

Outline

- 1) Log-likelihood and score
- 2) Fisher information
- 3) Cramér-Rao Lower Bound
- 4) Hammersley-Chapman-Robbins Ineq.

Log-likelihood, score

Assume \mathcal{P} has densities p_θ wrt μ , $\Theta \subseteq \mathbb{R}^d$

Common support: $\{x: p_\theta(x) > 0\}$ same $\forall \theta$

Recall $l(\theta; x) = \log p_\theta(x)$,

Thought of as random function of θ

Def The score is $\nabla l(\theta; x)$; plays a key role in many areas of statistics, esp. asymptotics.

Can think of as "local sufficient statistic":

$$\begin{aligned} p_{\theta_0 + \eta}(x) &= e^{l(\theta_0 + \eta; x)} \\ &\approx e^{\eta' \nabla l(\theta_0; x)} p_{\theta_0}(x) \quad \text{for } \eta \approx 0 \end{aligned}$$

Differential identities: (assuming enough regularity)

$$1 = \int_{\mathcal{X}} e^{l(\theta; x)} d\mu(x)$$

$$\frac{\partial}{\partial \theta_j} \Rightarrow 0 = \int \frac{\partial}{\partial \theta_j} l(\theta; x) e^{l(\theta; x)} d\mu(x)$$

$$\Rightarrow \mathbb{E}_\theta [\nabla l(\theta; x)] = 0$$

↑
only true if these are the same value of θ !

$$\begin{aligned} \frac{\partial}{\partial \theta_k} \Rightarrow 0 &= \int \left(\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_k} + \frac{\partial \ell}{\partial \theta_i} \cdot \frac{\partial \ell}{\partial \theta_k} \right) e^\ell d\mu \\ &= \mathbb{E}_\theta \left[\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_k} \right] + \mathbb{E}_\theta \left[\frac{\partial \ell}{\partial \theta_i} \frac{\partial \ell}{\partial \theta_k} \right] \end{aligned}$$

$$\Rightarrow J(\theta) = \text{Var}_\theta [\nabla \ell(\theta; x)] = \mathbb{E}_\theta [-\nabla^2 \ell(\theta; x)]$$

← same θ
← same θ

Called "Fisher Information"

[It is possible to extend this definition to certain cases where ℓ is not even differentiable, e.g. Laplace location family, but for our purposes we can just assume "sufficient regularity."]

Try with another statistic $\delta(x)$, let $g(\theta) = \mathbb{E}_\theta[\delta(x)]$ ("unbiased estimator")

$$g(\theta) = \int \delta e^\ell d\mu$$

$$\begin{aligned} \Rightarrow \nabla g(\theta) &= \int \delta \nabla \ell e^\ell d\mu = \mathbb{E}_\theta [\delta(x) \nabla \ell(\theta; x)] \\ &= \text{Cov}_\theta(\delta(x), \nabla \ell(\theta; x)) \end{aligned}$$

why?
 Since $\mathbb{E} \nabla \ell = 0$

Combining these results with Cauchy-Schwarz gives us the Cramér-Rao Lower Bound or Information Lower Bound:

1-param: $\text{Var}_\theta(\delta) \cdot \text{Var}_\theta(\dot{\ell}(\theta; X)) \geq \text{Cov}_\theta(\delta, \dot{\ell}(\theta; X))^2$
 $\Rightarrow \text{Var}_\theta(\delta) = \dot{g}(\theta)^2 / J(\theta)$

$\theta \in \mathbb{R}^d, g(\theta) \in \mathbb{R}$: $\text{Var}_\theta(\delta) \geq \nabla g(\theta)' J(\theta)^{-1} \nabla g(\theta)$

Interp: If $g(\theta)$ is estimand, no unbiased estimator can have smaller variance than $\nabla g(\theta)' J(\theta)^{-1} \nabla g(\theta)$

Ex.: (i.i.d. sample)

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_\theta^{(1)}(x) \quad \theta \in \Theta$$

$$X \sim p_\theta(x) = \prod_i p_\theta^{(1)}(x_i)$$

$$\text{Let } \ell_1(\theta; x_i) = \log p_\theta^{(1)}(x_i)$$

$$\ell(\theta; x) = \sum_i \ell_1(\theta; x_i)$$

$$J(\theta) = \text{Var}_\theta(\nabla \ell(\theta; x))$$

$$= \text{Var}_\theta(\sum_i \nabla \ell_1(\theta; x_i))$$

$$= n J_1(\theta)$$

where $J_1(\theta)$ is Fisher info in single observation

\Rightarrow Lower bound scales like n^{-1} (SD $\asymp n^{-1/2}$ for "regular" families)

Efficiency

CRLB is not nec. attainable.

We define the efficiency of an unbiased estimator as:

$$\text{eff}_{\theta}(\delta) = \frac{\text{CRLB}}{\text{Var}_{\theta}(\delta)} \quad (= \frac{1/J(\theta)}{\text{Var}_{\theta}(\delta)} \text{ if } g(\theta) = \theta \in \mathbb{R})$$

$$\text{eff}_{\theta}(\delta) \leq 1$$

We say $\delta(x)$ is efficient if $\text{eff}_{\theta}(\delta) = 1 \quad \forall \theta$

Depends on $\text{Corr}_{\theta}(\delta(x), \nabla \ell(\theta; x))$:

$$\begin{aligned} \text{eff}_{\theta}(\delta) &= \frac{\text{Cov}_{\theta}^2(\delta(x), \dot{\ell}(\theta; x))}{\text{Var}_{\theta}(\delta) \cdot \text{Var}_{\theta}(\dot{\ell}(\theta))} \\ &= \text{Corr}_{\theta}^2(\delta, \dot{\ell}(\theta)) \end{aligned}$$

$$\leq 1$$

$$\delta(x) \text{ is efficient} \Leftrightarrow \text{Corr}_{\theta}^2(\delta, \dot{\ell}(\theta)) = 1 \quad \forall \theta$$

Rarely achieved in finite samples but we can approach it asymptotically as $n \rightarrow \infty$

Hammersley - Chapman - Robbins Ineq.

CRLB requires differentiation under integral

Can make more general statement if we replace $\nabla \ell(\theta; x)$ with finite-difference:

$$\frac{p_{\theta+\varepsilon}(x)}{p_{\theta}(x)} - 1 = e^{\ell(\theta+\varepsilon; x) - \ell(\theta; x)} - 1$$

($\approx \varepsilon' \nabla \ell(\theta; x)$ small ε)

$$\mathbb{E}_{\theta} \left[\frac{p_{\theta+\varepsilon}}{p_{\theta}} - 1 \right] = \int \left(\frac{p_{\theta+\varepsilon}}{p_{\theta}} - 1 \right) p_{\theta} d\mu = 1 - 1 = 0$$

(assuming common support, or $p_{\theta+\varepsilon} \ll p_{\theta}$)

$$\text{Cov}_{\theta} \left(\delta, \frac{p_{\theta+\varepsilon}}{p_{\theta}} - 1 \right) = \int \delta \left(\frac{p_{\theta+\varepsilon}}{p_{\theta}} - 1 \right) p_{\theta} d\mu$$

$$= \mathbb{E}_{\theta+\varepsilon}(\delta) - \mathbb{E}_{\theta}(\delta)$$

$$= g(\theta+\varepsilon) - g(\theta)$$

$$\Rightarrow \text{Var}_{\theta}(\delta) \geq \frac{(g(\theta+\varepsilon) - g(\theta))^2}{\mathbb{E}_{\theta} \left[\left(\frac{p_{\theta+\varepsilon}}{p_{\theta}} - 1 \right)^2 \right]}$$

CRLB follows from $\varepsilon \rightarrow 0$, but \sup_{ε} gives better bound

Ex. Exponential Families

$$p_{\eta}(x) = e^{\eta' T(x) - A(\eta)} h(x)$$

$$\ell(\eta; x) = \eta' T(x) - A(\eta) + \log h(x)$$

$$\begin{aligned}\nabla \ell(\eta; x) &= T(x) - \nabla A(\eta) \\ &= T(x) - \mathbb{E}_{\eta} T(x)\end{aligned}$$

$$\text{Var}_{\eta}(\nabla \ell(\eta)) = \text{Var}_{\eta}(T(x)) = \nabla^2 A(\eta)$$

$$\nabla^2 \ell(\eta; x) = -\nabla^2 A(\eta)$$

$$\mathbb{E}_{\eta}[-\nabla^2 \ell(\eta; x)] = \nabla^2 A(\eta) \quad \checkmark$$

So any unbiased est. of η has

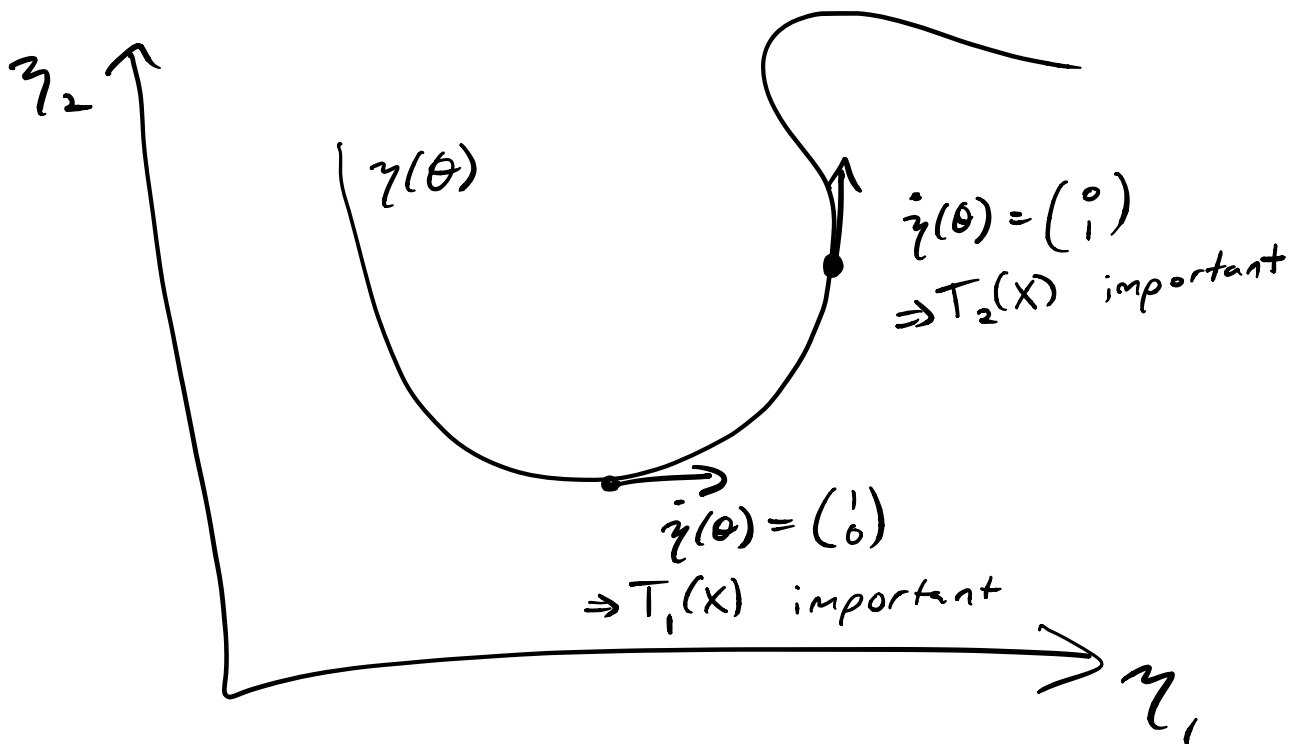
$$\text{Var}_{\eta}(\delta) \geq \nabla^2 A(\eta)^{-1}$$

Curved family: $p_{\theta}(x) = e^{\eta(\theta)'T(x) - B(\theta)} h(x)$, $\theta \in \mathbb{R}$
 $B(\theta) = A(\eta(\theta))$

$$l(\theta; x) = \eta(\theta)'T(x) - B(\theta) + \log h(x)$$

$$\begin{aligned} \dot{l}(\theta; x) &= \dot{\eta}(\theta)'T(x) - \dot{\eta}(\theta)'\nabla_{\eta}A(\eta(\theta)) \\ &= \dot{\eta}(\theta)'(T(x) - \nabla_{\eta}A(\eta(\theta))) \\ &= \dot{\eta}(\theta)'(T(x) - E_{\theta}T(x)) \end{aligned}$$

$\Rightarrow \dot{\eta}(\theta)'T(x)$ is "locally complete suff. stat."



Doubts about unbiasedness

The UMVUE might be very inefficient, or inadmissible, or just dumb, in cases where another approach makes much more sense.

Ex. $X \sim \text{Bin}(1000, \theta)$

Estimate $g(\theta) = \mathbb{P}_\theta(X \geq 500)$

UMVUE is $1\{X \geq 500\}$ (why?)

$\Rightarrow X = 500$? Conclude $g(\theta) = 100\%$
 $X = 499$? Conclude $g(\theta) = 0\%$

This is not epistemically reasonable!!

Could do much better with e.g. MLE or a Bayes estimator.

In fact, our theorem should make us suspicious of UMVUE's: every idiotic function of T is a UMVUE (of its own expectation)

Gaussian Sequence Model

$X_i \stackrel{\text{iid}}{\sim} N(\mu_i, 1) \quad i=1, \dots, d \quad \text{indep.}$

or $X \sim N_d(\mu, I_d) \quad \mu \in \mathbb{R}^d, \text{ estimate } \varphi^2 = \|\mu\|^2$

X is complete sufficient

$$\begin{aligned} \mathbb{E}_\mu \|X\|^2 &= \mathbb{E}_0 [\|\mu + X\|^2] \\ &= \|\mu\|^2 + \mathbb{E}_0 \|X\|^2 + 2\cancel{\mathbb{E}_0 [\mu^T X]}^0 \\ &= \|\mu\|^2 + d \end{aligned}$$

$$\Rightarrow \delta(x) = \|x\|^2 - d$$

If $\mu = 0$, $\delta(x) < 0$ about half the time!

$$(\|x\|^2 - d)_+ = \max(0, \|x\|^2 - d)$$

strictly dominates UMVU

Gets worse: Ex 4.7 in Keener

$X \sim \text{Truncated Poisson}(\theta)$

$$p_{\theta}(x) = \frac{\theta^x e^{-\theta}}{x! (1 - e^{-\theta})} \quad \begin{array}{l} x = 1, 2, \dots \\ \theta > 0 \end{array}$$

Estimate $g(\theta) = e^{-\theta}$ (mass lost to truncation)

Keener shows UMVUE is $\delta(X) = (-1)^{X+1}$

$$\frac{e^{-\theta}}{1 - e^{-\theta}} \left(\theta - \frac{\theta^2}{2} + \frac{\theta^3}{3!} - \dots \right) = \frac{e^{-\theta}}{1 - e^{-\theta}} \left[1 - \left(1 + (-\theta) + \frac{(-\theta)^2}{2!} + \dots \right) \right]$$

[Idiotic but we cannot improve using any unbiased estimator]

Sometimes insisting on unbiasedness leads us to absurd results.

Unbiasedness has bad reputation, but other methods have their problems too.]