

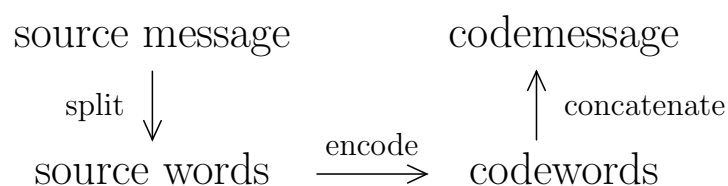
Interlude — Terminology for Coding Theory

An *alphabet* is a finite set of symbols, most commonly either $\{0, 1\}$ or the ASCII character set. A *word* over an alphabet A is a finite string of symbols from A (repetitions allowed). The *length* $l(w)$ of a word w is the number of symbols in it (including repetitions: so $l(00101) = 5$). There are $|A|^n$ different words of length n over A , since there are $|A|$ choices for the first symbol, for each of these there are $|A|$ choice for the second symbol, etc.

A *code* over an alphabet A is a finite set C of distinct words over A , called *codewords*. The *codemessage* corresponding to a sequence of codewords is the string obtained by concatenating these codewords; e.g., if $C = \{0, 010, 011\}$ then 00110110100 is the codemessage corresponding to the sequence $c_1c_3c_3c_2c_1$ of codewords, but 111 is *not* a codemessage.

To *use* a code in practice, one requires also:

- (1) a *source alphabet* A_s ;
- (2) a set of permissible *source messages* (usually, all possible words over A_s);
- (3) a finite set of distinct words over A_s , called *source words* (often, all possible words over A_s of some fixed length);
- (4) a procedure for dividing up an arbitrary source message into source words; and
- (5) a bijection from the set of source words to the set of codewords:



However, in much of coding theory the source messages/words are irrelevant and will be ignored.

Chapter 2 Coding for Data-Compression

2.1. Instantaneous and uniquely decipherable codes.

A word w_1 is a *prefix* of a word w_2 if w_1 occurs at the start of w_2 ; e.g., 011 is a prefix of 01101. A code C is *instantaneous* if no codeword is a prefix of any other codeword. C is *uniquely decipherable* if each codemessage corresponds to exactly one sequence of codewords. Instantaneous codes are uniquely decipherable:

Example 1. $C_1 = \{0, 11, 100, 101\}$.

C_1 is instantaneous. To decipher a codemessage, read from the left until a codeword is completed, then iterate with what is left. E.g.,

$$0 \mid 11 \mid 0 \mid 0 \mid 100 \mid 101 \mid 0 \mid 0 .$$

Example 2. $C_2 = \{0, 010, 011\}$.

C_2 is not instantaneous (e.g., 0 is a prefix of 010), but it *is* uniquely decipherable: in a codemessage, 111 is impossible, 11 can occur only as part of 011, and any other 1's must occur as part of 010; what is left is all 0's. E.g.,

$$\underline{0} \mid \underline{011} \mid \underline{011} \mid \underline{010} \mid \underline{011} \mid \underline{010} \mid \underline{0} \mid \underline{0} \mid \underline{010} \mid \underline{011} \mid \underline{0} .$$

Example 3. $C_3 = \{00, 01, 10, 010\}$.

C_3 is neither instantaneous nor uniquely decipherable. For example, the codemessage 01000010 has the two interpretations

$$\begin{array}{c} 01 \mid 00 \mid 00 \mid 10 \\ \text{and } 010 \mid 00 \mid 010 . \end{array}$$

Theorem 2.1. Let A be an alphabet of α symbols, and let m and $l_1 \leq \dots \leq l_m$ be positive integers. Then the following three statements are equivalent.

- (a) There exists an instantaneous code C over A with codewords c_1, \dots, c_m of lengths l_1, \dots, l_m . (L. G. Kraft, 1949)

[MS dissertation, MIT.]

- (b) There exists a uniquely decipherable code C' over A with codewords c'_1, \dots, c'_m of lengths l_1, \dots, l_m . (B. McMillan, 1956)

[b. 1915; Bell Labs.]

- (c) $\sum_{i=1}^m \alpha^{-l_i} \leq 1$.

Proof. (a) \implies (b) because every instantaneous code is uniquely decipherable.

(c) \implies (a): Choose c_1 to be an arbitrary word over A of length l_1 . For $2 \leq s \leq m$, suppose we have chosen c_1, \dots, c_{s-1} of lengths l_1, \dots, l_{s-1} so that none is a prefix of any other. We must now choose c_s .

There are α^{l_s} possible words of length l_s , but $\alpha^{l_s-l_i}$ of them have c_i as a prefix, for each i ($1 \leq i \leq s-1$). However,

$$\sum_{i=1}^{s-1} \alpha^{l_s-l_i} < \alpha^{l_s},$$

since

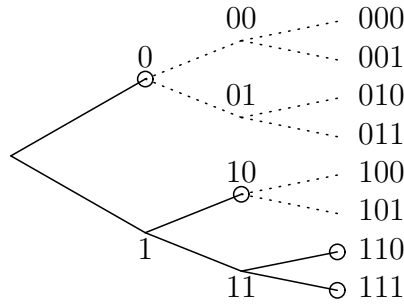
$$\sum_{i=1}^{s-1} \alpha^{-l_i} < \sum_{i=1}^m \alpha^{-l_i} \leq 1.$$

Thus there exists a word of length l_s that does not have any of c_1, \dots, c_{s-1} as a prefix. Choose c_s to be any such word. Then, for $1 \leq i \leq s-1$, c_i is not a prefix of c_s , and c_s is not a prefix of c_i since $l_s \geq l_i$. Continue until $s = m$.

One way of visualizing this proof is on a tree. Suppose, for example, $\alpha = 2$, $m = 4$, $(l_1, l_2, l_3, l_4) = (1, 2, 3, 3)$. Then

$$\sum_{i=1}^m \alpha^{-l_i} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{8} = 1,$$

and so there is an instantaneous code with these word lengths.



(b) \implies (c) (J. Karush, 1961): Write $x := \alpha^{-1}$ and let r be a positive integer. Then

$$\begin{aligned} \left(\sum_{i=1}^m \alpha^{-l_i} \right)^r &= \left(\sum_{i=1}^m x^{l_i} \right)^r = \left(\sum_{w \in C'} x^{l(w)} \right)^r \\ &= \sum_{w_1 \in C'} \sum_{w_2 \in C'} \dots \sum_{w_r \in C'} x^{l(w_1)} x^{l(w_2)} \dots x^{l(w_r)}, \end{aligned}$$

since each term in the sum is obtained by choosing $x^{l(w_1)}$ from the ‘first’ bracket, $x^{l(w_2)}$ from the ‘second’, \dots , $x^{l(w_r)}$ from the ‘last’, for some $w_1, w_2, \dots, w_r \in C'$. But the summand is x^i , where i is the length of the codemessage $w_1 w_2 \dots w_r$, and so the sum equals $\sum_{i=1}^{\infty} a_i x^i$, where a_i is the number of sequences of r codewords that give codemessages of length i . Note that two such sequences cannot give the same codemessage, since C' is uniquely decipherable, and so $a_i \leq \alpha^i$, the number of different strings of i symbols from A ; and of course $a_i = 0$ if $i > r l_m$. So

$$\left(\sum_{i=1}^m \alpha^{-l_i} \right)^r = \sum_{i=1}^{\infty} a_i x^i = \sum_{i=1}^{\infty} a_i \alpha^{-i} \leq \sum_{i=1}^{r l_m} \alpha^i \alpha^{-i} = r l_m.$$

Hence

$$\sum_{i=1}^m \alpha^{-l_i} \leq (r l_m)^{1/r} \rightarrow 1 \text{ as } r \rightarrow \infty.$$

Since r can be arbitrarily large, (c) follows. //

2.2. The noiseless coding theorem for memoryless sources. Suppose we have a source S that emits words w_1, \dots, w_m

with positive probabilities p_1, \dots, p_m , each word emitted being independent of all previous words (hence, ‘memoryless’). Suppose we have an alphabet A of α symbols, a uniquely decipherable code C over A with codewords c_1, \dots, c_m of lengths l_1, \dots, l_m , and an

$$\text{encoding map } \begin{array}{l} w_1 \rightarrow c_1 \\ w_2 \rightarrow c_2 \\ \vdots \\ w_m \rightarrow c_m \end{array}.$$

The *average length* of C is $l(C) := \sum_{i=1}^m p_i l(c_i) = \sum_{i=1}^m p_i l_i$. C is *optimal* or *compact* if $l(C)$ is as small as possible subject to all these conditions. This means that, on average over all source messages, the length of the corresponding codemessage is minimal among all those produced by uniquely decipherable codes over alphabet A .

Theorem 2.2. *The noiseless coding theorem* (C. E. Shannon, 1948). If C is an optimal code for the above source S over the alphabet A , then

$$\frac{H(S)}{\log_2 \alpha} \leq l(C) < \frac{H(S)}{\log_2 \alpha} + 1,$$

where $H(S) := -\sum_{i=1}^m p_i \log_2 p_i$, the entropy or uncertainty of the source S .

Proof. Let C be *any* uniquely decipherable code for S over A , with word lengths l_1, \dots, l_m . Let $\beta := \sum_{i=1}^m \alpha^{-l_i}$ and $q_i := \frac{\alpha^{-l_i}}{\beta}$ ($i = 1, \dots, m$), so that (q_1, \dots, q_m) is a probability distribution and $\log_2 q_i = -l_i \log_2 \alpha - \log_2 \beta$. By Lemma 1.3.1,

$$\begin{aligned} H(S) &= -\sum_{i=1}^m p_i \log_2 p_i \leq -\sum_{i=1}^m p_i \log_2 q_i \\ &= \left(\sum_{i=1}^m p_i l_i \right) \log_2 \alpha + \left(\sum_{i=1}^m p_i \right) \log_2 \beta \\ &= l(C) \log_2 \alpha + \log_2 \beta. \end{aligned}$$

But $\beta = \sum_{i=1}^m \alpha^{-l_i} \leq 1$ by Theorem 2.1, so $\log_2 \beta \leq 0$ and $\frac{H(S)}{\log_2 \alpha} \leq l(C)$, as required.

To prove the upper bound, it suffices to prove that *there exists* a uniquely decipherable code satisfying this bound, since then every optimal code must satisfy it. So, for each i , choose l_i to be the *least* integer such that $\frac{-\log_2 p_i}{\log_2 \alpha} \leq l_i$,

$$\text{i.e., } -\log_2 p_i \leq l_i \log_2 \alpha, \quad \text{i.e., } p_i^{-1} \leq \alpha^{l_i}, \quad \text{i.e., } \alpha^{-l_i} \leq p_i.$$

Then

$$\frac{-\log_2 p_i}{\log_2 \alpha} \leq l_i < \frac{-\log_2 p_i}{\log_2 \alpha} + 1.$$

Also, $\sum_{i=1}^m \alpha^{-l_i} \leq \sum_{i=1}^m p_i = 1$, and so there exists a uniquely decipherable code C with these codeword lengths by Theorem 2.1.

$$\text{Also, } l(C) = \sum_{i=1}^m p_i l_i < \sum_{i=1}^m \left(\frac{-p_i \log_2 p_i}{\log_2 \alpha} + p_i \right) = \frac{H(S)}{\log_2 \alpha} + 1. \quad //$$

In the binary case ($\alpha = 2$), R. G. Gallager proved in 1978 that $l(C) \leq H(S) + \max_i p_i + \log_2 \left(\frac{2 \log_2 e}{e} \right) \approx H(S) + \max_i p_i + 0.0861$.

2.3. Huffman's algorithm. This is an algorithm for constructing *instantaneous* optimal codes. We shall consider only *binary* codes, with alphabet $\{0, 1\}$. Let S be a memoryless source that emits words w_1, \dots, w_m with probabilities $p_1 \geq \dots \geq p_m > 0$. Let S' be a memoryless source that emits words w'_1, \dots, w'_{m-1} with probabilities $p_1, \dots, p_{m-2}, p_{m-1} + p_m$. Let C' be an instantaneous code for S' over $\{0, 1\}$ with encoding map

$$\begin{array}{l} w'_1 \rightarrow c'_1 \\ \vdots \\ w'_{m-1} \rightarrow c'_{m-1} \end{array}.$$

Let C be the code

$$\begin{array}{l} w_1 \rightarrow c'_1 \\ \vdots \\ w_{m-2} \rightarrow c'_{m-2} \\ w_{m-1} \rightarrow c'_{m-1}0 \\ w_m \rightarrow c'_{m-1}1 \end{array}$$

for S . Clearly C is instantaneous.

Lemma 2.3.1. $l(C) = l(C') + p_{m-1} + p_m$.

Proof.
$$l(C) = \sum_{i=1}^{m-2} p_i l(c'_i) + (p_{m-1} + p_m)(l(c'_{m-1}) + 1)$$

$$= l(C') + p_{m-1} + p_m. \quad //$$

Theorem 2.3. C is optimal for S if and only if C' is optimal for S' .

Proof. Let C^* be an *instantaneous* optimal code for S over $\{0, 1\}$ (which exists by Theorem 2.1), with encoding map

$$\begin{array}{c} w_1 \rightarrow c_1^* \\ \vdots \\ w_m \rightarrow c_m^* \end{array}.$$

Claim 1. If $p_i > p_j$ then $l(c_i^*) \leq l(c_j^*)$.

Proof. If $l(c_i^*) > l(c_j^*)$, form a new instantaneous code \hat{C} from C^* by interchanging c_i^* and c_j^* ; then

$$\begin{aligned} l(\hat{C}) &= l(C^*) - p_i l(c_i^*) - p_j l(c_j^*) + p_j l(c_i^*) + p_i l(c_j^*) \\ &= l(C^*) - (p_i - p_j)(l(c_i^*) - l(c_j^*)) \\ &< l(C^*), \quad \Rightarrow \Leftarrow. \end{aligned}$$

Claim 2. Among the codewords of maximum length in C^* , there are two that agree in all but the last digit.

Proof. If not, then we can delete the last digit of every codeword of maximum length and still have an instantaneous code, $\Rightarrow \Leftarrow$.

By Claims 1 and 2, we may assume w.l.o.g. that c_{m-1}^* and c_m^* agree in all but the last digit, since $l(C^*)$ does not change if we interchange two longest codewords, or two codewords used with the same probability. Say $c_{m-1}^* = c^*0$ and $c_m^* = c^*1$. Let $C^{*'} be the code$

$$\begin{array}{c} w'_1 \rightarrow c_1^* \\ \vdots \\ w'_{m-2} \rightarrow c_{m-2}^* \\ w'_{m-1} \rightarrow c^* \end{array}$$

Example 2.

	S	C		S'	C'		S''	C''		S'''	C'''		$S^{(4)}$	$C^{(4)}$	
w_1	0.3	00		w_1	0.3	00	w_1	0.3	00	w_{3456}	0.45	1	w_{12}	0.55	0
w_2	0.25	01		w_2	0.25	01	w_2	0.25	01	w_1	0.3	00	w_{3456}	0.45	1
w_3	0.2	11		w_3	0.2	11	w_{456}	0.25	10	w_2	0.25	01			
w_4	0.1	101		w_{56}	0.15	100	w_3	0.2	11						
w_5	0.1	1000		w_4	0.1	101									
w_6	0.05	1001													

Compression with Huffman's algorithm is limited because:

- (1) the encoding map must be included with the message;
- (2) the *only* form of redundancy it removes is that implied by the differing frequencies of source words; e.g., if the source words are 0 and 1 (only), then the source message 010101... is highly redundant, but Huffman's algorithm will not compress it at all.

There is also a computational disadvantage, in that the message must be scanned twice, once to choose the source words and discover their frequencies, and once to encode them.

Nevertheless, effective computer programs based on Huffman's algorithm exist and typically achieve compressions of around 45–50% on long text files and 75% on source code for computer programs.

2.4. The LZW algorithm. (A. Lempel and J. Ziv, 1977; T. A. Welch, 1984.) This algorithm avoids all the above problems. The codewords are all 2^n binary words of length n , for some n . The encoder and decoder start with the same code table, containing a numerical list of symbols in the source alphabet. At each stage, the encoder encodes the longest initial string of the remaining message for which a codeword exists in the table, and then adds that string plus the next symbol as a new source word in the table (unless the table is already full). The decoder decodes the incoming code and reconstructs the same code table one word behind the encoder.

Example with source alphabet $_$ (space), a, b.

Encoder

Decoder

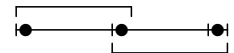
a|_|b|b|a|a|_|a|b|_|a|_a|ba|bab|_|b|a
1|0|2|2|1|3|1|2|0|8|6|13|4|1

1|0|2|2|1|3|1|2|0|8|6|13|4|1
a|_|b|b|a|a|_|a|b|_|a|_a|ba|bab|_|b|a

$_$	0	a_a	8
a	1	ab	9
b	2	b_	10
a_	3	_a	11
_b	4	a_ab	12
bb	5	bab	13
ba	6	bab_	14
aa	7	[_ba	15]

$_$	0	a_a	8
a	1	ab	9
b	2	b_	10
a_	3	_a	11
_b	4	a_ab	12
bb	5	bab	13
ba	6	bab_	14
aa	7	[_ba	15]

If the decoder receives a codeword that is not in the code table, then it must have been put into the encoder's code table at the previous word, and so corresponds to the previous source word received plus its first symbol.



The disadvantages of this method are:

- (1) It takes time for the code table to build up, so the first part of the message is expanded, not compressed. (Thus the method does not work with short messages.)
- (2) After a long time, the code table becomes full, and if the characteristics of the message change markedly after that, then compression is less effective.

Nevertheless, the algorithm achieves compressions of typically 55–60% on long text files and 85% on source code for computer programs.

A test on a long text file:

program	reduced to	compression	method
compress	41%	59%	LZW algorithm
gzip	28%	72%	Lempel–Ziv 1977
bzip2	23%	77%	Burrows–Wheeler block sorting + Huffman coding