

Základy teorie informace a kódování, entropie

Kódování: zobrazení $C: X \rightarrow D^*$ distribuívaná náhodná veličina

$$D^* = \bigcup_{k=1}^{\infty} D^k$$

D^* - množina konečných řetězců symbolů \in - ánní abecedy D

obraz $C(x)$ - kódové slovo \rightarrow jeho délka $l(x)$

Střední délka ~~kódu~~ ^{kódu}: $L(c)$ náhodná veličina X s $p(x)$ rozdělením

$$L(c) = \sum_{x \in X} l(x) p(x) = E l(x)$$

$L(c)$ nemůže být menší než entropie daného textu $H(X)$

\rightarrow pokud se rovnají, máme optimální kód

Nesingulární kódy

\rightarrow pokud je zobrazení C prosté = můžeme jednoznačně dekodovat daný kód symbolu
- nemusí ale platit pro všechny

$$x \neq x' \Rightarrow C(x) \neq C(x')$$

Jednoznačně dekodovatelné kódy

\rightarrow pokud je C^* nesingulární

$$[C^*(x_1 x_2 \dots x_n) = C(x_1) C(x_2) \dots C(x_n)]$$

Instanční (prefixový) kód

\rightarrow žádný kódové slovo není prefixem jiného

\rightarrow slova sestávají z kódových slov, mohou na ně být

Huffmanovo kódování

- agregace vždy nejmenší pravděpodobných \rightarrow rychlé přičítání hodnot

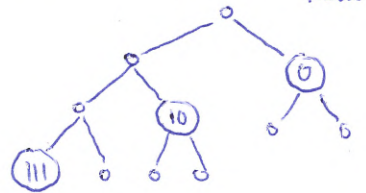
- je optimální: pro Huff. kód C^* a libovolně univ. dec. kód C' $L(C^*) \leq L(C')$

Kraft - McMillan inequality

prefixový kód musí splňovat podmínku: délky kódových slov splňují nerovnost:

$$\sum_i D^{-l_i} \leq 1$$

$|D\text{-číslo abeceda}|$ délky kódových slov



+ pro každou n-tici dělů splňující existující instancí kód

→ McMillan rozšířil nerovnost i na více dekodovatelných

Optimální kódy

$$L(c) \geq H_D(x) \quad \leftarrow \text{entropie} \quad - \text{platí pouze pro instanci (prefixový kód)!}$$

• rovnost nastává pokud $D^{-l_i} = p_i$ pro všechna i

$$H_D(x) \leq L(c^*) \leq H_D(x) + 1$$

\downarrow kódová slova \downarrow prefixový kód

→ optimální kód se vzdáli od ideálu max o 1

Entropie - míra neupřesnění

(pravděpodobnostní funkce $p(x) \rightarrow \forall x \in \mathcal{X} : p(x) = P(X=x)$)

$$\text{Entropie } H(x) = \boxed{H(x) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x)}$$

- platí pouze na rozdíl, ne na číselných hodnotách

• protože $H_b(x)$ může mít základ logaritmu $\log_b p(x)$

Entropie je střední hodnota míry neupřesnění (střední informace)

$$I(x) = -\log p(x) \quad \leftarrow \text{vždy nezáporná, protože 0 je jistých jeví}$$

$$H(x) = \mathbb{E}(I(x))$$

- je nejmenší pro konstantní rozdělení (nejistota)

= střední míra neupřesnění měřící veličiny H

optimální počet binárních slov = mezi $H(x)$ a $H(x) + 1$ (viz optimální kódy)

Schröderova entropie $H(X, Y)$ - entropie sdruženého rozdělení náhodných veličin

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y)$$

- analogicky se definuje pro náhodné vektory

Podmíněná entropie $H(Y|X)$ - podm. ent. náhodných veličin X, Y se sdruženým rozdělením $p(x, y)$

$$H(Y|X) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x)$$

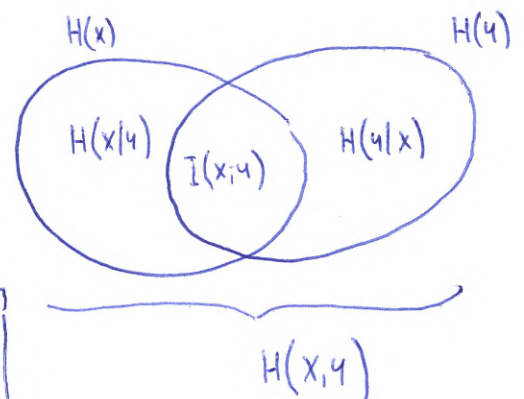
Rekurrenční pravidlo

$$H(X, Y) = H(X) + H(Y|X)$$

\Rightarrow která část informace je v Y navíc oproti X

Relativní entropie (Kullback - Leiblerova vzdálenost)

$$D(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$



Vzájemná informace

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

$$= D(p(x, y) \parallel p(x)p(y))$$

- míra informace, kterou sdílí dvě veličiny X a Y

\rightarrow vzdálenost od nezávislosti

$$I(X, Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y)$$