

Huffmanovo kódování

Huffman našel v r. 1951 optimální kód C_H .

Příklad: $\mathcal{A} = \{a, b\}$
 $p(a) = 0,9999$
 $p(b) = 0,0001$

Huffman: $C_H(a) = 0$
 $C_H(b) = 1$ $L(C_H) = \text{optimální}$
dílna
= 1

Shannon: dílny 1, 14
 $L(C_S) = 0,9999 + 0,001477$

algoritmus Huffmanova kódování:

- Seřadíme symboly sestupně dle pravděpodobnosti
- Stavíme dva nejmenší pravděpodobné symboly do jednoho a zastavíme stavěním pravděpodobnosti
- Pak navíme až dostaneme pouze D znaků a ty "okódujeme". $C_i = i \cdot D - 1$.
- Zpět kódujeme větven!

Príklad: hmotnosť kád pr
 $\vec{p} = (0,25; 0,25; 0,20; 0,15; 0,15)$

x_i	$p(x_i)$	x_i	$p(x_i)$	x_i	$p(x_i)$	x_i	$p(x_i)$	x_i	$p(x_i)$	hod
1	0,25	(4,5)	0,30	(2,3)	0,45	(1,4,5)	0,55	0,1	0,1	0,1
2	0,25	1	0,25	(4,5)	0,30	(2,3)	1	0,45	1,0	1,0
3	0,20	2,10	0,25	(1,0)	0,25				1,1	1,1
4	0,15	3,11	0,20						0,00	0,00
5	0,15								0,01	0,01

$$L(CH) = 2 \cdot 0,25 + 2 \cdot 0,25 + 2 \cdot 0,2$$

$$+ 3 \cdot 0,15 + 3 \cdot 0,15 = 2,3$$

$$H(X) =$$

seřiz oblaček: jsm n $\{1, 4, 5\}$ nebo $n \{2, 3\}$

↓

jsem 1 nebo $\{4, 5\}$

↓

jsem 4 nebo jsem 5

Entropie: dubn' hodnotn' pro prvn'ímým' jsm oblaček vedouc' k identifikaci oblačku (znak)

Pr'klad: $D=3$

u	$p(u)$	u	$p(u)$	$C(u)$
1	0,25	$\{3, 4, 5\}$	$0,5^0$	1
2	0,25	1	$0,25^1$	2
3	0,2	2	$0,25^2$	00
4	0,15			01
5	0,15			02

$$L(CH) = 1,5$$

Príklad: $A = \{a, b, c, d\}$

$$\vec{p} = (1/3, 1/3, 1/4, 1/12)$$

Huffmanov kód dá dĺžky
(2, 2, 2, 2) alebo (1, 2, 3, 3)

$$L(CH) =$$

Príklad: zpráva obsahuje 10^4 znakov = 10^4 bajtov
z abecedy $\{a, b, c, d, e\}$

z prahovou hodnotou

$$\vec{p} = (0,35; 0,17; 0,17; 0,16; 0,15)$$

$$L(CH) = 2,3$$

z komprimovaná zpráva: 2300 bajtov

$$10^4 \cdot L(CH) = 2875 \text{ bajtov}$$

$$H(X) = 2,23284 \text{ nosič an. 31.}$$

↓
keď sa použije blokový
kódovanie.

Kompresný poměr: 28,75%.

vyhody : optimalita
: snadna implementace
: nema patentova ochranu
: je instantni

nevychody : vyžaduje naiteni alého
: sdružení dat kvůli unáve
: četnosti (ale existuje adaptivní
: kódování)

: nutno připravit kódovací
: tabulku
: nejsou explicitně dány délky
: kódových slov

• nema ziskovani

Poznámka: U D-národního kódu
přiblížíme na konci D-značku.
Po b-krocích máme

$|X| - kD$ -- značka.

Abychom dostali D značku musíme
někdy na začátek přidat symboly
"Dummy" s nulovými pravděpodobnostmi.

úvahy o optimálním kódu

p_1, p_2, \dots, p_m pravděpodobnosti $p_1^2, p_2^2, \dots, p_m^2$
 l_1, l_2, \dots, l_m příslušné délky
kódových slov

Lemma 11 Každý optimální kód
splňuje vlastnost mono-lamie:

$$p_j > p_k \Rightarrow l_j \leq l_k$$

12) Pro instancí optimální kód platí,
že každý kód maximální
délky má souhlas.

(13) Existuje optimální kód tak, se kódová slova odpovídající dvoje nejmenšímu proudeč podstatně se liší v posledním bitu.

Důk: (1) Když ne tak přehrazení příslušných kódových slov skutečné délky jsou dva

$$p_1 > p_2 \Rightarrow \underbrace{p_1 l_1 + p_2 l_2}_{\text{optimální}} \leq \underbrace{p_1 l_2 + p_2 l_1}_{\text{přehrazení délky}}$$

$$\Rightarrow (p_1 - p_2)(l_2 - l_1) \geq 0 \Rightarrow l_2 \geq l_1$$

"swap"

(2) ať existuje kód s maximální délkou bez souhlasu



pak skutečným o větvi bez souhlasu kód vyčerpání (ať bude instantsní) - Spar.

(3) Každý kód max. délky má souhlasu. Převodíme otuž dvoje dvoje nejmenšímu proudeč podstatně

Tato výměna neovlivní střední délku.
(příklad: 1001 (sčet prw. posl. úroveň))

Def: Kanonický kód je optimální
instancí kód splňující (1), (2), (3)
výše.

Příklad: Které z následujících kódů
je Huffmanovy kódy

a) 20, 10, 111 ✓ pro $p = (1/2, 1/4, 1/4)$

b) 200, 01, 10, 1101

↑
nera srovnání, musí být
skočce

c) 201, 101 - není optimální

řivady k optimalitě Huffmanova kódu

$$\vec{p} = (p_1, p_2, \dots, p_m)$$

$$p_1, p_2, \dots, p_m$$

L_m^* - optimální délka

$$\vec{q} = (p_1, p_2, \dots, p_{m-1} + p_m)$$

L_{m-1}^* - optimální délka

Expanse kódu: $C_{m-1}^*(\vec{q})$ - kanonický kód pro \vec{q}

Přivedeme mu kód pro \vec{p} následujícími operacemi

$C_{m-1}^*(\vec{q})$ kódy délky				$C_m(\vec{p})$	
p_1	w_1'	l_1	\rightarrow	p_1	w_1'
p_2	w_2'	l_2	\rightarrow	p_2	w_2'
			\rightarrow		
p_{m-2}	w_{m-2}'	l_{m-2}	\rightarrow	p_{m-1}	w_{m-1}'
$p_{m-1} + p_m$	w_{m-1}'	l_{m-1}	\rightarrow	p_m	w_{m-1}'

Snováme si délky

$$(*) L(C_m(\vec{p}^*)) = L_{m-1}^* + p_{m-1} + p_m$$

Naspek, máme kanonický kód pro \vec{p}^* . Sestrojíme kód pro \vec{q}^* tak, si to udělá \log a skutečně souhlasí s minimální prahovou hodnotou a přirovnání ho prahové hodnotě $p_m + p_{m-1}$. Nový kód, $C_{m-1}(\vec{q}^*)$ má délku stejnou

$$\begin{aligned} \leftarrow & \rightarrow p_m + p_{m-1} \\ \rightarrow & \rightarrow p_m + p_{m-1} \end{aligned}$$

$$(**) L(C_{m-1}(\vec{q}^*)) = L_m^* - p_{m-1} - p_m$$

Sčítáme (*) a (**):

$$L(C_m(\vec{p}^*)) + L(C_{m-1}(\vec{q}^*)) = L_{m-1}^* + L_m^*$$

$$\Rightarrow \underbrace{(L_{m-1}^* - L(C_{m-1}(\vec{q}^*)))}_{\geq 0} + \underbrace{(L_m^* - L(C_m(\vec{p}^*)))}_{\geq 0} = 0$$

$$\Rightarrow \left. \begin{aligned} L(C_{m-1}(\vec{q}^*)) &= L_{m-1}^* \\ L(C_m(\vec{p}^*)) &= L_m^* \end{aligned} \right\} \Rightarrow \begin{aligned} C_m(\vec{p}^*) &\text{ je} \\ C_{m-1}(\vec{q}^*) &\text{ je optimální!} \end{aligned}$$

Tedy v Huffmanově kódování
se optimální strom v abec
směněch. Tedy máme:

Věta: Huffmanův kód je optimální!

Důk: Předchozí úvahy + kód pro
2-prvkovou abecdu je optimální!

Příklad: špatně úloha

Máme 6 různých úloh, jedno je zkušební a
to a pravidel.

$$\left(\frac{8}{23}, \frac{6}{23}, \frac{4}{23}, \frac{2}{23}, \frac{2}{23}, \frac{1}{23} \right)$$

a) Uspořádáme úlohy po úlohách. Jakby' je
přibližně počet úloh.

Benome to od největší' pro k nejmenší'

$$\begin{aligned} \sum p_i l_i &= 1 \cdot \frac{8}{23} + 2 \cdot \frac{6}{23} + 3 \cdot \frac{4}{23} + 4 \cdot \frac{2}{23} \\ &+ 5 \cdot \frac{2}{23} + 6 \cdot \frac{1}{23} = \frac{55}{23} = 2,39 \end{aligned}$$

b) Ypsi strategie. Kóme skúch!
Láve do mi mušine mišat vranky

Adiláme Huffmanov kód: Dvo
díky $(2, 2, 2, 3, 4, 4)$

$$\sum_{i=1}^6 p_i l_i = 2,35 \text{ - optimalní}$$

Shannon-Fano-Elias kód

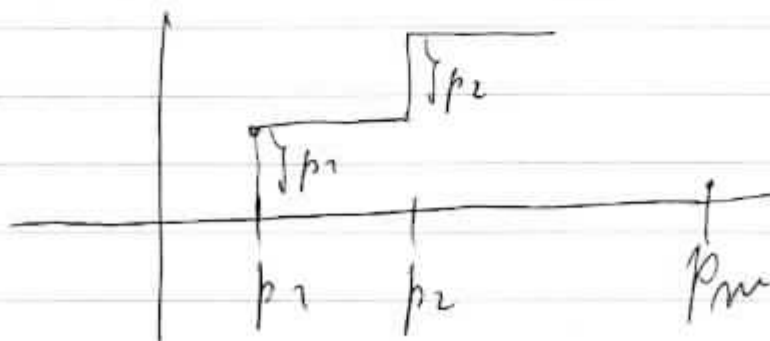
$$X = \{1, 2, \dots, m\}$$

$$p(a) > 0 \quad \forall a$$

Distribuční funkce $F(x) = \sum_{a \leq x} p(a)$

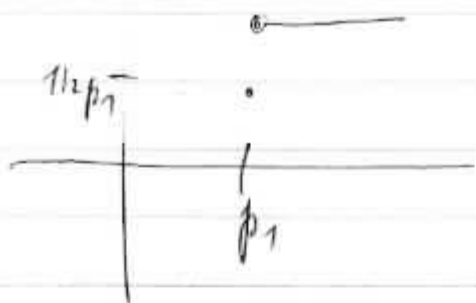
(zprava spajitá)

$$\dots \frac{1}{\dots}$$



Modifikace

$$\bar{F}(x) = \sum_{a \leq x} p(a) + \frac{1}{2} p(x)$$



$\bar{F}(x)$ rovinně dyodicky a vzhledem k
vzhledem $l(x)$ klesá.

Značím $[\bar{F}(x)]_{l(x)}$

$$\text{Pak } l(x) = \lceil \log \frac{1}{p(x)} \rceil + 1$$

$$\text{Pak } \log \frac{1}{p(x)} + 1 \leq l(x) \leq \log \frac{1}{p(x)} + 2$$

$$2 \cdot \frac{1}{p(x)} \leq 2^{l(x)} \leq 4 \cdot \frac{1}{p(x)}$$

$$\Rightarrow 2^{-l(x)} \leq \frac{p(x)}{2}$$

$$\text{Dále } \bar{F}(x) - [\bar{F}(x)]_{l(x)} <$$

$$< \frac{1}{2^{l(x)}} \leq \frac{p(x)}{2}$$

Tedy $[\bar{F}(x)]_{l(x)}$ existuje vnitřní
okružní odpovídajícímu x .

Tedy $l(x)$ bude stav k papíru x .
 Namísto dastaneme instancí kód
 neboť všechna pokračování kódu
 $z_1 \dots z_n$ jsou v intervalu

$$\left[0, z_1 \dots z_n; 0, z_1 \dots z_n + \frac{1}{2^n}\right)$$

Tyto intervaly jsou disjunktní.
 neboť mají délku $\frac{1}{2^{n(x)}} < \frac{p(x)}{2}$ a $[\bar{F}(x)]_{x(x)} \in \bar{F}(x)$
 Délka kódu: $\text{ale } z_{n(x)} \text{ je podmíněně}$

$$L = \sum_{x \in X} p(x) l(x) = \sum p(x) (\lceil \log \frac{1}{p(x)} \rceil + 1)$$

$$L \leq H(X) + 2$$

Problém:

x	$p(x)$	$\bar{F}(x)$	binární	$l(x)$	kód
1	0,25	0,125	0,001	3	001
2	0,5	0,5	0,10	2	10
3	0,125	0,8125	0,1101	4	1101
4	0,125	0,9375	0,1111	4	1111

0,75 — 0,8
 0,25 — 0,5
 0,125 = 2^{-3}

Střední délka je jedinečná a kritická
jakožto řada. Ale důležitá je
a relativní úspěšnost

Věta Necht' $l(x)$ je kódová délka
asociovaná s Shannonovým kódem.
Necht' $l'(x)$ je kódová délka
asociovaná s libovolným jiným
jednosměrně dekódovatelným kódem.

Pak

$$P[l(x) \geq l'(x) + c] \leq \frac{1}{2^{c-1}}$$

např.: pokud, že $l(x)$ je kódová S je
 $\frac{1}{24} = \frac{1}{16}$.

kódová 3 a $\text{prav.} \leq 1/2$.

kva se dvěma kódy a to bylo má
lepší kód.

Dikaz

$$P[l(x) \geq l'(x) + c] \leq$$

$$\leq P\left[\log \frac{1}{p(x)} \geq l'(x) + c - 1\right]$$

$$= P\left[p(x) \leq 2^{-l'(x) - c + 1}\right]$$

$$= \sum p(u)$$

$$\sum_{u|p(u) \leq 2^{-l'(u) - c + 1}} p(u)$$

$$\leq \sum_{u|p(u) \leq 2^{-l'(u) - c + 1}} 2^{-l'(u) - c + 1}$$

$$\leq \sum_{u|p(u) \leq 2^{-l'(u) - c + 1}} 2^{-l'(u) - c + 1}$$

$$\leq \sum_{u|p(u) \leq 2^{-l'(u) - c + 1}} 2^{-l'(u) - c + 1}$$

Kraftov nerovnost.

Dalším rozborem lze ukázat, že
Shannonův kód vyhovuje
všechnu podmínkám.

Entropie stochastických procesů

časová řada, diskrétní stochastický proces a
postupnost náhodných veličin

$$X_1, X_2, \dots$$

na prv. prostoru (Ω, \mathcal{A}, P)
Předpoklad: $X_i: \Omega \rightarrow \mathcal{X}$ - konečná množina

Definice. Stochastický proces $(X_i)_{i=1}^{\infty}$

se nazývá stacionární zkrátě

$$P[X_1 = \omega_1, X_2 = \omega_2, \dots, X_m = \omega_m] =$$

$$P[X_{1+l} = \omega_1, X_{2+l} = \omega_2, \dots, X_{m+l} = \omega_m]$$

$$\forall m, l = 1, 2, \dots \text{ a } \omega_1, \dots, \omega_m \in \mathcal{X}$$

- Každý nezávislý proces je stacionární
- Zorníkat parse na předchozí hodnotě
- Markovův řetězec

Def. Časová řada X_1, X_2, \dots
se nazývá Markovův řetězec
zkrátě

$$P[X_{m+1} = \omega_{m+1} \mid X_m = \omega_m, X_{m-1} = \omega_{m-1}, \dots, X_1 = \omega_1] =$$

$$= P[X_{m+1} = \omega_{m+1} \mid X_m = \omega_m] \quad \forall m, \omega_1 \rightarrow \omega_{m+1} \in \mathcal{X}$$

Príklad: $(X_n)_n$ - Markovian řada

$$P[X_3 = a_3, X_2 = a_2, X_1 = a_1] =$$

$$= P[X_3 = a_3 | X_2 = a_2, X_1 = a_1] \cdot P[X_2 = a_2, X_1 = a_1]$$

$$= P[X_3 = a_3 | X_2 = a_2] \cdot P[X_2 = a_2 | X_1 = a_1] \cdot P[X_1 = a_1]$$

Vedeme k neličovému pravidlu

$$p(a_1, \dots, a_n) = p(a_1) \cdot p(a_2 | a_1) \cdot p(a_3 | a_2) \cdot \dots \cdot p(a_n | a_{n-1})$$

Def Markovian řada je časově invariantní řada

$$P[X_{m+1} = a | X_m = b] = P[X_2 = b | X_1 = a]$$

$\forall m, \forall a, b \in \mathcal{X}$

Věta: Všechny Markovské řady jsou časově invariantní.

Průchodová matice

$$X = \{1, 2, \dots, m\}$$

$$P = (P_{ij})$$

$$P_{ij} = P[X_{m+1} = j \mid X_m = i] = P[X_2 = j \mid X_1 = i]$$

\vec{p}_m - prav. vektor pro X_m

$$\vec{p}_{m+1} = \vec{p}_m \cdot P$$

$$\Rightarrow \vec{p}_{m+1} = \vec{p}_1 \cdot P^m$$

Def: Stacionární roššení \vec{p} Markovského
 \forall řádku i pravdep. vektor

$$\vec{p} = \vec{p} \cdot P$$

$$i.e. \vec{p}_1 = \vec{p} \Rightarrow \vec{p}_m = \vec{p}$$

Stacionární roššení \vec{p} \forall řádku
soustavy rovnic

$$(P^T - E)\vec{p} = 0.$$

řádk. matice

Nutnou podmínkou pro existenci
nultního řádku \vec{p} je regularita $P^T - E$.

ani tehty nemusí písemně existovat - primál
a normalizace $\alpha_1 + \dots + \alpha_n = 1$.

Příklad: Řešete se dvěma stavy $X = \{1, 2\}$.

matice přechodu

$$P = \begin{bmatrix} 1-\alpha & \alpha \\ \beta & 1-\beta \end{bmatrix} \quad \alpha, \beta \geq 0.$$



$$P^T - E = \begin{bmatrix} -\alpha & \beta \\ \alpha & -\beta \end{bmatrix}$$

tedy

$$\begin{aligned} -\alpha p_1 + \beta p_2 &= 0 \\ p_1 + p_2 &= 1 \end{aligned}$$

$$p_1 = \frac{\beta}{\alpha + \beta} \quad | \quad p_2 = \frac{\alpha}{\alpha + \beta} \quad \text{stationární stav.}$$

Definice Entropie stochastického procesu

$$X = (X_1, X_2, \dots)$$

je definována jako

$$H(X) = \lim_{n \rightarrow \infty} \frac{H(X_1, X_2, \dots, X_n)}{n}$$

pokud tato limita existuje

Příklady: a) jsou stříj generují náhodně
znaky a všechny má znaky se stejnou
pravděpodobností. Nezávislé na sobě.

$$H(X_1, \dots, X_n) = \log m^n = n \cdot \log m$$

$$H(X) = \lim_{n \rightarrow \infty} \frac{n \cdot \log m}{n} = \log m$$

b) záměněm i.i.d

$$H(X) = \frac{1}{n} \lim_{n \rightarrow \infty} H(X_1, \dots, X_n) =$$

$$= \frac{1}{n} n \cdot H(X_1) = H(X_1).$$