

# Group Seminar

---

## Arithmetic Coding

WANG Yilei  
Sep 14 (2018)

# Outline

- ❑ Background
  - Source Coding
  - Huffman Coding
  - Limitation of Symbol Code
  - Beyond Symbol Code
- ❑ Arithmetic Coding
  - Introduction
  - Representation of Interval
  - Upper Bound
  - Integer Implementation
  - Rescaling while encoding
- ❑ Adaptive Arithmetic Coding
- ❑ Applications
  - Constant Weight Binary Code
  - Constant Weight Integer Code

# Background

## □ Source Coding

Source alphabet  $S = \{s_1, \dots, s_m\}$  and code alphabet  $A = \{a_1, \dots, a_n\}$ .

*We would like to transmit text written in the source alphabet, but our channel accepts only code alphabetic characters.*

*Common case in computer science:  $A = \{0, 1\}$ . We will always assume  $A = \{0, 1\}$  in the remaining discussion.*

# Background

## □ Source Coding

Source alphabet  $S = \{s_1, \dots, s_m\}$  and code alphabet  $A = \{a_1, \dots, a_n\}$ .

Encoding function: A one-to-one function  $\phi: S^+ \rightarrow A^+$ .

Encoding scheme: A list of productions  $\{s_i \rightarrow w_i\} (w_i \in A^+, i = 1, \dots, m)$ .

Each encoding scheme corresponds to an encoding function by concatenation.

*Example:  $S = \{a, b, c\}$ ,  $A = \{0, 1\}$ , and the scheme is*

$$a \rightarrow 0$$

$$b \rightarrow 10$$

$$c \rightarrow 11,$$

*then  $\phi(acbb) = 0111010$ .*

Symbol code: Source code with encoding function constructed by encoding scheme.

# Background

## □ Source Coding

Source alphabet  $S = \{s_1, \dots, s_m\}$  and code alphabet  $A = \{a_1, \dots, a_n\}$ .

Encoding function: A one-to-one function  $\phi: S^+ \rightarrow A^+$ .

Encoding scheme: A list of productions  $\{s_i \rightarrow w_i\} (w_i \in A^+, i = 1, \dots, m)$ .

Each encoding scheme corresponds to an encoding function by concatenation.

Symbol code: Source code with encoding function constructed by encoding scheme.

Prefix code: An encoding scheme that there are no  $i \neq j$  such that  $w_i$  is an initial segment, or prefix (from left to right) of  $w_j$ .

*Prefix code:*

$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 11$

*Non-prefix code:*

$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 01$

*Prefix code is quite easy to be decoded.*

# Background

## □ Source Coding

Source alphabet  $S = \{s_1, \dots, s_m\}$  and code alphabet  $A = \{a_1, \dots, a_n\}$ .

Encoding function: A one-to-one function  $\phi: S^+ \rightarrow A^+$ .

Encoding scheme: A list of productions  $\{s_i \rightarrow w_i\} (w_i \in A^+, i = 1, \dots, m)$ .

Each encoding scheme corresponds to an encoding function by concatenation.

Symbol code: Source code with encoding function constructed by encoding scheme.

Prefix code: An encoding scheme that there are no  $i \neq j$  such that  $w_i$  is an initial segment, or prefix (from left to right) of  $w_j$ .

## □ Huffman Coding

Suppose each source alphabet  $X$  is generated i.i.d. with  $\Pr(X = s_i) = f_i$ .

Average code word length:  $\bar{\ell} = \sum_{i=1}^m f_i \cdot |w_i|$ , where  $|w_i|$  is the length of  $w_i$ .

Entropy:  $h = H(X) = -\sum_{i=1}^m f_i \log f_i$ . Shannon's source coding theorem:  $h \leq \bar{\ell}$ .

Huffman coding: Optimal symbol code (also prefix code) with  $h \leq \bar{\ell} < h + 1$ .

# Background

## □ Limitation of Symbol Code

Consider a random source string with length  $k$ :  $\mathbf{X} = X_1 X_2 \cdots X_k$ .

The entropy of string  $\mathbf{X}$ :  $H(\mathbf{X}) = \sum_{i=1}^k H(X_i) = kh$ .

The average length of code string using Huffman coding:

$$E[|\phi(\mathbf{X})|] = E[|\phi(X_1) \cdots \phi(X_k)|] = E\left(\sum_{i=1}^k |\phi(X_i)|\right) = \sum_{i=1}^k E[\phi(X_i)] = k\bar{\ell}$$

The number of overhead bits is  $k(\bar{\ell} - h)$ , which is linear to the length of source string when  $\bar{\ell} \neq h$ .

*Example:  $S = \{a, b, c\}$  with frequencies  $\{0.7, 0.2, 0.1\}$ ,  $A = \{0,1\}$ .*

*Then  $h = -(0.7 \log 0.7 + 0.2 \log 0.2 + 0.1 \log 0.1) \approx 1.16$ .*

*The encoding scheme given by Huffman coding is  $\{a \rightarrow 0, b \rightarrow 10, c \rightarrow 11\}$ .*

*Then  $\bar{\ell} = 0.7 \times 1 + 0.2 \times 2 + 0.1 \times 2 = 1.40$ .*

*The number of overhead bits is about 0.24 per source alphabet.*

# Background

## □ Limitation of Symbol Code

Consider a random source string with length  $k$ :  $\mathbf{X} = X_1X_2 \cdots X_k$ .

The entropy of string  $\mathbf{X}$ :  $H(\mathbf{X}) = \sum_{i=1}^k H(X_i) = kh$ .

The average length of code string using Huffman coding:

$$E[|\phi(\mathbf{X})|] = E[|\phi(X_1) \cdots \phi(X_k)|] = E\left(\sum_{i=1}^k |\phi(X_i)|\right) = \sum_{i=1}^k E[|\phi(X_i)|] = k\bar{\ell}$$

The number of overhead bits is  $k(h - \bar{\ell})$ , which is linear to the length of source string when  $\bar{\ell} \neq h$ .

## □ Beyond Symbol Code

To reduce the number of overhead bits, we remove the restriction that encoding function  $\phi$  is always constructed by concatenation, i.e. we don't ensure  $\phi(xy) = \phi(x)\phi(y)$  for any  $x \in S^+, y \in S^+$  now.

Arithmetic coding achieves totally less than 1 bit overhead for any finite source string.



# Arithmetic Coding

## □ Introduction

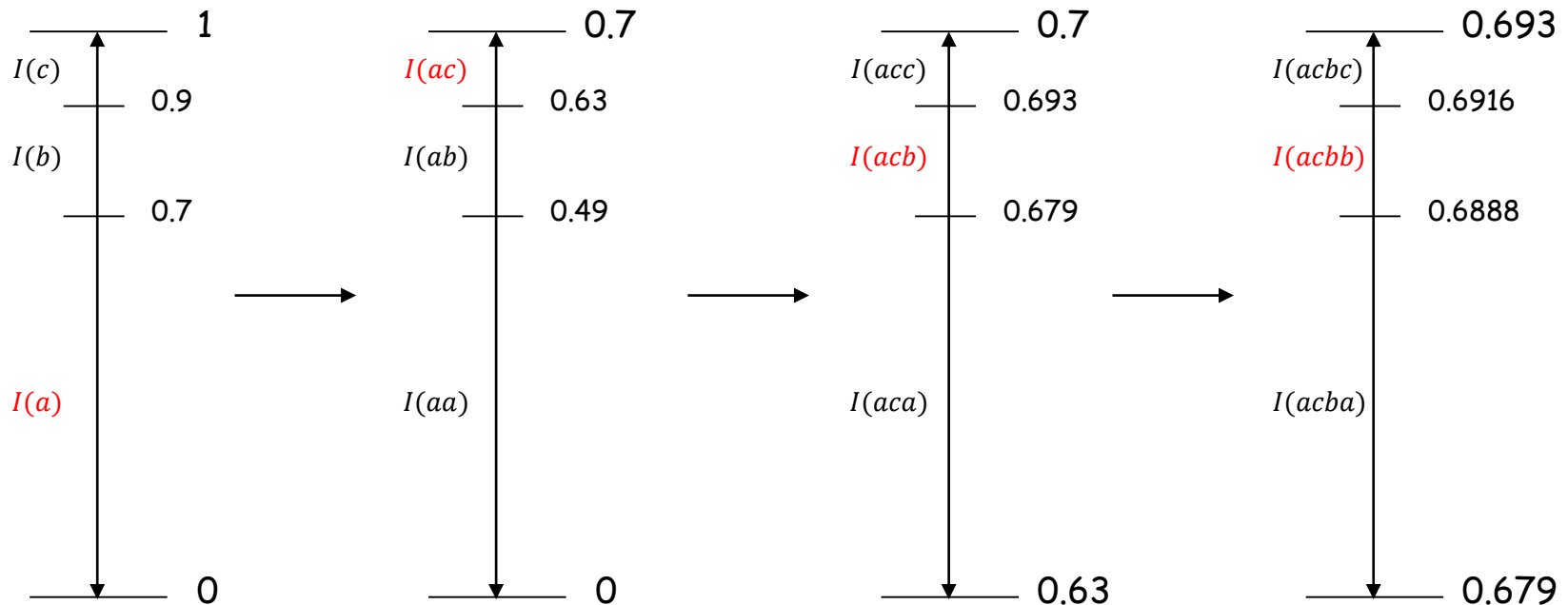
Each source string  $s \in S^*$  corresponds to an interval  $I(s) = [\alpha(s), \beta(s))$ .

Initial case: For empty string  $\epsilon$ ,  $I(\epsilon) = [0, 1)$ .

Recursion:

$$I(ss_i) = [\alpha(s) + \sum_{j=1}^{i-1} f_j (\beta(s) - \alpha(s)), \alpha(s) + \sum_{j=1}^i f_j (\beta(s) - \alpha(s))]$$

Example:  $S = \{a, b, c\}$  with frequencies  $\{0.7, 0.2, 0.1\}$ ,  $A = \{0, 1\}$ .  $I(acbb) = ?$



# Arithmetic Coding

## □ Introduction

Each source string  $s \in S^*$  corresponds to an interval  $I(s) = [\alpha(s), \beta(s))$ .

Initial case: For empty string  $\epsilon$ ,  $I(\epsilon) = [0, 1)$ .

Recursion:

$$I(ss_i) = [\alpha(s) + \sum_{j=1}^{i-1} f_j (\beta(s) - \alpha(s)), \alpha(s) + \sum_{j=1}^i f_j (\beta(s) - \alpha(s))]$$

Observation: Suppose  $X$  is a random source string with length  $k$ , and  $x$  is a fixed source string with length  $k$ . Then

- $\Pr(X = x) = \beta(x) - \alpha(x)$ ,
- $\alpha(ss_i) \geq \alpha(s)$  and  $\beta(ss_i) \leq \beta(s)$ .

## □ Representation of Interval

Let  $t$  be the integer such that  $2^{-t} \leq \beta - \alpha < 2^{-t+1}$ , i.e.  $t = \lceil -\log(\beta - \alpha) \rceil$ .

Let  $c = \lceil 2^t \alpha \rceil$ , then  $\alpha \leq c/2^t < \beta$ . We claim that the binary representation of  $c$  is the code string we need.

*Example:*  $\alpha = 0.6888, \beta = 0.6916, \beta - \alpha = 0.0028, t = 9, c = 353$ ,  
 $c/2^t = (0.101100001)_2$ , and the code string is 101100001.

# Arithmetic Coding

## □ Upper Bound

The average code length of length  $k$  random source string  $X$  is

$$\begin{aligned} E[|\phi(X)|] &= \sum_{s \in S^k} \Pr(X = s) t(s) \\ &= \sum_{s \in S^k} \Pr(X = s) [-\log(\beta(s) - \alpha(s))] \\ &< - \sum_{s \in S^k} \Pr(X = s) \log(\beta(s) - \alpha(s)) + \sum_{s \in S^k} \Pr(X = s) \\ &< - \sum_{s \in S^k} \Pr(X = s) \log \Pr(X = s) + 1 = H(X) + 1. \end{aligned}$$

## □ Integer Implementation

Arithmetic coding requires  $O(k)$  float number operations, and  $\beta(s) - \alpha(s)$  becomes smaller while the length of  $s$  grows, which needs unlimited registers, or the error will accumulate step by step.

There is an integer implementation of arithmetic coding, which needs only  $O(k)$  integer operations. With  $w = O(\log(k/(f_{\min} \epsilon)))$  registers, it achieves

$$E[|\phi(X)|] < H(X) + 1 + \epsilon.$$

# Arithmetic Coding

## □ Rescaling while encoding

It is not necessary to always wait until the entire source string is scanned and processed before you have a code for it. Instead, we are able to transmit bits while encoding.

Note that  $\alpha$  is non-decreasing while  $\beta$  is non-increasing while encoding. We keep a underflow count  $r$  with initial value 0. Each time  $\alpha$  or  $\beta$  is updated, we check whether

- $\alpha \geq 1/2$ . If so, we apply transformation  $x \rightarrow 2x - 1$  to  $\alpha$  and  $\beta$ , and transmit  $10^r$  (1 followed by  $r$  0's) to the decoder. Set  $r = 0$ .
- $\beta < 1/2$ . If so, we apply transformation  $x \rightarrow 2x$  to  $\alpha$  and  $\beta$ , and transmit  $01^r$  (0 followed by  $r$  1's) to the decoder. Set  $r = 0$ .
- $1/4 \leq \alpha < 1/2 \leq \beta < 3/4$ . If so, we apply transformation  $x \rightarrow 2x - 1/2$  to  $\alpha$  and  $\beta$ , and set  $r = r + 1$ .

*Example:  $S = \{a, b, c\}$  with frequencies  $\{0.7, 0.2, 0.1\}$ ,  $A = \{0, 1\}$ .  $I(acbb) = ?$*

# Adaptive Arithmetic Coding

Arithmetic coding is able to deal with the model without assumption that every source alphabet is generated independent and identically.

Recursion:

$$I(ss_i) = [\alpha(\mathbf{s}) + \sum_{j=1}^{i-1} f_j (\beta(\mathbf{s}) - \alpha(\mathbf{s})), \alpha(\mathbf{s}) + \sum_{j=1}^i f_j (\beta(\mathbf{s}) - \alpha(\mathbf{s}))]$$

Let  $\mathbf{s} = x_1 x_2 \dots x_k$ , and define the upper cumulative probabilities as follows:

$$Q_k(i | x_1, \dots, x_k) = \sum_{j=1}^i \Pr(X_{k+1} = s_j | X_1 = x_1, \dots, X_k = x_k)$$

where  $Q_k(0 | x_1, \dots, x_k) = 0$ . Then the recursion becomes

$$I(ss_i) = [\alpha(\mathbf{s}) + Q_k(i-1 | x_1, \dots, x_k)(\beta(\mathbf{s}) - \alpha(\mathbf{s})), \alpha(\mathbf{s}) + Q_k(i | x_1, \dots, x_k)(\beta(\mathbf{s}) - \alpha(\mathbf{s}))]$$

Note that both the analysis of upper bound and the integer implementation still work in adaptive arithmetic coding.

# Applications

## □ Constant Weight Binary Code

Let  $S = \{0, 1\}$ . Define constant weight binary code with length  $n$  and weight  $w$  as  $B(n, w) = \{s = s_1 s_2 \cdots s_n \in S^n \mid s_1 + s_2 + \cdots + s_n = w\}$ .

We want to encode  $B$  with least average code length.

Let  $p(s) = \Pr(\mathbf{X} = s) = \Pr(X_1 = s_1, X_2 = s_2, \dots, X_n = s_n)$ , then  $\sum_{s \in B} p(s) = 1$ .

Therefore,  $H(\mathbf{X}) = -\sum_{s \in B} p(s) \log p(s) \leq \log \binom{n}{w}$ , and two sides are equal if and only if  $p(s) = \binom{n}{w}^{-1}$  for all  $s \in B$ . If the distribution of  $B$  is unknown, we assume  $p(s) = \binom{n}{w}^{-1}$  so as to achieve worst case optimal.

*Example:*  $n = 5, w = 3$ .  $\phi(10011) = ?$

It is easy to see that for any  $s \in B$ ,  $|\phi(s)| = \lceil \log \binom{n}{w} \rceil$ . Therefore,  $E[|\phi(\mathbf{X})|] = \lceil H(\mathbf{X}) \rceil$ , which is optimal.

# Application

## □ Constant Weight Integer Code

Let  $S = \{0, 1, \dots, N\}$ , where  $N$  is a large positive integer. Define constant weight integer code with length  $n$  and weight  $w (\leq N)$  as  $C(n, w) = \{s = s_1 s_2 \dots s_n \in S^n \mid s_1 + s_2 + \dots + s_n = w\}$ .

We want to encode  $C$  with least average code length.

There is a bijection  $g: C(n, w) \rightarrow B(n + w - 1, w)$ . Let  $\phi'$  be the optimal encoding function from  $B(n + w - 1, w)$  to  $\{0, 1\}^+$ , then  $\phi = g \circ \phi'$  is an optimal encoding function from  $C(n, w)$  to  $\{0, 1\}^+$ .

We can also apply arithmetic coding directly to  $C$  if we notice the following equation holds:

$$\begin{aligned} Q_k(i \mid x_1, \dots, x_k) &= \sum_{j=1}^i \Pr(X_{k+1} = j \mid X_1 = x_1, \dots, X_k = x_k, X_1 + \dots + X_n = w) \\ &= 1 - \Pr(X_{k+1} \geq i + 1 \mid X_1 = x_1, \dots, X_k = x_k, X_1 + \dots + X_n = w) \\ &= 1 - \binom{w - x_1 - \dots - x_k + n - k - i - 2}{n - k - 1} / \binom{w - x_1 - \dots - x_k + n - k - 1}{n - k - 1} \end{aligned}$$

# References

- Johnson Jr P D, Harris G A, Hankerson D C. Introduction to information theory and data compression[M]. Chapman and Hall/CRC, 2003.
- MacKay D J C, Mac Kay D J C. Information theory, inference and learning algorithms[M]. Cambridge university press, 2003.
- Wikipedia contributors. "Arithmetic coding." Wikipedia, The Free Encyclopedia. Wikipedia, The Free Encyclopedia, 4 Sep. 2018. Web. 13 Sep. 2018.
- Jones C. An efficient coding system for long source sequences[J]. IEEE Transactions on Information Theory, 1981, 27(3): 280-291.
- Ramabadran T V. A coding scheme for m-out-of-n codes[J]. IEEE Transactions on Communications, 1990, 38(8): 1156-1163.