

Entropie zpráv (řádu 0 a vyšší), modelování. Statistické metody

komprese dat

→ data redundancy removal

- short codes for common events

Model: static, semi-adaptive, adaptive

Entropie S - finite set of source units, C - finite set of codewords $f: S \rightarrow C^+$

entropy of unit x_i $H_i = -\log_2 p_i$ - míra neuvěřitelnosti?

$$H_{\text{avg}}(S) = -\sum_{i=1}^n p_i \log_2 p_i \text{ bits} \quad \text{— průměrná entropie — abeceda délky } n \text{ (velikosti) entropie vzorku!}$$

entropy of source message
 $X = x_1 x_2 x_3 \dots x_k \in S^+$

$$H(X) = -\sum_{j=1}^k \log_2 p_j \quad \text{— } k \text{ je délka zprávy } X$$

length of message X $L(X) = \sum_{j=1}^k d_j \text{ bits}$

redundancy of code K for message X :

$$R(X) = L(X) - H(X) = \sum_{j=1}^k (d_j + \log_2 p_j) \text{ bits}$$

average length of a codeword: $L_{\text{avg}}(K) = \sum_{i=1}^n d_i p_i \text{ bits}$

average redundancy of code K : $R_{\text{avg}}(K) = L_{\text{avg}}(K) - H_{\text{avg}}(K)$

Empirical entropy \triangleright 0-th order ($T \in S^+$)

$$H_0(T) = -\sum_{a \in S} \frac{n_a^T}{n} \log_2 \frac{n_a^T}{n}$$

number of symbols a in message T

\triangleright k -th order entropy

$$H_k(T) = \frac{1}{n} \sum_{w \in S^k} |w_T| H_0(w_T)$$

symbols following each other in T

$$0 \leq H_k(T) \leq H_{k-1}(T) \leq \dots \leq H_0(T) \leq \log_2 |S|$$

$$H_0 = T = \text{abbbaababbaaababa}$$

a	9/17
b	8/17

$$H_0 = - \left(\frac{9}{17} \cdot \log_2 \frac{9}{17} + \frac{8}{17} \cdot \log_2 \frac{8}{17} \right) = 0,9975$$

$$H_1 = w_T$$

a	babbaabb
b	bbaabaaa

$$H_1 = \frac{1}{17} \cdot \left(8 \cdot H_0(\text{babbaabb}) + 8 \cdot H_0(\text{bbaabaaa}) \right)$$

$$= \frac{1}{17} \cdot \left(8 \left(\frac{5}{8} \cdot \log_2 \frac{5}{8} + \frac{3}{8} \cdot \log_2 \frac{3}{8} \right) + \dots \right) = 1,8983$$

Shannon - Fano

- dělení intervalů pravděpodobnosti na poloviny (w nejprůměji)
- získávání optimálního kódu, jednoduchá na implementaci

$$H_{\text{avg}} \leq L_{\text{avg}} \leq H_{\text{avg}} + 1$$

Huffman coding - optimální prefixový kód

- konstrukce bottom-up s využitím prioritní fronty

|| Sibling property - každý node má sourozence (kromě root uzlu)

- vždy se dají seřadit dle klesající pravděpodobnosti tak, že sourozenci budou v listu vedle sebe (luzí uzel na listu pravi, suchá pozice pro pravého)

$$\text{Huffman code} \Leftrightarrow \text{Sibling property}$$

Kraft - MacMillan Inequality

$$\sum_{i=1}^n 2^{-L_i} \leq 1$$

- každý jednoznačný kód s délkami kódů $L_1 \dots L_n$ musí splňovat tuto podmínku

- pokud je suma větší než 1, kód nelze jednoznačně dekodovat

Arithmetic coding

MI-SP-SP-14

(2)

- kódujeme do intervalu rozdeleného dle pravdepodobnosti jednotlivých symbolů
- konce se musí řídit EOF charakterem nebo explicitní dílkou
- kódují blokem entropii $H_{avg} = -\sum p(x_i) \log_2 p(x_i)$
- velmi dobře pro adaptivní model

Celostupňová implementace:

$$\text{Low} = \text{xxxx}00\dots \quad (0,0000)$$

$$\text{High} = \text{yyyy}99\dots \quad (0,9999)$$

- stejná hodnota jsou poslány na výstup a proměnné shiftůvek doleva

$$\text{Low} = \text{Low} + (\text{High} - \text{Low} + 1) \text{HighCP}(x) / \text{Totalfreq} - 1$$

$$\text{High} = \text{Low} + (\text{High} - \text{Low} + 1) \text{LowCP}(x) / \text{Totalfreq}$$

Underflow: $\overset{\text{Low}}{\text{High}} = 49\text{xxxx} \quad \& \quad \text{High} = 50\text{yyyy}$

$$\rightarrow \text{Low} = 4\text{xxxx}0, \text{High} = 5\text{yyyy}9, \text{counter}++$$

- poslední počet dle $\text{Low}[0] = \text{High}[0] = 4$

\rightarrow output 4 a counter 9

\rightarrow v opákném případě 5 a nulý

Adaptive arithmetic coding

- vyčíslování stromu - každý uzel symbol, frekvence a cumulative freq. k níž se přidá

$$a_8, 19, 40$$

- pravděpodobnější blíže ke kořeni

- strom balancovaný jako heap (lavičkový strom)

$$a_2, 12, 16$$

$$a_3, 12, 8$$

Range encoding

- here output ne jako binární číslo, ale jako číslo jiné báse
- menší počet normalizací, větší bitová operace