# INFORMATION THEORY & CODING
## Week 7 : Source Coding 3 - Huffman Code

Dr. Rui Wang

Department of Electrical and Electronic Engineering
Southern Univ. of Science and Technology (SUSTech)

Email: wang.r@sustech.edu.cn

November 8, 2022

# Huffman Codes

## Problem 5.1

Given source symbols and their probabilities of occurence, how to design an optimal source code (prefix code and the shortest on average)?

**Huffman Codes**

- Merge the $D$ symbols with the smallest probabilities, and generate one new symbol whose probability is the summation of the $D$ smallest probabilities.

- Assign the $D$ corresponding symbols with digits $0, 1, \ldots, D-1$, then go back to Step 1.

Repeat the above process until $D$ probabilities are merged into probability 1.

Huffman code:

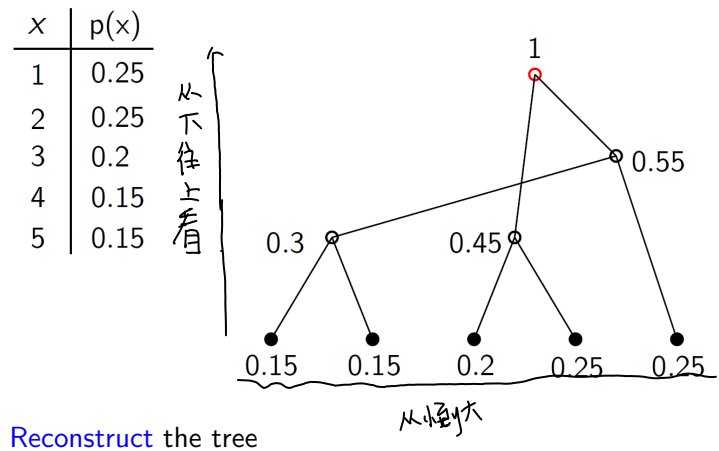Step I: make sure: $1+(D-1)k$ symbols

Step II: $D$ symbols with the least probabilities
$\Rightarrow$ one symbol with sum probabilities.

Step III: repeat step II until one symbol left.

Step IV: Assign codeword.

**Example 1**

| x | p(x) |
|---|------|
| 1 | 0.25 |
| 2 | 0.25 |
| 3 | 0.2  |
| 4 | 0.15 |
| 5 | 0.15 |

从下往上看

1

0.55

0.3

0.45

0.15  0.15  0.2  0.25  0.25

从短小

Reconstruct the tree

**Example 1**

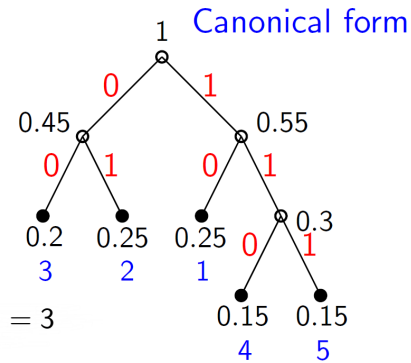| $x$ | $p(x)$ | $C(x)$ |
|-----|--------|--------|
| 1 | 0.25 | 10 |
| 2 | 0.25 | 01 |
| 3 | 0.2 | 00 |
| 4 | 0.15 | 110 |
| 5 | 0.15 | 111 |

**Validations:**

$\ell(1) = \ell(2) = \ell(3) = 2, \ell(4) = \ell(5) = 3$

$Ł = \sum \ell(x) p(x) = 2.3 \text{bits}$

$H_2(X) = - \sum p(x) \log_2 p(x) = 2.29 \text{bits}$

$L \geq H_2(X)$

Canonical form

不等的原因：分布不适配 $D^{-li}$

**Example 2**

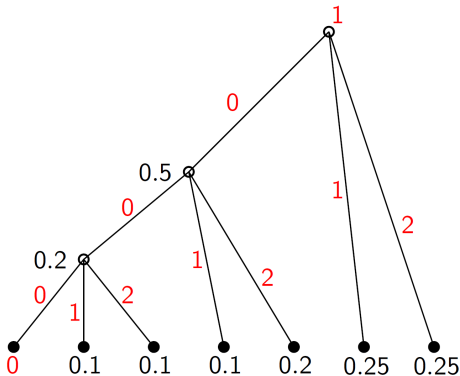| $x$ | $p(x)$ |
|---|---|
| 1 | 0.25 |
| 2 | 0.25 |
| 3 | 0.2 |
| 4 | 0.1 |
| 5 | 0.1 |
| 6 | 0.1 |
| Dummy | 0 |

$\mathcal{D} = \{0, 1, 2\}$

At one time, we merge $D$ symbols, and at each stage of the reduction, the number of symbols is reduced by $D - 1$. We want the total # of symbols to be $1 + k(D - 1)$. If not, we add dummy symbols with probability 0.

positive integer

D-ary

**Example 2** ($D \geq 3$)

| $x$ | p(x) | C(x) |
|---|---|---|
| 1 | 0.25 | 1 |
| 2 | 0.25 | 2 |
| 3 | 0.2 | 02 |
| 4 | 0.1 | 01 |
| 5 | 0.1 | 002 |
| 6 | 0.1 | 001 |
| Dummy | 0 | 000 |



**Validations:**

$L = \sum \ell(x)p(x) = 1.7$ ternary digits

$H_3(X) = -\sum p(x)\log_3 p(x) \approx 1.55$ ternary digits

南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

# Optimality of Huffman Codes

## Lemma 5.8.1

*For any distribution, the optimal prefix codes (with minimum expected length) should satisfy the following properties:*

1. *If $p_j > p_k$, then $\ell_j \leq \ell_k$.*
2. *The two longest codewords have the same length.*
3. *There exists an optimal prefix code, such that two of the longest codewords differ only in the last bit and correspond to the two least likely symbols.*

proof: Suppose $C$ with $p_j > p_k$ and $\ell_j > \ell_k$.

we should prove $C$ is not optimal.

$$L(C) = \sum_i p_i \ell_i = p_j \ell_j + p_k \ell_k + \sum_{i \neq j, k} p_i \ell_i$$

$C': j \to \ell_k; \ k \to \ell_j; \ C.$

$$L(C') = p_j \ell_k + p_k \ell_j + \sum_{i \neq j, k} p_i \ell_i$$

$$L(C) - L(C') = \underbrace{(p_j - p_k)}_{>0} \underbrace{(\ell_j - \ell_k)}_{>0} > 0$$

$L(C) > L(C') \Rightarrow C$ is not optimal.

proof: code $C$.

$C(j): b_1, b_2 \cdots b_m$

$C(k): b_1' b_2' \cdots b_m' b_{m+1}' \cdots b_{m+n}'$

$C'$

$C'(i) = \begin{cases} C(i), \, i \neq k \\ b_1' b_2' \cdots b_m', \, i = k. \end{cases}$ ⟶ $C'$ is prefix

⟶ $L(C') < L(C)$

⟶ 不为 $C'$ 中其它 codeword 的 prefix

and. $C'$ 中其它 code word 不会为它的 prefix

# Optimality of Huffman Codes

- 1. If $p_j > p_k$, then $\ell_j \leq \ell_k$.

### Proof.

Suppose that $C_m$ is an optimal code. Consider $C_m'$, with the codewords $j$ and $k$ of $C_m$ interchanged. Then

$$\underbrace{L\left(C_m'\right) - L\left(C_m\right)}_{\geq 0} = \sum p_i \ell_i' - \sum p_i \ell_i$$

$$= p_j \ell_k + p_k \ell_j - p_j \ell_j - p_k \ell_k$$

$$= \underbrace{\left(p_j - p_k\right)}_{>0} \left(\ell_k - \ell_j\right)$$

Thus, we must have $\ell_k \geq \ell_j$. $\qquad\square$

# Optimality of Huffman Codes

- 2. The two longest codewords have the same length.

## Optimality of Huffman Codes

- 3. There exists an optimal prefix code, such that two of the longest codewords differ only in the last bit and correspond to the two least likely symbols.

### Proof.

If there is a maximal-length codeword without a sibling, we can delete the last bit of the codeword and still preserve the prefix property. This reduces the average codeword length and contradicts the optimality of the code. Hence, every maximum-length codeword in any optimal code has a sibling. Now we can exchange the longest codewords s.t. the two lowest-probability source symbols are associated with two siblings on the tree, without changing the expected length. □

proof: Optimal prefix code $C$.

longest codeword. length $= m$.

$b_1 b_2 \cdots b_m \in C$.

$b_1' b_2' \cdots b_m' \in C$.

consider a new codeword. $b_1 b_3 \cdots b_{m-1} \overline{b_m}$

① $b_1 b_2 \cdots b_{m-1} \overline{b_m} \in C$. if $v \Rightarrow C$ is we wanted

② $b_1 b_2 \cdots b_{m-1} \overline{b_m} \notin C$

If ② is true, construct $C'$: replace $b_1' b_2' \cdots b_m'$ of $C$

by $b_1 b_2 \cdots b_{m-1} \overline{b_m}$.       $\Downarrow$

       optimal prefix code

$\Rightarrow$ canonical form.
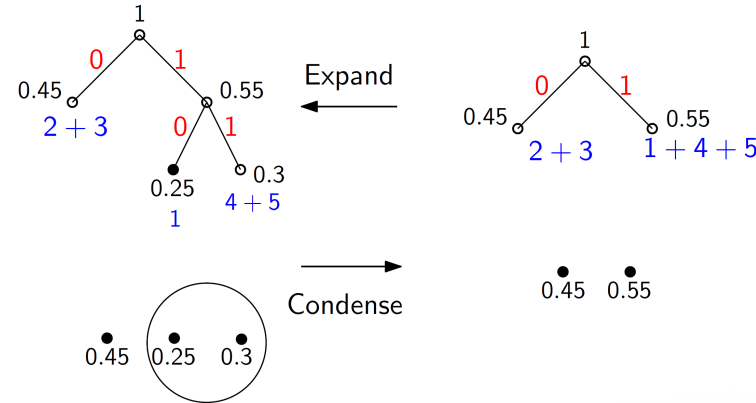
# Optimality of Huffman Codes

## Lemma 5.8.1

*For any distribution, the optimal prefix codes (with minimum expected length) should satisfy the following properties:*

1. *If $p_j > p_k$, then $\ell_j \leq \ell_k$.*
2. *The two longest codewords have the same length.*
3. *There exists an optimal prefix code, such that two of the longest codewords differ only in the last bit and correspond to the two least likely symbols.*

$\Rightarrow$ If $p_1 \geq p_2 \geq \cdots p_m$, then there exists an optimal code with $\ell_1 \leq \ell_2 \leq \cdots \ell_{m-1} = \ell_m$, and codewords $C(x_{m-1})$ and $C(x_m)$ differ only in the last bit. (canonical codes)

南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

- We prove the **optimality** of Huffman codes by **induction**. Assume binary code in the proof.



Expand ← / Condense

Right-hand handwritten notes:

$P_1, P_2 \cdots P_{m+1}$

$P_1 \geq P_2 \geq \cdots \geq P_{m+1}$

Huffman code for $(P_1 P_2 \cdots P_{m+1})$

$\Downarrow$ condense

Huffman code for $(P_1 P_2 \cdots P_{m-1}, P_m + P_{m+1})$

prefix code of $\underline{C\text{-form}}$ for $(P_1 P_2 \cdots P_{m+1})$

$\Downarrow$ condense.

prefix code for $(P_1, P_2 \cdots P_{m-1}, P_m + P_{m+1})$

W.L.O.G 2-ary code.

Step I: Huffman code is optimal for two-symbol dist
$|X| = 2$.

Step II: Suppose Huffman code is optimal for $|X| = m$
$m \geq 2$.

Step III: We should prove: Huffman code is optimal
for $|X| = m+1$, $P_1 \geq P_2 \geq \cdots \geq P_{m+1}$

$C_{m+1}$: Huffman code for $(P_1, P_2 \cdots P_m, P_{m+1})$

$C_m^*$: code condensed from $C_{m+1}$
$\Rightarrow$ Huffman code for $(P_1, P_2, \cdots, P_m + P_{m+1})$

据 Step II. it is optimal.

$C_{m+1}^*$: optimal prefix code of C-form for $(P_1, P_2, \cdots P_{m+1})$

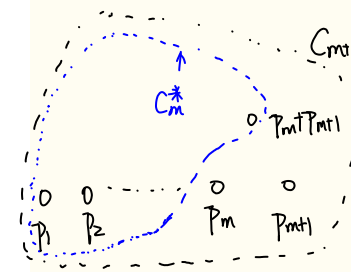$C_m$: condensed from $C_{m+1}^* \Rightarrow$ prefix code for $(P_1, P_2, \cdots P_m + P_{m+1})$

## Proof.

For $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ with $p_1 \geq p_2 \geq \cdots \geq p_m$, we define the Huffman reduction $\mathbf{p}' = (p_1, p_2, \ldots, p_{m-1+p_m})$ over an alphabet size of $m-1$. Let $C_{m-1}^*(\mathbf{P}')$ be an optimal Huffman code for $\mathbf{p}'$, and let $C_m^*(\mathbf{p})$ be the canonical optimal code for $\mathbf{p}$. □

**Key idea.**

expand $C_{m-1}^*$ to $C_m(\mathbf{p}) \Rightarrow L(C_m) = L(C_m^*)$

$$L(C_{m+1}) = \sum_{i=1}^{m+1} p_i l_i = \sum_{i=1}^{m+1} p_i l_i + p_m l_{m+1} + p_{m+1} l_{m+1}$$

$$l_m = l_{m+1}$$

$$L(C_m^*) = \sum_{i=1}^{m+1} p_i l_i + (p_m + p_{m+1})(l_{m+1} - 1)$$

$$L(C_{m+1}) - L(C_m^*) = p_m + p_{m+1}$$

$$L(C_{m+1}^*) - L(C_m) = p_m + p_{m+1}$$

$$\underbrace{L(C_{m+1}) - L(C_{m+1}^*)}_{\geq 0} = \underbrace{L(C_m^*) - L(C_m)}_{\leq 0}$$

$$\rightarrow L(C_{m+1}) = L(C_{m+1}^*) = L(C_m^*) = L(C_m)$$

# Optimality of Huffman Codes

### Proof.

For $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ with $p_1 \geq p_2 \geq \cdots \geq p_m$, we define the Huffman reduction $\mathbf{p}' = (p_1, p_2, \ldots, p_{m-1+p_m})$ over an alphabet size of $m-1$. Let $C_{m-1}^*(\mathbf{P}')$ be an optimal Huffman code for $\mathbf{p}'$, and let $C_m^*(\mathbf{p})$ be the canonical optimal code for $\mathbf{p}$. $\qquad\square$

| | $C_{m-1}^*(\mathbf{p}')$ | | $C_m(\mathbf{p})$ | |
|---|---|---|---|---|
| $p_1$ | $w_1'$ | $l_1'$ | $w_1 = w_1'$ | $l_1 = l_1'$ |
| $p_2$ | $w_2'$ | $l_2'$ | $w_2 = w_2'$ | $l_2 = l_2'$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |
| $p_{m-2}$ | $w_{m-2}'$ | $l_{m-2}'$ | $w_{m-2} = w_{m-2}'$ | $l_{m-2} = l_{m-2}'$ |
| $p_{m-1} + p_m$ | $w_{m-1}'$ | $l_{m-1}'$ | $w_{m-1} = w_{m-1}'0$ | $l_{m-1} = l_{m-1}' + 1$ |
| | | | $w_m = w_{m-1}'1$ | $l_m = l_{m-1}' + 1$ |

南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

# Optimality of Huffman Codes

## Proof.

For $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ with $p_1 \geq p_2 \geq \cdots \geq p_m$, we define the Huffman reduction $\mathbf{p}' = (p_1, p_2, \ldots, p_{m-1+p_m})$ over an alphabet size of $m - 1$. Let $C_{m-1}^*(\mathbf{P}')$ be an optimal Huffman code for $\mathbf{p}'$, and let $C_m^*(\mathbf{p})$ be the canonical optimal code for $\mathbf{p}$. $\square$

| $C_{m-1}(\mathbf{p}')$ | | | $C_m^*(\mathbf{p})$ | |
|---|---|---|---|---|
| $p_1$ | $w_1'$ | $l_1'$ | $w_1 = w_1'$ | $l_1 = l_1'$ |
| $p_2$ | $w_2'$ | $l_2'$ | $w_2 = w_2'$ | $l_2 = l_2'$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $p_{m-2}$ | $w_{m-2}'$ | $l_{m-2}'$ | $w_{m-2} = w_{m-2}'$ | $l_{m-2} = l_{m-2}'$ |
| $p_{m-1} + p_m$ | $w_{m-1}'$ | $l_{m-1}'$ | $w_{m-1} = w_{m-1}'0$ | $l_{m-1} = l_{m-1}' + 1$ |
| | | | $w_m = w_{m-1}'1$ | $l_m = l_{m-1}' + 1$ |

# Optimality of Huffman Codes

## Proof.

For $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ with $p_1 \geq p_2 \geq \cdots \geq p_m$, we define the Huffman reduction $\mathbf{p}' = (p_1, p_2, \ldots, p_{m-1+p_m})$ over an alphabet size of $m-1$. Let $C_{m-1}^*(\mathbf{P}')$ be an optimal Huffman code for $\mathbf{p}'$, and let $C_m^*(\mathbf{p})$ be the canonical optimal code for $\mathbf{p}$. □

expand $C_{m-1}^*(\mathbf{p}')$ to $C_m(\mathbf{p})$

$$L(\mathbf{p}) = L^*(\mathbf{p}') + p_{m-1} + p_m$$

condense $C_m^*(\mathbf{p})$ to $C_{m-1}(\mathbf{p}')$

$$L^*(\mathbf{p}) = L(\mathbf{p}') + p_{m-1} + p_m$$

# Optimality of Huffman Codes

## Proof.

For $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ with $p_1 \geq p_2 \geq \cdots \geq p_m$, we define the Huffman reduction $\mathbf{p}' = (p_1, p_2, \ldots, p_{m-1+p_m})$ over an alphabet size of $m - 1$. Let $C_{m-1}^*(\mathbf{P}')$ be an optimal Huffman code for $\mathbf{p}'$, and let $C_m^*(\mathbf{p})$ be the canonical optimal code for $\mathbf{p}$. □

$$L(\mathbf{p}) = L^*\left(\mathbf{p}'\right) + p_{m-1} + p_m$$
$$L^*(\mathbf{p}) = L\left(\mathbf{p}'\right) + p_{m-1} + p_m$$

$$(\underbrace{L\left(\mathbf{p}'\right) - L^*\left(\mathbf{p}'\right)}_{\geq 0}) + (\underbrace{L(\mathbf{p}) - L^*(\mathbf{p})}_{\geq 0}) = 0$$

# Optimality of Huffman Codes

## Proof.

For $\mathbf{p} = (p_1, p_2, \ldots, p_m)$ with $p_1 \geq p_2 \geq \cdots \geq p_m$, we define the Huffman reduction $\mathbf{p}' = (p_1, p_2, \ldots, p_{m-1+p_m})$ over an alphabet size of $m-1$. Let $C^*_{m-1}(\mathbf{P}')$ be an optimal Huffman code for $\mathbf{p}'$, and let $C^*_m(\mathbf{p})$ be the canonical optimal code for $\mathbf{p}$. $\qquad\square$

Thus, $L(\mathbf{p}) = L^*(\mathbf{p})$. Minimizing the expected length $L(C_m)$ is equivalent to minimizing $L(C_{m-1})$. The problem is reduced to one with $m-1$ symbols and probability masses $(p_1, p_2, \ldots, p_{m-1} + p_m)$. Proceeding this way, we finally reduce the problem to two symbols, in which case the optimal code is obvious.

# Optimality of Huffman Codes

## Theorem 5.8.1

*Huffman coding is* *optimal*, *that is, if $C^*$ is a Huffman code and $C'$ is any other uniquely decodable code,* $L(C^*) \leq L(C')$.

## Remark

Huffman coding is a *greedy algorithm* in which it merges the two least likely symbols at each step.

LOCAL OPT $\rightarrow$ GLOBAL OPT

# Textbook

Related Sections : 5.6 - 5.8