

# Review of source coding methods

---

## Source coding

### **Quantization (a lossy data compression)**

Analogue source is quantized into finite number of levels.

Introduces distortion - Data is lost and cannot be recovered.

Examples

- Pulse code modulation PCM

- Differential pulse code modulation DPCM

- Delta modulation DM

- Uniform quantization

...

The fundamental limit on the performance is given by the rate-distortion bound.

### **Noiseless coding (lossless data compression)**

Digital data (the results of quantization) are compressed. Compression attempts to represent data with as few bits as possible.

Examples

- Huffmann coding

- Lempel-Ziv coding

- Arithmetic coding

---

## Waveform sources

- ❑ Voice coding
  - Waveform coders
    - PCM, ADPCM, DM
  - Vocoders
- ❑ Music coding
  - mp3, ogg
- ❑ Video coding
  - mpeg2, mpeg4
- ❑ Still picture coding
  - jpeg

## Discrete sources

- ❑ Digital data compression
    - Constant length code words
    - Variable length code words
      - Huffman encoding
-

## Source coding theorem

- Foundation of source coding laid by Shannon
  - Tells the minimum source encoded word length from a discrete information source
  - A continuous information source must first be sampled, the sampling frequency should be equal to or larger than two times the signal bandwidth
  - This usually means that the samples are correlated, the information contains redundancy
  - Transmitting the redundancy part of the information means wasting of primary communication resources, transmit power and channel bandwidth
  - The task of source encoding is to remove the redundancy, if necessary it can be added in the receiver
-

## Self entropy

$$H(X) = - \sum_{i=1}^n P(x_i) \log P(x_i)$$

Entropy of a discrete source is maximum when the output symbols are equally probable

## Conditional entropy

$$H(X|Y) = - \sum_{j=1}^m \sum_{i=1}^n P(x_i, y_j) \log P(x_i | y_j)$$

In a communication channel the conditional entropy is called also equivocation



## Self information, conditional self information, mutual self information

$$l(x_i) = -\log P(x_i)$$

$$l(x_i | y_j) = -\log P(x_i | y_j)$$

$$l(x_i, y_j) = -\log P(x_i, y_j)$$

## Mutual information

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

## Differential entropy

$$h(X) = -\int_{-\infty}^{\infty} p(x_i) \log p(x_i) dx$$

## Entropy of a normally distributed variable

$$h(N(0, \sigma^2)) = \frac{1}{2} \log(2\pi e \sigma^2)$$

---

## Example of the source encoding theorem

- For a discrete memoryless information source (subsequent symbols are statistically independent) the average source encoded code word length must fulfil the condition

$$\mathbf{E}\{L_{cwd}\} \geq H_s \quad (3 - 1)$$

$L_{cwd}$  is the length of a source coded code-word

$H_s$  is the entropy of the source defined as

$$H_s = \sum_{k=1}^K p_k \log_2 \left( \frac{1}{p_k} \right) \quad (3 - 2)$$

$p_k$  is the occurrence probability of each of the  $K$  source symbols

---

For an information source with memory (subsequent symbols are not statistically independent) the source entropy is defined as

$$H_s = \lim_{k \rightarrow \infty} \left( \frac{1}{k} \sum_{i=1}^k H_s (s_i | s_1, s_2, \dots, s_{i-1}) \right) \quad (3 - 3)$$

$$H_s (s_i | s_1, s_2, \dots, s_{i-1}) = \sum_{s_1} \dots \sum_{s_i} P \{s_1, \dots, s_i\} \log_2 \left( \frac{1}{P \{s_i | s_1, \dots, s_{i-1}\}} \right) \quad (3 - 4)$$

For a continuous signal source the average entropy is infinite because an infinite codeword length is needed to describe a signal amplitude sample. Instead a differential entropy is defined as

$$H = \int p(x) \log_2 \left( \frac{1}{p(x)} \right) dx \quad (3 - 5)$$

The efficiency of the source encoder is defined as

$$\eta = \frac{H_s}{\mathbf{E}\{L_{c wd}\}} \quad (3 - 6)$$

---

### Example 3.1

A discrete source generates the symbols 0,1,2,3,4,5,6,7,8, and 9 with equal probability. The source entropy is

$$H_s = \sum_{k=1}^K p_k \log_2 \left( \frac{1}{p_k} \right) = 10 \cdot 0.1 \cdot \log_2 \left( \frac{1}{0.1} \right) = 3.32 \text{ bit/symbol}$$

If 4 bits/symbol are used in the transmission the source encoder efficiency is

$$\eta = \frac{H_s}{\mathbf{E}\{L_{c wd}\}} = \frac{3.32}{4} = 0.83$$

## Compression

Fundamental limit of the compression is given by the entropy of the source

If all the  $L$  source symbols are equally probable the average code length is

$$H(X) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i) \leq \log_2 L$$

The average code rate is

$R = \log_2 L$  if  $L$  is power of 2 or

$$R = \lfloor \log_2 L \rfloor + 1$$

We can improve (reduce) the code length by combining multiple symbols

$$\lfloor J \log_2 L \rfloor + 1 \Rightarrow R = \lfloor \log_2 L \rfloor + \frac{1}{J}$$

In case the symbols are not equally probable more likely symbols can be allocated shorter code words

Prefix code a code where end of each codeword can explicitly identified - the code is uniquely decodable and uniquely decodable

Letter	$P(a_i)$	Code I	Code II	Code III
$a_1$	1/2	1	0	0
$a_2$	1/4	00	10	01
$a_3$	1/8	01	110	011
$a_4$	1/8	10	111	111

## Huffmann coding

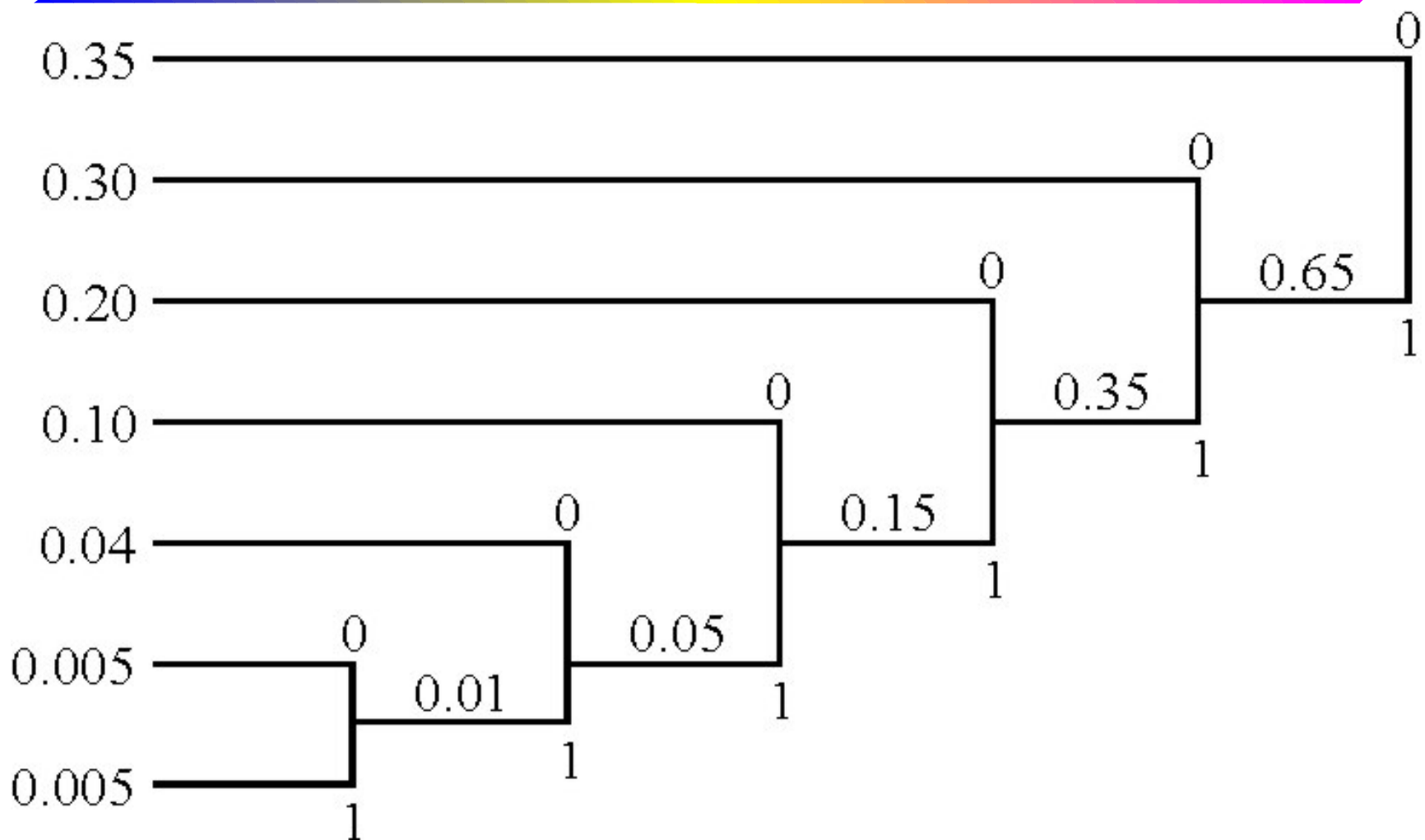
The algorithm minimizes the average number of digits required for describing the source.

Order the probabilities

$$p(x_1) \geq p(x_2) \geq \dots \geq p(x_L)$$

Encoding algorithm

1. Select two least probable symbols and tie them together
  2. Assign to the upper branch of the tied symbol 0 and to the lower branch 1
  3. Replace the two lowest probability symbols with the new combined symbol.  
The combined symbol probability is the sum of the probabilities of the initial symbols
  4. If there are any symbol left go back to step 1
-



**Example 1 Huffman coding**

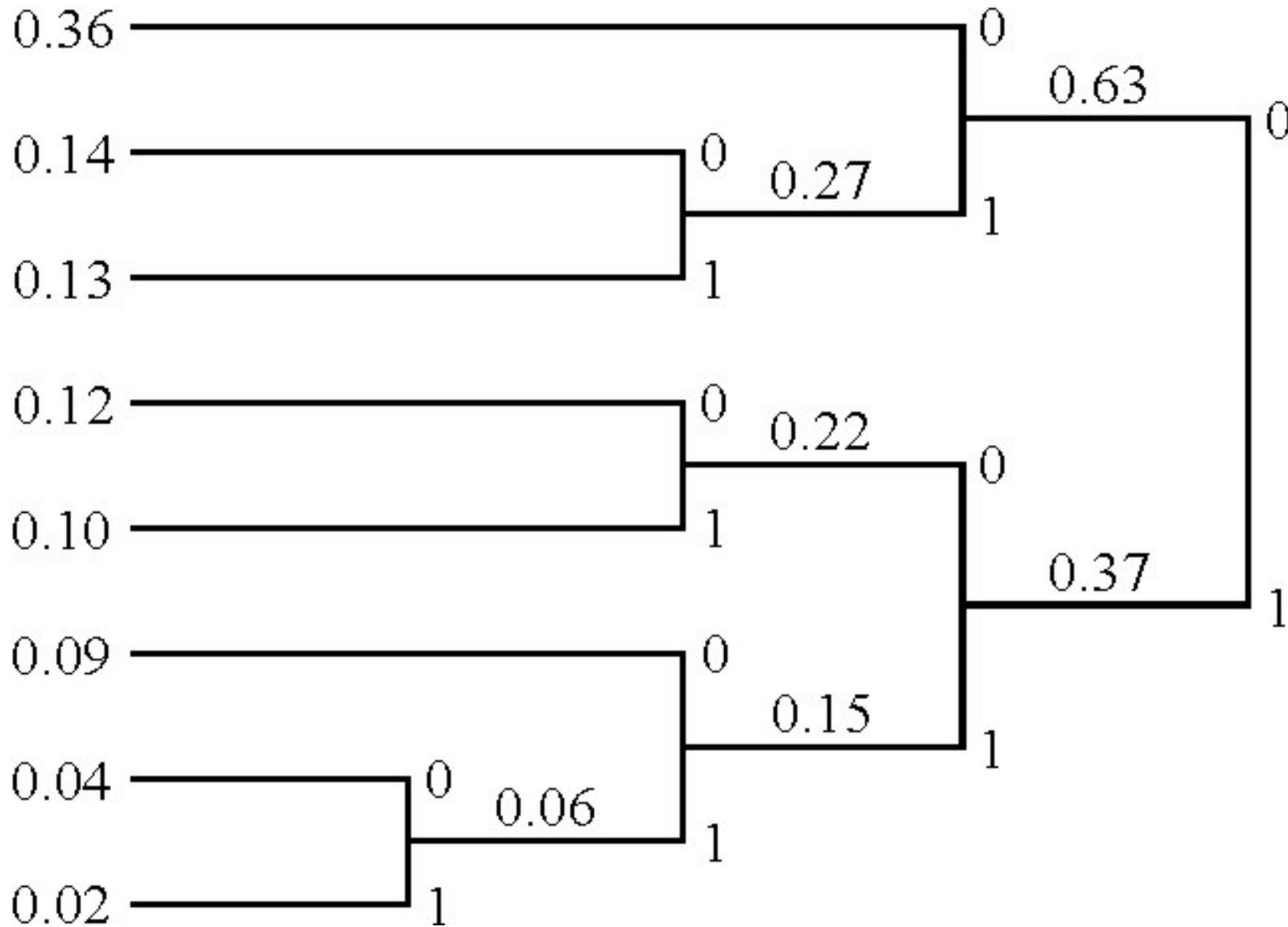


Input symbols	Pb	Self information $-\log_2(P(a_i))$	Assigned bit sequence	Length of codeword $l_i$
$a_1$	<b>0.35</b>	<b>1.51</b>	<b>0</b>	<b>1</b>
$a_2$	<b>0.30</b>	<b>1.74</b>	<b>01</b>	<b>2</b>
$a_3$	<b>0.20</b>	<b>2.32</b>	<b>011</b>	<b>3</b>
$a_4$	<b>0.10</b>	<b>3.32</b>	<b>0111</b>	<b>4</b>
$a_5$	<b>0.04</b>	<b>4.64</b>	<b>01111</b>	<b>5</b>
$a_6$	<b>0.005</b>	<b>7.64</b>	<b>011111</b>	<b>6</b>
$a_7$	<b>0.005</b>	<b>7.64</b>	<b>111111</b>	<b>6</b>

**Entropy**  $-\sum_{i=1}^7 P(a_i) \log_2(P(a_i)) = 2.11$

**Average codeword length**  $-\sum_{i=1}^7 P(a_i) l_i = 2.21$

---



**Example 2 Huffman coding**

---

Input symbols	Pb	Self information $-\log_2(P(a_i))$	Assigned bit sequence	Length of codeword $l_i$
$a_1$	<b>0.36</b>	<b>1.47</b>	<b>00</b>	<b>2</b>
$a_2$	<b>0.14</b>	<b>2.83</b>	<b>010</b>	<b>3</b>
$a_3$	<b>0.13</b>	<b>2.94</b>	<b>011</b>	<b>3</b>
$a_4$	<b>0.12</b>	<b>3.04</b>	<b>100</b>	<b>3</b>
$a_5$	<b>0.1</b>	<b>3.32</b>	<b>101</b>	<b>3</b>
$a_6$	<b>0.09</b>	<b>3.47</b>	<b>110</b>	<b>3</b>
$a_7$	<b>0.04</b>	<b>4.64</b>	<b>1110</b>	<b>4</b>
$a_8$	<b>0.02</b>	<b>5.64</b>	<b>1111</b>	<b>4</b>

**Entropy**  $-\sum_{i=1}^7 P(a_i) \log_2(P(a_i)) = 2.62$

**Average codeword length**  $-\sum_{i=1}^7 P(a_i) l_i = 2.7$

## Lempel Ziv coding

**Adaptive code - From the input data is created a codebook by parsing through the data vector.**

**Each element in the codebook is addressed as the address of the prefix of the bit sequence observed plus one new bit**

### Example

**Input sequence**

**10101101001001110101000011001110101100011011**

**extracted phrases**

**1 0 10 11 01 00 100 111 010 1000 011 001 110 101 10001 1011**



	<b>Dict. location</b>	<b>Dict. content</b>	<b>Codeword</b>
<b>1</b>	<b>0001</b>	<b>1</b>	<b>00001</b>
<b>2</b>	<b>0010</b>	<b>0</b>	<b>00000</b>
<b>3</b>	<b>0011</b>	<b>10</b>	<b>00010</b>
<b>4</b>	<b>0100</b>	<b>11</b>	<b>00011</b>
<b>5</b>	<b>0101</b>	<b>01</b>	<b>00101</b>
<b>6</b>	<b>0110</b>	<b>00</b>	<b>00100</b>
<b>7</b>	<b>0111</b>	<b>100</b>	<b>00110</b>
<b>8</b>	<b>1000</b>	<b>111</b>	<b>01001</b>
<b>9</b>	<b>1001</b>	<b>010</b>	<b>01010</b>
<b>10</b>	<b>1010</b>	<b>1000</b>	<b>01110</b>
<b>11</b>	<b>1011</b>	<b>011</b>	<b>01011</b>
<b>12</b>	<b>1100</b>	<b>001</b>	<b>01101</b>
<b>13</b>	<b>1101</b>	<b>110</b>	<b>01000</b>
<b>14</b>	<b>1110</b>	<b>101</b>	<b>00111</b>
<b>15</b>	<b>1111</b>	<b>10001</b>	<b>10101</b>
<b>16</b>		<b>1011</b>	<b>11101</b>

---

## Rate distortion theory

In practical source coding there are constraints that make the process imperfect

This causes distortion which should be kept on a tolerable level

Therefore we talk about source encoding with a fidelity criterion

Rate distortion theory deals with this situation

With a limited code word length there is a finite set of code words and e.g. analogue information cannot be exactly represented → lossy data compression

---