

Kompresi dat a neb

Zduajové kódování

notace: A - Typy kódů, střední délka
 koncína množina

A^* - množina koncinych
paslaupnosti prvků z A

tj:

$$A^* = A \cup (A \times A) \cup \dots \cup (A \times \dots \times A) \cup \dots$$

X - náhodná veličina s
koncíní množina hodnotami
 X

$$p(x) = P[X=x]$$

$D = \{0, 1, \dots, D-1\}$ D -národní abeceda

Definice: Zduajový kód C pro X

X zobrazení z X do D^*

Značení: $C(x)$ - kódové slovo pro $x \in X$
 $l(x)$ - délka kódu $C(x)$

Příklad $X = \{red, blue\}$

$D = \{0, 1\}$

$C(red) = 00$ $C(blue) = 11$

Definice: Střední délka, $L(C)$,
kódu C je

$$L(C) = \sum_{u \in X} p(u) l(u)$$

Pr: $D = 20,14$

$$P[X=1] = 1/2 \quad C(1) = 0$$

$$P[X=2] = 1/4 \quad C(2) = 10$$

$$P[X=3] = 1/8 \quad C(3) = 110$$

$$P[X=4] = 1/8 \quad C(4) = 111$$

$$L(C) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = 1,75$$

$$H(X) = 1,75$$

Každá posloupnost, řádových slov
(bez mezí) se dá dekodovat

0110 | 111 | 10 | 0110

1 3 4 2 1 3

\exists -li $|D| > |X|$

musíme jednoznačně přirodit každou slova všechny délky 1. Pak $\exists L(C)$ nejmenší možné

$$L(C) = \sum_{x \in X} p(x) \cdot 1 = 1$$

Většinou ale máme malou abecedu

Definice Extenze C^+ kódu C je
řazení $C^+ : X^+ \rightarrow D^+$ definované

$$C^+(c_1, c_2, \dots, c_m) = \underbrace{C(c_1) C(c_2) \dots C(c_m)}_{\text{spojení kódových slov}}$$

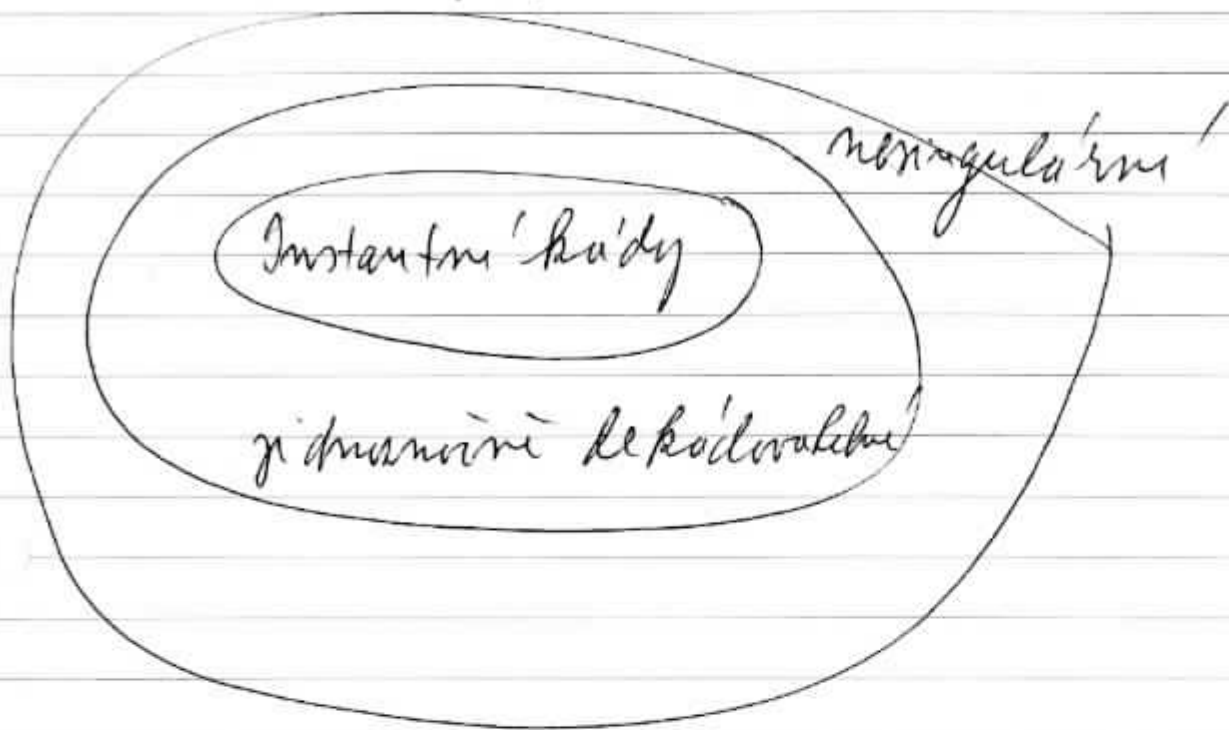
Definice (1) Kód C je nesingulární
 \exists -li C prosté řazení z X do D^+

$$\exists j: c = c' \Rightarrow C(c) \neq C(c')$$

(2) Kód C je jednoznačně dekodovatelný
jistě C^+ je prosté řazení

(3) Kód C je instantně jistě
žádné kódové slovo není předponou
jiného kódového slova.

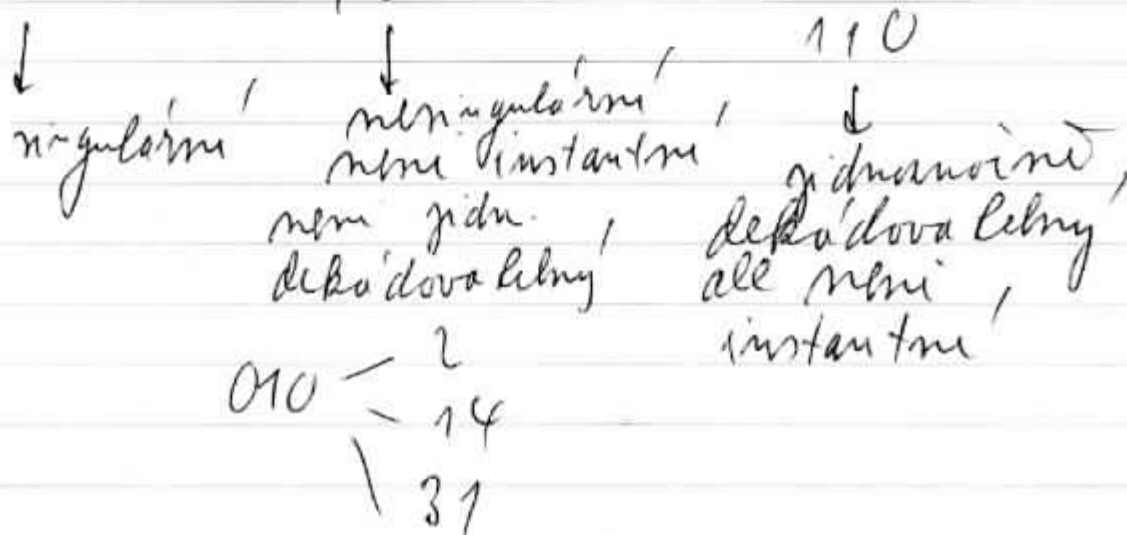
instantm' = prefix code



Pr'klad:

2C

1	0	0	10
2	0	010	00
3	0	01	11
4	0	10	



stromová implementace kódu



uzel = kód (mimo počátek)
délka = počet předků

Věta (Kraftova nerovnost)

Pro instantní kód c nad D -ární abecedou s délkami kódových slov

l_1, l_2, \dots, l_m

platí
$$\sum_{i=1}^m |D|^{-l_i} \leq 1.$$

Naopak, k daným l_1, \dots, l_m splňujícím nerovnost (6) existuje instantní kód s těmito délkami kódových slov.

Dikar c. 1

D-nármí strom

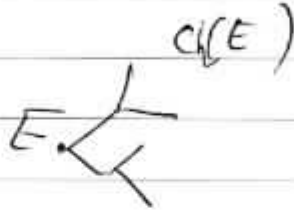
kód = usel

délka kódu = úroveň uslu

E - usel

$Ch(E)$ - dít uslu E

$Ch^l(E)$ - dít uslu E na úrovni l



Nějme E_1, E_2, \dots, E_m usly s délkami l_1, l_2, \dots, l_m

Užse předpokládá, že $l_1 \leq l_2 \leq \dots \leq l_m = l_{max}$

• Kód je instancí $\Leftrightarrow E_i \notin Ch(E_j) \forall i \neq j$

$\Leftrightarrow Ch(E_i) \cap Ch(E_j) = \emptyset \forall i \neq j$

k bodu z prvním mám spíše zdůvodnění sice
má domé úroveň \Rightarrow dva sice $l_{max} - l_i$ - mělo E_j y dít E_j

• $|Ch^{l_{max}}(E_i)| = D$; $i=1, \dots, m$

• je-li kód instancí, pak

$\sum_{i=1}^m |Ch^{l_{max}}(E_i)| \leq D^{l_{max}}$ (z disjunkčnosti)

$\Rightarrow \sum_{i=1}^m D^{l_{max}-l_i} \leq D^{l_{max}} \Rightarrow \sum_{i=1}^m D^{-l_i} \leq 1$

Mějme naopak délky $l_1 \leq l_2 \leq \dots \leq l_m$
 splňující $\sum_{i=1}^m D^{-l_i} \leq 1$. $\sum_{i=1}^m D^{l_{max}-l_i} \leq D^{l_{max}}$

Ustavíme inst. kód s těmito délkami.

• Vybereme jakýkoliv uzel E_1 s délkou l_1 .

$D^{l_{max}-l_1} \leq D^{l_{max}}$ a tedy existuje uzel na úrovni l_{max} , který není

v $Ch(E_1)$. Vezmeme jeho otce E_2 na úrovni l_2 . Musí platit $Ch(E_2) \cap Ch(E_1) = \emptyset$.
 Jelikož

$$|Ch(E_1) \cup Ch(E_2)| = D^{l_{max}-l_1} + D^{l_{max}-l_2} < D^{l_{max}}$$

existuje uzel E_3 na úrovni l_3 tak, že

$Ch(E_3), Ch(E_2), Ch(E_1)$ jsou disjunkční

atd.

poslední krok: $\sum_{i=1}^{m-1} D^{l_{max}-l_i} < D^{l_{max}}$

a najdeme E_m .



Důkaz č. 2

D-mákní rovnaj

$$\underbrace{y = y_1 \dots y_l}_{\text{ kód } y, y_i \in \{0, \dots, D-1\}} \longrightarrow \text{číslo } n \in \langle 0, 1 \rangle$$
$$y = \sum_{j=1}^l y_j \bar{D}^j = \frac{y_1}{D} + \frac{y_2}{D^2} + \dots + \frac{y_l}{D^l}$$

Všechna paknácování
téhoto rovnaj jsou menší než

$$0, y_1 y_2 \dots y_l + \underbrace{\sum_{j=l+1}^{\infty} (D-1) \bar{D}^j}$$

$$= (D-1) \frac{1}{D^{l+1}} \cdot \frac{1}{1-1/D} = \frac{1}{D^l}$$

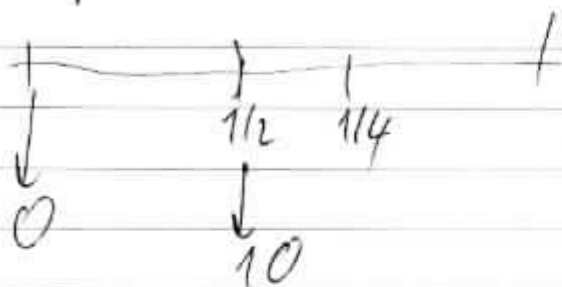
paknácování $\equiv \langle y, y + \frac{1}{D^l} \rangle$

Délka tohoto intervalu je $1/D^l$.

Máme nyní instancní kód s délkami l_1, l_2, \dots, l_m .

Jim odpovídají disjunktní intervaly v $\langle 0, 1 \rangle$ délky $\bar{D}^{l_1}, \bar{D}^{l_2}, \dots, \bar{D}^{l_m}$. Součet těchto délek je \leq menší než 1.

Naopak, máme-li délky
 $l_1 = l_2 = \dots = l_m$ splňující Kraftovu
 nerovnost, pak jsme nuceni přirodit
 intervaly od dubního konce délky l_i
 např. binární kód $l_1 = 1, l_2 = 2$



0, 10 - je instancí kód.

□

Příklad Kolik prosti má Manton?
 Máme kódová slova instancího
 kódu délky

1, 1, 2, 3, 2, 3

jak velika je abeceda?

Rěšení: $f(0) = \bar{0}^{-1} + \bar{0}^{-1} + \bar{0}^{-2} + \bar{0}^{-3} + \bar{0}^{-2} + \bar{0}^{-3}$
 $= 2\bar{0}^{-1} + 2\bar{0}^{-2} + 2\bar{0}^{-3} \leq 1$

$f(2) > 1$

$f(3) = 2\left(\frac{1}{3} + \frac{1}{9} + \frac{1}{27}\right) = \frac{26}{27} \leq 1$

Musíme vzít alespoň 3 znaky.

Pauzami: l_1, \dots, l_m - délky instantních kódů pro D

=) jsou dobré
i pro $D+1, D+2, \dots$

Překvapivý výsledek

Věta (McMillan, 1956)

Kódové délky l_1, l_2, \dots, l_m
řidmanovně dekadová kódy kódu
o D -násobné abecedě splňují nerovnost

$$\sum_{i=1}^m \frac{1}{D^{l_i}} \leq 1.$$

Důk: C^+ ulense kódu $C = C^+ \rightarrow D^*$

Kolik je nejvíce prvků z C^+ , které
 C^+ zahrnuje do D^m (kód délky m)
jelikož je C^+ prvek je takových
prvků nejvíce $|D^m| = D^m$

$u \in X \dots l(u) \dots$ délka kódů
 Filozofie přirození K

$$\left(\sum_{u \in X} D^{-l(u)} \right)^K =$$

$$\sum_{u_1 \in X} \sum_{u_2 \in X} \dots \sum_{u_K \in X} D^{-l(u_1)} \cdot D^{-l(u_2)} \dots D^{-l(u_K)}$$

$$= \sum_{(u_1, \dots, u_K) \in X^K} D^{-l(u_1 \cdot \dots \cdot u_K)} = \sum_{m=1}^{K \cdot l_{\max}} a(m) D^{-m}$$

↓
 počet posl. $u_1, \dots, u_K \in X^K$
 které mají rozšířeny
 kód délky m

víme $|a(m)| \leq D^m$

$$\hookrightarrow \leq \sum_{m=1}^{K \cdot l_{\max}} D^m \cdot D^{-m} = K \cdot l_{\max}$$

$$\Rightarrow \sum_{u \in X} D^{-l(u)} \leq (K \cdot l_{\max})^{1/K}$$

↓ $K \rightarrow \infty$

↑

$$\Rightarrow \sum_{u \in X} D^{-l(u)} \leq 1$$

Optimální kód

$$k \in \mathbb{R}$$

$T(k)$ = nejmenší celé číslo $\geq k$.

Tedy $k \leq T(k) < k+1$

Optimální kód je instancí kód
(jednoduše dekódovatelný) s nejmenší
střední délkou.

Existence: $\vec{p} = (p_1, p_2, \dots, p_m)$
 $p_1 \geq p_2 \geq \dots \geq p_m > 0$.

Je L je délka konkrétního instancího
kódu.

Pro jakýkoliv instancí kód jeho
algoritmicky kódové slovo má délku
větší než L (p_m je jeho délka
větší než L).

Je se tedy omezt na kódy s
 $l_{max} \leq \lceil 4/p_m \rceil$ a těch je konečné
množství.

Def: Instantní kód C^b optimální
jistě

$$L(C^b) \leq L(C)$$

pro každý instantní kód daneho
sdružení L^b optimální délka

mažeme dostatečně velké (velké D)
možeme dosáhnout $L^b = 1$.

čísly dávkou mezi pod klauz možeme
jit:

Věta (Shannonova věta o sdruženém
kódování)

Střední délka L libovolného instantního
 D -márního kódu pro náhodnou
veličinu X je větší nebo rovná

$H_D(X)$:

$$L \geq H_D(X).$$

Rovnost nastává $L = H_D(X)$ $p_i = 2^{-l_i}$

Diškas: máme předpoklad $p_i > 0 \forall i$.

l_1, l_2, \dots, l_m - délky kódových slov

$$L - H_D(X) = \sum_{i=1}^m p_i l_i - \sum_{i=1}^m p_i \log_D \frac{1}{p_i} =$$

$$= - \sum_{i=1}^m p_i \log_D D^{-l_i} + \sum_{i=1}^m p_i \log_D p_i$$

$$= \sum_{i=1}^m p_i \log_D \frac{p_i}{D^{-l_i}}$$

dě $c = \sum_{i=1}^m D^{-l_i} \leq 1$ - Kraft!

$$r_i = \frac{D^{-l_i}}{c} \Rightarrow \vec{r} = (r_1, \dots, r_m)$$

\vec{r} - pravděpodobnostní vektor.

$$L - H_D(X) = \sum_{i=1}^m p_i \log_D \frac{p_i}{r_i c} =$$

$$= \sum_{i=1}^m p_i \log_D \frac{p_i}{r_i} - \log_D c$$

$$= D(\vec{p} \parallel \vec{r}) + \log_D \frac{1}{c} \geq 0$$

≥ 0 .

$c \leq 1 \Rightarrow \geq 0$

$$\Rightarrow L \geq H_D(X).$$

Rovnost nastane $\Leftrightarrow \vec{p} = \vec{\bar{x}}$ a $c=1$

□

- Optimalní délka l_i noma entropie $L=1$ rozdělení p_i D -adické a Kraftova nerovnost p_i satisfikuje
 $p_i = D^{-l_i}$

- Jak blízko k entropii můžeme jít?

Shannonův kód:

$$l_i = \lceil \log_D \frac{1}{p_i} \rceil$$

Tyto délky splňují Kraftovu nerovnost:

$$\sum_{i=1}^m D^{-\lceil \log_D \frac{1}{p_i} \rceil} \leq \sum_{i=1}^m D^{-\log_D \frac{1}{p_i}} = \sum_{i=1}^m p_i = 1$$

Výše uvedená délka vyhovuje nerovnosti

$$\log_D \frac{1}{p_i} \leq l_i < \log_D \frac{1}{p_i} + 1$$

$$\sum_{i=1}^m p_i \log_D \frac{1}{p_i} \leq \sum_{i=1}^m l_i p_i < \sum_{i=1}^m p_i \log_D \frac{1}{p_i} + 1$$

Tedy $H_D(X) \leq L < H_D(X) + 1$.

Príklad

$$\vec{p} = (1/2, 1/4, 1/8, 1/16, 1/16)$$

Jaké jsou délky Shannonova kódu:

1, 2, 3, 4, 4.

slučm' délka

$$1/2 + 1/2 + \frac{3}{8} + \frac{4}{16} + \frac{4}{16} = \frac{16 + 6 + 8}{16} = \frac{30}{16} = \frac{15}{8}$$

$$H(\vec{p}) = 1,875$$

$$= 1,875$$

Optimální délka musí být jen lepší

$$HD(X) \leq L^* \leq HD(X) + 1$$

Věta: Pro optimální délku kódového slova L^* platí

$$HD(X) \leq L^* < HD(X) + 1.$$

• Neexistují firmule pro výpočet L^*

• binární kód \equiv série ano a ne sloček
Shannon: průměrný počet sloček k identifikaci abychom daného sdružení sdala omezen entropie a nekomplikovali sádky, že je pro entropie + 1 bit.

Příklad: Hrač A valí alykt ve všemine a hráč B identifikuje alykt série ano-ne sloček. V průměru potřebuje 38,5 sloček k určení alyklu. Valí optimální strategii. Jak velky je soubor alyklu?

Riešení: $L^* = 38,5$

$$\log |X| \geq H(X) \geq L^* - 1$$
$$= |X| \geq 2^{37,5} = 1,94 \cdot 10^{11}$$

Blokové kódování

• lze dosáhnout ~~entropie~~ ^{délky} libovolně blízké entropii

Mřížka: nezávislé X ale (X_1, X_2, \dots, X_n)
množina op. kódování \mathcal{X}



L_n^* ... střední délka optimálního kódu
pro (X_1, X_2, \dots, X_n) .

$$H(X_1, X_2, \dots, X_n) \leq \frac{L_n^*}{n} \leq H(X_1, \dots, X_n) + 1$$

$$n H(X) \leq L_n^* \leq n H(X) + 1$$

$$H(X) \leq \frac{L_n^*}{n} \leq H(X) + \frac{1}{n}$$

↓
střední délka na 1 symbol

Saunisi' a rychlost' entropie - postěje.

Člna se špatny' kód:

$X \dots x - p(x)$ rozdělení

$Y \dots y = x - q(x)$ rozdělení

Vynalíme Shannonov' kód dle $q(x)$:

$$l(x) = \lceil \log \frac{1}{q(x)} \rceil$$

Věta (špatny' kód) střední délka, L ,
kódu a délkami

$$l(x) = \lceil \log \frac{1}{q(x)} \rceil$$

splňuje nerovnost

$$H(X) + D(X||Y) \leq L < H(X) + D(X||Y) + 1$$

Důkaz: $L = \sum_{x \in X} p(x) l(x) =$

$$= \sum_{x \in X} p(x) \lceil \log \frac{1}{q(x)} \rceil \leq$$

$$\sum_{x \in X} p(x) \left[\log \frac{1}{q(x)} + 1 \right] =$$

$$= \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \cdot \frac{1}{p(x)} + 1$$

$$= D(X||Y) + H(X) + 1$$

druhá rovnost se dokáže stejně.

Příklad:

binární kódy

Symboly	$p(x)$	$q(x)$	$c_1(x)$	$c_2(x)$
1	1/2	1/2	0	0
2	1/4	1/8	10	100
3	1/8	1/8	110	101
4	1/16	1/8	1110	110
5	1/16	1/8	1111	111

c_1 je optimální kód pro p :

$$H(\vec{p}) = L(c_1) = 1,875$$

c_2 je optimální kód pro q :

$$H(\vec{q}) = 2 = L(c_2)$$

Uvažme kód c_2 pro distribuci \vec{p} .

Kód c_2 je Shannonův pro \vec{q} .

Kód c_1 je Shannonův pro \vec{p} .

distanční délka

$$1/2 + 3/4 + 3/8 + \frac{6}{16} = \frac{8+12+6+6}{16} =$$

$$= \frac{32}{16} = 2$$

Rozdíl od $H(\vec{p}) = L(c_1)$ je 0,125

c_2 je $D(p||q)$.