

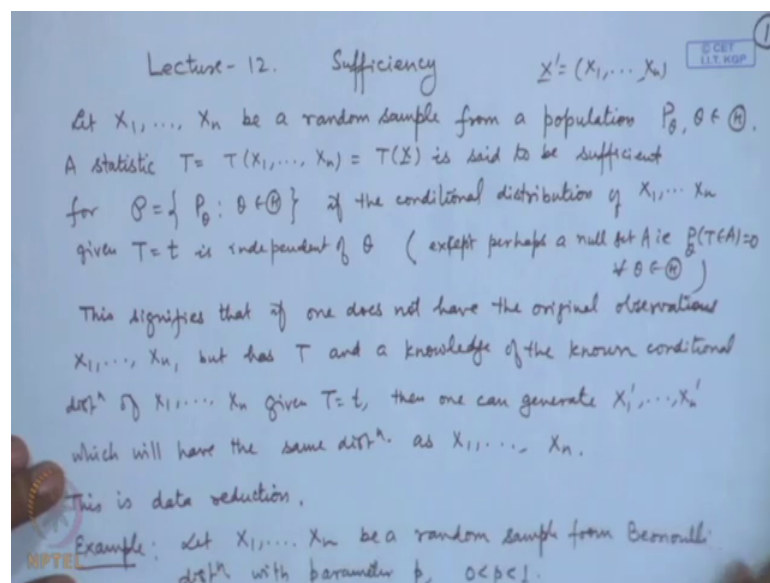
Statistical Inference
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture – 23
Sufficiency – I

Now, I start with a new concept that is called sufficiency. In the context of statistical inference, there is a concept which is useful to retain the necessary data without losing any information. What is the literal meaning of the word sufficiency? The literal meaning of the word sufficiency is that it is enough; sufficient means enough. So, usually we are dealing with the statistical model that we deal in the inference problem is that we say let X_1, X_2, \dots, X_n be a random sample meaning thereby that we have data on n observations or you can say n data points are available to us.

Now, in many of the practical problems, it becomes difficult to retain the data because, it may occupy lot of storage space whether it is on computer or it is in the form of hard copy of the data and then there is a danger of losing the data. It will be always interesting to say that let us keep the minimum things such that whatever information or whatever useful inferences, we want to make we are not suffering in that. That means, we do not lose any important part of it.

(Refer Slide Time: 01:55)



A formal specification of this concept is called sufficiency or sufficient statistic in the context of statistical inference.

So, let us introduce the formal definition of sufficiency as before we have a random sample. So, let X_1, X_2, \dots, X_n be a random sample from a population say $p(\theta)$ belonging to say script θ a statistic. So, a statistic we have already defined a statistic means a function of observations. So, T that is $T(X_1, X_2, \dots, X_n)$, which we also write as $T(X)$; that means we are denoting X as X_1, X_2, \dots, X_n . So, $T(X)$ is said to be sufficient. Now what do you mean by sufficient for what? So, we usually mention the word sufficient for the family of probability distributions.

In loose turns we also say sufficient for the parameter θ meaning thereby that whatever be the parameter under consideration, many times in the problems we will have 1 dimensional parameter 2 dimensional parameter etcetera. In that case, we will have to consider a specifically what parameter is being considered. So, the formal definition, I am writing for the family of probability distributions; meaning thereby that whatever parameters are under consideration this could be a scalar or a vector parameter.

So, this is said to be sufficient if the conditional distribution of X_1, X_2, \dots, X_n given T is equal to say small t is independent of θ . Of course, except perhaps a null set A that is on a set a where T takes probability 0. So, this is a exceptional case, but in general the distribution of the random sample given the statistic, if it is independent of the parameter then we say that this t is independent, then we say that this t is a sufficient statistic.

Now what is the physical interpretation of this definition that the distribution of X_1, X_2, \dots, X_n is free from θ and then we said is sufficient what does it mean? It means that now if the distribution is free from θ ; that means, the distribution of X_1, X_2, \dots, X_n given T is completely known. So, suppose we know T , we know the distribution of T now this conditional distribution of X_1, X_2, \dots, X_n given T , since it is free from θ then that is also known. Therefore, if I merge these 2 distributions that is the conditional distribution of X_1, X_2, \dots, X_n given T and the distribution of T I get the joint distribution of X_1, X_2, \dots, X_n and T from there I get the distribution of X_1, X_2, \dots, X_n .

It means that even if I may not have the initial X_1, X_2, \dots, X_n with us, but we can generate that distribution once again, because of the information or we can say the distribution of X_1, X_2, \dots, X_n given T being free from the parameter and T is known to us. This signifies

that if one does not have the original observations X_1, X_2, \dots, X_n , but has T and a knowledge of the known conditional distribution of X_1, X_2, \dots, X_n given T then one can generate say X_1^* , X_2^* , \dots , X_n^* , which will have the same distribution as X_1, X_2, \dots, X_n .

So, this is called data reduction as we will show later on that in most of the practical problems this sufficient statistics will become like 1 dimensional or 2 dimensional thing although, you may have any number of observations. So, this data reduction is helpful and we will show a statistically also that basing our decisions on the sufficient statistics is also useful. That means, if there is any inference made in the terms of estimation testing of hypothesis etcetera. If I am making inference based on the sufficient statistics we are better off.

So, let me explain this example say a binomial distribution example let me take. Suppose, I have X_1, X_2, \dots, X_n be a random sample from say Bernoulli distribution with parameter p a p lies between 0 and 1.

(Refer Slide Time: 09:09)

$T = \sum_{i=1}^n X_i \sim \text{Bin}(n, p)$ Consider the conditional distribution of X_1, \dots, X_n given $T=t$

$$P(X_1=x_1, \dots, X_n=x_n | T=t) = \frac{P(X_1=x_1, \dots, X_n=x_n, T=t)}{P(T=t)}$$

$$= \begin{cases} \frac{P(X_1=x_1, \dots, X_{n-1}=x_{n-1}, X_n=t - \sum_{i=1}^{n-1} x_i)}{P(T=t)}, & \text{if } t = \sum_{i=1}^n x_i \\ 0, & \text{if } t \neq \sum_{i=1}^n x_i \end{cases}$$

$$= \frac{p^{x_1} (1-p)^{1-x_1} \dots p^{x_{n-1}} (1-p)^{1-x_{n-1}} p^{t - \sum_{i=1}^{n-1} x_i} (1-p)^{1 - t + \sum_{i=1}^{n-1} x_i}}{\binom{n}{t} p^t (1-p)^{n-t}}$$

$$= \frac{p^t (1-p)^{n-t}}{\binom{n}{t} p^t (1-p)^{n-t}}$$

Let us consider say T is equal to $\sum_{i=1}^n X_i$ is equal to 1 to n . Let us look at the conditional distribution of consider the conditional distribution of X_1, X_2, \dots, X_n given T that is equal to X_1 is equal to X_1 and so on, X_n is equal to X_n given T is equal to t that is equal to probability of X_1 is equal to small x_1 and so on X_n is equal to small x_n T is equal to t divided by probability of T is equal to t .

Now, that is equal to since T is equal to $\sum X_i$, if $x_1 + x_2 + \dots + x_n$ is equal to t then only this probability will be calculated. In other cases, this will be simply equal to 0 so that is equal to probability of X_1 is equal to x_1 and so on, X_{n-1} is equal to x_{n-1} and X_n is equal to $t - \sum_{i=1}^{n-1} x_i$ is equal to $1 - \sum_{i=1}^{n-1} x_i$, if t is equal to $\sum_{i=1}^n x_i$ is equal to 1 to n otherwise this is 0.

Now here, we can make use of the fact that X_1, X_2, \dots, X_n are independently distributed Bernoulli random variables. So, if they are independent this probability of the joint occurrence will be equal to the product of these probabilities. So, this term let me write this term is anyway 0. So, this term is equal to probability of X_1 is equal to x_1 and so on X_{n-1} is equal to x_{n-1} probability of X_n is equal to $t - \sum_{i=1}^{n-1} x_i$ is equal to $1 - \sum_{i=1}^{n-1} x_i$ that is equal to p to the power X_1 $1 - p$ to the power $1 - X_1$ and so on; p to the power X_{n-1} $1 - p$ to the power $1 - X_{n-1}$ p to the power $t - \sum_{i=1}^{n-1} x_i$ is equal to $1 - \sum_{i=1}^{n-1} x_i$ $1 - p$ to the power $1 - t + \sum_{i=1}^{n-1} x_i$ is equal to $1 - \sum_{i=1}^{n-1} x_i$ and divided by probability t is equal to t .

Now, what is the distribution of t ? If X_1, X_2, \dots, X_n are Bernoullis independent then this will be binomial n, p . So, probability t is equal to t that will be equal to $\binom{n}{t} p^t (1-p)^{n-t}$. Now we can easily see these terms this p to the power terms if you add, you will get p to the power t . Similarly, if you add $1 - p$ exponents, you will get $n - \sum_{i=1}^n x_i$ so, that will cancel out with plus $\sum_{i=1}^n x_i$, you get $n - t$. So, you get it as p to the power t into $1 - p$ to the power $n - t$ divided by $\binom{n}{t} p^t (1-p)^{n-t}$.

(Refer Slide Time: 13:19)

$$\begin{aligned}
 & P(T=t) \\
 & \downarrow \\
 & P(X_1=x_1) \cdots P(X_{n-1}=x_{n-1}) P(X_n=t - \sum_{i=1}^{n-1} x_i) / P(T=t) \\
 & = \frac{p^{x_1} (1-p)^{1-x_1} \cdots p^{x_{n-1}} (1-p)^{1-x_{n-1}} \cdot p^{t - \sum_{i=1}^{n-1} x_i} (1-p)^{1 - t + \sum_{i=1}^{n-1} x_i}}{\binom{n}{t} p^t (1-p)^{n-t}} \\
 & = \frac{p^t (1-p)^{n-t}}{\binom{n}{t} p^t (1-p)^{n-t}} = \frac{1}{\binom{n}{t}}
 \end{aligned}$$

Now, this term simply cancels out. So, we get it as 1 by n c t. So, this conditional distribution then we can express as.

(Refer Slide Time: 13:31)

So $P(X_1=x_1, \dots, X_n=x_n | T=t) = \begin{cases} \frac{1}{\binom{n}{t}}, & t = \sum x_i \\ 0, & t \neq \sum x_i \end{cases}$

This is independent of p . So $T = \sum X_i$ is sufficient for $\{Ber(1, p), 0 < p < 1\}$.

Example: Let X_1, \dots, X_n be a random sample from $P(\lambda), \lambda > 0$
 $T = \sum X_i$

$$P(X_1=x_1, \dots, X_n=x_n | T=t) = \begin{cases} \frac{P(X_1=x_1, \dots, X_{n-1}=x_{n-1}, X_n=t - \sum_{i=1}^{n-1} x_i)}{P(T=t)} & t = \sum x_i \\ 0, & t \neq \sum x_i \end{cases}$$

Probability of X_1 is equal to small x_1 and so on X_n is equal to small x_n given T is equal to t that is equal to 1 by n c t for t is equal to $\sum x_i$ and it is equal to 0 , if t is not equal to $\sum x_i$. You look at this term there is no θ and no parameter appearing here p is not appearing here. So, T is equal to $\sum X_i$ is sufficient for the family of Bernoulli distributions.

We may also say it as that T sufficient for p here. Now note here, the physical significance of sufficiency, if we are observing X_1, X_2, \dots, X_n as n independent Bernoulli random variables. That means, they are observations related to success or failure in an Bernoulli and trials. For example, you are looking at a game of say dart and we are considering hitting a target and we make n aims at the target then what is important, whether individual hits whether this say second one hit correctly. Third one did not hit correctly, is it important information or out of n total attempts, how many are correct? That means, X that is some of exercise

Now, here you see in the concept of sufficiency exactly $\sum X_i$ is turning out to be sufficient. Therefore, this is the relevant information and whatever individual information about X_1, X_2, \dots, X_n is there that is not necessary to be written. In fact, now if we know this and we know the distribution of t that is binomial n, p , we can generate another random sample let us call it say X_1', X_2', \dots, X_n' , which will have Bernoulli $1, p$ distribution.

Let me explain through another example say, let X_1, X_2, \dots, X_n be a random sample from Poisson λ distribution, where λ is positive once again let us define T is equal to say $\sum X_i$. Now you can proceeding the same way like in the binomial case, we can consider X_1 is equal to X_1 and so on X_n is equal to X_n given T is equal t . So are going as before we get it as X_1 is equal to X_1 and so on X_{n-1} is equal to X_{n-1} minus 1 X_n is equal to t minus $\sum_{i=1}^{n-1} X_i$ is equal to 1 to n minus 1 divided by probability t is equal to t , if t is equal to $\sum_{i=1}^{n-1} X_i$ to n it is equal to 0, if t is not equal to $\sum_{i=1}^{n-1} X_i$ to n .

So once again, this term will be equal to $e^{-\lambda} \lambda^{X_i} / X_i!$ for i is equal to 1 to n minus 1.

(Refer Slide Time: 17:47)

This is independent of p . So $T = \sum X_i \sim \text{Bin}(n, p)$

$\{ \text{Ber}(1, p), 0 < p < 1 \}$

Example: let X_1, \dots, X_n be a random sample from $P(\lambda)$, $\lambda > 0$

$T = \sum X_i \sim P(n\lambda)$

$$P(X_1 = x_1, \dots, X_n = x_n | T = t) = \frac{P(X_1 = x_1, \dots, X_n = x_n, X_n = t - \sum_{i=1}^{n-1} x_i)}{P(T = t)}$$

$t = \sum_{i=1}^n x_i$

$$\prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \cdot \frac{e^{-\lambda} \lambda^{t - \sum_{i=1}^{n-1} x_i}}{(t - \sum_{i=1}^{n-1} x_i)!} \bigg/ \frac{e^{-n\lambda} (n\lambda)^t}{t!} = \dots$$

$t = \sum_{i=1}^n x_i$

$t \neq \sum_{i=1}^n x_i$

NPTL

And the last one is e to the power minus λ λ to the power t minus $\sum_{i=1}^n x_i$ 1 to n minus 1 divided by t minus $\sum_{i=1}^n x_i$ i is equal to 1 to n minus 1 factorial. Now this will be coming e to the power minus n minus 1 λ and then e to the power minus n λ and also we have in the denominator t . Now, this will follow Poisson n λ because, Poisson distribution is additive. So, if we are considering a random sample each one following Poisson λ then $\sum X_i$ will follow Poisson n λ . So, we can write e to the power minus n λ n λ to the power t by t factorial.

So, this e to the power minus n λ cancels out and if you look at λ to the power X_1 plus X_2 plus X_n minus 1 that cancels here, you get λ to the power t and in the denominator also we have λ to the power t here. So, what we get here? This t factorial will go in the numerator. So, let me write it here.

(Refer Slide Time: 19:01).

$$= \begin{cases} \frac{t!}{x_1! \dots x_n! (t - \sum_{i=1}^n x_i)!} & \text{if } t = \sum_{i=1}^n x_i \\ 0, & \text{if } t \neq \sum_{i=1}^n x_i \end{cases}$$

This is independent of λ . So $T = \sum X_i$ is sufficient for $\{P(\lambda): \lambda > 0\}$

Remarks: 1. Let T be sufficient for $\mathcal{P} = \{P_\theta: \theta \in \Omega\}$, and let U be a function of T . Then U is also sufficient for \mathcal{P} .

2. $P_\theta (X_1 = x_1, \dots, X_n = x_n \mid X_1 = t_1, \dots, X_n = t_n) = 1$ if $\sum x_i = \sum t_i$
 $= 0$ if $\sum x_i \neq \sum t_i$
 which is always free from the parameter.
 So $\underline{X} = (X_1, \dots, X_n)$ is always sufficient.

This is equal to t factorial divided by X_1 factorial X_2 factorial \dots X_n factorial t minus $\sum X_i$ is equal to 1 to n minus 1 factorial. If t is equal to $\sum X_i$ and it is equal to 0 , if t is not equal to $\sum X_i$ is equal to 1 to n

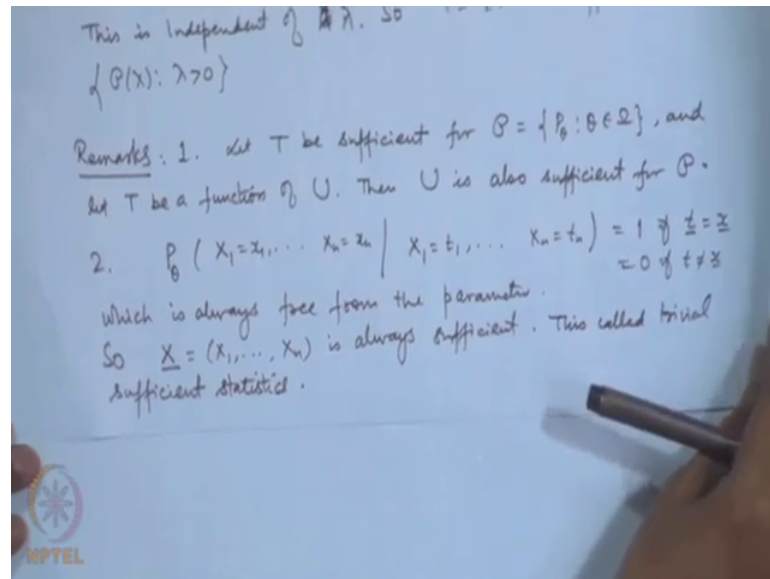
Once again, you notice here that this is independent of t . So, T is equal to $\sum X_i$, this is independent of λ sorry. So, T is equal to $\sum X_i$ is sufficient for the family of Poisson distributions, we may also say that $\sum X_i$ is sufficient for the parameter λ . Now we can make certain statements here if, I am considering conditional distribution of $X_1 X_2 \dots X_n$ given t and suppose t is a function of u . Then if I considered the conditional distribution of $X_1 X_2 \dots X_n$ given u then that will also be free from the parameter. Because, if that is not free from the parameter then the conditional distribution of $X_1 X_2 \dots X_n$ given t will also not be free from the parameter. Therefore, if t sufficient and t is a function of u then u is also sufficient and of course, if I have a 1 to 1 function of t then that will also be sufficient.

So, let me give some remarks here, let T be sufficient for a family of distributions and let T be a function of U then U is also sufficient for \mathcal{P} . Another point that you notice here that if I considered the conditional distribution of $X_1 X_2 \dots X_n$ given X_1 is equal to a small x_1 x_2 is equal to a small x_2 , x_n is equal to a small x_n . Then that is always independent of parameter, we can write conditional distribution of say X_1 is equal to X_1 X_n is equal to X_n given, say X_1 is equal to t_1 and so on, X_n is equal to t_n this is

equal to 1. If this t vector is same as x vector otherwise it is 0. So, this is naturally free from the parameter free from the parameter. So, the sample X is always sufficient.

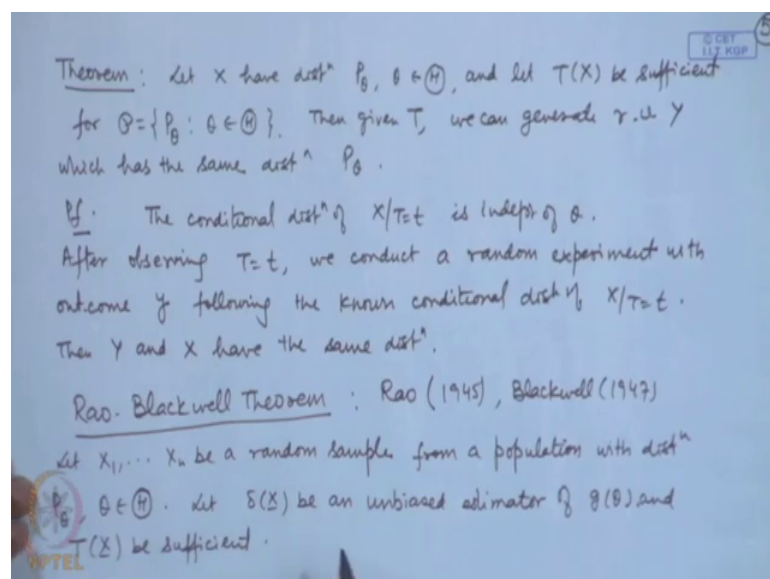
So, the full sample is always sufficient. We will be interested in getting some sort of reduction over there.

(Refer Slide Time: 22:59)



This is known as trivial sufficient statistics; trivial sufficient statistics.

(Refer Slide Time: 23:25)

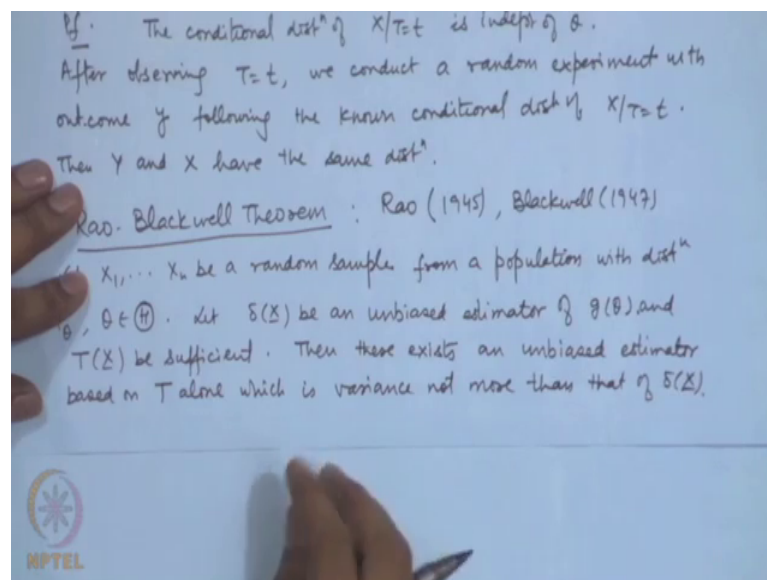


Let me formally prove that given a sufficient statistics, you can generate the original sample. So, let X have distribution say P_{θ} θ belonging to say script θ and let $T(X)$ be sufficient then given T . We can generate random variable y , which has the same distribution p_{θ} that is the same distribution of X . So, the conditional distribution of X given T is independent of θ . So, after observing T is equal to t , we conduct a random experiment with outcome, say y following the known conditional distribution of X given T then Y and X have the same distribution.

Now, another important significance you can say of sufficient statistics is that if we are considering any unbiased estimator. I can have another unbiased estimator, which is based on the sufficient statistics, and its variance will be less than or equal to the variance of the initial estimator this famous result is known as Rao Blackwell theorem.

It is named after the Indian statistician CR Rao, who proved this result in 1945 and David Blackwell 1947, let X_1, X_2, \dots, X_n be a random sample from a population with distribution p_{θ} θ belonging to say a script θ , let $\delta(X)$ be an unbiased estimator of parametric function say $g(\theta)$ and T be sufficient.

(Refer Slide Time: 27:05)



Then there exists an unbiased estimator based on T alone, which has variance not more than that of $\delta(X)$. Now this is a very very significant statement in a given problem, if I have a sufficient statistics then I can always base our unbiased estimators on that statistics. So, that I will do better than if I do not base it. That means, I will be utilizing

the full information in the sample for making my statistical inference the proof is in fact, not very difficult.

(Refer Slide Time: 28:15)

Proof. Let $h(t) = E[\delta(X) | T(X)=t]$ (this is independent of θ)
as T is sufficient. So $h(T)$ is a statistic

$$E_{\theta} h(T) = E_{\theta} E(\delta(X) | T=t) = E_{\theta} \delta(X) = g(\theta) \quad \forall \theta \in \mathcal{R}$$

So $h(T)$ is unbiased for $g(\theta)$.

$$\begin{aligned} \text{Var}_{\theta}(\delta(X)) &= E_{\theta} \{ \text{Var}(\delta(X) | T) \} + \text{Var}_{\theta} \{ E(\delta(X) | T) \} \\ &= \underbrace{E_{\theta} \{ \text{Var}(\delta(X) | T) \}}_{\geq 0} + \text{Var}_{\theta} \{ h(T) \} \end{aligned}$$

$\Rightarrow \text{Var}_{\theta} h(T) \leq \text{Var}_{\theta} \delta(X)$

Let us consider $h(t)$ to be expectation of $\delta(X)$ given T, X is equal to t . Since, we know the conditional distribution of X given T is independent of θ ; therefore, this expectation is going to be a function of t alone ok. This is independent of θ as T is sufficient. So, $h(t)$ is a statistic and I can consider it for my estimation purpose, let us consider expectation of $h(t)$. Now expectation of $h(t)$ is simply expectation of expectation $\delta(X)$ given t . Now this is nothing, but expectation of $\delta(X)$ that is equal to $g(\theta)$.

So, this new estimator that I have used $h(T)$ is unbiased. So, this is unbiased for $g(\theta)$ further. Let us consider say variance of $\delta(X)$. Now this variance of $\delta(X)$, I can express as expectation of variance $\delta(X)$ given T plus variance of expectation $\delta(X)$ given T . Now, this is equal to this quantity, if you see this is a non negative quantity and expectation of $\delta(X)$ given T be you are defined to be $h(T)$. So, this is equal to variance of $h(T)$. So, what we are getting? Variance of $\delta(X)$ is equal to variance of $h(T)$ plus a non negative quantity a non negative quantity. That means, variance of $h(T)$ is going to be less than or equal to variance of $\delta(X)$.

We will give applications of this result a little later.