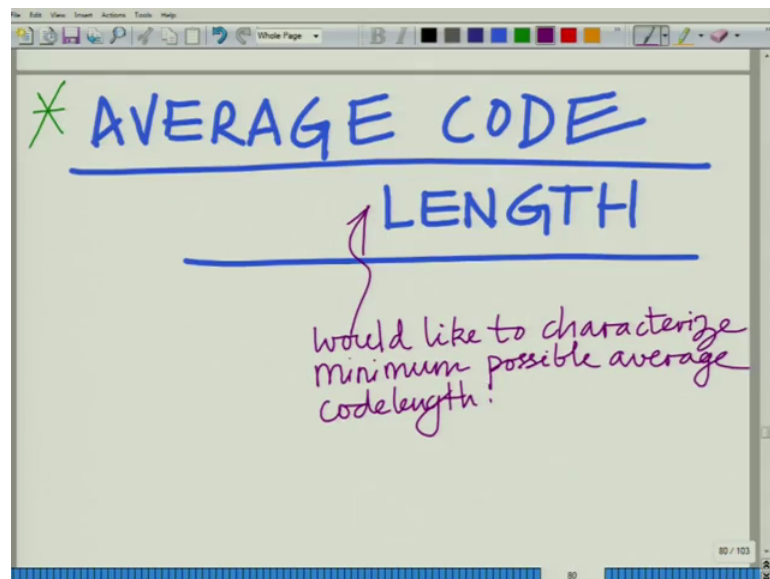


Principles of Communication Systems - Part II
Prof. Aditya K. Jagannatham
Department of Electrical Engineering
Indian Institute of Technology, Kanpur

Lecture - 44
Lower Bound on Average Code Length, Kullback-Leibler Divergence

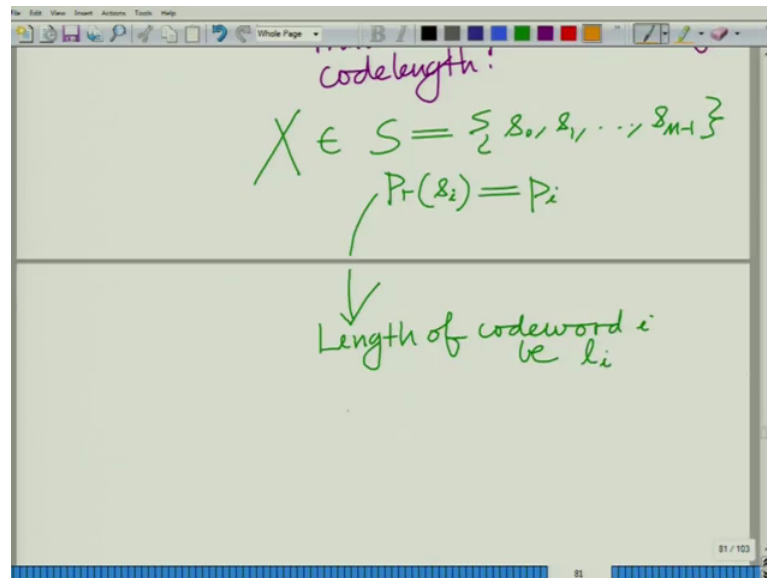
Hello. Welcome to another module in this massive open online course. So, we are looking at source coding, various aspects of source coding, different kinds of source codes in particular. We are interested in prefix free or instantaneous codes and in the previous module, we have looked at an important inequality that has to be satisfied by the lengths of the code words of a prefix free code that is given by the Kraft inequality. We will use this Kraft inequality to derive now a fundamental bound on the minimum possible average code length for a prefix free code, ok.

(Refer Slide Time: 00:50)



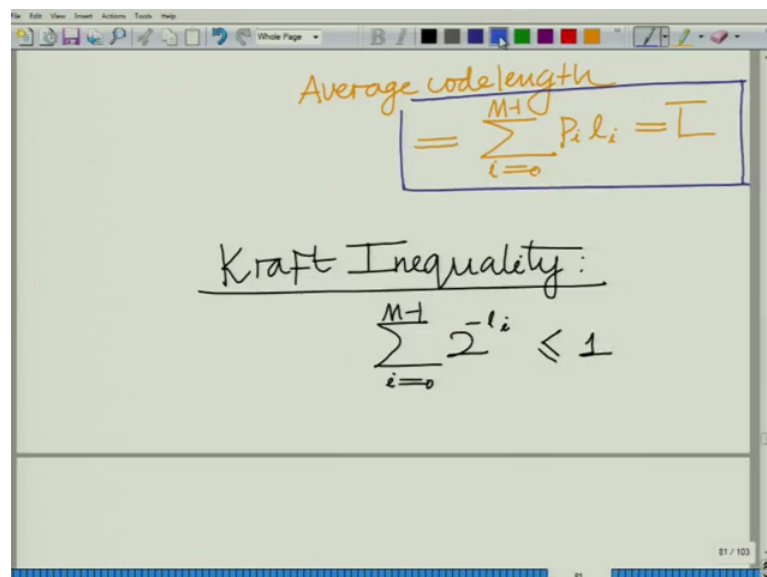
So, in this module, we would like to focus again on the average code length. As I have specified before, this is one of the fundamental aspects of a code that is we would like to characterize correct, we would like to characterize the minimum possible average code length. What is the minimum possible average code length? We would like to ask this question again towards this.

(Refer Slide Time: 01:54)



We would like to consider X drawn from a source with the alphabet s_0 up to s_{M-1} and we have probability of each symbol s_i is equal to probability of each alphabet s_i is purely symbol s_i is p_i . Now, the length we have seen, we have defined average length of the codeword.

(Refer Slide Time: 02:53)



Let the length of codeword be l_i , then we know that for the average code length or average codeword length is nothing, but expected value of l which is i equal to 0 to $M-1$ p_i , that is the summation of the lengths l_i weighted by the probabilities p_i ,

correct. This is expected value of, this is the average codeword length which we are denoting by \bar{l} . We have seen this before that the average codeword length is summation $p_i l_i$ equal to 0 to M minus 1.

Yesterday we have also defined, we have also seen the kraft inequality that the codeword lengths of any prefix free code have to satisfy and the kraft inequality, this is a fundamental inequality which we had derived from a binary representation. Let me just refresh your memory. So, I have i equal to 0 to M minus 1 to the power of minus l_i less than or equal to 1.

(Refer Slide Time: 4:20)

The image shows a whiteboard with the following handwritten content:

$$q_i = \frac{2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}}$$

Below this, it is noted that $q_i \geq 0$ because each $2^{-l_i} \geq 0$ and the denominator $\sum_{j=0}^{M-1} 2^{-l_j} \geq 0$.

Now, what we would like to define is, we would like to define this quantity q_i which is equal to 2 to the power of minus l_i by summation j equal to 0 to M minus 1 2 to the power of minus l_j . Just changing the index using a different index, instead of I am using j because I am using i equal to j equal to 0, j equal to 0 to M minus 1 2 to the power of minus l_j .

Now, what can we say about this q_i ? Now, we have defined this q_i . Now, if you can look at this q_i , the first thing you will observe is that q_i is greater than or equal to 0. All the quantities involved are positive. Q_i is greater than 0. Since each 2 to the power of minus l_i greater than equal to 0 summation j equal to 0 to M minus 1 2 to the power of minus l_j greater than equal to 0. So, q_i is less than equal to 0.

(Refer Slide Time: 05:43)

Further,
$$2^{-l_i} \leq \sum_{j=0}^{M-1} 2^{-l_j}$$
$$\Rightarrow \frac{2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}} \leq 1$$

$$\Rightarrow \frac{2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}} = q_i \leq 1$$

Further, we have 2 to the power of minus l_i , each 2 to the power of minus l_i . Remember all the quantities are positive. This is less than equal to summation j equal to 0 M minus 1 2 to the power of minus l_j because 2 to the power of minus l_i is one of the components in this summation, all right. So, all the quantities are positive and 2 to the power of minus l_i is in fact one of the quantities in the components in the summation on the right hand side. Therefore, 2 to the power of minus l_i is less than and equal to summation j equal to 0 to M minus 1 2 to the power of minus l_j which implies 2 to the power of minus l_i divided by summation j equal to 0 to M minus 1 2 to the power of minus l_j . This is less than or equal to let me just write this a little bit more clearly. This implies 2 to the power of minus l_i divided by summation j equal to 0 to M minus 1 2 to the power of minus l_j less than or equal to 1 or rather this is our, in fact nothing, but our q_j or q_i . This is q_i which is less than equal to 1.

(Refer Slide Time: 07:22)

The whiteboard shows the following handwritten work:

$$\Rightarrow \frac{2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}} = q_i \leq 1$$
$$0 \leq q_i \leq 1$$

Further,

$$\sum_{i=0}^{M-1} q_i = \sum_{i=0}^{M-1} \frac{2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}}$$

So, we have $q_i \geq 0$ less than or equal to q_i less than or equal to 1. Further summation i equal to 0 to M minus 1 q_i equals summation i equal to 0 to M minus 1 2 to the power of minus l_i divided by summation j equal to 0 to M minus 1 2 to the power of minus l_j , the denominator is a constant.

(Refer Slide Time: 07:59)

The whiteboard shows the following handwritten work:

$$= \frac{\sum_{i=0}^{M-1} 2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}} = 1$$
$$0 \leq q_i \leq 1$$
$$\sum_{i=0}^{M-1} q_i = 1$$

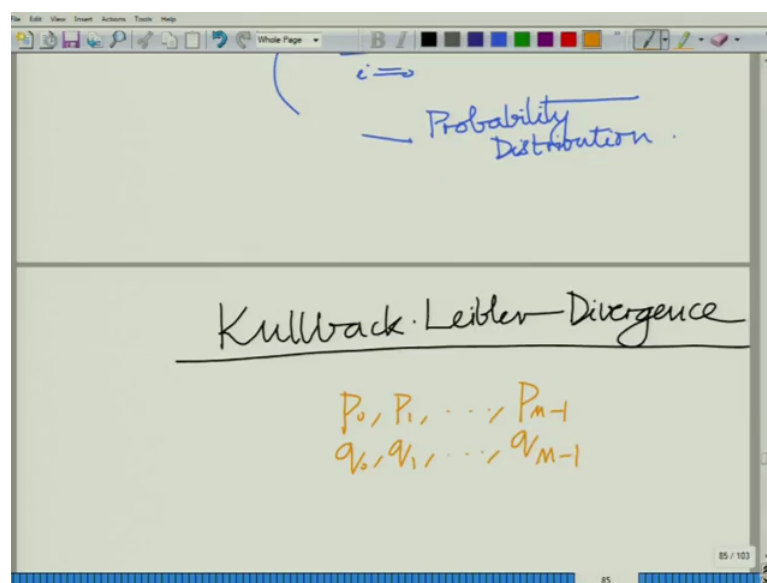
— Probability Distribution.

So, taking summation over the numerator, what we have is summation i equal to 0 to M minus 1 2 to the power of minus l_i divided by summation j equal to 0 to M minus 1 2 to the power of minus l_j . This is equal to 1.

So, what we have is each q_i is positive, that is it is non-negative. 0 is less than equal to q_i and q_i is less than equal to 1 and summation of all q_i is equal to 1 . So, naturally the q_i form a probability mass function on the probability distribution correct. So, q is correct. So, we have 0 less than or equal to each q_i less than or equal to 1 and we also have summation i equal to 0 to M minus 1 q_i equal to 1 . So, q_i is from a probability distribution.

Now, let us use the concept of the Kullback Leibler Divergence which we have seen before.

(Refer Slide Time: 09:17)



So, now we have a probability distribution P_i P_0 , we have a probability distribution P_0 P_1 P_m minus 1 q_0 q_1 q_m minus 1 and therefore, the KL divergence between these Kullback Leibler Divergence between these two probability.

(Refer Slide Time: 10:07)

Handwritten notes on a whiteboard showing the derivation of KL Divergence for Probability Mass Functions and Probability Density Functions.

Top section: KL Divergence For Probability Mass Functions

$$D(P||q) \geq 0$$
$$= \sum_{i=0}^{M-1} p_i \log_2 \left(\frac{p_i}{q_i} \right) \geq 0$$

Annotations: "Replacing integral by sum" (pointing to the summation), "Probability Density & Probability Masses" (pointing to the fraction $\frac{p_i}{q_i}$).

Bottom section: PDFs.

$$D(F_X(x)||g_X(x))$$
$$= \int_{-\infty}^{\infty} F_X(x) \log_2 \left(\frac{F_X(x)}{g_X(x)} \right) dx$$

So, we have two probability distributions, the KL divergence which is defined as D which is equal to summation i equal to 0 to M minus 1 $p_i \log$ to the base 2 p_i divided by q_i . This must be greater than or equal to 0. This is the KL divergence.

Remember we had looked at KL divergence between two probability density functions.

(Refer Slide Time: 10:47)

Handwritten notes on a whiteboard showing the derivation of KL Divergence for Probability Mass Functions and Probability Density Functions.

Top section: KL Divergence For Probability Mass Functions

$$D(P||q) \geq 0$$
$$= \sum_{i=0}^{M-1} p_i \log_2 \left(\frac{p_i}{q_i} \right) \geq 0$$

Bottom section: PDFs.

$$D(F_X(x)||g_X(x))$$
$$= \int_{-\infty}^{\infty} F_X(x) \log_2 \left(\frac{F_X(x)}{g_X(x)} \right) dx$$

We had defined it something like this, the KL divergence between, for instance if F and g which are two probability density functions, F of x and g of x , we had defined it for continuous probability density functions as two probability density functions correspond

into the random variable x as f of x g of x \log_2 to the base 2 f of x . This was a definition, where f of x and g of f are probability density functions.

Now, what we are doing is, we are again doing the same thing for discrete probability mass functions. P and q are probability mass functions correct, probability distributions on discrete symbols $s_0 s_1 s_{\text{minus } 1}$. So, what we have done is, we have taken the definition of KL divergence which we have defined for probability density functions and we have in fact now given the equivalent definition for probability mass functions which is obtained by of course replacing this integral by summation, that is replace continue integral which is the continuous sum replacing my integral by sum and of course, probability density functions by probability mass functions probability densities by probability masses. So, we have this is basically KL divergence for probability and you can say probability mass functions. So, that is what we have over here.

(Refer Slide Time: 13:12)

The image shows a handwritten derivation of the KL divergence $D(P||q) \geq 0$ for discrete probability mass functions. The derivation is as follows:

$$D(P||q) \geq 0$$

$$\Rightarrow \sum_{i=0}^{M-1} P_i \log_2 \left(\frac{P_i}{q_i} \right) \geq 0$$

$$\Rightarrow \underbrace{\sum_{i=0}^{M-1} P_i \log_2 P_i}_{-H(X)} + \sum_{i=0}^{M-1} P_i \log_2 \left(\frac{1}{q_i} \right) \geq 0$$

Below this, the value of q_i is defined as:

$$q_i = \frac{2^{-l_i}}{\sum_{j=0}^{M-1} 2^{-l_j}}$$

$$\Rightarrow \frac{1}{q_i} = \frac{\sum_{j=0}^{M-1} 2^{-l_j}}{2^{-l_i}}$$

Now, therefore, we have KL divergence greater than or equal to 0. Remember we said KL divergence and in fact, we had established proved using log concavity, right using concavity of the log function that KL divergence is always greater than equal to 0 and the same property of course we are done it in the case context of probability density functions, but the same is also valid for probability mass functions, ok.

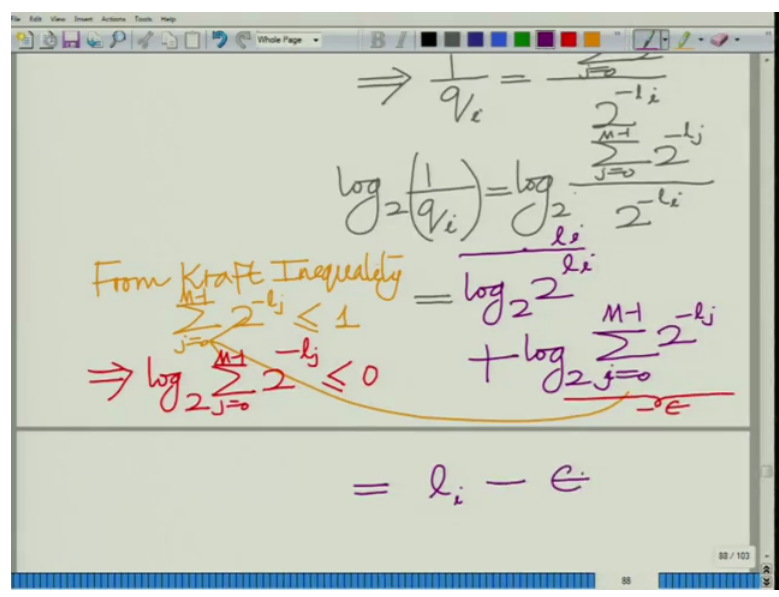
This implies that summation i equal to 0 to M minus 1 $P_i \log$ to the base 2 $P_i \log$ to the base 2 P_i greater than equal to 0 which implies summation i equal to 0 to M minus 1 P_i

$\log_2 \left(\frac{1}{q_i} \right) = -\sum_{j=0}^{M-1} \log_2 \left(\frac{2^{-l_j}}{2^{-l_i}} \right)$

over q_i greater than equal to 0. I can equivalently write in this fashion.

Now, if you look at this quantity $\log_2 \left(\frac{1}{q_i} \right)$, you will realize that this is nothing, but $-H(x)$ and now q_i remember equals to the power of 2^{-l_i} divided by $\sum_{j=0}^{M-1} 2^{-l_j}$ which implies $\frac{1}{q_i}$ is simply $\sum_{j=0}^{M-1} 2^{-l_j}$ divided by 2^{-l_i} .

(Refer Slide Time: 15:06)

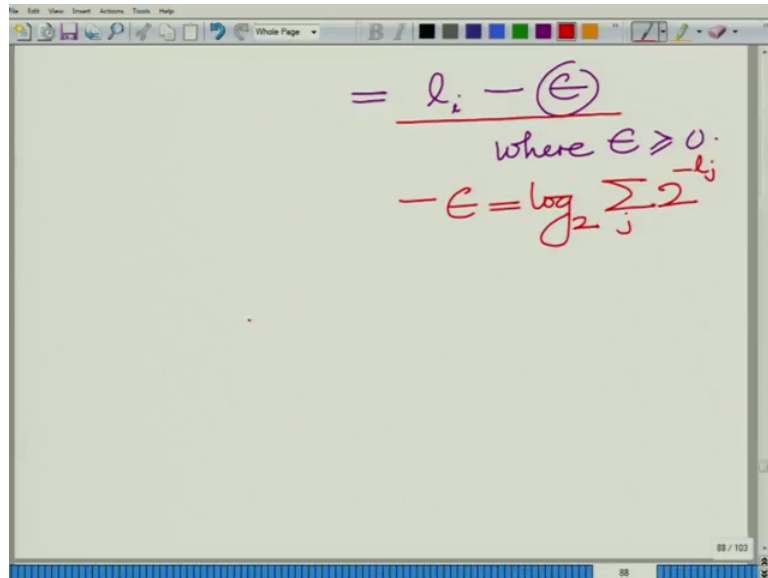


So, $\log_2 \left(\frac{1}{q_i} \right)$, this is nothing, but $\log_2 \left(\frac{1}{q_i} \right)$. Well, $\sum_{j=0}^{M-1} 2^{-l_j}$ divided by 2^{-l_i} can come to the numerators or that becomes 2^{-l_i} . So, this is simply $\log_2 2^{-l_i}$ plus $\log_2 \sum_{j=0}^{M-1} 2^{-l_j}$. Let me just write that also $\log_2 2^{-l_i}$ plus $\log_2 \sum_{j=0}^{M-1} 2^{-l_j}$.

Now, if you look at this, this quantity is nothing, but l_i . So, this will be l_i plus now if you look at this quantity here, we know from Kraft's inequality or from Kraft inequality, we know that $\sum_{j=0}^{M-1} 2^{-l_j} \leq 1$ which means implies $\log_2 \sum_{j=0}^{M-1} 2^{-l_j} \leq 0$. This has to be less than or equal to 0 correct from Kraft inequality. We know that $\sum_{j=0}^{M-1} 2^{-l_j}$ to the

power of 2^{-l_j} is less than or equal to 1. Therefore, if we take the logarithm of that quantity, the log has to be negative.

(Refer Slide Time: 17:46)



The image shows a whiteboard with handwritten mathematical expressions. The top expression is $= l_i - \epsilon$, where ϵ is circled. Below it, it says "where $\epsilon \geq 0$ ". The bottom expression is $-\epsilon = \log_2 \sum_j 2^{-l_j}$.

So, I can write this as minus epsilon, where epsilon is some positive quantity. So, that is the whole idea. So, I can write this as minus epsilon, where epsilon is some positive quantity and that is the reason for that is because, what is minus epsilon. Minus epsilon is log to the base 2 summation over j 2^{-l_j} and therefore, this is equal to l_i minus epsilon. So, now, if you substitute these quantities here, this is equal to l_i minus epsilon and now, what we have interestingly is if we call this equation as star.

(Refer Slide Time: 18:35)

From equation (*) above,

$$-H(X) + \sum_{i=0}^{M-1} P_i(l_i - \epsilon) \geq 0$$

$$\sum_{i=0}^{M-1} P_i l_i \geq H(X) + \sum_{i=0}^{M-1} P_i \epsilon$$

$$\underline{L} = H(X) + \epsilon \Rightarrow 0$$

$$\underline{L} \geq H(X)$$

If we call this equation as star, from star what we have from equation star above, what we have is minus H X plus minus H X plus minus H X plus, well P i into l i minus epsilon minus H X plus summation i equal to 0 to M minus 1 P i into l i minus epsilon is greater than or equal to 0 from KL divergence implies summation i equal to 0 to M minus 1 P i l i greater than or equal to H X plus summation i equal to 0 to M minus 1 P i into epsilon. Epsilon is constant that comes out summation P i which is equal to 1. So, this is simply H X plus epsilon implies this is greater than or equal to and epsilon is positive. Recall epsilon is greater than equal to 0 which implies this is greater than equal to H X, ok

(Refer Slide Time: 20:29)

$\bar{L} \geq H(X)$

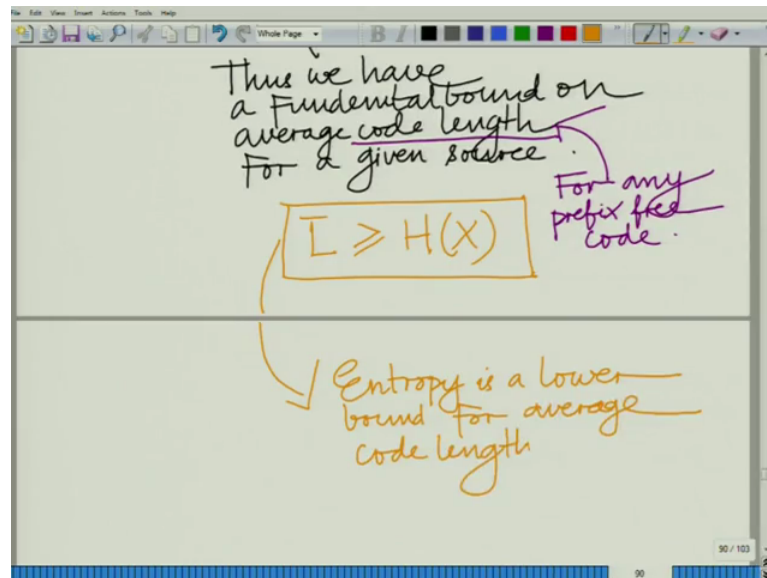
$$\sum_{i=0}^{M-1} P_i \log_2 \left(\frac{1}{P_i} \right)$$

Thus we have a fundamental bound on average code length for a given source.

For any prefix free code.

So, this is nothing, but \bar{L} . So, what we have is \bar{L} greater than or equal to $H(X)$ which is the entropy of the source, that is summation $P_i \log_2 \frac{1}{P_i}$ and there we have a very fundamental result, a fundamental bound on the average length of any prefix free code using the Kraft inequality. We have shown that the average length of any prefix code has to be lower bounded by the entropy of the source and this is a very fundamental elegant and interesting result. So, what this shows us is that, what this is telling us is that the average length of any prefix code, thus we have a fundamental bound on average code length for a given source and remember not any code, average code length, we have to qualify this. This is for any prefix free code, average code length for any prefix free code and we have shown that this is bounded by the entropy.

(Refer Slide Time: 22:17)



So, the average code length \bar{L} greater than equal to $H(X)$, that is average code length that is entropy is a lower bound. What this tells us that entropy is a lower bound for the average code length. So, this is very fundamental. What this says is no matter what prefix free code you design, all right this has to satisfy the Kraft inequality. What that tells us is that the average code length cannot be lower than the entropy. At most, it can be equal to entropy. It has to be greater than and equal to entropy.

So, the efficiency of a code can now be judged by how close is the average code length of the entropy. So, we have a convenient means to judge the efficiency of a code. Remember we said that lower the average length of the code, the more efficient it is and now, we have shown that you cannot arbitrarily reduce it to any non-zero quantity. This is lower bounded by the entropy. So, you can approach entropy. I mean one can desire or one can design a code to approach the entropy as closely as possible, but cannot of course make it lower than the entropy and therefore, the closeness of this average code length to the entropy you can characterize, can be used as a measure to characterize the efficiency of the designed prefix free code for a given source, all right.

So, now how closely can you approach this entropy and how to approach, it is something that we are going to see in the subsequent modules.

Thank you very much.