

HW now due on Wednesdays 10pm.

Day 15 - Huffman Code and Random Cryptosystem

Recap - Decryption matrix for rectangular transposition

The following matrix was obtained from the applet for breaking rectangular transposition, using the ciphertext in the first problem of HW4.

	1	2	3	4	5	6	7
1		17	17	21	34	22	24
2	25		20	20	19	27	35
3	16	19		18	22	34	25
4	18	21	33		19	19	21
5	24	32	18	21		20	17
6	26	24	24	22	17		19
7	22	21	15	32	23	20	

Rule: If the big entry is on row i , col j

then j follows i

in decryption perm.

Find the decrypting and encrypting permutations. . One row w/o big entry gives the last entry in perm.

. One col w/o big entry give the first perm entry.

$(a_1, a_2, a_3, a_4, a_5, a_6, a_7) = \text{decrypting perm.}$

. Row 6 has no big entry $\Rightarrow a_7 = 6$.

. To find a_6 , look at Col 6 for the big entry.

On 6th col, big entry is at row 3 (row 3, col 6)

$\Rightarrow a_6 = 3$

. To find $a_5 \Rightarrow$ look for big entry on col 3

$\Rightarrow a_5 = 4$

Decrypting perm: $(1, 5, 2, 7, 4, 3, 6)$

Encrypting perm. is the inverse of decrypting perm.

$$\text{decrypting perm} \parallel (1, 5, 2, 7, 4, 3, 6) = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 5 & 2 & 7 & 4 & 3 & 6 \end{pmatrix} (*)$$

To get the inverse, simply reverse the arrows in $(*)$

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 3 & 6 & 5 & 2 & 7 & 4 \end{pmatrix}$$

$(1, 3, 6, 5, 2, 7, 4)$: encrypting perm.

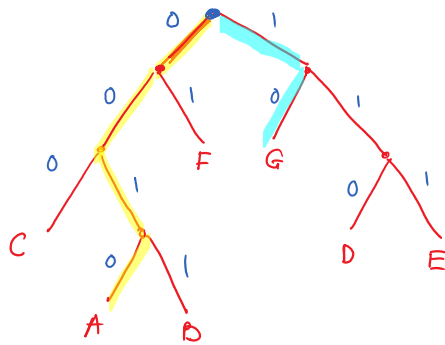
Example. Suppose a certain file contains only the letter with the following frequencies

frequency →

A	B	C	D	E	F	G
1	2	2	4	4	5	6

Construct the comma-free code that enables you to compress the file so that you can store it using the least number of bits.

- Sort frequencies in increasing order
- Group the smallest 2 entry, and replace w/ their sum.
- Re-sort, repeat ...



A → 0010 : need 4.

→ G → 10 : only 2 bits

D → 110

More common letter

||
less binary bits needed to encode.

In general, if $P(\alpha) = \frac{n_\alpha}{N}$, then code length for α is $\lceil \log_2(1/P(\alpha)) \rceil$

(frequency of α)

($\sum n_\alpha = \text{length of plaintext}$)

Letter	A	B	C	D	E	F	G
Frequency	1	2	2	4	4	5	6
Code	0010	0001	000	110	111	01	10
Bits	4	4	3	3	3	2	2

File length after encrypted is

$$\sum_{\alpha} \text{code length}(\alpha) \cdot \text{frequency}(\alpha) = 64$$

Average number of bits per letter is

$$\frac{\text{ciphertext length}}{\text{plaintext length}} = \frac{64}{24} \approx 2.66$$

Compare to the entropy of the file

$$\begin{aligned} \sum_{\alpha} P(\alpha) \log_2 \left(\frac{1}{P(\alpha)} \right) &= \frac{1}{24} \log_2 \left(\frac{1}{1/24} \right) + \frac{2}{24} \log_2 \left(\frac{1}{2/24} \right) + \frac{2}{24} \log_2 \left(\frac{1}{2/24} \right) \\ &\quad + \frac{4}{24} \log_2 \left(\frac{1}{4/24} \right) + \frac{4}{24} \log_2 \left(\frac{1}{4/24} \right) \\ &\quad + \frac{5}{24} \log_2 \left(\frac{1}{5/24} \right) + \frac{6}{24} \log_2 \left(\frac{1}{6/24} \right) \\ &\approx 2.62165 \end{aligned}$$

Theorem 1. Suppose the letter counts in the plaintext are n_1, n_2, \dots, n_k and let $N = n_1 + \dots + n_k$. Then the best possible code length (in terms of bits per letter) is

p/t length $H = \sum_{i=1}^k p_i \log_2 \left(\frac{1}{p_i} \right),$

where $p_i = n_i/N$ for all $1 \leq i \leq k$.

prob. of each letter i in p/t.

Fact: If p_i 's and q_i 's are probability distributions

$$p_1 + p_2 + \dots + p_k = 1 = q_1 + q_2 + \dots + q_k.$$

then $\sum_{i=1}^k p_i \log_2 \left(\frac{1}{p_i} \right) \leq \sum_{i=1}^k p_i \log_2 \left(\frac{1}{q_i} \right)$ due to the fact that $\log_2(\cdot)$ is a concave function.

* Take any binary tree with leaf heights h_1, h_2, \dots, h_k .

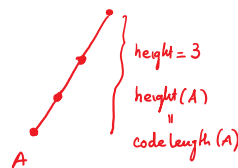
$$\text{File length after encryption} = \sum_{i=1}^k n_i h_i$$

$$= N \sum_{i=1}^k \frac{n_i}{N} \log_2 (2^{h_i})$$

$$= N \cdot \sum_{i=1}^k p_i \log_2 \left(\frac{1}{q_i} \right) \geq N \cdot \sum_{i=1}^k p_i \log_2 \left(\frac{1}{p_i} \right) = (\text{entropy of the p/t}) \cdot N.$$

$$\text{where } q_i = \frac{1}{2^{h_i}}$$

So $\frac{\text{file length after encryption}}{N} \geq \text{entropy of p/t}$



Theorem 2. The Huffman tree constructed from the probabilities p_1, p_2, \dots, p_k yields an expected code length that is within 1 bit of the entropy

$$H = \sum_{i=1}^k p_i \log_2 \left(\frac{1}{p_i} \right).$$

$$\text{Expected Codelength} = \sum_{i=1}^k p_i h_i = \sum_{i=1}^k p_i \lceil \log_2 \left(\frac{1}{p_i} \right) \rceil$$

So,

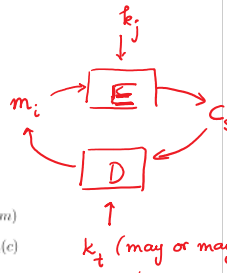
$$\underbrace{\sum_{i=1}^k p_i \log_2 \left(\frac{1}{p_i} \right)}_H \leq \sum_{i=1}^k p_i \lceil \log_2 \left(\frac{1}{p_i} \right) \rceil \leq \sum_{i=1}^k p_i \left(1 + \log_2 \left(\frac{1}{p_i} \right) \right)$$

$$H \leq \underbrace{\sum_{i=1}^k p_i}_1 + \underbrace{\sum_{i=1}^k p_i \log_2 \left(\frac{1}{p_i} \right)}_H$$

Random Crypto-systems

The set up of cryptography

- Message/plaintext space $\mathcal{M} = \{m_1, m_2, \dots, m_N\}$
- Key space $\mathcal{K} = \{k_1, k_2, \dots, k_S\}$
- Ciphertext space $\mathcal{C} = \{c_1, c_2, \dots, c_Q\} = \mathcal{C}$
- The encrypting function corresponding to the key k : $c = E_k(m)$
- The decrypting function corresponding to the key k : $m = D_k(c)$



In a random crypto-system:

- $M \in \mathcal{M}$ is the chosen message
- $K \in \mathcal{K}$ is the chosen key
- $C \in \mathcal{C}$ is the resulting ciphertext

random variables

$$P(M = m_i) = p_i$$

$$P(K = k_j) = q_j$$

$$C = E_{k_j}(m_i)$$

Remarks:

- We choose the key K independently of the message M .
- $C = E_K(M)$ so the ciphertext C is a random variable which depends on M and K .
- $M = D_K(C)$ so $H(K, C) = H(K, M)$
- $H(K|C)$ is the remaining uncertainty about the key after we intercept the ciphertext.

$$\begin{aligned} & \rightarrow P(M = m_i \cap K = k_j) \\ & = P(M = m_i) \cdot P(K = k_j) \\ & \text{for all } m_i \in \mathcal{M} \\ & \quad k_j \in \mathcal{K} \end{aligned}$$

$$H(K|C) = 0 \Leftrightarrow \text{the ciphertext determines the key : very bad!}$$

uncertainty about K, C
is the same as
the uncertainty about K, M

Choose a key by
spinning the wheel.

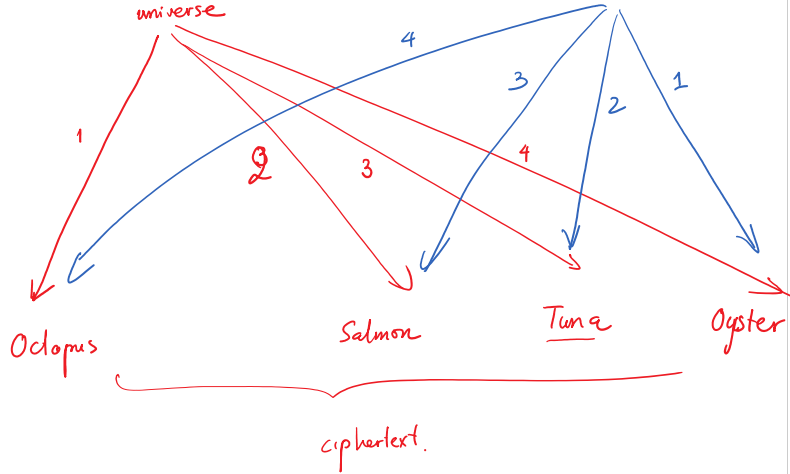


Example.

Thanos

destroy $\frac{1}{2}$
universe

chillax & watch
sun set : \mathcal{M}



ciphertext
'Tuna' — good : if key = k_2
— bad : if key = k_3

Theorem 3. For *ciphertext only* attack:

$$H(K|C) = H(K) + H(M) - H(C)$$

Theorem 4. For *known plaintext attack*:

$$H(K|C, M) = H(K) - H(C|M)$$

Definition 1. A cryptosystem is said to attain *perfect secrecy* if the ciphertext gives no information about the plaintext. That is, M, C are random variables, namely,

$$\mathbb{P}(M = m_i \mid C = c_j) = \mathbb{P}(M = m_i) \cdot \mathbb{P}(C = c_j)$$

for all $m_i \in \mathcal{M}$ and $c_j \in \mathcal{C}$.