

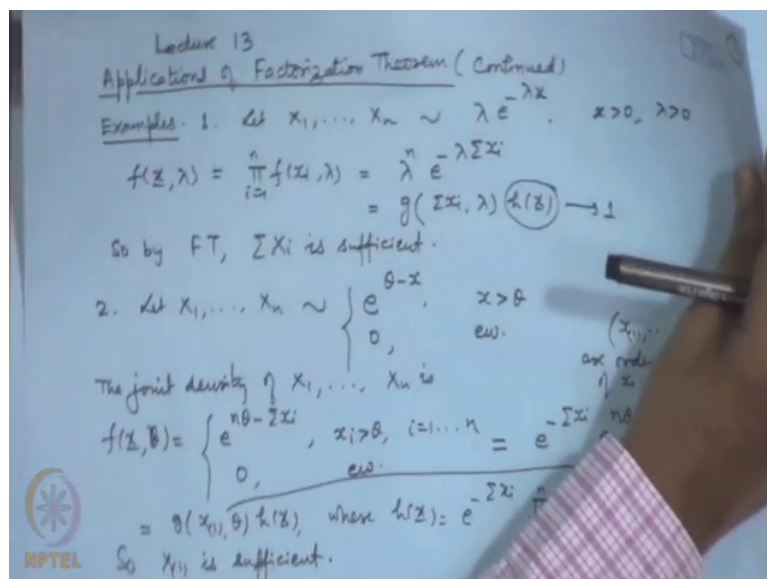
Statistical Inference
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture No. # 13
Sufficiency and Information

In the previous class, I have explained the concept of sufficiency. This concept is the concept which is called the principle of data relation. So, we have a random sample X_1, X_2, \dots, X_n , but if we have a sufficient statistic T , then that is sufficient that gives the complete information about the parameter which is contained in the sample. So, we need not retain it we have given one theorem, which is called factorization theorem and this is useful for deriving sufficient statistics in various probability models.

Yesterday, I have discussed the normal probability model and I have shown you that how, if we change the parameter space; that means, whether we have μ known or σ^2 known or both are unknown in each of the cases, the sufficient statistics changes. So, sufficiency is the property of the probability model under consideration. Let me explain it through a few more examples and we will use the concept of this factorization theorem here.

(Refer Slide Time: 01:29)



Let me start with exponential distribution, let X_1, X_2, \dots, X_n follow exponential distribution say with parameter λ . So, in the factorization theorem we need to write down the joint density of X_1, X_2, \dots, X_n that is equal to $\lambda^n e^{-\lambda \sum_{i=1}^n x_i}$. Now, this whole thing we can write as a function of $\sum_{i=1}^n x_i$ and λ and $h(x)$ this $h(x)$, I am taking to be one itself the constant.

So, you can see by factorization theorem by factorization theorem $\sum_{i=1}^n x_i$ is sufficient, let us consider another exponential model in which in place of a scale parameter we will have a location parameter. So, let us consider say X_1, X_2, \dots, X_n following exponential say, $\theta - x$ where x is greater than θ_0 elsewhere. Now, in this case the joint density of X_1, X_2, \dots, X_n is $f(x|\theta)$ that is equal to $e^{-n(\theta - \sum_{i=1}^n x_i)}$; however, this description of x_i greater than θ also place a role here.

Now, if we want to write it as a product here we will make use of the indicator function. So, we can write it like this $e^{-\sum_{i=1}^n x_i}$, $e^{-n\theta}$ indicator function of the set X_1 from θ to infinity and indicator function of other x_i 's from 2 to n from X_1 to infinity. So, what we can consider we can write it as $g(\sum_{i=1}^n x_i, \theta)$ into $h(x)$ where, $h(x)$ I am writing as $e^{-\sum_{i=1}^n x_i}$ into product i is equal to 2, 2^{-n} I of x , X_1 to infinity.

So, here this X_1, X_2, \dots, X_n they are denoting the order statistics of X_1, X_2, \dots, X_n . So, g is this function this is a function of X_1 and θ . So, we conclude that X_1 is sufficient. So, X_1 is sufficient. Now, note here when we had λ as the parameter and here we had a scale model the sufficient statistics was $\sum_{i=1}^n X_i$ although here, we are again we are dealing with the exponential distribution, but the nature of the parameter has changed therefore, the sufficient statistic is now, the minimum of the observations now, in a similar way let us take up the two parameter exponential distribution.

(Refer Slide Time: 05:18)

3. Let X_1, \dots, X_n be a random sample from a two parameter exponential distributions with pdf $f(x, \mu, \sigma) = \begin{cases} \frac{1}{\sigma} e^{-\frac{x-\mu}{\sigma}}, & x > \mu \\ 0, & \text{ew} \end{cases}$

The joint pdf of X_1, \dots, X_n is

$$f(x, \mu, \sigma) = \frac{1}{\sigma^n} e^{-\frac{\sum x_i}{\sigma}} e^{-\frac{\sum x_i}{\sigma}} \prod_{i=1}^n I_{(\mu, \infty)}(x_i)$$

$$= g(x_{(1)}, \sum x_i, \mu, \sigma) h(x)$$

So $(x_{(1)}, \sum x_i)$ is sufficient
or $(x_{(1)}, \bar{X})$ is sufficient.

Let us take X_1, X_2, \dots, X_n be a random sample from a two parameter exponential distribution say with density function $f(x, \mu, \sigma) = \frac{1}{\sigma} e^{-\frac{x-\mu}{\sigma}}$ for $x > \mu$ and it is equal to 0 otherwise. So, once again the joint probability density function of X_1, X_2, \dots, X_n this is now, $\frac{1}{\sigma^n} e^{-\frac{\sum x_i}{\sigma}}$ for $x_i > \mu$ and 0 otherwise.

And once again the condition that each of $x_i > \mu$ I can express in terms of the indicator function like x_1 is from μ to infinity and x_i is other x_i is they are from x_1 to infinity $i = 2$ to n . So, this portion I can write as $g(x_{(1)}, \sum x_i, \mu, \sigma)$ and this part is $h(X)$. So, here we concluded that $x_{(1)}$ and $\sum x_i$ is sufficient or we can also say $x_{(1)}$ and \bar{X} because this is a one to one function of this, this is sufficient.

I also want to mention here we have earlier considered, the maximum likelihood estimators. Now, let us remember one maximum likelihood estimators for each of these problems for example, in this case the maximum likelihood estimator for λ was $\frac{1}{\bar{X}}$ which is a function of $\sum x_i$ in this particular case. The maximum likelihood estimator was $x_{(1)}$ that is a minimum of observations and it is sufficient here similarly, here you see the maximum likelihood estimator for μ and σ , where $x_{(1)}$ and $\bar{X} - x_{(1)}$ respectively, which is again a one-to-one function of $x_{(1)}, \bar{X}$ that is the sufficient statistics.

So, we can observe that maximum likelihood estimator if it exists is actually, a function of the sufficient statistics the reason is obvious because in the factorization theorem, we are writing down the density as a function of the parameter and the sufficient statistics into a function, which is free from the parameter. Now, in the method of maximum likelihood estimator, we are maximizing the density function or the mass function with respect to the parameter.

Now, the part of the density which contains the parameter contains the variable only through the sufficient statistics therefore, the maximization problem will give a solution in terms of the sufficient statistics alone.

(Refer Slide Time: 09:15)

If maximum likelihood estimators exist, they are functions of sufficient statistics.

Examples: X_1, \dots, X_n a random sample from double exponential distⁿ. $\frac{1}{2} e^{-|x-\theta|}$, $x \in \mathbb{R}, \theta \in \mathbb{R}$

The joint pdf of X_1, \dots, X_n is

$$f(\underline{x}, \theta) = \frac{1}{2^n} e^{-\sum_{i=1}^n |x_i - \theta|} = \underbrace{\frac{1}{2^n}}_{h(\underline{x})} e^{-\sum_{i=1}^n |x_i - \theta|} = \underbrace{e^{-\sum_{i=1}^n |x_i - \theta|}}_{g(\underline{x}_1, \dots, \underline{x}_n, \theta)}$$

So (X_1, \dots, X_n) is sufficient

Order statistics
 $\hat{\theta}_{ML} = \text{Median}$ (\rightarrow a function of order statistics)

NPTEL

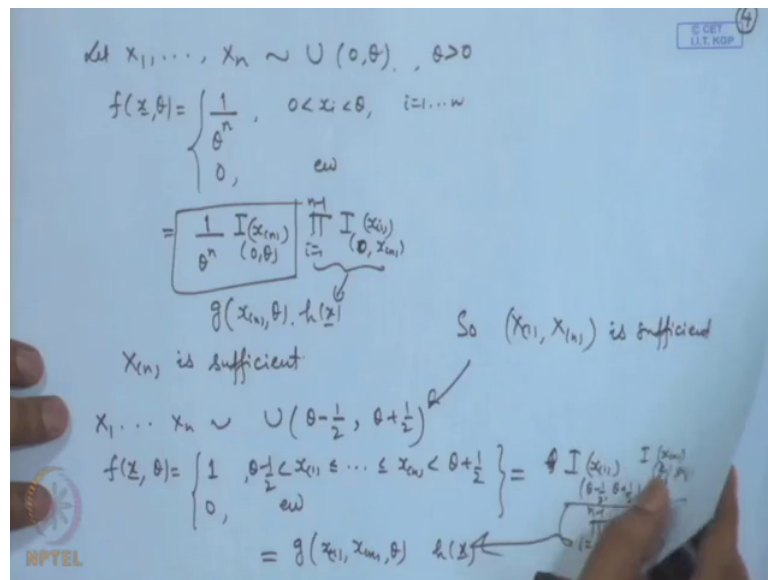
So, we have a general comment here that if maximum likelihood estimators exist, they are functions of sufficient statistics. Let us take some more examples here say for example, X_1, X_2, \dots, X_n a random sample from say double exponential distribution half e to the power minus x minus θ , where x is any real number and θ is any real number. In this case, if we consider the sufficiency. So, the joint distribution of X_1, X_2, \dots, X_n that is equal to one by two to power n e to the power minus sigma modulus x_i minus θ , i is equal to one to one. Now, here you observe I cannot reduce it further as a function of parameter and a another variable here because each of the x_i is appearing in the modulus sign and therefore, I cannot separate it out.

At the most I can consider the reduction as one by two to the power n to power minus sigma modulus of x_i order statistics minus theta. So, this function is now a function of the order statistics and theta and this you can. Now, call $h(x)$. So, we conclude that the order statistics X_1, X_2, \dots, X_n is sufficient this order statistics. Now, remember here for this problem what was the maximum likelihood estimator. The maximum likelihood estimator was median. Median of the observations and median is a function of this is a function of order statistics because if we have a odd number of observations say x_{2m+1} then x_{m+1} that is the middle of the observation was the median and if we have an even number of observations that is x_{2m} then any number between x_m and x_{m+1} and we can actually, consider say the middle of the two that is $x_m + x_{m+1} / 2$ as the maximum likelihood estimator.

So, this is a function of order of statistics in this case also. So, this statement is true in general another thing which I just now, pointed out that many times. When we are writing down the density function say in this case, we are we have to incorporate the region of the variable, which is dependent upon the parameter as a part of the joint density function because if we do not include it then we cannot derive the sufficient statistics. For example, if we have written only this part then there is no sufficient statistics here because $e^{-\theta x_i}$ can be separately written $e^{-\theta x_i}$ can be separately written.

However, this is not a complete description of the density unless we include the reason x_i greater than theta for all i and this is the way of including this a similar phenomena is observed in the uniform distributions also like in the uniform distribution.

(Refer Slide Time: 13:16)

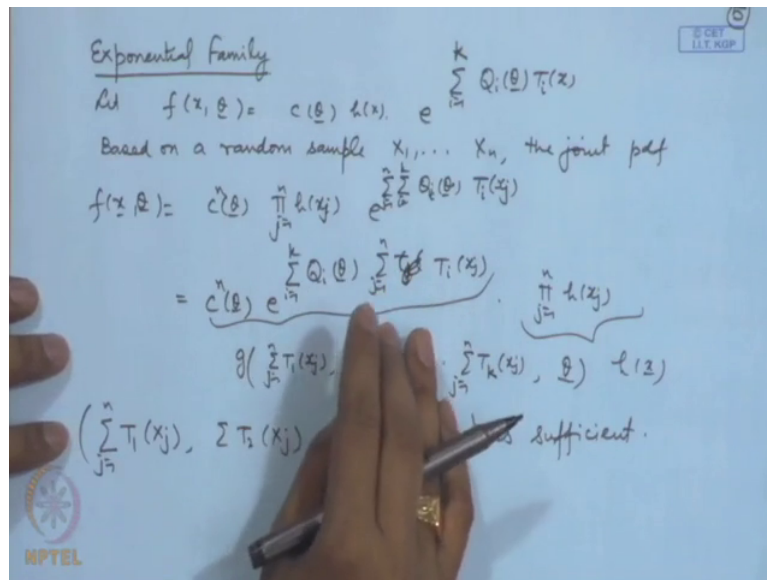


The range is dependent upon the range of variable is dependent on the parameter. So, let us consider say X_1, X_2, \dots, X_n say a random sample from uniform zero theta distribution. Now, in this case the joint density is equal to $1/\theta^n$ if each of these x_i is between 0 to θ and it is equal to 0 elsewhere. So, this part we will express as then one by $1/\theta^n$ and the other part you can write as $h(x)$. So, X_n is sufficient; however, if we consider say a uniform distribution which is on a two sided interval here we have taken one side as 0.

Suppose, we consider say from $\theta - 1/2$ to $\theta + 1/2$ in this case the joint distribution is simply 1 because $\theta + 1/2 - (\theta - 1/2) = 1$; however, each of the x_i 's they are between $\theta - 1/2$ and $\theta + 1/2$. So, this part then you can incorporate as indicator function of X_1 from $\theta - 1/2$ to $\theta + 1/2$ and the indicator function of X_n from $\theta - 1/2$ to $\theta + 1/2$. And the remaining order statistics lying between X_1 and X_n is equal to 1 . So, this you can see it is a function of X_1, X_n, θ into $h(x)$ this part is $h(x)$ and this part is a function X_1, X_n and θ . So, here X_1, X_n is sufficient although the parameter remains one dimensional here, but the order statistics contains two terms if you remember the maximum likelihood estimator.

The maximum likelihood estimator for this problem was any value between $X_{n-1/2}$ to $X_{1/2}$. So, which is a function of X_1, X_n . So, the statement that the maximum likelihood estimators, if they exist they are functions of the sufficient statistics is satisfied here also now this factorization theorem is very useful.

(Refer Slide Time: 17:18)



If we are considering a general distribution in an exponential family so, let us consider distributions in exponential family. So, if we are considering a k dimensional exponential family let $f(x, \theta)$ be equal to $c(\theta) h(x) e^{\sum_{i=1}^k Q_i(\theta) T_i(x)}$ is equal to say 1 to k this is called a k dimensional exponential family, provide the parameter space contains a k dimensional rectangle. So, based on a random sample X_1, X_2, \dots, X_n the joint probability density function we can write as $c(\theta)^n$ product of $h(x_j)$ e to the power $\sum_{i=1}^k Q_i(\theta) \sum_{j=1}^n T_i(x_j)$. So, let me change here I_j because I_j is being used here.

I is equal to 1 to k , x_j $\sum_{j=1}^n$ is equal to 1 to n . Now, this I can write as $c(\theta)^n$ product of $h(x_j)$ e to the power $\sum_{i=1}^k Q_i(\theta) \sum_{j=1}^n T_i(x_j)$ there is a mistake it is $T_i(x_j)$. So, $T_i(x_j)$. So, this becomes what I am doing is I am taking this summation inside. So, this becomes $\sum_{i=1}^k Q_i(\theta) \sum_{j=1}^n T_i(x_j)$ is equal to 1 to k into product of $h(x_j)$, j is equal to 1 to n . So, this part is now a function of $\sum_{j=1}^n T_1(x_j), \dots, \sum_{j=1}^n T_k(x_j)$ 1 to n and.

So, on sigma T k, x j, j is equal to 1 to n and the parameter and this part we can consider as h x. So, we conclude that sigma of T 1, x j sigma of T 2, x j and so on. Sigma of T k, x j, j is equal to 1 to n this is sufficient of course, when we write like this we assume that this q 1, q 2 etcetera are linearly, independent otherwise some of the terms can be merged together. Now, let me introduce the relationship between the Fisher's information measure and the concept of sufficiency.

(Refer Slide Time: 20:55)

Fisher's Information was defined as $X \sim f(x, \theta)$ (pdf, pmf) (under regularity conditions)

$$I_X(\theta) = E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \right]^2$$

If T is any statistic with density $\phi(t, \theta)$, then Fisher's Information in T is defined as $I_T(\theta) = E \left[\frac{\partial}{\partial \theta} \log \phi(T, \theta) \right]^2$ (under regularity conditions)

Relationship Between Sufficiency & Information

Let $T(X)$ have pdf (pmf) $\phi(t, \theta)$, $\frac{\partial}{\partial \theta} \phi(t, \theta)$ exists, (for any measurable set B in the space of T -values)

$$\frac{d}{d\theta} \int_B \phi(t, \theta) d\mu(t) = \int_B \frac{\partial \phi(t, \theta)}{\partial \theta} d\mu(t)$$

Then (i) $E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \mid T=t \right] = \frac{\partial \phi(t, \theta) / \phi(t, \theta)}{\partial \theta}$ a.e.

(ii) $I_X(\theta) \geq I_T(\theta)$ with equality holding if and only if T is sufficient.

Fisher's Information was defined as $X \sim f(x, \theta)$ (pdf, pmf) (under regularity conditions)

$$I_X(\theta) = E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \right]^2$$

If T is any statistic with density $\phi(t, \theta)$, then Fisher's Information in T is defined as $I_T(\theta) = E \left[\frac{\partial}{\partial \theta} \log \phi(T, \theta) \right]^2$ (under regularity conditions)

Relationship Between Sufficiency & Information

Theorem

Let $T(X)$ have pdf (pmf) $\phi(t, \theta)$, $\frac{\partial}{\partial \theta} \phi(t, \theta)$ exists, (for any measurable set B in the space of T -values)

$$\frac{d}{d\theta} \int_B \phi(t, \theta) d\mu(t) = \int_B \frac{\partial \phi(t, \theta)}{\partial \theta} d\mu(t)$$

$E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \mid T=t \right] = \frac{\partial \phi(t, \theta) / \phi(t, \theta)}{\partial \theta}$ a.e.

(ii) $I_X(\theta) \geq I_T(\theta)$ with equality holding if and only if T is sufficient.

So, if you remember the Fisher's information was defined as Fisher's information was defined as $I_X(\theta) = E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \right]^2$, here the

assumption is that the distribution of x is $f(x; \theta)$ and of course, this could be p.d.f or p.m.f with respect to a measure μ and we are making the assumption of regularity conditions that is differentiation under the integral sign is allowed.

So, this is under regularity conditions if the distribution of X is $f(x; \theta)$ then the information measure Fisher's information in X about θ is defined as expectation of $\frac{d}{d\theta} \log f(x; \theta)$ squared. Now, if t is any statistic and suppose the density of let me give the name as $\phi(t; \theta)$, then Fisher's information in T is defined as $I_T(\theta)$ is equal to expectation of $\frac{d}{d\theta} \log \phi(t; \theta)$ squared, once again we are making assumption about.

So, this could be p.d.f also or p.m.f and we should have the regularity conditions satisfied for ϕ also means, we should be able to differentiate the density with respect to that parameter, we should be able to differentiate under the integral sign. So, here also under regularity conditions we satisfy it under regularity conditions. So, we have the following result regarding relationship between sufficiency and information. Let $T(x)$ have p.d.f or p.m.f $\phi(t; \theta)$ $\frac{d}{d\theta} \phi(t; \theta)$ exists $\frac{d}{d\theta} \int \phi(t; \theta) d\mu_T$ over an image set b is equal to $\frac{d}{d\theta} \int \phi(t; \theta) d\mu_T$ for any measurable set b that is in the space of T values.

Then we have the following results first is that expectation of $\frac{d}{d\theta} \log f(x; \theta)$ given T is equal to t this conditional expectation is equal to $\frac{d}{d\theta} \log \phi(t; \theta)$ almost everywhere. Secondly, the information in the x is always greater than or equal to the information in any statistic T with equality holding if and only if T is sufficient. So, what we are saying suppose we have X_1, X_2, \dots, X_n as a sample and T is any statistic. Then in general the information content in a statistic will be less than or equal to the information contained in the full sample; however, if T is sufficient then it will be the same and this is a necessary and sufficient condition.

So, this is what I was mentioning from the, that is the utilization of the information or the content of the information in the concept of sufficiency that sufficient statistic contains all the information, which is available in the sample because we are saying $I_T(\theta)$ will become equal to $I_h(\theta)$. So, this is the physical meaning of the concept of sufficiency that if, we are considering this definition as the definition of information because this $I(x; \theta)$ we will call information in the sample. So, what we are saying is that there is no loss of information, if we consider a sufficient statistics. So, let me prove this theorem here.

(Refer Slide Time: 26:57)

Proof: Let X be the space of X -values & Y denote the space of T -values, $T: X \rightarrow Y$.
 Let \mathcal{B} be the σ -field of subsets of X & \mathcal{C} be the σ -field of subsets of Y .
 Let $B \in \mathcal{B}$, $C \in \mathcal{C}$, $B = T^{-1}(C) = \{x: T(x) \in C\}$.
 For any set $B \in \mathcal{B}$,

$$E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \mathbb{I}_B(X) \right] = \int_B \frac{\frac{\partial f(x, \theta)}{\partial \theta}}{f(x, \theta)} \cdot \frac{1}{f(x, \theta)} f(x, \theta) d\mu(x)$$

$$= \int_B \frac{\partial f(x, \theta)}{\partial \theta} d\mu(x) = \frac{d}{d\theta} \int_B f(x, \theta) d\mu(x) = \frac{d}{d\theta} P(X \in B)$$

$$= \frac{d}{d\theta} P(T \in C) = \frac{d}{d\theta} \int_C \phi(t, \theta) d\mu(t) = \int_C \frac{\partial \phi(t, \theta)}{\partial \theta} d\mu(t)$$

So, see we have say let x be the space of x values and say y denote the space of T values; that means, T is a function from x to say y naturally, we will be considering the sigma field's of subsets of x and similarly a sigma field of subsets of y also. So, let us use some notation say \mathcal{B} , be the sigma field of subsets of x and say \mathcal{C} be the sigma field of subsets of y which, we are considering here. So, now, let us consider B as a set C in \mathcal{C} then for that define say B is equal to $T^{-1}(C)$ that is the set of x such that $T(x)$ belongs to C . So, consider for any set B belonging to \mathcal{B} , let us consider expectation of $\frac{\partial}{\partial \theta} \log f(x, \theta)$ over the set B .

So, this is equal to $\frac{\partial f}{\partial \theta}$ divided by $f(x, \theta)$. Now, this is expectation. So, it becomes $\int_B \frac{\partial f}{\partial \theta} d\mu(x)$ over the set B . So, this $f(x, \theta)$ and this $f(x, \theta)$ cancels out, we are getting $\int_B \frac{\partial f}{\partial \theta} d\mu(x)$ over B . Now, this we can consider because we have made the assumption that, we can differentiate under the integral sign. So, this equal to $\frac{d}{d\theta} \int_B f(x, \theta) d\mu(x)$. Now, this the integral of the density of the random variable x over the set B . So, this is nothing, but probability of the set B .

Now, we have defined the set B to be the inverse function of or inverse image of T . So, B is the set where $T(x)$ belongs to C . So, this probability of x belonging to B is same as probability of T belonging to C therefore, we can write it as $\frac{d}{d\theta} \int_C \phi(t, \theta) d\mu(t)$ that is the density of a T with respect to the corresponding measure over the set C . Now, once again we have made the assumption that we can consider differentiation under the integral sign. So, this

becomes del phi by del theta d mu T over the set c now, we can divide and multiply by the density of t inside the integral sign.

(Refer Slide Time: 31:21)

$$= \int_C \frac{\partial \phi(t, \theta)}{\partial \theta} \cdot \frac{1}{\phi(t, \theta)} \phi(t, \theta) d\mu(t)$$

$$= \int_C \frac{\partial \log \phi(t, \theta)}{\partial \theta} \phi(t, \theta) d\mu(t) = E \left[\frac{\partial \log \phi(T, \theta)}{\partial \theta} I_C(T) \right]$$

By the definition of conditional expectation, we conclude that

$$E \left[\frac{\partial \log f(x, \theta)}{\partial \theta} \mid T=t \right] = \frac{\partial \phi(t, \theta)}{\partial \theta} / \frac{\phi(t, \theta)}{\phi(t, \theta)} \text{ a.e.}$$

(Remark: A function $g(t)$ is said to be $E(Y/T=t)$ if $E(Y I_B(T)) = E(g(T) I_B(T)) \rightarrow B \in \mathcal{G}$.)

So, we will get here this term as equal to del phi by del theta 1 by phi t theta phi T theta d mu T over set c. So, now, this becomes nothing, but the derivative of log of phi t theta d mu T over the set c this is nothing, but the expectation of this expectation of del log phi T theta by del theta indicator function of the set c look at the statement that, we have proved. Now, we started with expectation of del by del theta log of f x theta I b x.

We are showing that this term is now, equal to this term is equal to expectation of del log phi t theta by del theta I c T. Now, what is the relationship between x and T and b and c t is a function of x and b is the inverse image of the set c therefore, by the definition of the conditional expectation, we conclude that by the definition of conditional expectation we conclude that expectation of del by del theta log of f x theta given T is equal to T which is equal to del phi by del theta one by divided by phi of T theta, that is the statement given here of course, since we are obtaining this result from the expectation. So, we can say that this statement is true almost everywhere.

That means, the set we have this may not be true will have probability 0; that means, their set of values of small t for which this statement is not true then under the probability distribution of t that set will have probability 0. So, actually what we have used here is we have simply

used the definition of the conditional expectation. In fact, let me write here remark a function $g(t)$ is said to be conditional expectation of y given T if expectation of y I, b , T is equal to expectation of $g(t)$, I, d , t for all b more measurable sets of b . So, we have used this definition. So, what we have done is we have established a relationship in the, in the log likelihood or you can say the information content term in the density of the sufficient statistics and the original variable.

(Refer Slide Time: 35:17)

Consider $E \left[\frac{\partial \log \phi(T, \theta)}{\partial \theta} - \frac{\partial \log f(x, \theta)}{\partial \theta} \right]^2 \geq 0 \dots (1)$

The LHS is $= E \left(\frac{\partial \log \phi(T, \theta)}{\partial \theta} \right)^2 + E \left(\frac{\partial \log f(x, \theta)}{\partial \theta} \right)^2 - 2 E \left[\frac{\partial \log \phi(T, \theta)}{\partial \theta} \cdot \frac{\partial \log f(x, \theta)}{\partial \theta} \right] \dots (2)$

$E \left[\frac{\partial \log \phi(T, \theta)}{\partial \theta} \cdot \frac{\partial \log f(x, \theta)}{\partial \theta} \right] = E E \left[\frac{\partial \log \phi(T, \theta)}{\partial \theta} \cdot \frac{\partial \log f(x, \theta)}{\partial \theta} \mid T \right]$
 $= E \left(\frac{\partial \log \phi(T, \theta)}{\partial \theta} \right)^2$

So LHS of (1) is $I_T(\theta) + I_X(\theta) - 2 I_T(\theta)$
 $= I_X(\theta) - I_T(\theta)$

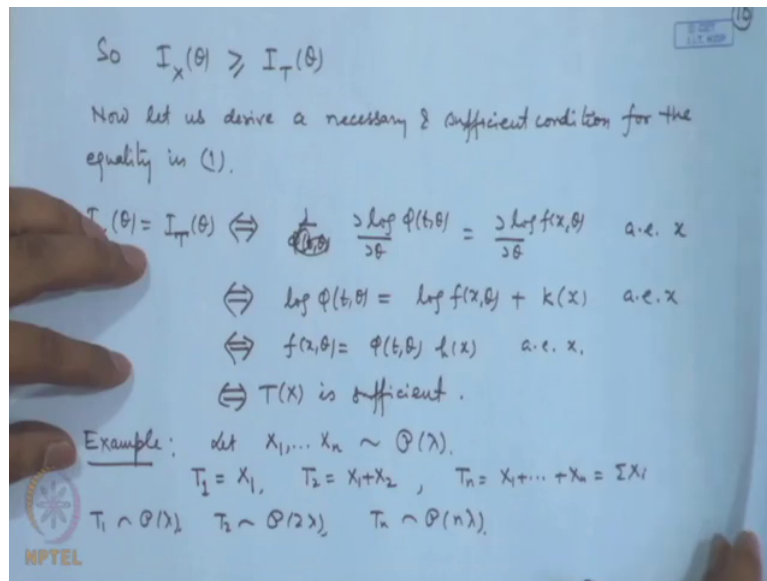
Let us look at the proof of the second part. So, consider here expectation of $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ this will be capital here because you are considering expectation. So, this term is going to be greater than or equal to 0 because this is a perfect square term here. Now, let us expand the left hand side the left hand side is equal to now you expand this. So, this is becoming expectation of $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ square plus expectation of $\frac{\partial \log f}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ square minus twice expectation and the product of these terms, that is $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$.

At this stage you notice here that expectation conditional expectation of $\frac{\partial \log f}{\partial \theta}$ given T is this term that is $\frac{\partial \log \phi}{\partial \theta}$ this term is nothing, but $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$. If I consider this expectation here I can write here it as expectation of expectation given T then this term becomes expectation of $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ into $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ you can express as expectation of expectation $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ del $\frac{\partial \log \phi}{\partial \theta}$ by $\frac{\partial \log f}{\partial \theta}$ given T . Now, if we use the relationship which we proved in the first

part that is this one then this conditional expectation becomes this term itself. So, this will become square of.

So, left hand side of one is then information in X, the first term is information in T plus information in X minus twice information in T that is equal to information in X minus information in T and the right hand side is it is greater than or equal to 0.

(Refer Slide Time: 38:54)



So, we conclude that so, $I_X(\theta)$ is greater than or equal to $I_T(\theta)$. So, information in a statistic is always less than or equal to the information in the full sample now let us consider when we will have equality. Now, let us derive a necessary and sufficient condition for the equality in 1 now when will there be equality if we are considering $I_X(\theta)$ is equal to $I_T(\theta)$ $I_T(\theta)$ equal to 0. So, that equal to 0 will come if we have equal to 0 here now this is an expectation of a non negative quantity if we say that expectation is 0 then the quantity itself must be 0 with probability one . So, $I_X(\theta)$ is equal to $I_T(\theta)$ is equivalent to saying that $\frac{\partial \log \phi(t, \theta)}{\partial \theta}$ or you can say $\frac{\partial \log \phi(t, \theta)}{\partial \theta}$ by $\frac{\partial \log f(x, \theta)}{\partial \theta}$ almost everywhere.

That means, the set of values of x where this is not true will have probability 0. Now, you integrate on both the sides. So, you will get $\log \phi(t, \theta)$ is equal to $\log f(x, \theta)$ plus a function of say x because this integration is with respect to theta. So, this is equivalent to

saying that if I consider $f(x|\theta)$, then it is equal to $\phi(t|\theta)$ into a function of x . Now, this is nothing, but factorization theorem. So, we are saying that T is sufficient.

So, the information in the statistic T is equivalent to equal to the information $n \times x$ if and only if the random variable the statistic T is sufficient. So, this Fisher's information is measure is extremely important concept. In fact, in the current physics or in the information theory this is widely used one can look at the references physics of Fisher's information there is currently a book, which has come out and it almost establishes entire physics theory on the Fisher's information measure.

Let me give an example of a calculation of the information, we will show that this statement is true. So, let me take up say let us consider say X_1, X_2, \dots, X_n following say Poisson λ distribution, let us take several statistics let us take say T_1 is equal to X_1 , T_2 is say $X_1 + X_2$ and say T_n is equal to $X_1 + X_2 + \dots + X_n$ that is $\sum x_i$ in the case of Poisson distribution, then we can easily derive the distribution. So, T_1 follows Poisson λ , T_2 will follow Poisson 2λ and T_n will follow Poisson $n\lambda$.

Let us, independently derive the information in T_1, T_2 and T_n and also let us derive the information in the full sample. What is information in full sample? X_1, X_2, \dots, X_n . So, let us derive all these things.

(Refer Slide Time: 43:39)

So $E[f(x, \lambda)] = \frac{e^{-\lambda} \lambda^x}{x!}, x = 0, 1, 2, \dots$

$\log f(x, \lambda) = -\lambda + x \log \lambda - \log x!$

$\frac{\partial}{\partial \lambda} \log f = -1 + \frac{x}{\lambda} = \frac{x - \lambda}{\lambda}$

$E\left[\frac{\partial}{\partial \lambda} \log f(x, \lambda)\right]^2 = \frac{E(x - \lambda)^2}{\lambda^2} = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$

$I_{T_1}(\lambda) = \frac{1}{\lambda}, I_{T_2}(\lambda) = \frac{n}{\lambda}, I_{T_n}(\lambda) = \frac{n}{\lambda}$

So, information in one of the x that is calculated. If I calculate the information in X_1 (Refer Slide Time: 43:48) and if I take n times that information is easily we can see an additive function. So, the density function or the probability mass function in the Poisson distribution is. So, log of this is $f(x, \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$. So, $\frac{\partial}{\partial \lambda} \log f(x, \lambda) = -1 + \frac{x}{\lambda}$, which we can write as $\frac{x - \lambda}{\lambda}$. So, expectation of $\frac{\partial}{\partial \lambda} \log f(x, \lambda)$ is equal to expectation of $\frac{x - \lambda}{\lambda}$ that is equal to expectation of x minus λ by λ .

Now, this is nothing, but the variance of x because in Poisson distribution expectation of x equal to λ . So, this is equal to $\frac{\lambda}{\lambda^2}$ that is equal to $\frac{1}{\lambda}$. So, if I consider the information in T_1 then that is equal to $\frac{1}{\lambda}$. If we consider the information in say x itself, then it is additive. So, it will become n by λ , if I consider information in T_2 that will be equal to $\frac{2}{\lambda}$ and if I consider information in T_n that is also equal to $\frac{n}{\lambda}$.

(Refer Slide Time: 43:39)

$$E \left[\frac{\partial}{\partial \lambda} \log f(x, \lambda) \right]^2 = \frac{E \left[\frac{(x - \lambda)^2}{\lambda^2} \right]}{\lambda^2} = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$$

$$I_{T_1}(\lambda) = \frac{1}{\lambda}, \quad I_X(\lambda) = \frac{n}{\lambda}, \quad I_{T_2}(\lambda) = \frac{2}{\lambda}, \quad I_{T_n}(\lambda) = \frac{n}{\lambda}$$
 So we observe that $I_X(\lambda) = I_{T_n}(\lambda)$ as T_n is sufficient.

Remark: Information is additive. Let X and Y be independent r.v.'s with distributions $f_1(x, \theta)$ & $f_2(y, \theta)$.

So, you can see this is less than this, this one is less than this; however, this one is equal to this and T_n that is $\sum x_i$ in the case of Poisson distribution, we have shown that it is sufficient statistics. So, we observe that information of X is same as information in T as T_n is sufficient. We write a comment here, that information is additive. So, suppose I am considering independent random variables let X and Y be independent random variables with distributions say $f_1(x, \theta)$ and $f_2(y, \theta)$.

(Refer Slide Time: 47:02)

The joint density of X & Y is

$$g(x, y, \theta) = f(x, \theta) f(y, \theta)$$

$$\log g = \log f(x, \theta) + \log f(y, \theta)$$

$$\frac{\partial}{\partial \theta} \log g = \frac{\partial}{\partial \theta} \log f_1(x, \theta) + \frac{\partial}{\partial \theta} \log f_2(y, \theta)$$

$$\frac{\partial^2}{\partial \theta^2} \log g = \frac{\partial^2}{\partial \theta^2} \log f_1(x, \theta) + \frac{\partial^2}{\partial \theta^2} \log f_2(y, \theta)$$

Exp. expectations

$$-E\left[\frac{\partial^2}{\partial \theta^2} \log g(x, y, \theta)\right] = -E\left[\frac{\partial^2}{\partial \theta^2} \log f_1(x, \theta)\right] - E\left[\frac{\partial^2}{\partial \theta^2} \log f_2(y, \theta)\right]$$

$$I_{(X+Y)}(\theta) = I_X(\theta) + I_Y(\theta)$$

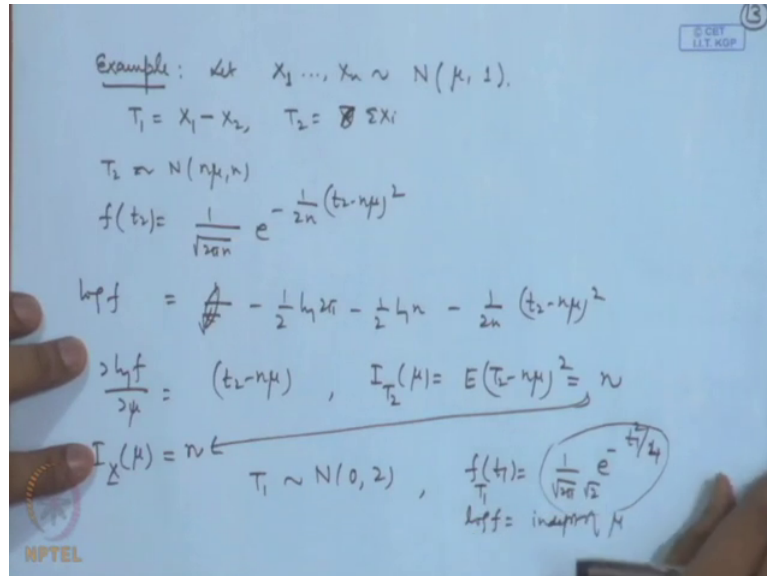
Then, let us consider information in x that is equal to expectation of del by del theta log of $f(x, \theta)$, theta whole square, which is also same as minus expectation of del 2 by del theta 2 log of $f(x, \theta)$, theta. We have seen this relationship similarly, information in y that is equal to information expectation of del by del theta log of $f(y, \theta)$ say, this is f_1 this is f_2 that we can also write as expectation of del 2 by del theta 2 log of $f_2(y, \theta)$. Information in x plus y so; that means, we will consider the joint distribution of the joint distribution of X and Y is because the distributions are independent, it is equal to the product of the $f(x, \theta)$ into $f(y, \theta)$.

So, if I take log of $f(x, \theta)$ into $f(y, \theta)$ it is equal to log of $f(x, \theta)$ plus log of $f(y, \theta)$. So, if I consider let me write this notation for this joint density here say g of x, y, θ , then log of g is equal to log of f plus log of f_1 plus log of f_2 . So, if I consider del by del theta log of g that is equal to del by del theta log of f_1 plus del by del theta log of f_2 . So, if I consider second order derivative del by del theta 2 log of g that is equal to del 2 by del theta 2 log of f_1 plus del 2 by del theta 2 log of f_2 .

So, if I take expectations here taking expectations. We get expectation of del 2 by del theta 2 log of $g(x, y, \theta)$ is equal to expectation of del 2 by del theta 2 log of $f_1(x, \theta)$ plus expectation of del 2 by del theta 2 log of $f_2(y, \theta)$. So, if I put a minus sign on both the sides then this is becoming information in x plus y and this is becoming information in x and this is information in y . So, we have proved that information is additive, if I am considering independent observations, then information in this total will be equal to the information in

one plus the information into the other one, but independence is used here. Let me explain the equality of sufficient statistic information by means of another example here.

(Refer Slide Time: 50:53)



Let us consider say x_1, x_2, \dots, x_n following say normal μ distribution. Let us consider say T_1 is equal to $x_1 - x_2$, T_2 is equal to \bar{x} or $\sum x_i$. So, what is the distribution of T_2 . T_2 will have normal $n\mu, n$. So, if we want to write down the distribution of this that is equal to $\frac{1}{\sqrt{2\pi n}} e^{-\frac{1}{2n}(t_2 - n\mu)^2}$.

That is equal to $\frac{1}{\sqrt{2\pi n}}$. So, if I take log of f I get $-\frac{1}{2} \log 2\pi n - \frac{1}{2} \log n - \frac{1}{2n} (t_2 - n\mu)^2$. So, $\frac{\partial \log f}{\partial \mu}$ that will be equal to $(t_2 - n\mu)$ because I get a $-n$ minus $2n$, here which will cancel out. So, if I consider information in T_2 that will be equal to expectation of $(t_2 - n\mu)^2$ that is equal to n .

Now, consider the information in the normal distribution we have already calculated it was equal to n in 1 of the variables it was one. So, in n variables it is n here which is matching here information in x about μ that was also equal to n . So, you can see that these two things are same let us look at the information in T_1 . Now, T_1 here is normal $0, 2$. So, the density function of T_1 will be free from μ . So, this is simply a constant because it is simply $\frac{1}{\sqrt{4\pi}} e^{-\frac{t_1^2}{4}}$.

You can see there is no μ occurring here. So, if I consider log of this, this is independent of μ and therefore, if I consider derivative with respect to μ that is going to be 0. So, information in T_1 is simply 0. Now, we will define this concept little later, if the information about the parameter is 0 in the statistic, it will be called ancillary statistic if the information is full; that means, whatever information in the whole sample is there and if that is equal then it is called a sufficient statistic.

So, this concept of information is very, very significant it actually tells the kind of statistic that we are considering and therefore, for what purpose it should be used. Now, I have also considered the cases that we can consider more than 1 sufficient statistics. So, we need to distinguish between different sufficient statistics in the sense that. What is the maximum reduction of the data that is possible? That is called the concept of minimal sufficiency. So, in the following lecture I will be starting the concept of minimal sufficiency. How to derive it? How to characterize the concept of minimal sufficiency? So, these are the topics that I will be covering in the next lecture.