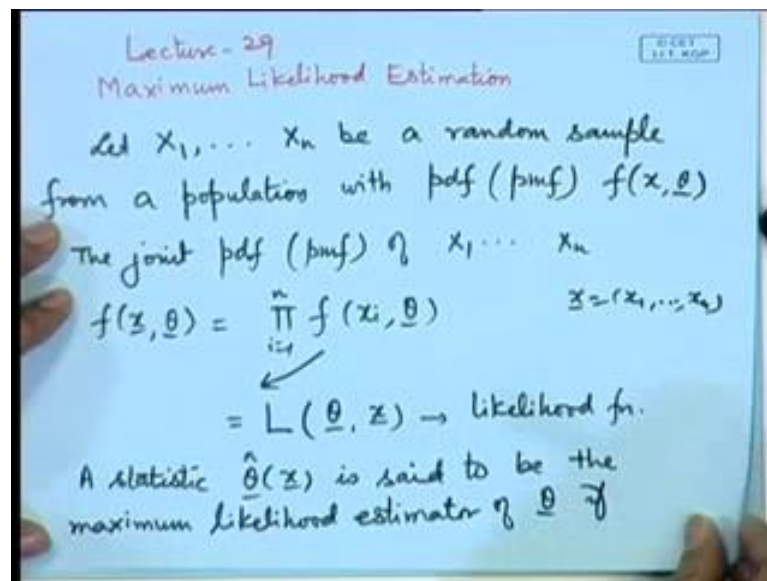


Probability and Statistics
Prof. Somesh Kumar
Department of Mathematics
Indian Institute of Technology, Kharagpur

Lecture - 57
Examples on MLE – I

In the last lecture I introduced method of maximum likelihood estimation, besides that we had also discussed the method of least square estimation and the method of moments. So, now, I continue the discussion on the method of maximum likelihood estimation which is actually the last among the three mention methods that I have given. So, let me repeat the definition of the maximum likelihood definition that how it is obtained, so firstly we define likelihood function.

(Refer Slide Time: 00:53)



So, if I have X_1, X_2, \dots, X_n let X_1, X_2, \dots, X_n be a random sample from a population with either pdf or pmf say $f(x, \theta)$. So, what we write the joint pdf or pmf of X_1, X_2, \dots, X_n . So, we write it as say $f(x, \theta) = \prod_{i=1}^n f(x_i, \theta)$ where i is equal to 1 to n ; we call this now see when X_1, X_2, \dots, X_n is observed to be some value x_1, x_2, \dots, x_n then we call this as the likelihood function of θ and we may put x here just to denote its dependence on x also this is called the likelihood function.

So, what we consider a statistic $\hat{\theta}(x)$ is said to be the maximum likelihood estimator of θ if $L(\hat{\theta}(x), x) \geq L(\theta, x)$ for all θ belonging to the parametric space.

(Refer Slide Time: 02:55)

$L(\hat{\theta}(x), x) \geq L(\theta, x) \quad \forall \theta \in \Theta$

MLE

Example 1. Let $X_1, \dots, X_n \sim \text{Ber}(1, p)$, $0 \leq p \leq 1$.

$$L(p, x) = \prod_{i=1}^n \left\{ p^{x_i} (1-p)^{1-x_i} \right\} \quad x_i = 0, 1$$

$$= p^{\sum x_i} (1-p)^{n - \sum x_i}$$

$$\ell(p) = \log L = \sum x_i \ln p + (n - \sum x_i) \ln(1-p)$$

$$\frac{d\ell}{dp} = \frac{\sum x_i}{p} - \frac{n - \sum x_i}{1-p} = 0$$

So, in short we use the terminology MLE for maximum likelihood estimator; a popular technique to obtain the optimizing value of θ is to take logarithm of L and then consider the derivative of L with respect to the parametric functions and then equate to 0 and solve those equations.

So, let me explain the method through some examples let X_1, X_2, \dots, X_n follows a Bernoulli distribution. So, here p is the unknown parameter, we want to find out the maximum likelihood estimator for p . So, we write down the likelihood function here that will be equal to p to the power x_i , $1 - p$ to the power $1 - x_i$ product i is equal to 1 to n . If you look at the inside quantity this is the probability mass function of x_i . So, here each of the x_i is can take value 0 or 1.

So, this can be simplified as p to the power $\sum x_i$, $1 - p$ to the power $n - \sum x_i$. So, we have to differentiate this with respect to p and solve it. So, what we can do we can consider \log of L , which I call as a small l of p that is equal to $\sum x_i \log p$ plus $n - \sum x_i \log(1 - p)$. So, $d l / d p$ is equal to $\sum x_i / p$ minus $n - \sum x_i / (1 - p)$ that is equal to 0 now then this gives.

(Refer Slide Time: 05:25)

$$L(p) = \prod_{i=1}^n \{ p^{x_i} (1-p)^{1-x_i} \}$$
$$= p^{\sum x_i} (1-p)^{n - \sum x_i}$$
$$= \sum x_i \ln p + (n - \sum x_i) \ln(1-p)$$
$$\frac{\sum x_i}{p} - \frac{n - \sum x_i}{1-p} = 0 \text{ gives } \hat{p} = \frac{\sum x_i}{n} = \frac{x}{n}$$

$(p \in (0, 1))$
 $x_i \in \{0, 1\}$

So, after some simplification we will get \hat{p} is equal to so in fact you can solve it let us combine the term, so it will give me $1 - p \sum x_i - n p + p \sum x_i$ divided by. So, this we can put p is belonging to the interval 0 to 1.

We can take the case p is equal to 0 and p is equal to 1 separately. So, this becomes $\sum x_i - n p$ is equal to 0; that means, \hat{p} is equal to $\sum x_i / n$, which I can call say x / n ; where x is the sigma of x_i . If you compare it with the method of moment estimator actually it is the same thing here; that means, the maximum likelihood estimator of the probability of success is actually the number of successes in n trials divided by the number of trials, that is the proportion of the number of successes which is a very logical estimator and it is coming through the method of maximum likelihood estimation.

(Refer Slide Time: 06:42)

2. X_1, \dots, X_n (i.i.d.) $\theta(\lambda)$, $\lambda > 0$

$$L(\lambda, \mathbf{x}) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}, \quad x_i = 0, 1, 2, \dots$$
$$= \frac{e^{-n\lambda} \lambda^{\sum x_i}}{\prod_{i=1}^n (x_i!)}$$

$\ell(\lambda) = \ln L(\lambda, \mathbf{x}) = -n\lambda + \sum x_i \ln \lambda - \ln(\prod x_i!)$

$\frac{d\ell}{d\lambda} = 0 \Rightarrow -n + \frac{\sum x_i}{\lambda} = 0$

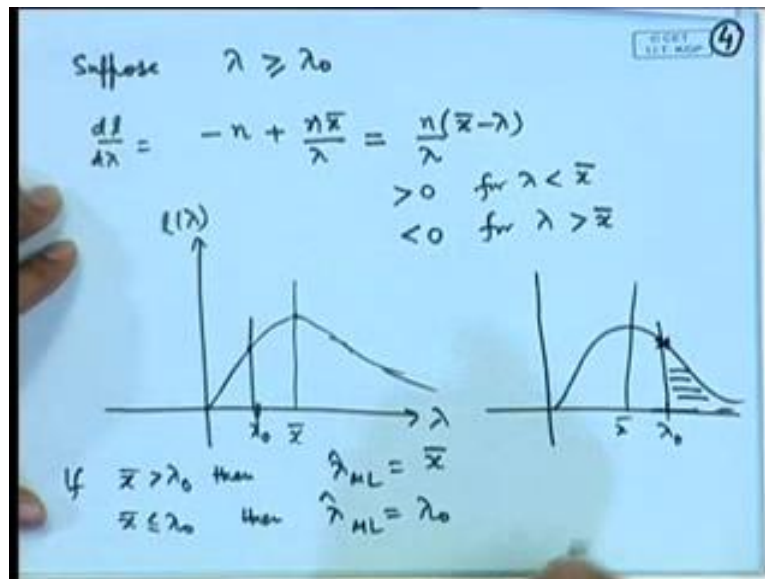
$$\Rightarrow \hat{\lambda}_{MLE} = \frac{\sum x_i}{n} = \bar{x}$$

\bar{x} is MLE for λ .

Let us take say X_1, X_2, \dots, X_n to be a random sample from Poisson distribution with parameter λ ; let us write down the likelihood function that is equal to $e^{-n\lambda} \lambda^{\sum x_i}$ divided by $\prod_{i=1}^n x_i!$. Here each of the x_i can take values 0, 1, 2 and so on and can be positive. So, this is equal to $e^{-n\lambda} \lambda^{\sum x_i}$ divided by $\prod_{i=1}^n x_i!$.

So, we can write the log likelihood function $-n\lambda + \sum x_i \ln \lambda - \ln(\prod x_i!)$. So, the likelihood equation $\frac{d\ell}{d\lambda} = 0$, that is $-n + \frac{\sum x_i}{\lambda} = 0$. So, the solution of this gives $\hat{\lambda}_{MLE} = \frac{\sum x_i}{n} = \bar{x}$. So, \bar{x} is the maximum likelihood estimator for λ . Note here that in these 2 cases the maximum likelihood estimator and the method of moments estimator are same; however, as we will see it is not a rule, in many cases they will not be the same let me take an example of that kind.

(Refer Slide Time: 08:46)



Another thing that we can observe here the maximization is done over the parameter space. So, in case there is a modification for example in the previous exercise, suppose we know that lambda is say greater than or equal to lambda naught; that means, we know that the rate of arrival in a Poisson process is bigger than a prescribed quantity. In that case if we look at x bar this is actually the maximization over the full parameter space and we may have a situation where x bar is actually less than lambda naught, in that case this does not satisfy the property that likelihood function is maximized over the parameter space.

So, what we do? Then consider the behavior of the likelihood function, what we are getting is $d\ell/d\lambda$ is equal to $-n + n\bar{x}/\lambda$. So, we write it in $n(\bar{x} - \lambda)/\lambda$ that we can write as $n(\bar{x} - \lambda)/\lambda$ that is $n\bar{x}/\lambda - n$. So, you can see here that if lambda is less than x bar then this is positive; it is less than 0 for lambda greater than x bar. So, we look at the behavior of the function it is increasing up to x bar and then decreasing after x bar.

So, let us plot it as a function of lambda the likelihood function; that means, on this side I have log likelihood, it is increasing and then it is decreasing this is the point x bar. Now suppose lambda naught is here, if lambda naught is here then you consider the maximum value x bar is satisfying this condition and therefore, x bar remains the MLE; that is if x bar is greater than lambda naught then lambda hat ML is x bar. However, we may have a situation where this is x bar and lambda naught is say here, in that case you see the parameter space is this; that means, the maximum value that is occurring is at actually

lambda naught; that means, if \bar{x} is less than or equal to lambda naught then lambda hat ML is equal to lambda naught.

(Refer Slide Time: 11:38)

So the MLE of λ is modified as

$$\hat{\lambda}_{ML} = \bar{x} \text{ if } \bar{x} > \lambda_0$$

$$= \lambda_0 \text{ if } \bar{x} \leq \lambda_0$$

X_1, \dots, X_n i.i.d. $N(\mu, \sigma^2)$

$$L(\mu, \sigma^2, \mathbf{z}) = \prod_{i=1}^n \left[\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} \right]$$

$$= \frac{1}{\sigma^n (2\pi)^{n/2}} e^{-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2}$$

$$\ell(\mu, \sigma^2) = \log L = -n \log \sigma - \frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

So, our estimator is then modified, so the maximum likelihood estimator of lambda is modified as let me write it as lambda hat ML this is equal to \bar{x} ; if \bar{x} is greater than lambda naught it is equal to lambda naught, if \bar{x} is less than or equal to lambda naught. Now this brings into focus another important property or you can say another important aspect of the maximum likelihood estimator, which would have been missed by the method of moments; because in the method of moments the restriction on the parameter space does not play any role there we simply look at the moment which is not affected by the restriction on the parameter space and so the method of moment estimator remains as \bar{x} . Whereas here you see the effect of the restriction on the parameter space is getting reflected on the maximum likelihood estimator which is actually a reasonable thing, because if \bar{x} is less than or equal to lambda naught which is going outside the parameter space, we should not take \bar{x} as an estimate for lambda.

Let us take another popular example X_1, X_2, \dots, X_n follows normal μ, σ^2 . If we consider the likelihood function here then it is a function of 2 parameters μ and σ^2 . So, $L(\mu, \sigma^2, \bar{x}, \mathbf{X})$ that is equal to product of $\frac{1}{\sigma \sqrt{2\pi}}$ e to the power minus $\frac{1}{2\sigma^2} (x_i - \mu)^2$. Now this term we

will simplify this can be written as $\frac{1}{\sigma^2} \sum (x_i - \mu)$ to the power n , $e^{-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2}$.

So, the log likelihood function minus $n \log \sigma$, minus $n \log \frac{1}{\sqrt{2\pi}}$, minus $-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2$. Now note here that I have written it as minus $n \log \sigma$; now one may ask that if parameter is sigma square then whether we should write it as sigma. Now answer is that both are because maximum likelihood estimation is invariant under the transformation of the parameters, suppose I obtain the MLE of sigma square in place of sigma and then I want for sigma, then I should simply take the square root of that. So, I may write it like this also it does not make any difference.

(Refer Slide Time: 14:50)

The likelihood equations are

$$\frac{\partial l}{\partial \mu} = 0 \Rightarrow \frac{1}{\sigma^2} \sum (x_i - \mu) = 0 \Rightarrow \hat{\mu}_{ML} = \bar{x}$$

$$\frac{\partial l}{\partial \sigma^2} = 0 \Rightarrow -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (x_i - \mu)^2 = 0$$

$$\Rightarrow \sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$$

$$\Rightarrow \hat{\sigma}_{ML}^2 = \frac{1}{n} \sum (x_i - \bar{x})^2 = \left(\frac{n-1}{n}\right) S^2$$

we know that $\mu \geq 0$

$$\frac{\partial l}{\partial \mu} = \frac{n(\bar{x} - \mu)}{\sigma^2} > 0 \text{ for } \mu < \bar{x}$$

$$< 0 \text{ for } \mu > \bar{x}$$

So, the likelihood equations in this case will be $\frac{\partial l}{\partial \mu} = 0$; that means, $\frac{1}{\sigma^2} \sum (x_i - \mu) = 0$, which will give $\hat{\mu}$ is equal to \bar{x} and $\frac{\partial l}{\partial \sigma^2} = 0$ that gives $-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (x_i - \mu)^2 = 0$, this gives σ^2 is equal to $\frac{1}{n} \sum (x_i - \mu)^2$; that means, a equation for sigma square involves μ . So, since we have already got the solution for μ we can substitute; so $\hat{\sigma}^2$ that is a maximum likelihood estimator becomes that is $\frac{n-1}{n} S^2$.

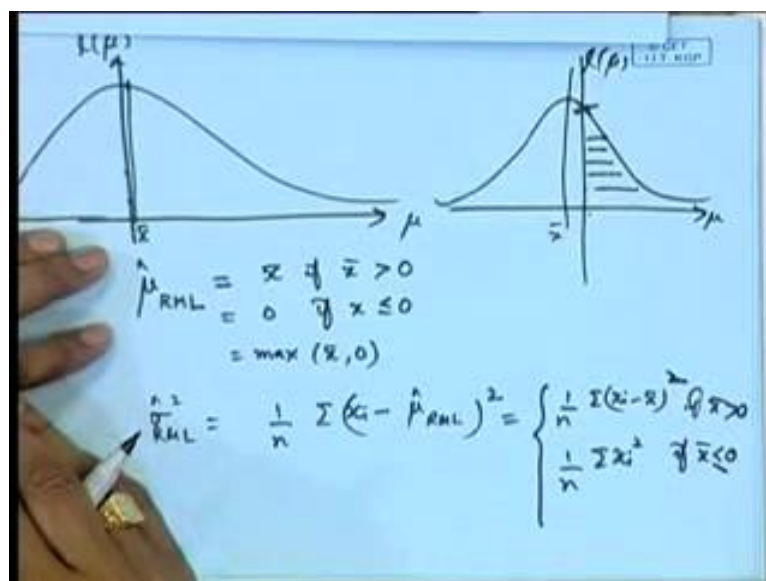
You can note here that it is similar to the method of moment's estimator in this case. So, in many situations the method of moment's estimators and the maximum likelihood

estimates concede, but there are other cases where they may not concede. We already have seen the example that if the parameter space gets modified then the maximum likelihood estimator gets modified. For example, in this particular situation suppose I consider suppose we know that μ is greater than or equal to 0 now this type of situation occurs when through some experience we already know that the mean is actually non negative although, the variable may be normally distributed, but because of the way the experiment has been framed or any other reason the parameter space is restricted; that means, we know that the mean is greater than or equal to 0.

Now, you see here we have the maximum likelihood estimator for μ as \bar{x} and for sigma square it is n minus 1 by n s square. Now if we see that \bar{x} by observation gives us a negative value then it will become unreasonable estimators for μ , which is taken to be greater than or equal to 0. So, we can analyze it in a proper way by looking at the behavior of the log likelihood function which we have actually maximized. So, if we look at with respect to μ we have $\frac{\partial l}{\partial \mu}$ as equal to. So, $\frac{\partial l}{\partial \mu}$ that function we got it as n times \bar{x} minus μ by sigma square.

So, you can see that it is greater than 0 for μ less than \bar{x} and it is less than 0 for μ greater than \bar{x} ; that means, the function is increasing up to \bar{x} and decreasing beyond \bar{x} bar.

(Refer Slide Time: 18:06)



That means the shape of the likelihood function as a function of μ is something like; so this could be the maximizing choice. Now if \bar{x} is bigger than 0 then we may take \bar{x} , but we may have a situation where \bar{x} may be less than 0. So, if the parameter space is this then the maximum is occurring at 0 itself.

So, the modified or you can say the restricted maximum likelihood estimator becomes \bar{x} if \bar{x} is greater than 0 it is equal to 0, if \bar{x} is less than or equal to 0; which we can actually call as maximum of \bar{x} and 0. Now is there any effect on the estimator for sigma square, the answer is yes because the maximum likelihood estimator for sigma square was obtained by substituting the estimator for μ in this second likelihood equation, so we will get sigma hat square ML as $\frac{1}{n} \sum (x_i - \hat{\mu}_{ML})^2$.

So, this we can write as when \bar{x} is positive then this is the old 1 that is $\frac{1}{n} \sum (x_i - \bar{x})^2$ if \bar{x} is positive; and it is equal to $\frac{1}{n} \sum x_i^2$ if \bar{x} is negative or less than or equal to 0. We will show later on that when the parameter space is restricted, this restricted maximum likelihood estimator has a greater performance as compared to the usual maximum likelihood estimator that we obtained before; we will define the criteria of better a little later.

(Refer Slide Time: 20:35)

4. $X_1, \dots, X_n \sim U[0, \theta], \theta > 0$

$$L(\theta, \mathbf{x}) = \begin{cases} \frac{1}{\theta^n} & , 0 \leq x_i \leq \theta, i=1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$$\ell = \log L = -n \ln \theta \quad \cdot \quad \frac{d\ell}{d\theta} = -\frac{n}{\theta} = 0 \text{ absurd}$$

$x_{(n)} = \max\{x_1, \dots, x_n\}$

Then L is maximized when θ is minimized
i.e. $\hat{\theta}_{ML} = x_{(n)}$.

$\mu_1^* = E(X) = \frac{\theta}{2} \quad \hat{\theta}_{MME} = 2\bar{x}$.

Let us take another example where this process of argument does not seem to work; let us consider say a random sample from uniform distribution on the interval say 0 to theta

where θ is a positive real number. So, here the likelihood function is simply $\frac{1}{\theta^n}$ for $0 < x_i \leq \theta$, for $i = 1, 2, \dots, n$ it is 0 otherwise. Now you see if you follow the previous procedure of taking logarithm and differentiating what I will get? If I take \log of L then that will give me $-\ln L = -n \ln \theta$ and if I differentiate let me call it is small l so $\frac{d}{d\theta} l$ that will give me $-\frac{n}{\theta}$. So, if I put this equal to 0 this is actually observed. So, why this is happening?

This is happening because we are not taking care of the region, when we are differentiating and putting equal to 0 basically we are trying to find out the minimum and maximum over a range of parameter. Here the range of the parameter value is dependent upon x_i so that is not been taken care by this kind of process. So, let us look at a direct argument, our aim is to maximize the likelihood function $\frac{1}{\theta^n}$ with respect to θ ; since θ is in the denominator it corresponds to the minimization with respect to θ , what is a minimum value of θ ?

The minimum value of θ will be actually since $E \theta$ is greater than each of the x_i the minimum value that θ can take will be the maximum of x_1, x_2, \dots, x_n . So, if I call x_n as the maximum of x_1, x_2, \dots, x_n that is largest order statistic then L is maximized when θ is minimized that is $\hat{\theta}_{ML} = X_n$. So, this is the maximum likelihood estimator for θ . Suppose I want to find out the method of moments estimator here, what is that? Consider the mean of this distribution that is the first moment that is $\theta/2$.

So, $\hat{\theta}$ method of moments estimator that will be twice and in place of μ_1' will put \bar{X} . So, you can see the situation here the method of moment estimator and the maximum likelihood estimator are totally different. In fact, they are here it is the maximum of the observation and here it is average 2 times the average of the observations; further this example illustrates that the method of taking logarithm and differentiating does not always work.

(Refer Slide Time: 24:00)

5. $X_1, \dots, X_n \sim \text{Exp}(\mu, \sigma)$ (9)

$$f(x, \mu, \sigma) = \frac{1}{\sigma} e^{-\frac{(x-\mu)}{\sigma}}, \quad x \geq \mu, \quad \sigma > 0, \mu \in \mathbb{R}$$

The likelihood function is

$$L(\mu, \sigma, \mathbf{x}) = \frac{1}{\sigma^n} e^{-\frac{1}{\sigma} \sum_{i=1}^n (x_i - \mu)}, \quad x_i \geq \mu, \quad i=1, \dots, n$$

$$= \frac{1}{\sigma^n} e^{-\frac{n}{\sigma}(\bar{x} - \mu)} = \frac{1}{\sigma^n} e^{\frac{n}{\sigma}(\mu - \bar{x})}, \quad x_i \geq \mu, \quad i=1, \dots, n$$

$$\ell = \ln L = -n \ln \sigma + \frac{n}{\sigma}(\mu - \bar{x})$$

This is maximized with respect to μ , when

$$\hat{\mu}_{ML} = X_{(1)} = \min\{X_1, \dots, X_n\}$$

Let us take another examples say X_1, X_2, X_n follow exponential distribution let me consider say 2 parameter exponential distribution μ, σ ; that means, the density function is $\frac{1}{\sigma} e^{-\frac{(x-\mu)}{\sigma}}$ that is the density of; here x is greater than or equal to μ σ is positive and μ is any real number. The likelihood function will be $\frac{1}{\sigma^n} e^{-\frac{1}{\sigma} \sum_{i=1}^n (x_i - \mu)}$ each x_i is greater than or equal to μ which we can write as $\frac{1}{\sigma^n} e^{-\frac{n}{\sigma}(\bar{x} - \mu)}$ or $\frac{1}{\sigma^n} e^{\frac{n}{\sigma}(\mu - \bar{x})}$ to the power n , e to the power n by σ μ minus \bar{x} .

Now, if we want to maximize this function with respect to μ , the likelihood equation and then derivative will not give a result like in the uniform case because the log likelihood function $\ln L = -n \ln \sigma + \frac{n}{\sigma}(\mu - \bar{x})$; you can see here if I differentiate with respect to μ and put equal to 0, I get an absurd result so; however, with respect to σ we can do that, but for finding out the maximum likelihood estimator with respect to μ we can use a direct argument, this is an increasing function of μ . So, the maximum value will be attained when μ attains its maximum value.

Since μ is always less than or equal to x_i , the maximization will occur when μ is the minimum of the x_i . So, this is maximized with respect to μ when $\hat{\mu}_{ML}$ is equal to $X_{(1)}$ that is the minimum of the observations. So, we can substitute this here and get the estimator for σ also.

(Refer Slide Time: 27:00)

Consider the derivative of this with respect to sigma that will give minus n by sigma, minus n by sigma square, mu minus x bar this is equal to 0. So, this gives sigma hat is equal to x bar minus mu hat that is equal to x bar minus x 1.

So, sigma hat ML is equal to X bar minus X 1, let us see what is the method of moments estimator in this case; to obtain the method of moments estimator we have to consider the moments of this distribution, now here it is a 2 parameter distribution I will have to find out first 2 moments. So, mu 1 prime is equal to mu plus sigma, mu 2 prime now for that we can consider certain transformation expectation of X minus mu square that is equal to twice sigma square, so expectation of X square minus 2 mu X plus mu square is equal to twice sigma square.

So, expectation of X square is equal to 2 sigma square, minus mu square plus twice mu expectation of X that is mu plus sigma that is equal to twice sigma square plus mu square plus 2 mu sigma. So, the second moment is twice sigma square plus mu square plus 2 mu sigma. Now if we consider mu 1 prime square minus mu 2 prime rather I consider mu 2 prime minus mu 1 prime square then that gives me sigma square and therefore, mu becomes mu 1 prime minus square root of mu 2 prime minus mu 1 prime square.

So, the method of moments estimators for mu and sigma square will be obtained as sigma hat square MME that is equal to 1 by n sigma x i minus x bar whole square and mu hat MME will be equal to X bar minus a square root 1 by n sigma x i minus x bar

whole square. Note here that the maximum likelihood estimators and the method of moment estimators are quite different from each other and again therefore, the question arises that which of this is better.