Fisher Information

Assume $X \sim f(x|\theta)$ with $\theta \in \Theta \subset \mathbb{R}$ polfor pmf Define $I_{x}(\Theta) = E_{\Theta}\left(\left(\frac{\partial}{\partial \Theta}\log f(x|\Theta)\right)^{2}\right)$ the derivative of the log-likelihood function evaluated at the true value of O. Fisher information is meaningful for families of distributions which are regular: 1) Fixed support: fx: f(x10) > 0} is the same for all Θ . (2) $\frac{\partial}{\partial A} \log f(x|\theta)$ must exist and be finite for all x and Θ . (3) If $E_{a}W(X) < \infty \forall \theta$, then $\left(\frac{\partial}{\partial \theta}\right)^{\kappa} E_{\theta} W(x) = \left(\frac{\partial}{\partial \theta}\right)^{\kappa} \int W(x) f(x|\theta) dx$ $= \left(W(x) \left(\frac{\partial}{\partial \Theta} \right)^{\kappa} f(x | \Theta) dx \right).$



= $I_{\times}(\Theta)$.

(3) If
$$X = (X_1, X_2, ..., X_n)$$
 and
 $X_1, X_2, ..., X_n$ are independent rv^3s ,
then $I_X(\Theta) = I_X(\Theta) + I_X(\Theta) + ... + I_X(\Theta)$.
Proof:
 $f(X|\Theta) = \prod_{i=1}^n f_i(X_i|\Theta)$
where $f_i(\cdot|\Theta)$ is the pdf (pmf) of X_i .
 $\frac{\partial}{\partial \Theta} \log f(X|\Theta) = \sum_{i=1}^n \frac{\partial}{\partial \Theta} \log f_i(X_i|\Theta)$
and the rv^3s in the sum are independent.
Thus
 $Var\left[\frac{\partial}{\partial \Theta} \log f(X|\Theta)\right] = \sum_{i=1}^n Var\left[\frac{\partial}{\partial \Theta} \log f_i(X_i|\Theta)\right]$
so that $I_X(\Theta) = \sum_{i=1}^n I_X(\Theta)$ by (2).
(4) If $X_1, X_2, ..., X_n$ are iid
and $X = (X_1, ..., X_n)_3$ then $I_X(\Theta) = I_X(\Theta)$.

(5) $I_{x}(\theta) = E_{A}\left(-\frac{\partial^{2}}{\partial\theta^{2}}\log f(X|\theta)\right)$ (Alternate formula for Fisher information) Proof: Abbreviate $\int f(x|\theta) dx$ as $\int f_{1}$ etc. $I = \int f$. Apply $\frac{\partial}{\partial \Theta}$ to both sides. $0 = \frac{\partial}{\partial \partial f} = \int \frac{\partial f}{\partial f} = \int \frac{\partial f}{\partial f} \cdot f$ $= \left(\left(\frac{\partial}{\partial \theta} \log f \right) \cdot f \right)$ Apply à again. $0 = \frac{\partial}{\partial \theta} \left(\left(\frac{\partial}{\partial \theta} \log f \right) \cdot f \right)$ $= \left(\frac{\partial}{\partial \Theta} \left[\left(\frac{\partial}{\partial \Theta} \log f \right) \cdot f \right] \right)$ $= \int \left(\frac{\partial^2}{\partial \theta^2} \log f\right) \cdot f + \int \left(\frac{\partial}{\partial \theta} \log f\right) \cdot \frac{\partial f}{\partial \theta}$

Noting that $\frac{\partial f}{\partial \Theta} = \frac{\partial f}{\frac{\partial \Theta}{\Gamma}} \cdot f$ $= \left(\frac{\partial}{\partial A} \log f\right) \cdot f$

this becomes

 $O = \int \left(\frac{d^2}{d\theta^2} \log f\right) \cdot f + \int \left(\frac{d}{d\theta} \log f\right)^2 \cdot f$

or

 $O = E\left(\frac{\partial}{\partial \theta^2} \log f(X|\theta)\right) + I_X(\theta)$

QED

Example: Fisher information for a Poisson sample. Observe $X = (X_1, \dots, X_n)$ iid Poisson (λ) . Find $I_X(\lambda)$. We know $I_X(\lambda) = nI_{X_1}(\lambda)$.

We shall calculate $I_{X_1}(\lambda)$ in three ways. Let $X = X_1$. Preliminaries:

$$f(x \mid \lambda) = \frac{\lambda^{x} e^{-\lambda}}{x!}$$
$$\log f(x \mid \lambda) = x \log \lambda - \lambda - \log x!$$
$$\frac{\partial}{\partial \lambda} \log f(x \mid \lambda) = \frac{x}{\lambda} - 1$$
$$-\frac{\partial^{2}}{\partial \lambda^{2}} \log f(x \mid \lambda) = \frac{x}{\lambda^{2}}$$

Method # 1:

$$I_X(\lambda) = E_\lambda \left[\left(\frac{\partial}{\partial \lambda} \log f(X \mid \lambda) \right)^2 \right] = E_\lambda \left[\left(\frac{X}{\lambda} - 1 \right)^2 \right]$$
$$= \operatorname{Var}_\lambda \left(\frac{X}{\lambda} \right) \qquad \text{since } E \left(\frac{X}{\lambda} \right) = \frac{EX}{\lambda} = \frac{\lambda}{\lambda} = 1$$
$$= \frac{\operatorname{Var}(X)}{\lambda^2} = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$$

Method # 2 (almost the same):

$$I_X(\lambda) = \operatorname{Var}_{\lambda}\left(\frac{\partial}{\partial\lambda}\log f(X|\lambda)\right) = \operatorname{Var}\left(\frac{X}{\lambda} - 1\right)$$
$$= \operatorname{Var}\left(\frac{X}{\lambda}\right) = \frac{1}{\lambda} \quad \text{as in Method } \#1$$

Method #3:

$$I_X(\lambda) = E_\lambda \left(-\frac{\partial^2}{\partial \lambda^2} \log f(X \mid \lambda) \right) = E_\lambda \left(\frac{X}{\lambda^2} \right) = \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$$

• Thus $I_{\underline{X}}(\lambda) = nI_{X_1}(\lambda) = \frac{n}{\lambda}$.

Example: Fisher information for Cauchy location family. Suppose X1, X2,..., Xn iid with $pdf f(x|\theta) = \frac{1}{\pi(1+(x-\theta)^2)}$ Let $X = (X_1, \dots, X_n), X \sim f(X|\theta).$ Find $I_{X}(\theta)$. $I_{X}(\Theta) = n I_{X}(\Theta) = n I_{X}(\Theta).$ $\frac{\partial}{\partial \Theta} \log f(\chi|\Theta) = \frac{\partial f}{\partial \Theta}$ $\frac{-1}{\pi (1 + (\chi - \theta)^2)^2} \cdot 2(\chi - \theta)(-1)$ $\pi\left(1+(\chi-\theta)^2\right)$ $\frac{2(x-\theta)}{(1+(x-\theta)^2)}$

 $I_{X}(\theta) = E\left[\left(\frac{\partial}{\partial \theta} \log f(X|\theta)\right)^{2}\right]$ $= E \left(\frac{2(X-\theta)}{(+(X-\theta)^2)^2} \right)^2$ $= \int_{-\infty}^{\infty} \left(\frac{2(x-\theta)}{(1+(x-\theta)^2)^2}\right)^2 \frac{1}{\pi(1+(x-\theta)^2)} dx$ $= \frac{4}{\pi} \int_{-\infty}^{\infty} \frac{(\chi - \theta)^2}{(1 + (\chi - \theta)^2)^3} d\chi$ Let $u = x - \theta$, du = dx. $=\frac{4}{\pi}\int_{-\infty}^{\infty}\frac{u^2}{(1+u^2)^3}du$ $=\frac{8}{\pi}\int_{0}^{\infty}\frac{u^{2}}{(1+1)^{2}}du$ substitute $x = \frac{1}{1+u^2}, \quad 1-x = \frac{u^2}{1+u^2}$ $dx = \frac{-2u}{(1+u^2)^2} du \quad (don't use)$

 $\mu = \left(\frac{1}{x} - 1\right)^{1/2}$ $du = \frac{1}{2} \left(\frac{1}{x} - 1 \right)^{-1/2} \left(-\frac{1}{x^2} \right) dx$ Then $\frac{\delta}{\pi} \int_{n}^{\infty} \frac{u^2}{(1+u^2)^3} du$ $= \frac{8}{\pi} \int_{0}^{\infty} \frac{u^{2}}{1+u^{2}} \left(\frac{1}{1+u^{2}}\right)^{2} du$ $= \frac{8}{\pi} \int_{1}^{0} (1-\chi) \chi^{2} \cdot \frac{1}{2} (\frac{1}{\chi} - 1)^{-1/2} (\frac{1}{\chi}) d\chi$ $= 4 \int_{-\infty}^{1} x^{+1/2} (1-x)^{1/2} dx$ $=\frac{4}{\pi}\int_{0}^{1}\chi^{3/2-1}(1-\chi)^{3/2-1}d\chi \text{ (Beta (integral))}$ $= \frac{4}{\pi} \frac{\Gamma(3/2) \Gamma(3/2)}{\Gamma(3/2+3/2)} = \frac{4}{\pi} \frac{(\frac{1}{2}\sqrt{\pi})^2}{21}$ $=\frac{1}{2}$. Therefore $I_X(\theta) = \frac{n}{2}$.

Uses of Fisher Information:

- Asymptotic distribution of MLE's
- Cramér-Rao Inequality (Information Inequality)

Asymptotic distribution of MLE's

• IID Case:

If $f(x | \theta)$ is a regular one-parameter family of pdf's (or pmf's) and $\hat{\theta}_n = \hat{\theta}_n(X_n)$ is the MLE based on $X_n = (X_1, \dots, X_n)$ where n is large and X_1, \dots, X_n are iid from $f(x | \theta)$, then

$$\widehat{ heta}_n \sim \operatorname{approx} N\left(heta \,, \, rac{1}{nI(heta)}
ight) \qquad ext{where } I(heta) \equiv I_{X_1}(heta)$$

and θ is the true value. Note that $nI(\theta) = I_{\pmb{X}_n}(\theta).$ More formally,

$$\frac{\widehat{\theta}_n - \theta}{\sqrt{\frac{1}{nI(\theta)}}} = \sqrt{nI(\theta)}(\widehat{\theta}_n - \theta) \stackrel{d}{\longrightarrow} N(0, 1) \quad \text{as } n \to \infty.$$

• More general case:

(Assuming various regularity conditions) If $f(\mathbf{x}|\theta)$ is a one-parameter family of joint pdf's (or joint pmf's) for data $X_n = (X_1, \ldots, X_n)$ where *n* is large (think of a large data set arising from a regression or time series model) and $\hat{\theta}_n = \hat{\theta}_n(X_n)$ is the MLE, then

$$\widehat{\theta}_n \sim \operatorname{approx} N\left(\theta \,, \, \frac{1}{I_{\boldsymbol{X}_n}(\theta)} \right) \qquad (\theta \equiv \operatorname{true value}).$$

Estimation of the Fisher Information

If θ is unknown, then so is $I_X(\theta)$.

Two estimates \hat{I} of the Fisher information $I_X(\theta)$ are

$$\widehat{I}_1 = I_X(\widehat{\theta})$$
 and $\widehat{I}_2 = -\frac{\partial^2}{\partial \theta^2} \log f(X \mid \theta) \Big|_{\theta = \widehat{\theta}}$

where $\hat{\theta}$ is the MLE of θ based on the data X.

 \hat{I}_1 is the obvious plug-in estimator. It can be difficult to compute when $I_X(\theta)$ does not have a known closed form.

The estimator \hat{I}_2 is suggested by the formula

$$I_X(\theta) = E\left(-\frac{\partial^2}{\partial\theta^2}\log f(X \mid \theta)\right)$$

It is often easy to compute, and is required in many Newton-Raphson style algorithms for finding the MLE (so that it is already available without extra computation).

The two estimates \hat{I}_1 and \hat{I}_2 are often referred to as the "expected" and "observed" Fisher information, respectively.

As $n \to \infty$, both estimators are consistent (after normalization) for $I_{X_n}(\theta)$ under various regularity conditions.

For example: in the iid case: \hat{I}_1/n , \hat{I}_2/n , and $I_{X_n}(\theta)/n$ all converge to $I(\theta) \equiv I_{X_1}(\theta)$.

Approximate Confidence Intervals for θ

Choose 0 < α < 1 (say, α = 0.05). Let z^* be such that

 $P(-z^* < Z < z^*) = 1 - \alpha$ where $Z \sim N(0, 1)$.

When n is large,

$$\begin{split} &\sqrt{I_X(\theta)}(\hat{\theta} - \theta) \sim \text{approx } N(0, 1) \quad \text{so that} \\ &P\left\{-z^* < \sqrt{I_X(\theta)}(\hat{\theta} - \theta) < z^*\right\} \approx 1 - \alpha \quad \text{or equivalently} \\ &P\left\{\hat{\theta} - z^*\sqrt{\frac{1}{I_X(\theta)}} < \theta < \hat{\theta} + z^*\sqrt{\frac{1}{I_X(\theta)}}\right\} \approx 1 - \alpha \,. \end{split}$$

This approximation continues to hold when $I_X(\theta)$ is replaced by an estimate \hat{I} (either \hat{I}_1 or \hat{I}_2):

$$P\left\{\widehat{\theta}-z^*\sqrt{\frac{1}{\widehat{I}}}<\theta<\widehat{\theta}+z^*\sqrt{\frac{1}{\widehat{I}}}\right\}\approx 1-\alpha.$$

Thus

$$\left(\widehat{ heta}-z^*\sqrt{rac{1}{\widehat{I}}}\,,\,\widehat{ heta}+z^*\sqrt{rac{1}{\widehat{I}}}
ight)$$

is an approximate $1 - \alpha$ confidence interval for θ . (Here $\hat{\theta}$ is the MLE and \hat{I} an estimate of the Fisher information.)

 $\frac{(\operatorname{ramer}-\operatorname{Rao}\ \operatorname{Inequality}\ (X \sim P_{\theta}, \theta \in \Theta \subset \mathbb{R})}{\operatorname{If}\ f(\chi|\theta)\ \text{is a regular one-parameter family,}}$ $E_{\theta}W(\chi) = \mathcal{T}(\theta)\ \text{for all }\theta, \text{ and}$ $\mathcal{T}(\theta)\ \text{is differentiable,}$ $\operatorname{then}\ \operatorname{Var}_{\theta}(W(\chi)) \geq \frac{(\mathcal{T}'(\theta))^2}{I_{\chi}(\theta)}.$

Proof:

 $\frac{\text{Preliminary Facts}}{\left(Ov(X,Y) \right)^{2}} \leq (Var X)(Var Y).$ This is a special case of the Cauchy-Schwarz inequality. It is better known to statisticians as $\rho^{2} \leq 1$ where $\rho = \frac{Cov(X,Y)}{Var X \cdot Var Y}$ is the Correlation between X and Y. $\frac{Ov(X,Y) = E \times Y}{Over Y} = E \times Y = O = E \times Y = 0$ This follows from the well known formula $Cov(X,Y) = E \times Y - (E \times XEY).$

Body of Proof: From (A) we have $\left[\operatorname{Cov}(W(X), \frac{1}{2}\log f(X|\theta))\right]^{2} \leq \left(\operatorname{Var}W(X)\right)\left(\operatorname{Var}(\frac{1}{2}\log f(X|\theta))\right)$ $I_{X}(\theta)$ From B we have $Cov(W(X), \frac{\partial}{\partial \theta} \log f(X|\theta)) = E[W(X), \frac{\partial}{\partial \theta} \log f(X|\theta)]$ (since $E \stackrel{\partial}{\partial \theta} \log f(X|\theta) = 0$) $= \left(W(\chi) \left(\frac{\partial}{\partial \Theta} \log f(\chi | \Theta) \right) f(\chi | \Theta) d\chi \right)$ $\Rightarrow = \frac{1}{9+1}$ $= \int W(\chi) \frac{\partial f(\chi, \theta)}{\partial \theta} d\chi$ $= \frac{\partial}{\partial \theta} \int W(\chi) f(\chi|\theta) d\chi$ (since $f(\chi|\theta)$ is a regular family) $=\frac{\partial}{\partial\Theta}E_{\Theta}W(X) = \mathcal{T}'(\Theta)$. Thus $[\gamma'(\Theta)]^2 \leq (Var_{\Theta} W(X)) I_X(\Theta)$.

Addendum:

Equality in (A) is achieved iff Y = a X + b for some constants a, b. Moreover, if EY = 0, then E(aX+b)=0forces b = -a EX so that Y = a(X - EX) for some constant a. Applying this to the proof of CRLB with X = W(X), $Y = \frac{\partial}{\partial \Theta} \log f(X|\Theta)$ tells us that $\operatorname{Var}_{\Theta} W(X) = (\gamma'(\Theta))^2$ $I_{\times}(\Theta)$ $iff \quad \frac{\partial}{\partial \Theta} \log f(X|\Theta) = A(\Theta) \left[W(X) - \gamma(\Theta) \right] \quad (*)$ for some function $a(\theta)$. Note: (*) is true only when

 $f(\underline{x}|\theta)$ is a lpef and $W(\underline{x}) = cT(\underline{x}) + d$ for some c, dwhere $T(\underline{x})$ is the natural suff. stat of the lpef. Asymptotic Efficiency

Given: A sequence of estimators

$$W_n = W_n(X_1, X_2, ..., X_n)$$
.
 X_n



Example: Observe X_1, \ldots, X_n iid Poisson (λ) .

• Estimation of $\tau(\lambda) = \lambda$

 $E\bar{X} = \lambda$

Does \bar{X} achieve the CRLB? Yes!

$$Var(\bar{X}) = \frac{Var(X_1)}{n} = \frac{\lambda}{n}$$
$$CRLB = \frac{(\tau'(\lambda))^2}{I_X(\lambda)} = \frac{1}{n/\lambda} = \frac{\lambda}{n}$$

Alternative: Check condition for exact attainment of CRLB.

$$\frac{\partial}{\partial \lambda} \log f(\boldsymbol{X} | \lambda) = \sum_{i=1}^{n} \frac{\partial}{\partial \lambda} \log f(X_i | \lambda) = \sum_{i} \left(\frac{X_i}{\lambda} - 1 \right)$$
$$= \frac{n}{\lambda} \left(\bar{X} - \lambda \right)$$

Note: Since \overline{X} attains the CRLB (for all λ), it must be the best unbiased estimator of λ .

Showing that an estimator attains the CRLB is one way to show it is best unbiased. (But see later remark.)

• Estimation of $au(\lambda) = \lambda^2$

Define
$$W = \frac{T(T-1)}{n^2}$$
 where $T = \sum_{i=1}^n X_i$.

 $EW = \lambda^2$ (see calculations below) and W is a function of the CSS T. Thus W is best unbiased for λ^2 .

Does W achieve the CRLB? No!!!

$$CRLB = \frac{(\tau'(\lambda))^2}{I_X(\lambda)} = \frac{(2\lambda)^2}{n/\lambda} = \frac{4\lambda^3}{n}$$
$$Var(W) = \frac{4\lambda^3}{n} + \frac{2\lambda^2}{n^2} \text{ (see calculations below)}$$

Alternative: Show condition for achievement of CRLB **fails**. As shown earlier:

$$\frac{\partial}{\partial \lambda} \log f(\mathbf{X} | \lambda) = \sum_{i} \left(\frac{X_i}{\lambda} - 1 \right) = \frac{T}{\lambda} - n$$

The CRLB is attained **iff** there exists $a(\lambda)$ such that

$$\frac{T}{\lambda} - n = a(\lambda) \left(\frac{T(T-1)}{n^2} - \lambda^2 \right) \,.$$

But the left side is linear in T, and the right is quadratic in T, so that no multiplier $a(\lambda)$ can make them equal for all possible values of T = 0, 1, 2, ...

Remark: This situation is not unusual. The best unbiased estimator often fails to achieve the CRLB.

But W is asymptotically efficient:

$$\lim_{n \to \infty} \frac{\operatorname{Var}(W)}{CRLB} = \lim_{n \to \infty} \frac{\frac{4\lambda^3}{n} + \frac{2\lambda^2}{n^2}}{\frac{4\lambda^3}{n}} = \lim_{n \to \infty} \left(1 + \frac{1}{2n\lambda}\right) = 1$$

• Calculations:

Suppose $Y \sim \text{Poisson}(\xi)$.

The factorial moments of the Poisson follow a simple pattern:

$$EY = = \xi EY(Y - 1) = \xi^{2} EY(Y - 1)(Y - 2) = \xi^{3} EY(Y - 1)(Y - 2)(Y - 3) = \xi^{4}$$
etc.

Proof of one case:

$$EY(Y-1)(Y-2) = \sum_{i=0}^{\infty} i(i-1)(i-2)\frac{\xi^{i}e^{-\xi}}{i!}$$
$$= \xi^{3} \sum_{i=3}^{\infty} \frac{\xi^{i-3}e^{-\xi}}{(i-3)!} = \xi^{3} \sum_{j=0}^{\infty} \frac{\xi^{j}e^{-\xi}}{j!} = \xi^{3}$$

From the factorial moments, we can calculate everything else. For example:

$$Var(Y(Y-1)) = E[\{Y(Y-1)\}^2] - [EY(Y-1)]^2$$

= $E[Y^2(Y-1)^2] - [\xi^2]^2$
= $E[\langle Y \rangle_4 + 4 \langle Y \rangle_3 + 2 \langle Y \rangle_2] - \xi^4$
= $[\xi^4 + 4\xi^3 + 2\xi^2] - \xi^4 = 4\xi^3 + 2\xi^2$
where $\langle Y \rangle_k \equiv Y(Y-1) \cdots (Y-k+1)$.

In our case $T\sim {\rm Poisson}(n\lambda)$ so that substituting $\xi=n\lambda$ in the above results leads to

$$ET(T-1) = (n\lambda)^2 = n^2\lambda^2$$

Var[T(T-1)] = 4(n\lambda)^3 + 2(n\lambda)^2 = 4n^3\lambda^3 + 2n^2\lambda^2
so that $W = T(T-1)/n^2$ satisfies:

$$EW = \lambda^{2}$$

Var(W) = $\frac{4\lambda^{3}}{n} + \frac{2\lambda^{2}}{n^{2}}$

An asymptotically inefficient estimator
Example: Let
$$X_{1,}X_{2,}...,X_{n}$$
 be iid
with $pdf f(x|\alpha) = \frac{x^{\alpha-1}e^{-x}}{\Gamma(\alpha)}$ for $x > 0$.
For this pdf , $Ex = Var X = \alpha$.
Clearly $E \overline{x} = \alpha$. Thus \overline{x} is MOM
estimator of α .
Is it asymptotically efficient? No. (Verified
below.)
Note: This is Ipef with natural sufficient
statistic $T = \sum_{i=1}^{n} \log X_i$. Since T is complete,
 $E(\overline{x}|T)$ is the UMVUE of α . Since \overline{x}
is not a function of T , we know
 $Var(\overline{x}) > Var[E(\overline{x}|T)]$.
But $Var[E(\overline{x}|T)] \ge CRLB$. Thus, without
calculation, we know that \overline{x} cannot achieve
the CRLB for any value of n . We now show
it does not achieve it asymptotically either.

$$Var \overline{X} = \frac{Var X_{l}}{n} = \frac{\alpha}{n} .$$

$$I_{X_{n}}(\alpha) = n I_{X_{l}}(\alpha)$$

$$= n \left[\frac{\Gamma''(\alpha) \Gamma(\alpha) - (\Gamma'(\alpha))^{2}}{(\Gamma(\alpha))^{2}} \right]$$
by a routine calculation
$$CRLB = \frac{1}{n I_{X_{l}}(\alpha)} .$$

Thus

$$\frac{\operatorname{Var}\overline{X}}{\operatorname{CRLB}} = \propto I_{X_{1}}(\alpha) \quad \text{which does } \underline{\operatorname{not}} \\ \text{depend on n.} \\ \text{Since } \overline{X} \text{ does } \underline{\operatorname{not}} \text{ achieve } \operatorname{CRLB} \text{ for any } n, \\ \text{we know } \alpha I_{X_{1}}(\alpha) > 1. \text{ Thus} \\ \lim_{n \to \infty} \operatorname{Var}\overline{X} = \alpha I_{X_{1}}(\alpha) > 1 \\ \operatorname{not} \operatorname{var}\overline{X} = \alpha I_{X_{1}}(\alpha) > 1 \\ \text{so that } \overline{X} \text{ is not asymp. eff.} \\ \end{array}$$

