

# Fisher information and Cramér-Rao Lower Bound

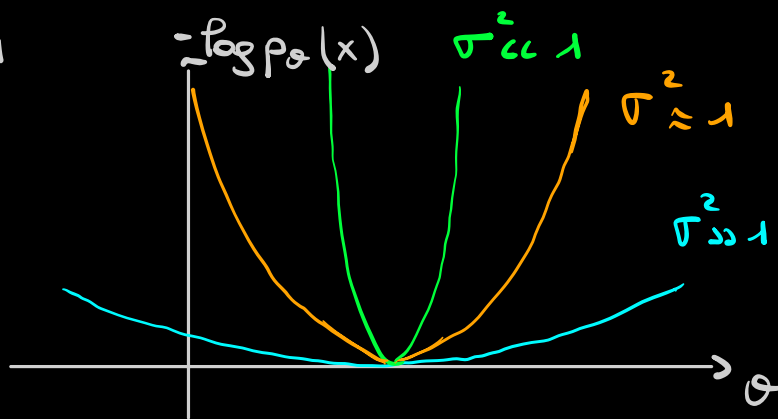
Example (Gaussians with given variance).  $X_j \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, \sigma^2)$

$\Theta = \mathbb{R}$ ,  $\sigma^2$  is known

Consider the estimator  $\hat{\theta} = \frac{1}{n} \sum X_j \rightarrow \text{UMVU!}$

$$V_{\theta}(\hat{\theta}) = \frac{1}{n} \cdot \sigma^2 = \frac{1}{n} \left( \frac{1}{\sigma^2} \right)^{-1}$$

Idea: the curvature of the (log) likelihood around the estimator



$$-\frac{\partial^2}{\partial \theta^2} \log p_{\theta}(x) = \frac{1}{\sigma^2}$$

indicates how "localized" the parameter  $\theta$  is. ( $\Rightarrow$  how much variance/uncertainty our estimate might have).

Averaging over  $x \sim P_{\theta}$  we obtain

$$I(\theta) = -\mathbb{E}_{\theta} \left( \frac{\partial^2}{\partial \theta^2} \log p_{\theta}(x) \right)$$

We expect (the inverse of) this quantity to encode how much variance we may expect our estimator (at least).

Idea: We want to capture the variability of a parametric model around  $\theta$ .

Assumption 1: Let  $(X, \mathcal{F}, \{P_\theta\}_{\theta \in \Theta})$  be a dominated statistical model with densities  $p_\theta(x)$  wrt. a dominating meas.  $\mu$ .

We assume that  $\Theta \subseteq \mathbb{R}^k$  is open and that  $\mathbb{E}$  denotes expectation wrt  $\mu$ .

$$p_\theta(x) > 0 \quad \forall x \in X \quad \forall \theta \in \Theta$$

to take logs

Ex: Gaussians:  $\Theta = (\mu, \sigma^2)$   $X_j \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$\Rightarrow p_\theta(x) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_j - \mu)^2}{2\sigma^2}\right) > 0 \quad \forall x \in \mathbb{R}^n$$

Assumption 2: It holds that for all mble  $f$  with  $\mathbb{E}(f^2(x)) < \infty$

$$\nabla_\theta \mathbb{E}(p_\theta(x) f(x)) = \mathbb{E}(\nabla_\theta p_\theta(x) f(x))$$

(for example:  $\nabla_\theta p_\theta(x)$  exists and  $\mathbb{E}(\|\nabla_\theta p_\theta(x)\|^2) < \infty \quad \forall \theta \in \Theta$ )

Lemma: Under A1, A2 we have

$$\mathbb{E}_\theta(\nabla_\theta \log p_\theta(x)) = 0 \quad \forall \theta \in \Theta$$

$$\text{Proof: } \mathbb{E}_\theta(\nabla_\theta \log p_\theta(x)) = \mathbb{E}_\theta\left(\frac{1}{p_\theta(x)} \cdot \nabla_\theta p_\theta(x)\right) =$$

$$= \int \frac{1}{\cancel{p_\theta(x)}} \nabla_\theta p_\theta(x) \cdot \cancel{p_\theta(x)} \mu(dx)$$

$$= \int \nabla_\theta p_\theta(x) \mu(dx) = \nabla_\theta \int p_\theta(x) \mu(dx)$$

$$= \nabla_\theta (1) = 0$$

□

Assumption 3:  $\mathbb{E}_\theta (\| \nabla_\theta \log p_\theta(x) \|^2) < \infty \quad \forall \theta \in \Theta$

Def: Let A1,2,3 hold. We call

$$I(\theta)_{ij} = \mathbb{E}_\theta \left( \frac{\partial}{\partial \theta_i} \log p_\theta(x) \cdot \frac{\partial}{\partial \theta_j} \log p_\theta(x) \right) \quad i, j \in \{1, \dots, k\}$$

the Fisher information matrix of the dominated model.

$\{p_\theta(x)\}_{\theta \in \Theta}$ . Further, we call  $s_\theta(x) = \nabla_\theta \log p_\theta(x)$  the score function.

Note: Alternatively, using chain rule, we can write

$$I(\theta)_{ij} = \mathbb{E}_\theta \left( \frac{1}{p_\theta(x)^2} \frac{\partial}{\partial \theta_i} p_\theta(x) \frac{\partial}{\partial \theta_j} p_\theta(x) \right)$$

Note: The matrix  $I(\theta)_{ij}$  is symmetric positive semidefinite since it is the covariance matrix  $C$  of the score:

$$\begin{aligned} \langle C \cdot \sigma, \sigma \rangle &= \sum_{i,j=1}^k \text{Cov}(Y_i, Y_j) \sigma_i \sigma_j \\ &= \mathbb{E} \left( \sum_{i,j=1}^k \sigma_i (Y_i - \mathbb{E} Y_i) (Y_j - \mathbb{E} Y_j) \sigma_j \right) \\ &= \mathbb{E} \left( \left( \sum_{i=1}^k \sigma_i (Y_i - \mathbb{E} Y_i) \right)^2 \right) \geq 0 \end{aligned}$$

under A1.3

✓

Lemma: For every  $i, j \in \{1, \dots, k\}$   $\mathbb{E}_\theta \left( \frac{1}{p_\theta(x)} \frac{\partial^2}{\partial \theta_i \partial \theta_j} p_\theta(x) \right) = 0$

Proof:  $\mathbb{E}_\theta \left( \frac{1}{p_\theta(x)} \frac{\partial^2}{\partial \theta_i \partial \theta_j} p_\theta(x) \right) = \int \frac{1}{p_\theta(x)} \frac{\partial^2}{\partial \theta_i \partial \theta_j} p_\theta(x) p_\theta(x) \mu(dx)$

$$= \frac{\partial^2}{\partial \theta_i \partial \theta_j} \int p_\theta(x) \mu(dx) = \frac{\partial^2}{\partial \theta_i \partial \theta_j} 1 = 0 \quad \square$$

Proposition: Under A1-3

$$I(\theta)_{ij} = -\mathbb{E}_\theta \left( \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p_\theta(x) \right).$$

Proof: Use the notation  $d_{ij}^2 = \frac{\partial^2}{\partial \theta_i \partial \theta_j}$ .

$$\begin{aligned} d_{ij}^2 \log p_\theta(x) &= d_i \left( d_j \log p_\theta(x) \right) = d_i \left( \frac{1}{p_\theta(x)} d_j p_\theta(x) \right) \\ &= \frac{1}{p_\theta(x)} d_{ij}^2 p_\theta(x) - \frac{1}{p_\theta(x)^2} d_i p_\theta(x) d_j p_\theta(x) \end{aligned}$$

$$\begin{aligned} \Rightarrow \mathbb{E}_\theta \left( d_{ij}^2 \log p_\theta(x) \right) &= \mathbb{E}_\theta \left( \frac{1}{p_\theta(x)} d_{ij}^2 p_\theta(x) \right) - \mathbb{E}_\theta \left( \frac{1}{p_\theta(x)^2} d_i p_\theta(x) d_j p_\theta(x) \right) \\ &= -I(\theta)_{ij} \quad \square \end{aligned}$$

Ex (Gaussian)  $\underline{X} = (X_1, \dots, X_n)$   $X_j \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$   $\theta = \mu$ .

$$\log p_\theta(x) = \sum_{j=1}^n -\frac{1}{2} \log 2\pi\sigma^2 - (x_j - \mu)^2 \frac{1}{2\sigma^2}$$

$$\Rightarrow \partial_\theta \log p_\theta(x) = + \sum_{j=1}^n \frac{x_j - \mu}{\sigma^2} \Rightarrow -\partial_\theta^2 \log p_\theta(x) = n \cdot \frac{1}{\sigma^2}$$

$$\Rightarrow \mathbb{E}_\theta \left( \left( \partial_\theta \log p_\theta(x) \right)^2 \right) = \frac{1}{\sigma^4} \sum_{j=1}^n \mathbb{E}_\theta \left( (x_j - \mu)^2 \right) = n \frac{1}{\sigma^2}$$

Now consider the statistical model  $(\mathcal{X}^n, \mathcal{F}^{\otimes n}, \{\mathbb{P}_\theta^{\otimes n}\}_{\theta \in \Theta})$ .

Theorem: Let A1-3 hold and consider the  $\mu_1 \otimes \mu_2$  dominated stat. model  $(\mathcal{X}_1 \times \mathcal{X}_2, \mathcal{F}_1 \otimes \mathcal{F}_2, (\mathbb{P}_\theta^{(1)} \otimes \mathbb{P}_\theta^{(2)})_{\theta \in \Theta})$  then

$$I^{(1+2)}(\theta) = I^{(1)}(\theta) + I^{(2)}(\theta)$$

$$\text{where } I_c(\theta) = \mathbb{E}_\theta^{(c)} \left( \frac{\partial}{\partial \theta_i} \log p_\theta^{(c)}(x) \frac{\partial}{\partial \theta_j} \log p_\theta^{(c)}(x) \right).$$

Corollary: for  $(\mathcal{X}^n, \mathcal{F}^{\otimes n}, (\mathbb{P}_\theta^{\otimes n})_\theta)$   $I_n(\theta) = n I_1(\theta)$

$\Rightarrow$  compute only  $I_1(\theta)$   $\searrow$  fisher of  $\mathbb{P}_\theta$   
 $\searrow$  fisher of  $\mathbb{P}_\theta^{\otimes n}$

Ex (Gaussian unknown mean)  $X_j \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$   $\theta = \mu$ .

$$I_1(\theta) = -\mathbb{E}_\theta \left( \frac{\partial^2}{\partial \theta^2} \log p_\theta(x) \right) = \frac{1}{\sigma^2} \Rightarrow I_n(\theta) = n \cdot \frac{1}{\sigma^2}$$

Ex (Gaussian unknown variance)  $X_j \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$   $\theta = \sigma^2$

$$\log p_\theta(x) = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{(x-\mu)^2}{2\sigma^2} \Rightarrow s_\theta(x) = -\frac{1}{2\sigma^2} + \frac{(x-\mu)^2}{2\sigma^4}$$

$$I_1(\theta) = -\mathbb{E}_\theta \left( \frac{\partial^2}{\partial \theta^2} \log p_\theta(x) \right) = -\mathbb{E}_\theta \left( \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} (x-\mu)^2 \right)$$

$$= -\left( \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} \sigma^2 \right) = \frac{1}{2\sigma^4} \Rightarrow I_n(\theta) = \frac{n}{2\sigma^4}$$

Exercise:  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ ,  $\theta = (\mu, \sigma^2)$   $I(\theta) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}$

Lemma: Let  $X \perp Y$  random vectors, then  $\text{Cov}(X+Y) = \text{Cov}(X) + \text{Cov}(Y)$

Proof of theorem:

$$\begin{aligned} I(\vartheta) &= \text{Cov}(\nabla_{\vartheta} \log p_{\vartheta}(x)) = \text{Cov}(\nabla_{\vartheta} (\log p_{\vartheta}^{(1)}(x) + \log p_{\vartheta}^{(2)}(x))) \\ &\stackrel{\text{Lemma}}{=} \text{Cov}(\nabla \log p_{\vartheta}^{(1)}(x)) + \text{Cov}(\nabla \log p_{\vartheta}^{(2)}(x)) = I_1(x) + I_2(x) \quad \square \end{aligned}$$

Theorem (Cramer-Rao) Let  $T$  be an unbiased estimator of  $g(\vartheta)$  with  $\mathbb{E}_{\vartheta}(T^2) < \infty \forall \vartheta \in \Theta$ . Then, assuming  $I$  is finite  $\forall \vartheta$ ,

$$W_{\vartheta}(T) \geq \nabla_{\vartheta} g(\vartheta)^T (I(\vartheta))^{-1} \nabla_{\vartheta} g(\vartheta)$$

Def: We call  $\nabla_{\vartheta} g(\vartheta)^T (I(\vartheta))^{-1} \nabla_{\vartheta} g(\vartheta)$  the Cramer-Rao lower bound for estimating  $g(\vartheta)$

Note: in 1d  $W_{\vartheta}(T) \geq \frac{(\partial_{\vartheta} g(\vartheta))^2}{I(\vartheta)}$

if  $g(\vartheta) = \vartheta$   $W_{\vartheta}(T) \geq I(\vartheta)^{-1}$

Proof: We have

$$\begin{aligned} \nabla_{\vartheta} g(\vartheta) &= \nabla_{\vartheta} \mathbb{E}_{\vartheta}(T) = \nabla_{\vartheta} \int T(x) p_{\vartheta}(x) \mu(dx) \\ &= \int T(x) \nabla_{\vartheta} p_{\vartheta}(x) \frac{1}{p_{\vartheta}(x)} p_{\vartheta}(x) \mu(dx) \end{aligned}$$

$$= \mathbb{E}_{\vartheta}(T(x) \nabla_{\vartheta} \log p_{\vartheta}(x)) - \mathbb{E}_{\vartheta}(T(x)) \mathbb{E}_{\vartheta}(\nabla_{\vartheta} \log p_{\vartheta}(x))$$

$$= \mathbb{E}_{\vartheta}(\nabla_{\vartheta} \log p_{\vartheta}(x) (T(x) - \mathbb{E}_{\vartheta}(T(x))))$$

$$\Rightarrow \langle x, \nabla_{\theta} \mathbb{E}_{\theta}(T) \rangle^2 = \mathbb{E}_{\theta}(\langle x, \nabla_{\theta} \log p_{\theta}(x) (T(x) - \mathbb{E}_{\theta}(T(x))) \rangle)$$

$$\stackrel{c.s.}{\leq} \mathbb{V}_{\theta}(T) \mathbb{E}(\langle x, \nabla_{\theta} \log p_{\theta}(x) \rangle^2) =$$

$$= \mathbb{V}_{\theta}(T) x^T I(\theta) x.$$

→ done in the 1d case: simplify  $x$  and get:  $\mathbb{V}(T) = \frac{(\partial_{\theta} g(\theta))^2}{I(\theta)}$ .

$$\Rightarrow \mathbb{V}_{\theta}(T) = \frac{\langle x, \nabla_{\theta} \mathbb{E}_{\theta}(T) \rangle^2}{x^T I(\theta) x}$$

optimize,

$$\xrightarrow{\text{over } x} \mathbb{V}_{\theta}(T) = \nabla_{\theta} g(\theta)^T I(\theta)^{-1} \nabla_{\theta} g(\theta)$$

Lemma: Let  $a \in \mathbb{R}^n$ ,  $A$  symmetric positive definite, then

$$\sup_{x \in \mathbb{R}^n} \frac{\langle x, a \rangle^2}{\langle Ax, x \rangle} = \langle A^{-1} a, a \rangle$$

PP: Let  $BB = A$ ,  $B$  symmetric  $\rightarrow$

$$\frac{\langle x, a \rangle}{\langle Ax, x \rangle} = \frac{\langle B^{-1} y, a \rangle}{\langle AB^{-1} y, B^{-1} y \rangle} = \frac{\langle y, B^{-1} a \rangle}{\|y\|^2} = \left( \frac{\langle y, b \rangle}{\|y\|} \right)^2 = \|b\|^2 = \langle (B^2)^{-1} a, a \rangle = \langle A^{-1} a, a \rangle \quad \square$$

→ maximized in  $y = b$

Ex (Poisson):  $\mathbf{X} = (X_1, \dots, X_n)$   $X_j \stackrel{iid}{\sim} \text{Pois}(\theta)$   $\theta > 0$

$$\log p_{\theta}(x) = -n\theta + x \log \theta - \log x!$$

$$\Rightarrow s_{\theta}(x) = \nabla_{\theta} \log p_{\theta}(x) = -1 + \frac{x}{\theta}$$

$$\Rightarrow I_1(\theta) = \mathbb{V}_{\theta}(s_{\theta}(x)) = \mathbb{V}_{\theta}\left(\frac{x}{\theta}\right) = \frac{1}{\theta^2} \mathbb{V}_{\theta}(x) = \frac{1}{\theta}$$

$$\text{For } \left. \begin{array}{l} \hat{\theta} = \frac{1}{n} \sum X_j \\ g(\theta) = \theta \end{array} \right\} \mathbb{V}_{\theta}(\hat{\theta}) = \frac{1}{n} \theta = \frac{1}{n \cdot \frac{1}{\theta}} = \frac{(\partial_{\theta} g(\theta))^2}{I_n(\theta)}$$

CRLB.

$$\text{Now for } g(\theta) = e^{-\theta}, \quad T = \left(\frac{n-1}{n}\right)^{\sum X_j} \quad (\text{UMVU})$$

$$\begin{aligned} \mathbb{E}_{\theta}(T^2) &= \sum_{j=0}^{\infty} \left(1 - \frac{1}{n}\right)^{2j} \frac{(n\theta)^j}{j!} e^{-n\theta} \\ &= e^{-n\theta} \sum_{j=0}^{\infty} \frac{1}{j!} \left(\frac{(n-1)^2}{n} \theta\right)^j = e^{-n\theta} e^{\frac{(n-1)^2}{n} \theta} \\ &= \exp\left(\frac{(1-2n)\theta}{n}\right) \end{aligned}$$

$$\begin{aligned} \Rightarrow \mathbb{V}_{\theta}(T) &= \mathbb{E}_{\theta} T^2 - (\mathbb{E}_{\theta} T)^2 = \mathbb{E}_{\theta} T^2 - e^{-2\theta} \\ &= e^{-2\theta} (e^{\theta/n} - 1) \rightarrow e^{-2\theta} \frac{\theta}{n} \end{aligned}$$

$$\text{Since } g'(\theta) = -e^{-\theta} \quad \frac{g'(\theta)^2}{n I_1(\theta)} = \frac{\theta e^{-2\theta}}{n}$$

$\Rightarrow$  UMVU estimator does not reach the CRLB.