

lec 12: Huffman Coding

Recap

* Extended Kraft inequality

* Shannon code

$$l(x) = \left\lceil \log \frac{1}{p(x)} \right\rceil$$

* $L \geq H_D(x)$

$$H_D(x) \leq L \leq H_D(x) + 1$$

* Wrong code

TODAY

* { KI for uniquely decodable codes

* { Huffman Coding

* { Shannon-Fano-Elias Code

Thm (McMillan) The codewords of any uniquely decodable (UD) code must satisfy the K.I.

$$\sum_{x \in \mathcal{X}_0} D^{-l(x)} \leq 1 \rightarrow \textcircled{1} \quad (|\mathcal{X}_0| < \infty)$$

Conversely, given $\{l(x), x \in \mathcal{X}_0\}$ satisfying $\textcircled{1}$, an UD can be found.

Pf: Consider $\left(\sum_{x \in \mathcal{X}_0} D^{-l(x)} \right)^k = \sum_{z_1} \dots \sum_{z_k} D^{-(l(z_1) + \dots + l(z_k))}$

$$= \sum_{z} D^{-l(z)} = \sum_{j=1}^{kl_{\max}} a_j D^{-j}$$

where $a_j = \#$ of strings of the k codewords from the source code with length j .

Note that no two strings of k codewords can result in the same coded string of length j source symbols

Hence $a_j \leq D^j$

$$\therefore \left(\sum_{x \in \mathcal{X}_0} D^{-l(x)} \right)^k = \sum_{j=1}^{kl_{\max}} a_j D^{-j} \leq \sum_{j=1}^{kl_{\max}} D^j D^{-j}$$

$$\therefore \sum_{x \in \mathcal{X}_0} D^{-l(x)} \leq \left(kl_{\max} \right)^{\frac{1}{k}} = kl_{\max}^{\frac{1}{k}} \text{ for any } k \in \mathbb{N}.$$

This must be true for any k and hence true as $k \rightarrow \infty$.

$$\therefore \sum_{x \in X_0} D^{-l(x)} \leq \lim_{k \rightarrow \infty} (k \cdot l_{\max})^{1/k} = 1$$

$$\Rightarrow \sum_{x \in X_0} D^{-l(x)} \leq 1$$

Converse holds, since one can find a prefix code if $\textcircled{1}$ is satisfied and every prefix code is a uniquely decodable code.

Corollary: The above result also holds if X_0 is countably infinite $|X_0| = \infty$.

Pf: let $|X_0| = \infty$, Take a finite subset $S \subseteq X_0$ $|S| < \infty$. let \mathcal{C} be a u.d. code for X_0 , and \mathcal{C}_S be a u.d. code the subset of \mathcal{C} obtained by restricting \mathcal{C} to the set S .

The resultant code \mathcal{C}_S also should be u.d.

$$\Rightarrow \sum_{x \in S} D^{-l(x)} \leq 1$$

let $S = \{x_1, x_2, \dots, x_n\}$ i.e. $|S| = n$.

$$\Rightarrow \lim_{n \rightarrow \infty} \sum_{i=1}^n D^{-l(x_i)} \leq 1$$

The converse holds true, cause one can construct a prefix free code if $\sum_{x \in X_0} D^{-l(x)} \leq 1$, ($|X_0| = \infty$) and prefix free code is u.d.



DATE _____

The Huffman code

Begin with some examples.

$|D| = 2$

(a)

$C(x)$	$x \in X_0$	$p(x)$	A	B	C	D	
01	1	0.25					
10	2	0.25					
11	3	0.2					
000	4	0.15					
001	5	0.15					

$L(c) = (0.7)2 + (0.3)3 = 2.3 \text{ bits}$

$|D| = 3$

(b)

$C(x)$	$x \in X_0$	$p(x)$	A	B
1	1	0.25		
2	2	0.25		
00	3	0.2		
01	4	0.15		
02	5	0.15		

$L(c) = (0.5)1 + (0.5)2 = 1.5 \text{ ternary digits/Source symbol}$

5 5 \rightarrow 3 \rightarrow 1

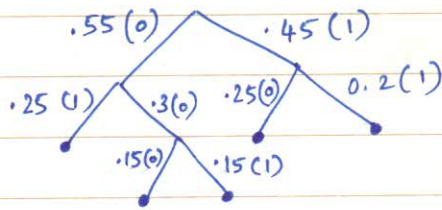
(c) $|D| = 3$

$C(x)$	$x \in X_0$	$p(x)$	A
1	1	0.25	
2	2	0.25	
00	3	0.2	
02	4	0.1	
010	5	0.1	
011	6	0.1	
012	Dummy	0	

$L(c) = (0.5)1 + (0.3)2 + (0.2)3 = 1.7 \text{ ternary digits/Source symbol}$

7 \rightarrow 5 \rightarrow 3 \rightarrow 1

Added a dummy so that you end up with 3 nodes.
 You would want $m = 1 \pmod{q-1}$
 where $q = |D|$ and $m = |X_0|$

Example @Optimality of the Huffman code.

(focus on the binary $\mathbb{D} = \{0, 1\}$ case)
finite x_i

Lemma: Given any $P(x)$, \exists an optimal prefix code satisfying:

1) the lengths l_i satisfy

$$P_j > P_i \Rightarrow l_j \leq l_i$$

2) the two longest codewords have same length

3) 2 of the longest codewords differ only in the last bit and correspond to the 2 least likely symbols.

① Suppose \mathcal{C}_m is optimal

\mathcal{C}_m	\mathcal{C}_m'
$P_i > P_j$	P_i, P_j
(say) $l_i > l_j$	$l_i' = l_j, l_j' = l_i$ rest remain

$$\begin{aligned} L(\mathcal{C}_m) - L(\mathcal{C}_m') &= P_i l_i + P_j l_j \\ &\quad - P_i l_j - P_j l_i \\ &= (P_i - P_j)(l_i - l_j) > 0 \end{aligned}$$

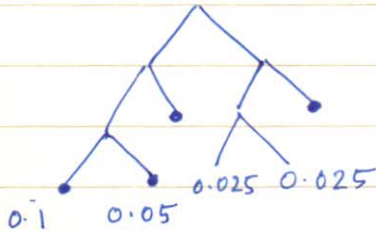
$\Rightarrow \mathcal{C}_m'$ is strictly better code than \mathcal{C}_m , \therefore

\mathcal{C}_m cannot be optimal (Contradiction)

② If not true, the picture would look like



③



Clearly if there are 25 codewords of the same length, these correspond to the smallest 28 probabilities.

Shuffling these codewords will not change the average code length $L(C_m)$. Hence WLOG we can assume ③ holds.

Claim: The Huffman procedure yields optimal code.

$$L(C_m) - L(C_{m-1}) = l_{\max} (P_m + P_{m-1}) - (l_{\max} - 1) (P_m + P_{m-1})$$

$$= P_m + P_{m-1} = \text{constant}$$

$P_1, P_2, \dots, P_{m-1}, P_m$
 $P_1, P_2, \dots, P_{m-2}, (P_m + P_{m-1})$

But C_{m-1} is a code for $(m-1)$ -symbol source.

cause there exists a prefix free code that is optimal with properties ①, ②, ③

\therefore To optimize $L(C_m)$ it is necessary & sufficient that $L(C_{m-1})$ is optimal.

For C_2 there is only one optimal code, therefore the claim follows.