

Chapter 3

Lossless Data Compression

Po-Ning Chen, Professor

Department of Electrical and Computer Engineering

National Chiao Tung University

Hsin Chu, Taiwan 30010, R.O.C.

Principle of Data Compression

I: 3-1

- Average codeword length

E.g.

$$\left\{ \begin{array}{l} P_X(x = \text{outcome}_A) = 0.5; \\ P_X(x = \text{outcome}_B) = 0.25; \\ P_X(x = \text{outcome}_C) = 0.25. \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} \text{code}(\text{outcome}_A) = 0; \\ \text{code}(\text{outcome}_B) = 10; \\ \text{code}(\text{outcome}_C) = 11. \end{array} \right.$$

Then the average codeword length is

$$\begin{aligned} & \text{len}(0) \cdot P_X(A) + \text{len}(10) \cdot P_X(B) + \text{len}(11) \cdot P_X(C) \\ &= 1 \cdot 0.5 + 2 \cdot 0.25 + 2 \cdot 0.25 \\ &= 1.5 \text{ bits.} \end{aligned}$$

- Categories of codes
 - Variable-length codes
 - Fixed-length codes (often treated as a subclass of variable-length codes)
 - * Segmentation is normally considered an implicit part of the codewords.

Principle of Data Compression

I: 3-2

Example of segmentation of fixed-length codes.

E.g. To encode the final grades of a class with 100 students.

Assume that there are three grade levels: *A*, *B* and *C*.

- Without segmentation

$$\lceil \log_2 3^{100} \rceil = 159 \text{ bits.}$$

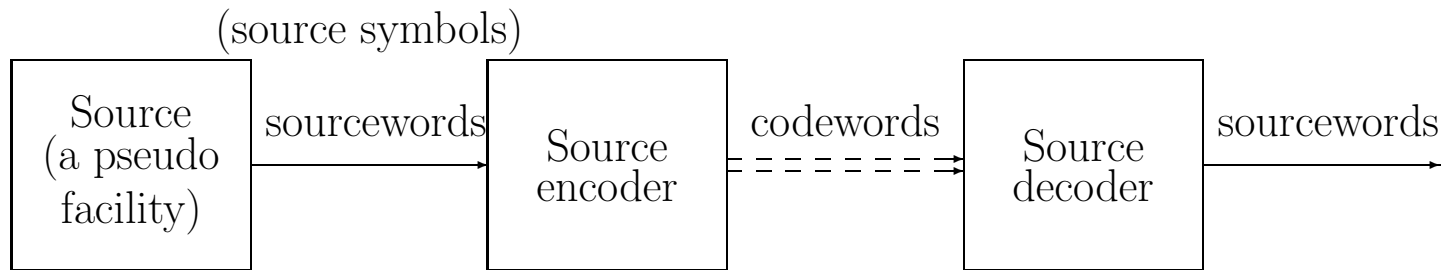
- With segmentation length of 10 students

$$10 \times \lceil \log_2 3^{10} \rceil = 160 \text{ bits.}$$

Principle of Data Compression

I: 3-3

- Fixed-length codes
 - Block codes
 - * Encoding of the next segment is independent of the previous segments
 - Fixed-length tree codes
 - * Encoding of the next segment retains and uses some knowledge of earlier segments
- Block diagram of a data compression system



Block diagram of a data compression system.

Key Difference in Data Compress Schemes

I: 3-4

- Block codes for *asymptotic* lossless data compression
 - Asymptotic in blocklength n
- Variable-length codes for *completely* lossless data compression

Block Codes for DMS

Definition 3.1 (Discrete memoryless source) A discrete memoryless source (DMS) $\{X_n\}_{n=1}^{\infty}$ consists of a sequence of independent and identically distributed (i.i.d.) random variables, X_1, X_2, X_3, \dots , all taking values in a common finite alphabet \mathcal{X} . In particular, if $P_X(\cdot)$ is the common distribution or probability mass function (pmf) of the X_i 's, then

$$P_{X^n}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P_X(x_i).$$

Block Codes for DMS

I: 3-5

Definition 3.2 An (n, M) block code¹ of blocklength n and size M for a discrete source $\{X_n\}_{n=1}^{\infty}$ is a set $\{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\} \subseteq \mathcal{X}^n$ consisting of M reproduction (or reconstruction) words, where each reproduction word is a sourceword (an n -tuple of source symbols).

- Note that \mathbf{c}_i is not a “codeword” but a “reproduction word.” It is an n -tuple of source symbols.
- One can binary-index the reproduction words in $\{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\}$ using $k \triangleq \lceil \log_2 M \rceil$ bits.
- As such k -bit words in $\{0, 1\}^k$ are usually stored for retrieval at a later date, the (n, M) block code can be represented by an encoder-decoder pair of functions (f, g) :
 - the encoding function $f : \mathcal{X}^n \rightarrow \{0, 1\}^k$ maps each sourceword x^n to a k -bit word $f(x^n)$, which we call a *codeword*.
 - the decoding function $g : \{0, 1\}^k \rightarrow \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_M\}$ is a retrieving operation that produces the reproduction words.

¹In the literature, both (n, M) and (M, n) have been used to denote a block code with blocklength n and size M . For example, R. W. Yeung adopts the former one, while T. M. Cover and J. A. Thomas use the latter. We use the (n, M) notation since $M = M_n$ is a function of n in general.

Block Codes for DMS

I: 3-6

- Since the above codewords are binary-valued, such a block code is called a *binary code*.
- More generally, a *D-ary block code* (where $D > 1$ is an integer) would use an encoding function $f : \mathcal{X}^n \rightarrow \{0, 1, \dots, D - 1\}^k$ where each codeword $f(x^n)$ contains k D -ary code symbols.
- Furthermore, since the behavior of block codes is investigated for sufficiently large n and M (tending to infinity), it is legitimate to replace $\lceil \log_2 M \rceil$ by $\log_2 M$ for the case of binary codes. With this convention, the *data compression rate* or *code rate* is

$$\text{bits required per source symbol} = \frac{k}{n} = \frac{1}{n} \log_2 M.$$

Similarly, for D -ary codes, the rate is

$$D\text{-ary code symbols required per source symbol} = \frac{k}{n} = \frac{1}{n} \log_D M.$$

For computational convenience, *nats* (under the natural logarithm) can be used instead of *bits* or *D-ary code symbols*; in this case, the code rate becomes:

$$\text{nats required per source symbol} = \frac{1}{n} \log M.$$

Block Codes for DMS

I: 3-7

- (Weakly) δ -typical set

$$\mathcal{F}_n(\delta) \triangleq \left\{ x^n \in \mathcal{X}^n : \left| -\frac{1}{n} \sum_{i=1}^n \log P_X(x_i) - H(X) \right| < \delta \right\}.$$

E.g. $n = 2$ and $\delta = 0.3$ and $\mathcal{X} = \{A, B, C, D\}$.

$$\text{The source distribution is } \begin{cases} P_X(A) = 0.4 \\ P_X(B) = 0.3 \\ P_X(C) = 0.2 \\ P_X(D) = 0.1 \end{cases}$$

The entropy equals:

$$0.4 \log \frac{1}{0.4} + 0.3 \log \frac{1}{0.3} + 0.2 \log \frac{1}{0.2} + 0.1 \log \frac{1}{0.1} = 1.27985 \text{ nats}$$

Then for $x_1^2 = (A, A)$,

$$\begin{aligned} \left| -\frac{1}{2} \sum_{i=1}^2 \log P_X(x_i) - H(X) \right| &= \left| -\frac{1}{2} (\log P_X(A) + \log P_X(A)) - 1.27985 \right| \\ &= \left| -\frac{1}{2} (\log 0.4 + \log 0.4) - 1.27985 \right| = 0.364 \end{aligned}$$

Block Codes for DMS

I: 3-8

Source	$\left -\frac{1}{2} \sum_{i=1}^2 \log P_X(x_i) - H(X) \right $
<i>AA</i>	<u>0.364</u> nats $\notin \mathcal{F}_2(0.3)$
<i>AB</i>	0.220 nats $\in \mathcal{F}_2(0.3)$
<i>AC</i>	0.017 nats $\in \mathcal{F}_2(0.3)$
<i>AD</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$
<i>BA</i>	0.220 nats $\in \mathcal{F}_2(0.3)$
<i>BB</i>	0.076 nats $\in \mathcal{F}_2(0.3)$
<i>BC</i>	0.127 nats $\in \mathcal{F}_2(0.3)$
<i>BD</i>	<u>0.473</u> nats $\notin \mathcal{F}_2(0.3)$
<i>CA</i>	0.017 nats $\in \mathcal{F}_2(0.3)$
<i>CB</i>	0.127 nats $\in \mathcal{F}_2(0.3)$
<i>CC</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$
<i>CD</i>	<u>0.676</u> nats $\notin \mathcal{F}_2(0.3)$
<i>DA</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$
<i>DB</i>	<u>0.473</u> nats $\notin \mathcal{F}_2(0.3)$
<i>DC</i>	<u>0.676</u> nats $\notin \mathcal{F}_2(0.3)$
<i>DD</i>	<u>1.023</u> nats $\notin \mathcal{F}_2(0.3)$

$$\Rightarrow \mathcal{F}_2(0.3) = \{AB, AC, BA, BB, BC, CA, CB\}.$$

Block Codes for DMS

I: 3-9

Source	$-\frac{1}{2} \sum_{i=1}^2 \log P_X(x_i) - H(X)$	codeword
<i>AA</i>	<u>0.364</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>AB</i>	0.220 nats $\in \mathcal{F}_2(0.3)$	001
<i>AC</i>	0.017 nats $\in \mathcal{F}_2(0.3)$	010
<i>AD</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>BA</i>	0.220 nats $\in \mathcal{F}_2(0.3)$	011
<i>BB</i>	0.076 nats $\in \mathcal{F}_2(0.3)$	100
<i>BC</i>	0.127 nats $\in \mathcal{F}_2(0.3)$	101
<i>BD</i>	<u>0.473</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>CA</i>	0.017 nats $\in \mathcal{F}_2(0.3)$	110
<i>CB</i>	0.127 nats $\in \mathcal{F}_2(0.3)$	111
<i>CC</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>CD</i>	<u>0.676</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>DA</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>DB</i>	<u>0.473</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>DC</i>	<u>0.676</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>
<i>DD</i>	<u>1.023</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>

We can therefore encode the **seven** outcomes in $\mathcal{F}_2(0.3)$ by **seven** distinct codewords, and encode all the remaining **nine** outcomes outside $\mathcal{F}_2(0.3)$ by a **single** codeword.

Theorem 3.3 (Shannon-McMillan) (Asymptotic equipartition property or AEP or Entropy stability property) If $\{X_n\}_{n=1}^{\infty}$ is a DMS with entropy $H(X)$, then

$$-\frac{1}{n} \log_2 P_{X^n}(X_1, \dots, X_n) \rightarrow H(X) \quad \text{in probability.}$$

In other words, for any $\delta > 0$,

$$\lim_{n \rightarrow \infty} \Pr \left\{ \left| -\frac{1}{n} \log_2 P_{X^n}(X_1, \dots, X_n) - H(X) \right| > \delta \right\} = 0.$$

- Almost all the source sequences in $\mathcal{F}_n(\delta)$ are nearly *equiprobable* or *equally surprising* (cf. Property 1 of Theorem 3.4); Hence, Theorem 3.3 is named AEP.

E.g. The probabilities of the elements in

$$\mathcal{F}_2(0.3) = \{AB, AC, BA, BB, BC, CA, CB\}$$

are respectively 0.12, 0.08, 0.12, 0.09, 0.06, 0.08 and 0.06.

The sum of these seven probability masses are 0.61.

Block Codes for DMS

I: 3-11

Source	$-\frac{1}{2} \sum_{i=1}^2 \log P_X(x_i) - H(X)$	codeword	reconstructed sequence
<i>AA</i>	<u>0.364</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>AB</i>	0.220 nats $\in \mathcal{F}_2(0.3)$	001	AB
<i>AC</i>	0.017 nats $\in \mathcal{F}_2(0.3)$	010	AC
<i>AD</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>BA</i>	0.220 nats $\in \mathcal{F}_2(0.3)$	011	BA
<i>BB</i>	0.076 nats $\in \mathcal{F}_2(0.3)$	100	BB
<i>BC</i>	0.127 nats $\in \mathcal{F}_2(0.3)$	101	BC
<i>BD</i>	<u>0.473</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>CA</i>	0.017 nats $\in \mathcal{F}_2(0.3)$	110	CA
<i>CB</i>	0.127 nats $\in \mathcal{F}_2(0.3)$	111	CB
<i>CC</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>CD</i>	<u>0.676</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>DA</i>	<u>0.330</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>DB</i>	<u>0.473</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>DC</i>	<u>0.676</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous
<i>DD</i>	<u>1.023</u> nats $\notin \mathcal{F}_2(0.3)$	<u>000</u>	ambiguous

Block Codes for DMS

I: 3-12

Theorem 3.4 (Consequence of the AEP) Given a DMS $\{X_n\}_{n=1}^{\infty}$ with entropy $H(X)$ and any δ greater than zero, then the weakly δ -typical set $\mathcal{F}_n(\delta)$ satisfies the following.

1. If $x^n \in \mathcal{F}_n(\delta)$, then

$$2^{-n(H(X)+\delta)} \leq P_{X^n}(x^n) \leq 2^{-n(H(X)-\delta)}.$$

2. $P_{X^n}(\mathcal{F}_n^c(\delta)) < \delta$ for sufficiently large n , where the superscript “c” denotes the complementary set operation.
3. $|\mathcal{F}_n(\delta)| > (1 - \delta)2^{n(H(X)-\delta)}$ for sufficiently large n , and $|\mathcal{F}_n(\delta)| \leq 2^{n(H(X)+\delta)}$ for every n , where $|\mathcal{F}_n(\delta)|$ denotes the number of elements in $\mathcal{F}_n(\delta)$.

Proof:

- Property 1 is an immediate consequence of the definition of $\mathcal{F}_n(\delta)$. I.e.,

$$\mathcal{F}_n(\delta) \triangleq \left\{ x^n \in \mathcal{X}^n : \left| -\frac{1}{n} \sum_{i=1}^n \log_2 P_X(x_i) - H(X) \right| < \delta \right\}.$$

Block Codes for DMS

I: 3-13

Thus,

$$\begin{aligned} \left| -\frac{1}{n} \sum_{i=1}^n \log_2 P_X(x_i) - H(X) \right| < \delta &\Leftrightarrow \left| -\frac{1}{n} \log_2 P_{X^n}(x^n) - H(X) \right| < \delta \\ &\Leftrightarrow H(X) - \delta < -\frac{1}{n} \log_2 P_{X^n}(x^n) < H(X) + \delta. \end{aligned}$$

- Property 2 is a direct consequence of the AEP. We nevertheless provide a direct proof of Property 2. Observe that by Chebyshev's inequality,

$$P_{X^n}(\mathcal{F}_n^c(\delta)) = P_{X^n} \left\{ x^n \in \mathcal{X}^n : \left| -\frac{1}{n} \log_2 P_{X^n}(x^n) - H(X) \right| > \delta \right\} \leq \frac{\sigma_X^2}{n\delta^2} < \delta,$$

for $n > \sigma_X^2/\delta^3$, where the variance

$$\begin{aligned} \sigma_X^2 &\triangleq \text{Var}[-\log_2 P_X(X)] = \sum_{x \in \mathcal{X}} P_X(x) [\log_2 P_X(x)]^2 - (H(X))^2 \\ &\leq E[(\log_2 P_X(X))^2] = \sum_{x \in \mathcal{X}} P_X(x) (\log_2 P_X(x))^2 \leq \sum_{x \in \mathcal{X}} \max_{0 \leq p \leq 1} p (\log_2 p)^2 \quad (p^* = \frac{1}{e^2}) \\ &= \sum_{x \in \mathcal{X}} \frac{4}{e^2} [\log_2(e)]^2 = \frac{4}{e^2} [\log_2(e)]^2 \times |\mathcal{X}| < \infty \quad \text{for finite alphabet} \end{aligned}$$

is a constant independent of n .

Block Codes for DMS

I: 3-14

- To prove Property 3, we have from Property 1 that

$$1 \geq \sum_{x^n \in \mathcal{F}_n(\delta)} P_{X^n}(x^n) \geq \sum_{x^n \in \mathcal{F}_n(\delta)} 2^{-n(H(X)+\delta)} = |\mathcal{F}_n(\delta)| 2^{-n(H(X)+\delta)},$$

and, using Properties 1 and 2, we have that

$$1 - \delta < 1 - \frac{\sigma_X^2}{n\delta^2} \leq \sum_{x^n \in \mathcal{F}_n(\delta)} P_{X^n}(x^n) \leq \sum_{x^n \in \mathcal{F}_n(\delta)} 2^{-n(H(X)-\delta)} = |\mathcal{F}_n(\delta)| 2^{-n(H(X)-\delta)},$$

for $n \geq \sigma_X^2/\delta^3$.

□

Shannon's Source Coding Theorem

I: 3-15

Theorem 3.5 (Shannon's source coding theorem) Given integer $D > 1$, consider a discrete memoryless source $\{X_n\}_{n=1}^{\infty}$ with entropy $H_D(X)$. Then the following hold.

- *Forward part (achievability):* For any $0 < \varepsilon < 1$, there exists $0 < \delta < \varepsilon$ and a sequence of D -ary block codes $\{\mathcal{C}_n = (n, M_n)\}_{n=1}^{\infty}$ with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_D M_n \leq H_D(X) + \delta \quad (3.2.1)$$

satisfying

$$P_e(\mathcal{C}_n) < \varepsilon \quad (3.2.2)$$

for all sufficiently large n , where $P_e(\mathcal{C}_n)$ denotes the probability of decoding error for block code \mathcal{C}_n .

- *Strong converse part:* For any $0 < \varepsilon < 1$, any sequence of D -ary block codes $\{\mathcal{C}_n = (n, M_n)\}_{n=1}^{\infty}$ with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_D M_n < H_D(X) \quad (3.2.3)$$

satisfies

$$P_e(\mathcal{C}_n) > 1 - \varepsilon$$

for all n sufficiently large.

Shannon's Source Coding Theorem

I: 3-16

Note

- Since ε can be made arbitrarily small, (3.2.2) is equivalent to

$$\limsup_{n \rightarrow \infty} P_e(\mathcal{C}_n) = 0.$$

- In parallel, (3.2.3) is equivalent to

$$\limsup_{n \rightarrow \infty} P_e(\mathcal{C}_n) = 1.$$

Keys of the proof of the forward part

- Only need to prove the existence of such block code.
- The code chosen is indeed the weakly δ -typical set.

Shannon's Source Coding Theorem

I: 3-17

Proof:

Forward Part:

- Without loss of generality, we will prove the result for the case of binary codes (i.e., $D = 2$). Also recall that subscript D in $H_D(X)$ will be dropped (i.e., omitted) specifically when $D = 2$.
- Given $0 < \varepsilon < 1$, fix δ such that $0 < \delta < \varepsilon$ and choose $n > 2/\delta$.
- Binary-index the sourcewords in $\mathcal{F}_n(\delta/2)$ with the following encoding map:

$$\begin{cases} x^n \rightarrow \text{binary index of } x^n, & \text{if } x^n \in \mathcal{F}_n(\delta/2); \\ x^n \rightarrow \text{all-zero codeword,} & \text{if } x^n \notin \mathcal{F}_n(\delta/2). \end{cases}$$

Then by the Shannon-McMillan AEP theorem, we obtain that

$$M_n = |\mathcal{F}_n(\delta/2)| + 1 \leq 2^{n(H(X)+\delta/2)} + 1 < 2 \cdot 2^{n(H(X)+\delta/2)} < 2^{n(H(X)+\delta)},$$

for $n > 2/\delta$. Hence, a sequence of $\mathcal{C}_n = (n, M_n)$ block code satisfying (3.2.1) is established.

- It remains to show that the error probability for this sequence of (n, M_n) block code can be made smaller than ε for all sufficiently large n .

Shannon's Source Coding Theorem

I: 3-18

By the Shannon-McMillan AEP theorem,

$$P_{X^n}(\mathcal{F}_n^c(\delta/2)) < \frac{\delta}{2} \quad \text{for all sufficiently large } n.$$

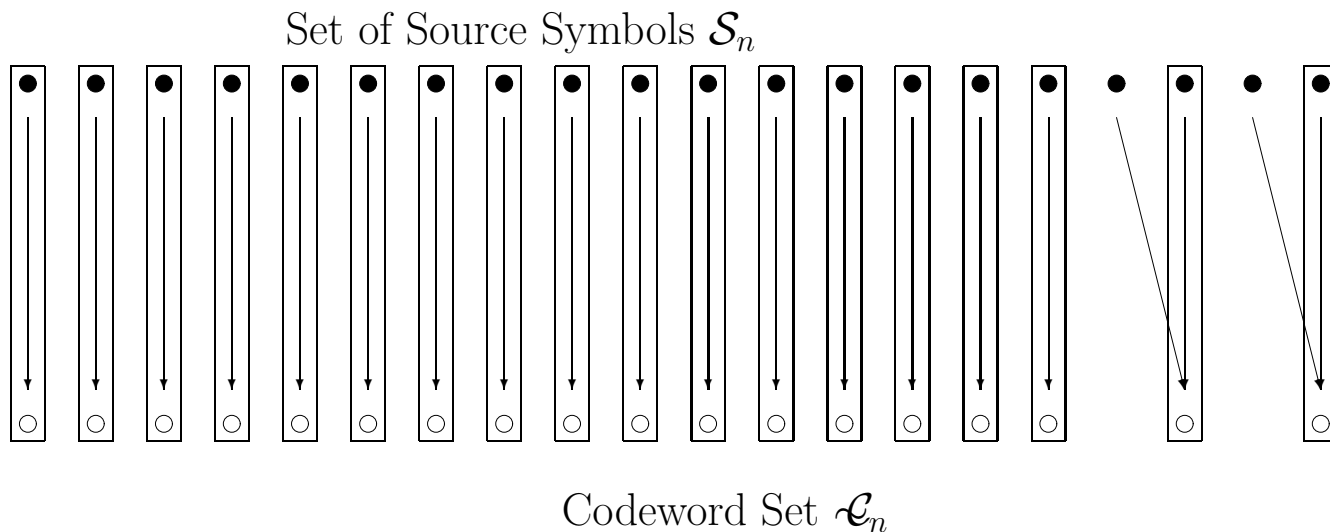
Consequently, for those n satisfying the above inequality, and being bigger than $2/\delta$,

$$P_e(\mathcal{C}_n) \leq P_{X^n}(\mathcal{F}_n^c(\delta/2)) < \delta \leq \varepsilon.$$

(See slide I: 3-11 to confirm that only the “ambiguous” sequences outside the typical set contribute to the probability of error.)

Shannon's Source Coding Theorem

I: 3-19



Strong Converse Part:

- Fix any sequence of block codes $\{\mathcal{C}_n\}_{n=1}^{\infty}$ with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_2 |\mathcal{C}_n| < H(X).$$

Let \mathcal{S}_n be the set of source symbols that can be correctly decoded through \mathcal{C}_n -coding system. (A quick example is depicted above.) Then $|\mathcal{S}_n| = |\mathcal{C}_n|$.

Shannon's Source Coding Theorem

I: 3-20

- By choosing δ small enough with $\varepsilon/2 > \delta > 0$, and also by definition of limsup operation, we have

$$(\exists N_0)(\forall n > N_0) \quad \frac{1}{n} \log_2 |\mathcal{S}_n| = \frac{1}{n} \log_2 |\mathcal{C}_n| < H(X) - 2\delta,$$

which implies

$$|\mathcal{S}_n| < 2^{n(H(X)-2\delta)}.$$

Shannon's Source Coding Theorem

I: 3-21

- Furthermore, from Property 2 of the Consequence of the AEP, we obtain that

$$(\exists N_1)(\forall n > N_1) \quad P_{X^n}(\mathcal{F}_n^c(\delta)) < \delta.$$

Consequently, for $n > N \triangleq \max\{N_0, N_1, \log_2(2/\varepsilon)/\delta\}$, the probability of correctly block decoding satisfies

$$\begin{aligned} 1 - P_e(\mathcal{C}_n) &= \sum_{x^n \in \mathcal{S}_n} P_{X^n}(x^n) \\ &= \sum_{x^n \in \mathcal{S}_n \cap \mathcal{F}_n^c(\delta)} P_{X^n}(x^n) + \sum_{x^n \in \mathcal{S}_n \cap \mathcal{F}_n(\delta)} P_{X^n}(x^n) \\ &\leq P_{X^n}(\mathcal{F}_n^c(\delta)) + |\mathcal{S}_n \cap \mathcal{F}_n(\delta)| \cdot \max_{x^n \in \mathcal{F}_n(\delta)} P_{X^n}(x^n) \\ &< \delta + |\mathcal{S}_n| \cdot \max_{x^n \in \mathcal{F}_n(\delta)} P_{X^n}(x^n) \\ &< \frac{\varepsilon}{2} + 2^{n(H(X)-2\delta)} \cdot 2^{-n(H(X)-\delta)} \\ &< \frac{\varepsilon}{2} + 2^{-n\delta} \\ &< \varepsilon, \end{aligned}$$

which is equivalent to $P_e(\mathcal{C}_n) > 1 - \varepsilon$ for $n > N$.

□

Code Rates for Data Compression

I: 3-22

Notes

- Ultimate data compression rate

$$R \triangleq \limsup_{n \rightarrow \infty} \frac{1}{n} \log M_n \text{ nats per source symbol.}$$

- Shannon's source coding theorem

– Arbitrary good performance can be achieved by **extending the block-length**.

$$(\forall \varepsilon > 0 \text{ and } 0 < \delta < \varepsilon)(\exists \mathcal{C}_n) \text{ such that } \frac{1}{n} \log M_n < H(X) + \delta \text{ and } P_e(\mathcal{C}_n) < \varepsilon.$$

So $R = \limsup_{n \rightarrow \infty} \frac{1}{n} \log M_n$ can be made smaller than $H(X) + \delta$ for arbitrarily small δ .

In other words, at rate $R < H(X) + \delta$ for arbitrarily small $\delta > 0$, the error probability can be made *arbitrarily zero* ($< \varepsilon$).

- How about further making $R < H(X)$? Answer:

$$\left(\forall \{ \mathcal{C}_n \}_{n \geq 1} \text{ with } \limsup_{n \rightarrow \infty} \frac{1}{n} \log |\mathcal{C}_n| < H(X) \right) P_e(\mathcal{C}_n) \rightarrow 1.$$

Summary of Shannon's Source Coding Theorem

I: 3-24

- Notes to the strong converse theorem
 - It is named the *strong converse theorem* because the result is very *strong*.
 - * All code sequences with $R < H(X)$ have error probability approaching 1!
 - Of course, you can always design a lousy code with error probability approaching 1. Here, what the theorem truly claims is that all designs are *lousy*.
 - The strong converse theorem applies to all stationary-ergodic sources.
- A weak converse statement (than the strong converse) is:
 - For general sources, such as non-stationary non-ergodic sources, we can find some code sequence with $R < H(X)$ whose error probability is only bounded away from zero, and does not approach 1 at all.

Block Codes for Stationary Ergodic Sources

I: 3-25

- Recall that the merit of the *stationary ergodic* assumption is on its validity of *law of large numbers*.
- However, in order to extend the Shannon's source coding theorem to stationary ergodic sources, we need to generalize the information measure for such sources.

Definition 3.7 (Entropy rate) The *entropy rate* for a source $\{X_n\}_{n=1}^{\infty}$ is denoted by $H(\mathcal{X})$ and defined by

$$H(\mathcal{X}) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} H(X^n)$$

provided the limit exists, where $X^n = (X_1, \dots, X_n)$.

- **Comment:** The limit of $\lim_{n \rightarrow \infty} \frac{1}{n} H(X^n)$ exists for all stationary sources.

Lemma 3.8 For a stationary source $\{X_n\}_{n=1}^{\infty}$, the conditional entropy

$$H(X_n|X_{n-1}, \dots, X_1)$$

is non-increasing in n and also bounded from below by zero. Hence by Lemma A.20, the limit

$$\lim_{n \rightarrow \infty} H(X_n|X_{n-1}, \dots, X_1)$$

exists.

Proof: We have

$$H(X_n|X_{n-1}, \dots, X_1) \leq H(X_n|X_{n-1}, \dots, X_2) \tag{3.2.4}$$

$$\begin{aligned} &= H(X_n, \dots, X_2) - H(X_{n-1}, \dots, X_2) \\ &= H(X_{n-1}, \dots, X_1) - H(X_{n-2}, \dots, X_1) \tag{3.2.5} \\ &= H(X_{n-1}|X_{n-2}, \dots, X_1) \end{aligned}$$

where (3.2.4) follows since conditioning never increases entropy, and (3.2.5) holds because of the stationarity assumption. Finally, recall that each conditional entropy $H(X_n|X_{n-1}, \dots, X_1)$ is non-negative. □

Block Codes for Stationary Ergodic Sources

I: 3-27

Lemma 3.9 (Cesaro-mean theorem) If $a_n \rightarrow a$ as $n \rightarrow \infty$ and $b_n = (1/n) \sum_{i=1}^n a_i$, then $b_n \rightarrow a$ as $n \rightarrow \infty$.

Proof: $a_n \rightarrow a$ implies that for any $\varepsilon > 0$, there exists N such that for all $n > N$, $|a_n - a| < \varepsilon$. Then

$$\begin{aligned} |b_n - a| &= \left| \frac{1}{n} \sum_{i=1}^n (a_i - a) \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n |a_i - a| \\ &= \frac{1}{n} \sum_{i=1}^N |a_i - a| + \frac{1}{n} \sum_{i=N+1}^n |a_i - a| \\ &\leq \frac{1}{n} \sum_{i=1}^N |a_i - a| + \frac{n - N}{n} \varepsilon. \end{aligned}$$

Hence, $\lim_{n \rightarrow \infty} |b_n - a| \leq \varepsilon$. Since ε can be made arbitrarily small, the lemma holds. \square

Block Codes for Stationary Ergodic Sources

I: 3-28

Theorem 3.10 For a stationary source $\{X_n\}_{n=1}^{\infty}$, its entropy rate always exists and is equal to

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1).$$

Proof: The result directly follows by writing

$$\frac{1}{n}H(X^n) = \frac{1}{n} \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1) \quad (\text{chain-rule for entropy})$$

and applying the Cesaro-mean theorem. □

Observation 3.11 It can also be shown that for a stationary source, $(1/n)H(X^n)$ is non-increasing in n and $(1/n)H(X^n) \geq H(X_n | X_{n-1}, \dots, X_1)$ for all $n \geq 1$. (The proof is left as an exercise. See Problem 3.)

Practices of Finding the Entropy Rate

I: 3-29

- I.i.d. source

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X^n) = H(X)$$

since $H(X^n) = n \times H(X)$ for every n .

- First-order stationary Markov source

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X^n) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1) = H(X_2 | X_1),$$

where

$$H(X_2 | X_1) \triangleq - \sum_{x_1 \in \mathcal{X}} \sum_{x_2 \in \mathcal{X}} \pi(x_1) P_{X_2 | X_1}(x_2 | x_1) \cdot \log P_{X_2 | X_1}(x_2 | x_1),$$

and $\pi(\cdot)$ is the stationary distribution for the Markov source.

- In addition, if the Markov source is also *binary*,

$$\lim_{n \rightarrow \infty} \frac{1}{n} H(X^n) = \frac{\beta}{\alpha + \beta} H_b(\alpha) + \frac{\alpha}{\alpha + \beta} H_b(\beta),$$

where $H_b(\alpha) \triangleq -\alpha \log \alpha - (1 - \alpha) \log(1 - \alpha)$ is the binary entropy function, and $P_{X_2 | X_1}(0|1) = \alpha$ and $P_{X_2 | X_1}(1|0) = \beta$

Theorem 3.12 (Generalized AEP or Shannon-McMillan-Breiman Theorem [12]) If $\{X_n\}_{n=1}^{\infty}$ is a stationary ergodic source, then

$$-\frac{1}{n} \log_2 P_{X^n}(X_1, \dots, X_n) \xrightarrow{a.s.} H(\mathcal{X}).$$

Theorem 3.13 (Shannon's source coding theorem for stationary ergodic sources) Given integer $D > 1$, let $\{X_n\}_{n=1}^{\infty}$ be a stationary ergodic source with entropy rate (in base D)

$$H_D(\mathcal{X}) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} H_D(X^n).$$

Then the following hold.

- *Forward part (achievability):* For any $0 < \varepsilon < 1$, there exists δ with $0 < \delta < \varepsilon$ and a sequence of D -ary block codes $\{\mathcal{C}_n = (n, M_n)\}_{n=1}^{\infty}$ with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_D M_n < H_D(\mathcal{X}) + \delta,$$

and probability of decoding error satisfied

$$P_e(\mathcal{C}_n) < \varepsilon$$

for all sufficiently large n .

Shannon's Source Coding Theorem Revisited

I: 3-31

- *Strong converse part:* For any $0 < \varepsilon < 1$, any sequence of D -ary block codes $\{\mathcal{C}_n = (n, M_n)\}_{n=1}^{\infty}$ with

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log_D M_n < H_D(\mathcal{X})$$

satisfies

$$P_e(\mathcal{C}_n) > 1 - \varepsilon$$

for all n sufficiently large.

Problems of Ergodicity Assumption

I: 3-32

- A discrete memoryless (i.i.d.) source is stationary and ergodic.
- In general, it is hard to check whether a process is ergodic or not.
- A stationary process is a mixture of several stationary ergodic processes (i.e., its n -fold distribution can be written as the mean of the n -fold distributions of stationary ergodic processes) if, and only if, it is *not ergodic*.
 - For example, let P and Q be two distributions on a finite alphabet \mathcal{X} such that the process $\{X_n\}_{n=1}^{\infty}$ is i.i.d. with distribution P and the process $\{Y_n\}_{n=1}^{\infty}$ is i.i.d. with distribution Q . Flip a biased coin (with Heads probability equal to θ , $0 < \theta < 1$) *once* and let

$$Z_i = \begin{cases} X_i & \text{if Heads} \\ Y_i & \text{if Tails} \end{cases}$$

for $i = 1, 2, \dots$. Then the resulting process $\{Z_i\}_{i=1}^{\infty}$ has its n -fold distribution as a mixture of the n -fold distributions of $\{X_n\}_{n=1}^{\infty}$ and $\{Y_n\}_{n=1}^{\infty}$:

$$P_{Z^n}(a^n) = \theta P_{X^n}(a^n) + (1 - \theta) P_{Y^n}(a^n)$$

for all $a^n \in \mathcal{X}^n$, $n = 1, 2, \dots$. Then the process $\{Z_i\}_{i=1}^{\infty}$ is stationary but *not ergodic*.

Problems of Ergodicity Assumption

I: 3-33

- A specific case that ergodicity can be verified is that of the Markov sources.

Observation

1. An irreducible finite-state Markov source is ergodic.
 - Note that irreducibility can be verified in terms of the transition probability matrix. For example, all the entries in transition probability matrix are non-zero.
2. The generalized AEP theorem holds for irreducible stationary Markov sources. For example, if the Markov source is of the first-order, then

$$-\frac{1}{n} \log P_{X^n}(X^n) \xrightarrow{a.s.} \lim_{n \rightarrow \infty} \frac{1}{n} H(X^n) = H(X_2|X_1).$$

Redundancy for Lossless Data Compression

I: 3-34

- A source can be compressed only when it has redundancy.
 - A very important concept is that **the output of a perfect lossless data compressor should be (asymptotic) i.i.d. with (asymptotic) uniform marginal distribution**. Because if it were not so, there would be redundancy in the output and hence the compressor cannot be claimed **perfect**.
- This arises the need to define the redundancy of a source.
- Categories of redundancy
 - intra-sourcword redundancy
 - * due to non-uniform marginal distribution
 - inter-sourcword redundancy
 - * due to the source memory

Definition (Redundancy)

1. Source redundancy due to the *non-uniformity of the source marginal distribution* ρ_d :

$$\rho_d \triangleq \log |\mathcal{X}| - H(X_1).$$

2. Source redundancy due to the *source memory* ρ_m :

$$\rho_m \triangleq H(X_1) - H(\mathcal{X}).$$

3. Hence, the source total redundancy ρ_t is given by:

$$\rho_t \triangleq \rho_d + \rho_m = \log |\mathcal{X}| - H(\mathcal{X}).$$

Redundancy for Lossless Data Compression

I: 3-36

E.g.

Source	ρ_d	ρ_m	ρ_t
i.i.d. uniform	0	0	0
i.i.d. non-uniform	$\log_2 \mathcal{X} - H(X_1)$	0	ρ_d
1st-order symmetric Markov	0	$H(X_1) - H(X_2 X_1)$	ρ_m
1st-order non-symmetric Markov	$\log_2 \mathcal{X} - H(X_1)$	$H(X_1) - H(X_2 X_1)$	$\rho_d + \rho_m$

- A first-order Markov process is symmetric if for any x_1 and \hat{x}_1 ,

$$\{a : a = P_{X_2|X_1}(y|x_1) \text{ for some } y\} = \{a : a = P_{X_2|X_1}(y|\hat{x}_1) \text{ for some } y\}.$$

Variable-Length Code for Lossless Data Compression I: 3-37

- Non-singular codes
 - To encode all sourcewords with distinct variable-length codewords
- Uniquely decodable codes
 - Concatenation of codewords (without punctuation mechanism) can be uniquely decodable.

E.g., a non-singular but non-uniquely decodable code

code of A = 0,

code of B = 1,

code of C = 00,

code of D = 01,

code of E = 10,

code of F = 11.

The code is not uniquely decodable because the codeword sequence, 01, can be reconstructed as either AB or D .

Variable-Length Code for Lossless Data Compression I: 3-38

Definition 3.14 Consider a discrete source $\{X_n\}_{n=1}^{\infty}$ with finite alphabet \mathcal{X} along with a D -ary code alphabet $\mathcal{B} = \{0, 1, \dots, D - 1\}$, where $D > 1$ is an integer. Fix integer $n \geq 1$, then a D -ary n -th order variable-length code (VLC) is a map

$$f : \mathcal{X}^n \rightarrow \mathcal{B}^*$$

mapping (fixed-length) sourcewords of length n to D -ary codewords in \mathcal{B}^* of variable lengths, where \mathcal{B}^* denotes the set of all finite-length strings from \mathcal{B} (i.e., $c \in \mathcal{B}^* \Leftrightarrow \exists$ integer $l \geq 1$ such that $c \in \mathcal{B}^l$).

The *codebook* \mathcal{C} of a VLC is the set of all codewords:

$$\mathcal{C} = f(\mathcal{X}^n) = \{f(x^n) \in \mathcal{B}^* : x^n \in \mathcal{X}^n\}.$$

Variable-Length Code for Lossless Data Compression I: 3-39

Definition 3.15 Let \mathcal{C} be a D -ary n -th order VLC code

$$f : \mathcal{X}^n \rightarrow \{0, 1, \dots, D - 1\}^*$$

for a discrete source $\{X_n\}_{n=1}^{\infty}$ with alphabet \mathcal{X} and distribution $P_{X^n}(x^n)$, $x^n \in \mathcal{X}^n$. Setting $\ell(\mathbf{c}_{x^n})$ as the length of the codeword $\mathbf{c}_{x^n} = f(x^n)$ associated with sourceword x^n , then the *average codeword length* for \mathcal{C} is given by

$$\bar{\ell} \triangleq \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \ell(\mathbf{c}_{x^n})$$

and its *average code rate* (in D -ary code symbols/source symbol) is given by

$$\bar{R}_n \triangleq \frac{\bar{\ell}}{n} = \frac{1}{n} \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \ell(\mathbf{c}_{x^n}).$$

Variable-Length Code for Lossless Data Compression I: 3-40

Theorem 3.16 (Kraft inequality for uniquely decodable codes) Let \mathcal{C} be a uniquely decodable D -ary n -th order VLC for a discrete source $\{X_n\}_{n=1}^{\infty}$ with alphabet \mathcal{X} . Let the $M = |\mathcal{X}|^n$ codewords of \mathcal{C} have lengths $\ell_1, \ell_2, \dots, \ell_M$, respectively. Then the following inequality must hold

$$\sum_{m=1}^M D^{-\ell_m} \leq 1.$$

Proof: Suppose that we use the codebook \mathcal{C} to encode N sourcewords ($x_i^n \in \mathcal{X}^n$, $i = 1, \dots, N$) arriving in a sequence; this yields a concatenated codeword sequence

$$\mathbf{c}_1 \mathbf{c}_2 \mathbf{c}_3 \dots \mathbf{c}_N.$$

Let the lengths of the codewords be respectively denoted by

$$\ell(\mathbf{c}_1), \ell(\mathbf{c}_2), \dots, \ell(\mathbf{c}_N).$$

Consider

$$\left(\sum_{\mathbf{c}_1 \in \mathcal{C}} \sum_{\mathbf{c}_2 \in \mathcal{C}} \dots \sum_{\mathbf{c}_N \in \mathcal{C}} D^{-[\ell(\mathbf{c}_1) + \ell(\mathbf{c}_2) + \dots + \ell(\mathbf{c}_N)]} \right).$$

It is obvious that the above expression is equal to

$$\left(\sum_{\mathbf{c} \in \mathcal{C}} D^{-\ell(\mathbf{c})} \right)^N = \left(\sum_{m=1}^M D^{-\ell_m} \right)^N.$$

Variable-Length Code for Lossless Data Compression I: 3-41

(Note that $|\mathcal{C}| = M$.) On the other hand, all the code sequences with length

$$i = \ell(\mathbf{c}_1) + \ell(\mathbf{c}_2) + \cdots + \ell(\mathbf{c}_N)$$

contribute equally to the sum of the identity, which is D^{-i} . Let A_i denote the number of N -codeword sequences that have length i . Then the above identity can be re-written as

$$\left(\sum_{m=1}^M D^{-\ell_m} \right)^N = \sum_{i=1}^{LN} A_i D^{-i}, \quad \text{where } L \triangleq \max_{\mathbf{c} \in \mathcal{C}} \ell(\mathbf{c}).$$

Since \mathcal{C} is by assumption a uniquely decodable code, the codeword sequence must be unambiguously decodable. Observe that a code sequence with length i has at most D^i unambiguous combinations. Therefore, $A_i \leq D^i$, and

$$\left(\sum_{m=1}^M D^{-\ell_m} \right)^N = \sum_{i=1}^{LN} A_i D^{-i} \leq \sum_{i=1}^{LN} D^i D^{-i} = LN,$$

which implies that

$$\sum_{m=1}^M D^{-\ell_m} \leq (LN)^{1/N}.$$

The proof is completed by noting that the above inequality holds for every N , and the upper bound $(LN)^{1/N}$ goes to 1 as N goes to infinity. \square

Source Coding Theorem for Variable-Length Code

I: 3-42

Theorem 3.17 The average rate of every uniquely decodable D -ary n -th order VLC for a discrete memoryless source $\{X_n\}_{n=1}^{\infty}$ is lower-bounded by the source entropy $H_D(X)$ (measured in D -ary code symbols/source symbol).

Proof: Consider a uniquely decodable D -ary n -th order VLC code for the source $\{X_n\}_{n=1}^{\infty}$

$$f : \mathcal{X}^n \rightarrow \{0, 1, \dots, D - 1\}^*$$

and let $\ell(\mathbf{c}_{x^n})$ denote the length of the codeword $\mathbf{c}_{x^n} = f(x^n)$ for sourceword x^n .

Hence,

$$\begin{aligned}
 \bar{R}_n - H_D(X) &= \frac{1}{n} \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \ell(\mathbf{c}_{x^n}) - \frac{1}{n} H_D(X^n) \\
 &= \frac{1}{n} \left[\sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \ell(\mathbf{c}_{x^n}) - \sum_{x^n \in \mathcal{X}^n} (-P_{X^n}(x^n) \log_D P_{X^n}(x^n)) \right] \\
 &= \frac{1}{n} \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \log_D \frac{P_{X^n}(x^n)}{D^{-\ell(\mathbf{c}_{x^n})}} \\
 &\geq \frac{1}{n} \left[\sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \right] \log_D \frac{[\sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n)]}{[\sum_{x^n \in \mathcal{X}^n} D^{-\ell(\mathbf{c}_{x^n})}]} \\
 &\quad \text{(log-sum inequality)} \\
 &= -\frac{1}{n} \log \left[\sum_{x^n \in \mathcal{X}^n} D^{-\ell(\mathbf{c}_{x^n})} \right] \\
 &\geq 0
 \end{aligned}$$

where the last inequality follows from the Kraft inequality for uniquely decodable codes and the fact that the logarithm is a strictly increasing function. \square

Summaries for Uniquely Decodability

I: 3-44

1. Uniquely decodability \Rightarrow the Kraft inequality holds.
2. Uniquely decodability \Rightarrow average code rate of VLCs for memoryless sources is lower bounded by the source entropy.

Exercise 3.18

1. Find a non-singular and also non-uniquely decodable code that violates the Kraft inequality. (Hint: Slide I: 3-37.)
2. Find a non-singular and also non-uniquely decodable code that beats the entropy lower bound. (Hint: Same as the previous one.)

A Special Case of Uniquely Decodable Codes

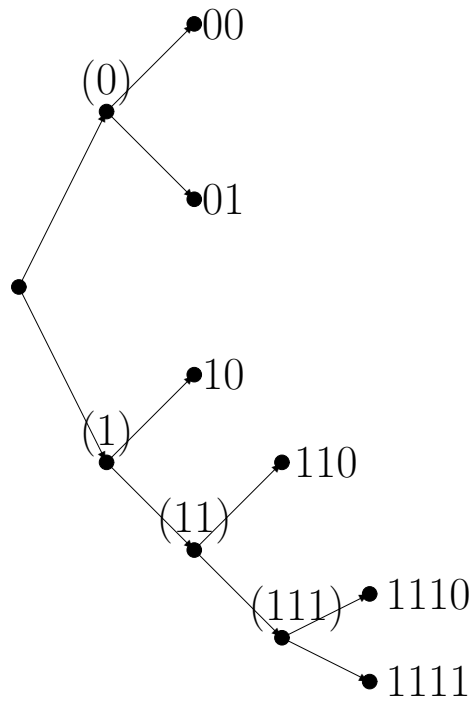
I: 3-45

- Prefix codes or instantaneous codes
 - Note that a uniquely decodable code may not necessarily be decoded instantaneously.

Definition 3.19 A code is called a *prefix code* or an *instantaneous code* if no codeword is a prefix of any other codeword.

Tree Representation of Prefix Codes

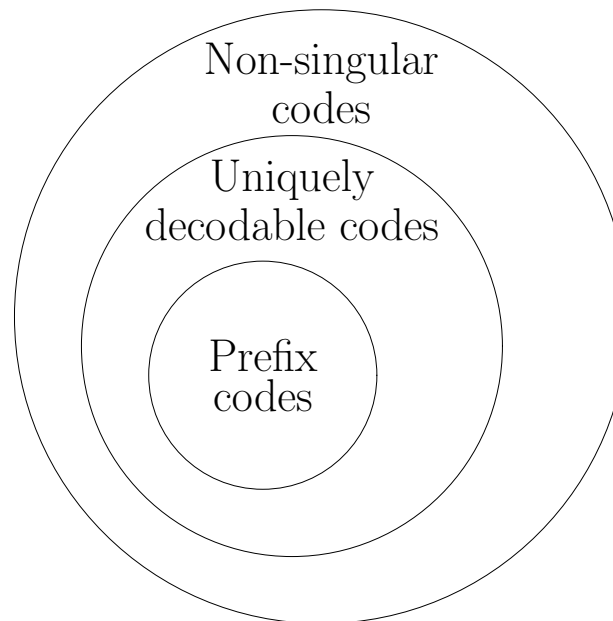
I: 3-46



The codewords are those residing on the leaves, which in this case are 00, 01, 10, 110, 1110 and 1111.

Classification of Variable-Length Codes

I: 3-47



Prefix Code to Kraft Inequality

I: 3-48

Theorem 3.20 (Kraft inequality for prefix codes) There exists a D -ary n th-order prefix code for a discrete source $\{X_n\}_{n=1}^{\infty}$ with alphabet \mathcal{X} if, and only if, the codewords of length ℓ_m , $m = 1, \dots, M$, satisfy the Kraft inequality, where $M = |\mathcal{X}|^n$.

Proof: Without loss of generality, we provide the proof for the case of $D = 2$ (binary codes).

1. [**The forward part**] *Prefix codes satisfy the Kraft inequality.*

The codewords of a prefix code can always be put on a tree. Pick up a length

$$\ell_{\max} \triangleq \max_{1 \leq m \leq M} \ell_m.$$

- A tree has originally $2^{\ell_{\max}}$ nodes on level ℓ_{\max} .
- Each codeword of length ℓ_m obstructs $2^{\ell_{\max} - \ell_m}$ nodes on level ℓ_{\max} .
 - In other words, when any node is chosen as a codeword, all its children will be excluded from being codewords.
 - There are exactly $2^{\ell_{\max} - \ell_m}$ excluded nodes on level ℓ_{\max} of the tree. We therefore say that each codeword of length ℓ_m obstructs $2^{\ell_{\max} - \ell_m}$ nodes on level ℓ_{\max} .

Prefix Code to Kraft Inequality

I: 3-49

- Note that no two codewords obstruct the same nodes on level ℓ_{\max} . Hence the number of totally obstructed codewords on level ℓ_{\max} should be less than $2^{\ell_{\max}}$, i.e.,

$$\sum_{m=1}^M 2^{\ell_{\max} - \ell_m} \leq 2^{\ell_{\max}},$$

which immediately implies the Kraft inequality:

$$\sum_{m=1}^M 2^{-\ell_m} \leq 1.$$

This part can also be proven by stating the fact that a prefix code is a uniquely decodable code. The objective of adding this proof is to illustrate the characteristics of a tree-like prefix code.

Prefix Code to Kraft Inequality

I: 3-50

2. [**The converse part**] *Kraft inequality implies the existence of a prefix code.*

Suppose that $\ell_1, \ell_2, \dots, \ell_M$ satisfy the Kraft inequality. We will show that there exists a binary tree with M selected nodes where the i^{th} node resides on level ℓ_i .

- Let n_i be the number of nodes (among the M nodes) residing on level i (namely, n_i is the number of codewords with length i or $n_i = |\{m : \ell_m = i\}|$), and let

$$\ell_{\max} \triangleq \max_{1 \leq m \leq M} \ell_m.$$

- Then from the Kraft inequality, we have

$$n_1 2^{-1} + n_2 2^{-2} + \dots + n_{\ell_{\max}} 2^{-\ell_{\max}} \leq 1.$$

- The above inequality can be re-written in a form that is more suitable for this proof as:

$$n_1 2^{-1} \leq 1$$

$$n_1 2^{-1} + n_2 2^{-2} \leq 1$$

...

$$n_1 2^{-1} + n_2 2^{-2} + \dots + n_{\ell_{\max}} 2^{-\ell_{\max}} \leq 1.$$

Prefix Code to Kraft Inequality

I: 3-51

Hence,

$$\begin{aligned}n_1 &\leq 2 \\n_2 &\leq 2^2 - n_1 2^1 \\&\dots \\n_{\ell_{\max}} &\leq 2^{\ell_{\max}} - n_1 2^{\ell_{\max}-1} - \dots - n_{\ell_{\max}-1} 2^1,\end{aligned}$$

which can be interpreted in terms of a tree model as:

- the 1st inequality says that the number of codewords of length 1 is less than the available number of nodes on the 1st level, which is 2.
- The 2nd inequality says that the number of codewords of length 2 is less than the total number of nodes on the 2nd level, which is 2^2 , minus the number of nodes obstructed by the 1st level nodes already occupied by codewords.
- The succeeding inequalities demonstrate the availability of a sufficient number of nodes at each level after the nodes blocked by shorter length codewords have been removed.
- Because this is true at every codeword length up to the maximum codeword length, the assertion of the theorem is proved. □

Source Coding Theorem for Variable-Length Codes I: 3-52

Corollary 3.21 A uniquely decodable D -ary n -th order code can always be replaced by a D -ary n -th order prefix code with the same average codeword length (and hence the same average code rate).

Proof: Uniquely decodability

\Rightarrow the Kraft inequality holds.

\Rightarrow [**The converse part**] *Kraft inequality implies the existence of a prefix code.*

Theorem 3.22 Consider a discrete memoryless source $\{X_n\}_{n=1}^{\infty}$.

1. For any D -ary n -th order prefix code for the source, the average code rate is no less than the source entropy $H_D(X)$.
2. There must exist a D -ary n -th order prefix code for the source whose average code rate is no greater than $H_D(X) + \frac{1}{n}$, namely,

$$\bar{R}_n \triangleq \frac{1}{n} \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \ell(\mathbf{c}_{x^n}) \leq H_D(X) + \frac{1}{n}, \quad (3.3.1)$$

where \mathbf{c}_{x^n} is the codeword for sourceword x^n , and $\ell(\mathbf{c}_{x^n})$ is the length of codeword \mathbf{c}_{x^n} .

Source Coding Theorem for Variable-Length Codes I: 3-53

Proof: A prefix code is uniquely decodable, and hence it directly follows from Theorem 3.17 that its average code rate is no less than the source entropy.

To prove the second part, we can design a prefix code satisfying both (3.3.1) and the Kraft inequality, which immediately implies the existence of the desired code by Theorem 3.20.

- Choose the codeword length for sourceword x^n as

$$\ell(\mathbf{c}_{x^n}) = \lfloor -\log_D P_{X^n}(x^n) \rfloor + 1. \quad (3.3.2)$$

Then

$$D^{-\ell(\mathbf{c}_{x^n})} \leq P_{X^n}(x^n).$$

- Summing both sides over all source symbols, we obtain

$$\sum_{x^n \in \mathcal{X}^n} D^{-\ell(\mathbf{c}_{x^n})} \leq 1,$$

which is exactly the Kraft inequality.

Source Coding Theorem for Variable-Length Codes

I: 3-54

- On the other hand, (3.3.2) implies

$$\ell(\mathbf{c}_{x^n}) \leq -\log_D P_{X^n}(x^n) + 1,$$

which in turn implies

$$\begin{aligned} \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \ell(\mathbf{c}_{x^n}) &\leq \sum_{x^n \in \mathcal{X}^n} [-P_{X^n}(x^n) \log_D P_{X^n}(x^n)] + \sum_{x^n \in \mathcal{X}^n} P_{X^n}(x^n) \\ &= H_D(X^n) + 1 = nH_D(X) + 1, \end{aligned}$$

where the last equality holds since the source is memoryless. □

Prefix Codes for Block Sourcewords

I: 3-55

E.g. A memoryless source with source alphabet

$$\{A, B, C\}$$

and probability distribution

$$P_X(A) = 0.8, \quad P_X(B) = P_X(C) = 0.1$$

has entropy being equal to

$$-0.8 \cdot \log_2 0.8 - 0.1 \cdot \log_2 0.1 - 0.1 \cdot \log_2 0.1 = 0.92 \text{ bits.}$$

- One of the best binary first-order or single-letter encoding (with $n = 1$) prefix codes for this source is given by

$$\mathbf{c}(A) = 0, \mathbf{c}(B) = 10 \text{ and } \mathbf{c}(C) = 11,$$

where $\mathbf{c}(\cdot)$ is the encoding function.

- Then the resultant average code rate for this code is

$$0.8 \times 1 + 0.2 \times 2 = 1.2 \text{ bits} \geq 0.92 \text{ bits.}$$

The optimal variable-length code for a specific source X has average codeword length <i>strictly larger</i> than the source entropy.
--

Prefix Codes for Block Sourcewords

I: 3-56

- Now if we consider a second-order (with $n = 2$) prefix code by encoding two consecutive source symbols at a time, the new source alphabet becomes

$$\{AA, AB, AC, BA, BB, BC, CA, CB, CC\},$$

and the resultant probability distribution is calculated by

$$(\forall x_1, x_2 \in \{A, B, C\}) \quad P_{X^2}(x_1, x_2) = P_X(x_1)P_X(x_2)$$

as the source is memoryless. Then one of the best binary prefix codes for the source is given by

$$\mathbf{c}(AA) = 0$$

$$\mathbf{c}(AB) = 100$$

$$\mathbf{c}(AC) = 101$$

$$\mathbf{c}(BA) = 110$$

$$\mathbf{c}(BB) = 111100$$

$$\mathbf{c}(BC) = 111101$$

$$\mathbf{c}(CA) = 1110$$

$$\mathbf{c}(CB) = 111110$$

$$\mathbf{c}(CC) = 111111.$$

Prefix Codes for Block Sourcewords

I: 3-57

- The average code rate of this code now becomes

$$\frac{0.64(1 \times 1) + 0.08(3 \times 3 + 4 \times 1) + 0.01(6 \times 4)}{2} = 0.96 \text{ bits,}$$

which is closer to the source entropy of 0.92 bits.

- As n increases, the average code rate will be brought closer to the source entropy.

Prefix Codes for Block Sourcewords

I: 3-58

Theorem 3.23 (Lossless variable-length source coding theorem) Fix integer $D > 1$ and consider a discrete memoryless source $\{X_n\}_{n=1}^{\infty}$ with distribution P_X and entropy $H_D(X)$ (measured in D -ary units). Then the following hold.

- *Forward part (achievability):* For any $\varepsilon > 0$, there exists a D -ary n -th order prefix (hence uniquely decodable) code

$$f : \mathcal{X}^n \rightarrow \{0, 1, \dots, D - 1\}^*$$

for the source with an average code rate \bar{R}_n satisfying

$$\bar{R}_n \leq H_D(X) + \varepsilon$$

for n sufficiently large.

- *Converse part:* Every uniquely decodable code

$$f : \mathcal{X}^n \rightarrow \{0, 1, \dots, D - 1\}^*$$

for the source has an average code rate $\bar{R}_n \geq H_D(X)$.

Proof: The forward part follows directly from Theorem 3.22 by choosing n large enough such that $1/n < \varepsilon$, and the converse part is already given by Theorem 3.17 (cf. Slide I: 3-42). □

Final Note on Prefix Codes

I: 3-59

Observation 3.24 Theorem 3.23 actually also holds for the class of *stationary sources* by replacing the source entropy $H_D(X)$ with the source entropy rate

$$H_D(\mathcal{X}) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} H_D(X^n),$$

measured in D -ary units. The proof is very similar to the proofs of Theorems 3.17 and 3.22 with slight modifications (such as using the fact that $\frac{1}{n} H_D(X^n)$ is non-increasing with n for stationary sources).

Huffman Codes

I: 3-60

Lemma 3.25 Let \mathcal{C} be an optimal binary prefix code with codeword lengths ℓ_i , $i = 1, \dots, M$, for a source with alphabet $\mathcal{X} = \{a_1, \dots, a_M\}$ and symbol probabilities p_1, \dots, p_M . We assume, without loss of generality, that

$$p_1 \geq p_2 \geq p_3 \geq \dots \geq p_M,$$

and that any group of source symbols with identical probability is listed in order of increasing codeword length (i.e., if $p_i = p_{i+1} = \dots = p_{i+s}$, then $\ell_i \leq \ell_{i+1} \leq \dots \leq \ell_{i+s}$). Then the following properties hold.

1. Higher probability source symbols have shorter codewords: $p_i > p_j$ implies $\ell_i \leq \ell_j$, for $i, j = 1, \dots, M$.
2. The two least probable source symbols have codewords of equal length: $\ell_{M-1} = \ell_M$.
3. Among the codewords of length ℓ_M , two of the codewords are identical except in the last digit.

Huffman Codes

I: 3-61

Proof:

- 1) If $p_i > p_j$ and $\ell_i > \ell_j$, then it is possible to construct a better code \mathcal{C}' by interchanging (“swapping”) codewords i and j of \mathcal{C} , since

$$\begin{aligned}\bar{\ell}(\mathcal{C}') - \bar{\ell}(\mathcal{C}) &= p_i \ell_j + p_j \ell_i - (p_i \ell_i + p_j \ell_j) \\ &= (p_i - p_j)(\ell_j - \ell_i) \\ &< 0.\end{aligned}$$

Hence code \mathcal{C}' is better than code \mathcal{C} , contradicting the fact that \mathcal{C} is optimal.

- 2) We first know that $\ell_{M-1} \leq \ell_M$, since:

- If $p_{M-1} > p_M$, then $\ell_{M-1} \leq \ell_M$ by result 1) above.
- If $p_{M-1} = p_M$, then $\ell_{M-1} \leq \ell_M$ by our assumption about the ordering of codewords for source symbols with identical probability.

Now, if $\ell_{M-1} < \ell_M$, we may delete the last digit of codeword M , and the deletion cannot result in another codeword since \mathcal{C} is a prefix code. Thus the deletion forms a new prefix code with a better average codeword length than \mathcal{C} , contradicting the fact that \mathcal{C} is optimal. Hence, we must have that $\ell_{M-1} = \ell_M$.

Huffman Codes

I: 3-62

- 3) Among the codewords of length ℓ_M , if no two codewords agree in all digits except the last, then we may delete the last digit in all such codewords to obtain a better codeword. □

Huffman Codes

I: 3-63

Lemma 3.26 (Huffman) Consider a source with alphabet $\mathcal{X} = \{a_1, \dots, a_M\}$ and symbol probabilities p_1, \dots, p_M such that

$$p_1 \geq p_2 \geq \dots \geq p_M.$$

Consider the *reduced source* alphabet \mathcal{Y} obtained from \mathcal{X} by combining the two least likely source symbols a_{M-1} and a_M into an equivalent symbol $a_{M-1,M}$ with probability $p_{M-1} + p_M$. Suppose that \mathcal{C}' , given by $f' : \mathcal{Y} \rightarrow \{0, 1\}^*$, is an optimal code for the reduced source \mathcal{Y} . We now construct a code \mathcal{C} , $f : \mathcal{X} \rightarrow \{0, 1\}^*$, for the original source \mathcal{X} as follows:

- The codewords for symbols a_1, a_2, \dots, a_{M-2} are exactly the same as the corresponding codewords in \mathcal{C}' :

$$f(a_1) = f'(a_1), f(a_2) = f'(a_2), \dots, f(a_{M-2}) = f'(a_{M-2}).$$

- The codewords associated with symbols a_{M-1} and a_M are formed by appending a “0” and a “1”, respectively, to the codeword $f'(a_{M-1,M})$ associated with the letter $a_{M-1,M}$ in \mathcal{C}' :

$$f(a_{M-1}) = [f'(a_{M-1,M})0] \quad \text{and} \quad f(a_M) = [f'(a_{M-1,M})1].$$

Then code \mathcal{C} is optimal for the original source \mathcal{X} .

Huffman Codes

I: 3-64

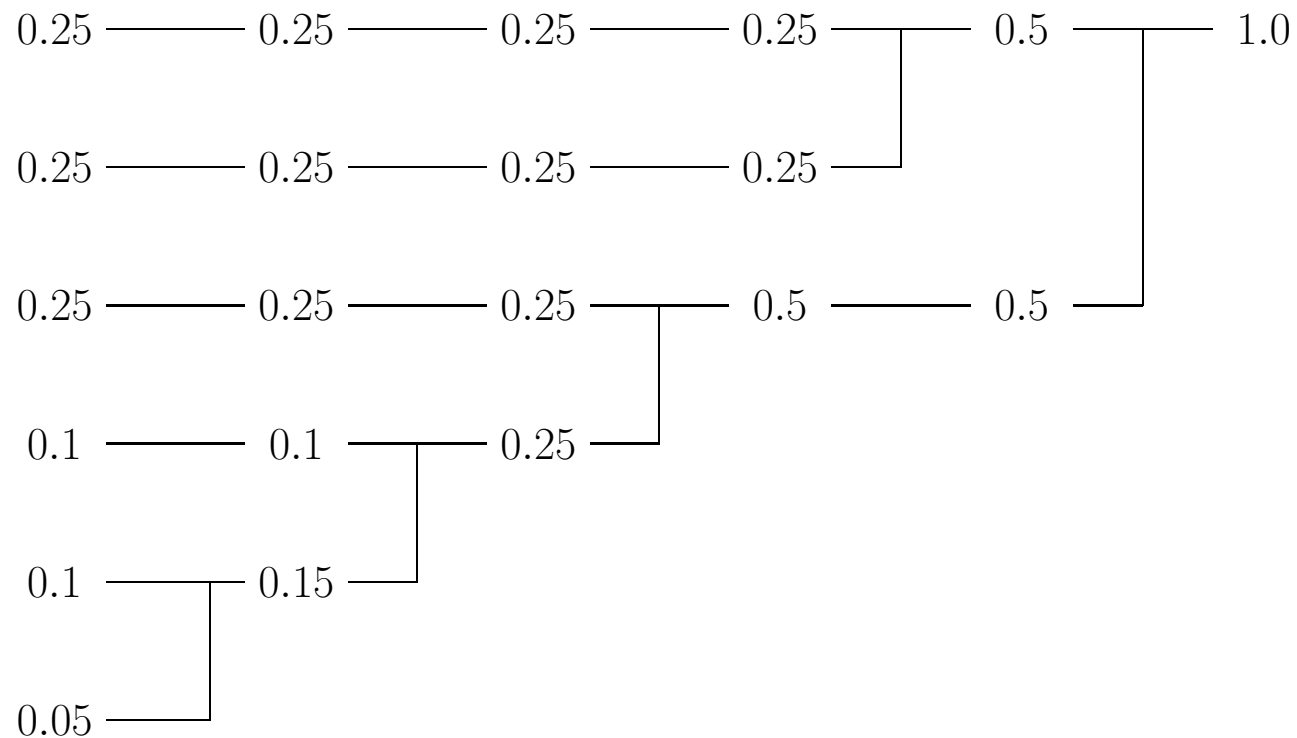
Huffman encoding algorithm:

1. Repeatedly apply the above lemma until one is left with a reduced source with two symbols. An optimal binary prefix code for this source consists of the codewords 0 and 1.
2. Then proceed backwards, constructing (as outlined in the above lemma) optimal codes for each reduced source until one arrives at the original source.

Huffman Codes

Example 3.27 Consider a source with alphabet $\{1, 2, 3, 4, 5, 6\}$ and symbol probabilities 0.25, 0.25, 0.25, 0.1, 0.1 and 0.05, respectively.

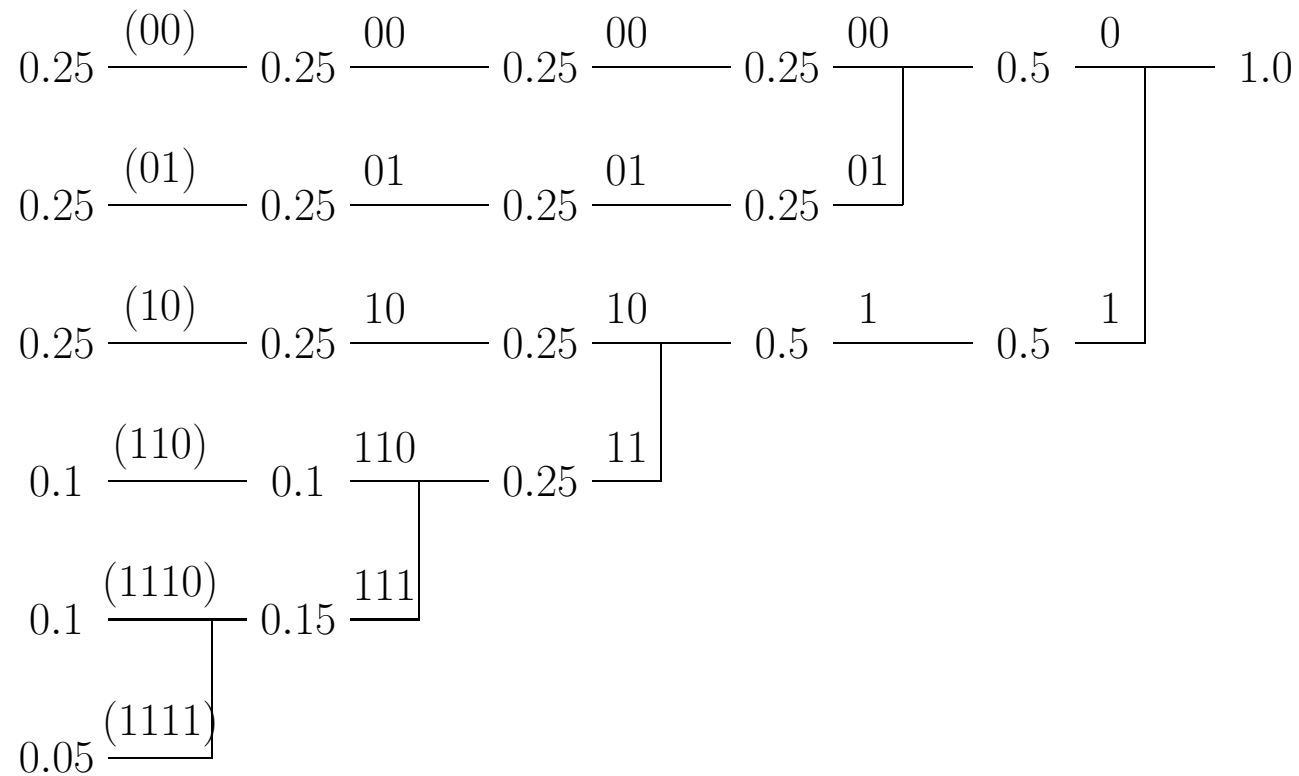
• Step 1:



Huffman Codes

I: 3-66

- Step 2:



By following the Huffman encoding procedure as shown in above figures, we obtain the Huffman code as

00, 01, 10, 110, 1110, 1111.

Huffman Codes

I: 3-67

Observation 3.28

- Huffman codes are not unique for a given source distribution;
 - e.g., by inverting all the code bits of a Huffman code, one gets another Huffman code,
 - or by resolving ties in different ways in the Huffman algorithm, one also obtains different Huffman codes.
- One can obtain optimal codes that are not Huffman codes;
 - e.g., by interchanging two codewords of the same length of a Huffman code, one can get another non-Huffman (but optimal) code.
 - Furthermore, one can construct an optimal *suffix* code (i.e., a code in which no codeword can be a suffix of another codeword) from a Huffman code by reversing the Huffman codewords.
 - Binary Huffman codes always satisfy the Kraft inequality with equality (their code tree is “saturated”).

Huffman Codes

I: 3-68

- Any n -th order binary Huffman code $f : \mathcal{X}^n \rightarrow \{0, 1\}^*$ for a stationary source $\{X_n\}_{n=1}^\infty$ with finite alphabet \mathcal{X} satisfies:

$$H(\mathcal{X}) \leq \frac{1}{n}H(X^n) \leq \bar{R}_n < \frac{1}{n}H(X^n) + \frac{1}{n}.$$

Thus, as n increases to infinity, $\bar{R}_n \rightarrow H(\mathcal{X})$ but the complexity as well as encoding-decoding delay grows exponentially with n .

- Finally, note that *non-binary* (i.e., for $D > 2$) Huffman codes can also be constructed in a mostly similar way as for the case of binary Huffman codes by designing a D -ary tree and iteratively applying Lemma 3.26, where now the D least likely source symbols are combined at each stage.
 - The only difference from the case of binary Huffman codes is that we have to ensure that we are ultimately left with D symbols at the last stage of the algorithm to guarantee the code's optimality.
 - This is remedied by expanding the original source alphabet \mathcal{X} by adding “dummy” symbols (each with zero probability) so that the alphabet size of the expanded source $|\mathcal{X}'|$ is the smallest positive integer greater than or equal to $|\mathcal{X}|$ with

$$|\mathcal{X}'| = 1 \pmod{D - 1}.$$

Huffman Codes

I: 3-69

– For example, if $|\mathcal{X}| = 6$ and $D = 3$ (ternary codes), we obtain that

$$|\mathcal{X}'| = 7,$$

meaning that we need to enlarge the original source \mathcal{X} by adding one dummy (zero-probability) source symbol.

Shannon-Fano-Elias Codes

I: 3-70

Assume $\mathcal{X} = \{1, \dots, M\}$ and $P_X(x) > 0$ for all $x \in \mathcal{X}$.

Note that it is not guaranteed that

$$P_X(1) \geq P_X(2) \geq \dots \geq P_X(M).$$

Define

$$F(x) \triangleq \sum_{a \leq x} P_X(a),$$

and

$$\bar{F}(x) \triangleq \sum_{a < x} P_X(a) + \frac{1}{2}P_X(x).$$

Encoder: For any $x \in \mathcal{X}$, express $\bar{F}(x)$ in decimal binary form, say

$$\bar{F}(x) = .c_1c_2 \dots c_k \dots,$$

and take the first k (fractional) bits as the codeword of source symbol x , i.e.,

$$(c_1, c_2, \dots, c_k),$$

where

$$k \triangleq \lceil \log_2(1/P_X(x)) \rceil + 1.$$

Shannon-Fano-Elias Codes

I: 3-71

Decoder: Given codeword (c_1, \dots, c_k) , compute the cumulative sum of $F(\cdot)$ starting from the smallest element in $\{1, 2, \dots, M\}$ until the first x satisfying

$$F(x) \geq .c_1 \dots c_k.$$

Then x should be the original source symbol.

Shannon-Fano-Elias Codes

I: 3-72

Proof of decodability: For any number $a \in [0, 1]$, let $[a]_k$ denote the operation that chops the binary representation of a after k bits (i.e., removing the $(k + 1)^{\text{th}}$ bit, the $(k + 2)^{\text{th}}$ bit, etc). Then

$$\bar{F}(x) - [\bar{F}(x)]_k < \frac{1}{2^k}.$$

Since $k = \lceil \log_2(1/P_X(x)) \rceil + 1$,

$$\begin{aligned} \frac{1}{2^k} &\leq \frac{1}{2} P_X(x) \\ &= \left[\sum_{a < x} P_X(a) + \frac{P_X(x)}{2} \right] - \sum_{a \leq x-1} P_X(a) \\ &= \bar{F}(x) - F(x-1). \end{aligned}$$

Hence,

$$F(x-1) = \left[F(x-1) + \frac{1}{2^k} \right] - \frac{1}{2^k} \leq \bar{F}(x) - \frac{1}{2^k} < [\bar{F}(x)]_k.$$

In addition,

$$F(x) > \bar{F}(x) \geq [\bar{F}(x)]_k.$$

Consequently, x is the first element satisfying

$$F(x) \geq .c_1 c_2 \dots c_k.$$

Shannon-Fano-Elias Codes

I: 3-73

Average codeword length:

$$\begin{aligned}\bar{\ell} &= \sum_{x \in \mathcal{X}} P_X(x) \left\lceil \log_2 \frac{1}{P_X(x)} \right\rceil + 1 \\ &< \sum_{x \in \mathcal{X}} P_X(x) \log_2 \frac{1}{P_X(x)} + 2 \\ &= (H(X) + 2) \text{ bits.}\end{aligned}$$

Observation 3.29 The Shannon-Fano-Elias code is a prefix code.

Shannon-Fano-Elias Codes

I: 3-74

Appendix

- *Fano code*: The *Fano code* is generated according to the following algorithm:
 1. Arrange the symbols in order of nonincreasing probability.
 2. Divide the list of ordered symbols into two parts, with the total probability of the left part being as close to the total probability of the right part as possible.
 3. Assign the binary digit 0 to the left part of the list, and the digit 1 to the right part.
 4. Recursively apply Step 2 and Step 3 to each of the two parts, subdividing into further parts and adding bits to the codewords until each symbol is the single member of a part.

Note that effectively this algorithm constructs a tree. Hence, the Fano code is a prefix code.

Shannon-Fano-Elias Codes

I: 3-75

Exemplified construction of Fano code for alphabet of size 5 and probability distributions 0.35, 0.25, 0.15, 0.15 and 0.1, respectively.

p_1	p_2	p_3	p_4	p_5
0.35	0.25	0.15	0.15	0.1
0.6		0.4		
0		1		
0.35	0.25	0.15	0.15	0.1
0	1	0.15	0.25	
			1	
			0.15	0.1
			0	1
00	01	10	110	111

Remarks

- In the literature, the Fano code is usually known as the *Shannon–Fano code*, even though it is an invention of Professor Robert Fano from MIT and not of Shannon. (Another example for such mis-naming is Stein’s Lemma, which is not the invention of Prof. Charles M. Stein.)

Shannon-Fano-Elias Codes

I: 3-76

– To make things even worse, there exists another code that is also known as *Shannon-Fano code*, but actually should be called *Shannon code* because it was proposed by Shannon.

* Under the premise that the symbols must be in order of nonincreasing probability, this *Shannon code* is exactly the Shannon-Fano-Elias code just introduced with

$$k \triangleq \lceil \log_2(1/P_X(x)) \rceil.$$

Note that without symbols being arranged in order of nonincreasing probability, a larger k should be used:

$$k \triangleq \lceil \log_2(1/P_X(x)) \rceil + 1.$$

The advantage of taking larger k is that the symbols are only required to be ordered lexicographically.

– The Kraft Inequality is still satisfied for one less k :

$$\sum_{x \in \mathcal{X}} 2^{-\lceil \log_2 \frac{1}{P_X(x)} \rceil} \leq \sum_{x \in \mathcal{X}} 2^{-\log_2 \frac{1}{P_X(x)}} = \sum_{i=1}^M P_X(x) = 1.$$

So the existence of such prefix code (with k being one less) is priori known. Shannon code substantiates this prognosis.

Shannon-Fano-Elias Codes

I: 3-77

- Shannon code performs similarly to the Fano code, but Fano code is in general slightly better.
- The idea of Shannon-Fano-Elias code has been additionally credited to the late Professor Peter Elias from MIT (hence the name *Shannon-Fano-Elias coding*), but actually Elias denied this. The concept has probably come from Shannon himself during a talk that he gave at MIT.

Universal Lossless Variable-Length Codes

I: 3-78

- The Huffman codes and Shannon-Fano-Elias codes can be constructed only when the source statistics is known.
- If the source statistics is unknown, is it possible to find a code whose average codeword length is arbitrarily close to entropy? Yes, if “asymptotic achievability” is allowed.

Adaptive Huffman Codes

I: 3-79

- Let the source alphabet be $\mathcal{X} \triangleq \{a_1, \dots, a_M\}$.

- Define

$$N(a_i|x^n) \triangleq \text{number of } a_i \text{ appearances in } x_1, x_2, \dots, x_n.$$

- Then the (current) relative frequency of a_i is

$$\frac{N(a_i|x^n)}{n}.$$

- Let $\mathbf{c}_n(a_i)$ denote the Huffman codeword of source symbol a_i with respect to distribution

$$\left[\frac{N(a_1|x^n)}{n}, \frac{N(a_2|x^n)}{n}, \dots, \frac{N(a_J|x^n)}{n} \right].$$

Adaptive Huffman Codes

I: 3-80

Now suppose that $x_{n+1} = a_j$.

1. The codeword $\mathbf{c}_n(a_j)$ is outputted.
2. Update the relative frequency for each source outcome according to:

$$\frac{N(a_j|x^{n+1})}{n+1} = \frac{n \times [N(a_j|x^n)/n] + 1}{n+1}$$

and

$$\frac{N(a_i|x^{n+1})}{n+1} = \frac{n \times [N(a_i|x^n)/n]}{n+1} \quad \text{for } i \neq j.$$

Adaptive Huffman Codes

I: 3-81

Definition 3.30 (Sibling property) A prefix code is said to have the *sibling property* if its codetree satisfies:

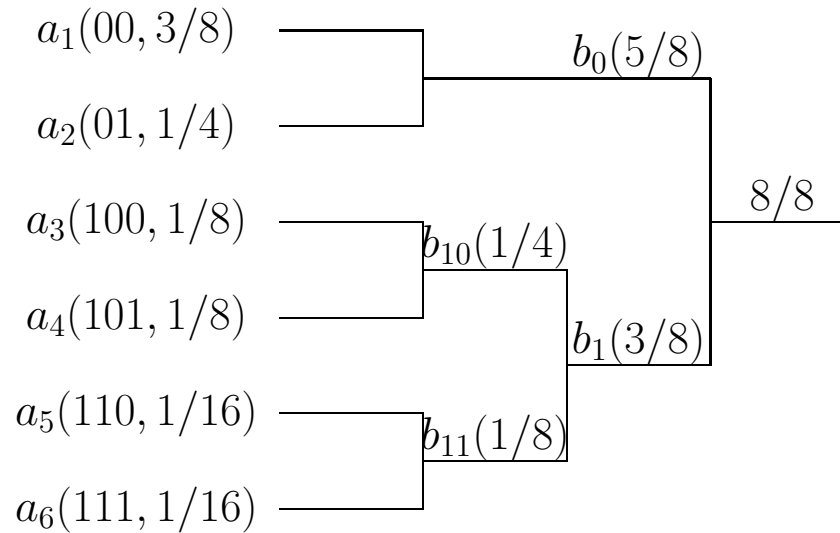
1. every node in the code-tree (except for the root node) has a sibling (i.e., the code-tree is saturated), and
2. the node can be listed in non-decreasing order of probabilities with each node being adjacent to its sibling.

Observation 3.31 A prefix code is a Huffman code if, and only if, it satisfies the sibling property.

Adaptive Huffman Codes

I: 3-82

E.g.



$$\begin{aligned}
 & \underbrace{b_0 \left(\frac{5}{8} \right) \geq b_1 \left(\frac{3}{8} \right)}_{\text{sibling pair}} \geq \underbrace{a_1 \left(\frac{3}{8} \right) \geq a_2 \left(\frac{1}{4} \right)}_{\text{sibling pair}} \\
 & \geq \underbrace{b_{10} \left(\frac{1}{4} \right) \geq b_{11} \left(\frac{1}{8} \right)}_{\text{sibling pair}} \geq \underbrace{a_3 \left(\frac{1}{8} \right) \geq a_4 \left(\frac{1}{8} \right)}_{\text{sibling pair}} \geq \underbrace{a_5 \left(\frac{1}{16} \right) \geq a_6 \left(\frac{1}{16} \right)}_{\text{sibling pair}}
 \end{aligned}$$

Adaptive Huffman Codes

I: 3-83

E.g. (cont.)

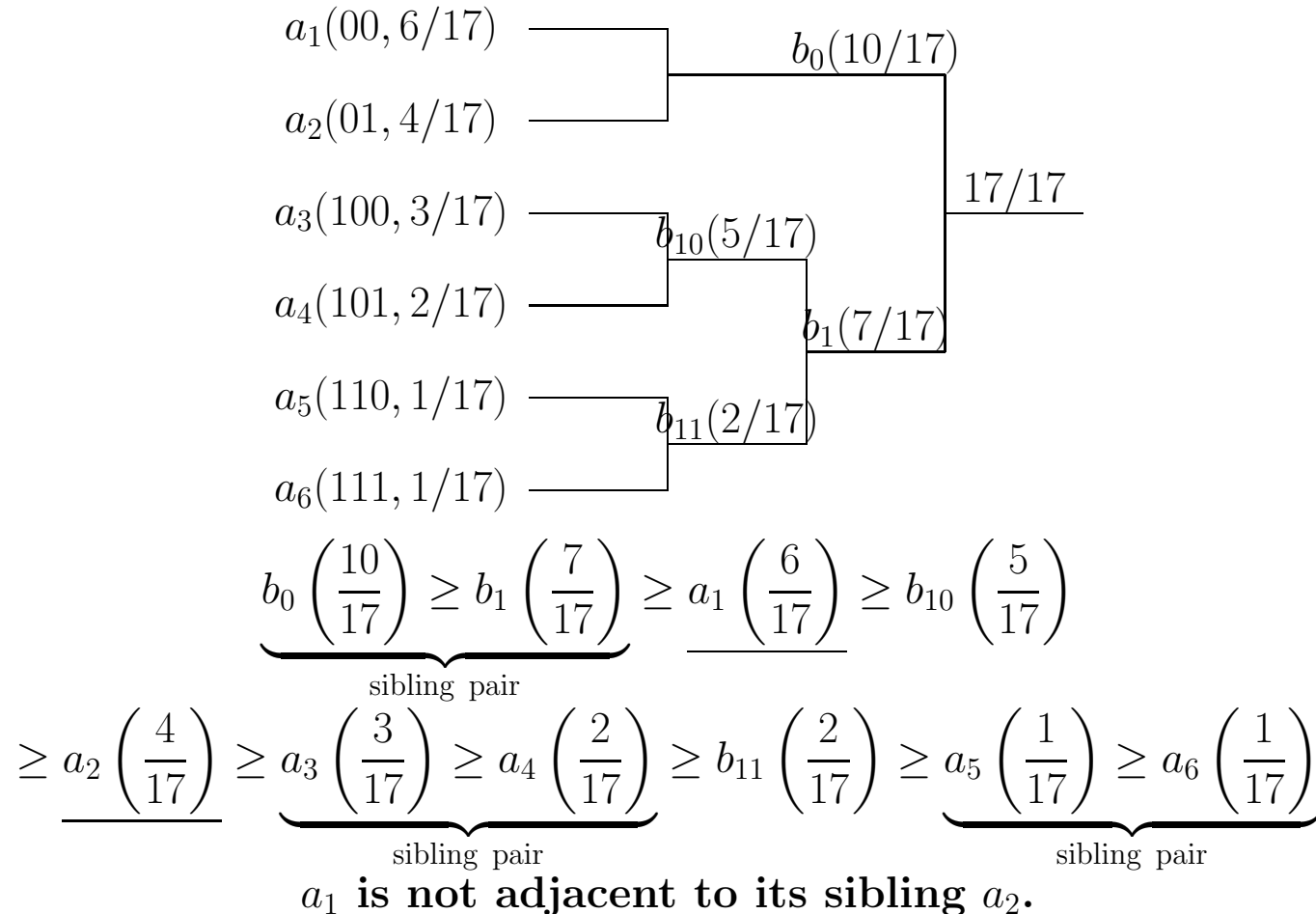
- If the next observation ($n = 17$) is a_3 , then its codeword 100 is outputted.
- The estimated distribution is updated as

$$\begin{aligned}P_{\hat{X}}^{(17)}(a_1) &= \frac{16 \times (3/8)}{17} = \frac{6}{17}, & P_{\hat{X}}^{(17)}(a_2) &= \frac{16 \times (1/4)}{17} = \frac{4}{17} \\P_{\hat{X}}^{(17)}(a_3) &= \frac{16 \times (1/8) + 1}{17} = \frac{3}{17}, & P_{\hat{X}}^{(17)}(a_4) &= \frac{16 \times (1/8)}{17} = \frac{2}{17} \\P_{\hat{X}}^{(17)}(a_5) &= \frac{16 \times [1/(16)]}{17} = \frac{1}{17}, & P_{\hat{X}}^{(17)}(a_6) &= \frac{16 \times [1/(16)]}{17} = \frac{1}{17}.\end{aligned}$$

The sibling property is no longer true; hence, the Huffman codetree needs to be updated.

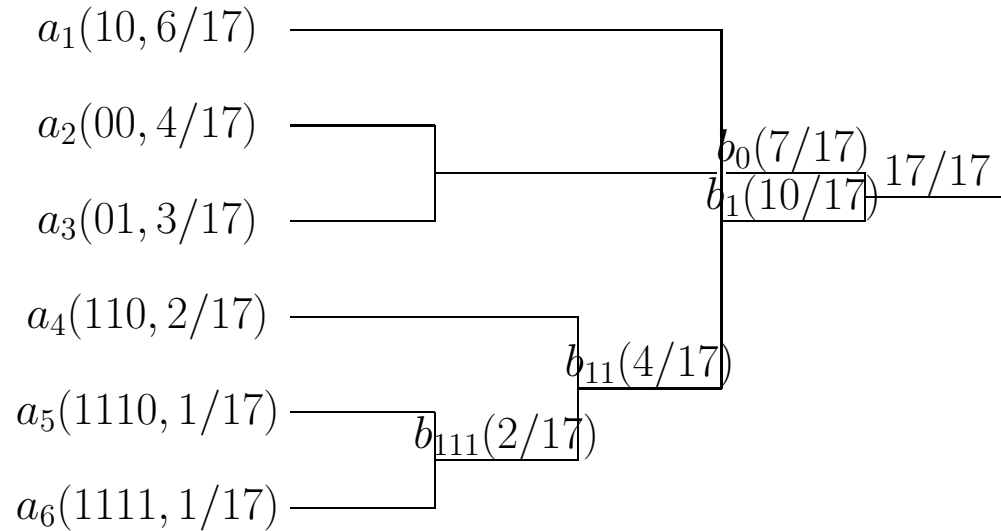
Adaptive Huffman Codes

I: 3-84



Adaptive Huffman Codes

E.g. (cont.) The updated Huffman codetree.



$$\begin{aligned}
 & \underbrace{b_1 \left(\frac{10}{17} \right) \geq b_0 \left(\frac{7}{17} \right)}_{\text{sibling pair}} \geq \underbrace{a_1 \left(\frac{6}{17} \right) \geq b_{11} \left(\frac{4}{17} \right)}_{\text{sibling pair}} \\
 \geq & \underbrace{a_2 \left(\frac{4}{17} \right) \geq a_3 \left(\frac{3}{17} \right)}_{\text{sibling pair}} \geq \underbrace{a_4 \left(\frac{2}{17} \right) \geq b_{111} \left(\frac{2}{17} \right)}_{\text{sibling pair}} \geq \underbrace{a_5 \left(\frac{1}{17} \right) \geq a_6 \left(\frac{1}{17} \right)}_{\text{sibling pair}}
 \end{aligned}$$

Lempel-Ziv Codes

I: 3-86

Encoder:

1. Parse the input sequence into strings that have never appeared before.
2. Let L be the number of distinct strings of the parsed source. Then we need $\log_2 L$ bits to index these strings (starting from one). The codeword of each string is the index of its prefix concatenated with the last bit in its source string.

E.g.

- The input sequence is 1011010100010;
- Step 1:
 - The algorithm first eats the first letter 1 and finds that it never appears before. So 1 is the *first string*.
 - Then the algorithm eats the second letter 0 and finds that it never appears before, and hence, put it to be the *next string*.
 - The algorithm eats the next letter 1, and finds that this string has appeared. Hence, it eats another letter 1 and yields a new string 11.
 - By repeating these procedures, the source sequence is parsed into strings as

1, 0, 11, 01, 010, 00, 10.

Lempel-Ziv Codes

I: 3-87

- Step 2:

- $L = 8$. So the indices will be:

parsed source : 1 0 11 01 010 00 10
index : 001 010 011 100 101 110 111 ·

- E.g., the codeword of source string 010 will be the index of 01, i.e. 100, concatenated with the last bit of the source string, i.e. 0.

- The resultant codeword string is:

$(000, 1)(000, 0)(001, 1)(010, 1)(100, 0)(010, 0)(001, 0)$

or equivalently,

0001000000110101100001000010.

Theorem 3.32 The Lempel-Ziv algorithm asymptotically achieves the entropy rate of any stationary ergodic source (with unknown statistics).

Notes on Lempel-Ziv Codes

I: 3-88

- An example of variations of Lempel-Ziv algorithm is the “compress” Unix command.
- The conventional Lempel-Ziv encoder requires two passes: the first pass to decide L , and the second pass to generate real codewords.
- The algorithm can be modified so that it requires only one pass over the source string.
- Also note that the above algorithm uses an equal number of bits— $\log_2 L$ —to all the location index, which can also be relaxed by proper modifications.

Key Notes

I: 3-89

- Average per-source-symbol codeword length versus per-source-symbol entropy
 - Average per-source-symbol codeword length is exactly the code rate for fixed-length codes.
- Categories of codes
 - Fixed-length codes (and their relation with segmentation or blocking)
 - * Block codes
 - * Fixed-length tree codes
 - Variable-length codes
 - * Non-singular codes
 - * Uniquely decodable codes
 - * Prefix codes
- AEP theorem
- Weakly δ -typical set and Shannon-McMillan theorem
- Shannon's source coding theorem and its converse theorem for DMS

Key Notes

I: 3-90

- Entropy rate and the proof of its existence for stationary sources
- Generalized AEP
- Shannon's source coding theorem and its converse theorem for stationary-ergodic sources
- Redundancy of sources
- Kraft inequality and its relation to uniquely decodable codes, as well as prefix codes
- Source coding theorem for variable-length codes
- Huffman codes and adaptive Huffman codes
- Shannon-Fano-Elias codes
- Lempel-Ziv codes