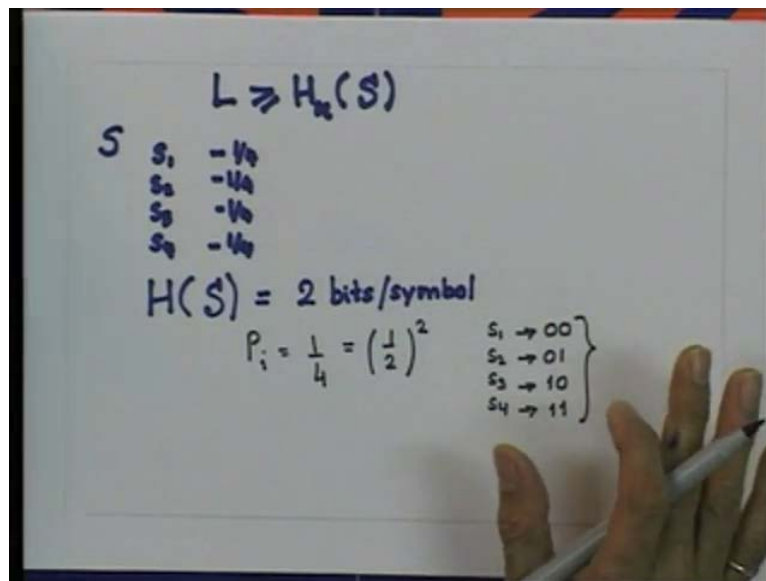


Information Theory And Coding
Prof. S. N. Merchant
Electrical Engineering
Indian institute of Technology, Bombay

Lecture – 10
Shannon's First Theorem

In the last class, we derived a very important result in information theory, which states that the average length of a code can never be greater than the entropy of a source.

(Refer Slide Time: 01:01)

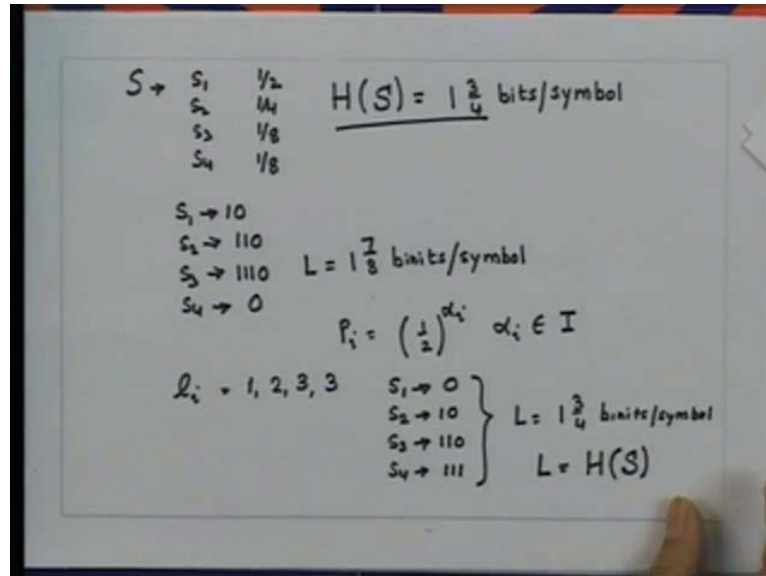


So, what we derived was average length has to be always greater than equal to the entropy of the source measured in array units. Now, to appreciate the importance of this result let us re visit some of the examples which we had studied in earlier in our course. So, one example which we which we which we had considered was that we had a source s consisting of 4 symbols s_1, s_2, s_3, s_4 each of this symbols are equiprobable. So, we know that the entropy of the source measured with the base 2 is equal to 2 bits per symbol.

Now, if I try to design any code a binary instantaneous code than my length of that code can never be less than 2 bits per symbol. And in this case each of this probabilities P_i which is equal to $1/4$ is of the form half raised to two. So, what it implies that I can design a compact code with 4 code words each of length 2 and 1 such code is. So, now the length of this code is also 2 binitis per symbol. So, there exists no uniquely decodable

code for this source with smaller average code word length, that is the length of 2 bits per symbol. Now, let us look at one more example to understand the derivation of that important result that is L is greater than or equal to entropy of a source.

(Refer Slide Time: 03:58)



So, if I take this same source s with the same 4 symbols, but with different probabilities given as half one fourth, one eighth, one eighth then the entropy for this source will turn out to be $1 \frac{3}{4}$ bits per symbol. Earlier in our course we had designed an instantaneous code for the source and that code was given as S_1 is 10, S_2 is 110, S_3 is 1110 and S_4 is 0. This code is a uniquely decodable code in fact it is an instantaneous code. Now, to calculate the length for this code it will turn out to be equal to $1 \frac{7}{8}$ bits per symbol.

So, even in this case we find the length of the code greater than the entropy of the source, but now each P_i in this case also is of the form P_i is of the form half raised to α_i where α_i belongs to integer. Therefore, it is possible to achieve the lower bound of $1 \frac{3}{4}$ bits per symbol and this is done by setting L_i equal to 1, 2, 3 and 3 respectively for this 4 symbol S_1, S_2, S_3 and S_4 . So, the code would be $S_1 \rightarrow 0, S_2 \rightarrow 10, S_3 \rightarrow 110$ and $S_4 \rightarrow 111$, if the length the average length for this code will come out to be $1 \frac{3}{4}$ bits per symbol. So, again in this case we find that the length average length turns out to be equal to the entropy of the source. And as a final example to explain the importance of L is equal to $H(S)$ let us consider a source s with seven symbols.

(Refer Slide Time: 07:13)

The whiteboard contains the following handwritten text:

S

S_1	$-\frac{1}{3}$
S_2	$-\frac{1}{3}$
S_3	$-\frac{1}{9}$
S_4	$-\frac{1}{9}$
S_5	$-\frac{1}{27}$
S_6	$-\frac{1}{27}$
S_7	$-\frac{1}{27}$

$H_3(S) = \frac{13}{9}$ trinary units/symbol

$P_i = \left(\frac{1}{3}\right)^{\alpha_i} \quad \alpha_i \in \mathbb{I}$

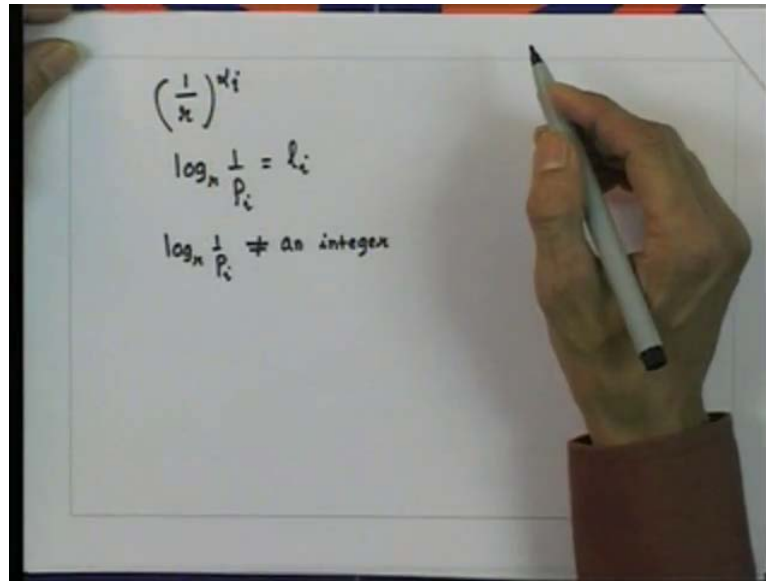
S_1	$\rightarrow 0$
S_2	$\rightarrow 1$
S_3	$\rightarrow 20$
S_4	$\rightarrow 21$
S_5	$\rightarrow 220$
S_6	$\rightarrow 221$
S_7	$\rightarrow 222$

$L = \sum_{i=1}^7 P_i l_i = \frac{13}{9}$ trinary symbols / Source symbol

So I have a source s consisting of seven symbols $S_1, S_2, S_3, S_4, S_5, S_6, S_7$ each with symbol probabilities given as one third, 1 third, 1 by 9, 1 by 9, 1 by 9, 1 by 7. Now, the entropy for this source will take in the 3 array units we calculate that will turn out to be 13 by 9 trinary units per symbol. And again we find that the symbol probabilities are of the form one third α_i where α_i belongs to integer.

And we can once we have the lengths α_i are nothing but the lengths for the code and we can design the instantaneous code as follows S_1 is 0 S_2 1. Now, if you calculate the length for this code as $P_i l_i$ i equal to 1 to 7 will turn out to be 13 by 9 trinary symbols per source symbol. So, far what we have seen is that we have looked at the coding problem for the 0 memory source with symbol probabilities of the form 1 by r^{α_i} .

(Refer Slide Time: 09:56)



So, if I have $\log_r 1/P_i$ is equal to l_i . If I have my P_i of this form then I can write $\log_r 1/P_i$ is equal to l_i where I choose my l_i equal to α_i . Now, the next question arises is that if this condition is not satisfied then how do I choose my length. Now, what it means that if $\log_r 1/P_i$ is not an integer then how do I choose my length.

It might seem reasonable that a compact code could be formed by choosing l_i as the first integer which is greater than $\log_r 1/P_i$. Now, this tempting conjecture is not valid if we select l_i in this manner it is not necessary that we will get a compact code, but selecting l_i in this manner where l_i is the integer, which is just larger than $\log_r 1/P_i$ can lead to some important results. So, let us select l_i , therefore as the unique integer satisfying this condition.

(Refer Slide Time: 12:01)

The whiteboard contains the following handwritten mathematical steps:

$$\log_r \frac{1}{P_i} \leq l_i < \log_r \frac{1}{P_i} + 1 \rightarrow \textcircled{1}$$

$$\frac{1}{P_i} \leq r^{l_i}$$

$$\Rightarrow P_i \geq r^{-l_i} \rightarrow \textcircled{2}$$

$$\sum_{i=1}^Q P_i \geq \sum_{i=1}^Q r^{-l_i}$$

$$1 \geq \sum_{i=1}^Q r^{-l_i}$$

$$P_i \sum_{i=1}^Q \log_r \frac{1}{P_i} \leq \sum_{i=1}^Q P_i l_i < \sum_{i=1}^Q P_i \log_r \frac{1}{P_i} + \sum_{i=1}^Q P_i$$

$$\underline{H_x(S) \leq L} \rightarrow \textcircled{3}$$

So, what we will do is we will select l_i as a integer which is just larger than this value it means that it follows this inequality. So, if I choose a set of l_i which follow this inequality then the next question is it possible for me to design or synthesize an instantaneous code which uses this set of l_i . So, that question can be answered that question can be answered if I can test Kraft's inequality.

So, taking exponential of the left inequality of equation 1 we will get $1/P_i$ is less than equal to r^{l_i} which implies that P_i is greater than r^{-l_i} . Now, summing over all i we obtain if I sum this all I assuming that the source is of size Q then I can write this relationship and this is 1 therefore, I get this relationship. Now, this relationship is you know that we have discussed this earlier and we have shown that if this condition is satisfied then it is possible for us to get an instantaneous code for that source. So, this choosing l_i according to this relationship given by 1 is acceptable for synthesis of an instantaneous code.

Now, equation 1 defines an acceptable set of l_i for an instantaneous code multiplying equation 1 by P_i and summing up over all, i we will find if I sum this up $P_i \sum_{i=1}^Q \log_r \frac{1}{P_i}$ is equal to $\sum_{i=1}^Q P_i l_i$ $P_i \sum_{i=1}^Q \log_r \frac{1}{P_i} + \sum_{i=1}^Q P_i$ these are summed over all i we get the relationship as this is the entropy of the source measured in array units this average length of the code, another very important result which we have derived.

Now, there is a difference between this result and the result which is. So, earlier another result which we saw earlier was a $H_r S$ is less than equal to L there is a important difference between this relationship and this relationship. This relationship expresses a bound for the average land of a code independent of any particular coding scheme the bound requires only that the code be instantaneous, whereas equation 3 on the other hand is derived by assuming the coding method given by equation 1. So, if I use the coding method described by equation 1 then I get this relationship and this relationship provides both a lower and upper bound on the average length of the code.

Now, since this relationship is valid for any 0 memory source we may apply it to the n th extension of the source. Let us see basically what happens if I apply this kind of coding scheme to an n th extension of a source. So, let us assume that I have a source s which is a 0 memory source and I look at its n th extension of this source. If I follow the strategy for coding which is given by this equation then this equation is also valid for the n th extension because n th extension of a 0 memory source is again a 0 memory source.

(Refer Slide Time: 18:09)

$S \quad S^n$
 $H_x(S^n) \leq L_n < H_x(S^n) + 1 \rightarrow (4)$
 $L_n = \sum_{i=1}^{Q^n} P(\sigma_i) \lambda_i \quad \lambda_i \rightarrow \text{length} \rightarrow \sigma_i$
 $\frac{L_n}{n} \rightarrow \text{average length for code symbols used per single source symbol from } S$
 $H_x(S^n) = n H_x(S)$
 $\frac{H_x(S^n)}{n} \leq \frac{L_n}{n} < \frac{H_x(S^n)}{n} + \frac{1}{n} \rightarrow S-(a)$
 $\lim_{n \rightarrow \infty} \frac{L_n}{n} = H_x(S) \rightarrow S-(b)$

So, the relationship which I get is $H_r S^n$ is less than equal to a L_n less than $H_r S^n$ plus 1, where L_n is the average length which is defined as probability of σ_i λ_i , i equal to 1 to Q^n where Q^n is the size of the n th extension source σ_i are the symbols of the n th extension of the source s^n λ_i are the length.

So, λ_i corresponds to the length of the code word corresponding to symbols from n th extension that is \bar{S} . This is the average length and L_n by n will give me the average length for code symbols used per single source symbol from S . Now, we have also seen that $H_r S_n$ that is entropy of the n th extension of a 0 memory source is equal to n times the entropy of the original source. So, based on this relationship we can write this equation. What this equation says is that it is possible to make L_n by n as close as we wish to $H_r S$ by coding the n th extension of s rather than original source that is S because S keeps on increasing 1 by n tend towards 0 . So, in that case L_n by n will tend towards $H_r S$ that is entropy of the original source.

Now, sp limit of n tending to infinity would be L_n by n is equal to $H_r S$. This is a very important relationship which we have derived. This equation is known as Shannon's first theorem or the noiseless coding theorem. It is one of the two major theorems of information theory equation 5 a tells us that we can make the average number of r array code symbols per source symbol as small as possible, but not smaller than the entropy of the source measured in r array units.

So, the prize which we pay for decreasing L_n by n quantity is the increasing complexity of the coding scheme. Now, all this discussion which we have done pertains to a 0 memory source and its extension. The next question arises is are this results also valid for a source which is not a 0 memory source, but say a Markov source. Let us extend our discussion to a Markov source.

(Refer Slide Time: 23:59)

Handwritten mathematical derivations on a whiteboard:

- Diagram showing $S \rightarrow \bar{S}$ with arrows pointing to the respective terms.
- Equation: $L = \sum_{i=1}^q P_i \lambda_i$
- Equation: $H_r(\bar{S}) \leq L$
- Equation: $H_r(S) \leq H_r(\bar{S}) \leq L$
- Equation (boxed): $H_r(S) \leq L$
- Equation: $S - S^N \quad \{s_1, s_2, s_3, \dots, s_N\} \equiv v_i$
- Equation: $H(V) \leq L_n < H(V) + 1$
- Equation: $N H_r(S) \leq L_n < N H_r(S) + 1$
- Equation: $H_r(S) \leq \frac{L_n}{N} < H_r(S) + \frac{1}{N}$
- Equation: $\ell = 1 \rightarrow q^N$
- Equation: $H(V) \hat{=} N H_r(S)$

Let me assume that I have a Markov source S and have its adjoint \bar{S} . We have seen the definition of a adjoint of a Markov source. Adjoint of a Markov source is a source with the source with the same source symbol as the source alphabet as the source S . And the symbol probabilities of the symbols in \bar{S} is the same as the first order symbol probabilities of source S . So, that is a definition of \bar{S} . Now, the process of encoding the symbols in the source S and the symbols which are identical to the source in the source alphabet S into an instantaneous block code is identical, because the source symbols are same. And the probability of the symbols are also the same.

In that case what it means that the average length which is defined as $\sum_{i=1}^Q p_i L_i$ is equal to 1 to Q this average length is also identical for both S and \bar{S} . \bar{S} however is a 0 memory source and we may apply the earlier derived result to obtain $H_r \bar{S}$ is less than equal to L . Now, we also have seen that entropy of the original source is always less than or equal to the entropy of its adjoint. And this is less than equal to L . So, what follows is that $H_r S$ is less than or equal to L .

So, again even for a Markov source we have shown that the average length of a code when we code the symbols individually in the source is greater or equal to the entropy of the source. Let us extend this result to an n th extension of a Markov source. So, I have a source S which is a Markov source and I am looking at the coding of this source in groups of n elements. So, that means I am looking at s_n where i code $S_{i-1} S_i S_{i+1}$ up to S_{i+n} and as 1 unit.

Now, if I start coding the source original source s in blocks of n source symbols then I can say that this source symbol form a new source symbol which is the i . And now i will range from 1 to 2 raise to n . Now, let us follow the same strategy which we followed earlier for a 0 memory source. Now, for a 0 memory source and its n th extension we had derived a result saying that L_n is greater or equal to this is result which we had derived for the 0 memory source.

And it is a n th extension, now based on this same thing is valid for a Markov source I can say that L_n is greater than or equal to H of v less than H of v plus 1. So, the same result is valid when I code this source in groups of n source symbols. Now, we have also seen that by definition H_v is equal to n times $H_n S$, where $H_n S$ is the average information per symbol. So, if you go by this definition which we had seen earlier in our

lectures I can write this expression as $n H_n \geq L_n$. Now, if I divide both the sides by n what I get is this result. Now, this is the average length of code symbols per source symbol of S . Now if I take the limit as n tends to infinity.

(Refer Slide Time: 31:14)

Handwritten notes on a whiteboard showing the derivation of entropy and a table for Shannon coding assignment.

Equation: $\lim_{N \rightarrow \infty} \frac{L_N}{N} = H_\infty(S) \rightarrow \text{Entropy of a source } S$

Equation: $\eta = \frac{H_\infty(S)}{L}$

Shannon-coding assignment

Inequality: $\log_2 \frac{1}{P_i} \leq l_i < \log_2 \frac{1}{P_i} + 1$

S_i	P_i	$\log_2 \frac{1}{P_i}$	l_i	code	$\frac{l_i}{n}$	\mathcal{B}
S_1	$2/3$	0.58	1	0	0	0
S_2	$2/9$	2.17	3	100	10	10
S_3	$1/9$	3.17	4	1010	11	11

$H(S) = 1.22 \text{ bits/symbol}$ $L_A = 1.78 \text{ bits/symbol}$ $L_B = 1.33 \text{ bits/symbol}$
 $H(S) < L_A < H(S) + 1$

I will get limit n tending to infinity L_n by n is equal to $H_\infty(S)$. And this is the entropy of a source which is a Markov source. So, again we have derived the result that if I want my average length to be as close as possible to the entropy of the source, then I have to use the extensions of the source. Now, another parameter we can define which relates the entropy and the average length of the source and that is known as efficiency of a code. Efficiency of a code is defined as entropy of source divided by the average length of the code for that source. So, this is by definition and efficiency of a code, what is desired is to have as large as possible the value for this etc.

Now, a question that arises is that we have used the coding technique which is based on this inequality. So, this inequality provides a some method of choosing the word length l_i by encoding the symbols from S_n and taking n sufficiently large the quantity L_n by n can be made as small as possible, but it cannot be smaller than the entropy of the source. The question is suppose if n is not very large n or capital N , if this quantity is not very large the theorem does not tell us what value of L what is the average length we shall obtain in that case. It does not guarantee that L or L_n by n will be smallest for that fixed n .

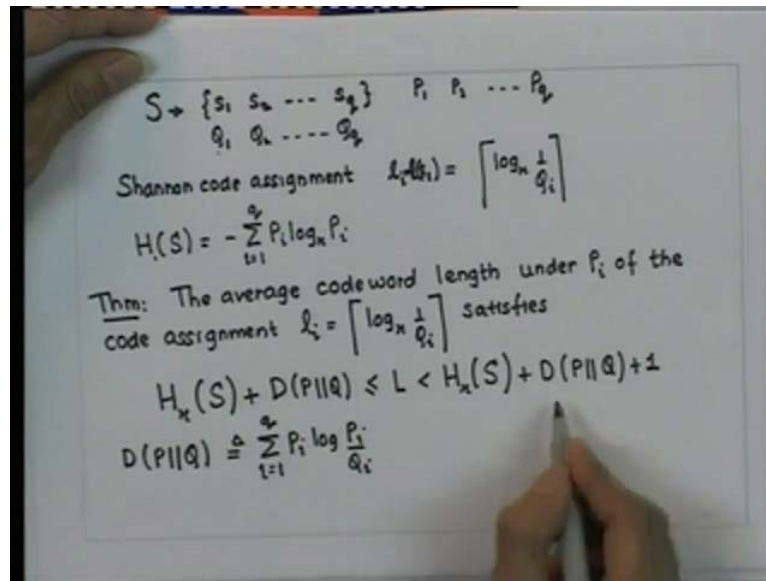
Let us try to understand this with a simple illustration. Suppose, I have a source consisting of 3 symbols S_1 S_2 S_3 and let us assume the P_i is for this is given as two third, $\frac{2}{9}$, $\frac{1}{9}$ \log_2 of $\frac{1}{P_i}$. If I assume that I am going to design a binary code than is equal to 0.58, 2.17, 3.17. So, my L_i which follows in equality will be 1 3 4. So, I can design a code which is 0 1 0 0 1 0 1 0.

So, I have designed a code a which is based on this coding strategy. Now, if you calculate the entropy for this source turns out to be H_s is equal to 1.22 bits per symbol. And if I calculate the length for this code is 1.78 bits per symbol. If you look at 1.78 it satisfies this inequality of H_s this inequality is satisfied by this code a , but unfortunately this coding strategy does not give me a compact code, a code with the length average length of the code to be as small as possible. So, if I take another code in this case if I have another code say B which is 0 1 0 1 one which is again an instantaneous code. And you can calculate the length for this code turns out to be 1.33 bits per symbol.

So, what it means that this procedure does not guarantee a compact code. Now, this code B 1.33 is very close to 1.22. So, what it also implies that by going for an extension of the source I may not able to achieve much. So, what this example demonstrate is that using this coding strategy, it is not essential that we always get a compact code. We could get a compact code when n is very large then when we consider n th extension of a source.

Otherwise we have seen in this example that follow another strategy I could get average length of the code which is smaller than the strategy for given by this inequality. Now, this strategy of coding is known as Shannon coding strategy or Shannon coding assignment for the lengths. Another question which arises is that whenever the design using this strategy based on some probability of symbols given here, but in reality if this probability is not correct than what is the effect on the average length. So, let us answer this question.

(Refer Slide Time: 39:28)



So, what we say is that we have a source s with source symbols given as S_1, S_2 , up to S_q . The real probability of these symbols are P_1, P_2 up to P_q , but for some reason this is not available and the estimate of these probabilities are available. So, let us call those estimates given as Q_1, Q_2 up to Q_q . So, as far as we are concerned since this is known to us we will be designing our code based on this information.

So, when we do that what it means that if we use Shannon code assignment. Then my lengths for the code word for each source symbol will be decided by log of this is a symbol which we use for finding out the first integer larger than this value. So, our designing of L_i will be based on Q_i . Now, the two function or the two probabilities are given by this. So, the entropy of the source is $P_i \log P_i$. So, this is a real entropy of the source.

And if we design our code properly than we expect that every length of our code should be as close as possible to $H_x(S)$, but since this is not known to us and we will be designing our code based on the length given by these probabilities. Then the average length of the code which we will be getting in a practical case would be very much different from the entropy of the source.

How much is the difference is in the entropy and the average length based on this is of interest to me. And this difference is given in the form of a theorem. The theorem says that the average code word length under real probability of real probabilities given by P_i

of the code assignment L is equal to \log of 1 by Q_i . This is what is available to me satisfies the following relationship where by definition P of Q is by definition given as $P_i \log$ of $P_i Q_i$ is equal to 1 to q . So, this is what we want to prove.

(Refer Slide Time: 44:33)

The image shows a handwritten mathematical proof on a piece of paper. The proof starts with the definition of the average code length L as a sum over all source symbols i of $P_i \lceil \log_r \frac{1}{Q_i} \rceil$. It then shows that this is less than the sum of $P_i (\log_r \frac{1}{Q_i} + 1)$. This is further simplified to $\sum P_i \log_r \frac{P_i}{Q_i} + 1$, which is equal to $\sum P_i \log_r \frac{P_i}{Q_i} + \sum P_i \log_r \frac{1}{P_i} + 1$. This is then identified as $D(P||Q) + H_r(S) + 1$. Finally, it concludes that $L < H_r(S) + D(P||Q) + 1$.

$$\begin{aligned} \text{Proof: } L &= \sum_{i=1}^n P_i \left\lceil \log_r \frac{1}{Q_i} \right\rceil \\ &< \sum_{i=1}^n P_i \left(\log_r \frac{1}{Q_i} + 1 \right) \\ &= \sum P_i \log_r \frac{P_i}{Q_i} + 1 \\ &= \sum_{i=1}^n P_i \log_r \frac{P_i}{Q_i} + \sum P_i \log_r \frac{1}{P_i} + 1 \\ &= D(P||Q) + H_r(S) + 1 \\ L &< H_r(S) + D(P||Q) + 1 \end{aligned}$$

Let us look at the proof of it the average length of the code which we will get is L is equal to P_i is the real probability of occurrence of a source symbol. And the design length for this source symbol S_i would be given by this quantity because Q_i are available to us. And this quantity is less than $P_i \log$ of 1 by Q_i plus 1 this is equal to $P_i \log$ of $P_i Q_i$ plus 1 by P_i . This is equal to summation of $P_i \log_r P_i Q_i$ plus $P_i \log$. This by definition is equal to d of plus $H_r(S)$ plus 1 . So, we have shown that L is less than $H_r(S)$ plus. Similarly, we can show the other side of the inequality not very difficult to do that.

(Refer Slide Time: 47:12)

$$\begin{aligned}
 L &= \sum_i P_i \left\lceil \log_x \frac{1}{Q_i} \right\rceil \\
 &\geq \sum_i P_i \log_x \frac{1}{Q_i} \\
 &= \sum_i P_i \log_x \frac{P_i}{Q_i} \cdot \frac{1}{P_i} \\
 &= D(P||Q) + H_x(S) \\
 H_x(S) + D(P||Q) &\leq L < H_x(S) + D(P||Q) + 1 \\
 D(P||Q) &\triangleq \sum_{i=1}^n P_i \log \frac{P_i}{Q_i} \leftarrow \text{Relative Entropy or Kullback Leibler distance between two probabilities (distribution) density}
 \end{aligned}$$

So, we write L is equal to a again $P_i \log$ of 1 by Q_i and this is greater or equal to summation where $i P_i \log r 1$ by Q_i is equal to this can be shown as equal to. So, finally, we get result as $H_r S$ plus d of Q where d of $P Q$ is by definition given as $P_i \log P_i Q_i$. So, this is known as this is termed as relative entropy or it is also known Kullback Leibler distance between two probability distributions. More appropriate would be to say to two probabilities density.

So, the penalty which we pay for the wrong choice for the probabilities of the source symbol is in terms of the relative entropy. So, our average length is going to increase by this quantity. So, we have looked at the mathematical relationship between the average length of a code and the entropy of the source for which the code has been designed. We also look at some of the methods to synthesize an instantaneous code. Now, in the next class we will have a look at different coding strategies which can be adapted in particular scenario to obtain average code word length to be as cool as possible to the entropy of a source without going further extension of the source.