

# Discovery of Huffman codes

Inna Pivkina\*

## 1 Introduction

The story of the invention of Huffman codes is described in an article by Gary Stix in the September 1991 issue of *Scientific American*, pp. 54, 58 ([4]). The following is excerpted from the article.

PROFILE: DAVID A. HUFFMAN  
Encoding the “Neatness” of Ones and Zeroes

Large networks of IBM computers use it. So do high-definition television, modems and a popular electronic device that takes the brain work out of programming a videocassette recorder. All these digital wonders rely on the results of a 40-year-old term paper by a modest Massachusetts Institute of Technology graduate student—a data compression scheme known as Huffman encoding.

In 1951 David A. Huffman and his classmates in an electrical engineering graduate course on information theory were given the choice of a term paper or a final exam. For the term paper, Huffman’s professor, Robert M. Fano, had assigned what at first appeared to be a simple problem. Students were asked to find the most efficient method of representing numbers, letters or other symbols using a binary code. Besides being a nimble intellectual exercise, finding such a code would enable information to be compressed for transmission over a computer network or for storage in a computer’s memory.

Huffman worked on the problem for months, developing a number of approaches, but none that he could prove to be the most efficient. Finally, he despaired of ever reaching a solution and decided to start studying for the final. Just as he was throwing his notes in the garbage, the solution came to him. “It was the most singular moment of my life,” Huffman says. “There was the absolute lightning of sudden realization.”

That epiphany added Huffman to the legion of largely anonymous engineers whose innovative thinking forms the technical underpinnings for the accoutrements of modern living - in his case, from facsimile machines to modems and a myriad of other devices. “Huffman code is one of the fundamental ideas that people in computer science and data communications are using all the time,” says Donald E. Knuth of Stanford University, who is the author of the multivolume series *The Art of Computer Programming*.

Huffman says he might never have tried his hand at the problem - much less solved it at the age of 25 - if he had known that Fano, his professor, and Claude E. Shannon, the creator of information theory, had struggled with it. “It was my luck to be there at the right time and also not have my professor discourage me by telling me that other good people had struggled with this problem,” he says.

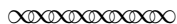
---

\*Department of Computer Science; New Mexico State University; Las Cruces, NM 88003; [ipivkina@cs.nmsu.edu](mailto:ipivkina@cs.nmsu.edu).

## 2 Shannon-Fano Coding 1944

In this section we will study work of Shannon and Fano, in particular, definition of a unit of information, how to determine the amount of information, and a recording procedure (Shannon-Fano coding) they developed for improving the efficiency of transmission by reducing the number of digits required to transmit a message. Shannon-Fano coding was developed independently by Shannon and Fano in 1944. A greedy strategy can be used to produce a Shannon-Fano coding. The strategy uses a top-down approach and does not necessarily produce the optimal code.

Let us start with the notion of a unit of information. It is defined in section I from Fano [1].



### I Definition of the Unit of Information

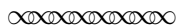
In order to define, in an appropriate and useful manner, a unit of information, we must first consider in some detail the nature of those processes in our experience which are generally recognized as conveying information. A very simple example of such processes is a yes-or-no answer to some specific question. A slightly more involved process is the indication of one object in a group of  $N$  objects, and, in general, the selection of one choice from a group of  $N$  specific choices. The word “specific” is underlined because such a qualification appears to be essential to these information-conveying processes. It means that the receiver is conscious of all possible choices, as is, of course, the transmitter (that is, the individual or the machine which is supplying the information). For instance, saying “yes” or “no” to a person who has not asked a question obviously does not convey any information. Similarly, the reception of a code number which is supposed to represent a particular message does not convey any information unless there is available a code book containing all the messages with the corresponding code numbers.

Considering next more complex processes, such as writing or speaking, we observe that these processes consist of orderly sequences of selections from a number of specific choices, namely, the letters of the alphabet or the corresponding sounds. Furthermore, there are indications that the signals transmitted by the nervous system are of a discrete rather than of a continuous nature, and might also be considered as sequences of selections. If this were the case, all information received through the senses could be analyzed in terms of selections. The above discussion indicates that the operation of selection forms the bases of a number of processes recognized as conveying information, and that it is likely to be of fundamental importance in all such processes. We may expect, therefore, that a unit of information, defined in terms of a selection, will provide a useful basis for a quantitative study of communication systems.

Considering more closely this operation of selection, we observe that different informational value is naturally attached to the selection of the same choice, depending on how likely the receiver considered the selection of that particular choice to be. For example, we would say that little information is given by the selection of a choice which the receiver was almost sure would be selected. It seems appropriate, therefore, in order to avoid difficulty at this early stage, to use in our definition the particular case of equally likely choices - that is, the case in which the receiver has no reason to expect that one choice will be selected rather than any other. In addition, our natural concept of information indicates that the information conveyed by a selection increases with the number of choices from which the selection is made, although the exact functional relation between these two quantities is not immediately clear.

On the basis of the above considerations, it seems reasonable to define as the unit of information the simplest possible selection, namely the selection between two equally likely choices, called, hereafter, the “elementary selection”. For completeness, we must add to this definition the postulate, consistent

with our intuition, that  $N$  independent selections of this type constitute  $N$  units of information. By independent selections we mean, of course, selections which do not affect one another. We shall adopt for this unit the convenient name of “bit” (from “binary digit”), suggested by Shannon. We shall also refer to a selection between two choices (not necessarily equally likely) as a “binary selection”, and to a selection from  $N$  choices, as an  $N$ -order selection. When the choices are, a priori, equally likely, we shall refer to the selection as an “equally likely selection”.



**Exercise 2.1.** Fano writes about specific questions and specific choices. Explain in your own words what he means by specific. Give an example of a specific question or choice. Explain why it is specific. Give an example of a question or choice which is not specific. Explain why it is not specific.

**Exercise 2.2.** Give an example of the simplest possible selection (or the “elementary selection”). Give an example of a selection which is not the simplest possible one.

**Exercise 2.3.** How much information does one flip of a coin convey? How much information is conveyed in  $N$  flips of a coin?

Once the unit of information is defined, it can be used to measure information conveyed by selections and messages. Fano talks about selection from several equally likely choices in section II of his work.



## II Selection from $N$ Equally Likely Choices

Consider now the selection of one among a number,  $N$ , of equally likely choices. In order to determine the amount of information corresponding to such a selection, we must reduce this more complex operation to a series of independent elementary operations. The required number of these elementary selections will be, by definition, the measure in bits of the information given by such an  $N$ -order selection.

Let us assume for the moment that  $N$  is a power of two. In addition (just to make the operation of selection more physical), let us think of the  $N$  choices as  $N$  objects arranged in a row, as indicated in Figure 1.

These  $N$  objects are first divided in two equal groups, so that the object to be selected is just as likely to be in one group as in the other. Then the indication of the group containing the desired object is equivalent to one elementary selection, and, therefore, to one bit. The next step consists of dividing each group into two equal subgroups, so that the object to be selected is again just as likely to be in either subgroup. Then one additional elementary selection, that is a total of two elementary selections, will suffice to indicate the desired subgroup (of the possible four subgroups). This process of successive subdivisions and corresponding elementary selections is carried out until the desired object is isolated from the others. Two subdivisions are required for  $N = 4$ , three from  $N = 8$ , and, in general, a number of subdivisions equal to  $\log_2 N$ , in the case of an  $N$ -order selection.

The same process can be carried out in a purely mathematical form by assigning order numbers from 0 to  $N - 1$  to the  $N$  choices. The numbers are then expressed in the binary system, as shown in Figure 1, the number of binary digits (0 or 1) required being equal to  $\log_2 N$ . These digits represent

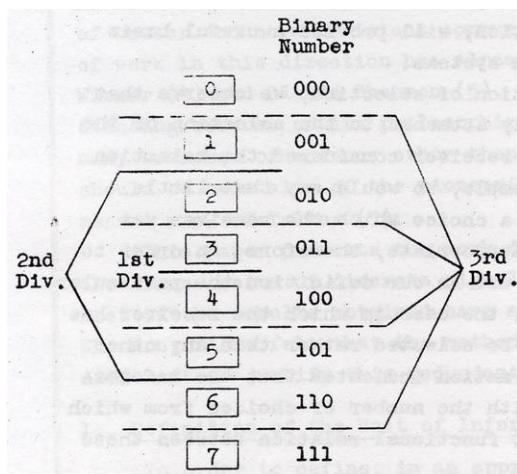


Fig. 1 Selection procedure for equally likely choices.

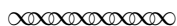
an equal number of elementary selections and, moreover, correspond in order to the successive divisions mentioned above. In conclusion, an  $N$ -order, equally likely selection conveys an amount of information

$$H_N = \log_2 N \quad (1)$$

The above result is strictly correct only if  $N$  is a power of two, in which case  $H_N$  is an integer. If  $N$  is not a power of two, then the number of elementary selections required to specify the desired choice will be equal to the logarithm of either the next lower or the next higher power of two, depending on the particular choice selected. Consider, for instance, the case of  $N = 3$ . The three choices, expressed as binary numbers, are then

$$00; 01; 10.$$

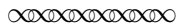
If the binary digits are read in order from left to right, it is clear that the first two numbers require two binary selections - that is, two digits, while the third number requires only the first digit, 1, in order to be distinguished from the other two. In other words, the number of elementary selections required when  $N$  is not a power of two is equal to either one of the two integers closest to  $\log_2 N$ . It follows that the corresponding amount of information must lie between these two limits, although the significance of a non-integral value of  $H$  is not clear at this point. It will be shown in the next section that Eq.(1) is still correct when  $N$  is not a power of two, provided  $H_N$  is considered as an average value over a large number of selections.



**Exercise 2.4.** Assume that we select one letter from the English alphabet at random. How much information does this selection convey?

**Exercise 2.5.** Assume that we have the following equally likely choices expressed as binary numbers: 000, 001, 010, 100, 101, and 110. If the binary digits are read in order from left to right, how many selections does each of the numbers require? Will your answer change if the digits are read from right to left?

Now that we know the amount of information conveyed by one selection, how can we determine the amount of information contained in a sequence of selections, or in a message? It is explained in the following quote from section III from Fano [1].



### III Messages and Average Amount of Information

We have determined in the preceding section the amount of information conveyed by a single selection from  $N$  equally likely choices. In general, however, we have to deal with not one but long series of such selections, which we call messages. This is the case, for instance, in the transmission of written intelligence. Another example is provided by the communication system known as pulse-code modulation, in which audio waves are sampled at equal time intervals and then each sample is quantized, that is approximated by the closest of a number  $N$  of amplitude levels.

Let us consider, then, a message consisting of a sequence of  $n$  successive  $N$ -order selections. We shall assume, at first, that these selections are independent and equally likely. In this simpler case, all the different sequences which can be formed equal in number to

$$S = N^n, \quad (2)$$

are equally likely to occur. For instance, in the case of  $N = 2$  (the two choices being represented by numbers 0 and 1) and  $n = 3$ , the possible sequences would be 000, 001, 010, 100, 011, 101, 110, 111. The total number of these sequences is  $S = 8$  and the probability of each sequence is  $1/8$ . In general, therefore, the ensemble of the possible sequences may be considered as forming a set of  $S$  equally likely choices, with the result that the selection of any particular sequence yields an amount of information

$$H_S = \log_2 S = n \log_2 N. \quad (3)$$

In words,  $n$  independent equally likely selections give  $n$  times as much information as a single selection of the same type. This result is certainly not surprising, since it is just a generalization of the postulate, stated in Section II, which forms an integral part of the definition of information.

It is often more convenient, in dealing with long messages, to use a quantity representing the average amount of information per  $N$ -order selection, rather than the total information corresponding to the whole message. We define this quantity in the most general case as the total information conveyed by a very long message divided by the number of selections in the message, and we shall indicate it with the symbol  $H_N$ , where  $N$  is the order of each selection. It is clear that when all the selections in the message are equally likely and independent and, in addition,  $N$  is a power of two, the quantity  $H_N$  is just equal to the information actually given by each selection, that is

$$H_N = \frac{1}{n} \log_2 S = \log_2 N. \quad (4)$$

We shall show now that this equation is correct also when  $N$  is not a power of two, in which case  $H_N$  has to be actually an average value taken over a sufficiently long sequence of selections.

The number  $S$  of different and equally likely sequences which can be formed with  $n$  independent and equally likely selections is still given by Eq.(2), even when  $N$  is not a power of two. On the contrary, the number of elementary selections required to specify any one particular sequence must be written now in the form

$$B_S = \log_2 S + d, \quad (5)$$

where  $d$  is a number, smaller in magnitude than unity, which makes  $B_S$  an integer and which depends on the particular sequence selected. The average amount of information per  $N$ -order selection is then, by definition,

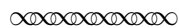
$$H_N = \lim_{n \rightarrow \infty} \frac{1}{n} (\log_2 S + d). \quad (6)$$

Since  $N$  is a constant and since the magnitude of  $d$  is smaller than unity while  $n$  approaches infinity, this equation together with Eq.(2) yields

$$H_N = \log_2 N. \tag{7}$$

We shall consider now the more complex case in which the selections, although still independent, are not equally likely. In this case, too, we wish to compute the average amount of information per selection. For this purpose, we consider again the ensemble of all the messages consisting of independent selections and we look for a way of indicating any one particular message by means of elementary selections. If we were to proceed as before, and divide the ensemble of messages in two equal groups, the selection of the group containing the desired message would no longer be a selection between equally likely choices, since the sequences themselves are not equally likely. The proper procedure is now, of course, to make equal for each group not the number of messages in it but the probability of its containing the desired message. Then the selection of the desired group will be a selection between equally likely choices. This procedure of division and selection is repeated over and over again until the desired message has been separated from the others. The successive selections of groups and subgroups will then form a sequence of independent elementary selections.

One may observe, however, that it will not generally be possible to form groups equally likely to contain the desired message, because shifting any one of the messages from one group to the other will change, by finite amounts, the probabilities corresponding to the two groups. On the other hand, if the length of the messages is increased indefinitely, the accuracy with which the probabilities of the two groups can be made equal becomes better and better since the probability of each individual message approaches zero. Even so, when the resulting subgroups include only a few messages after a large number of divisions, it may become impossible to keep the probabilities of such subgroups as closely equal as desired unless we proceed from the beginning in an appropriate manner as indicated below. The messages are first arranged in order of their probabilities, which can be easily computed if the probabilities of the choices are known. The divisions in groups and subgroups are then made successively without changing the order of the messages, as illustrated in Figure 2. In this manner, the smaller subgroups will contain messages with equal or almost equal probabilities, so that further subdivisions can be performed satisfactory. It is clear that when the above procedure is followed, the number of binary selections required to separate any message from the others varies from message to message. Messages with a high probability of being selected require less binary selections than those with lower probabilities. This fact is in agreement with the intuitive notion that the selection of a little-probable message conveys more information than the selection of a more-probable one. Certainly, the occurrence of an event which we know a priori to have a 99 per cent probability is hardly surprising or, in our terminology, yields very little information, while the occurrence of an event which has a probability of only 1 per cent yields considerably more information.



In the above, Figure 2 illustrates a method which can be used to code messages. Assume that we are given an ensemble of messages and probabilities of all the messages in the ensemble. First, we sort the messages in decreasing order of their probabilities. Then, we divide the sequence of messages into two consecutive groups in such a way that the probability of a message being in the first group is as nearly equal to the probability of a message being in the second group as possible. After the 1st division, we have two groups each of which consists of one or more messages. If a part consists of more than one message, we divide it into two parts again in the same manner - probabilities of the two parts should be as nearly equal as possible. We continue divisions until

Probabilities of Groups Obtained by Successive Divisions									
I Div.	II Div.	III Div.	IV Div.	V Div.	VI Div.	Message	P(i)	Recoded Message	$P(i)B_g(i)$
0.49						00	0.49	0	0.49
0.51						01	0.14	100	0.42
		0.14				10	0.14	101	0.42
	0.28					02	0.07	1100	0.28
	0.23					20	0.07	1101	0.28
			0.07			11	0.04	1110	0.16
			0.07			12	0.02	11110	0.10
		0.14				21	0.02	111110	0.12
		0.09				22	0.01	111111	0.06
			0.04						
			0.05						
				0.02					
				0.03					
					0.02				
					0.01				
									$(B_g)_{av.} = 2.33$

Fig. 2

every part contains one message. Each division of a group of messages contributes one binary digit to codes of the messages in the group. When divided into two subgroups, messages in one subgroup have 0 added to their codes; messages in the other subgroup have 1 added to their codes. This method of encoding messages was developed independently by Shannon [3]. We will call it Shannon-Fano coding. The following quote from [3] gives one more description of Shannon-Fano coding: "... arrange the messages ... in order of decreasing probability. Divide this series into two groups of as nearly equal probability as possible. If the message is in the first group its binary digit will be 0, otherwise 1. The groups are similarly divided into subsets of nearly equal probability and the particular subset determines the second binary digit. This process is continued until each subset contains only one message."

Let us see how Shannon-Fano coding is done in the example from Figure 2. We start with messages 00, 01, 10, 02, 20, 11, 12, 21, and 22; their probabilities are listed in column "P(i)". After the 1st division, the two groups of messages are {00} and {01, 10, 02, 20, 11, 12, 21, 22} with probabilities 0.49 and 0.51, respectively. Message 00 is assigned code 0; messages in the second group have their codes begin with 1. Since the second group consists of more than one message it is divided again in two groups, {01, 10} and {02, 20, 11, 12, 21, 22}, with probabilities 0.28 and 0.23, respectively. Messages 01 and 10 have 0 added as 2nd binary digit in their codes; their codes begin with "10". Messages 02, 20, 11, 12, 21, and 22 have 1 added as 2nd binary digit in their codes. Divisions continue until each group contains one message only. Codes obtained for each message are listed in column "Recoded message".

Successive divisions of messages and obtained message codes can be represented as a code tree as shown in Figure 3. Each leaf is labeled with a message and its probability. Each internal node is labeled with the sum of probabilities of the leaves in its subtree. The root is labeled with 1 -

the sum of probabilities of all the messages. Each edge from a node to its child is labeled with a digit (0 or 1) which appears in codes of leaves of the subtree rooted at the child. The code of a message can be obtained by reading labels of edges on the path from the root of the tree to the leaf containing the message. For instance, path from the root of the tree in Figure 3 to the leaf containing message “20” consists of edges labeled 1, 1, 0, and 1, which indicates that the code of message “20” is “1101”.

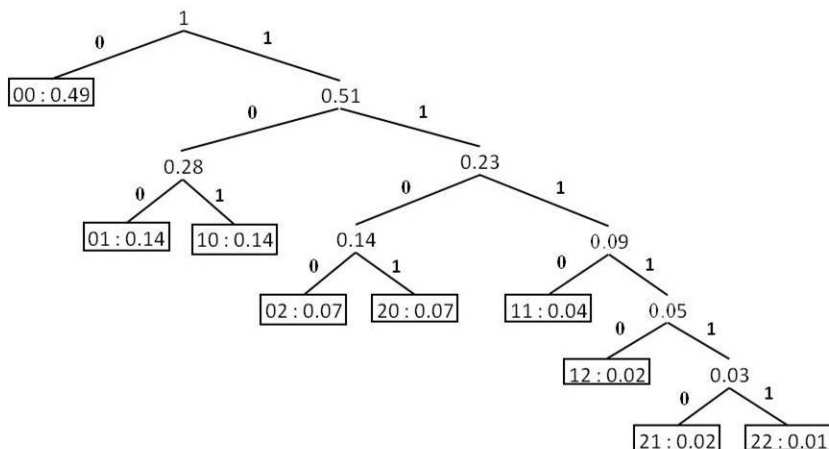


Fig. 3 Tree corresponding to the coding scheme in Figure 2.

Shannon-Fano coding is an example of a greedy algorithm. Greedy algorithms are algorithms that at every step make a choice that looks best at the moment. In Shannon-Fano coding, at every step (division) a choice (division of messages into two consecutive groups) is made that looks best at the moment. Best in a sense that the probabilities of the two parts should be as nearly equal as possible. Shannon-Fano coding uses top-down approach, dividing groups of messages into subgroups. The code tree is built top-down, from the root (that corresponds to all the messages) to the leaves (each leaf corresponds to a single message). At each step, a group of messages is divided into two subgroups and two children corresponding to the subgroups are added to the node that corresponds to the original group.

**Exercise 2.6.** Assume that we have messages A, B, C, D, E, F, G. The probability of A is 0.12, the probability of B is 0.34, the probability of C is 0.05, the probability of D is 0.22, the probability of E is 0.15, the probability of F is 0.02, and the probability of G is 0.10. Use Shannon-Fano coding to produce recoded messages for each of A, B, C, D, E, F, G. Represent obtained coding scheme as a tree.

**Exercise 2.7.** Assume that there are  $n$  messages. For each message  $i$ , we know the probability of the message,  $P(i)$ . Write a pseudocode of the procedure described by Fano (Shannon-Fano coding). What is the running time of the procedure?

**Exercise 2.8.** Implement Fano’s method. Given a set of messages with their probabilities, your program should output codes for each message.

Sometimes, instead of probabilities we know frequencies of messages - how often a message appears in a certain text. In this case, we can compute probabilities by assuming that they are proportional to corresponding frequencies. In other words, the probability of a message is set to be the frequency of the message divided by the sum of frequencies of all messages. However, in the



Shannon-Fano coding procedure we do not need to compute probabilities, we can use frequencies instead.

**Exercise 2.9.** Argue that if you use frequencies instead of probabilities in Shannon-Fano coding, you get the same codes for messages as you would get when using probabilities.

**Exercise 2.10.** Use Fano’s procedure to produce encodings for the six symbols with their frequencies shown in Figure 4. Draw a code tree that corresponds to the encoding you obtained.

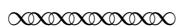
Symbol	Frequency
A	56
B	40
C	33
D	75
E	6
F	24

Fig. 4

Shannon-Fano coding method does not always produce the best (optimal) encoding - encoding that reduces the number of digits required to transmit a message the most. Huffman approach is better as it always finds an optimal encoding.

### 3 Huffman encoding: I

Let us look at the work of Huffman. The Introduction from Huffman’s paper “A Method for the Construction of Minimum-Redundancy Codes” ([2]) is quoted below. It defines the “optimum code” and presents basic restrictions that a code must satisfy to be optimum.



#### Introduction

One important method of transmitting messages is to transmit in their place sequences of symbols. If there are more messages which might be sent than there are kinds of symbols available, then some of the messages must use more than one symbol. If it is assumed that each symbol requires the same time for transmission, then the time for transmission (length) of a message is directly proportional to the number of symbols associated with it. In this paper, the symbol or sequence of symbols associated with a given message will be called the “message code.” The entire number of messages which might be transmitted will be called the “message ensemble.” The mutual agreement between the transmitter and the receiver about the meaning of the code for each message of the ensemble will be called the “ensemble code.”

Probably the most familiar ensemble code was stated in the phrase “one if by land and two if by sea”, and the message codes were “one” and “two.”

In order to formalize the requirements of an ensemble code, the coding symbols will be represented by numbers. Thus, if there are  $D$  different types of symbols to be used in coding, they will be represented by the digits 0, 1, 2, ...,  $(D - 1)$ . For example, a ternary code will be constructed using the three digits 0,1, and 2 as coding symbols.

The number of messages in the ensemble will be called  $N$ . Let  $P(i)$  be the probability of the  $i$ th message. Then

$$\sum_{i=1}^N P(i) = 1 \quad (1)$$

The length of a message,  $L(i)$ , is the number of coding digits assigned to it. Therefore, the average message length is:

$$L_{av} = \sum_{i=1}^N P(i)L(i) \quad (2)$$

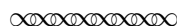
The term “redundancy” has been defined by Shannon as a property of codes. A “minimum-redundancy code” will be defined here as an ensemble code which, for a message ensemble consisting of a finite number of members,  $N$ , and for a given number of coding digits,  $D$ , yields the lowest possible average message length. In order to avoid the use of the lengthy term “minimum-redundancy,” this term will be replaced here by “optimum.” It will be understood then that, in this paper, “optimum code” means “minimum-redundancy code.”

The following basic restrictions will be imposed on an ensemble code:

- (a) No two messages will consist of identical arrangements of coding digits.
- (b) The message codes will be constructed in such a way that no additional indication is necessary to specify where a message code begins and ends once the starting point of a sequence of messages is known.

Restriction (b) necessitates that no message be coded in such a way that its code appears, digit for digit, as the first part of any message code of greater length. Thus, 01, 102, 111, and 202 are valid message codes for an ensemble of four members. For instance, a sequence of these messages 1111022020101111102 can be broken up into the individual messages 111-102-202-01-01-111-102. All the receiver needs to know is the ensemble code. However, if the ensemble has individual message codes including 11, 111, 102, and 02, then when a message sequence starts with the digits 11, it is not immediately certain whether the message 11 has been received or whether it is only the first two digits of the message 111. Moreover, even if the sequence turns out to be 11102, it is still not certain whether 111-02 or 11-102 was transmitted. In this example, change of one of the two message codes 111 or 11 indicated.

C.E. Shannon and R. M. Fano have developed ensemble coding procedures for the purpose of proving that the average number of binary digits required per message approaches from above the average amount of information per message. Their coding procedures are not optimum, but approach the optimum behavior when  $N$  approaches infinity. Some work has been done by Kraft toward deriving a coding method which gives an average code length as close to possible to the ideal when the ensemble contains a finite number of members. However, up to the present time, no definite procedure has been suggested for the construction of such a code to the knowledge of the author. It is the purpose of this paper to derive such a procedure.



**Exercise 3.1.** In the Introduction Huffman defined a notion of “message ensemble”. What is the message ensemble in the example of Figure 2? What is the average message length of the encoding in Figure 2?

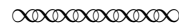
**Exercise 3.2.** What is the average message length of the encoding in Exercise 2.10?

**Exercise 3.3.** Does the coding in Figure 2 satisfy restrictions (a) and (b) from the Introduction of Huffman’s paper? Will an encoding produced by Shannon-Fano method always satisfy restrictions (a) and (b), or not?

**Exercise 3.4.** Recall that Huffman defines a code to be “optimum” if it yields the lowest possible average message length. Let a message ensemble consists of three messages: A, B, and C. Assume that only two symbols, 0 and 1, can be used in coding. Give an example of a code for the ensemble which is not optimum. Prove that it is not optimum. You may choose any values for probabilities of the messages. (Hint: In order to prove that the code is not optimum, give another code that has a smaller average message length.)

## 4 Huffman encoding: II

In the following excerpts from sections Derived Coding Requirements and Optimum Binary Code from Huffman’s paper ([2]) Huffman derives more requirements that an optimum code must satisfy, proposes a procedure to find an optimum code and proves that it is correct (that the code found by the procedure is always optimum).



### Derived Coding Requirements

For an optimum code, the length of a given message code can never be less than the length of a more probable message code. If this requirement were not met, then a reduction in average message length could be obtained by interchanging the codes for the two messages in question in such a way that the shorter code becomes associated with the more probable message. Also, if there are several messages with the same probability, then it is possible that the codes for these messages may differ in length. However, the codes for these messages may be interchanged in any way without affecting the average code length for the message ensemble. Therefore, it may be assumed that the messages in the ensemble have been ordered in a fashion such that

$$P(1) \geq P(2) \geq \dots \geq P(N - 1) \geq P(N) \quad (3)$$

And that, in addition, for an optimum code, the condition

$$L(1) \leq L(2) \leq \dots \leq L(N - 1) \leq L(N) \quad (4)$$

holds. This requirement is assumed to be satisfied throughout the following discussion.

It might be imagined that an ensemble code could assign  $q$  more digits to the  $N$ th message than to the  $(N - 1)$ th message. However, the first  $L(N - 1)$  digits of the  $N$ th message must not be used as the code for any other message. Thus the additional  $q$  digits would serve no useful purpose and would unnecessarily increase  $L_{av}$ . Therefore, for an optimal code it is necessary that  $L(N)$  be equal to  $L(N - 1)$ .

The  $k$ th prefix of a message code will be defined as the first  $k$  digits of that message code. Basic restriction (b) could then be restated as: No message shall be coded in such a way that its code is a prefix of any other message, or that any of its prefixes are used elsewhere as a message code.

Imagine an optimum code in which no two of the messages coded with length  $L(N)$  have identical prefixes of order  $L(N) - 1$ . Since an optimum code has been assumed, then none of these messages of length  $L(N)$  can have codes or prefixes of any order which correspond to other codes. It would then be possible to drop the last digit of all of this group of messages and thereby reduce the value of  $L_{av}$ . Therefore, in an optimum code, it is necessary that at least two (and no more than  $D$ ) of the codes with length  $L(N)$  have identical prefixes of order  $L(N) - 1$ .

One additional requirement can be made for an optimum code. Assume that there exists a combination of the  $D$  different types of coding digits which is less than  $L(N)$  digits in length and which is not used as a message code or which is not a prefix of a message code. Then this combination of digits could be used to replace the code for the  $N$ th message with a consequent reduction of  $L_{av}$ . Therefore, all possible sequences of  $L(N) - 1$  digits must be used either as message codes, or must have one of their prefixes used as message codes.

The derived restrictions for an optimum code are summarized in condensed form below and considered in addition to restrictions (a) and (b) given in the first part of this paper:

$$(c) \quad L(1) \leq L(2) \leq \dots \leq L(N - 1) = L(N). \quad (5)$$

(d) At least two and not more than  $D$  of the messages with code length  $L(N)$  have codes which are alike except for their final digits.

(e) Each possible sequence of  $L(N) - 1$  digits must be used either as a message code or must have one of its prefixes used as a message code.

### Optimum Binary Code

For ease of development of the optimum coding procedure, let us now restrict ourselves to the problem of binary coding. Later this procedure will be extended to the general case of  $D$  digits.

Restriction (c) makes it necessary that the two least probable messages have codes of equal length. Restriction (d) places the requirement that, for  $D$  equal to two, there be only two of the messages with coded length  $L(N)$  which are identical except for their last digits. The final digits of these two codes will be one of the two binary digits, 0 and 1. It will be necessary to assign these two message codes to the  $N$ th and the  $(N - 1)$ st messages since at this point it is not known whether or not other codes of length  $L(N)$  exist. Once this has been done, these two messages are equivalent to a single composite message. Its code (as yet undetermined) will be the common prefixes of order  $L(N) - 1$  of these two messages. Its probability will be the sum of the probabilities of the two messages from which it was created. The ensemble containing this composite message in the place of its two component messages will be called the first auxiliary message ensemble.

This newly created ensemble contains one less message than the original. Its members should be rearranged if necessary so that the messages are again ordered according to their probabilities. It may be considered exactly as the original ensemble was. The codes for each of the two least probable messages in this new ensemble are required to be identical except in their final digits; 0 and 1 are assigned as these digits, one for each of the two messages. Each new auxiliary ensemble contains one less message than the preceding ensemble. Each auxiliary ensemble represents the original ensemble with full use made of the accumulated necessary coding requirements.

The procedure is applied again and again until the number of members in the most recently formed auxiliary message ensemble is reduced to two. One of each of the binary digits is assigned to each of these two composite messages. These messages are then combined to form a single composite message with probability unity, and the coding is complete. Now let us examine Table I. The left-hand column

Table I (Huffman)

OPTIMUM BINARY CODING PROCEDURE

Message Probabilities												
Original Message Ensemble	Auxiliary Message Ensembles											
	1	2	3	4	5	6	7	8	9	10	11	12
0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20	0.20
0.18	0.18	0.18	0.18	0.18	0.18	0.18	0.18	0.18	0.18	0.18	0.18	0.18
0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06
0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06	0.06
0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
*0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03
0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01

contains the ordered message probabilities of the ensemble to be coded.  $N$  is equal to 13. Since each combination of two messages (indicated by a bracket) is accompanied by the assigning of a new digit to each, then the total number of digits which should be assigned to each original message is the same as the number of combinations indicated for that message. For example, the message marked \*, or a composite of which it is a part, is combined with others five times, and therefore should be assigned a code length of five digits.

When there is no alternative in choosing the two least probable messages, then it is clear that the requirements, established as necessary, are also sufficient for deriving an optimum code. There may arise situations in which a choice may be made between two or more groupings of least likely messages. Such a case arises, for example, in the fourth auxiliary ensemble of Table I. Either of the messages of probability 0.08 could have been combined with that of probability 0.06. However, it is possible to rearrange codes in any manner among equally likely messages without affecting the average code length, and so a choice of either of the alternatives could have been made. Therefore, the procedure given is always sufficient to establish an optimum binary code.

The lengths of all the encoded messages derived from Table I are given in Table II.

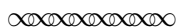
Having now determined proper lengths of code for each message, the problem of specifying the actual digits remains. Many alternatives exist. Since the combining of messages into their composites is similar to the successive confluences of trickles, rivulets, brooks, and creeks into a final large river, the procedure thus far described might be considered analogous to the placing of signs by a water-borne insect at each of these junctions as he journeys downstream. It should be remembered that the code which we desire is that one which the insect must remember in order to work his way back up stream. Since the placing of the signs need not follow the same rule, such as "zero-right-returning", at each junction, it can be seen that there are at least  $2^{12}$  different ways of assigning code digits for our example.

The code in Table II was obtained by using the digit 0 for the upper message and the digit 1 for the lower message of any bracket. It is important to note in Table I that coding restriction (e) is

Table II (Huffman)  
RESULTS OF OPTIMUM BINARY CODING PROCEDURE

$i$	$P(i)$	$L(i)$	$P(i)L(i)$	$Code$
1	0.20	2	0.40	10
2	0.18	3	0.54	000
3	0.10	3	0.30	011
4	0.10	3	0.30	110
5	0.10	3	0.30	111
6	0.06	4	0.24	0101
7	0.06	5	0.30	00100
8	0.04	5	0.20	00101
9	0.04	5	0.20	01000
10	0.04	5	0.20	01001
11	0.04	5	0.20	00110
12	0.03	6	0.18	001110
13	0.01	6	0.06	001111
			$L_{av} = 3.42$	

automatically met as long as two messages (and not one) are placed in each bracket.



Huffman algorithm is another example of a greedy algorithm as it looks for two least probable messages at every step. Huffman algorithm uses bottom-up approach. At every step it combines two messages into a single composite message. The code tree is built bottom-up, from the leaves to the root. At each step, two nodes that correspond to two least probable messages are made children of a new node that corresponds to the composite message.

Both, Shannon-Fano and Huffman algorithms are greedy algorithms that try to solve the same problem. One uses top-down approach while the other one uses bottom-up approach. However, the greedy strategy of Shannon-Fano does not always produce an optimal solution. Huffman approach, on the other hand, always finds an optimal solution.

**Exercise 4.1.** Draw a code tree that corresponds to Tables I and II in Huffman's paper.

**Exercise 4.2.** Use the Huffman algorithm for messages and probabilities from Figure 2 (Fano). What codes do you get for the messages? Is the encoding better or worse than the one obtained by Fano?

**Exercise 4.3.** Use Huffman algorithm for symbols and frequencies from Figure 4. What code tree do you get? What are codes for the symbols? Is the obtained encoding better or worse than the one produced earlier by Shannon-Fano approach?

**Exercise 4.4.** Assume that message ensemble consists of  $n$  messages:  $1, 2, \dots, n$ . Frequency of message  $i$  is equal to Fibonacci number  $F_i$  ( $1 \leq i \leq n$ ). Recall that  $F_1 = 1, F_2 = 1, F_k = F_{k-1} + F_{k-2}$ . What will a code tree produced by Huffman algorithm look like?

**Exercise 4.5.** Use Shannon-Fano approach for messages and probabilities from Table II (Huffman). What codes do you get for the messages? Is the encoding better or worse than the one obtained by Huffman?

**Exercise 4.6.** Give your own example when Shannon-Fano procedure produces an optimal code. You need to specify what the messages are, what their probabilities (or frequencies) are, compute Shannon-Fano coding, and compare it with optimal code obtained by using Huffman procedure.

**Exercise 4.7.** Give your own example when Shannon-Fano procedure does not produce an optimal code. You need to specify what the messages are, what their probabilities (or frequencies) are, compute Shannon-Fano coding, and compare it with the optimal code obtained by using Huffman procedure.

**Exercise 4.8.** Given a message ensemble and the ensemble code, how can you show that the code is optimal? Is verification of a code being optimal easier or not than finding an optimal code?

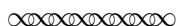
**Exercise 4.9.** You are given a code tree for a message ensemble. If you swap two leaves in the tree and update labels of internal nodes so that they are equal to the sum of probabilities of their descendant leaves, you get a different code for the ensemble. Consider the following conjecture: if every code obtained by swapping two leaves in the code tree is no better (is the same or worse) than the code of the original tree, then the code of the original tree is an optimum code. Prove the conjecture or give a counterexample.

**Exercise 4.10.** Write a program that would produce a Huffman encoding given messages and their frequencies.

**Exercise 4.11.** In the procedure described by Fano, the divisions in groups and subgroups are made without changing the order of the messages. Changing the order of messages may allow us to divide the group into subgroups with less difference between probabilities of subgroups. Assume that at every step of the Fano procedure we get the best possible split. Will this modified procedure yield an optimum code?

## 5 Huffman encoding: III

Huffman describes a generalization of his method to the case when three or more digits are used for coding. It is done in Generalization of the Method section from Huffman's paper ([2]).



### Generalization of the Method

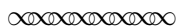
Optimum coding of an ensemble of messages using three or more types of digits is similar to the binary coding procedure. A table of auxiliary message ensembles similar to Table I will be used. Brackets indicating messages combined to form composite messages will be used in the same way as was done in Table I. However, in order to satisfy restriction (e), it will be required that all these brackets, with the possible exception of one combining the least probable messages of the original ensemble, always combine a number of messages equal to  $D$ .

It will be noted that the terminating auxiliary ensemble always has one unity probability message. Each preceding ensemble is increased in number by  $D - 1$  until the first auxiliary ensemble is reached. Therefore, if  $N_1$  is the number of messages in the first auxiliary ensemble, then  $(N_1 - 1)/(D - 1)$  must be an integer. However  $N_1 = N - n_0 + 1$ , where  $n_0$  is the number of the least probable messages combined in a bracket in the original ensemble. Therefore,  $n_0$  (which, of course, is at least two and no more than  $D$ ) must be of such a value that  $(N - n_0)/(D - 1)$  is an integer.

Table III (Huffman)  
OPTIMUM CODING PROCEDURE FOR  $D=4$

Message Probabilities		$L(i)$	Code
<i>Original Message Ensemble</i>	<i>Auxiliary Ensembles</i>		
0.22	0.22	1	1
0.20	0.20	1	2
0.18	0.18	1	3
0.15	0.15	2	00
0.10	0.10	2	01
0.08	0.08	2	02
0.05	0.07	3	030
0.02	0.07	3	031

In Table III an example is considered using an ensemble of eight messages which is to be coded with four digits;  $n_0$  is found to be 2. The code listed in the table is obtained by assigning the four digits 0, 1, 2, and 3, in order, to each of the brackets.



**Exercise 5.1.** Implement Huffman generalized method to produce encoding which uses a given number of digits.

## References

- [1] R.M. Fano. The transmission of information. Technical Report 65, Research Laboratory of Electronics, M.I.T., Cambridge, Mass., 1949.
- [2] David A. Huffman. A method for the construction of minimum-redundancy codes. *Proceedings of the Institute of Radio Engineers*, 40(9):1098–1101, September 1952.
- [3] C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423, 623–656, July, October 1948.
- [4] Gary Stix. Profile: David A. Huffman. *Scientific American*, 265(3):54, 58, September 1991.



## Notes to the Instructor

The story of the invention of Huffman codes is a great story that demonstrates that students can do better than professors. David Huffman was a student in an electrical engineering course in 1951. His professor, Robert Fano, offered students a choice of taking a final exam or writing a term paper. Huffman did not want to take the final so he started working on the term paper. The topic of the paper was to find the most efficient (optimal) code. What professor Fano did not tell his students was the fact that it was an open problem and that he was working on the problem himself. Huffman spent a lot of time on the problem and was ready to give up when the solution suddenly came to him. The code he discovered was optimal, that is, it had the lowest possible average message length. The method that Fano had developed for this problem did not always produce an optimal code. Therefore, Huffman did better than his professor. Later Huffman said that likely he would not have even attempted the problem if he had known that his professor was struggling with it.

The project uses excerpts from Fano's work ([1]) and from Huffman's paper ([2]) where they present their encodings. Both Fano and Huffman used greedy strategies to find the codes. However, Fano's greedy algorithm would not always produce an optimal code while Huffman's greedy algorithm would always find an optimal solution. The purpose of the project is for students to learn greedy algorithms, prefix-free codes, Huffman encoding, binary tree representations of codes, and the basics of information theory (unit and amount of information). The project demonstrates that greedy strategy could be applied in different ways to the same problem, sometimes producing an optimal solution and sometimes not.

The project is designed for a junior level Data Structures and Algorithms course. It can be used when covering the topic of greedy algorithms. The project uses quotes from Fano's and Huffman's papers ([1], [2]) to look into greedy algorithms, prefix-free codes, Huffman encoding, binary tree representations of codes, unit and amount of information.

The project is divided into several sections. Each section, except for Introduction, contains a reading assignment (read selected quotes from the original sources and some additional explanations) and a list of exercises on the material from the reading. Exercises are of different types, some requiring simple understanding of a definition or an algorithm, and some requiring proofs. Exercises 2.8, 4.10, and 5.1 are programming exercises and require good programming skills.

The material naturally splits in two parts. The first part (sections 1 and 2) is based on Fano's paper ([1]). It talks about measuring information and Fano encoding. Students may need substantial guidance with this part. The second part (sections 3, 4, and 5) is based on Huffman's paper ([2]). Sections 3 and 4 discuss Huffman encoding and compare it to Fano's. Section 5 can be skipped as it describes a generalization of Huffman's method to the case when three or more digits are used for coding. Each of the two parts can be used as a one week homework assignment which can be done individually or in small groups. For instance, the first assignment may ask students to read sections 1 and 2 and do exercises 2.1-2.10. The second assignment may ask students to read sections 3 and 4 and do exercises 3.1-3.4 and 4.1-4.11 (section 5 is skipped).