

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

Module No. # 01
Lecture No. # 07
Properties of MLEs

In the last lectures, I have derived the form of the maximum likelihood estimators for various probability models. I have demonstrated how the role of likelihood is there, in determining the final form of the estimator. Suppose there is a prior information then it has an effect. Now in today's lecture I will spend some time, on discussing important properties of the maximum likelihood estimators. First of all we note that, see various problems we have done and in most of those problems you have got a value of the maximum likelihood estimator; that means, there is a function which is corresponding to the estimator. However, that is not necessarily the case; sometimes we may have a non-uniqueness.

(Refer Slide Time: 01:07)

Lecture - 7. Properties of MLEs.

Non-uniqueness of the MLE
Let $X_1, \dots, X_n \sim U[\theta-a, \theta+a]$ $\theta \in \mathbb{R}, a > 0$
Here a is a known constant.
The likelihood function is
$$L(\theta, \underline{x}) = \frac{1}{(2a)^n} \quad \theta-a \leq x_{(1)} \leq \dots \leq x_{(n)} \leq \theta+a$$
$$= 0, \quad \text{elsewhere.}$$

 L is maximum when $\theta-a \leq x_{(1)}$ or $\theta \leq x_{(1)}+a$
and $\theta+a \geq x_{(n)}$ or $\theta \geq x_{(n)}-a$
So any value θ between $x_{(n)}-a$ to $x_{(1)}+a$ is a MLE of θ
We may choose the midpoint, i.e. $\frac{x_{(1)}+x_{(n)}}{2}$ as the MLE.

So, let me give an example of that, non-uniqueness of the M L E. Let me discuss one example of that, let we have a sample x_1, x_2, \dots, x_n , from a uniform distribution on the interval

theta minus a to theta plus a, where theta is a real number, and a is a positive number. Here, a is a known constant. So, the problem is here to estimate, the parameter theta of this uniform distribution; that means, the spread is from theta minus a to theta plus a. Since theta is unknown we do not know the starting and the end point of the spread. So, let us consider the likelihood function. The likelihood function is. Now here the density function is $\frac{1}{2a}$. So, if I consider the joint distribution of x_1, x_2, \dots, x_n , it will become $\frac{1}{2a}$ to the power n.

And each of the x_i 's lies from theta minus a to theta plus a, therefore we can summarize this information in the form, that theta minus a less or equal to x_1 and so on, less than or equal to x_n , less than or equal to theta plus a. Let me put it here close interval, because I am including the end points here it is equal to 0 elsewhere. Naturally, you can see that the maximum value of the likelihood function is $\frac{1}{2a}$ to the power n, because this is a constant value here at other points it is 0. Now, this is satisfied when this inequality holds true. Therefore, let us see the optimal range of theta for which this value is attained. So, L is maximum, when theta minus a is less than or equal to x_1 , or you can say that, theta is less than or equal to x_1 plus a and theta plus a is greater than or equal to x_n , or theta is greater than or equal to x_n minus a. Naturally if I choose any value of theta in the interval x_n minus a to x_1 plus a, that will be the maximum likelihood estimator.

So, any value of between x_n minus a to x_1 plus a is, a maximum likelihood estimator of theta. So, this is an example where the maximum likelihood estimator is not unique. However, we may choose the, for example the midpoint of this, that will be the $\frac{x_1 + x_n}{2}$. We may choose the midpoint that is $\frac{x_1 + x_n}{2}$, as the M L E. now another feature which we noticed in the various problems that we have done, that in most of the cases we got it a very nice function, for example we go it as \bar{x} , one by n sigma x_i minus \bar{x} whole square, the median, the largest or the smallest etcetera. In most of these cases the maximum likelihood estimator is in a closed form, and also a mathematically elegant form, but even that is not necessary. Let me take another example where we do not get a nice analytic form.

(Refer Slide Time: 05:32)

MLE need not be in a nice analytic form.

Let $X_1, \dots, X_n \sim N(\theta, \theta^2)$, $\theta > 0$

$$L(\theta, \underline{x}) = \frac{1}{(\theta\sqrt{2\pi})^n} e^{-\frac{1}{2\theta^2} \sum (x_i - \theta)^2}, \quad x_i \in \mathbb{R}, \theta > 0$$

$$\begin{aligned} \ell(\theta) &= \log L(\theta, \underline{x}) = -n \log \theta - \frac{n}{2} \log 2\pi - \frac{\sum (x_i - \theta)^2}{2\theta^2} \\ &= -\frac{n}{\theta} + \frac{\sum (x_i - \theta)}{\theta^2} + \frac{\sum (x_i - \theta)^2}{\theta^2} \\ &= \frac{1}{\theta^3} \left[\sum (x_i - \theta)^2 + \theta \sum (x_i - \theta) - n\theta^2 \right] \\ &= \frac{1}{\theta^3} \left[\sum x_i^2 - 2n\theta\bar{x} + n\theta\bar{x} - n\theta^2 - n\theta^2 \right] \\ &= \frac{1}{\theta^3} \left[\sum x_i^2 - n\theta\bar{x} - n\theta^2 \right] \end{aligned}$$

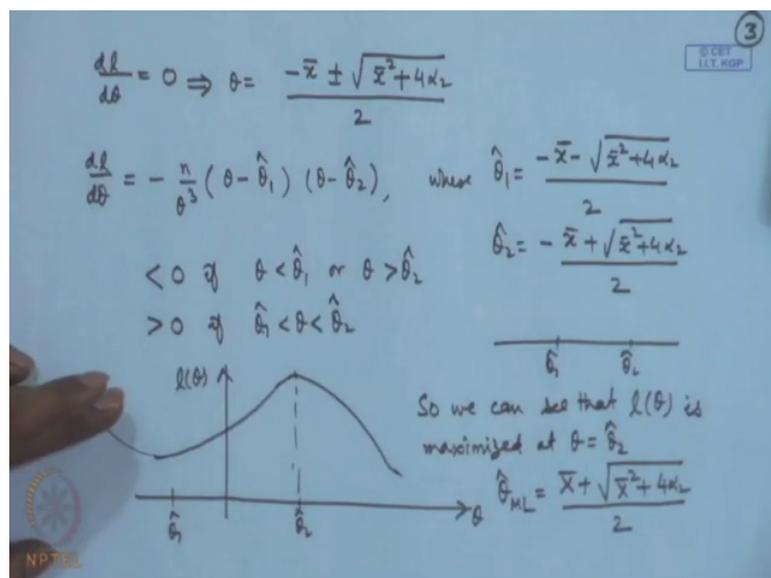
M L E need not be in a nice analytic form. Let me take one example here, let $x_1 \times 2 \times n$ be a random sample from a normal distribution, with mean θ and variance θ^2 . So, naturally here θ is a positive parameter here, the likelihood function. So, here your variance is actually the square of the mean, so there is an inter relationship here. So, the problem reduces to one parameter. So, if we consider the likelihood function, it is a joint distribution 1 by $\theta \sqrt{2\pi}$ to the power n , minus e to the power minus 1 by $2\theta^2$, $\sum x_i - \theta$ whole square, and here each of the x_i is on the real line whereas, θ is positive. So, if we consider the log likelihood, that is equal to minus $n \log \theta$ minus $\frac{n}{2} \log 2\pi$ minus $\frac{\sum x_i - \theta$ whole square divided by $2\theta^2$. There is naturally a difference from the situation when we had considered $\mu \sigma^2$ square, because then we had two parameters, and we had considered the maximization with respect to both of them.

Now since μ has been replaced by θ , so this is a consolidated function of θ that is coming here, and we have to maximize with respect to θ , but nevertheless this is a differentiable function, and therefore, we can think of the usual calculus procedure. Let us look at $d\ell$ by $d\theta$. So, that is equal to minus n by θ , plus $\sum x_i - \theta$ divided by θ^2 , plus $\sum x_i - \theta$ whole square divided by θ^3 , because the derivative of 1 by θ^2 will be minus 2 by θ^3 , so that simplifies it to this. This term we can write as $\sum x_i - \theta$ square, plus $\theta \sum x_i - \theta$, minus n

theta square divided by theta cube. Now we can expand these terms here, you will get 1 by theta cube, sigma x i square minus twice. Now, you get 2 theta x i, when you put summation here it becomes sigma x i, which we can write as n x bar.

So, minus 2 n theta x bar, then you have theta sigma x i, again which we can write as n x bar, so this becomes n theta x bar, minus theta square and this is summation here. So, minus n theta square, minus n theta square, and of course there was another term here which we missed here, theta square here. So, n theta square will come here with a plus sign, plus n theta square. So, naturally this term cancel's out, and here one of the n theta x bar cancels out. We get here 1 by theta cube, sigma x i square, minus n theta x bar, minus n theta square. Now we can consider the, if we put this equal to 0 then this is nothing, but a quadratic equation in theta, which will have two roots, because the denominator is theta cube which is always positive.

(Refer Slide Time: 09:59)



So, we can look at this d l by d theta is equal to 0, gives theta is equal to. So, we can write this equation, this sigma x i square term is coming we can take out n here. So, this we can write as minus n by theta cube, theta square plus theta x bar minus 1 by n sigma x i square, I will use the notation alpha 2. This alpha 2 notation we have introduced in the method of moments, this is the second sample moment. This alpha 2 is 1 by n sigma x i square. So, if we put this equal to 0, we can straight forwardly apply the b square minus 4 a c formula. So, we get theta is equal to minus x bar, plus minus square root x bar square, plus 4 alpha 2 divided

by 2. So, naturally there are two solutions, and we have to see the increasing and the decreasing nature of this. So, we can express $d \ln l$ by $d \theta$ as, equal to $-\frac{n}{\theta^3} (\theta - \theta_1)$; let me call it θ_1 head into $\theta - \theta_2$ head. Where I am taking θ_1 head to be the solution with the negative, that is $-\bar{x} - \sqrt{\bar{x}^2 + 4\alpha^2}$ divided by 2. And θ_2 head is equal to $-\bar{x} + \sqrt{\bar{x}^2 + 4\alpha^2}$ by 2.

Now let us look at the sign scheme of this. This term will be negative, if θ is less than θ_1 head or θ is greater than θ_2 head. Because if θ is less than θ_1 head, this term is negative. Here we can see that θ_1 head is less than θ_2 head. So, if we are considering θ_1 head and θ_2 head here. So, if θ is below θ_1 head and θ_2 head, then both of these terms are negative, their product is positive, so this entire term $d \ln l$ by $d \theta$ will become negative. Similarly, if θ is greater than θ_2 head, then this term is positive as well as this term is positive. So, the overall term will become negative. And this will become positive if θ_1 head is less than θ is less than θ_2 head. Therefore, we can look at the behavior of the likelihood function as θ varies.

Of course, we can actually plot it here; θ_1 head will be somewhere here, because both the terms are negative here. Minus \bar{x} of course, minus \bar{x} could be, because \bar{x} can be negative or positive, but this term is certainly negative and it is bigger. So, naturally I think this will become negative, whereas θ_2 head is going to be positive. The function is decreasing before θ_1 head. So, something like this, and then it will increase between θ_1 head to θ_2 head, and from θ_2 onwards again it will start decreasing. So, you can see that, at θ is equal to θ_2 head we get a maximizing value. So, we can see that, $\ln l$ is maximized at θ is equal to θ_2 head. And another point which we notice here, that this is actually a positive value. And from our model that we have considered here, θ should be positive.

So, it is natural that our maximum likelihood estimator confirms to that range here, and it is happening here. So, the maximum likelihood estimator is $\bar{x} + \sqrt{\bar{x}^2 + 4\alpha^2}$ by 2. Naturally, you can see that the form of the maximum likelihood estimator is not in a nice analytic form. In fact, you are getting square roots. So, once again taking expectation etcetera, checking whether it is unbiased and all those things will be quite complicated. So, the statement that maximum likelihood estimator need not be in a nice analytic form. We may have even more difficult situation; that is we may not be able to solve

the likelihood equation. In this case although solution is coming it is not in a good form, but there may be a situation, where we may not be able to solve it explicitly. So, let me give an example of that situation also.

(Refer Slide Time: 15:43)

MLE may not be in a closed form
 Let $X_1, \dots, X_n \sim \text{Gamma}(r, \lambda)$
 Case: λ is known and r is unknown ($\lambda=1$)

$$L(r, x) = \prod_{i=1}^n \left[\frac{\lambda^r}{\Gamma(r)} e^{-\lambda x_i} x_i^{r-1} \right]$$

$$= \frac{\lambda^{nr}}{(\Gamma(r))^n} e^{-\lambda \sum x_i} (\prod x_i)^{r-1}$$

$$\ell(r) = \log L = nr \log \lambda - n \log \Gamma(r) - \lambda \sum x_i + (r-1) \log \prod x_i$$

M L E may not be in a closed form. Let us consider, say a random sample from a gamma distribution, with parameter say r and λ . Now there can be two cases as we have seen earlier. r could be known λ may be unknown, and r could be unknown and λ may be known or both may be unknown. Let us consider the case, λ is known, and r is unknown. If λ is known since it is occurring as a scale parameter we may take it to be 1. Now, let us consider the likelihood function. So, likelihood function will be a function of r now, λ to the power r by $\Gamma(r)$, e to the power minus λx_i x_i to the power r minus 1, product i is equal to 1 to n . So, this if you take λ to the power $n r$ by $\Gamma(r)$ to the power n , e to the power minus λ , $\sum x_i$ product x_i to the power r minus 1. Let us take the log of this; that is equal to $n r \log \lambda$, plus minus $n \log \Gamma(r)$, minus $\lambda \sum x_i$, plus r minus one log of product x_i .

(Refer Slide Time: 15:43)

Let $X_1, \dots, X_n \sim \text{Gamma}(\gamma, \lambda)$

Case: λ is known and γ is unknown ($\lambda=1$)

$$L(\gamma, \underline{x}) = \prod_{i=1}^n \left[\frac{\lambda^\gamma}{\Gamma(\gamma)} e^{-\lambda x_i} x_i^{\gamma-1} \right]$$

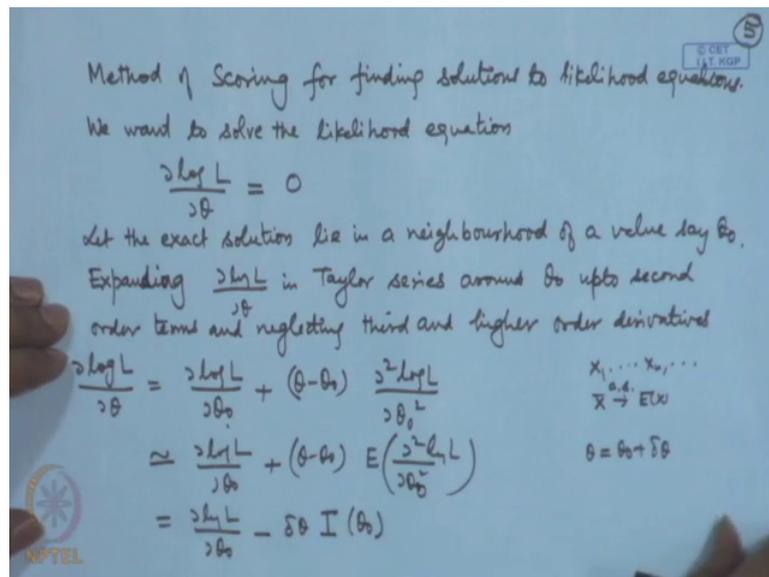
$$= \frac{\lambda^{nr}}{(\Gamma(\gamma))^n} e^{-\lambda \sum x_i} (\prod x_i)^{\gamma-1}$$

$$\log L = nr \log \lambda - n \log \Gamma(\gamma) - \lambda \sum x_i + (\gamma-1) \sum \log x_i$$

$$\frac{d \log L}{d \gamma} = -\frac{n}{\Gamma(\gamma)} \Gamma'(\gamma) + \sum \log x_i = 0 \quad \text{Euler's digamma fn.}$$

Now, let us put lambda is equal to 1 here then this term becomes much simpler. This particular term vanishes here. You get minus n log of gamma r, minus sigma x i plus r minus 1, sigma log of x i. Now if we treat it as a function of r then derivative of this with respect to r, will give me minus n by gamma r into gamma prime r. This is known as digamma function, minus this will become 0 plus sigma log of x i. So, this if we put 0, gamma prime r by gamma r, this is known as Euler's Digamma Function. So, it is not a very nice analytic function, and you cannot solve that this is equal to what will be the solution. You will have to use a numerical method such as; say Newton Raphson method or any other numerical method to solve this non-linear equation. Now in the cases when the explicit solution of the likelihood equation is not possible, a modification to the Newton Raphson method was suggested by fisher, and this is known as the method of scoring.

(Refer Slide Time: 19:22)



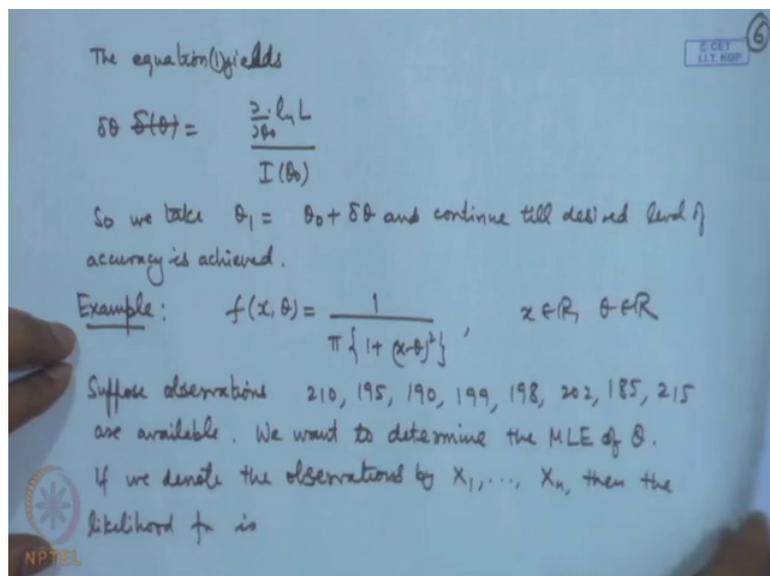
The method of scoring for finding solutions to likelihood equations. So, the method is briefly as follows; we are actually looking at, del by del theta log of l is equal to 0, we are trying to solve this equation. So, in case the solution is existing there is no problem. However, there may be cases when the solution exists, but we are not able to get exact analytic form. So what we will do, consider we want to solve the likelihood equation, del log l by del theta is equal to 0. And suppose the exact solution lies in the neighborhood of theta naught, or you can say it lies, it could be theta naught or we can assume that it is near about theta naught. Let the exact solution lie in a neighborhood of a value, say theta naught. Now once again here we will use the, techniques of analysis to determine. So, for example, if we want to find out the roots of a an in general non-linear equation, then what do we do.

We study the behavior of the function for example, if we are saying f x is equal to 0, then we look at the behavior of the function, we try to locate the roots where they may be lying and then we apply any numerical method, because generally the initial approximation is important. For example, Newton Raphson method, in one initial approximation is required. If we are using say bisection method, then two initial approximations are required, such that both of them are on the either side of the solution. So, similarly here we guess the initial roots, say theta naught. And let us consider expanding, say del log l by del theta in taylor series around theta naught. Now once we say Taylor series expansion, we are making the assumption that the derivatives of this exist. So, let us consider, only up to second order, up to

second order terms and neglecting third and higher order derivatives. So, basically it means the derivative evaluated at theta naught, similarly here it means the derivative evaluated at theta naught, the second derivative. Now this term, what fisher suggested we approximate by; that means, in place of this term we have written expectation.

Now there are certain justifications for this, for example this likelihood function is the joint distribution, so when we are taking log it is becoming summation here. So, this becomes summation term here. Now we know by the laws of large numbers that, if I have x_1, x_2, \dots, x_n a sequence of i i d random variables then \bar{x} converges to expectation \bar{x} almost surely, that is with probability one; that means, if n is large enough this approximation is all right. So, we have replaced this term by its expectation and this expectation of this with a minus sign, is known as the fisher's information. So, this is equal to $\text{del log l by del theta naught, minus}$. Now this theta point we consider in the neighborhood of theta, so let us write it as delta theta . So, this is $\text{delta theta and minus expectation}$ this is called fisher's information, at the point theta naught.

(Refer Slide Time: 25:03)



So, now from here what we get, the equation yields; let me call it equation number one. The equation one yields $\text{delta theta is equal to, delta by delta theta naught log l divided by i theta naught}$, because what we are going to do we are having $\text{delta log l by delta theta}$ is equal to 0. That means we are putting this term is equal to 0 here. So, if we put 0, then we can simplify this and we get the first approximation delta theta , as $\text{del by del theta naught log l by i theta}$

naught. So, we take the next, it rate as theta naught plus this delta theta and continue. So, theta 2 will then again become, where delta theta will be evaluated at theta 1 and so on, so continue till desired level of accuracy is achieved. So, this modified method it is known as Fisher's Newton Raphson method or fisher's scoring method. I will explain it through one example. Let us consider say, we have observations from a Cauchy distribution, where x is any real number and theta is any real parameter. Suppose observations 210, 195, 191, 99, 198, 202, 195 and 2 1 5 8 observations are available, and we want to determine the maximum likelihood estimator of theta, based on this sample. Now, in general if we write x_1, x_2, \dots, x_n , then what will be the likelihood function in the case of.

(Refer Slide Time: 28:13)

$$L(\theta, \mathbf{x}) = \prod_{i=1}^n \frac{1}{\pi \sqrt{1 + (x_i - \theta)^2}}$$

$$l(\theta) = \log L = \sum_{i=1}^n \log \frac{1}{\pi \sqrt{1 + (x_i - \theta)^2}}$$

$$= -n \log \pi - \sum_{i=1}^n \log \{1 + (x_i - \theta)^2\}$$
 The likelihood equation is $\frac{dl}{d\theta} = 0$

$$\Rightarrow 2 \sum_{i=1}^n \frac{(x_i - \theta)}{1 + (x_i - \theta)^2} = 0 \quad \dots (2)$$
 For $n=1$, we get $\hat{\theta} = x_1$, however for $n \geq 2$ it is a non-linear equation of degree $2n-1$.

if we denote the observations by say x_1, x_2, \dots, x_n , then the likelihood function is, product i is equal to 1 to n , 1 by $\pi \sqrt{1 + x_i - \theta \text{ square}}$. So, if we take log of this, we will get $\sum \log$ of 1 by $\pi \sqrt{1 + x_i - \theta \text{ square}}$, which we can write as $-n \log \pi$, minus $\sum \log$ of $1 + x_i - \theta \text{ whole square}$. So, if we look at the likelihood equation, $d l$ by $d \theta$ is equal to 0 , this is equivalent to. This term will yield 0 if you differentiate and here you will get 1 by $1 + x_i - \theta \text{ square}$ and then derivative of that, that will be 2 times $x_i - \theta$ with a minus sign. So, we will get $2 \sum x_i - \theta$, divided by $1 + x_i - \theta \text{ square}$, i is equal to 1 to n . Now naturally you can see here, this equation you cannot solve for n greater than or equal to 2 . For n is equal to

1, this will give simply theta is equal to x 1. For n is equal to 1, we get theta head is equal to x 1. However, for n greater than or equal to 2, it is a non-linear equation.

In fact, even if you write two terms here, then you will get x 1 minus theta by one plus x 1 minus theta square, plus x 2 minus theta divided by 1 plus x 2 minus theta square. And; obviously, that equation we will be having terms up to theta cube in the numerator. So, in general if I am writing n terms here, then n in each of the terms you will get square in the denominator. So, if you multiply n minus 1 of them, you will get 2 n minus 2 and then numerator. So, this will give me a non-linear equation of degree 2 n minus 1. So, naturally we cannot solve this theoretically. Let us apply the method of scoring in this problem. Now method of scoring involves, as we have seen just now that we should calculate the term called i theta naught, i theta naught is obtained as minus expectation del two log l by del theta naught square. And that is also equal to, let me write it here i theta is equal to minus expectation del 2 log l by del theta square. This is also equal to expectation of del log l by del theta whole square. So, one can do it in either way. Let us look at the calculation for this part.

(Refer Slide Time: 32:10)

We therefore apply the method of scoring here.

$$\log f = -\ln \pi - \ln (1 + (x-\theta)^2)$$

$$\frac{\partial \log f}{\partial \theta} = \frac{2(x-\theta)}{1+(x-\theta)^2}$$

$$E\left(\frac{\partial \log f}{\partial \theta}\right)^2 = 4 E \frac{(x-\theta)^2}{\{1+(x-\theta)^2\}^2} = \frac{4}{\pi} \int_{-\infty}^{\infty} \frac{(x-\theta)^2}{\{1+(x-\theta)^2\}^2} dx = \frac{4}{\pi} \int_{-\infty}^{\infty} \frac{y^2}{(1+y^2)^2} dy$$

$$\frac{8}{\pi} \int_0^{\infty} \frac{y^2}{(1+y^2)^2} dy = \frac{8}{\pi} \int_0^{\pi/2} \frac{\tan^2 \theta \sec^2 \theta}{(\sec^2 \theta)^2} d\theta = \frac{8}{\pi} \int_0^{\pi/2} \sin^2 \theta \cos^2 \theta d\theta = \frac{1}{2}$$

$$I(\theta) = \frac{n}{2} \quad \text{So } \delta(\theta) = \frac{4}{n} \cdot \sum_{i=1}^n \frac{(x_i - \theta)}{1 + (x_i - \theta)^2}$$

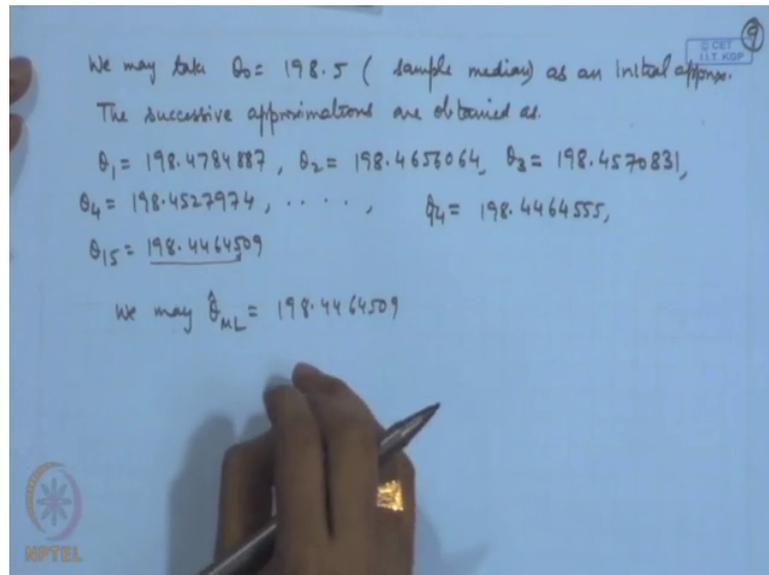
So, we therefore, apply the method of scoring here. Now for cauchy distribution, log of f is equal to minus log of pi, minus log of 1 plus x minus theta square. So, del log f by del theta, is equal to twice x minus theta divided by 1 plus x minus theta square. So, expectation of del log f by del theta whole square; that is equal to four times expectation of x minus theta square, divided by 1 plus x minus theta square, whole square. So, we evaluate this, that is

equal to four times integral x minus θ square, divided by 1 plus x minus θ whole square, and whole square of that, multiplied by the density function of x . And the density function of that is 1 by π 1 plus x minus θ square, so you will get power three here, and 1 by π , I will write here this is from minus infinity to infinity. Now you can easily transform this by putting x minus θ is equal to y . So, you get this as equal to four by π , y square divided by 1 plus y square whole cube dy , easily you can see that, this is an even function.

So, this becomes eight by π , 0 to infinity y square divided by 1 plus y square whole cube dy . Now this type of integral is standard we can substitute something like y is equal to $\tan \theta$. So, this will give me eight by π , 0 to π by 2 , \tan square θ \sec square θ divided by \sec cube θ , \sec square θ whole cube, $d\theta$ and that is equal to 8 by π , 0 to π by 2 \sin square θ \cos square θ $d\theta$, and that is equal to half. So, this you can see it is free from θ . So, information at the point θ naught, that will become n by 2 . Now the function that you need to calculate for the scoring method is $\Delta \theta$. $\Delta \theta$ is equal to $\frac{d}{d\theta}$ \log of θ divided by i of θ naught.

So, if we look at this term here, $\Delta \theta$, that will be equal to four by n $\sum x_i$ minus θ , divided by 1 plus x_i minus θ whole square, i is equal to 1 to n . Now the question is that what should be the initial approximation. Now, in the Cauchy distribution the sample mean is inconsistent, because we have seen the distribution of the sample mean is the same as that of the initial observation x_i , each x_i . However, we can see that sample median will be a consistent estimator here. Now from the given data of this, the middle observation will turn out to be. Two middle observations are there that is 198 and 199 , because there are if you arrange it in the ascending or descending order, then these are the two middle observations. So, the midpoint of that can be considered as the initial approximation.

(Refer Slide Time: 36:37)

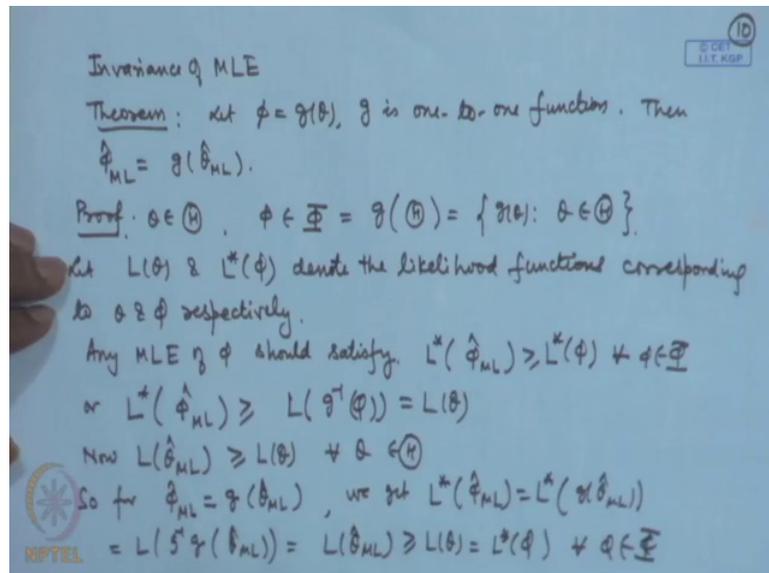


So, we may take θ_0 is equal to 198 point 5. The sample median as an initial approximation. Now, for each θ . So, if I take θ_0 here 198 point 5, here we substitute n is equal to 8 and we have the data available to us, x_i 's in the form of these values 210, 195 etcetera. So, if we substitute these values from the initial approximation we can get the successive approximations. The successive approximations are obtained as θ_1 is equal to 198 point 4784887, you have done up to seven decimal places. Next approximation give 198 point 4656064, θ_3 is 198 point 4570831, θ_4 is equal to 198 point 4527974 and so on. If we look at 14 that is 198 point 4464555, and θ_{15} is equal to 198 point 4464509. So, it is accurate up to five decimal places, so we may take the solution as 198 point 4464509. Of course, you can see that, this is not much different from the sample median, because the sample median was 198 point 5.

So, this method of scoring can be applied to various cases, whenever we are getting a non-linear equation, for which the solution is not in a tractable form. 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_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, $1/\theta$ is 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: 39:51)

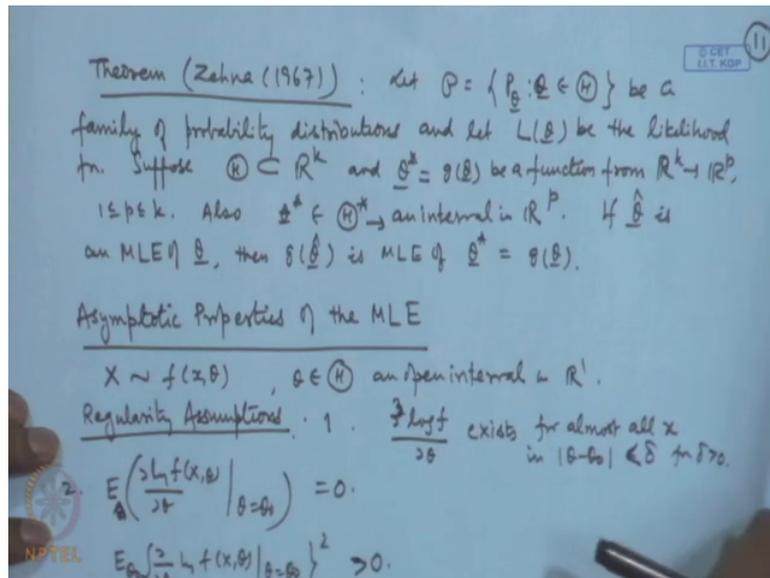


Now, this property I will be proving in two forms; invariance of MLE. Firstly, I will prove it for the one-to-one functions, and then actually I will give the general proof, which is true for any function. Let ϕ be $g(\theta)$, where g is a one-to-one function. Then $\hat{\phi}_{ML}$ is equal to g of $\hat{\theta}_{ML}$. Let us look at the proof of this, so suppose my parameter space for θ is Θ , and ϕ we wrote as Φ and this is actually the $g(\theta)$ space, that is the set of all $g(\theta)$ values as θ varies over Θ . Let $L(\theta)$ and $L^*(\phi)$ denote, the likelihood functions corresponding to θ and ϕ respectively. Essentially they are that same function, because $L(\theta)$ is obtained as the joint distribution, written at the point θ . Now in that you substitute, because $g(\theta)$ is equal to ϕ . So, if you substitute in terms of ϕ here in that function, then that will be a function of ϕ and we denote it 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, $L^*(\hat{\phi}_{ML}) \geq L^*(\phi)$ for all ϕ belonging to Φ , or $L^*(\hat{\phi}_{ML}) \geq L(g^{-1}(\phi))$, that is equal to $L(\theta)$. Now $L(\hat{\theta}_{ML})$, it is greater than or equal to $L(\theta)$, for all θ belonging to Θ . So, for $\hat{\phi}_{ML}$ equal to g of $\hat{\theta}_{ML}$, we get $L^*(\hat{\phi}_{ML}) = L^*(g(\hat{\theta}_{ML})) = L(g^{-1}(g(\hat{\theta}_{ML}))) = L(\hat{\theta}_{ML}) \geq L(\theta) = L^*(\phi)$ for all $\phi \in \Phi$.

is greater than or equal to $L(\theta)$, that is equal to $L(\phi)$ for all ϕ . So, we can say that $\hat{\phi}$ is the MLE of θ . Now naturally this result is true, we have proved for g being a one-to-one function. However, even if we have any function the same invariance property can be used. A justification for this was provided by Zehna in 1967.

(Refer Slide Time: 44:04)



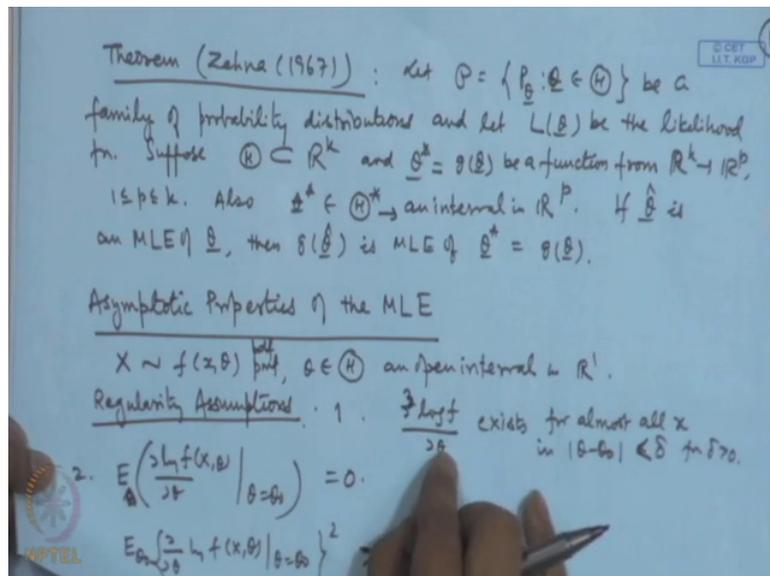
I will state the result without proof here, let $\{P_\theta : \theta \in \Theta\}$ be a family of probability distributions, and let $L(\theta)$ be the likelihood function. Suppose Θ is a subset of k dimensional space, and θ^* be a function from \mathbb{R}^k to \mathbb{R}^p , where p will be less than or equal to k . Also we assume that the range of θ^* 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 $g(\theta)$. If there are several values which lead to the same value, because now the function can be a many-to-one 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, $L(\theta)$ is $L(g^{-1}(\phi))$ which we call $L(\phi)$. Whereas in the case of many-to-one 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 skipping the proof here, for the details we can look at the paper by Zehna in 1967. I will now give some asymptotic properties of the maximum likelihood estimators. So, let me call it large sample properties are asymptotic properties of the maximum likelihood estimators. 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 the usually call Cramer–Rao or Frechet-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 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 exist around θ_0 , then if we assume this derivative existing in an interval or in a neighborhood of θ_0 then it is enough, less than δ for some δ positive.

(Refer Slide Time: 49:55)

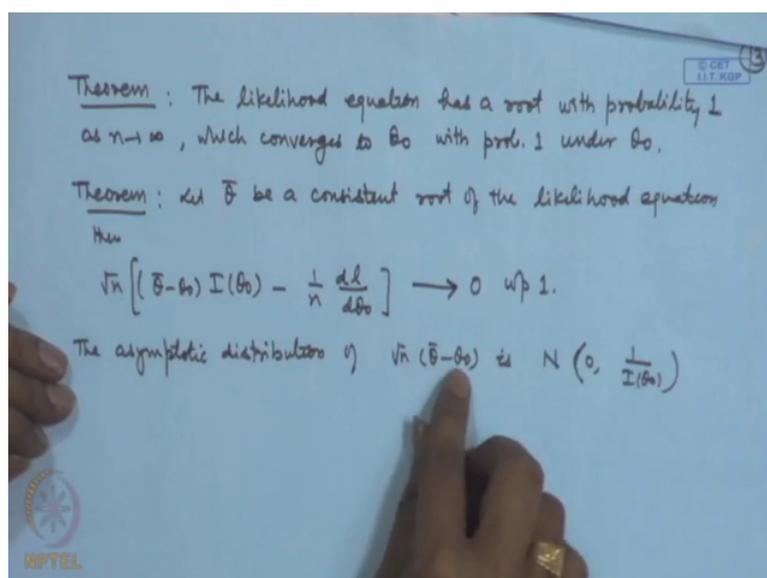


The second condition is that expectation of $\frac{\partial}{\partial \theta} \log f$ by $\frac{\partial}{\partial \theta} \theta$, at θ_0 is equal to θ_0 naught, that is equal to 0. Expectation of θ_0 naught $\frac{\partial}{\partial \theta} \log f$ by $\frac{\partial}{\partial \theta} \theta$ log of f x θ , at θ_0 is equal to θ_0 naught, whole square that is positive. The third condition is that the third

derivative of the density or the mass function is bounded in a neighborhood of θ_0 , and this function itself is having finite expectation, a bounded expectation at any point in the interval $\theta_0 - \delta$, $\theta_0 + \delta$. So, if we are considering x_1, \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 $\frac{d l}{d \theta} = \sum_{i=1}^n f'(x_i, \theta)$. Now, here prime means derivative with respect to θ , this is the likelihood equation. Let us look at the conditions once again, whatever assumptions we have made, here f can be p.d.f or p.m.f.

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. If we look at this one expectation of $\left(\frac{d \log l}{d \theta}\right)^2$. This is called 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 theorem without proof.

(Refer Slide Time: 53:26)



The first result says, that the likelihood equation has a root with probability one as n tends to infinity, and the root converges to θ_0 with probability one under θ_0 . So, this says that the likelihood equation has a root with probability one, and the root converges to θ_0 with probability one, so this is a very important result. And the second one says that the, let $\bar{\theta}_n$ be a consistent root of the likelihood equation, then $\sqrt{n}(\bar{\theta}_n - \theta_0) \xrightarrow{d} N(0, I^{-1}(\theta_0))$. The information measure we have defined, $I(\theta) = -E[\partial^2 \log L(\theta)]$. This converges to 0 with probability one. Now what does it mean, that $\sqrt{n}(\bar{\theta}_n - \theta_0) \xrightarrow{d} N(0, I^{-1}(\theta_0))$. This result is something like, that the root of the likelihood equation that is a maximum likelihood estimator is asymptotically efficient, because what term you are getting $\bar{\theta}_n - \theta_0$.

See if you take it to this side $\sqrt{n}(\bar{\theta}_n - \theta_0) \xrightarrow{d} N(0, I^{-1}(\theta_0))$. And the asymptotic distribution of $\sqrt{n}(\bar{\theta}_n - \theta_0)$ is normal $N(0, I^{-1}(\theta_0))$; that is as n becomes large, the distribution of $\sqrt{n}(\bar{\theta}_n - \theta_0)$ converges to that of a normal distribution with mean 0, and variance $I^{-1}(\theta_0)$, where $I(\theta_0)$ is the Fisher's Information Measure. The proofs of these results are not 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.

Cauchy 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 one; that means, it converges to the true value with probability one. Moreover the asymptotic distribution is normal, and it is also second order efficiency in the sense of Rao. So, that makes the use of maximum likelihood estimators a fairly important practice in the statistical theory. See this property which I have written at the end, that is $\sqrt{n}(\bar{\theta}_n - \theta_0) \xrightarrow{d} N(0, I^{-1}(\theta_0))$, or say $\sqrt{n}(\bar{\theta}_n - \theta_0)$ is asymptotically normal. This is also known as consistent asymptotic normal property or Chain 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 Asymptotically Normal (BAN) estimators, so Cramér-Rao. So, under certain conditions such estimator exists. And maximum likelihood estimators are more likely to satisfy these properties. I will be completing this discussion now, and we move over to, once again the unbiased estimators. We have defined unbiased estimators, but I had mentioned that there can be many unbiased estimators. So, which one we should choose. We will introduce that criteria in the next lecture.