

Uvažme lineární kód $C: X \rightarrow \Sigma^*$. $L(C)$ vs $H(X)$?

$$H(X) - 2H(X) - 2 \leq L(C)$$

minimální kód C bez prefixů kód tak, že ke každému stavu přidáme jeho

délku zdejším a oddělime 01

$$\text{např. } C(x) = 0011 \rightsquigarrow$$

zavázka
↓
11000001 $C(x)$

4 bitů

a každý bit je zdejším

Generování náhodných čísel dle zadané distribuce vs entropie

komprese: je dána distr. objektů, vyjádřete nejkratší reprezentaci

• OI reprezentace \approx schéma kódů min.

další otázka: chceme sampleovat z výše uvedené distribuce, jaké potřebujeme kódování min?

Pr:
$$X = \begin{cases} a & \text{s pstr. } 1/2 \\ b & \text{s pstr. } 1/4 \\ c & \text{s pstr. } 1/4 \end{cases}$$

0 \rightarrow a

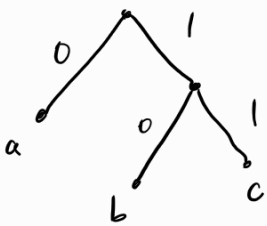
10 \rightarrow b

11 \rightarrow c

\rightarrow průměrně 1.5 bitů

$$H(X) = \frac{1}{2} \lg 2 + \frac{1}{4} \lg 4 + \frac{1}{4} \lg 4 = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = 1.5 \dots \text{náhoda?}$$

• obrázek



průměrně?

• významná tabulka jako algoritmus mapující X na prvky binárního stromu

• vlastnosti

- i) strom by měl být úplný (\approx žádné větvy s 1 potomkem)
- ii) počet dorážení listů v hloubce k je 2^{-k}
- iii) očekávaný # přejetých bitů \Leftrightarrow očekávaná hloubka stromu

• také nějaké zdání, že nezávislost distribuce \Leftrightarrow # bitů ke sampleování rovná bit-ů

• pro vzhled a značení jeho hloubka $h(n)$ a hloubka stromu $h(t)$ entropie

Lemma: Pro lib. úplný strom (tj. # potomků je vždy 0 nebo 2) uvádíme první rozdělání t.j. na listy v hloubce k dává 2^{-k} . Ukážete, že $E[h(T)] = H(\text{první dist. na stráně})$

Důk:

$$E[h(T)] = \sum_{l \in \text{listy}} h(l) 2^{-h(l)}$$

$$H(\text{dist. na listech}) = - \sum_{l \in \text{listy}} 2^{-h(l)} \log 2^{-h(l)} = E[h(\text{listech})] \quad \square$$

algebra entropií

Věta: Pro lib. alg. sampler X , očekávaný # bitů je aspoň entropie X .

↑
t.j. $h(T)$ (což je n.v.)

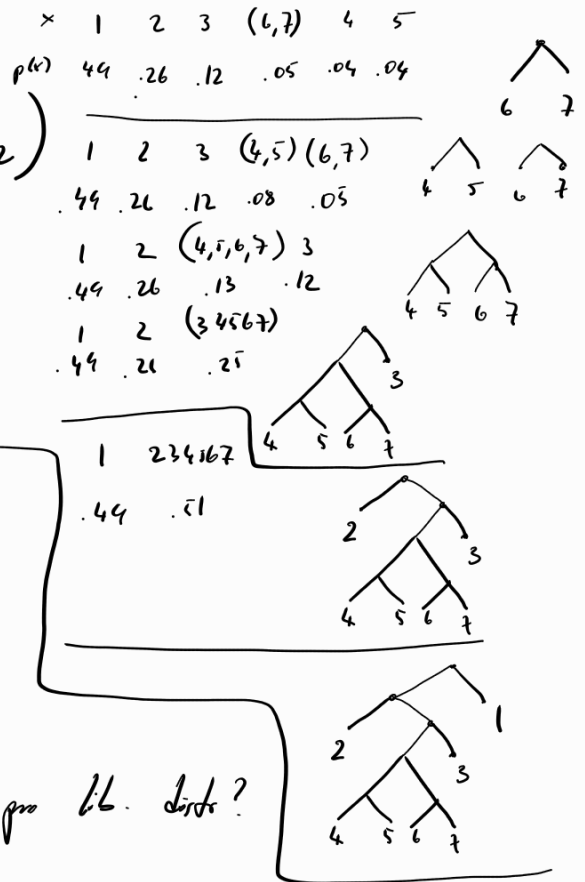
- Důk:
- využití indukce lze rekonstruovat stromem T jako výše
 - ovšem listy stromu jsou $\{1, 2, \dots\}$ (pokud je strom nekonečný, Y taký)
 - uvádíme n.v. Y na listech T t.j. $\Pr[Y=y] = 2^{-h(y)}$ (h je hloubka)
 - lemmata výše $E[h(T)] = H(Y)$
 - n.v. X je vlastně Y , protože chování algoritmu závisí na náhodných listech
- $\Rightarrow H(X) \leq H(Y) \quad \square$

— pokud první jsou vždy mochný den, máme rovnost (inspekce dělení)

· TODO: lepší shod

Uvažme n.v.

$$X = \begin{pmatrix} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 & x_7 \\ 0.49 & 0.26 & 0.12 & 0.04 & 0.04 & 0.03 & 0.02 \end{pmatrix}$$



a) Nalezněte Huffmanův kód pro X.

b) Jaká je přím. délka kód. slova vs entropie? 2.02 2.01

Které z násled. kódů nemohou být Huffmanův kód pro lib. zdroj?

a) $\{0, 10, 11\}$ ano, pro $(\frac{1}{2}, \frac{1}{4}, \frac{1}{4})$

b) $\{00, 01, 10, 110\}$ ne, takže má jednodušší strom

c) $\{01, 10\}$ ne, 1 potřebe

$$x_i = \{1, 0\}$$

Máme n vprávních, a každý z nich je kombinací/vnější s číseli $p_i > \frac{1}{2}$. Pomocí ano-ne dotazů chceme zjistit který z nich jsou vadné. Lib. ano-ne otázka je povolena.

a) Najděte vhodný dotaz odhad na # potřebných dotazů.

b) Dříve tomu, že používáme nejdelší opt. sekci dotazů co do přím. délky. Stále se nám,

že jsou ve sítě, kdy používání zpráva tu největší probl. dotazů. Slavní papírky, na co se ptá poslední dotaz. Může jít jiná sítě rodinné?

c) Uvězte horní odhad na # vyřezávaných dotazů

, použijte Huffmanův kód a polskou hr. odhadování \leftarrow \checkmark
 rozumnost

a) x_i : je itij ikon vobitij? $H(x_1, x_2, \dots, x_n) = \sum H(x_i) = \sum H(p_i)$

b)

pr. dle Hufmannova kódu spájajúce od najvyššieho post. nahor
 - každé post. dle je pta, jstli je post. íkán rostlý

c) $\lceil \sum H(p_i) \rceil + 1$ (= Huffman)

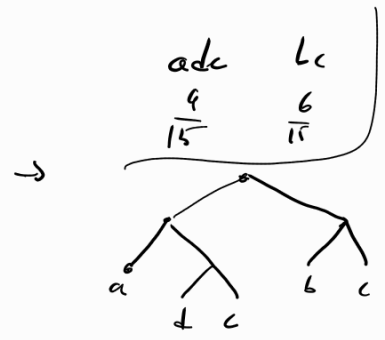
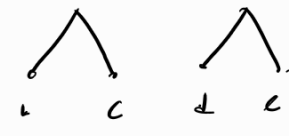
Uvažme Huffmanin kód pro post. distr. $(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{2}{5}, \frac{2}{5})$. Ukaže, že je opt. pro

$(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})$

a	b	c	(d,e)
$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{2}{5}$



a	(b,c)	(d,e)
$\frac{1}{5}$	$\frac{2}{5}$	$\frac{2}{5}$



uvážme lib kód

$E[L(C)] = \sum_{x \in \{a,b,c,d,e\}} p(x) \cdot \text{délka-kódu}(x) = \sum \frac{1}{5} \cdot \text{délka} = \frac{\sum \text{délka}}{5}$

\sum délka musí být celí čísla

$E[\text{Huffman kód}] = 2 \cdot \frac{2}{5} + 3 \cdot \frac{1}{5} = 1 \frac{2}{5}$

lepší kód musí být v číselní ≤ 11

ale to je méně než $H[(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5})]$

Sarce odbyl ten ruk, že, nelze bezstranně komprimovat X pod prům. délka slova $H(X)$ a že opt. kód má prům. délku $\leq H(X) + 1$. Nejde o X t.č. opt. kód má prům. délku slova nekonečně blíže k $H(X) + 1$.

stačí vzít $X = \begin{cases} 1 \\ 0 \end{cases}$ s post. $1-\epsilon$ pro $\epsilon \rightarrow 0$

$H(X) \rightarrow 0$, ale prům. 1 bit je potřeba