

In this lecture we are going to devise methods to find estimators of population parameters satisfying certain criterion. A population parameter is estimated using what is called a sample statistic.

A sample statistic is a function of a random sample. Let us consider a population $f_X(\cdot; \theta)$, where θ denotes the population parameter.

If θ is known then, the distribution of X is completely known. In fact the idea of estimation comes in when θ is not known. In fact θ can be a vector or a scalar. If $X \sim \text{Poisson}(\lambda)$, then $\theta = \lambda$, a real scalar, whereas if $X \sim N(\mu, \sigma^2)$, then $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$ with $\theta_1 = \mu$ & $\theta_2 = \sigma^2$.

For the moment assume $\theta \in \mathbb{R}$. Then having a random sample from x_1, \dots, x_n of size n from $f_X(\cdot, \theta)$, we construct a statistic symbolically given as

$$T = l(x_1, x_2, \dots, x_n)$$

We say that the random variable T is an unbiased estimator of θ if

$$E(T) = \theta.$$

Now the expectation of T is calculated based on the sampling distribution of T .

The expression $E(T) - \theta$ is called the bias in the estimation

if $E(T) - \theta \neq 0$.

If $\theta \in \mathbb{R}^k$ say i.e. $\theta = (\theta_1, \dots, \theta_k)$, then we can have

k different ^{statistics} estimators, $T_i = l_i(x_1, \dots, x_n)$, $i=1, \dots, k$

which can acts as estimator of θ .

Note that \bar{X} , the sample mean & S^2 the sample variance have been proved in the last chapter as unbiased estimators of population ^{mean} variance and population variance respectively.

(1)

Such an approach to get an approximation for θ is often called point estimate estimation. The statistic T when viewed as a random variable is called an estimator and the specific value of T is called an estimate. Sometimes $\hat{\theta}$ is denoted as an estimate of θ . There are several approaches to find an estimator. We would also sometimes refer $\hat{\theta}$ as an estimator of θ corresponding to the estimate $\hat{\theta}$.

Section 1: Method of Moments

Let x_1, \dots, x_n be a random sample of size n drawn from a population $f_x(x; \theta_1, \dots, \theta_k)$, where $\theta \in \mathbb{R}^k$, with $\theta = (\theta_1, \dots, \theta_k)$. Let us recall the r -th moment of X , i.e

$$\mu'_r = E[X^r]$$

$$\therefore \mu'_r = \int x^r f_x(x; \theta_1, \dots, \theta_k) dx$$

\therefore We can write

$$\mu'_r = \mu'_r(\theta_1, \dots, \theta_k)$$

The r -th sample moment

$$M'_r = \frac{x_1^r + x_2^r + \dots + x_n^r}{n} = \frac{1}{n} \sum_{i=1}^n x_i^r$$

To find $\theta = (\theta_1, \dots, \theta_k)$, set up the k -equations, in $\theta_1, \dots, \theta_k$

$$\mu'_j(\theta_1, \dots, \theta_k) = M'_j, \quad j=1, \dots, k$$

Suppose $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ are the solutions to the equations. If say $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$ is unique, then we can consider $\hat{\theta}_1, \dots, \hat{\theta}_k$ to be the estimates of $\theta_1, \dots, \theta_k$.

E.g. 1: x_1, \dots, x_n be a random sample from an exponential population

$$f_x(x, \theta) = \theta e^{-\theta x}; \quad \theta > 0, \quad x \in (0, \infty).$$

Thus we will solve ~~for one equation~~

$$\therefore M'_1 = \mu'_1 = \mu'_1(\theta) = \frac{1}{\theta}$$

$\therefore \hat{\theta} = \frac{1}{M'_1}$ So if we take a particular sample
 x_1, \dots, x_n , then $\hat{\theta} = \frac{n}{x_1 + \dots + x_n}$

So $\hat{\theta}$ is estimated by $\frac{1}{\bar{x}} = \frac{n}{x_1 + \dots + x_n}$

$$\therefore \hat{\theta} = l(\bar{x}) (x_1, \dots, x_n) = \frac{n}{x_1 + \dots + x_n}$$

$\hat{\theta}$ is the estimator.

E.g. 2: x_1, \dots, x_n be a random sample of size n , drawn from a normal population with mean μ and variance σ^2
 $\therefore \theta = (\mu, \sigma^2)$.

$$\therefore M'_1 = \mu'_1 = \theta/\mu$$

$$M'_2 = \mu'_2 = \sigma^2 + \mu^2$$

μ is estimated by $M'_1 = \bar{x}$, i.e. \bar{x} is the unbiased estimator of μ .

σ^2 is estimated by the statistic

$$\therefore \hat{\theta}_2 = \frac{M'_2 - \mu^2}{n} = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{n}{n} \bar{x}^2$$

$$\begin{aligned} \text{Now } \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 &= \frac{1}{n} \left[\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2 \right] \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x} \underbrace{\frac{1}{n} \sum_{i=1}^n x_i}_{\bar{x}} + \frac{n\bar{x}^2}{n} \\ &= \frac{1}{n} \sum_{i=1}^n x_i^2 - 2\bar{x}^2 + \bar{x}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \end{aligned}$$

(3)

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$\therefore \hat{\theta}$ is the estimator of σ^2 and note that through the method of moments we do not have S^2 as the estimator of σ^2 . Thus $\hat{\theta}$ is not an un-biased estimator of σ^2 . Thus method of moments need not give us unbiased estimators always.

Section 2: Maximum Likelihood Estimator

The technique of maximum likelihood function is based on using the fact that x_1, \dots, x_n are iid random variables.

Let us first look at the simple case where $\theta \in \mathbb{R}$. Let us draw a sample of size n , from a population $f_X(x; \theta)$.

The likelihood function in this case is defined

as

$$L(\theta) = f_{X_1}(x_1, \theta) \cdot f_{X_2}(x_2, \theta) \cdots f_{X_n}(x_n, \theta)$$

$$L(x, \theta) = L(\theta) = f_{X_1} \dots f_{X_n}(x_1, \dots, x_n, \theta) \rightarrow (\text{By independence})$$

$$L(\theta) = f_{X_1}(x_1, \theta), f_{X_2}(x_2, \theta), \dots, f_{X_n}(x_n, \theta)$$

$$\therefore L(\theta) = \prod_{i=1}^n f_{X_i}(x_i, \theta) *$$

We compute $\max_{\theta \in A} L(\theta)$, where A is the set in which θ is restricted to belong. We find that value of θ , which maximizes

$L(\theta) = L(\theta; x_1, \dots, x_n)$, for a chosen sample. Then if $\hat{\theta}$ is the maximizer for a chosen sample i.e. $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$, then $\hat{\theta}(x_1, \dots, x_n)$ is called an maximum likelihood estimator estimate of θ for the chosen sample values, $X_1 = x_1, \dots, X_n = x_n$. So the estimator of θ is $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$.

thus $\hat{\theta}$ is called the maximum likelihood estimator of θ

Noting that $\ln(\text{base. loge})$ is an increasing function. We know that if $L(\theta) > 0, \forall \theta \in A$, then if $\hat{\theta}$ is the minimizer, we have

$$L(\theta) \leq L(\hat{\theta}), \quad \forall \theta \in A$$

$$\Rightarrow \ln L(\theta) \leq \ln L(\hat{\theta}) \quad \forall \theta \in A$$

Thus $\hat{\theta}$ is also the maximizer of $\ln L(\theta)$, over A . In most cases the set A is an open interval. Thus to compute $\hat{\theta}$ we can use the equation

$$\frac{\partial}{\partial \theta} \ln L(\theta) = 0$$

$$\frac{\partial}{\partial \theta} \ln L(\theta) = 0$$

$$\boxed{\frac{\partial}{\partial \theta} \ln L(\theta) = 0} \rightarrow (\ast)$$

Then also check the ex second order condition $\left. \frac{\partial^2}{\partial \theta^2} \ln L(\theta) \right|_{\theta=\hat{\theta}} < 0$

where $\left. \frac{\partial}{\partial \theta} \ln L(\theta) \right|_{\theta=\hat{\theta}} = 0$, $\hat{\theta}$ solves (\ast) . Such a $\hat{\theta}$ will provide an estimator or maximum likelihood estimator using the expression

$$\hat{\theta} = \hat{\theta}(x_1, \dots, x_n).$$

Now suppose $\theta \in \mathbb{R}^k$, and we draw a random sample of size n x_1, \dots, x_n , from a population described by $f_x(x; \theta) = f_{x_1}(x_1, \theta_1, \dots, \theta_k)$.

In this case we have the likelihood function as

$$\begin{aligned} L(x, \theta_1, \dots, \theta_n) &= f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta_1, \dots, \theta_k) \\ &= f_{x_1}(x_1, \theta_1, \dots, \theta_k) \cdot f_2(x_2, \theta_1, \dots, \theta_k) \cdots f_n(x_n, \theta_1, \dots, \theta_k) \\ &= \prod_{i=1}^n f_{x_i}(x_i, \theta_1, \dots, \theta_k). \end{aligned}$$

as before
 \therefore Similarly, the maximum likelihood function ~~estimate~~ estimation problem is given as.

$$\max_{\theta \in K} L(x, \theta_1, \dots, \theta_k), \quad \text{where } K \subset \mathbb{R}^k$$

$$\max_{\theta \in K} \ln L(x, \theta_1, \dots, \theta_k), \quad \text{where } K \subset \mathbb{R}^k$$

(5)

Eng Now if $\hat{\theta}$ is the maximizer, then as before we have $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_n)$

$$\hat{\theta}_i = \hat{\theta}_i(x_1, \dots, x_n), \quad i=1, \dots, n$$

So there are n , maximum likelihood estimators. We shall now provide two examples. In the first case where $\theta \in \mathbb{R}$, and in the second case $\theta \in \mathbb{R}^2$.

E.g. 3: Let x_1, \dots, x_n be a random sample drawn from an exponential population $f_X(\cdot, \theta)$, and is given as

$$f_X(x_i, \theta) = \theta e^{-\theta x_i}, \quad x_i > 0, \quad \theta > 0.$$

$$\therefore L(\theta, x) = L(\theta) = \prod_{i=1}^n f_{X_i}(x_i, \theta)$$

$$= \prod_{i=1}^n \theta e^{-\theta x_i}$$

$$= \prod_{i=1}^n (\theta)^n e^{-\theta \sum_{i=1}^n x_i}$$

The maximum likelihood estimation problem (MLE)
 $\max_{\theta > 0} \ln L(\theta)$ (P1)

$$\therefore \ln L(\theta) = n \ln \theta + (-\theta \sum_{i=1}^n x_i)$$

$$\text{Thus } \frac{\partial}{\partial \theta} \ln L(\theta) = n \cdot \frac{1}{\theta} - \sum_{i=1}^n x_i = 0$$

$$\therefore \frac{n}{\theta} = \sum_{i=1}^n x_i$$

$$\therefore \frac{1}{\theta} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

$$\text{or } \theta = \frac{1}{\bar{x}}$$

$$\text{Now } \begin{aligned} \frac{\partial^2}{\partial \theta^2} L(\theta) &= -\frac{n}{\theta^2} && \left. \right\} \text{ Thus } \theta = \frac{1}{\bar{x}} \text{ maximizes (strictly)} \\ &= -n \bar{x}^2 < 0 && \text{problem P1).} \end{aligned}$$

$\therefore \hat{\theta} = \frac{1}{\bar{x}}$ is the mle estimator of θ

(6) m.l.e: Short form for "maximum likelihood estimator!"

Now for the case when there are more than one parameter, i.e. $\theta \in \mathbb{R}^k$, we first have to find $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ satisfying

$$\frac{\partial \ln L(\theta_1, \dots, \theta_k)}{\partial \theta_j} = 0 \quad \text{for } j=1, \dots, k.$$

Once we find a $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$, we compute the Hessian matrix of $\nabla^2 \ln L(\theta_1, \dots, \theta_k)$, and see if this negative definite at $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$. In the above expression we took the liberty of writing $L(x, \theta_1, \dots, \theta_k) = L(\theta_1, \dots, \theta_k)$.

When $\theta \in \mathbb{R}^2$, i.e. $\theta = (\theta_1, \theta_2)$, then we have to check the following. First find $(\hat{\theta}_1, \hat{\theta}_2)$ satisfying

$$\left. \begin{aligned} \frac{\partial \ln L(\theta_1, \theta_2)}{\partial \theta_1} &= 0 \\ \frac{\partial \ln L(\theta_1, \theta_2)}{\partial \theta_2} &= 0 \end{aligned} \right\}$$

Then the Hessian matrix of $\ln L(\theta_1, \theta_2)$ is given as

$$\nabla^2 \ln L(\theta_1, \theta_2) = \begin{bmatrix} \frac{\partial^2 \ln L}{\partial \theta_1^2}, & \frac{\partial^2 \ln L}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2 \ln L}{\partial \theta_2 \partial \theta_1}, & \frac{\partial^2 \ln L}{\partial \theta_2^2} \end{bmatrix}$$

To show that $\nabla^2 \ln L(\theta_1, \theta_2)$ is positive definite we have
 $\frac{\partial^2 \ln L}{\partial \theta_1^2} < 0$ and $\frac{\partial^2 \ln L}{\partial \theta_2^2} < 0$ and $\det[\nabla^2 \ln L(\theta_1, \theta_2)] > 0$ evaluated at $(\hat{\theta}_1, \hat{\theta}_2)$. Or find the eigen-values which both has to be negative.

Note: If we have a diagonal matrix, which is 2×2 , then to check whether it is negative definite, just see if the diagonal elements are negative as they are the eigenvalues.

E.g.: Let us consider a random sample from a normal population, with mean μ and variance σ^2 .

$$\begin{aligned} \therefore L(x, \mu, \sigma^2) &= \prod_{i=1}^n f_{x_i}(x_i, \mu, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sigma(\sqrt{2\pi})} e^{-\frac{1}{2\sigma^2}(x_i-\mu)^2} \\ &= \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2} \end{aligned}$$

where $\sigma > 0, -\infty < \mu < +\infty$

$$\therefore \ln(L(x, \mu, \sigma^2)) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2$$

\therefore The maximum likelihood estimation problem; ie MLE problem is

$$\max_{\begin{array}{l} \sigma > 0 \\ -\infty < \mu < +\infty \end{array}} \ln(L(x, \mu, \sigma^2)) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i-\mu)^2$$

$$\begin{aligned} \therefore \frac{\partial}{\partial \mu} \ln L(x, \mu, \sigma^2) &= 0 \quad \& \frac{\partial}{\partial \sigma^2} \ln L(x, \mu, \sigma^2) = 0 \\ &\Downarrow \\ \therefore \frac{1}{\sigma^2} \sum_{i=1}^n (x_i-\mu) &= 0 \quad \& -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i-\mu)^2 = 0 \\ &\Downarrow \\ &- \frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i-\mu)^2 = 0 \end{aligned}$$

$$\therefore \begin{cases} \hat{\mu} = \frac{1}{n} \sum x_i = \bar{x} \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \end{cases} \rightarrow \text{The MLE estimators?}$$

To confirm that these are MLE estimators we have to establish the fact that these are maximizers. The simplest way to realize this is that $\ln L(x, \mu, \sigma^2)$ is a convex function of (μ, σ^2) and hence the critical point is a global maximizer. But this is beyond the scope of this course. Hence we check the usual way stated in the previous page.

The Hessian matrix $\nabla^2 \ln L(x, \mu, \sigma^2)$ evaluated at $(\hat{\mu}, \hat{\sigma}^2)$ is given as

$$\nabla^2 \ln L(x, \mu, \sigma^2) = \begin{pmatrix} -\frac{n}{\sigma^4} & 0 \\ 0 & \frac{1}{\sigma^4} \left[\frac{n}{2} + n^2 \right] \end{pmatrix}$$

$$\text{Thus } -\frac{n}{\sigma^4} < 0 \quad \text{and} \quad \frac{1}{\sigma^4} \left[\frac{n}{2} + n^2 \right] > 0, \text{ hence}$$

$\nabla^2 \ln L(x, \mu, \sigma^2)$ is negative definite. Thus $\hat{\mu}$ & $\hat{\sigma}^2$ are the maximizers of the mle problem.

Detail computation of $\nabla^2 \ln L(x, \mu, \sigma^2)$:

$$\frac{\partial}{\partial \mu} \frac{\partial}{\partial \mu} \ln L(x, \mu, \sigma^2) = -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu)$$

$$\therefore \frac{\partial^2}{\partial \mu^2} \ln L(x, \mu, \sigma^2) = \sum_{i=1}^n \frac{1}{\sigma^4} = \frac{n}{\sigma^4}$$

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \frac{\partial}{\partial \mu} \ln L(x, \mu, \sigma^2) &= \frac{\partial}{\partial \sigma^2} \ln L(x, \mu, \hat{\sigma}^2) \\ &= -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) \\ &= -\frac{1}{\hat{\sigma}^4} \left[\sum_{i=1}^n (x_i - \bar{x}) \right] \\ &= -\frac{1}{\hat{\sigma}^4} \left(\sum_{i=1}^n x_i - n\bar{x} \right) = 0 \end{aligned}$$

$$\frac{\partial}{\partial \sigma^2} \frac{\partial}{\partial \sigma^2} \ln L(x, \mu, \sigma^2) = \frac{n}{2} \frac{1}{\sigma^4} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

$$\therefore \frac{\partial}{\partial \mu} \frac{\partial}{\partial \sigma^2} \ln L(x, \mu, \hat{\sigma}^2) = \frac{1}{2\hat{\sigma}^2} \cdot 2 = \left(\sum_{i=1}^n (x_i - \bar{x}) \right) = 0$$

$$\frac{\partial}{\partial \sigma^2} \left(\frac{\partial}{\partial \sigma^2} \ln L(x, \mu, \hat{\sigma}^2) \right) = \frac{n}{2} \frac{1}{\hat{\sigma}^4} = \frac{n}{\hat{\sigma}^4} \left[\frac{n}{2} + n^2 \right]$$

(we use $\sum_{i=1}^n (x_i - \bar{x})^2 = n\hat{\sigma}^2$) (9)

We would like to note that second order condition only tells us that the ~~more~~ estimators that we have got are ~~really~~ only local minimizer albeit strict one. To establish that they are global we need to use ~~convexity~~ ^{concavity}. The second order check also establishes the concavity of functions here.

A crash course on Concavity Let I be an interval in \mathbb{R} and $f: I \rightarrow \mathbb{R}$.

The function $f: I \rightarrow \mathbb{R}$ is called concave if for any $x, y \in I$ & $\lambda \in [0, 1]$

$$f(\lambda y + (1-\lambda)x) \geq \lambda f(y) + (1-\lambda) f(x). \quad \rightarrow (A)$$

Note: If $x = y$ are in I ; then for any $\lambda \in [0, 1]$, $\lambda y + (1-\lambda)x \in I$

• Thus if f is differentiable, then from (A), by Taylor's Theorem.

$$f(x) + \lambda f'(x)(y-x) + o(\lambda) \geq \lambda f(y) + (1-\lambda) f(x), \quad \forall \lambda \in (0, 1)$$

where $\frac{o(\lambda)}{\lambda} \rightarrow 0$ as $\lambda \rightarrow 0$.

$$\therefore \lambda f'(x)(y-x) + o(\lambda) \geq \lambda (f(y) - f(x))$$

$$\therefore \text{As } \lambda \rightarrow 0, \quad f'(x)(y-x) \geq f(y) - f(x)$$

If $f'(x) = 0 \Rightarrow f(x) \geq f(y) \Rightarrow x$ is a global maximum maximizer

Let $f: I \times I \rightarrow \mathbb{R}$, then f is concave if for any $(x_1, x_2) \in I \times I$

and $(y_1, y_2) \in I_1 \times I_2$ & $\lambda \in (0, 1)$

$$f(\lambda y_1 + (1-\lambda)x_1, \lambda y_2 + (1-\lambda)x_2) \geq \lambda f(y_1, y_2) + (1-\lambda) f(x_1, x_2)$$

Again using Taylor's Theorem in 2-dimensions we have $\forall (x_1, x_2), (y_1, y_2)$

$$\frac{\partial f}{\partial x_1}(y_1 - x_1) + \frac{\partial f}{\partial x_2}(y_2 - x_2) \geq f(y_1, y_2) - f(x_1, x_2).$$

$$\therefore \frac{\partial f}{\partial x_1} = 0 \text{ & } \frac{\partial f}{\partial x_2} = 0, \quad f(x_1, x_2) \geq f(y_1, y_2), \quad \forall (y_1, y_2) \in I \times I$$

$\Rightarrow (x_1, x_2)$ is the global maximizer.

E.g.: $-\alpha(x-\mu)^2 = f(x),$

$\alpha > 0, \mu \in \mathbb{R}$

is concave in x

$$f(x) = -\ln x$$

$x > 0, i.e. x \in \mathbb{R}^+$

is concave in x .

(10)

If ~~the~~ $f_{xx} < 0$, then

f is concave on I if

$f''(x) \leq 0, \forall x \in I$

f is concave on $I \times I$ if

$\nabla^2 f$ is negative definite on $I \times I$

In our study in the example the log likelihood functions can be shown to be concave.

Section 3 Fisher Information & Cramer-Rao Inequality

The maximum likelihood estimation technique, is very popular, and the reason for this is follows. Given a sample observation,

$x_1 = x_1, \dots, x_n = x_n$ the likelihood function $L(x, \theta)$

provides us the "likelihood" or the frequency of the occurrence of the observation x_1, x_2, \dots, x_n . Since the parameter θ is unknown we ask the question: What is the distribution for which the observations x_1, \dots, x_n occur the maximum number of times? This can be done by finding the θ , which maximizes the joint pdf, since it is the parameter θ , which specifies the distribution. We shall now see how the likelihood function can be used to build some important measures related to estimators.

We shall begin with the notion of Fisher Information, but let us first write down the two basic assumptions we need. We will ~~first~~ study only the single parameter case.

Assumptions

a) The pdf $f_x(x, \theta)$ is twice continuously differentiable as a function of the parameter θ .

b) We can differentiate ~~twice~~ with respect to θ under the integral sign for the integral $\int f_x(x, \theta) dx$. In fact we will assume that we can differentiate twice under the integral sign with respect to θ .

c) The parameter θ belongs to an open interval.

The Fisher information function is built on the idea of a Fisher score function.

Let x_1, \dots, x_n be a random sample from a population $f_x(x, \theta)$.

The Fisher score function for the i th observation is given as the r.v. $F(x_i, \theta)$

whose value for the observation $x=x_i$ is given as $\bar{F}(x_i, \theta) = \frac{\partial}{\partial \theta} \ln f_x(x_i, \theta)$

$$\therefore F(x_i, \theta) = \frac{1}{f_x(x_i, \theta)} \cdot \frac{\partial}{\partial \theta} f_x(x_i, \theta)$$

$$\therefore E[F(x_i, \theta)] = \int_{-\infty}^{\infty} \frac{1}{f_x(x_i, \theta)} \frac{\partial}{\partial \theta} f_x(x_i, \theta) f_x(x_i, \theta) dx$$

$$= \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f_x(x_i, \theta) dx$$

$$= \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f_x(x_i, \theta) dx = \frac{\partial}{\partial \theta} (1) = 0.$$

$$\therefore E[\bar{F}(x_i, \theta)] = 0$$

The total Fisher score function for a sample of size n is given as the r.v.

$$\bar{F}_n(\theta) = \sum_{i=1}^n \cancel{F(x_i, \theta)}$$

$$\boxed{\bar{F}_n(\theta) = \sum_{i=1}^n \bar{F}(x_i, \theta)}$$

$$\boxed{\text{Note } E(\bar{F}_n(\theta)) = 0}$$

The first Fisher information about the i -th observation is denoted as $I(\theta)$ and is given as

$$\boxed{I(\theta) = \text{Var}(F(x_i, \theta))}$$

$$\begin{aligned}
 I_i(\theta) &= \text{Var}(F(x_i, \theta)) \\
 &= \sqrt{E[(F(x_i, \theta))^2]} \quad (\because E(F(x_i, \theta)) = 0) \\
 &= \sqrt{E\left[\left(\frac{\partial}{\partial \theta} \ln f_x(x_i, \theta)\right)^2\right]} \\
 &= \int_{-\infty}^{\infty} \left[\frac{\partial}{\partial \theta} \ln f_x(x_i, \theta) \right]^2 f_x(x_i, \theta) dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{f_x(x_i, \theta)} \cdot \frac{1}{f_x^2(x_i, \theta)} \left[\frac{\partial}{\partial \theta} f_x(x_i, \theta) \right]^2 f_x(x_i, \theta) dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{f_x(x_i, \theta)} \left[\frac{\partial}{\partial \theta} f_x(x_i, \theta) \right]^2 dx
 \end{aligned}$$

Usually in texts one writes

$$I_i(\theta) = E\left[\left(\frac{\partial}{\partial \theta} \ln f_x(x_i, \theta)\right)^2\right]$$

For a random sample of size n , the Fisher information is given as.

E.g1: Consider a random sample drawn from a population $f_x(x, \theta)$, given as drawn from an exponential population; $f_x(x, \theta)$, given as

$$\therefore f_x(x, \theta) = \theta e^{-\theta x}; \quad x > 0, \theta > 0$$

Thus

$$\begin{aligned}
 F(x_i, \theta) &= \frac{\partial}{\partial \theta} \ln(\theta e^{-\theta x_i}) \\
 &= \frac{1}{\theta e^{-\theta x_i}} \left[\frac{\partial}{\partial \theta} (\theta e^{-\theta x_i}) \right] \\
 &= \frac{1}{\theta e^{-\theta x_i}} \left[e^{-\theta x_i} + \theta e^{-\theta x_i} (-x_i) \right]
 \end{aligned}$$

(13)

$$\begin{aligned}
\therefore \bar{F}(x_i, \theta) &= \frac{1}{\theta} \left[1 - e^{-\theta x_i} \right] \\
&= \frac{(1 - \theta x_i)}{\theta} \\
\therefore \bar{F}(x_i, \theta) &= \left[\frac{1 - \theta x_i}{\theta} \right] \\
\\
\therefore I(\theta) &= E \left[(\bar{F}(x_i, \theta))^2 \right] \\
&= \int_0^\infty \frac{(1 - \theta x_i)^2}{\theta} \cdot \theta e^{-\theta x_i} dx \\
&= \int_0^\infty \frac{(1 - \theta x_i)^2}{\theta^2} \cdot \theta e^{-\theta x_i} dx \\
&= \frac{1}{\theta} \int_0^\infty (1 - \theta x_i)^2 e^{-\theta x_i} dx \\
&= \frac{1}{\theta} \int_0^\infty (1 - 2\theta x_i + \theta^2 x_i^2) \cdot e^{-\theta x_i} dx \\
&= \frac{1}{\theta} \int_0^\infty e^{-\theta x_i} dx - \frac{2}{\theta} \int_0^\infty x_i \cdot \theta e^{-\theta x_i} dx \\
&\quad + \cancel{\frac{1}{\theta} \int_0^\infty \theta^2 x_i^2 e^{-\theta x_i} dx} \\
&= \frac{1}{\theta^2} \int_0^\infty \theta e^{-\theta x_i} dx - \frac{2}{\theta} E(x_i) + \cancel{\frac{1}{\theta} E(x_i^2)} \\
&= \frac{1}{\theta^2} - \frac{2}{\theta} \cdot \frac{1}{\theta} + \cancel{\frac{1}{\theta} \left[\text{Var}(x_i) + (E(x_i))^2 \right]} \\
&= \frac{1}{\theta^2} - \frac{2}{\theta^2} + \cancel{\frac{1}{\theta} \left[\frac{1}{\theta^2} + \frac{1}{\theta^2} \right]} \\
&= \frac{1}{\theta^2} - \frac{2}{\theta^2} + \cancel{\frac{1}{\theta^2} \left[\frac{2}{\theta^2} \right]} = -\frac{1}{\theta^2} + \cancel{\frac{2}{\theta^4}} \\
&= \cancel{\frac{2}{\theta^2}} - \frac{1}{\theta^2} = \frac{1}{\theta^2} \left[\frac{2}{\theta^2} - 1 \right]
\end{aligned}$$

(14)

So the Fisher information for the whole sample is given as

$$I_n(\theta) = \text{Var}(F_n(\theta)) \\ = \text{Var}(\theta F(x_1, \theta) + \dots + F(x_n, \theta))$$

Since x_1, \dots, x_n are independent, $F(x_1, \theta), \dots, F(x_n, \theta)$ are independent. Thus

$$I_n(\theta) = \text{Var}(F_{\theta}(x_1, \theta)) + \dots + \text{Var}(F(x_n, \theta)) \\ = \sum I_i(\theta)$$

$$\therefore I_n(\theta) = \sum_{i=1}^n I_i(\theta)$$

But each x_i 's are having the same distribution.

$$\therefore \boxed{I_n(\theta) = n I_i(\theta)}$$

The Cramer-Rao Inequality

A random sampling is called regular if its Fisher information is continuous, strictly positive and bounded for all θ in the given range. For example if we consider a random sample from x_1, \dots, x_n from an exponential distribution, then

$$I_n(\theta) = n I_i(\theta) \\ = \frac{n}{\theta^2} > 0$$

$I_n(\theta)$ is of course continuous on $\theta > 0$. However, $I_n(\theta)$ is not bounded above, as a function of θ . But as $I_n(\theta)$ is continuous and positive we can consider the sampling from an exponential distribution is regular.

Theorem II.1: Cramer-Rao Bound (Actually a lower bound)

Let $T = \hat{\theta}(x_1, \dots, x_n)$ be a statistic which an estimator of θ of the population $f(x, \theta)$. Let the bias $b_n(\theta) = E[\hat{\theta}] - \theta$ be continuously differentiable. Then

$$\text{Var}(\hat{\theta}) \geq \frac{(1 + b'_n(\theta))^2}{I_n(\theta)}$$

Proof:

$$\theta + b_n(\theta) = E[\hat{\theta}] = \int_{\mathbb{R}^n} \hat{\theta}_n(x_1, \dots, x_n) f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) dx_1 \dots dx_n$$

Now differentiation under the integral sign we have.

$$\begin{aligned} 1 + b'_n(\theta) &= \int_{\mathbb{R}^n} \hat{\theta}_n(x_1, \dots, x_n) \frac{\partial f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta)}{\partial \theta} dx_1 \dots dx_n \\ &= \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) \frac{f_x(x_1, \dots, x_n, \theta)}{f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta)} \frac{\partial}{\partial \theta} f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) dx_1 \dots dx_n \\ &= \int_{\mathbb{R}^n} \hat{\theta}(x_1, \dots, x_n) F_n(\theta) f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) dx_1 \dots dx_n \end{aligned}$$

$$\begin{aligned} F_n(\theta) &= \sum_{i=1}^n F(x_i, \theta) = \sum_{i=1}^n \frac{1}{f_x(x_i, \theta)} \frac{\partial}{\partial \theta} f_x(x_i, \theta) \\ &= \sum_{i=1}^n \frac{1}{f_{x_i}(x_i, \theta)} \frac{\partial}{\partial \theta} f_{x_1, \dots, x_n}(x_1, \dots, x_n, \theta) \end{aligned}$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln f_x(x_i, \theta)$$

$$= \frac{\partial}{\partial \theta} \left[\sum_{i=1}^n \ln f_{x_i}(x_i, \theta) \right]$$

$$= \frac{\partial}{\partial \theta} \left[\ln \cancel{f_x(x)} \ln \prod_{i=1}^n f_{x_i}(x_i, \theta) \right] = \frac{\partial}{\partial \theta} \ln \cancel{f_x(x)} \downarrow$$

(By independence).

$$\therefore F_n(\theta) = \frac{1}{\int_{x_1 x_2 \dots x_n} f_{x_1 \dots x_n}(x_1 \dots x_n, \theta)} \frac{\partial}{\partial \theta} f_{x_1 \dots x_n}(x_1 \dots x_n, \theta)$$

$$\therefore 1 + b'_n(\theta) = E[\hat{\theta} F_n(\theta)] = \text{cov}(\hat{\theta}, F_n(\theta))$$

Note that $\text{cov}(\hat{\theta}, F_n(\theta)) = E[\hat{\theta} F_n(\theta)] = -E[\hat{\theta}] E[F_n(\theta)]$

$$= E[\hat{\theta} F_n(\theta)] \quad (\because E[F_n(\theta)] = 0 \text{ as derived earlier})$$

\therefore let $P_n^2 = \text{correlation coefficient of } \hat{\theta}, F_n(\theta) \text{ and } F_n(\theta)$

$$\therefore P_n^2 = \frac{\text{cov}(\hat{\theta}, F_n(\theta))}{\text{Var}(\hat{\theta}) \text{Var}(F_n(\theta))}$$

But $P_n^2 \leq 1$ as $P_n \in [-1, +1]$

$$\Rightarrow \frac{\text{cov}(\hat{\theta}, F_n(\theta))}{\text{Var}(\hat{\theta}) \text{Var}(F_n(\theta))} \leq 1.$$

Hence

$$\text{Var}(\hat{\theta}) \geq \frac{\text{cov}(\hat{\theta}, F_n(\theta))}{\text{Var}(F_n(\theta))}$$

$$\Rightarrow \boxed{\text{Var}(\hat{\theta}) = \frac{(1 + b'_n(\theta))^2}{I_n(\theta)}} \quad \square$$

This is called the Cramer-Rao Lower bound. Now if $\hat{\theta}$ is an unbiased estimator, then $b'_n(\theta) = 0 \Rightarrow$ for any θ in the given range. Thus we have

$$\boxed{\text{Var}(\hat{\theta}) \geq \frac{1}{I_n(\theta)}}$$

This is called the Cramer-Rao Inequality. An unbiased estimator $\hat{\theta}$ is called efficient if

$$\boxed{\text{Var}(\hat{\theta}) = \frac{1}{I_n(\theta)}}$$

Let us finish our discussion with an example by drawing a sample of size n ; from a normal population $N(\mu, \sigma^2)$, where μ is unknown but σ^2 is ~~not~~ known.

We shall show that $\bar{X}_n = \bar{X} = \frac{x_1 + \dots + x_n}{n}$ is indeed an efficient estimator of μ . \bar{X}_n is known to be unbiased as we have

already shown. Now $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$. We have

$$f_X(x, \mu) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

$$\therefore \ln f_{X_i}(x_i, \mu) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

In the random-variable form we have

$$\ln f_{X_i}(\bar{x}_i, \mu) = \frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$F(x_i, \mu) = \frac{x_i - \mu}{\sigma^2}$$

$$\begin{aligned} I_n(\theta) &= n I_i(\theta) = n \text{Var}(F(x_i, \mu)) = n E[(F(x_i, \mu))^2] \\ &= n E[\cancel{(F(x_i, \mu))^2}] \\ &= n E\left[\frac{(x_i - \mu)^2}{\sigma^4}\right] \\ &= \frac{n}{\sigma^4} E[(x_i - \mu)^2] \\ &= \frac{n}{\sigma^4} \cdot \text{Var}(x_i) \\ &= \frac{n}{\sigma^4} \cdot \sigma^2 = \frac{n}{\sigma^2} \end{aligned}$$

Now $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n} = \frac{1}{I_n(\theta)}$, proving that

\bar{X}_n is indeed an efficient estimator.

The part on Fisher Information is partially based on Chapter 1, of the book titled: Mathematical Statistics: Asymptotic minimax theory, by A. Korostelev and O. Korosteleva, AMS-2011.