

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

Lecture – 15
Lower Bounds for Variance- I

So, now we will take up another topic that is for the Lower Bounds for the Variance. Now, what is this concept? Earlier, we have seen that unbiasedness is a desirable property or desirable criteria to use an estimator. However, we have also seen the example that in a given problem, there can be several unbiased estimators. Now, if there are several unbiased estimators which one to choose, then we can decide some additional criteria such as variance. The one which has smaller variance will be considered to be more stable in some sense.

Now therefore, we need to have an estimate of that what could be the variance or what could be the minimum variance. So this gives the idea or you can say this led to the development of methods for finding out lower bounds for the variance of an unbiased estimator.

(Refer Slide Time: 01:26)

Lecture-8. Lower Bounds for Variance-1 SECRET 1

In this section we will discuss various methods for determining the lower bounds on the variance of unbiased estimators.

Wolfowitz's Regularity Conditions.

Let X_1, \dots, X_n be a random sample from a distribution having pdf (pmf) $f(x, \theta)$ w.r.t measure μ . An estimator $\delta(X)$ is to be considered for θ .

(i) θ lies in an open interval H of the real line

(ii) $\frac{\partial f(x, \theta)}{\partial \theta}$ exists $\forall \theta \in H, \forall x$

(iii) $\int \left(\frac{\partial}{\partial \theta} \log f(x, \theta) \right)^2 d\mu(x) > 0 \forall \theta \in H$

under the integral sign for any θ such that the above integral exists.

$E \left[\frac{\partial}{\partial \theta} \log f(X, \theta) \right]^2 > 0 \forall \theta \in H$

MPTEL

So, in this section, we will discuss various methods for determining the lower bounds on the variance of unbiased estimators. As we have seen in the case of maximum likelihood estimation in the last results that I gave that variance asymptotic variance of the

maximum likelihood estimator was 1 by the information. Now, this is asymptotic variance. So, if the maximum likelihood estimator is the best in some sense then its variance will not be below $1/I(\theta_0)$, that means, the Fisher information measure.

The question comes that whether similar result we can give for finite samples. Now, this is the precisely the question that was posed to Indian statisticians C.R. Rao in his class in 1943 at Indian Statistical Institute, and he started working out for finite samples. And it led to the famous lower bound by Rao. However, at the same time the result was also proved by Fisher in 1943, by Cramer in 1946, therefore, it is now popularly called Fisher, Rao, Cramer inequality.

Now, once again in order to prove this, we need certain regularity conditions they are known by the name Wolfowitz's wits regularity conditions named after the statistician Jacob Wolfowitz. So, as before we have a random sample be a random sample from a distribution having say PDF. And of course, it could be pmf $f(x, \theta)$ with respect to say measure μ . So, we assume the usual conditions for the existence of a density function or the mass function etcetera.

Now, an estimator $\delta(x)$ is to be considered for the parameter θ . We make the assumptions that θ lies in an open interval of the real line, the derivative of the density or the mass function exist, and of course for all x are for almost all x . The integral I used a more general notation, because if it is discrete this will be replaced by summation, I have written here $d\mu$, so that takes care of both the cases. So, this is a n fold integral or summation. This can be differentiated under the integral sign for any δ such that this is an integral function that means this integral exists that means, for any integral function its expectation should be or its integral should be differentiable. So, that the above integral exists. This is positive for all θ . Once again this is related to the Fisher information measure.

(Refer Slide Time: 06:46)

Frechet-Rao-Cramer Inequality : Under assumptions (i) - (iv) of (1943) (1945) (1946) SECRET I.I.T. KGP 2

$E_{\theta} \delta(X) = \theta + b(\theta)$, then

$$\text{Var}_{\theta}(\delta) \geq \frac{\{1 + b'(\theta)\}^2}{n E \left[\frac{\partial}{\partial \theta} \log f(X, \theta) \right]^2} \quad \dots (1)$$

Proof: $E_{\theta} \delta(X) = \theta + b(\theta) \quad \forall \theta \in \Theta$

$$\Rightarrow \int \delta(x) \prod_{i=1}^n f(x_i, \theta) d\mu(x) = \theta + b(\theta) \quad \forall \theta \in \Theta$$

Differentiating under the integral sign w.r.t θ , we get

$$\int \delta(x) \left\{ \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i, \theta) \right\} \prod_{i=1}^n f(x_i, \theta) d\mu(x) = 1 + b'(\theta)$$

$\Rightarrow E[\delta(X) S(X, \theta)] = 1 + b'(\theta) \quad \forall \theta \in \Theta$

Under these conditions we have the following inequality I will call it fresh Fisher Rao Cramer inequality. Rao Cramer inequality, because Fisher's is paper appeared in 1943, Rao's paper appeared in 1945, Cramer's paper appeared in 1946. So, they all seem to have done it independently under assumptions 1 to 4. If expectation of delta x is equal to theta plus b theta, then variance of delta is greater than or equal to 1 plus b prime theta whole square divided by n times expectation del by del theta log of f x theta whole square.

Firstly, let us look at the proof of this. So, what we are doing is that for an estimator delta, we are providing the lower bound for the variance. This right hand side you can see it is not dependent upon the twice of the estimator that we have chosen, that means, any estimator of any unbiased estimator of theta plus b theta will have the minimum variance which will be greater than or equal to this, because this is the lower bound. So, it may be attend or it may not be attend. Let us look at the proof of this result first of all. So, expectation of delta x is equal to theta plus b theta. Now, this is of course, true these statements are true for all theta.

Now, we are assuming that we can differentiate under the integral sign. So, this is delta product f of x i, theta d mu x. Now, this denotes d mu x 1, d mu x 2, d mu x n this is equal to theta plus b theta for all theta differentiating under the integral sign. Let me again emphasize that this integral is a generalized lebesgue integral. That means if we are

dealing with the discrete distributions, then this will be replaced by the summation. So, this is $\delta x_1, x_2, \dots, x_n$ product of f of x_i , $\theta d\mu \times 1 d\mu \times 2 d\mu \times \dots \times n$. So, this is the n fold integral.

So, if you differentiate with respect to θ , we will get δx . Now, derivative of the product that you can easily write as $\sum \frac{\partial}{\partial \theta} \log f(x_i, \theta)$ multiplied by product $f(x_i, \theta)$ that is equal to $1 + b'(\theta)$. Now, we use some notation this term I call say $S(x, \theta)$. Then I am getting δx into $S(x, \theta)$ into the joint distribution of x_1, x_2, \dots, x_n and $d\mu \times 1 d\mu \times 2 d\mu \times \dots$. So, this we can write as expectation of δx into $S(x, \theta)$ it is equal to $1 + b'(\theta)$.

Now, what we can see that this term, if we look at this we have made the assumption here that for any function δ for which this integral exists this can be differentiated. So, if we look at this particular term that is $S(x, \theta)$ then expectation of $S(x, \theta)$ can also be differentiated under the integral sign. If we look at that then this is going to be 0. Let us see this. Let me give this 2 here.

(Refer Slide Time: 11:43)

Now we have $\int \prod_{i=1}^n f(x_i, \theta) d\mu(x) = 1 \quad \forall \theta \in \Theta$ (3)

So once again, differentiating (3) under the integral sign, we get

$$\int \left\{ \sum_{i=1}^n \frac{\frac{\partial}{\partial \theta} f(x_i, \theta)}{f(x_i, \theta)} \right\} \prod_{i=1}^n f(x_i, \theta) d\mu(x) = 0 \quad \forall \theta \in \Theta$$

$$\Rightarrow \int \left\{ \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i, \theta) \right\} \prod_{i=1}^n f(x_i, \theta) d\mu(x) = 0 \quad \forall \theta \in \Theta$$

$$\Rightarrow \int \delta(x, \theta) \prod_{i=1}^n f(x_i, \theta) d\mu(x) = 0 \quad \forall \theta \in \Theta$$

$$\Rightarrow E_{\theta} \delta(x, \theta) = 0 \quad \forall \theta \in \Theta$$

Using this in (2), we can write

$$\text{Cov}(\delta(x), S(x, \theta)) = 1 + b'(\theta)$$

Now, we have the integral of the distribution of x_1, x_2, \dots, x_n equal to 1 by the property of the distribution that the integral or the summation should be equal to 1 over the whole range. So, once again if we differentiate let me call it relation 3 under the integral sign, we get $\sum \frac{\partial}{\partial \theta} \log f(x_i, \theta)$ into product of $f(x_i, \theta)$. See, if you differentiate one particular term, then other will be there. So, we can keep that also and

then divide by that. So, this becomes $\sum \frac{\partial}{\partial \theta} f(x_i, \theta)$, θ by $f(x_i, \theta)$, θ product $f(x_i, \theta) d\mu(x)$ is equal to 0.

Now, this term I can write as $\frac{\partial}{\partial \theta} \log$ of $f(x_i, \theta)$. Now, compare this, here we defined $S(x, \theta)$ to be $\sum \frac{\partial}{\partial \theta} \log$ of $f(x_i, \theta)$ and this is the term. So, what we have got here we have got integral of $S(x, \theta)$, $S(x, \theta)$ product $f(x_i, \theta) d\mu(x)$ is equal to 0, that means, expectation of $S(x, \theta)$ is 0. If expectation of a random variable is 0, then expectation of that random variable equal to the covariance term. So, we can say that using this in 2, we can write that covariance between $\delta(x)$ and $S(x, \theta)$ is equal to $1 + b'(\theta)$. Now, this relation we square it.

(Refer Slide Time: 15:03)

Squaring the above relation, we get

$$\{1 + b'(\theta)\}^2 = \text{Cov}^2(\delta(x), S(x, \theta))$$

$$\leq \text{Var}(\delta(x)) \text{Var}(S(x, \theta)) \quad \dots (4)$$

(Using Cauchy-Schwarz Inequality)

Further, $\text{Var}_\theta S(x, \theta) = \text{Var} \left[\sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i, \theta) \right]$

$$= n \text{Var} \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \right]$$

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

$$= I_x(\theta)$$

So (4) gives, $\text{Var} \delta(x) \geq \frac{\{1 + b'(\theta)\}^2}{n E \left\{ \frac{\partial}{\partial \theta} \log f(x, \theta) \right\}^2} = \frac{\{1 + b'(\theta)\}^2}{I(\theta)}$

Squaring the above relation we get $1 + b'(\theta)$ square is equal to covariance square $\delta(x)$ $S(x, \theta)$. Now, covariance square this is less than or equal to the variance of $\delta(x)$ into variance of $S(x, \theta)$ if we use Cauchy-Schwarz inequality. So, this is less than or equal to variance of $\delta(x)$ into variance of $S(x, \theta)$, this is true in general let me say it here using Cauchy-Schwarz inequality. Now, once again since expectation of $S(x, \theta)$ is 0, variance is nothing but expectation of $S(x, \theta)$ square or we can also say that variance of $S(x, \theta)$. Now, that is equal to variance of $\sum \frac{\partial}{\partial \theta} \log$ of $f(x_i, \theta)$, θ . Now, this is variance of a sum.

Now each term in the sum involves each x_i ; x_i 's are independent and identically distributed random variables. So, this becomes nothing but the n times we can say

variance of δ by δ θ \log of say $f \times 1 \theta$. Since expectation of δ by δ θ $\log f \times \theta$ is 0, this is nothing but expectation of δ by δ θ $\log f \times \theta$ square. So, this is equal to n times expectation by δ by δ θ \log of $f \times \theta$ square.

So, if we are using the notation $I \theta$ for this term, then this is nothing but the Fisher's information in the sample. We can say Fisher's information contained in the full sample. So, this we can then write 4. Here we are having variance $\delta \times$ greater than or equal to 1 plus v prime θ whole square divided by this, and that term is this variance of $\delta \times$ greater than or equal to 1 plus b prime θ . See this will be whole square here divided by n times expectation δ by δ θ \log of $f \times 1 \theta$ whole square which we can also write as 1 plus b prime θ square by $I \theta$ in the sample. This means the random sample is x_1, x_2, \dots, x_n . So, this is exactly the statement of the Cauchy-Schwarz of the Fisher, Rao, Cramer inequality.

Now, we can look at the various ramifications of this. First of all in the assumption we have taken the δ estimator to have expectation θ plus $b \theta$. Suppose, our parameter of interest is θ and δ is an unbiased estimator then $b \theta$ will be 0. If $b \theta$ is here then this term will vanish. So, the lower bound will come as simply 1 by the information or 1 by n times expectation δ by δ θ \log of $f \times \theta$.

(Refer Slide Time: 19:35)

Cor. 4 $\delta(X)$ is unbiased for θ , then

$$\text{Var}_\theta(\delta(X)) \geq \frac{1}{n \left[\frac{\partial}{\partial \theta} \log f(X, \theta) \right]^2} = \frac{1}{I_\theta(\theta)} \rightarrow \text{Fisher's Information in } (x_1, \dots, x_n) \text{ about } \theta.$$

Remarks: 1. The equality in FRC inequality is achieved iff $\delta(X)$ & $S(X, \theta)$ are linearly related with prob. 1, i.e. \exists functions $\alpha(\theta)$ & $\beta(\theta) \Rightarrow$

$$\delta(X) + \alpha(\theta) S(X, \theta) = \beta(\theta) \text{ with prob. 1.}$$

2. Under the regularity conditions

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

So, we have the following case. As a corollary I write, if $\delta \times$ is unbiased for θ , then variance of $\delta \times$ is greater than or equal to 1 by n times expectation δ by δ

theta log of f x 1 theta whole square that is 1 by I theta, this term as I defined Fisher's information in x 1, x 2, x n about theta.

Another point that let us see the Rao Cramer inequality that we have proved the proof used Cauchy-Schwarz inequality. Now, Cauchy-Schwarz inequality has a condition for the equality also when is that true. Equality is true when delta and s are that means, they are linearly related you can say that S is a linear function of delta or delta is a linear function of S. Since, here the random variables are involved we have to say that there are linear functions with probability 1.

So, we can say as a remark the equality in FRC inequality is achieved if and only if delta x and S x, theta are linearly related with probability 1, that is their exist functions say alpha theta and say beta theta such that we can say delta x plus alpha theta S x, theta is equal to say beta theta with probability 1. Now, another point I have been using that expectation of del by del theta log f x theta square. And earlier I wrote this also as minus expectation del 2 by del theta square log f x theta. Now, that is true provided the regularity conditions are satisfied.

So, let me prove that also here. Under the regularity conditions, under the regularity conditions expectation of del by del theta log of f x theta square is equal to minus expectation del 2 by del theta 2 log of f x theta.

(Refer Slide Time: 23:23)

$$\frac{\partial}{\partial \theta} \log f(x, \theta) = \frac{\partial}{\partial \theta} \left(\frac{f'(x, \theta)}{f(x, \theta)} \right) = \frac{f''(x, \theta) f(x, \theta) - (f'(x, \theta))^2}{f^2(x, \theta)}$$

$$E \left[\frac{\partial}{\partial \theta} \log f(x, \theta) \right]^2 = \int \underbrace{f''(x, \theta) f(x, \theta)}_{> 0} dP(x) - \int \left\{ \frac{f'(x, \theta)}{f(x, \theta)} \right\}^2 f(x, \theta) dP(x)$$

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

Examples: 1. $X \sim \text{Bin}(n, p)$, n is known, $0 \leq p \leq 1$.
 We estimate p here.
 $E \left(\frac{X}{n} \right) = p$. So $\frac{X}{n}$ is unbiased for p . $\text{Var} \left(\frac{X}{n} \right) = \frac{p(1-p)}{n}$.
 $f(x, p) = \binom{n}{x} p^x (1-p)^{n-x}$. $\log f(x, p) = \log \binom{n}{x} + x \log p + (n-x) \log(1-p)$
 $\frac{\partial \log f(x, p)}{\partial p} = \frac{x}{p} - \frac{n-x}{1-p} = \frac{x - np}{p(1-p)}$

So, let us look at the proof of this. Expectation of see we have to consider the second derivative here. So, let us write this $\frac{d^2}{d\theta^2} \log f(x|\theta)$ that is equal to $\frac{d}{d\theta}$ of first derivative. Now, the first derivative is nothing but f' by f . So, if you differentiate this, you will get second derivative here, multiplied by f minus derivative of this and this. So that becomes square divided by $f(x|\theta)^2$. So, if we consider expectation of this that is equal to integral of $f''(x|\theta) f(x|\theta) dx$. So, this will be cancelled out because when we multiply by $f(x|\theta)$. And $f(x|\theta)^2$ that will cancel out minus second term will become $f'(x|\theta) f(x|\theta)$.

Now, this term is 0, because of the assumption because integral $f(x|\theta) dx$ is equal to one. So, you differentiate under the integral sign. So, this becomes 0. So, this is nothing, but minus expectation of $\frac{d}{d\theta} \log f(x|\theta)$ by $\frac{d}{d\theta}$ whole square. So, these are two alternative ways of evaluating this Fisher's information measure. Now, let me give examples of the situations where the lower bound is attained, and also the examples where the lower bound is not attained. Certainly, whenever the lower bound will be attained, the unbiased estimator will become minimum variance unbiased estimator because it is attaining the lower bound.

So, there cannot be an either and by the estimator which will have the variance is smaller than this form. So, this is one nice way of proving that a given estimator is minimum variance unbiased estimator. However, in the case when it is not attained, then it is difficult to prove the minimum variance unbiased estimator using this approach for that we will take up another case or another approach here.

So, let me start with the some of the standard distributions let us consider say binomial distribution with parameters n and p , where n is known. So, the parameter is actually, and p takes any value between 0 and 1. So, we have to consider the estimation of p here. Now, easily you can see that $\frac{x}{n}$ is an unbiased estimator of p $\frac{x}{n}$ is unbiased for p . And also let us look at what is variance of $\frac{x}{n}$ variance of this is simply $p(1-p)/n$.

Now, let us look at the lower bound here if it is unbiased then the lower bound is simply equal to $1/p(1-p)$ by the information measure. So, here we can calculate this density function is $n \binom{n-1}{x} p^x (1-p)^{n-x}$. So, we take log of this that

is equal to \log of $m c x$ plus $x \log p$ plus n minus $x \log 1$ minus p . So, derivative of this with respect to p will give x by p minus n minus x by 1 minus p which we can write as x minus $n p$ divided by p into 1 minus p .

So, in order to apply the lower bound, we calculate the information. And the information term is equal to n times expectation $\frac{\partial \log f}{\partial \theta} \frac{\partial \log f}{\partial \theta}$. Since, in this case, we have only one observation, so n will not be there we simply calculate this. So, we have already evaluated the derivative $\frac{\partial \log f}{\partial p}$ by $\frac{\partial \log f}{\partial p}$. Now, we square it and then take the expectation, so that gives us expectation.

(Refer Slide Time: 28:36)

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

So the FRC lower bound for the variance of an unbiased estimator of p is $\frac{p(1-p)}{n}$ which equals $V\left(\frac{x}{n}\right)$ here.

So $\frac{x}{n}$ is UMVUE of p .
(uniformly minimum variance unbiased estimator).

2. Let $X_1, \dots, X_n \sim P(\lambda)$, $\lambda > 0$.
We want to estimate λ .

$f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$, $x = 0, 1, \dots, \infty$

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

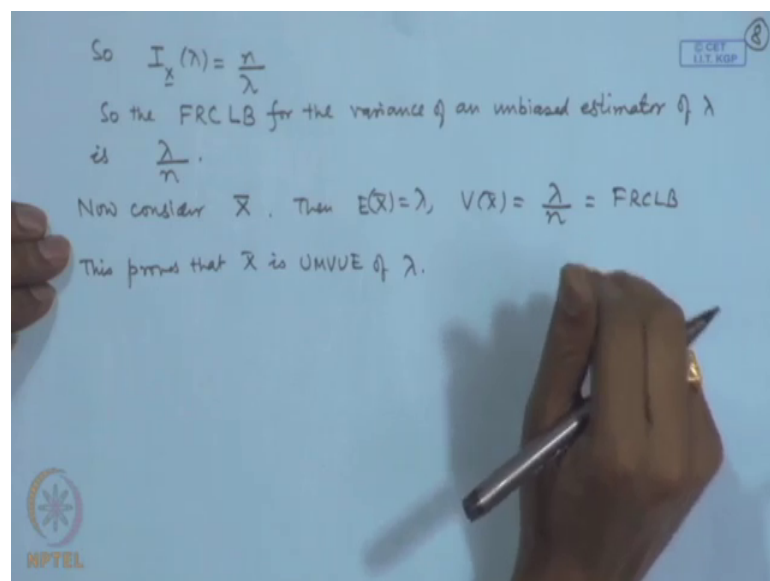
$\frac{\partial \log f}{\partial \lambda} = -1 + \frac{x}{\lambda} = \frac{x - \lambda}{\lambda}$, $E \left(\frac{\partial \log f}{\partial \lambda} \right)^2 = \frac{E(x - \lambda)^2}{\lambda^2} = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$

$\frac{\partial \log f}{\partial p}$ whole square that is equal to expectation of x minus $n p$ square divided by p square and 2 1 minus p square. Now, this is nothing but the variance of x that is $n p$ into 1 minus p in a binomial distribution. So, you get it as n by p into 1 minus p . So, the FRC lower bound for the variance of an unbiased estimator of p is p into 1 minus p by n . Now, in this particular case, you observe here variance of x by n was equal to p into 1 minus p by n which equals variance of x by n here. So, x by n is uniformly minimum variance unbiased estimator of p , so that is uniformly minimum variance unbiased estimator. So, you can see here the method is quite useful in actually proving that a given estimator is UMVUE.

Now let us take say Poisson example. So, suppose we have a random sample from Poisson distribution with the parameter λ . So, naturally we want to estimate

lambda. Now, let us consider the density function e to the power minus lambda, lambda to the power x by x factorial log of f that is equal to minus lambda plus x log of lambda minus log of x factorial. So, if we consider the derivative of this with respect to lambda, then we get minus 1 plus x by lambda that we can write as x minus lambda by lambda. So, expectation of $\frac{\partial \log f}{\partial \lambda}$ square that will be equal to expectation of $\frac{x - \lambda}{\lambda}$ square by lambda square. Now, in the Poisson distribution case, expectation of x is lambda. Therefore, this is nothing but the variance and this is also lambda. So, this is lambda by lambda square that is equal to $\frac{1}{\lambda}$.

(Refer Slide Time: 32:19)



That gives us so you get here the information as n by lambda. So, the FRC lower bound for the variance of an unbiased estimator of lambda is lambda by n . Now, consider say \bar{x} then expectation of \bar{x} is lambda what is variance of \bar{x} variance is equal to lambda by n which is equal to this FRC lower bound. This proves that \bar{x} is UMVUE of lambda. In this particular case, in the Poisson example I had given several unbiased estimators. For example, S^2 I had given $x_1 + x_2$ by 2 I had considered each x_i 's also unbiased for lambda. But you can see that among all of them, \bar{x} will be preferred, because this is the uniformly minimum variance unbiased estimator.

That is all for today's lecture.