

Kompresia dát

Zdroje redundancie:

- priveľa znakov
- frekvencie výskytu znakov sú rôzne
- podmienené pravdepodobnosti rôznych postupností znakov sú rôzne
- dá sa vyjadriť vzorcom alebo algoritmom (malá kolmogorovská zložitosť)
- možnosti interpolácie a aproximácie

$$\text{Kompresný pomer } k = \frac{\textit{dĺžka kompresovaného súboru}}{\textit{dĺžka pôvodného súboru}}$$

Niekedy sa ako kompresný pomer uvádza hodnota $1 - k$.

Vlastnosti kódov

- Jednoznačná dekódovateľnosť (Každú postupnosť vieme jednoznačne rozdeliť na kódové slová.)
- Prefixová vlastnosť (Žiadne kódové slovo nie je prefixom iného kódového slova.)
- Vlastnosť okamžitej rozhodnuteľnosti (Pri prečítaní posledného znaku poznáme koniec kódového slova.)

Príklad:

	kód 1	kód 2
A	0	0
B	10	01
C	110	011
D	1110	0111

Kód 1 má prefixovú vlastnosť.

Kód 2 je jednoznačne dekódovateľný, ale nemá prefixovú vlastnosť, ani vlastnosť okamžitej rozhodnuteľnosti.

Prefixová vlastnosť implikuje obe ostatné vlastnosti.

Kompresia rozptýlených záznamov

- Opakovací znak a počet opakovaní
- Bitová mapa
- Separácia konštánt

Príklad: $d_1 d_2 000000 d_3 d_4 d_5 d_6 000 d_7 d_8 d_9 000$

Opakovací znak: $d_1 d_2 \Re 6 d_3 d_4 d_5 d_6 \Re 3 d_7 d_8 d_9 \Re 3$

Bitová mapa: 110000001111000111000 : $d_1 d_2 d_3 d_4 d_5 d_6 d_7 d_8 d_9$

Separácia: $\underbrace{26699}_{(12)}: d_1 d_2 d_3 d_4 d_5 d_6 d_7 d_8 d_9$

Distribučný
vektor

Na nepárnych miestach počet dátových položiek. Na párnych miestach počet konštánt. Súčet znamená celkový počet predchádzajúcich.

Nájdienie L -tej položky z pôvodných dát v kompresovanom súbore: Nájdeme najmenšie i také, že $L \leq v_{i-1} + v_i$. Ak i je párne konštanta. Ak i je nepárne $p = L - v_{i-1}$.

Diferenčné metódy kompresie

Princíp pamätá sa len rozdiel od predchádzajúceho záznamu. Metóda je vhodná na menné zoznamy, bibliografické údaje, slovníky, tabuľky hodnôt funkcií. Možno ju kombinovať so slovníkovou metódou. Nevýhodou je citlivosť na chyby.

Pr.:	Michalčík	0	Michalčík	(0,3)	Michalč
	Michna	4	na	(4,0)	na
	Mikelka	2	kelka	(2,4)	kel
	Mikletič	3	letič	(3,2)	let
	Mikulcová	3	ulcová	(3,7)	ul
	Mikulič	4	lič	(4,2)	l
	Mikuličová	7	ová	(7,7)	
	Mikulová	5	ová	(5,7)	

Pridaním slovníka 15 najbežnejších koncoviek mien: er, ič, ík, ka, ná, ný, ová, ... zlepší sa kompresný pomer.

Frekvenčná analýza – Huffmanové kódy

Princíp minimalizácia očakávanej dĺžky kódu. Nech p_i je pravdepodobnosť výskytu i -tého kódového znaku a l_i je dĺžka. Potom očakávaná dĺžka kódu H je daná vzorcom:

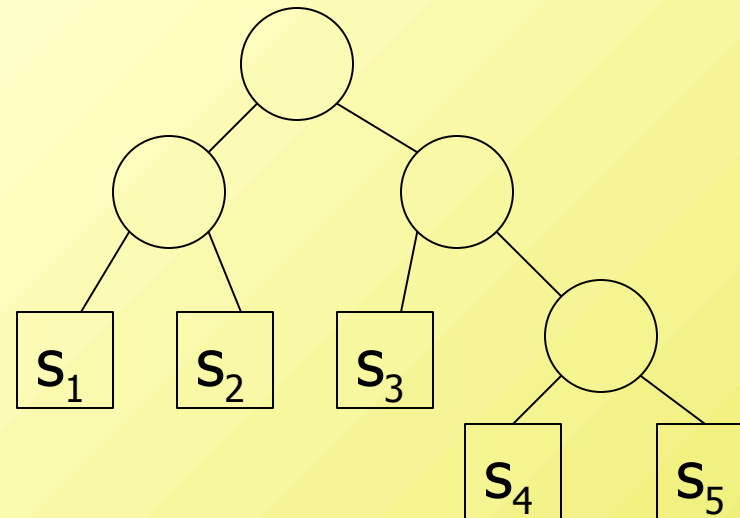
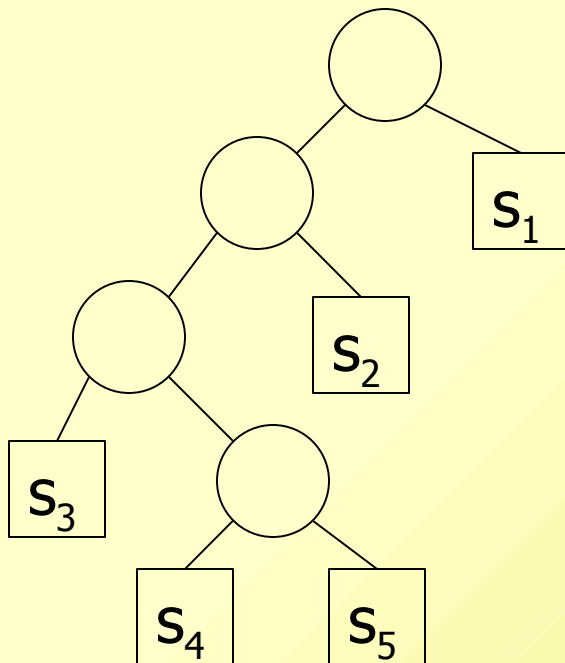
$$H = \sum_{i=1}^n p_i l_i, \quad \text{kde } n \text{ je počet znakov.}$$

Algoritmus konštrukcie Huffmanovho kódu (Knuth):
Utriedíme znaky podľa p_i . Zlúčime znaky dva znaky s najmenšími hodnotami p_i . Riešime problém pre $n-1$ hodnôt.

Algoritmus nie je jednoznačný. V niektorých prípadoch sa môžeme rozhodnúť, ktoré uzly budeme spájať. Možno je výhodne v takomto prípade minimalizovať súčet vnútorných ciest. Dostaneme tak kódy, ktoré sa dĺžkou najmenej odlišujú.

Konštrukcia Huffmanovho kódu

znak, p_i ;	znak, p_i ;	znak, p_i ;	znak, p_i ;	znak, p_i ;	Kód
s_1 0.4	0.4	0.4	s_{3452} 0.6	s_{34521} 1	1
s_2 0.2	0.2	s_{345} 0.4	s_1 0.4		01
s_3 0.2	0.2	s_2 0.2			000
s_4 0.1	s_{45} 0.2				0010
s_5 0.1					0011

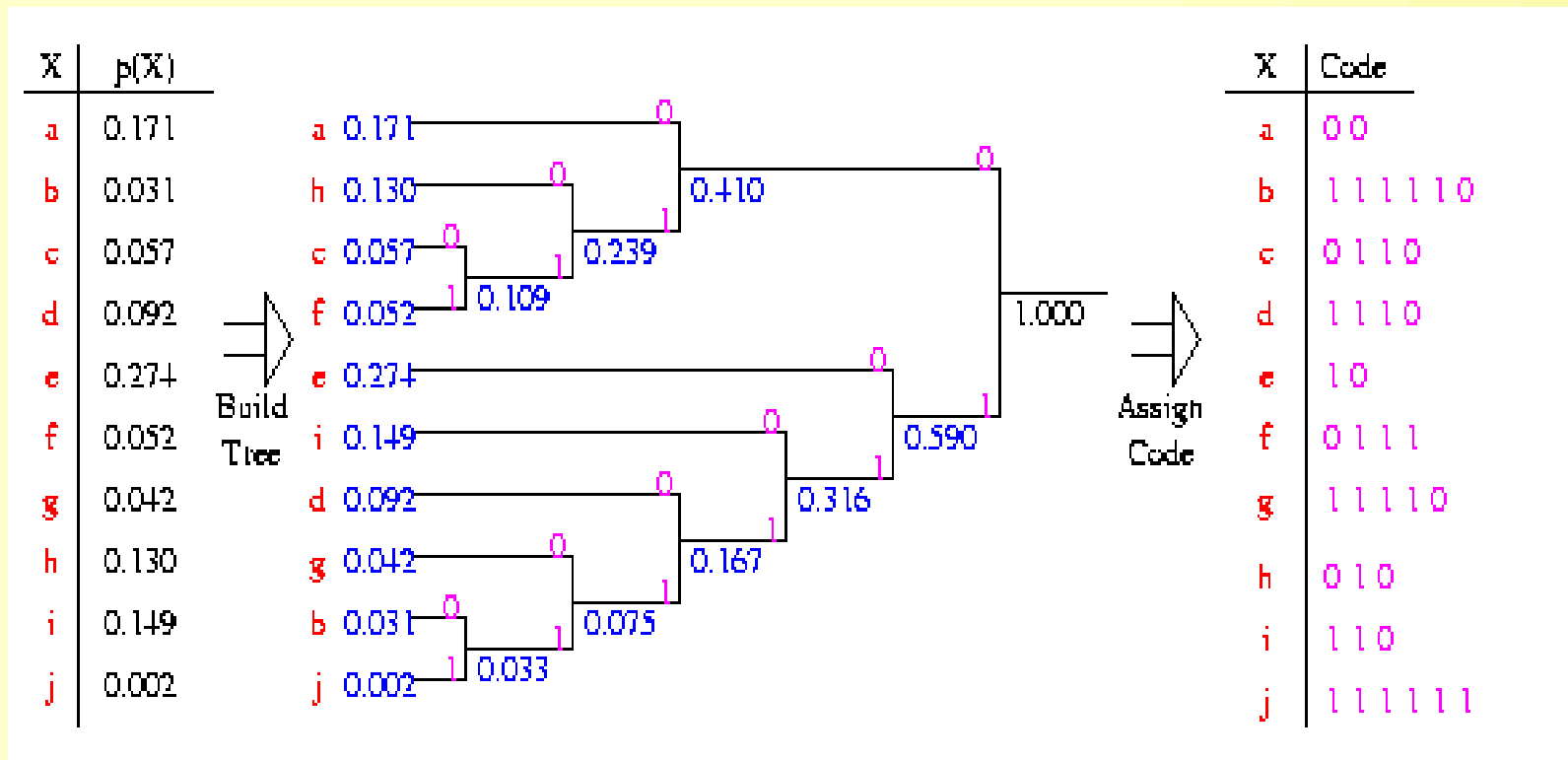


00
01
10
110
111

Iný strom a kód

Kompresia dát

Väčší príklad



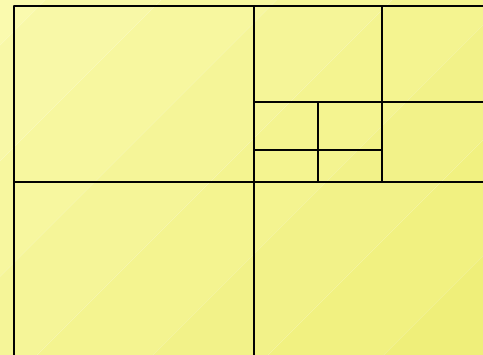
Run length coding

Niekedy sa stáva, že pravdepodobnosti výskytu znakov sú približne rovnaké, ale je veľká pravdepodobnosť, že nasledujúci údaj bude ako predošlý. Vznikajú dlhé sekvencie (runs) rovnakých znakov. (Je to typické pre grafické dáta.) Kódovanie pomocou znaku opakovania a počtu opakovaní.

Príklad: aaaabbbbccccabccccbbbaaaa $l=25$
 a \Re 4b \Re 4c \Re 3abc \Re 5b \Re 3a \Re 4 $l=20$

Účinnosť tejto metódy je tým väčšia, čím sú sekvencie dlhšie.

Quad trees a oct trees
v počítačovej grafike
sú založené na
podobnom princípe.



Kompresia s využitím kódového slovníka

Daný text t máme kompresovať pomocou slovníka n fráz $d = \{p_i: 1 \leq i \leq n\}$. Predpokladáme, že slovník obsahuje všetky elementárne (jednoznakové) frázy. Nech $|t| = N$ označuje dĺžku vstupného textu a t_j j -tý symbol vstupného textu.

Označme $B(j) = \{p_i: \forall (0 \leq k < |p_i|) t_{j+k} = p_{ik}\}$, kde p_{ik} je k -tý symbol frázy p_i . Nech ω_i je cena frázy p_i . (Obvykle $\omega_i = 1$ pre všetky i .)

Cenu $F(1)$ minimálneho pokrytia textu t dostaneme, ako riešenie systému rekuretných rovníc:

$$F(N + 1) = 0$$

$$F(j) = \min_{p_i \in B(j)} \{F(j + |p_i|) + \omega_i\}$$

$F(j)$ je cena pokrytia textu od pozície j (včítane) do konca.

Na pokrytie použijeme tie frázy, pomocou ktorých sa $F(1)$ vypočítalo.

Príklad - kompresia pomocou slovníka

Slovník: { A, AR, ARG, E, EN, G, GU, M, ME, MENT, N, R, U, UM, T }
 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Text: ARGUMENT

$B(1)=\{1,2,3\}$ $B(4)=\{13,14\}$ $B(7)=\{11\}$

$B(2)=\{12\}$ $B(5)=\{8,9,10\}$ $B(8)=\{15\}$

$B(3)=\{6,7\}$ $B(6)=\{4,5\}$

$F(9)=0$, $F(8)=1$, $F(7)=2$, $F(6)=2$, $F(5)=1$, $F(4)=2$,
 15 15,8 15,14 10 10, 13

$F(3)=2$, $F(2)=3$, $F(1)=3$

10,7 12,10,7 10,7,2 (10,13,3)

AR-GU-MENT 2,7,10; ARG-U-MENT 3,13,10.

Napriek zdanlivej zložitosti, kompresovať podľa statického slovníka sa dá v „lineárnom“ čase. Konštrukcia optimálneho slovníka, k danému textu je NP ťažký problém.

Adaptívne metódy kompresie

Hoci optimálne pokrývanie sa dá robiť efektívne, vyžaduje najprv prečítanie celého vstupu. Konštrukcia slovníka je tak, či tak ťažký problém. Nápad vzdať sa optimálnosti, použiť nejaký „pažravý“ (greedy) algoritmus a robiť všetko za letu (on the fly).

```
Algoritmus: while not koniec do  
            begin zakóduj čo najdlhší úsek vstupu;  
                modifikuj slovník  
            end;
```

Výhodou týchto metód je, že pri dekódovaní sa použije ten istý algoritmus modifikácie slovníka a slovník nie je potrebné prenášať a uchovávať s dátami.

Ziv – Lempel – Welch

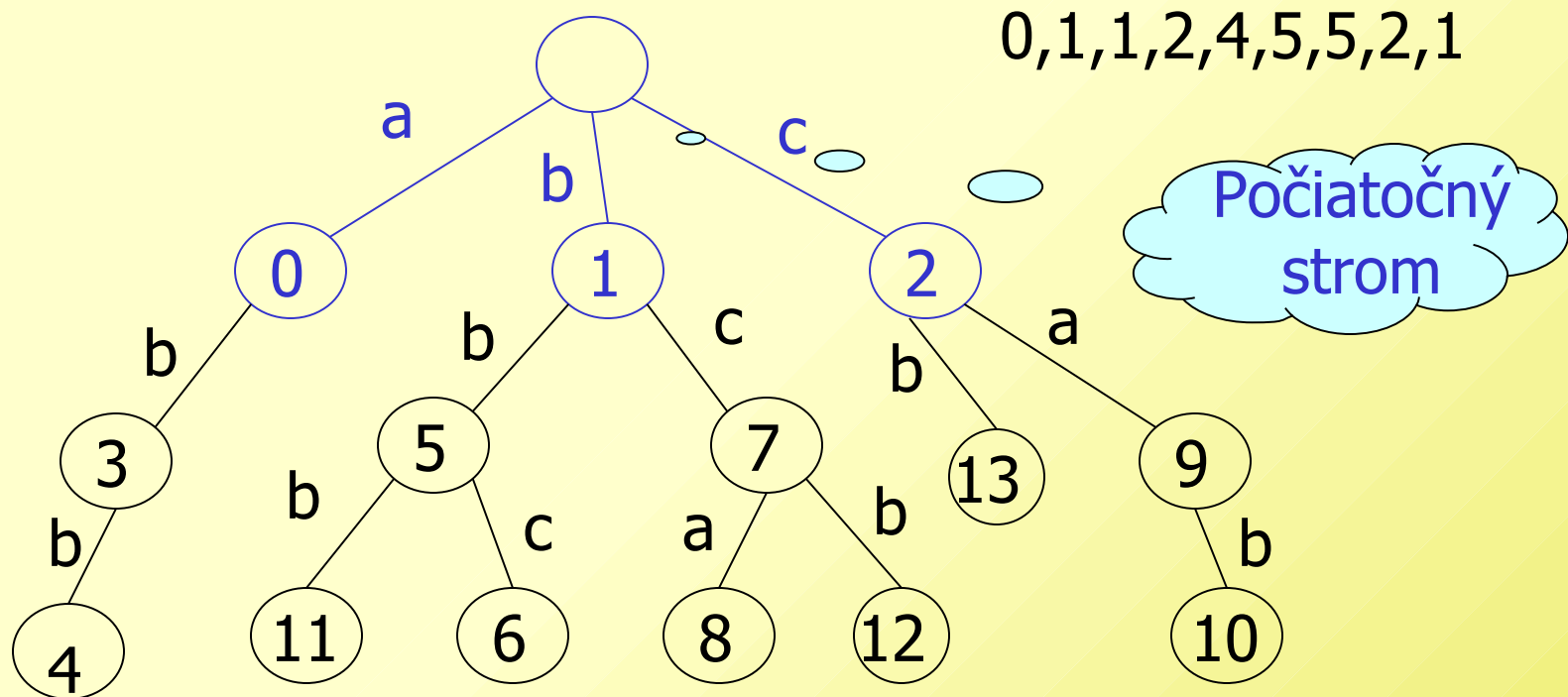
Slovník je strom zostavený z podreťazcov prečítaného úseku textu. Maximálna výška tohoto stromu h a najväčší možný počet vrcholov n sú parametre metódy.

Vrcholy sú označené číslami kódmi a hrany symbolmi vstupnej abecedy. Na počiatku sú všetky symboly listy v hĺbke 1. Tieto postupne rozširujeme o suffixy prečítaného textu (dĺžky menšej alebo rovnej h) pokiaľ strom nemá n vrcholov. Uzol stromu je kód pre cestu od vrcholu k nemu. Ak sa strom naplní nemusíme ho už ďalej modifikovať, ale môžeme ho reorganizovať a prečíslovať vrcholy.

Kódovanie sa robí tak, že nájdeme najdlhšiu cestu od koreňa k uzlu stromu, ktorej hrany sa zhodujú so vstupom.

Príklad konštrukcie adaptívneho slovníku

$\Sigma = \{a, b, c\}$, $h = 4$, $n = 16$, $t = abbcabbbbbbcb$



Najjednoduchší spôsob LZ kódovania

Inicializuj slovník; /* znaky abecedy */

Nastav sa na začiatok vstupného slova W ;

repeat

Z danej pozície najdi najdlhší úsek u slova W , ktorý je v slovníku;

Zakóduj u indexom v slovníku;

Nastav sa nasledujúcu pozíciu vstupného slova $W[\text{next}]$;

Ak v slovníku je ešte miesto pridaj $uW[\text{next}]$ do slovníka;

until koniec vstupu;

Posledný príkaz môžeme rôzne modifikovať. Napr. obmedziť dĺžku pridávaného slova, alebo pridávať všetky sufixy $uW[\text{next}]$. Aj činnosť v prípade, že slovník je plný môže byť rôzna. Môžeme pokračovať bez pridávania, alebo vyprázdniť všetko až po prvú úroveň a začať znova, alebo vynechať len nejakú časť.

Aritmetické kódovanie

Huffmanovo kódovanie nie je optimálne pretože dĺžka kódu musí byť celé číslo. Idea kódovať pomocou intervalov podintervalov $<0, 1)$. Celý text jeden interval.

Symbol	pravdepodobnosť	interval
a	0.2	$<0, 0.2)$
e	0.3	$<0.2, 0.5)$
i	0.1	$<0.5, 0.6)$
o	0.2	$<0.6, 0.8)$
u	0.1	$<0.8, 0.9)$
!	0.1	$<0.9, 1)$

Text je eaii!.

Pri kompresii i dekompresii začíname s intervalom $<0, 1)$.
Vždy po načítaní symbolu zúžime interval.

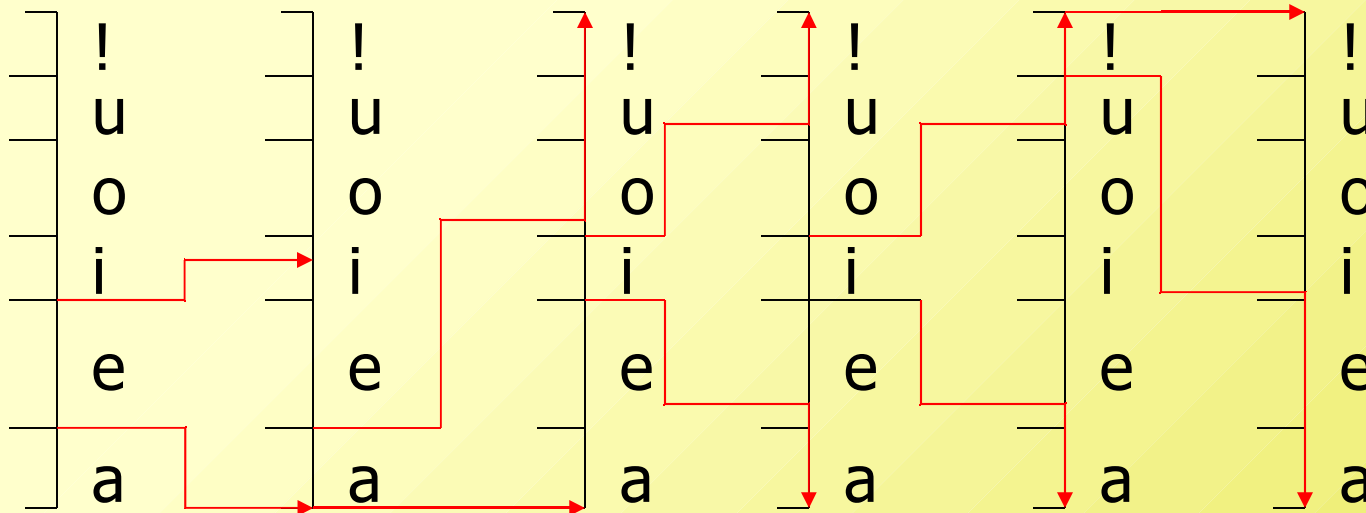
Aritmetické kódovanie - pokračovanie

Proces kódovania (kompresie)

! <0, 1)
e <0.2, 0.5)
a <0.2, 0.26)
i <0.23, 0.236)
i <0.233, 0.2336)
! <0.23354, 0.2336)

Ak posledný znak je znak konca textu, môžeme kompresovaný text zakódovať ľubovoľným číslom z výsledného intervalu.

Proces dekódovania (dekompresie)



Technické detaily aritmetického kódovania

- Presná aritmetika (integer, fixed)
- kumulatívne frekvencie namiesto pravdepodobností
- underflows - scaling

Metóda bude adaptívna, keď po každom prečítaní symbolu, aktualizujeme kumulatívne frekvencie, či pravdepodobnosti.

Metóda je vhodná na:

- adaptívnu kompresiu textu
- kompresiu postupností celých čísiel
- kompresiu čierno-bielých obrázkov

Transformácie „zmenšujúce entropiu“

Entropia H je očakávaný počet bitov potrebný na zakódovanie znaku a_i s pravdepodobnosťou výskytu $p(a_i)$.

$$H = - \sum_{i=1}^n p(a_i) \log_2 p(a_i)$$

Dobré kódy sú blízke tejto hranici:

- Aritmetické kódovanie
- Huffman

Pred kompresiou zakódujeme text nejakou transformáciou, ktorá zmenší jeho entropiu a tým umožní jeho lepšiu kompresiu.

Presun na začiatok – MTF

Najprv znaky zakódujeme do slovníka v nejakom poradí. Znaky kódujeme poradovým číslom. Použitý znak presunieme na začiatok.

Príklad: a b a b a a b c c b b c c c c b d b c c

0	1	...	2
0 1 2 3	0 1 2 3	...	0 1 2 3
a b c d	b a c d		c b a d

Výsledok: 0 1 1 1 1 0 1 2 0 1 0 1 0 0 0 1 3 1 2 0

Frekvencie znakov: a b c d 0 1 2 3
 4 7 8 1 8 9 2 1

Entropia: 1.74 1.6

Extrém: a a a a a b b b b b c c c c c d d d d d

Výsledok: 0 0 0 0 0 1 0 0 0 0 2 0 0 0 0 3 0 0 0 0

Entropia: 2 0.85

Burrows Wheelerova transformácia

Algoritmus 1 (kódovanie):

Z reťazca r dĺžky n vytvoríme všetkých n jeho cyklických posuvov o jeden znak vľavo a tieto utriedíme. Výstup L sú posledné znaky každého reťazca a index (pozícia) pôvodného reťazca v utriedenej postupnosti.

Algoritmus 2 (inverzná transformácia – dekódovanie):

Znaky zakódovaného reťazca L utriedime. Dostaneme reťazec F . Vytvoríme dekódovací vektor $T[i] = j$, práve vtedy ak symbol $F[i]$ je na j tej pozícii $L[j]$. Pritom predpokladáme, že triedenie bolo stabilné. Priradíme i počiatočný index a dekodujeme algoritmom:

```
repeat output(F[i]); i := T[i] until EOF;
```

Príklad

Kódovanie:

0	abrakadabra
1	brakadabraa
2	rakadabraab
3	akadabraabr
4	kadabraabra
5	adabraabrak
6	dabraabraka
7	abraabrakad
8	braabrakada
9	raabrakadab
10	aabrakadabr

A

0	aabrakadab	r
1	abraabraka	d
2	abrakadabr	a
3	adabraabra	k
4	akadabraab	r
5	braabrakad	a
6	brakadabra	a
7	dabraabrak	a
8	kadabraabr	a
9	raabrakada	b
10	rakadabraa	b

Výstup po transformácii: **2 ; rdakraaaabb**

Dekódovanie

0 1 2 3 4 5 6 7 8 9 10
L = rdakraaaabb
F = aaaaabbdkrr
T = 2 5 6 7 8 9 10 13 0 4

Že dekódovanie funguje.
Posuňme riadky matice A
cyklicky o jedno doprava
dostaneme maticu A^s .
Bude začínať stĺpcami LF.
Prvý symbol $A[T[i]]$ je
druhý symbol $A[i]$, pretože
 $A^s[T[i]] = A[i]$.

$l = 2$
 $l = T[2] = 6$
 $l = T[6] = 10$
 $l = T[10] = 4$
 $l = T[4] = 8$
 $l = T[8] = 3$
 $l = T[3] = 7$
 $l = T[7] = 1$
 $l = T[1] = 5$
 $l = T[5] = 9$
 $l = T[9] = 0$

T[l]
a
b
r
a
k
a
d
a
b
r
a

Použitie

Burrows Wheelerova transformácia nekompresuje, ale triedením presúva rovnaké symboly do zhŕukov. Umožňuje teda lepšiu kompresiu.

BZip, BZip2:

- Rozdelí text na rozumne veľké bloky pevnej dĺžky
 - Čím väčší blok, tým lepšia kompresia
- Na blok aplikuje Burrows Wheelerovu transformáciu
- Potom aplikuje presun na začiatok
- Kompresuje aritmetickým kódovaním

Množstvo vylepšení a programátorských trikov

- Triedenie len na základe prefixu
- Netreba skutočne robiť cyklické posuvy
- Umiernený posun na začiatok