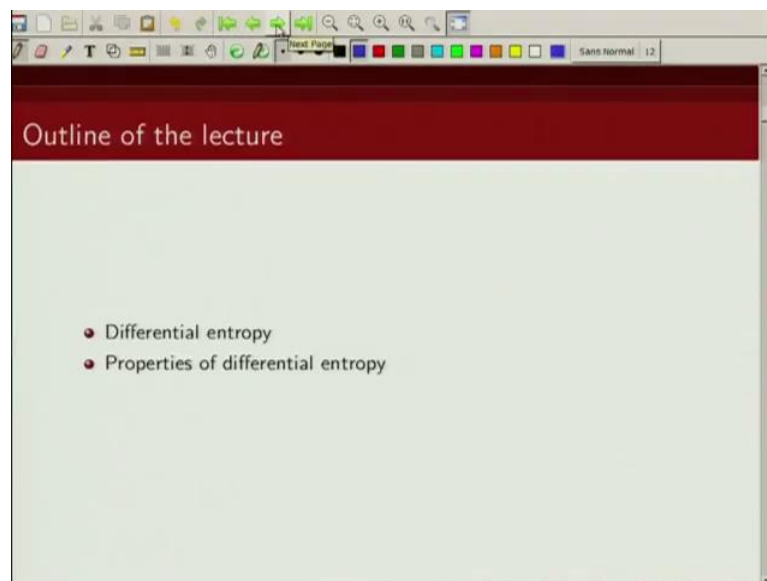**An Introduction to Information Theory**
**Prof. Adrish Banerjee**
**Department of Electronics and Communication Engineering**
**Indian Institute of Technology, Kanpur**

**Lecture – 11**
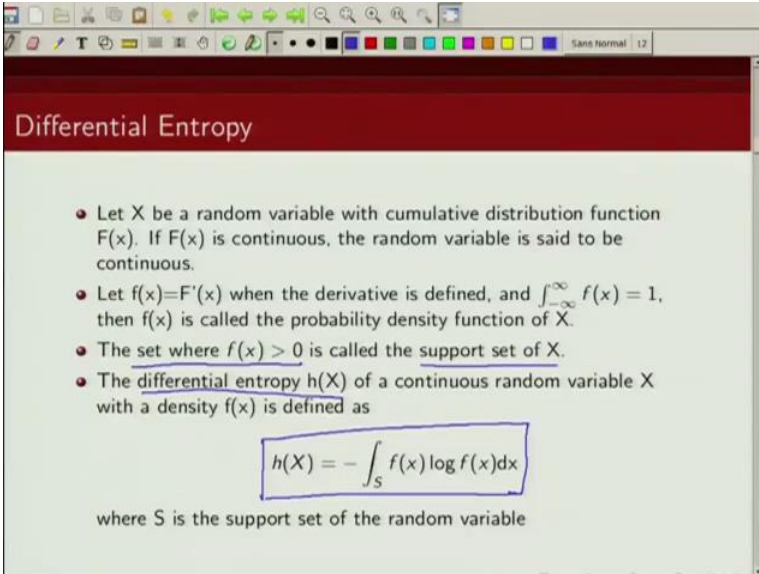**Differential Entropy**

Welcome to the course on An Introduction to Information Theory. So far we have been talking about discrete random variable and entropy associated with discrete random variable. Now, today we are going to talk about entropy for continuous random variable and that is known as differential entropy. So, it's the entropy of continuous random variable.

(Refer Slide Time: 00:41)



In this class, we will define differential entropy and we will prove some properties of differential entropy.

Let x be a random variable whose cumulative distribution function is given by F of x. If F of x is continuous, then we say that random variable is continuous. So, if derivative of F of x exists, n is defined which has the following property that integration from minus infinity to infinity is 1, then F of x is known as the probability density function of x. The set over which F of x is greater than zero is known as support set of x. These definitions are very similar to the definitions we had for discrete case and if you recall for the discrete random variable, we define the support set as the set where probability p of x is greater than 0. Now, differential entropy of a continuous random variable is defined as minus of integration of this probability density function log of probability density function, integration over the support set of this random variable.

So, this is how we define differential entropy and if you recall for discrete random variable, we define entropy as minus of summation p of x log of p of x. This is very similar to that.

(Refer Slide Time: 02:34)



Now, let us take an example. Let X is a normal distributed random variable. So, the density function phi of x is given by pi sigma square exponential minus x square by 2 sigma square. So, this is a zero mean variance sigma square normal distributed random variable whose density function can use the term f of x or I am using here phi of x is given by this expression. So, differential entropy is defined in this particular fashion. This can be written as phi of x, now if you take log of phi of x, we get log of this term and log exponential will be just this term. We get minus phi of x and log of phi is given by this expression. So, this and log of 1 by under root of 2 pi sigma square and now this can be written as minus phi of x minus x square 2 sigma square d of x and minus of minus actual log of 2 pi sigma square phi of x d x.

So, this minus minus becomes plus, this integration phi x x square 2 sigma square d x, it is nothing but expected value of x square by 2 sigma square. This 1st term that you see here is nothing but expected value of x square divided by 2 sigma square and there is no x here, so, this is a constant and if you integrate over phi of x we will get 1. So, this term is nothing but given by this. Now, the expected value of x square is given by variance of x plus expected value of x whole square and the expected value of x is 0. This term is 0 and variance is given by sigma square. This is variance sigma square and expected value of x square is given by sigma square.

So, sigma square sigma square cancels out, what we will get is half and I can also write it as half is half of natural log of e and this is of the form half of natural log of e plus half of natural log of 2 pi sigma square. This is log a plus log b and this can be written as log a times b and then this term be written like this. Now, note that we are talking about natural logs. So, this is in nats, if we take log to the base 2, then the units will be bits. The differential entropy of a 0 mean Gaussian random variable with variance sigma square is given by this expression half of log 2 pi e sigma square.

(Refer Slide Time: 06:28)



Now, if x 1 x 2 x 3 are sequence of random variables drawn independent identically distributed according to distribution f of x then, 1 minus n log of f of X 1 X 2 X 3 X n, this converges to differential entropy in probability. The proof of this is very similar to the proof that we did in the discrete case, this follows from the weak law of large number. So, this can be written as product of f x i and when you take log of product term, you get summation and then you will get summation of minus 1 by n log of f x i and that by weak law of large numbers is basically expected value of this which is equal to differential entropy and this converges to differential entropy in probability.

So, we are skipping this proof because this is very similar to what we have done earlier now for an epsilon greater than 0 for any n we can define typical set with respect to this

density function as follows. Typical set is defined as x 1 x 2 x 3 x n belonging to this set such that minus 1 by n log of f x 1 x 2 x 3, basically it is close to the true differential entropy and the difference is within epsilon.

(Refer Slide Time: 08:35)



So, very similar to how we defined typical set, we for the discrete case and for the continuous random variable also we can define typical set. Now, in case of discrete random variable, we define how many such typical sequence is there, the similar analogous thing for continuous random variable is what we call volume of a set. So, volume of a set is defined as analogous to the number of typical sequence here we have the volume of the typical set. Now, typical set has following properties, the probability of typical sequence for very large n that is greater than one minus epsilon, now again this follows from AEP property that we know that this converges to differential entropy in probability and this establishes the proof that probability of this typical set is greater than 1 minus epsilon.

(Refer Slide Time: 09:43)



The next property which says the f of x lies between this and this. Now, this can be proved from the definition of typical set. So, if you go back and look at the definition of typical set, this empirical differential entropy, the difference between that and the two differential entropy. The absolute difference should be less than or equal to epsilon.

(Refer Slide Time: 10:48)

So, depending on whether this term is greater or this term is greater because we are taking absolute difference and we will get either this inequality or we will get this inequality. Again the proof is similar to the proofs we have done many times before so, I am just skipping the details of the proof. The next property that we are going to show is the volume of the typical set is upper bounded by this.

Now, we know that if you integrate this density function over the support set, you will get probability 1 and if we integrate over typical set then since the typical set is subset of this S of n. So, we get here this greater than equal to because we are integrating over a smaller set and that is why this greater than equal to sign comes. We know that this density function is lower bounded by this. So, if we plug in this lower bound on f of x 1 x 2 x n here, what we get is this and that is why we are writing it as greater than equal to. We can take this out and if we integrate it over the typical set, what we get is a volume of the typical set and volume of the typical set multiplied by this is less than equal to 1.

So, from here we get volume of typical set is upper bounded by this 2 raise power n times differential entropy plus epsilon. Similarly, we can prove that volume of a typical set is lower bounded by this or just follows from we know probability of this typical set is greater than equal to 1 minus epsilon. So, 1 minus epsilon is less than equal to this probability of typical set which is given by this expression and now we know the upper bound on this f of x 1 x 2 x 3 x n, this follows from this property. So, we know the upper bound on f of x.

(Refer Slide Time: 13:24)



Now, if we plug in the upper bound of f of x 1 x 2 x n here, we get here less than equal to and subsequently you take this out and integrate it over the typical set what we get is volume of typical set multiplied by this and this is greater than equal to 1 minus epsilon. So, from here we get volume of typical set is at least given by this.

(Refer Slide Time: 14:02)

Now, what is the relation between differential entropy and discrete? So, we are comparing the differential entropy of a continuous random variable with the entropy of a discrete random variable that we get by discretizing that continuous random variable. This result says that entropy of n bit quantization of a continuous random variable is approximately given by differential entropy of X plus n. Let us consider a random variable X whose density function is given by f of x. Now, as we said we want to compare the entropy of the continuous random variable with the discretized version of this continuous random variable. So, we divide the range of X into bins of length delta and from mean value theorem, we know that there exist an x i such that f of x i times delta is equal to this.
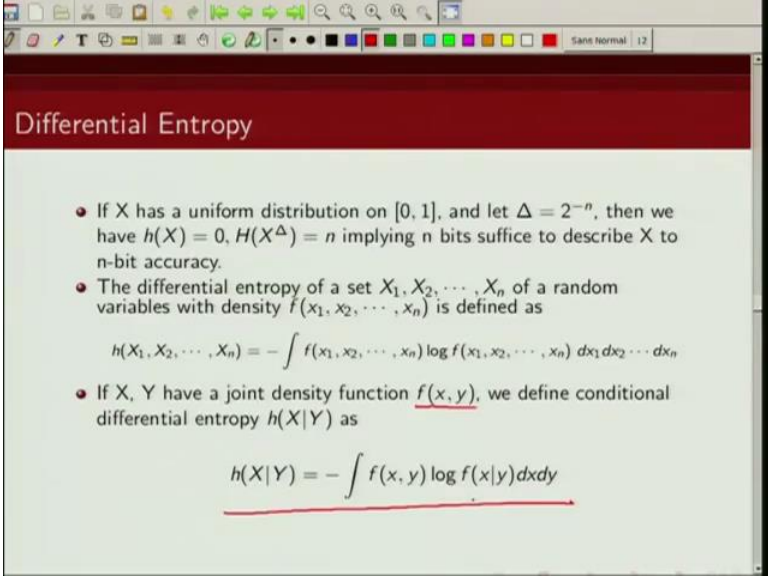
(Refer Slide Time: 15:37)



Then we are considering a quantized random variable and calling it x delta which is equal to X i if X lies between i delta and i plus 1 delta and the probability that X delta is equal to x is then given by this f of x d x integrate over this interval i delta to i plus 1 delta which by mean value theorem is given by f x i into delta. So, let us compute the entropy of this quantized random variable X delta. This is given by this expression which it follows from the definition of entropy, now p i is given by f x i delta. So, we plug in the value of p i here, we plug in the value of p i here and here we have log of a times b. This can be written as log of a plus log of b. So, then log of f x i delta can be written as

log of f x i and log of delta and this comes here, this comes here.

So, what we have here is this term plus this term, now this term does not depend on x i s and if you sum f x delta I, this will basically be 1 and we get minus log of delta and this term. Thus, the small delta approximates to the differential entropy.

(Refer Slide Time: 17:37)



We can write the entropy of the quantized random variable to be equal to differential entropy minus log of delta when delta is very small and if we assume X to have uniform distribution between 0 and 1 and we take our delta to be 2 raise power minus n, then in this case differential entropy is given by log of a which is a log of 1 which is 0 and log of delta minus log of delta will be n. So, H of delta will be h of X plus h of X minus log of delta minus log of delta is n, h of X is zero and then H of X delta will be n which implies that basically n bits are sufficient to describe X within n-bit accuracy.

Now, here we have made use of the fact that what is the differential entropy of a uniformly distributed random variable. So, we have a random variable which is uniformly distributed random variable. So, we have a random variable which is uniformly distributed between zero to a, its differential entropy from the definition h of X is given by minus our support group f of x log of f of x d of x, now if it is uniformly distributed between 0 to a, it looks like this. So, 0 to a is uniformly distributed, this will

be 1 by a, this is my f of x, this will be minus 1 by a log of 1 by a d x and this integration will be found 0 to a. So, this can be written as 1 by a integration log of a d x 0 to a, this is 1 by a into a into log of a.

So, this will be log of a, in this case a is 1, then log of 1 is 0 and that is why I wrote differential entropy for X which is uniformly distributed is between 0 to 1, its differential entropy is 0. One point I just wanted to make which is different for differential entropy compared to the entropy for discrete random variable is, for discrete random variable the entropy is greater than equal to zero, here the entropy can be less than zero. For example, if a is a fraction, let us say a is half then entropy in this case will be log of half which is minus log of 2 which is a negative quantity. So, differential entropy can be negative. However, volume of a typical set is positive.

Next, we are going to define differential entropy for random variable X 1 X 2 X 3 X n, so, differential entropy for a set x 1 x 2 x n whose joint density is given by this and can be similarly defined like this. So, it is a very straight forward extension of definition of differential entropy with single random variable. Similarly, we can also define conditional differential entropy. So, if X and Y have joint density function given by this, then we can define conditional differential entropy in this fashion. Again this is, if you look at the form of this expression, it is very similar to the form that we had for the discrete case.

(Refer Slide Time: 21:55)



Similarly, we can define mutual information. So, mutual information between two random variable with joint density function given by this is nothing but divergence between f x y and the marginal's f x and f y and this can be written as differential entropy of X minus conditional differential entropy of X given Y or differential entropy of Y minus differential entropy of Y given X.

(Refer Slide Time: 22:56)

Similarly, we could define the divergence between two densities f and g. This is defined as expected value f of log of f by g. This is the divergence between f and g. Now, let us consider an example, we are considering a multivariate normal distribution. So, X 1 X 2 X 3 X n have a multivariate normal distribution with mean given by mu and covariance matrix given by X show that the joint differential entropy is given by this expression where this is nothing but determinate of K, the covariance matrix. So, since X 1 X 2 X 3 X n follows a multivariate normal distribution such that density function is given by this expression.

(Refer Slide Time: 23:46)



Next, we apply the definition of joint differential entropy. So, this will be minus f of x log of f of x d x. When we take log of this function, we get log of this plus this term minus half x minus mu transpose K inverse x minus mu. So, this is what I am writing here and I have minus half x minus mu transpose K inverse x minus mu that is one term and then we had this minus of 1 by under root 2 pi n determinant of K. So, we have these two terms. Now, second term if I integrate over f of x, I will just get this and let us look at the first term which is nothing but expected value of this. So, in matrix form I can write this in this particular form, the summation over i and j, x i minus mu i K inverse i j x j minus mu j. So, after I can keep it here and I can combine these two terms, this can be written as expected value of this multiplied by K inverse of this.

(Refer Slide Time: 25:25).



Next, I can write this as K j i and K i j inverse. So, this will be non-zero only for those terms where i is equal to j and that is basically this and this is going to be i. So, summation over all j, this will give me n and this one I am getting n by 2 and then this term is coming from earlier, so we get this. Now, n by 2 I can write as n by 2 natural log of e and then I have this term which is half of natural log of 2 pi n and determinant of K. This I can also write as half of natural log of e raise power n plus this term. Now, I have half log of this particular term here. So, log of a plus log of b kind of form and this can be written as log of a b. So, this can be combined into this, now if I write unit in terms of log to the base 2, I get my differential entropy to be this. For a multivariate normal distribution with mean mu and co variance matrix K, the differential entropy is given by this.

(Refer Slide Time: 27:08)



Let us take an example. So, X and y are multivariate Gaussian distributed mean 0, co variance matrix given by K which is this. We can write h of X and h of Y which is nothing but half of log 2 pi sigma square and similarly you can write joint differential entropy which is given by half of log 2 pi square and determinant of k which is given by this and now we can write mutual information also as h of X plus h of Y minus joint entropy and this comes out to be this. Now, if rho 0 which is basically these terms of 0, you can see X and Y are independent and in that case the mutual information between X and Y is going to be 0 because X and Y are independent random variables and if rho is plus minus 1 that means they are perfectly correlated then you can see mutual information is going to be infinite.

(Refer Slide Time: 28:39)



Now, similar to the discrete random variable case, here also the divergence is greater than equal to 0. So, divergence between two density function f and g is greater than equal to 0, again we can prove it in the similar fashion, it is minus of divergence between f and g which is given by this expression, now log is a concave function if you recall log is a concave function, log is like this and from Jensen's inequality, the expected value of a function is less than equal to function evaluated at expected value and if it is a concave function then Jensen's inequality says this. So, this is like expected value of this log function. Since, log is a concave function then by Jensen's inequality, the log of expected value of the X should be more than expected value of the function.

From Jensen's inequality, we know this relation holds if f of x is concave, if this is a concave function then this expected value of log will be less than equal to log expected value and this basically is then integration of g over the support set is 1. So, log of 1 will be 0 and what we have proved is minus of divergence is less than equal to 0 and if multiplied by minus 1 both sides, we get the divergence between two densities f and g is greater than equal to 0. Now, similarly we could define chain rule also for continuous random variable. So, differential entropy between X 1 X 2 X 3 X n can be written using chain rule in this particular fashion and since we know conditioning cannot increase entropy, again this can be very easily proved. We have just now shown divergence is

greater than equal to 0 and we know mutual information is divergence between the joint densities and the marginal's.

(Refer Slide Time: 31:09)



So, mutual information should be greater than equal to 0 and from there we can prove that conditioning cannot increase entropy and from the chain rule which is given here if we apply this condition that conditioning cannot increase entropy, we get this condition and of course, this joint differential entropy is equal to this integral differential entropy when these X i s are independent. So, I may have given the proof, this follows from the chain rule and the fact that conditioning cannot increase entropy only when X and Y are independent that this is equal.

Now, we can make use of this result to prove what is known as Hadamards inequality. So, if we consider a multivariate random variable with 0 mean and covariance matrix given by K, now if we compute its differential entropy, this will be given by this expression half of log 2 pi e raise power n determinant of x and this h of X I, this half of log of 2 pi e K i i. So, if we compute this differential entropy here, we will get terms of the form, this is half log 2 pi e, this K i i form summation over all n, we will get something of this form whereas, this differential entropy is the form half log of 2 pi e raise to power n determinant of this..

So, summation of log, this can be written as half log of product from i equal to 1 to n K i i and you have this 2 pi e raise to power n term here and if you compare these two forms you have half log of 2 pi e n and you have half log of two pi e n and since this is less than equal to this, we will get the condition that determinant of K is less than equal to this particular term. So, this follows from two results, one is this one and second is the differential entropy of this multivariate Gaussian random variable.

(Refer Slide Time: 35:29)



## Properties of differential entropy

- $h(X + c) = h(X)$
- Proof: Let $Y = X + c$. Then $f_Y(y) = f_X(y - c)$ and $S_Y = \{x - c : x \in S_X\}$. Letting $x = y - c$, we have

$$
\begin{aligned}
h(X) &= -\int f_X(x) \log f_X(x)\, dx \\
&= -\int f_X(y - c) \log f_X(y - c)\, dy \\
&= -\int f_Y(y) \log f_Y(y)\, dy \\
&= h(Y) = h(X + c)
\end{aligned}
$$

We will prove some more properties of differential entropy. So, translation does not change differential entropy. So, differential entropy of X plus c is same as h of X. This is straight forward to prove and we have a new random variable Y which is X plus c and we can write the density function of y in terms of density function of x.

Now, we can write density differential entropy of X in this particular fashion and X is nothing but y minus c. So, this is given by this and this is nothing of f of y and this is equal to differential entropy of Y which is nothing but differential entropy of y is X plus c. and this is differential entropy of X plus c. So, what we have proved is differential entropy of X is same as differential entropy of X plus c. So, translation does not change differential entropy.

(Refer Slide Time: 36:47)



This is the effect of scaling. So, differential entropy of a of a times X is given by differential entropy of X plus log of absolute value of a. Here, y is given by a of X, we can write the density function of y in terms of density function of x and from the definition, differential entropy of y is given by this expression, now we plug in the densities of y in terms of densities of x which is this and this we simplify log of 1 by absolute value of a plus log of f of x of this and y by a is nothing but my x and this can be simplified into two terms, one is this particular term and second term is log of a.

So, we can see basically this will be log of 1 by a, log of 1 by a is minus log of a and minus minus that becomes plus and when you integrate it over the density function, this density function will become one. So, you got log of a integration of this density function will give you ,once it is log of a and the next term that you will get is this log of f of x y a which is nothing but f of x log of f of x. This is nothing but differential entropy of X, this is log of absolute value of a. So, the effect of scaling is as follows, h of differential entropy of a of X is differential entropy of X plus log of absolute value of a.

Now, the next result that we are going to show is as follows. So, if you have a random vector X with zero mean and variance given by this, then its differential entropy is upper bounded by half log of 2 pi e raise power n determinant of this covariance matrix and equality happens when X is multivariate normal distributed random variable.

So, if X is the random variable with zero mean and covariance matrix K then, we are saying it is the multivariate Gaussian distributed random variable which will have the maximum differential entropy and if you recall its counterpart in case of discrete random variable, it is the uniform distributed random variable which will have the maximum entropy. Now, how do we prove it? We have been given that mean is 0 and covariance is fixed K. So, let g of X is a density satisfying this condition, this is basically the covariance. Let phi of x is zero mean multivariate distribution which has the same second order moment which is given by K. Now, let us compute the divergence between the density g and this zero mean multivariate normal distributed phi K. So, this divergence between g and phi k can be written from the definition as g log of g by phi K.

This can be written as minus differential entropy of g minus integration of g log phi K, now log phi K is of quadratic form, log phi K is of quadratic form and we have been given that x and phi K has the same covariance that is they have the same. So, integration

of x i x j phi K d x is same as integration of x i x j g x d x and then this term g of log phi K will be same as phi K log phi K because this has a quadratic form log of phi x will be a quadratic form x i x j form and we know that integration of g x x i x j, this is same as integration of x i x j phi K. So, this is equal to this and that is why you are able to write it like this.

Now, this is nothing but plus of differential entropy of this zero mean multivariate distribution and this is greater than equal to 0 that means, h of phi k is greater than equal to h of g. So, in other words it is the multivariate normal distribution random variable which will have the maximum differential entropy. So, with this we will conclude our discussion on differential entropy. In the next class, we will talk about Gaussian channel.

Thank you