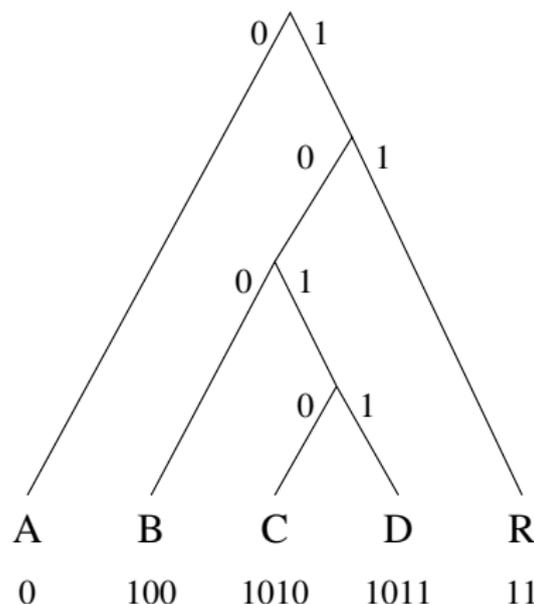
Codierung:

01001101010010110100110

AB R AC AD AB R A

Binärer Präfixcode lässt sich durch Binärbaum (Trie) darstellen:

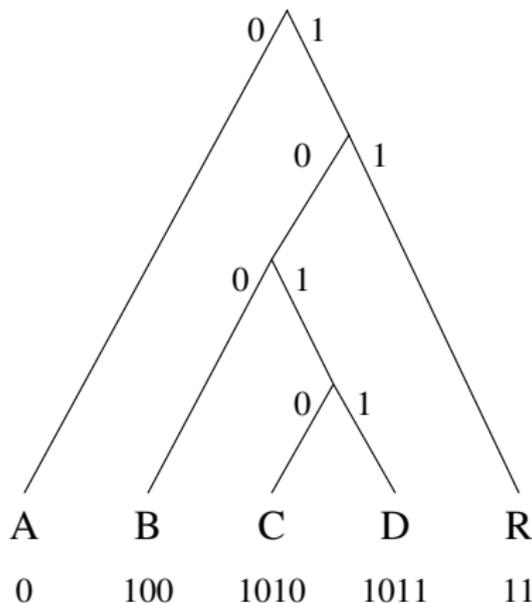
- **Blätter** = zu codierende Zeichen
- Zur **Codierung eines Zeichens**:
laufe von Wurzel zum Blatt des Zeichens und gebe Kantenlabel aus (0 bei Verzweigung nach links; 1, nach rechts)
- übliche Bezeichnung: **Trie**
(dazu später mehr)



Dekodierung (top-down)

Gegeben: Bitsequenz und Präfixcode als Trie

- 1 Beginne bei Wurzel im Baum
- 2 Wiederhole bis zum Ende der Bitsequenz
 - 1 Lies nächstes Bit. Bei 0 verzweige nach links im Baum, sonst nach rechts.
 - 2 Falls erreichter Knoten ein Blatt ist, gib das Zeichen des Blatts aus und springe zurück zur Wurzel.



Der *Huffman-Code* ist ein präfix-freier Code mit variabler Länge, der die unterschiedlichen Zeichen-Häufigkeiten ausnutzt und systematisch erzeugt werden kann.

Beispiel

“A SIMPLE STRING TO BE ENCODED USING
A MINIMAL NUMBER OF BITS”

Dazugehörige Häufigkeits-Tabelle

	_	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
k	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26
count[k]	11	3	3	1	2	5	1	2	0	6	0	0	2	4	5	3	1	0	2	4	3	2	0	0	0	0	0

Erzeugung des Huffman-Codes

Schritt 1: Zähle Häufigkeit der Zeichen in der zu codierenden Zeichenfolge

Das folgende Programm ermittelt die Buchstaben-Häufigkeiten einer Zeichenfolge `a` und trägt diese in ein Feld `count` ein.

```
for ( i = 0 ; i <= 26 ; i++ ) count[i] = 0;
for ( i = 0 ; i < a.length ; i++ ) count[index(a[i])]++;
```

Dabei liefert `index(c)` den Index eines Zeichens `c`:

`index(A)=1, index(B)=2, ...; index(_):=0`

Schritt 2: Aufbau des Tries (entsprechend den Häufigkeiten)

Für jedes Zeichen mit Häufigkeit $\neq 0$, erzeuge einen Baum jeweils mit dem Zeichen als einzigem Knoten.

Schritt 3: Kombiniere iterativ die Bäume mit kleinsten Häufigkeiten

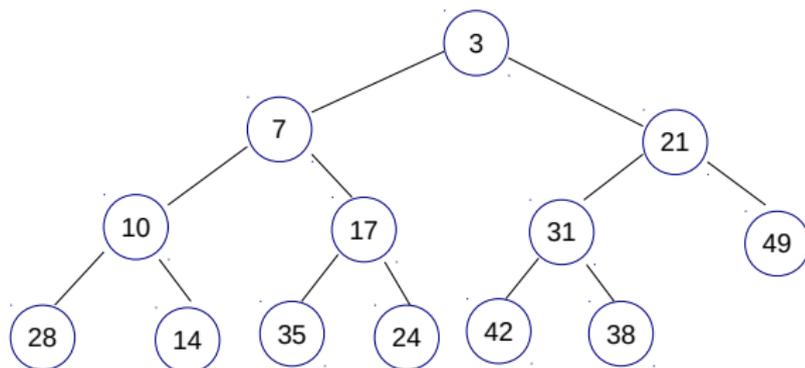
- Wähle zwei Bäume mit minimalen Häufigkeiten (Die Häufigkeit eines Baums ist die Summe der Häufigkeiten aller seiner Zeichen. Während der Erzeugung speichern wir Häufigkeit jedes Baums ab.)
- Verwende hierfür als Zwischenstruktur einen *MinHeap*
- Erzeuge neuen Knoten, der diese beiden Bäume als Nachfolger hat
- Iteriere Schritt 3 bis alle Knoten miteinander zu einem einzigen Baum verbunden sind.

Am Ende sind Zeichen mit geringen Häufigkeiten weit unten im Baum; Zeichen mit großen Häufigkeiten nah an der Wurzel.

- Datenstruktur Heap („Halde“) als Spezialfall eines binären Baumes
- **Ein Heap ist ein Binärbaum**, der
 - die Heapeigenschaft hat (Kinder sind größer/kleiner als der Vater)
 - bis auf die letzte Ebene vollständig besetzt ist
 - höchstens eine Lücke in der letzten Ebene hat, die ganz rechts liegt
- Heaps i. A. genutzt, um den Datentyp *Vorrangwarteschlange* (engl. „priority queue“) zu implementieren
- Folgende Operationen sollen dabei unterstützt werden:
 - Lesen des Objektes mit der kleinsten/größten Priorität
 - Löschen des Objektes mit der kleinsten/größten Priorität
 - Einfügen eines neuen Elements mit beliebiger Priorität

Aus der Definition folgt, dass

- ein Teilbaum eines Heaps wiederum ein Heap ist
- in einem beliebigen Pfad im Heap die Elemente sortiert sind
- Je nach dem, ob bei der Sortierung die Eigenschaft *größer* oder *kleiner* betrachtet wird, spricht man von *minHeap* oder *maxHeap*



Typische Operationen auf Priority Queues:

- isempty, insert, readmin, delmin

Für das **Einfügen**:

- hänge das einzufügende Element an den freien Knoten der vorletzten Ebene
- stelle die Heap-Eigenschaft wieder her

Für das **Entnehmen** eines minimalen/maximalen Elements:

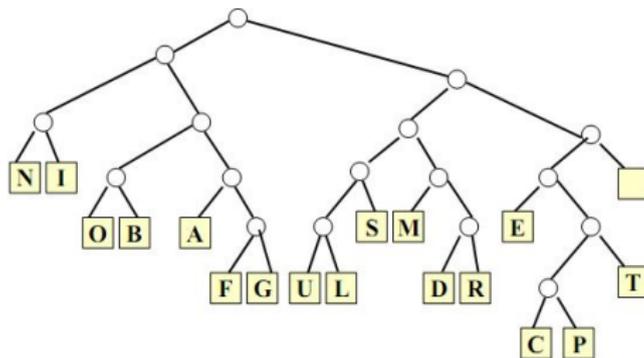
- lösche das Wurzelement
- füge das letzte Blatt an der Wurzel ein
- stelle die Heap-Eigenschaft wieder her (*heapify*)

Aufwand:

- Min-Lesen $O(1)$
- Min-Löschen $O(\log n)$
- Einfügen $O(\log n)$

Trie für die Huffman-Codierung von "A SIMPLE STRING ..."

	_	A	B	C	D	E	F	G	I	L	M	N	O	P	R	S	T	U
k	0	1	2	3	4	5	6	7	9	12	13	14	15	16	18	19	20	21
count[k]	11	3	3	1	2	5	1	2	6	2	4	5	3	1	2	4	3	2

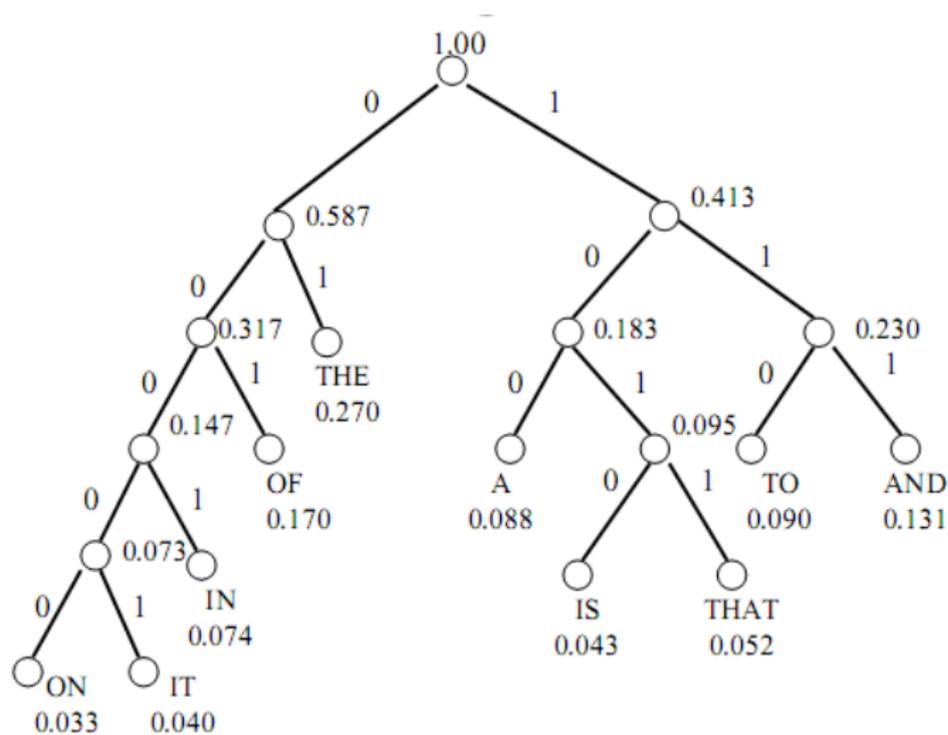


Ableitung des eigentlichen Codes aus dem Trie: Kantenlabel 0 nach links und 1 nach rechts, lese Code eines Zeichens auf Pfad von Wurzel zu Zeichen ab.

Huffman-Codierung für Wörter der Englischen Sprache

Codierungs-Einheit	Wahrscheinlichkeit des Auftretens	Code-Wert	Code-Länge
the	0.270	01	2
of	0.170	001	3
and	0.137	111	3
to	0.099	110	3
a	0.088	100	3
in	0.074	0001	4
that	0.052	1011	4
is	0.043	1010	4
it	0.040	00001	5
on	0.033	00000	5

Huffman-Codierungs-Baum für englische Wörter



Häufigkeitsverteilung deutscher Bigramme

<i>en</i>	4.47	<i>nw</i>	0.55	<i>ab</i>	0.25	<i>nm</i>	0.16
<i>er</i>	3.40	<i>us</i>	0.54	<i>il</i>	0.25	<i>pe</i>	0.16
<i>ch</i>	2.80	<i>nn</i>	0.53	<i>mm</i>	0.25	<i>rl</i>	0.16
<i>nd</i>	2.58	<i>nt</i>	0.52	<i>nz</i>	0.25	<i>sm</i>	0.16
<i>ei</i>	2.26	<i>ta</i>	0.51	<i>sg</i>	0.25	<i>sp</i>	0.16
<i>de</i>	2.14	<i>eg</i>	0.50	<i>sw</i>	0.25	<i>th</i>	0.16
<i>in</i>	2.04	<i>eh</i>	0.50	<i>rn</i>	0.24	<i>wo</i>	0.16
<i>es</i>	1.81	<i>zu</i>	0.50	<i>ro</i>	0.24	<i>af</i>	0.15
<i>te</i>	1.78	<i>al</i>	0.49	<i>ea</i>	0.23	<i>lu</i>	0.15
<i>ie</i>	1.76	<i>ed</i>	0.48	<i>fr</i>	0.23	<i>mu</i>	0.15
<i>un</i>	1.73	<i>ru</i>	0.48	<i>sd</i>	0.23	<i>no</i>	0.15
<i>ge</i>	1.68	<i>rs</i>	0.47	<i>tt</i>	0.23	<i>nv</i>	0.15
<i>st</i>	1.24	<i>ig</i>	0.45	<i>tw</i>	0.23	<i>rf</i>	0.15
<i>ic</i>	1.19	<i>ts</i>	0.45	<i>gr</i>	0.22	<i>ut</i>	0.15
<i>he</i>	1.17	<i>ma</i>	0.43	<i>tz</i>	0.22	<i>br</i>	0.14
<i>ne</i>	1.17	<i>sa</i>	0.43	<i>fe</i>	0.21	<i>ez</i>	0.14
<i>se</i>	1.17	<i>wa</i>	0.43	<i>gt</i>	0.21	<i>ho</i>	0.14
<i>ng</i>	1.07	<i>ac</i>	0.42	<i>rh</i>	0.21	<i>ka</i>	0.14
<i>re</i>	1.07	<i>eu</i>	0.42	<i>ds</i>	0.20	<i>os</i>	0.14
<i>au</i>	1.04	<i>so</i>	0.41	<i>du</i>	0.20	<i>bl</i>	0.13
<i>di</i>	1.02	<i>ar</i>	0.40	<i>mi</i>	0.20	<i>dw</i>	0.13
<i>be</i>	0.96	<i>tu</i>	0.40	<i>nb</i>	0.20	<i>ep</i>	0.13
<i>ss</i>	0.94	<i>ck</i>	0.37	<i>nk</i>	0.20	<i>hm</i>	0.13

Häufigkeitsverteilung deutscher Trigramme

<i>ein</i>	1.22	<i>ese</i>	0.27	<i>hre</i>	0.18	<i>nne</i>	0.14	<i>auc</i>	0.11
<i>ich</i>	1.11	<i>auf</i>	0.26	<i>hei</i>	0.18	<i>nes</i>	0.14	<i>als</i>	0.11
<i>nde</i>	0.89	<i>ben</i>	0.26	<i>lei</i>	0.18	<i>ond</i>	0.14	<i>alt</i>	0.11
<i>die</i>	0.87	<i>ber</i>	0.26	<i>nei</i>	0.18	<i>oen</i>	0.14	<i>eic</i>	0.11
<i>und</i>	0.87	<i>eit</i>	0.26	<i>nau</i>	0.18	<i>sdi</i>	0.14	<i>esc</i>	0.11
<i>der</i>	0.86	<i>ent</i>	0.26	<i>sge</i>	0.18	<i>sun</i>	0.14	<i>enh</i>	0.11
<i>che</i>	0.75	<i>est</i>	0.26	<i>tte</i>	0.18	<i>von</i>	0.14	<i>eil</i>	0.11
<i>end</i>	0.75	<i>sei</i>	0.26	<i>wei</i>	0.18	<i>bei</i>	0.13	<i>fen</i>	0.11
<i>gen</i>	0.71	<i>and</i>	0.25	<i>abe</i>	0.17	<i>chl</i>	0.13	<i>gan</i>	0.11
<i>sch</i>	0.66	<i>ess</i>	0.25	<i>chd</i>	0.17	<i>chn</i>	0.13	<i>hte</i>	0.11
<i>cht</i>	0.61	<i>ann</i>	0.24	<i>des</i>	0.17	<i>chw</i>	0.13	<i>iea</i>	0.11
<i>den</i>	0.57	<i>esi</i>	0.24	<i>nte</i>	0.17	<i>ech</i>	0.13	<i>ieb</i>	0.11
<i>ine</i>	0.53	<i>ges</i>	0.24	<i>rge</i>	0.17	<i>edi</i>	0.13	<i>nli</i>	0.11
<i>nge</i>	0.52	<i>nsc</i>	0.24	<i>tes</i>	0.17	<i>enk</i>	0.13	<i>rda</i>	0.11
<i>nun</i>	0.48	<i>nwi</i>	0.24	<i>uns</i>	0.17	<i>eun</i>	0.13	<i>rsc</i>	0.11
<i>ung</i>	0.48	<i>tei</i>	0.24	<i>vor</i>	0.17	<i>enz</i>	0.13	<i>std</i>	0.11
<i>das</i>	0.47	<i>eni</i>	0.23	<i>dem</i>	0.16	<i>hau</i>	0.13	<i>sst</i>	0.11
<i>hen</i>	0.47	<i>ige</i>	0.23	<i>hin</i>	0.16	<i>ite</i>	0.13	<i>tre</i>	0.11
<i>ind</i>	0.46	<i>aen</i>	0.22	<i>her</i>	0.16	<i>ief</i>	0.13	<i>uss</i>	0.11
<i>enw</i>	0.45	<i>era</i>	0.22	<i>lle</i>	0.16	<i>imm</i>	0.13	<i>all</i>	0.10
<i>ens</i>	0.44	<i>ern</i>	0.22	<i>nan</i>	0.16	<i>ihr</i>	0.13	<i>aft</i>	0.10
<i>ies</i>	0.44	<i>rde</i>	0.22	<i>tda</i>	0.16	<i>iss</i>	0.13	<i>bes</i>	0.10
<i>ste</i>	0.44	<i>ren</i>	0.22	<i>tel</i>	0.16	<i>kei</i>	0.13	<i>dei</i>	0.10
<i>ten</i>	0.44	<i>tun</i>	0.22	<i>ueb</i>	0.16	<i>mei</i>	0.13	<i>erf</i>	0.10
<i>ere</i>	0.43	<i>ing</i>	0.21	<i>ang</i>	0.15	<i>nsi</i>	0.13	<i>ess</i>	0.10
<i>lic</i>	0.42	<i>sta</i>	0.21	<i>cha</i>	0.15	<i>nem</i>	0.13	<i>esw</i>	0.10

Künstliche Sprache (nach Kumpfmüller)

Einergruppen (Buchstabhäufigkeit)

EME GKNEET ERS TITBL BTZENFNDGBGD EAI E LASZ
BETEATR IASMIRCH EGEOM

Zweiergruppen (Paarhäufigkeit)

AUSZ KEINU WONDINGLİN DUFNRN ISAR STEISBERER ITEHM
ANORER

Dreiergruppen

PLANZEUNDGES PHIN INE UNDEN ÜBBEICHT GES AUF ES SO
UNG GAN DICH WANDERSO

Vierergruppen

ICH FOLGEMÄSZIG BIS STEHEN DISPONIN SEELE NAMEN

- Welche Arten von Redundanz können Lauflängencodierung und Huffmancodierung ausnutzen bzw. nicht ausnutzen?
- Welches Verfahren sollte wann benutzt werden?
- Welches Verfahren ist gut für Text?
- Sind Kombinationen, d.h. die Hintereinanderausführung von Kompressions-Verfahren sinnvoll?

- Verschiedene Kompressionsverfahren arbeiten unterschiedlich gut bei Daten unterschiedlichen Typs (Bilder, numerische Daten, Text).
- **Schlussfolgerung: Wir sollten uns für weitere Verfahren interessieren!**

Die LZ-Familie von Kompressionsalgorithmen

LZ77, LZ78, LZW, ...

Nutze unterschiedliche Häufigkeit von *Zeichenfolgen* aus.

Grundlegende Idee: verwende ein “Wörterbuch” von Zeichenfolgen, so dass häufige Zeichenfolgen durch Verweise auf das Wörterbuch platzsparend codiert werden können. (LZ78,LZW: explizites Wörterbuch; LZ77: implizit, “sliding window”)

Text-Beispiel: RABARBABARBARABARBARBARENBART¹

Grundsätzlicher Konflikt: Je grösser das Wörterbuch, desto

- + mehr und längere Zeichenfolgen können durch Verweise codiert werden
- mehr Platzbedarf pro Verweis
- mehr Platzbedarf für Wörterbuch
- höhere Laufzeit und Platzbedarf für Kompression und Dekompression

¹in kompressionsfreundlich-reformierter Rechtschreibung

Historischer Exkurs (LZ-Familie)

- | | |
|------|---|
| 1977 | Abraham Lempel und Jacob Ziv erfinden den Kompressionsalgorithmus LZ77 |
| 1983 | Terry A. Welch (Sperry Corporation - später Unisys) patentiert LZW (Variante von LZ78) |
| 1987 | CompuServe veröffentlicht GIF als freie und offene Spezifikation (Version 87a) |
| 1989 | Vorstellung von GIF 89a |
| 1993 | Unisys informiert CompuServe über die Verwendung ihres patentierten LZW-Algorithmus in GIF |
| 1994 | Unisys gibt öffentlich bekannt, Gebühren für die Verwendung des LZW-Algorithmus einzufordern |
| 1995 | die PNG Gruppe wird gegründet, erste PNG Bilder werden ins Netz gestellt, die PNG-Spezifikation 0.92 steht im W3C |
| 1997 | PNG-Unterstützung in Netscape 4.04 und IE 4.0 |

- Michael Burrows und David Wheeler, Mai 1994
- **Grundidee:** Daten vorbehandeln, um sie besonders gut komprimieren zu können.
- Kompression in 3 Schritten
 - Burrows-Wheeler-Transformation (BWT):
produziere lange Runs von Zeichen
 - Move-To-Front (MTF):
produziere daraus lange runs von Nullen
 - dann erst komprimieren (Huffman-Codierung)
- Die Kompression ist sehr gut geeignet für Text, da Kontext berücksichtigt wird.
- Verwendet in `bzip`

Beispiel für Burrows-Wheeler-Transformation

Vorwärtstransformation von HelloCello

	0	1	2	3	4	5	6	7	8	9
0	H	e	l	l	o	C	e	l	l	o
1	e	l	l	o	C	e	l	l	o	H
2	l	l	o	C	e	l	l	o	H	e
3	l	o	C	e	l	l	o	H	e	l
4	o	C	e	l	l	o	H	e	l	l
5	C	e	l	l	o	H	e	l	l	o
6	e	l	l	o	H	e	l	l	o	C
7	l	l	o	H	e	l	l	o	C	e
8	l	o	H	e	l	l	o	C	e	l
9	o	H	e	l	l	o	C	e	l	l

HelloCello

	0	1	2	3	4	5	6	7	8	9
0	C	e	l	l	o	H	e	l	l	o
1	H	e	l	l	o	C	e	l	l	o
2	e	l	l	o	C	e	l	l	o	H
3	e	l	l	o	H	e	l	l	o	C
4	l	l	o	C	e	l	l	o	H	e
5	l	l	o	H	e	l	l	o	C	e
6	l	o	C	e	l	l	o	H	e	l
7	l	o	H	e	l	l	o	C	e	l
8	o	C	e	l	l	o	H	e	l	l
9	o	H	e	l	l	o	C	e	l	l

/

L

ooHCeellll

- die Zeilen links über “HelloCello” sind zyklische Rotationen des Texts
- die Zeilen rechts sind lexikographisch sortiert
- Betrachten wir die 1. und letzte Spalte:
 - letzte Spalte: 1-Char Prefix der jeweiligen Zeile
 - ab 1. Spalte: Suffixe des jeweils letzten Buchstabens der Zeile
- für Suffixe mit häufig auftretendem Präfix tauchen die Präfix-Buchstaben gruppiert in der letzten Zeile auf
- (nicht immer perfekt gruppiert, aber in “guten” Runs)
- BWT hat eine inverse Funktion . . .

- Erstelle eine Zuordnung von Zeichen in codierter Position (letzte Spalte der Codierungstabelle) zur sortierten Position (erste Spalte der Codierungstabelle)
- Position in letzter Spalte entspricht Zeilennummer
- Position in erster Spalte übernimmt die Position aus der letzten Spalte
- Bestimme Startposition
 - Suche Zeile, welche uncodierten Text enthält
 - Startposition ist Zeilennummer in letzter Spalte
- Beginne bei Startposition
- While noch nicht alle Positionen besucht:
 - Lies Ergebnis aus sortierter Liste (*Zeile3*)
 - Gehe zu Position in Spalte codiert (*Zeile4*)

Zuordnung: Beispiel

codiert (letzte Spalte)	o	o	H	C	e	e	l	l	l	l
Position (letzte Spalte)	0	1	2	3	4	5	6	7	8	9
Sortiert (erste Spalte)	C	H	e	e	l	l	l	l	o	o
Position in Spalte codiert	3	2	4	5	6	7	8	9	0	1
Ausgabe		H	e		l		l		o	

Die Burrows-Wheeler Transformation

- Ein Eingabeblock der Länge N wird als quadratische Matrix dargestellt, die alle Rotationen des Eingabeblocks enthält.
- Die Zeilen der Matrix werden alphabetisch sortiert.
- Die letzte Spalte und die Zeilennummer des Originalblocks werden ausgegeben.

In der Ausgabe werden Zeichen mit ähnlichem Kontext zusammensortiert
⇒ lange Runs gleicher Buchstaben.

Daraus lässt sich der Originalblock wieder rekonstruieren
(wird hier nicht bewiesen, wurde nur demonstriert).

- *Berechnung der Korrektheit* via equational Reasoning: “Functional Pearls: Inverting the Burrows-Wheeler Transform” (Bird, Mu, 2004)
- wichtig ist folgende Eigenschaft eines:
 - Sortieralgorithmus auf
 - Matrix von Rotationen eines Strings
- Die Matrix nach Zeilen sortiert ist äquivalent zu
- n -mal die komplette Matrix 1-Char rechts rotieren und *nur nach dem 1. Char* stabil zu sortieren
- mit diesem Argument kann die ursprüngliche sortierte Matrix (aus der Vorwärts-Variante) aus der letzten Spalte allein wieder hergestellt werden

- Gruppierung typischerweise nicht perfekt (im Vergleich zum Sortieren)
- aber reversibel (sonst wäre BWT fuer Kompressionsverfahren nicht nützlich)
- Eingaben sollten relativ lang sein, mehrere Kilobyte, um genügend lange Runs zu erzeugen
- BWT erfordert Prefix-Suffix Muster

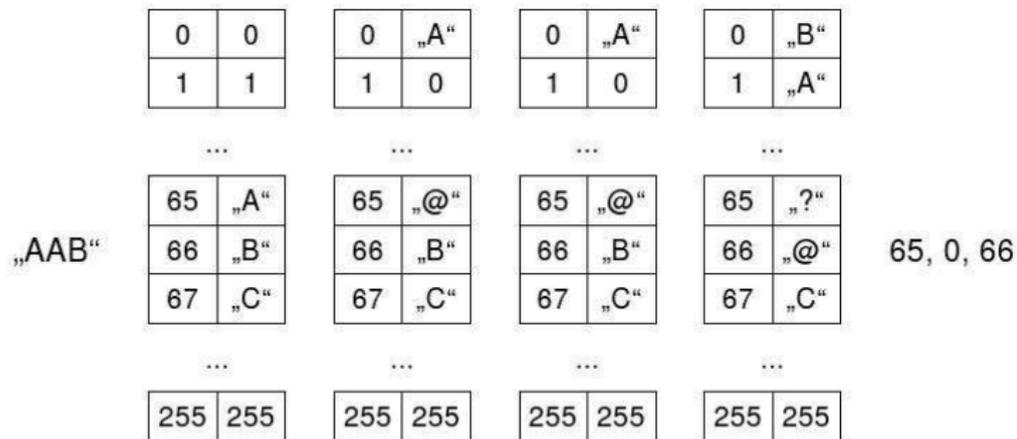
Idee: Codiere Runs von Buchstaben mit möglichst kleinen Zahlen

- Beginne mit Array `...A ...Z ...`, wobei 'A' = 65, 'B' = 66, ...
- das 1. Auftreten eines Buchstaben im Run codiert den (momentanen) Index des Buchstabens
- dieser Buchstabe wandert an den Anfang des Array
- weiteres Auftreten des Buchstaben wird durch '0' codiert

also:

- 1 Lese Buchstabe
- 2 finde Index zu Buchstabe und schreibe Index
- 3 Verschiebe Buchstabe an 0'te Position
- 4 Weiter mit Schritt 1, falls Eingabe nicht zu Ende

MTF: Move-To-Front-Coding



Eingabe: A A B B C C C A A B C C C
 Ausgabe: 65 0 66 0 67 0 0 2 0 2 2 0 0

- 1 Beginne mit Array `...A...Z...`, wobei `'A' = 65`, `'B' = 66`, ...
- 2 die Eingabe ist eine Liste der Indices
- 3 Schreibe Buchstaben welcher durch Index (zB. `65 → 'A'`) codiert
- 4 Verschiebe Buchstaben an den Anfang des Arrays
- 5 weiter mit Schritt 2

Wozu das Ganze

- BWT erzeugt durch Sortierung lange Läufe gleicher Zeichen
- MTF erzeugt daraus kleine (Index-)Zahlen für häufig wiederholte Zeichen, genauer: Lange Folgen von '0', erlaubt gute Lauflängen-Kompression
- Huffman-Kodierer ermöglicht kurze Bitfolgen für häufige Zeichen

bzip-Funktionen:

$$\mathbf{BWT\text{-}Kompression}(x) = \text{Huffmann}(\text{MTF}(\text{BWT}(x)))$$

$$\mathbf{BWT\text{-}Dekompression}(y) = \text{BTW}^{-1}(\text{MTF}^{-1}(\text{Huffmann}^{-1}(y)))$$

[= BWT-Kompression⁻¹(y)]

Definition (Shannon)

Entscheidungsinformation: Anzahl optimal gewählter binärer Entscheidungen zur Ermittlung eines Zeichens in einem Zeichenvorrat

Entropie: der mittlere Informationsgehalt eines Textes

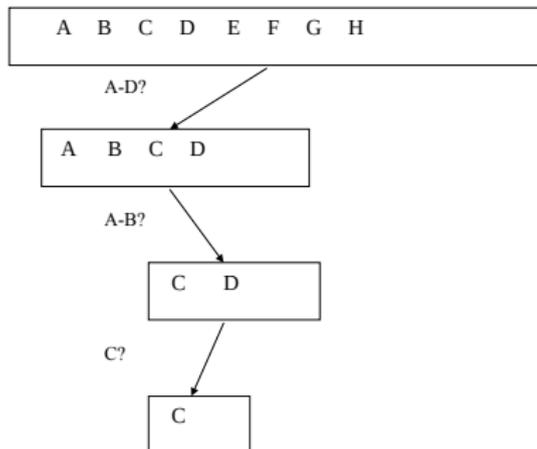
Redundanz: Anteil einer Nachricht, der keine Information enthält

- gegeben sei ein Alphabet $\mathcal{A} = \{A \dots Z\}$ (oder ein beliebiger anderer Zeichenvorrat)
- für jeden Buchstaben $x \in \mathcal{A}$ sei die Wahrscheinlichkeit, dass dieser Buchstabe auftritt: $p(x)$ mit $\sum_x p(x) = 1$.
- der Informationsgehalt sei $I(x) = -\log_2 p(x)$
- hierbei meint *Information* im engeren Sinne *Entscheidungsinformation*, also die Anzahl optimal gewählter binärer Entscheidungen zur Ermittlung eines Zeichens innerhalb eines Zeichenvorrats

Entscheidungsinformation - Beispiel

Gegeben seien 8 Zeichen. Nach maximal wieviel Schritten ist ein Zeichen gefunden?

Entscheidungsbaum:



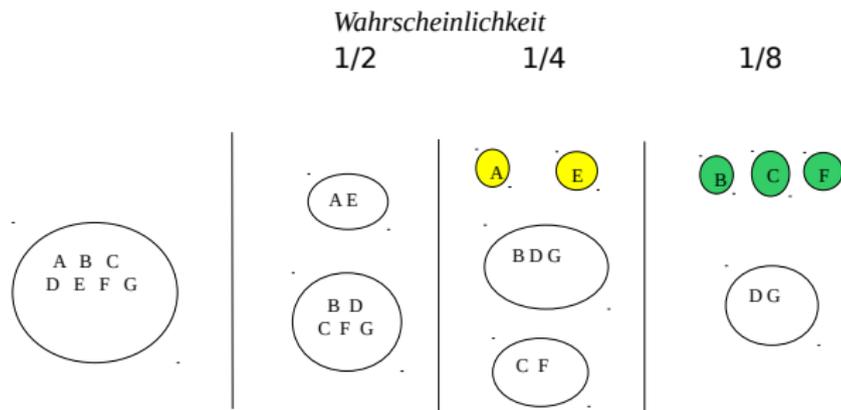
Entscheidungsinformation - Beispiel

Wichtig für den allgemeinen Fall ist die Aufteilung eines Alphabets nicht in gleich große, sondern **gleich wahrscheinliche** Mengen von Zeichen.

Das i -te Zeichen ist nach k_i Alternativentscheidungen isoliert,

seine Wahrscheinlichkeit ist $p_i = (1/2)^{k_i}$,

sein Informationsgehalt $k_i = \text{ld}(1/p_i)$ bit.



Buchstaben	p_i	Codierung
A	1/4	00
E	1/4	01
F	1/8	100
C	1/8	101
B	1/8	110
D	1/16	1110
G	1/16	1111

Mittlerer Entscheidungsgehalt pro Zeichen (Entropie):

$$\begin{aligned} H &= p_1 l_1 + p_2 l_2 + \dots + p_n l_n \\ &= \sum p_i \log_2(1/p_i) \text{ bit} \\ &= 2/4 + 2/4 + 3/8 \dots = 2,625 \end{aligned}$$

- die *Entropie* eines Zeichens im Text ist:
$$H_1 = \sum_x p(x) I(x) = - \sum_x p(x) \log_2 p(x)$$
- für Wörter der Länge n ergibt sich:
$$H_n = - \sum_{w \in \mathcal{A}^n} p(w) \log_2 p(w)$$
- mit $H = \lim_{n \rightarrow \infty} H_n/n$
- falls die Zeichen stoch. unabhängig sind: $H_n = nH_1$, also $H = H_1$

- besitzt in einer Codierung einer Nachricht das i -te Zeichen die Wortlänge N_i , so ist $L = \sum p_i N_i$ die **mittlere Wortlänge**.
- sind alle Elemente des Zeichenvorrats genau gleichwahrscheinlich, gilt $L = H$
- im allgemeinen gilt das *Shannonsche Codierungstheorem*
 - $H \leq L$
 - jede Nachricht kann so codiert werden, dass die Differenz $L-H$ beliebig klein wird
- Die Differenz $L-H$ heißt *Code-Redundanz*, die Größe $1-H/L$ *relative Code-Redundanz*

Informationsgehalt Schriftsprache Deutsch

- 30 Buchstaben (inkl. Zwischenraum), also $I = \log_2 30 = 4,9$ bit
- mittlerer Informationsgehalt (*Entropie*) unter Berücksichtigung von Bigrammen $H = 1,6$ bit
- **Redundanz:** $4,9 - 1,6 \text{ bit} = 3,3 \text{ bit}$
- Text ist auch dann noch lesbar, wenn jeder zweite Buchstabe fehlt

Bei reduzierter Redundanz wird das Lesen sehr viel mühsamer

BEI REDUZIERTER REDUNDANZ WIRD DAS LESEN SEHR VIEL
MÜHSAMER

BEI REDUZIERTER REDUNDANZ WIRD DAS LESEN SEHR VIEL MÜHSA-
MER

BE RE UZ ER ER ED ND NZ IR DA LE EN EH VI LM HS ME

(nach Breuer)

Was bedeutet das alles?

- “zufällige” Texte haben die höchste Entropie und lassen sich nicht komprimieren
- falls häufig gleiche Zeichen hintereinander stehen, so sind diese Zeichen im Text stochastisch voneinander abhängig
- damit reduziert sich ihre Entropie und auch die des gesamten Textes
- Kompressionsverfahren nutzen dies, um möglichst nahe an das theoretische Optimum für Kompression zu kommen
- warum nicht “optimal”: die Kompressionsverfahren müssen die Abhängigkeiten sehen können