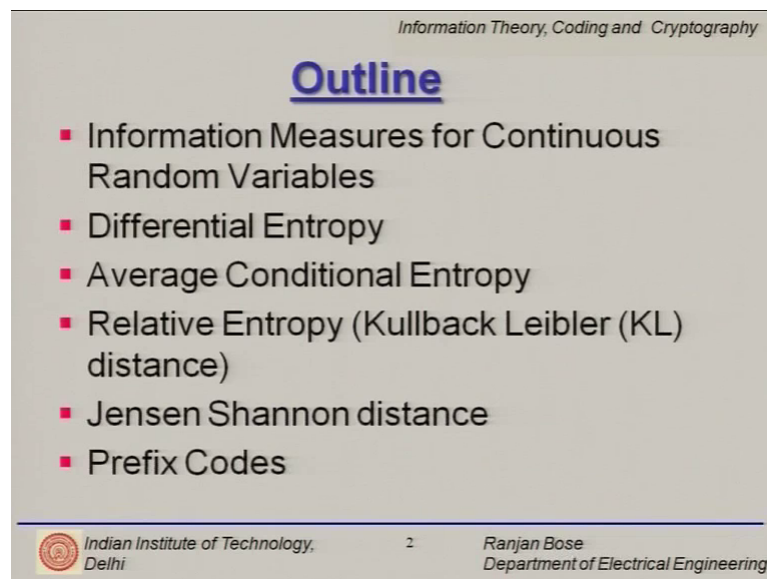


Information Theory, Coding and Cryptography
Dr. Ranjan Bose
Department of Electrical Engineering
Indian Institute of Technology, Delhi

Module - 03
Source Coding
Lecture – 03

(Refer Slide Time: 00:24)



Information Theory, Coding and Cryptography

Outline

- Information Measures for Continuous Random Variables
- Differential Entropy
- Average Conditional Entropy
- Relative Entropy (Kullback Leibler (KL) distance)
- Jensen Shannon distance
- Prefix Codes

Indian Institute of Technology, Delhi 2 Ranjan Bose Department of Electrical Engineering

Hello and welcome to module 3 of source coding. Let us look at the outline of today's talk. Today we are going to consider a very interesting thing, information measures for continuous random variables. So, far we have only dealt with discrete random variables. So, what does it mean? We will look at that. Then we will formulate the notion of differential entropy, we will follow it up with Average Conditional Entropy. We look at Relative Entropy which is also known as Kullback Leibler Distance. We will look at Jensen Shannon's distance and finally, we will look at Prefix Codes ok.

Let us start once again. We will cut off the first three minutes all right. All of you are settled in ok, shall start with the regular formality hello and welcome to module three of source coding. Let us look at the outline of today's talk. We will start with information measures for continuous random variable as opposed to discrete random variables that we talked about in the previous module. Then we will formulate the notion of Differential Entropy. We will look at Average Conditional Entropy for Continuous

Random Variables. Then we will discuss something called Relative Entropy which is a kind of distance measure called the Kullback Leibler Distance. We will then look at Jensen Shannon's distance and finally, we will introduce the notion of Prefix Codes. So, this is our general outline. But first let us start with a quick recap of what we have done already.

(Refer Slide Time: 02:38)


Information Theory, Coding and Cryptography

Self Information

- Consider a discrete random variable X with possible outcomes $x_i, i = 1, 2, \dots, n$.
- The **self information** of the event $X = x_i$ is defined as

$$I(x_i) = \log\left(\frac{1}{P(x_i)}\right) = -\log P(x_i).$$

- When the base of the logarithm is 2 the units of $I(x)$ are in **bits**
- When the base is e , the units are in **nats** (natural units).

 Indian Institute of Technology, Delhi4Ranjan Bose
Department of Electrical Engineering

So, we have already looked at Average Mutual Information. We talked about entropy and its relation to say self-information. We went on to discuss conditional entropy, joint entropy and so and so forth.

So, very quick look at what self information was. So, if you remember we talked about a discrete random variable X with possible outcomes x_i equal to 1, 2, 3 up to n and self information was defined as $I(x_i) = \log \frac{1}{P(x_i)}$ and when the base of the log was 2, the units was in bits, but please note this n does not have to be finite; what if i goes from 1, 2, 3, 4 up to n . Let us look at an example.

(Refer Slide Time: 03:41)

$$P(x_i) = \frac{1}{2^i}$$
$$\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots, \dots$$
$$\sum_{i=1}^{\infty} P(x_i) = 1$$
$$I(x_i) = \log \frac{1}{P(x_i)}$$
$$H(X) = \sum_{i=1}^{\infty} P(x_i) \log \frac{1}{P(x_i)}$$
$$= 2 \text{ bits}$$

ETSC, IIT DELHI

So, suppose I have $P(x_i)$ equal to $1/2^i$. So, the probabilities look like $1/2, 1/4, 1/8$ and so and so forth, but I do not stop at n . I go right up to infinity. Now if you do a basic sanity check overall, you will see that it adds up to 1. So, this is indeed a probability measure and if it is a probability measure, then I should be able to talk about the self-information and what do I do? I plug into the formula. So, if I want to do $I(x_i)$, I will have for a particular case $\log 1/P(x_i)$, but if I want to have the notion of $H(X)$ which is the average self-information, then I have summation of $P(x_i) \log 1/P(x_i)$ and this i will go from 1 to infinity.

Now, if I just plug in the values of the probabilities here and I solve it you can do so. It is a pretty straightforward answer because $\log 1/2^1, \log 1/2^2$ and so and so forth with probabilities multiplied here you will get up a summation and it adds up to 2 bits.

So, please note that even though there are infinite number of possible outcomes here, the net average self-information is bounded. So, even though the formula does not allow you to sum beyond n , you can always have an answer up to a summation up to infinity.

(Refer Slide Time: 05:42)

Information Theory, Coding and Cryptography

Mutual Information

- **Mutual information**
$$I(x_i; y_j) = \log \left(\frac{P(x_i | y_j)}{P(x_i)} \right)$$
- **Observe that**
- $$\frac{P(x_i | y_j)}{P(x_i)} = \frac{P(x_i | y_j)P(y_j)}{P(x_i)P(y_j)} = \frac{P(x_i, y_j)}{P(x_i)P(y_j)} = \frac{P(y_j | x_i)}{P(y_j)}$$
- **Therefore**
- $$I(x_i; y_j) = \log \left(\frac{P(x_i | y_j)}{P(x_i)} \right) = \log \left(\frac{P(y_j | x_i)}{P(y_j)} \right) = I(y_j; x_i)$$

Indian Institute of Technology, Delhi 5 Ranjan Bose
Department of Electrical Engineering

We then looked at mutual information and mutual information between x_i and y_j was defined as $\log \frac{P(x_i | y_j)}{P(x_i)}$ and we made a very interesting observation that $I(x_i; y_j)$ is equal to $I(y_j; x_i)$. So, it is symmetric in nature.

(Refer Slide Time: 06:05)

Information Theory, Coding and Cryptography

Average Mutual Information

- **Definition** The **average mutual information** between two random variables X and Y is given by

$$I(X; Y) = \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) I(x_i; y_j) = \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i)P(y_j)}$$

$I(X; Y) \geq 0$, with equality if and only if X and Y are **statistically independent**.

- **Average mutual information cannot be negative !**

Indian Institute of Technology, Delhi 6 Ranjan Bose
Department of Electrical Engineering

We then went on to define the notion of a average mutual information which is just a average it over the joint probabilities of $P(x_i, y_j)$ and I get capital I, the mutual information X semicolon Y is defined as follows and we know that $I(X; Y)$ is non-negative greater than or equal to 0 and the equality is achieved, if and only if X and

X and Y are statistically independent. Please note, these are discrete random variables. What we have to do today is to look at continuous random variables and whether there is a notion of mutual information for continuous random variables or not.


(Refer Slide Time: 06:56)

Information Theory, Coding and Cryptography

Information Measure for Continuous Random Variables

- The definitions of mutual information for discrete random variables can be directly extended to continuous random variables
- Let X and Y be random variables with joint probability density function (pdf) $p(x, y)$ and marginal pdfs $p(x)$ and $p(y)$.
- The average mutual information between X and Y is defined as follows

$$I(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x)p(y|x) \log \frac{p(y|x)p(x)}{p(x)p(y)} dx dy$$



Indian Institute of Technology,
Delhi

7

Ranjan Bose
Department of Electrical Engineering

So, let us look at this notion of information measure for continuous random variables. Now the definition of mutual information for discrete random variables can directly be extended to continuous random variables, but this is talking about the definition. We will talk about the meaning, whether the meaning can be extended or not in a later slide.


So, let X and Y be random variables with joint probability density functions $p(x, y)$ and marginal pdf's $p(x)$ and $p(y)$, fair enough? We are talking about continuous random variables here. And we then define the average mutual information between x and y as $I(x, y)$ is equal to integration double integration over this is $p(x, y)$ and $\log p(y|x)$ into $p(x)$ over $p(x)$ into $p(y)$ $dx dy$. So, this is the definition of continuous random variable the information measure for that.

(Refer Slide Time: 08:00)

Information Theory, Coding and Cryptography

Information Measure for Continuous Random Variables

- It should be pointed out that the definition of average mutual information can be carried over from discrete random variables to continuous random variables, but the concept and **physical interpretation cannot**.
- The reason is that the information content in a continuous random variable is actually **infinite**, and we require infinite number of bits to represent a continuous random variable precisely.
- The self information and hence the entropy is **infinite**.
- To get around the problem we define a quantity called the **differential entropy**.

 Indian Institute of Technology,
Delhi8Ranjan Bose
Department of Electrical Engineering

So, let us put a word of caution. Even though, we could extend the definition; the physical interpretation probably cannot be stretched. So, we should point out that the definition of average mutual information can be carried over from discrete random variables to continuous random variables, but the concept and the physical interpretation cannot.

What was the physical interpretation for a discrete random variable? Well, when we defined an average mutual information between X and Y , the basic physical interpretation was having observed Y , what can you say about X ? In general on an average, what can you say? Maybe you can say something? Maybe you can say nothing. So, X occurrence of X communicate something about occurrence of Y and vice versa and that is basically captured by $I(X; Y)$.

Now, you would like to say the same thing about continuous random variable, but unfortunately that is not the case. The reason is that the information contained in a continuous random variable is actually infinite. I mean what is the best way to look at it? Take a sample, represent it correctly. How many decimal points do you need? You can sample it and say it as 2.309921729, but you keep going. It is a continuous random variable. It is a point on the real line. So, you keep going and you really need infinite number of bits even to represent a single sample value, let alone the entire function.


So, the information content truly is infinite and therefore, we cannot really go on dealing with infinite information all the time. Let alone compare and what one communicates about the other random variable. So, we have just seen that the self-information entropy is infinite and we have to get around this problem and we define a new quantity called differential entropy. So, what each one of the random variables encompass infinite information? What about the difference? Maybe the difference is not infinite.

(Refer Slide Time: 10:47)

Information Theory, Coding and Cryptography

Differential Entropy

- The **Differential Entropy** of a continuous random variable X is defined as
$$h(X) = - \int_{-\infty}^{\infty} p(x) \log p(x)$$
- Again, it should be understood that there is **no physical meaning** attached to the above quantity

 *Indian Institute of Technology, Delhi* 9 *Ranjan Bose*
Department of Electrical Engineering

So let us define the differential entropy for a continuous random variable. So, the differential entropy is defined as h of X , mind you h is small h lowercase h as opposed to uppercase h for discrete random variables and is defined as an integration minus infinity to infinity $p(x) \log p(x)$ with a negative sign; as usual if the base of the log is 2, then the units are in bits. Same word of caution, there is no physical meaning attached to it.

If you remember in the earlier lectures, we tossed a coin and if it was a fair coin, we said that the average self-information for that source tossing a fair coin was 1 bit and it made sense because you needed 1 bit to represent either head or a tail. So, there is a strong physical interpretation; however, for a continuous random variable, we have no such luck.

Student: (Refer Time: 11:54).


Yes we will talk about this differential part very shortly, it will come.

(Refer Slide Time: 12:00)

Information Theory, Coding and Cryptography

Some properties of Differential Entropy

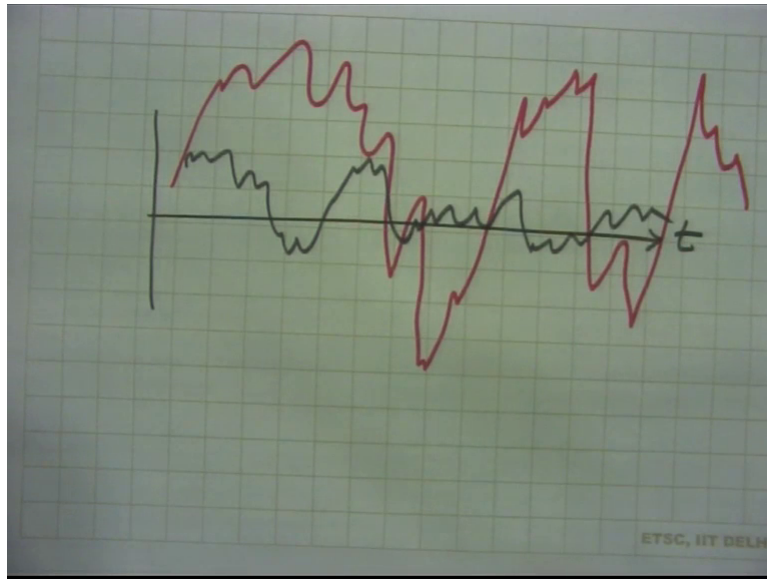
- Chain Rule:
$$h(X_1, X_2, \dots, X_n) = \sum_{i=1}^n h(X_i | X_1, X_2, \dots, X_{i-1})$$
- Translation does not alter the differential entropy
$$h(X + c) = h(X),$$
- And
$$h(aX) = h(X) + \log|a|.$$

 Indian Institute of Technology, Delhi 10 Ranjan Bose
Department of Electrical Engineering

Now let us look at the properties of differential entropy. So, we talk about this chain rule where $h(X_1, X_2, \dots, X_n)$; this K should be dot, dot, dot, dot up to X_n and it can be represented in a conditional form as $h(X_i | X_1, X_2, \dots, X_{i-1})$. The other property is that the translation of the random variable X does not change the differential entropy and I am relieved because if I add a constant, it really does not add to the randomness of X and the information content is primarily the measure for randomness.

So, translation it turns out; does not alter the differential entropy. And if you just multiply by a scalar, so the differential entropy of not X , but a times X is just a d c shift. How do you visualize this practically? Imagine a continuous random variable. So, let us plot this and try to get a physical interpretation for this.

(Refer Slide Time: 13:36)



Suppose, this is just one capture of my random process because I have put a time axis here and suppose I pass it through an amplifier and multiply it with a . And now I get another capture of the same thing ok. So, if you can see just the physical observation tells you that the second random variable which is a times multiplied with the first one has a higher level of randomness.

The variance has gone up and consequently the information content has to go up, but how does it go up? It gives you a basic dc shift here. So, if you look at $h(aX)$ equal to $h(X)$ plus \log absolute value of a is a . There is a strong interpretation attached to the scaling factor.

(Refer Slide Time: 14:54)


Information Theory, Coding and Cryptography

Relative Entropy

- An interesting question to ask is how similar (or different) are two probability distributions?
- Relative entropy is used as a measure of distance between two distributions.
- The **Relative Entropy** or **Kullback Leibler (KL) distance** between two probability mass functions $p(x)$ and $q(x)$ is defined as

$$D(p \parallel q) = \sum_{x \in X} p(x) \log \left(\frac{p(x)}{q(x)} \right)$$

- It can be interpreted as the expected value of $\log \left(\frac{p(x)}{q(x)} \right)$

 Indian Institute of Technology, Delhi 11 Ranjan Bose
Department of Electrical Engineering

Now, another interesting question to ask is, how similar or if you are a pessimist; how different are two probability distributions fine? So, we talk about a notion of relative entropy as a measure of distance between two distributions, distance measure. Well I tell you, how different to probability distributions can be. If I ask you a question, I give you Gaussian and I say, how different is it with respect to another Gaussian or how different is a distribution with respect to a Rayleigh distribution? These questions are valid and we would like to have an answer to that.

So, we talk about this notion of relative entropy or by the two guys who defined Kullback and Leibler. This is called the Kullback Leibler distance between two probability mass functions; $p(x)$ and $q(x)$. It is defined as $D(p \parallel q)$ is nothing, but $p(x) \log p(x)$ by $q(x)$ ok. This you can see is nothing, but an expected value of $\log p(x)$ over $q(x)$. So, what they have done is taken the two quantities $p(x)$ and $q(x)$, if you just think hard enough it is nothing, but $p(x) \log p(x)$ minus $p(x) \log q(x)$, but if you remember $p(x) \log p(x)$ is a measure of the average self-information. So, therefore, the notion of relative entropy, how is ones information relative to the other?

Now.

Student: Sir, so what is the physical significance of $p(x) \log p(x)$? It is actually a $p(x) \log p(x)$ there is a self-information and minus $p(x) \log p(x)$.

(Refer Slide Time: 17:21)

$$D(p||q) = \sum p(x) \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$
$$= \sum p(x) \log p(x) - \sum p(x) \log q(x)$$

The image shows a handwritten derivation of the Kullback-Leibler divergence formula on a grid background. The first line is $D(p||q) = \sum p(x) \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$. The second line is $= \sum p(x) \log p(x) - \sum p(x) \log q(x)$. A horizontal double-headed arrow is drawn under the first term of the second line. A vertical arrow points down from the second term of the second line to the label $q(x)$.

Yes. So, let us repeat the question. The question being asked is we have defined the Kullback Leibler distance as follows and with a little bit of imagination, you can see that this is nothing, but fair enough. Now, this has a physical interpretation. It tells me something about the average self-information of p of x ; however, with a little bit of distortion because this is not q of x , otherwise I was really finding out the difference between the self-information, average self-information of p x and q x . So, they have kind of distorted it with a purpose and we will look at why it tells you because I am looking at a similarity measure between the two.

Now, why would they define it like this? It is very simple. Suppose p x and q x are identical, then you have \log of 1 and clearly it is 0 and we are relieved to find the two distributions which are identical, their distance is 0. The more different they are, the larger should be the value of this distance and this notion; so, this is the basic definition. Logarithmic measure is required because right from the beginning, we have argued by a \log measure for information is the only logical way to go and they have found out the difference between $\log p$ x and $\log q$ of x , but that just would not do. So, they have averaged it to give you the relative information measure.

(Refer Slide Time: 19:37)


Information Theory, Coding and Cryptography

Relative Entropy

- Does Kullback Leibler distance follow symmetry property?
- Is $D(p \parallel q) = D(q \parallel p)$?
- Lets check for $\sum p(x) \log \frac{p(x)}{q(x)} = \sum q(x) \log \frac{q(x)}{p(x)}$
- We find

$$\sum p(x) \log p(x) - \sum p(x) \log q(x) \neq \sum q(x) \log q(x) - \sum q(x) \log p(x)$$

Does NOT follow symmetry property

 Indian Institute of Technology,
Delhi12Ranjan Bose
Department of Electrical Engineering

As we have seen there are clearly problems with this and our next slide will tell you, what the problems are. There are many problems. First is, if we talk about this to be a distance that is KL distance Lullback Leibler distance, then distance has 3 properties. Number 1 non-negative; so greater than or equal to 0; Number 2, symmetry property; distance between a to b is the same as b to a and triangle inequality right. The sum of the two sides of a triangle should be greater than the third side right.

So, first thing you can easily verify that the Kullback leibler distance is indeed non negative, but let us look at the other two properties. So, the question we ask ourselves is, Does the KL distance really follow the symmetry property? Just now couple of minutes back we had this discussion and we saw that there was some asymmetry. So, why do not we test it out? Question we are asking is $D p \text{ parallel } q$ equal to $D q \text{ parallel } p$? So, we plug in the values and we check for whether this is the definition of $D p \text{ parallel } q$. Is it really equal to $D q \text{ parallel } p$? And if we expand it out, it does not take much effort to see that it is in really not true.

So, the first conclusion is this Kullback Leibler distance, the distance is a misnomer. It is a wrong thing to call it a distance, even though the proposals have called it a distance; it does not follow the symmetry property of a distance. What does it mean? If I say, how different is distribution p from distribution q? My answer will differ with respect to how different is a distribution q with respect to p. Nonetheless it is used in practical life.


(Refer Slide Time: 21:41)

Information Theory, Coding and Cryptography

Relative Entropy

- Does Kullback Leibler distance follow Triangle Inequality?
- Is $D(p \parallel q) + D(q \parallel r) \geq D(p \parallel r)$?
- Lets check for $\sum p(x) \log \frac{p(x)}{q(x)} + \sum q(x) \log \frac{q(x)}{r(x)} \geq \sum p(x) \log \frac{p(x)}{r(x)}$
- We find $\sum (-p(x) + q(x)) \log \frac{q(x)}{r(x)} \geq 0$
- This relation does not hold if $p(x) > q(x)$.

Does NOT follow Triangle Inequality

 Indian Institute of Technology,
Delhi13Ranjan Bose
Department of Electrical Engineering

Now, we look at the second property of the distance. Does it follow the triangle inequality? Does it satisfy that condition? What does it mean? Let us take 3 distributions. So, we talk about distance between p and q, distance between q and r and distance between p and r. And the question we are asking is the distance p parallel q plus q parallel r, sum of the two sides of the triangle greater than or equal to the third side and again we plug in this values and indeed we find that for simply p x greater than q of x, we would see that this relation does not hold; it is not universally true.

So, we get that this relative entropy which is nothing, but the Kullback Leibler distance does not follow the triangle inequality. Let us look at an example

(Refer Slide Time: 22:47)


Information Theory, Coding and Cryptography

Example

- Consider a Gaussian distribution $p(x)$ with mean and variance given by (μ_1, σ_1^2) ,
- Consider another Gaussian distribution $q(x)$ with mean and variance given by (μ_2, σ_2^2)
- We find the KL distance between two Gaussian distributions as

$$D(p \parallel q) = \frac{1}{2} \left[\frac{\sigma_1^2}{\sigma_2^2} + \left(\frac{\mu_2 - \mu_1}{\sigma_2} \right)^2 - 1 - \log_2 \left(\frac{\sigma_1^2}{\sigma_2^2} \right) \right]$$

- The distance becomes zero when the two distributions are identical, i.e., $\mu_1 = \mu_2$ and $\sigma_1^2 = \sigma_2^2$
- It is interesting to note that when $\mu_1 \neq \mu_2$ the distance is minimum for $\sigma_1^2 = \sigma_2^2$

 Indian Institute of Technology, Delhi 14 Ranjan Bose
Department of Electrical Engineering

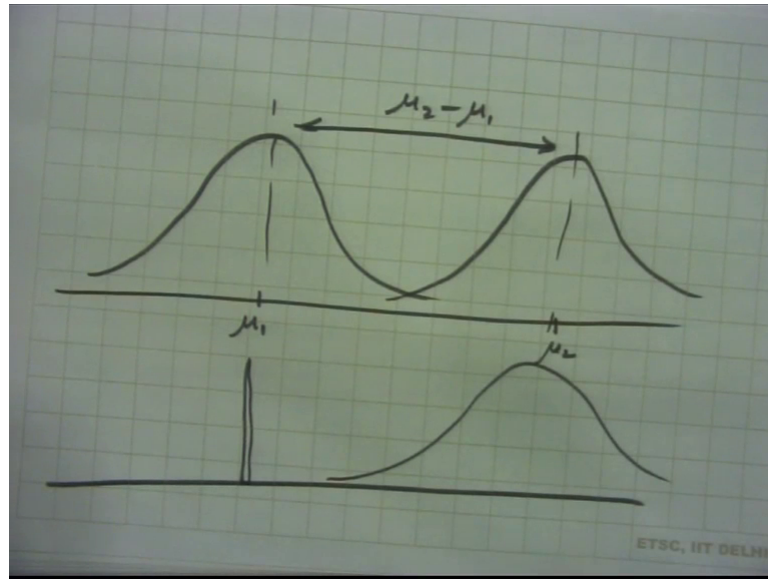
So, we are curious to find out, how similar two Gaussian random variables; call two Gaussian distributions are? So, what are my two distributions; $p(x)$ and $q(x)$, then it is an interesting exercise, $p(x)$ has mean μ_1 and variance σ_1^2 ; $q(x)$ has a mean μ_2 and a variance σ_2^2 . So, how different are they? We simply find out $D(p \parallel q)$, you plug in the value and you can do a little bit of math and you get this expression.

Student: (Refer Time: 23:30) by doing you summation instead of integration using (Refer Time: 23:34) presents a plays a continuous case.

Right question being asked is, why use using summation way and integration? We have defined as a mass function, you can use integration right. As we have defined in the earlier cases. So, if you look at the two Gaussian distributions $p(x)$ and $q(x)$, you get the relative entropy $D(p \parallel q)$ as follows. So, as expected it is a function of σ_1^2 , σ_2^2 , μ_1 , μ_2 and so and so forth. You will note that $D(p \parallel q)$ is not the same as $D(q \parallel p)$. Nonetheless some interesting observations can be seen.

So, when does this distance become zero? Well obviously, when $\mu_1 = \mu_2$ and $\sigma_1^2 = \sigma_2^2$, this distribution will be 0 that is they are identical. But if you say that suppose $\sigma_1^2 = \sigma_2^2$, but $\mu_1 \neq \mu_2$ right, then the distance is minimum. What does it mean? It means that the Gaussian spread is the same.

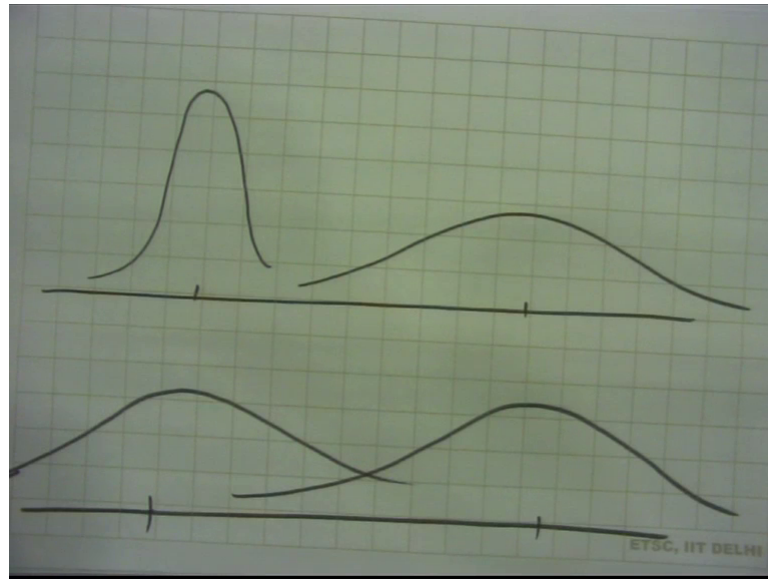
(Refer Slide Time: 25:00)



So, if you look at these two distributions. So, variance is the same, but they are apart right. In this case KL distance tells us that this is minimum and it is only a function of μ_1 and μ_2 fine. But it is minimized when the variances are the same, but if you go ahead and say that if either σ_1^2 tends to 0 or σ_2^2 tends to 0, then you see; then we are in a fix because the distribution the distance becomes infinite.

So, this is the case, when either 1 becomes in the limiting case, a delta function. So, they are very different. You can visibly see that this, these are different and somewhere in the middle if you have one distribution like this, the other one is like this; variances are different, means are different. They are clearly different.

(Refer Slide Time: 26:22)



On the other side if I put variances are the same, means are different; then the distance scale distance is less. So, these two probability mass functions, if the distributions are much more similar to each other than these two distributions. Not only it is visibly correct, you can also calculate it from the expression for KL distance ok. So, there is a strong physical connotation attached to how similar or different two distributions are.

(Refer Slide Time: 27:20)


Information Theory, Coding and Cryptography

Average Mutual Information

- The average mutual information can be seen as the relative entropy between the joint distribution, $p(x, y)$, and the product distribution, $p(x)p(y)$, i.e.,

$$I(X; Y) = D(p(x, y) \| p(x)p(y))$$

- We note that, in general, $D(p \| q) \neq D(q \| p)$
- Thus, even though the relative entropy is a distance measure, it does not follow the symmetry property of distances.
- To overcome this, another measure, called the **Jensen Shannon distance**, is sometimes used

 Indian Institute of Technology, Delhi16Ranjan Bose
Department of Electrical Engineering

Now, we move on to the average mutual information for continuous random variable and it is a very elegant way to define it. It is $I(X; Y)$ is nothing, but the distance

between the joint distribution and the product of the distribution. So, it is a relative entropy between the joint distribution $p \times y$ and the product of the distribution $p \times p_y$. Clearly if we have $p \times y$ is a same as $p \times p_y$, we get a direct obvious conclusion right. If they are independent, then the we have a one notion. But in general we know that this distance measure is not symmetric. So, Jensen Shannon came up with an alternate measure of this distance which is symmetric and let us define it also.

(Refer Slide Time: 28:30)

Information Theory, Coding and Cryptography


Jensen Shannon Distance

- The **Jensen Shannon distance** between two probability mass functions $p(x)$ and $q(x)$ is defined as

$$JSD(p \parallel q) = \frac{1}{2} D(p \parallel m) + \frac{1}{2} D(q \parallel m)$$

where $m = \frac{1}{2}(p+q)$

- If the base of the logarithm is 2, then, $0 \leq JSD(p \parallel q) \leq 1$
- Jensen Shannon distance is sometimes referred to as **Jensen Shannon divergence** or **Information Radius** in literature.

 Indian Institute of Technology, Delhi
17
Ranjan Bose
Department of Electrical Engineering

So, we are now looking at overcoming the shortcoming of the Kullback Leibler distance. We talk about the Jensen Shannon distance. Again it is between two probability mass functions $p \times x$ and $q \times x$, but it is defined as JSD start standing for Jensen Shannon distance is half of D, relative entropy $p \parallel m$ plus half $D \ q \parallel m$ where m is an intermediate point between p and q .

So, this you can check for yourself is symmetric. So, $JSD \ p \parallel q$ is equal to $JSD \ q \parallel p$. So, Jensen Shannon distance is also referred to as Jensen Shannon divergence or information radius and literature. So, those terms are used interchangeably. And this value is limited between 0 and 1 under the condition that the base of the log is 2.

So, why are we doing all of this? There has to be some practical utility for all of these mathematical exercises. So, far we have built in some tools to understand, what is the best way to represent symbols?


(Refer Slide Time: 30:00)

Information Theory, Coding and Cryptography

Efficient Representation of Symbols

- Lets explore **efficient representation** (efficient coding) of symbols generated by a source.
- The primary motivation is the **compression of data** due to efficient representation of the symbols.
- Suppose a discrete **memoryless** source (DMS) outputs a symbol every t seconds.
- Each symbol is selected from a **finite set of symbols** $x_i, i = 1, 2, \dots, L$, occurring with probabilities $P(x_i), i = 1, 2, \dots, L$.
- The **entropy** of this DMS in bits per source symbols is

$$H(X) = - \sum_{j=1}^L P(x_j) \log_2 P(x_j) \leq \log_2 L.$$

 Indian Institute of Technology,
Delhi18Ranjan Bose
Department of Electrical Engineering

So, in the next few slides, we will explore; what are the ways to efficiently represent in other words efficiently, code symbols generated by a source? In earlier lectures, we talked about, what could be a source? It could be a man tossing a fair coin and shouting out 1 0 0 1 1 0 or he could be tossing 2 independent coins and saying 1 0 0 0 0 0 1 and so and so forth or he could be tossing unfair coins or he could be typing his SMS, that is a source or a monkey typing of keyboard, that is the source; all of these are sources.

They generate symbols it could be a b x y z p q or 1 0 0 1 or it could be voltages all of them are symbols for me and I need to represent them efficiently. What is the primary motivation? Data compression, efficient representation leads to compression of data. Suppose we have a discrete memory less source and its outputs assemble every t seconds. So, each symbol is selected from a finite set of symbols.


So, again we make this assumption that the set is finite, we will move to infinite sets also and we can always define the average self-information or entropy of this discrete memoryless source from the theory we have developed so far. And you can always show that this is upper bounded by \log to the base 2 L . So, let us get some definitions in order.

(Refer Slide Time: 32:00)

Information Theory, Coding and Cryptography

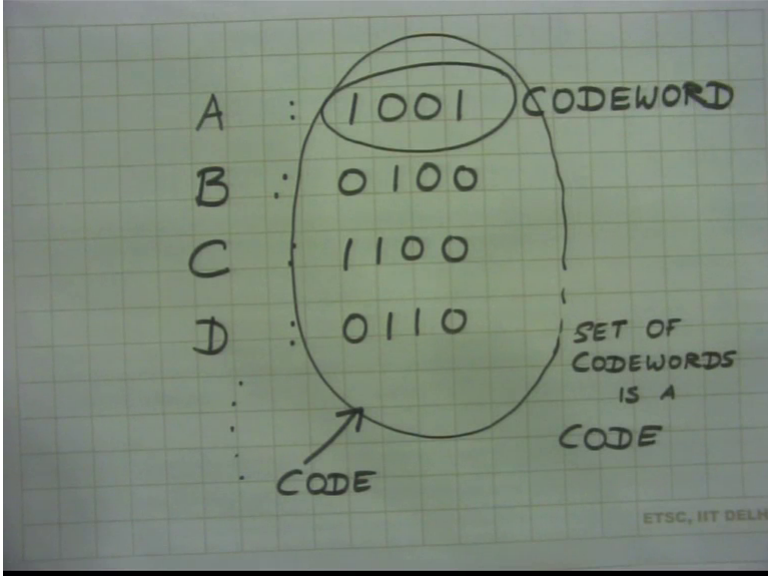
What is a code?

- A **code** is a set of vectors called **codewords**
- As we saw earlier, to encode the letters of the English alphabet, we need $R = \lfloor \log_2 26 \rfloor + 1 = 5$ bits.
- The fixed length code for the English alphabet suggests that each of the letters in the alphabet is equally important (probable) and hence each one requires 5 bits for representation.
- However, we know that some of the letters are less common (x, q, z etc.) while some others are more frequently used (s, t, e etc.).
- It appears that allotting equal number of bits to both the frequently used letters as well as not so commonly used letters is *not* an efficient way of representation (coding).

 Indian Institute of Technology, Delhi19Ranjan Bose
Department of Electrical Engineering

We will represent symbols. So, suppose we want to represent the letters of the English alphabet, so A can be represented as 1 0 0 1, B can be represented as 0 1 0 0 and so and so forth. So, a vector which represents a symbol is actually called a codeword. So, if A is represented as 1 0 0 1, then the codeword for A is 1 0 0 1 ok.

(Refer Slide Time: 32:44)



A : 1001 CODEWORD
B : 0100
C : 1100
D : 0110
...
CODE

SET OF CODEWORDS IS A CODE

ETSC, IIT DELHI

So, let us look at it in a simple manner. So, A it is represented as B and so and so forth. So, this is a codeword; however, the set the set of codewords is called a code. So, this is a code. In this example a code is a set of codewords, 1 codeword for A, 1 for B, 1 for C

and 1 for D. So, code is a set of vectors called codewords. So, we have seen in an earlier example that we can encode the letters from an English alphabet and we need a certain number of bits. Here I have shown in my earlier example, a 4 bit representation, I can have a 5 bit representation clearly. If there are 26 characters in the English alphabet, then I have to have minimum 5 number of bits, otherwise I cannot have unique representation.

And we have also observed in our previous lectures that certain number of alphabets, certain alphabets are more frequent a, e, s, t; some are less frequent x, q, z, j and so, it does not make sense to represent all the alphabets with equal number of bits. What is the rationale behind it? Bits are expensive. To transmit a bit, I need power, I need bandwidth, I need time. All three are very precious quantities for me in my modern communication systems. If x does not appear frequently, why should I allocate certain number of bits? Or if A appears more frequently, maybe I should use fewer number of bits to represent it.

(Refer Slide Time: 35:17)

Information Theory, Coding and Cryptography

What is a variable length code?

- It appears that allotting equal number of bits to both the frequently used letters as well as not so commonly used letters is *not* an efficient way of representation (coding).
- Intuitively, we should represent the more frequently occurring letters by fewer numbers of bits and represent the less frequently occurring letters by larger number of bits.
- In this manner, if we have to encode a whole page of written text, we might end up using fewer number of bits overall.
- When the source symbols are not equally probable, a more efficient method is to use a **Variable Length Code (VLC)**

Indian Institute of Technology, Delhi 20 Ranjan Bose
Department of Electrical Engineering

This brings us to the notion of variable length code. Remember a code is a set of codewords. So, in the example that I have shown earlier, I only had fixed length code where A, B, C, D all of them had four bits, but now maybe we should be looking at and the notion of variable length code. What is the necessity for this? When the source symbols are not equally probable, it makes sense to use fewer number of bits to represent more frequently occurring symbols and vice versa and therefore, we would use the notion of variable length codes.

(Refer Slide Time: 36:05)

Information Theory, Coding and Cryptography


Example

- Suppose we have only the **first eight letters** of the English alphabet (A – H) in our vocabulary.
- The **fixed length code** for this set of letters would be

Letter	Codeword	Letter	Codeword
A	000	E	100
B	001	F	101
C	010	G	110
D	011	H	111

- A **variable length code** for the same set of letters can be

Letter	Codeword	Letter	Codeword
A	00	E	101
B	010	F	110
C	011	G	1110
D	100	H	1111

 Indian Institute of Technology, Delhi 21 Ranjan Bose
Department of Electrical Engineering

Let us look at an example. Suppose we have only the first 8 letters of the English alphabet A to H in our vocabulary. So, you go only you going to use A, B, C, D, E, F, G and H. H it is a very convenient example, log to the base 8, log to the base 2 of 8 gives me 3. So, conveniently I can have fixed length code, 3 bits per symbol right from 0 0 0 0 1 up to 1 1 1 and I have got the first fixed length code. It is no brainer ok.

At the same time I would say, hey how about representing them with unequal number of bits. Why do not we do a variable length code? So, I have an example. A is 0 0, B is 0 1 0, C is 0 1 1 so and so forth. I run out of certain number of bits and then I have to use 4 bits also, but I am not too bad. I am using 2 bits and 4 bits and 3 bits. So, maybe I will come out.

(Refer Slide Time: 37:14)


Information Theory, Coding and Cryptography

Example

- Suppose we have to code the series of letters:
"A BAD CAB".
- The fixed length and the variable length representation of the pseudo sentence would be

Fixed Length Code	000 001 000 011 010 000 001	Total bits = 21
Variable Length Code	00 010 00 100 011 00 010	Total bits = 18

- Note that the variable length code uses **fewer numbers of bits**
- This is because the letters appearing more frequently in the pseudo sentence are represented with fewer numbers of bits.

 *Indian Institute of Technology, Delhi* 22 *Ranjan Bose*
Department of Electrical Engineering

So, how about looking at a practical example? Since I have only 8 letters in my alphabet, let us make a small sentence. A BAD CAB, it only uses only 4 first 4 letters. So, if you take a bad cab and use the first code, what is a code? Is a set of codewords. So, the table 1 is a code, table two is also a code. Table 1 is a fixed length code, table 2 is a variable length code.

So, if you use the fixed length code, you have how many characters? 1, 2, 3, 4, 5, 6, 7, 8; so 7 into 3, 21 bits is what I expect from the fixed length code. But if you look at the variable length code, I have got fewer number of bits; looks like we have a winner right. We have been able to save 3 bits percentage wise that is not too bad.

(Refer Slide Time: 38:22)

Information Theory, Coding and Cryptography


Example

- We look at yet another variable length code for the first 8 letters of the English alphabet

Letter	Codeword	Letter	Codeword
A	0	E	10
B	1	F	11
C	00	G	000
D	01	H	111

- This second variable length code appears to be more efficient in terms of representation of the letters.

Variable Length Code 1	00 010 00 100 011 00 010	Total bits = 18
Variable Length Code 2	0 1001 0001	Total bits = 9

 Indian Institute of Technology, Delhi23Ranjan Bose
Department of Electrical Engineering

So, let us look at another variable length code for the first 8 letters of the English alphabet because this time I am really excited, I am going to squeeze out fewer number of bits will be used to represent ABC and D so and so forth.

Student: So, if it always (Refer Time: 38:40) variable length code will reduce the number of bits may be in some cases you can also increase.

So the question will asked is, does variable number of code, the variable length of codewords always reduce the representation? So, answer depends on a, how efficient is our code and b, what is the frequency with which the letters are appearing. So, we will give an example where a variable length code can actually lead to expansion and not a compression. So, that can also happen, but we talk about on an average. Yes, on an average a variable length code.

So, the most 3, most important 3 words are on an average. What are you saying on an average? Yes on an average a variable length code, if designed properly will be able to compress, but once in a while for a very special set of input characters, it can lead to expansion. But hey we carry out our communication over millions of bits and it averages out. So, we turn out a winner.

(Refer Slide Time: 40:00)


Information Theory, Coding and Cryptography

Uniquely decodable !

Letter	Codeword	Letter	Codeword
A	0	E	10
B	1	F	11
C	00	G	000
D	01	H	111

Variable Length Code 2	0 1001 0001	Total bits = 9
-------------------------------	-------------	-----------------------

- However there is a problem with VLC2.
- Consider the sequence of bits 0 1001 0001 which is used to represent **A BAD CAB**
- We could regroup the bits in a different manner to have [0] [10][0][1] [0][0][01] which translates to **A EAB AAD**
- or we can decode the vector as [0] [1][0][0][1] [0][0][0][1] which stands for **A BAAB AAAB !**
- **Not uniquely decodable**

 Indian Institute of Technology, Delhi24Ranjan Bose
Department of Electrical Engineering

So, if you complete this example, if this much more constricted code A only 1 bit, B only 1 bit, C 2 bits and so and so forth and only when I am first I am run out of 1 and 2 bits, I go to 3 and 3 bits. And then I again encode a bad cab, I just have a 9 bits, but the problem is the decoding part. You cannot decode it.

At least not uniquely because if you look at it a bad cab is truly 0 1 0 0 1 0 0 0 1, but it can be broken up into different sets and it becomes A BAD AAD or A BAAB AAAB. So, clearly it is not uniquely decodable. Since I do not know a priori, how long are the codewords; since they are variable length? I do not know at the decoding end. What was actually sent? I am in trouble unless I have a smart way to overcome this deficiency.


(Refer Slide Time: 41:13)

Information Theory, Coding and Cryptography

Revisit VLC1 and VLC2

VLC1				VLC2			
Letter	Codeword	Letter	Codeword	Letter	Codeword	Letter	Codeword
A	00	E	101	A	0	E	10
B	010	F	110	B	1	F	11
C	011	G	1110	C	00	G	000
D	100	H	1111	D	01	H	111

- **VLC2:** We have no clue where the codeword of one letter (symbol) ends and where the next one begins, since the lengths of the codewords are variable.
- However, this problem does not exist with the **VLC1**.
- *It can be seen that no codeword forms the prefix of any other codeword.*

 Indian Institute of Technology,
Delhi25Ranjan Bose
Department of Electrical Engineering

So, let us revisit variable length code 1 and variable length code 2. If you remember variable length code 1, did give us a reduction, but not by much and variable length code 2 give us a major reduction, but it could not be uniquely decoded. So, you look at variable length code 1, can it be uniquely decoded? Well the answer is yes. The answer is yes because if we make a simple observation that no codeword is a prefix of any other codeword.

So, variable length code 1 has a very very unique characteristic. What is it? No codeword is a prefix that is no other codeword starts with any other codeword. No codeword is a prefix of any other codeword and hence decoding is absolutely unique and instantaneous. The moment you find a valid codeword has come, you declare the result because there is no point in looking further because no codeword is a prefix of any other codeword.

(Refer Slide Time: 42:37)


Information Theory, Coding and Cryptography

Prefix Codes

- Observation: No codeword of VLC2 forms the prefix of any other codeword.
- This is called the **prefix condition**.
- So, as soon as a sequence of bits corresponding to any one of the possible codewords is detected, we can declare that symbol decoded.

A **Prefix Code** is one in which no codeword forms the prefix of any other codeword.

- Such codes are also called **Instantaneous Codes**.

 Indian Institute of Technology,
Delhi26Ranjan Bose
Department of Electrical Engineering

This interesting property is called the prefix condition and the decoding strategy is very simple as soon as a sequence of bits corresponding to any one valid codeword is detected, we declared the resemble being decoded.

So, we now formally define, what is a prefix code? A prefix code is one, in which no codeword forms the prefix of any other codeword and since I can instantaneously declare the results as we go along. These are also called instantaneous codes. So, let us summarize what we have learned today.

We started with information measure for continuous random variables; we made a distinction between discrete random variables, continuous random variables. We could extend the definition, but not the physical interpretation because the average self-information of contained in a continuous random variable is actually infinite.

(Refer Slide Time: 43:48)

Information Theory, Coding and Cryptography

Summary

- Information Measures for Continuous Random Variables
- Differential Entropy
- Average Conditional Entropy
- Relative Entropy (Kullback Leibler (KL) distance)
- Jensen Shannon distance
- Prefix Codes

Indian Institute of Technology, Delhi 27 Ranjan Bose
Department of Electrical Engineering

We then talked about differential entropy and we observe that even though, you cannot really talk about in real terms; the information content of a continuous random variable. You can talk in terms of a differential mode. So, X is a continuous random variable, Y is a continuous random variable. So, $h(X) - h(Y)$ may not be infinite. Even though $h(X)$ and $h(Y)$, each of them are infinite and hence this word name differential. So, it only makes sense any meaning, when it is taken as a difference of $2h(X)$ and $h(Y)$ ok; hence the name differential.

We talked about the average conditional entropy for continuous random variable, then we have raised the very interesting question. How do you say two probability distributions are similar or different? What is the similarity measure? And we talked about relative entropy is also called the Kullback Leibler distance, we also observed that it is a distance measure, but it is a pseudonym, it is a misnomer only the non-negativity is satisfied. It does not follow the triangle inequality nor does it follow the symmetry property of a distance measure.

So, to overcome that we talked about the Jensen Shannon distance which is symmetric and finally, we introduced the notion of prefix codes; that is the first step towards efficient representation of symbols and ultimately we have look at ways to compress data, speech, images what have you right. So, that is one of the fundamental

contributions of source coding. It let us you calculate the theoretical limits to which I can compress my data and no further. That is where we will go in our next module.

Thank you.