

# **Case Study of Shannon Fano Coding**

**Bachelor of Technology  
In  
Electronics & Communication Engineering**

**By  
NEKKANTI GUNA SAI KIRAN  
ENROLLMENT No. : BT20ECE075**

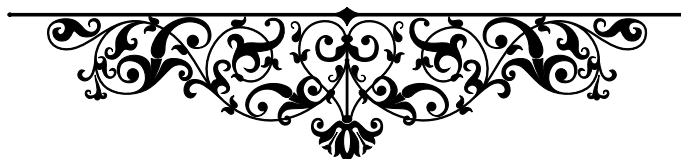
**Under the Guidance of  
Dr. Rashmi Pandhare**



**Indian Institute of Information Technology, Nagpur**



## Coding Techniques



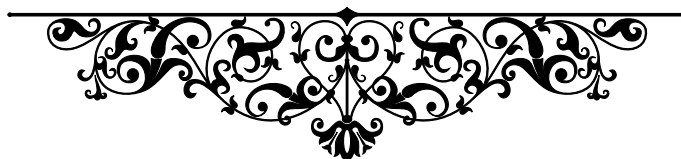
# I

## Details

- **Name :** Nekkanti Guna Sai Kiran
- **Enrollment No :** BT20ECE075
- **Lab Experiment:** 04
- **Course Title :** Coding Techniques
- **Lab Name :** Case Study On Shannon fano Coding In Information theory
- **Date :** 16/02/2023



## Case Study on Shannon Fano Coding



# I

## Introduction to Lossless data Compression

### 2.1.1 Information Theory And Source Coding

The most significant feature of the communication system is its unpredictability or uncertainty. The transmitter transmits at random any one of the pre-specified messages. The probability of transmitting each individual message is known. Thus our quest for an amount of information is virtually a search for a parameter associated with a probability scheme. The parameter should indicate a relative measure of uncertainty relevant to the occurrence of each message in the message ensemble.

The principle of improbability (which is one of the basic principles of the media world)—“if a dog bites a man, it’s no news, but if a man bites a dog, it’s a news”—helps us in this regard.

Hence there should be some sort of inverse relationship between the probability of an event and the amount of information associated with it. The more the probability of an event, the less is the amount of information associated with it, vice versa.

$$I(x_j) = f(1/p(X_j)),$$

Where  $X_j$  is an event with a probability  $p(X_j)$  the amount of information associated with it is  $I(x_j)$ .

Now let there be another event  $Y_k$  such that  $X_j$   $Y_k$  are independent. Hence probability of the joint event is  $p(X_j, Y_k) = p(X_j) p(Y_k)$  with associated information content ,

$$I(X_j, Y_k) = f(1/p(X_j, Y_k)) = f(1/p(X_j) \cdot 1/p(Y_k))$$

The total information  $I(X_j, Y_k)$  must be equal to the sum of individual information  $I(X_j) + I(Y_k)$ , where  $I(Y_k) = f(1/p(Y_k))$ .

Thus it can be seen that function  $f$  must be a function which converts the operation of multiplication into addition.

## LOGARITHM

It is one such function. Thus, the basic equation defining the amount of information (or self-information) is,  $I(X_j) = \log(1/P(X_j)) = -\log(P(X_j))$

when base is 2 (or not mentioned) the unit is bit, when base is  $e$  the unit is nat, when base is 10 the unit is decit or Hertley.

## ENTROPY:

Entropy is defined as the average information per individual message.

Let there be  $L$  different messages  $m_1, m_2, \dots, m_L$ , with their respective probabilities of occurrences be  $p_1, p_2, \dots, p_L$ . Let us assume that in a long time interval,  $M$  messages have been generated. Let  $M$  be very large so that  $M \gg L$ . The total amount of information in all  $M$  messages.

The number of messages  $m_1 = M \cdot p_1$ , the amount of information

in message  $m_1 = \log_2(1/P(X_i))$ , thus the total amount of information in all  $m_1$  messages =  $M \cdot p_1 \cdot \log_2(1/P(X_i))$ .

So, the total amount of information in all  $L$  messages will then be  $I = M \cdot (P_1) \cdot \log(1/P_1) + M \cdot (P_2) \cdot \log(1/P_2) + \dots + M \cdot (P_L) \cdot \log(1/P_L)$

So, the average information per message, or entropy, will then be  $H = I/M = (P_1) \cdot \log(1/P_1) + (P_2) \cdot \log(1/P_2) + \dots + (P_L) \cdot \log(1/P_L)$ ;  
Hence,  $H(X) = - \sum_{i=1}^L P(X_i) \log_2(P(X_i))$  , summation  $i=1$  to  $L$

The Entropy of a source in bits/symbol is given by  $H(X) = - \sum_{i=1}^L P(X_i) \log_2(P(X_i))$  , summation  $i=1$  to  $L$  Where  $X_i$  are the symbols with probabilities  $P(X_i)$  ,  $i=1,2,3 \dots L$   
The equality holds when the symbols are equally likely.

There are two types of code possible:

- 1) Fixed Length Code : All code words are of equal length
- 2) Variable Length Code: All code words are not of equal length. In such cases, it is important for the formation of uniquely decodable code that all the code words satisfy the PREFIX CONDITION, which states that “no code word forms the prefix of any other code word”.

The necessary sufficient condition for the existence of a binary code with code words having lengths  $n_1, n_2, \dots, n_L$  that satisfy the prefix condition is,

$$2^{\text{power}(-n_k)} \leq 1, \text{ summation } k = 1 \dots L,$$

## Source Coding Theorem:

Let  $X$  be ensemble of letters from a discrete memory less source with finite Entropy  $H(X)$  output symbols  $X_i$  with probabilities  $P(X_i)$ ,  $i=1,2,3,\dots,L$

It is possible to construct a code that satisfies the prefix condition has an average length  $R$  that satisfies the following inequality,

$H(x) \leq R < H(x)+1$ , the efficiency of the prefix code is defined as

$$= H(x)/R,$$

where,  $H(X) = - \sum_{i=1}^L P(X_i) \log_2(P(X_i))$

$$R = - \sum_{i=1}^L n_i \log_2(P(X_i))$$

Here  $n_i$  denotes the length of  $i$ th code word

The source coding theorem tells us that for any prefix code used to represent the symbols from a source, the minimum number of bits required to represent the source symbols on an average must be at least equal to the entropy of the source. If we have found a prefix code that satisfies  $R=H(x)$  for a certain source  $X$ , we must abandon further search because we can not do any better. The theorem also tells us that a source with higher entropy (uncertainty) requires on an average, more number of bits to represent the source symbols in terms of a prefix code.



**Proof:**

Lower bound:

First consider the lower bound of the inequality. For codewords that have length  $n_k$ ,  $1 \leq k \leq L$ , the difference  $H(x) - R$  can be expressed as

- $H(x) - R = \sum_{k=1}^L P(X_k) \log_2(1/P(X_k)) - \sum_{k=1}^L P(X_k) n_k$ , summation
- $H(x) - R = \sum_{k=1}^L P(X_k) \log_2(2^{-n_k}/P(X_k))$ , summation
- $H(x) - R = \sum_{k=1}^L P(X_k) (\log_2 e - n_k) (2^{-n_k}/P(X_k))^{-1}$ , summation
- $H(x) - R = \sum_{k=1}^L P(X_k) (\log_2 e - n_k) (2^{-n_k} - 1)$ , summation
- $H(x) - R \geq 0$  (using KRAFT'S INEQUALITY)
- $H(x) \geq R$

Upper bound:

Let us select a code word length  $n_k$  such that

$$2^{\lceil -n_k \rceil} P(X_k) < 2^{-n_k + 1}$$

First consider,  $2^{\lceil -n_k \rceil} P(X_k)$

- $2^{\lceil -n_k \rceil} P(X_k) = 1$ , summation  $k=1 \dots L$
- $\log_2(P(X_k)) < (-n_k + 1)$
- $n_k < 1 - \log_2(P(X_k))$
- $P(X_k) n_k < P(X_k) + P(X_k) \log_2(P(X_k))$ , summation  $k=1 \dots L$
- $R < H(x) + 1$

## II

### Shannon Fano Coding

#### 2.2.1 Shannon Fano Coding Algorithm

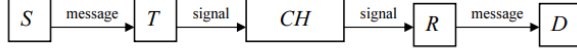
Shannon coding is yet another coding algorithm, which follows PREFIX Condition can achieve compression performance close to that of Huffman coding.

With his paper “The Mathematical Theory of Communication” (1948), Shannon offered precise results about the resources needed for optimal coding and for error-free communication. This 3 paper was immediately followed by many works of application to fields as radio, television and telephony. Shannon’s theory was later mathematically axiomatized (Khinchin 1957).

According to Shannon (1948; see also Shannon and Weaver 1949), a general communication system consists of five parts:

- A **source S**, which generates the message to be received at the destination.
- A **transmitter T**, which turns the message generated at the source into a signal to be transmitted. In the cases in which the information is encoded, encoding is also implemented by this system.
- A **channel CH**, that is, the medium used to transmit the signal from the transmitter to the receiver.

- A **receiver R**, which reconstructs the message from the signal.
- A **destination D**, which receives the message.



The source  $S$  is a system with a range of possible states  $S_1, \dots, S_n$  usually called *letters*, whose respective probabilities of occurrence are  $p(s_1), \dots, p(s_n)$ .

The amount of information generated at the source by the occurrence of  $S_i$  is defined as:

$$I(s_i) = \log(1/p(s_i)) = -\log p(s_i) \quad (1)$$

Since  $S$  produces sequences of states, usually called *messages*, the *entropy of the source*  $S$  is defined as the average amount of information produced at the source:

$$H(S) = \sum_{i=1}^n p(s_i) \log(1/p(s_i)) = -\sum_{i=1}^n p(s_i) \log p(s_i) \quad (2)$$

Analogously, the destination  $D$  is a system with a range of possible states  $d_1, \dots, d_m$ , with respective probabilities  $p(d_1), \dots, p(d_m)$ . The amount of information  $I(d_j)$  received at the destination by the occurrence of  $d_j$  is defined as:

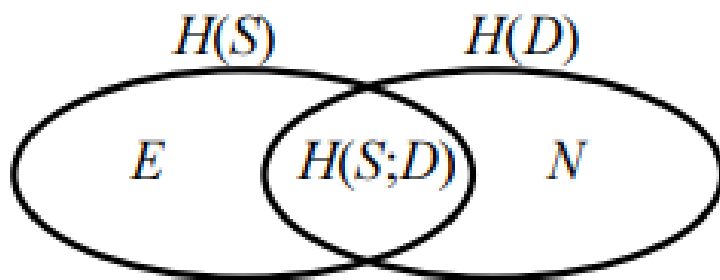
$$I(d_j) = \log(1/p(d_j)) = -\log p(d_j) \quad (3)$$

And the entropy of the destination  $D$  is defined as the average amount of information received at the destination:

$$H(D) = \sum_{j=1}^m p(d_j) \log(1/p(d_j)) = -\sum_{j=1}^m p(d_j) \log p(d_j)$$

In his original paper, Shannon (1948, p. 349) explains the convenience of the use of a logarithmic function in the definition of the entropies: it is practically useful because many important parameters in engineering vary linearly with the logarithm of the number of possibilities; it is intuitive because we use to measure magnitudes by linear comparison with unities of measurement; it is mathematically more suitable because many limiting operations in terms of the logarithm are simpler than in terms of the number of possibilities. In turn, the choice of a logarithmic base amounts to a choice of a unit for measuring information. If the base 2 is used, the resulting unit is called ‘bit’ –a contraction of binary unit. With these definitions, one bit is the amount of information obtained when one of two equally likely alternatives is specified.

The relationship between the entropies of the source  $H(S)$  and of the destination  $H(D)$  can be represented in the following diagram:



where,

- **H(S;D)** is the mutual information: the average amount of information generated at the source S and received at the destination D
- **E** is the equivocation: the average amount of information generated at S but not received at D.
- **N**, is the noise: the average amount of information received at D but not generated at S,

As the diagram clearly shows, the mutual information can be computed as:

$$\mathbf{H(S;D)=H(S)-E=H(D)-N}$$

Equivocation E and noise N are measures of the dependence between the source S and the destination D

- If S and D are completely independent, the values of E and N are maximum (  $E=H(S)$  and  $N=H(D)$ ), and the value of  $H(S;D)$  is minimum ( $H(S;D) = 0$  ).
- If the dependence between S and D is maximum, the values of E and N are minimum ( $E=N=0$ ), and the value of  $H(S;D)$  is maximum (  $H(S;D)=H(S)=H(D)$  ) .

The values of E and N are functions not only of the source and the destination, but also of the communication channel CH. The introduction of the communication channel leads directly to the possibility of errors in the process of transmission: the channel CH is defined by the matrix

where  $P(d/s)$  is the conditional probability of the occurrence of  $d$  in the destination  $D$  given that  $s$  occurred in the source  $S$ , and the elements in any row add up to 1. On this basis,  $E$  and  $N$  can be computed as:

$$N = \sum_{i=1}^n p(s_i) \sum_{j=1}^m p(d_j/s_i) \log(1/p(d_j/s_i))$$

$$E = \sum_{j=1}^m p(d_j) \sum_{i=1}^n p(s_i/d_j) \log(1/p(s_i/d_j))$$

where  $p(s_i/d_j) = p(d_j/s_i)p(s_i)/p(d_j)$ . The *channel capacity*  $C$  is defined as:

$$C = \max_{p(s_i)} H(S; D)$$

where the maximum is taken over all the possible distributions  $p(s)$  at the source.  $C$  is the largest average amount of information that can be transmitted over the communication channel  $CH$

On the other hand, in the early 1940s, it was thought that the increase of the rate in the information transmission over a communication channel would always increase the probability of error. The Second Theorem, or Noisy-Channel Coding Theorem, surprised the communication theory community by proving that that assumption was not true as long as the communication rate was maintained below the channel capacity. The channel capacity is equal to the maximum rate at which the information can be sent over the channel and recovered at the destination with a vanishingly low probability of error.

The formal simplicity of Shannon's theory might suggest that the interpretation of the involved concepts raises no difficulty. As we will see in the following sections, this is not the case at all.

### III

## Algorithm

Shannon coding is yet another coding algorithm, which follows PREFIX Condition can achieve compression performance close to that of Huffman coding.

Following are the steps to obtain the codewords by using Shannon Coding Algorithm:

- Arrange the symbols along with their Probability in decreasing order of Probabilities.
- Select all those symbols on one side whose probabilities is closest to half of the sum of probabilities of all symbols, remaining symbols on the other side assign the values 0 1 (first bit of their codeword) respectively to all the symbols in each of the two groups.
- Select the group having 0 assigned to each of its symbol in step 2 repeat step 2 for this group. Repeat the same task for the other group (whose symbols have been assigned 1 in the step 2)
- Repeat step 3 till a group remains to exist like above having more than one symbol.



Symbol	Stage I	Step I	Step II	Step III	Step IV	Step V
S <sub>0</sub>	0.30	0	0			00
S <sub>1</sub>	0.25	0	1			01
S <sub>2</sub>	0.20	1	0			10
S <sub>3</sub>	0.12	1	1	0		110
S <sub>4</sub>	0.08	1	1	1	0	1110
S <sub>5</sub>	0.05	1	1	1	1	1111

$$H(X) = 2.36 \text{ b/symbol}$$

$$L = 2.38 \text{ b/symbol}$$

$$\eta = H(X)/L = 0.99$$

The redundancy  $\gamma$  is defined as

$$\gamma = 1 - \eta$$

## IV

### Conclusion

Despite of its formal precision and its great many applications, Shannon's theory still offers an active terrain of debate when the interpretation of its main concepts is the task at issue. In this article we have tried to analyze certain points that still remain obscure or matter of discussion, and whose elucidation contribute to the assessment of the different interpretative proposals about the concept of information. Moreover, the present argumentation might shed light on the problems related with the so-called 'quantum information theory', in particular as formulated by Benjamin Schumacher (1995) on the basis of an explicit analogy with the first Shannon coding theorem. Furthermore, the discussion about the interpretation of the concepts involved in Shannon's theory would turn out to be particularly relevant if, as some believe, there were not two kinds of information classical and quantum, but only information encoded in different ways

## V

### ACKNOWLEDGMENT

This work was supported by Indian Institute of Information Technology Nagpur to promote research and innovation on Shannon Fano Coding Under the guidance of Dr. Rashmi Pandhare

## VI

### REFERENCES

- Aczel, J. and Forte, B. (1986), “ Generalized entropies and the maximum entropy principle, In: Bayesian Entropy and Bayesian Methods in Applied Statistics, edited by J. H. Justice, Cambridge University Press, Cambridge, pp. 95-100.
- Bar-Hillel, Y. and Carnap, R. (1952), “An outline of a theory of semantic information,” Tech. Rep. No., 247, Research Lab. of Electronics, MIT.
- Berger, T. (1971), Rate Distortion Theory, Englewood Cliffs, N.J.: Prentice-Hall.
- Brillouin, L. (1962) Science and Information Theory, Academic Press, New York.