

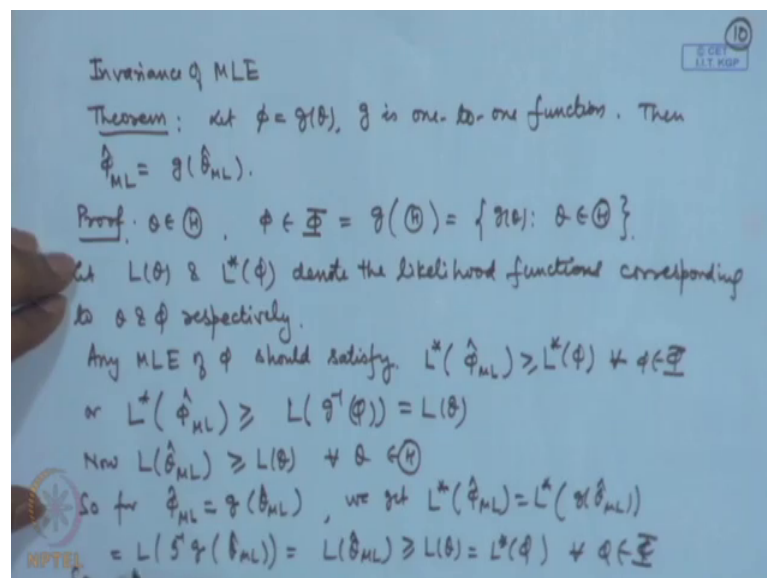
Statistical Inference
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture – 14
Properties of MLES - II

From now, there are certain other properties of the maximum likelihood estimators like invariance which make it very attractive. What is the meaning of invariance? Suppose, we are able to obtain as a natural parameters θ , θ_1 , θ_2 etcetera say suppose we have a one parameter problem and we have θ . So, we obtain the MLE of θ ; however, suppose in the given problem it may be required that θ^2 is the quantity of interest 1 by θ is a quantity of interest, \log of θ is a quantity of interest, in that case we can substitute the maximum likelihood estimator in that function.

In general, if we are considering function g of θ , then $g(\hat{\theta}_{ML})$ will be the actual MLE of $g(\theta)$.

(Refer Slide Time: 01:13)



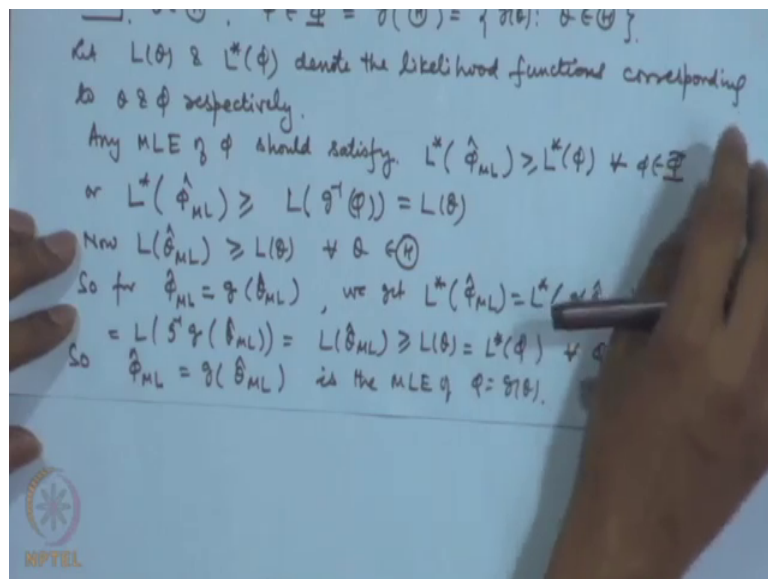
Now, this property I will be proving in two forms; invariance of MLE. Firstly, I will prove it for the one-one functions and then, actually I will give the general proof which is true for any function.

Let, ϕ be g of θ where g is a one to one function, then $\hat{\phi}_{ML}$ is equal to g of $\hat{\theta}_{ML}$. Let us look at the proof of this. So, suppose my parameter space for θ is capital script Θ and ϕ we wrote as capital Φ and this is actually the g of θ space, that is the set of all g of θ values as θ varies over a script Θ . Let $L(\theta)$ and $L^*(\phi)$ denote the likelihood functions corresponding to θ and ϕ respectively. Essentially, they are the same function because $L(\theta)$ is obtained as the joint distribution written at the point θ .

Now, in that you substitute because g of θ is equal to ϕ . So, if we substitute in terms of ϕ here in that function, then that will be a function of ϕ and we denoted by L^* . So, they are actually same functions, but written in as functions of different variables. Now, any maximum likelihood estimator of ϕ should satisfy $L^*(\hat{\phi}_{ML}) \geq L^*(\phi) \forall \phi \in \Phi$ or $L^*(\hat{\phi}_{ML}) \geq L(g^{-1}(\phi)) = L(\theta)$. Now, $L(\hat{\theta}_{ML}) \geq L(\theta)$ for all $\theta \in \Theta$. So for $\hat{\phi}_{ML} = g(\hat{\theta}_{ML})$, we get $L^*(\hat{\phi}_{ML}) = L(g^{-1}(\hat{\phi}_{ML})) = L(\hat{\theta}_{ML}) \geq L(\theta) = L^*(\phi) \forall \phi \in \Phi$. So $\hat{\phi}_{ML} = g(\hat{\theta}_{ML})$ is the MLE of $\phi = g(\theta)$.

So, for $\hat{\phi}_{ML}$ equal to g of $\hat{\theta}_{ML}$, we get $L^*(\hat{\phi}_{ML})$ is equal to $L^*(g$ of $\hat{\theta}_{ML})$ that is equal to $L(g^{-1}$ of g of $\hat{\theta}_{ML})$ that is equal to L of $\hat{\theta}_{ML}$ which is greater than or equal to $L(\theta)$ that is equal to $L^*(\phi)$ for all ϕ .

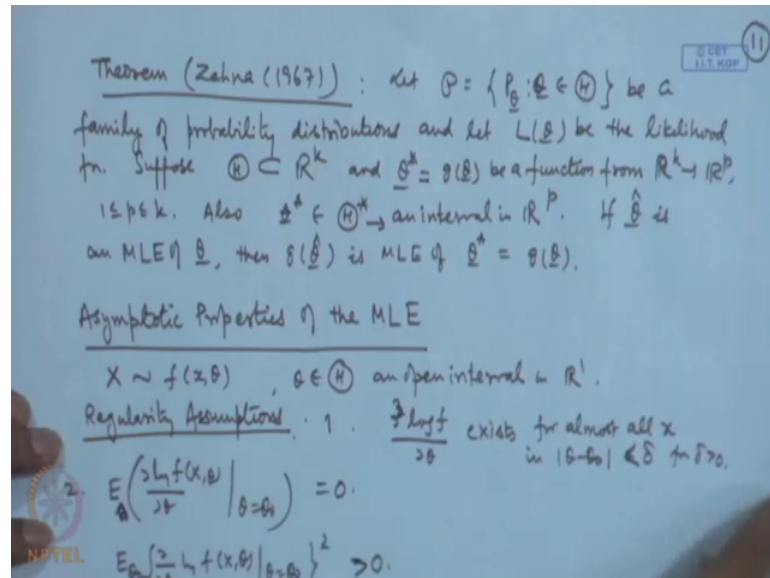
(Refer Slide Time: 04:51)



So, you can say that $\hat{\phi}_{ML}$ is equal to g of $\hat{\theta}_{ML}$ is the MLE of ϕ is equal to g of $\hat{\theta}_{ML}$. Now, naturally this result is true we have proved for g being a one-one

function. However, even if we have any function the same invariance property can be used as a justification for this was provided by Zahna 1967, I will state the result without proof here.

(Refer Slide Time: 05:26)



Let p_θ belonging to script θ be a family of probability distributions and let L_θ be the likelihood function. Suppose, θ is a subset of k dimensional Euclidean space and θ^* be a function from \mathbb{R}^k to \mathbb{R}^p where p will be less than or equal to k . Also, we assume that the range of θ^* , this is an interval in \mathbb{R}^p . So, if $\hat{\theta}$ is an MLE of θ , then $g(\hat{\theta})$ is MLE of θ^* . The proof of this requires that for every value of a g_θ , if there are several values which will lead to the same value because now we the function can be a an even function, then the likelihood function for θ^* is defined as the maximum of all those values.

In the case of one to one transformation, what we have done is L_θ is L of $g^{-1}(\phi)$ which we call $L^* \phi$ whereas, in the case of an even function there can be several values corresponding to one value of ϕ , in this case θ^* there can be several values. So, what we will do that for all those values we take the maximum of the likelihood function and then we maximize that.

So, when we associate maximum for each inverse image, what will happen is that, we are actually creating a one to one function and therefore, this theorem is once again applicable. However, I am just keeping the proof here for the details we can look at the

paper by Zahna in 1967. I will now give some more asymptotic properties of the maximum likelihood estimators. So, let me call it large sample properties or asymptotic properties of the maximum likelihood estimator.

Now, these properties are true under certain conditions which we usually call regularity conditions. Now, these conditions were initially given by Cramer and these are usually called Cramer Rao or, Fisher Cramer Rao regularity conditions. So, I will just call it regularity conditions.

So, in general we are considering a class of probability distributions. Now, they may have probability densities or probability mass functions. So, let me write that the density function or the mass function, this is a general notation I am using θ belongs to a script θ , this is an open interval in real line, then we have the following regularity assumptions. The assumptions are as follows.

That the up to the third order derivative exists and this should be for all θ ; however, it is enough if we assume it in a neighborhood of the solution. Suppose, we know that the solution exists around θ_0 , then if we assume this derivative existing and in interval or in a neighborhood of θ_0 , then it is enough; less than δ for some δ positive.

The second condition is that expectation of $\frac{\partial}{\partial \theta} \log f$ by $\frac{\partial}{\partial \theta}$ at θ_0 is equal to θ_0 that is equal to 0, expectation of $\theta_0 \frac{\partial}{\partial \theta} \log f$ by $\frac{\partial}{\partial \theta}$ at θ_0 is equal to θ_0^2 that is positive.

(Refer Slide Time: 11:11)

3. $\left| \frac{\partial^3 \ln f}{\partial \theta^3} \right| < M(x) \quad \forall \theta \in (\theta_0 - \delta, \theta_0 + \delta)$

$E_{\theta} M(x) < K \quad \forall \theta \in (\theta_0 - \delta, \theta_0 + \delta)$

Let X_1, \dots, X_n be i.i.d. as X and $X_i = x_i$ be observed values then

$L(\theta, \mathbf{x}) = \sum_{i=1}^n \log f(x_i, \theta)$

The likelihood equation is

$\frac{dL}{d\theta} = \sum_{i=1}^n \frac{f'(x_i, \theta)}{f(x_i, \theta)} = 0$ is the likelihood equation

Under these assumptions, we have the following large sample results for the MLE:

The third condition is that the third derivative of the density or the mass function is bounded in a neighborhood of theta naught and this function itself is having finite expectation a bounded expectation at any point in the interval theta naught minus delta to theta naught plus delta. So, if we are considering X_1, X_2, \dots, X_n independent and identically distributed as x and observed values are then the likelihood equation is the log likelihood is and the likelihood equation is $dL/d\theta = \sum_{i=1}^n f'(x_i, \theta)/f(x_i, \theta) = 0$. Now, here prime means derivative with respect to theta, this is the likelihood equation.

(Refer Slide Time: 13:02)

Theorem (Zehna (1967)): Let $\mathcal{P} = \{P_{\theta} : \theta \in \mathcal{H}\}$ be a family of probability distributions and let $L(\theta)$ be the likelihood fn. Suppose $\mathcal{H} \subset \mathbb{R}^k$ and $\underline{g} = g(\theta)$ be a function from $\mathbb{R}^k \rightarrow \mathbb{R}^p$, $1 \leq p \leq k$. Also $\mathcal{H}^* \subset \mathcal{H} \rightarrow$ an interval in \mathbb{R}^p . If $\hat{\theta}$ is an MLE of θ , then $g(\hat{\theta})$ is MLE of $\underline{g} = g(\theta)$.

Asymptotic Properties of the MLE

$X \sim f(x, \theta)$ pdf, $\theta \in \mathcal{H}$ an open interval in \mathbb{R}^1 .

Regularity Assumptions:

- $\frac{\partial \log f}{\partial \theta}$ exists for almost all x in $|\theta - \theta_0| \leq \delta$ for $\delta > 0$.
- $E_{\theta} \left(\frac{\partial \log f(x, \theta)}{\partial \theta} \right) = 0$.
- $E_{\theta_0} \left(\frac{\partial^2 \log f(x, \theta)}{\partial \theta^2} \right) \Big|_{\theta = \theta_0} < 0$.

Let us look at the conditions once again whatever assumptions we have made here f can be pdf or pmf θ belongs to the parameter space which is an open interval in the real line, we are assuming the derivative up to the third order exist. And, the assumption is at least for an interval in the neighborhood of the solution and then expectation of the first order derivative at θ_0 is equal to zero. The expectation of first derivative square that should be positive. In fact, we have defined earlier this as the information function.

(Refer Slide Time: 13:48)

The image shows a whiteboard with handwritten mathematical notes. At the top, it states $\frac{\partial \log L}{\partial \theta} = 0$. Below this, it says "the exact solution lies in a neighbourhood of a value say θ_0 ". The next line says "using $\frac{\partial \log L}{\partial \theta}$ in Taylor series around θ_0 upto second order terms and neglecting third and higher order derivatives".

The main derivation is as follows:

$$\frac{\partial \log L}{\partial \theta} = \frac{\partial \log L}{\partial \theta_0} + (\theta - \theta_0) \frac{\partial^2 \log L}{\partial \theta_0^2} + \dots$$

$$\approx \frac{\partial \log L}{\partial \theta_0} + (\theta - \theta_0) E\left(\frac{\partial^2 \log L}{\partial \theta_0^2}\right)$$

$$0 = \frac{\partial \log L}{\partial \theta_0} - \delta \theta I(\theta_0) \dots (1)$$

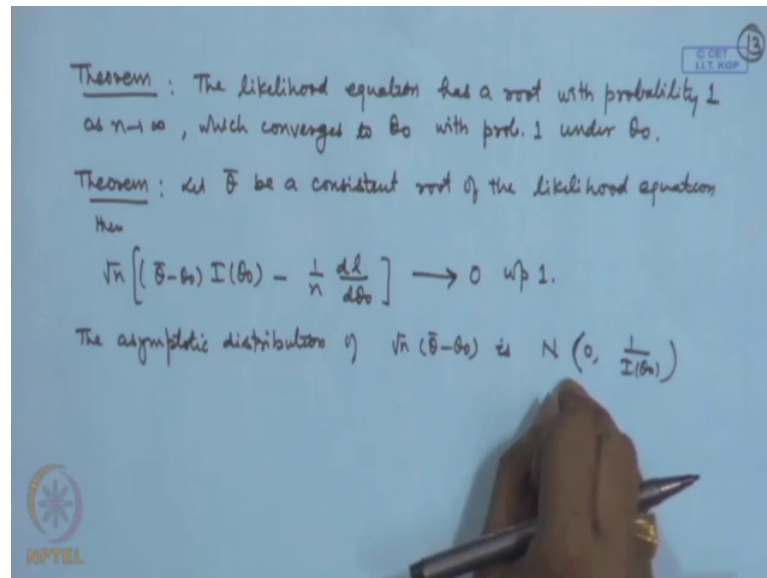
On the right side of the whiteboard, there are additional notes:

- X_1, \dots, X_n, \dots
- $X \rightarrow E(X)$
- $\theta = \theta_0 + \delta \theta$
- $I(\theta) = -E\left(\frac{\partial^2 \log L}{\partial \theta^2}\right)$
- $= E\left(\frac{\partial \log L}{\partial \theta}\right)^2$
- Fisher's information measure

If we look at this one expectation of $\frac{\partial \log L}{\partial \theta}$ whole square, this is called Fisher's information major Fisher's information measure. We will talk more about it somewhat later third assumption is that the third order derivative is bounded by an integrable function. Now, under these conditions, we have under these assumptions we have the following large sample results for the maximum likelihood estimator.

I will state it in the form of theorems without proof.

(Refer Slide Time: 14:47)



The first result says that the likelihood equation has a root with probability 1 as n tends to infinity and the root converges to θ_0 with probability 1 under θ_0 . So, this says that the likelihood equation has a root with probability 1 and the root converges to θ_0 with probability 1. So, this is a very important result and the second one says that let $\bar{\theta}$ be a consistent root of the likelihood equation, then
$$\sqrt{n} \left[(\bar{\theta} - \theta_0) I(\theta_0) - \frac{1}{n} \frac{d^2 l}{d\theta^2} \right] \rightarrow 0 \text{ w.p. 1.}$$
 The asymptotic distribution of $\sqrt{n}(\bar{\theta} - \theta_0)$ is $N(0, \frac{1}{I(\theta_0)})$.

Now, what does it mean that $\sqrt{n}(\bar{\theta} - \theta_0) I(\theta_0) - \frac{1}{n} \frac{d^2 l}{d\theta^2} \rightarrow 0$ with probability 1. This result is something like that the root of the likelihood equation that is the maximum likelihood estimator is asymptotically efficient because what term we are getting $\bar{\theta} - \theta_0$, you see if you take it to this side $\frac{1}{n} \frac{d^2 l}{d\theta^2}$ divided by $I(\theta_0)$. And, the asymptotic distribution of $\sqrt{n}(\bar{\theta} - \theta_0)$ is normal $N(0, \frac{1}{I(\theta_0)})$ that is as n becomes large, the distribution of $\bar{\theta} - \theta_0$ converges to that of a normal distribution with mean 0 and variance $\frac{1}{I(\theta_0)}$, where $I(\theta_0)$ is the Fisher's information measure.

The proofs of these results are not very difficult actually they use the laws of large numbers and the central limit theorem at various points. I am skipping the proof here and what is more important is that this is true under fairly general conditions for example, the

assumptions that I have stated now these assumptions are true for say binomial distribution, say for Poisson distribution, say for normal distribution, say for gamma distribution; that means, there is a large class of distributions particularly the distributions in the exponential family which satisfy these conditions Gaussian distribution is not in the exponential family, but even that also satisfies this condition. So, there is a fairly large class of distributions and densities which will actually satisfy this property.

So, under fairly general conditions we can say that the likelihood equation has a solution, the solution is consistent with probability 1; that means, it converges to the true value with probability one moreover the asymptotic distribution is normal and it is also second order efficient in the sense of Rao. So, that makes the use of likelihood maximum likelihood estimators a fairly important practice in the statistical theory.

See this property which I have written at the end that is a square root n $\bar{\theta}$ minus θ_0 or say square root and $\bar{\theta}$ minus θ_0 is asymptotically normal. This is also known as consistent as asymptotic normal property or kind estimator. So, if an estimator is consistent as well as its asymptotic distribution is normal. So, naturally these are having some desirable properties, we also say best asymptotically normal estimator that is best estimator so, can be. So, under certain conditions such estimator have exist and maximum likelihood estimators are more likely to satisfy these properties. I will be completing this discussion now.