S&DS 242/542: Theory of Statistics Lecture 15: Fisher information and the Cramer-Rao bound

#### Recap: Asymptotic normality of the MLE

Let  $X_1, \ldots, X_n \stackrel{ID}{\sim} f(x \mid \theta)$  with a single parameter  $\theta \in \mathbb{R}$ . Under regularity conditions, the MLE  $\hat{\theta}$  satisfies

$$\sqrt{n}(\hat{\theta} - \theta) o \mathcal{N}(0, \frac{1}{I(\theta)})$$

in distribution as  $n \to \infty$ , where  $I(\theta)$  is the Fisher information, defined by

$$I(\theta) = \mathsf{Var}_{\theta} \left[ \frac{\partial}{\partial \theta} \log f(X \mid \theta) \right] = -\mathbb{E}_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log f(X \mid \theta) \right]$$

Thus the distribution of  $\hat{\theta}$  for large *n* is approximately  $\mathcal{N}(\theta, \frac{1}{nI(\theta)})$ , and an asymptotic  $(1 - \alpha)$ -confidence interval for  $\theta$  is

$$\hat{\theta} \pm z^{(1-\alpha/2)} \sqrt{\frac{1}{nI(\hat{\theta})}}$$

## Expected log-likelihood function



Recall that as  $n \to \infty$ , if the true parameter is  $\theta^*$ , then the average log-likelihood approaches the *expected* log-likelihood

$$\frac{1}{n}\ell_n(\theta) = \frac{1}{n}\sum_{i=1}^n \log f(X_i \mid \theta) \to L(\theta) = \mathbb{E}_{\theta^*}[\log f(X \mid \theta)]$$

in probability as  $n \to \infty$ . The expected log-likelihood  $L(\theta)$  is maximized at the true parameter  $\theta^*$ .

### Expected score



The derivative of  $L(\theta)$  is

$$L'(\theta) = \frac{\partial}{\partial \theta} \mathbb{E}_{\theta^*} [\log f(X \mid \theta)] = \mathbb{E}_{\theta^*} [\underbrace{\frac{\partial}{\partial \theta} \log f(X \mid \theta)}_{=\text{score}}]$$

So  $L'(\theta)$  is the expected score. Since  $L(\theta)$  is maximized at  $\theta^*$ , this also explains why the expected score is 0 at the true parameter  $\theta^*$ .

## Geometry of the Fisher information



The second derivative of  $L(\theta)$  is

$$L''( heta) = rac{\partial^2}{\partial heta^2} \mathbb{E}_{ heta^*}[\log f(X \mid heta)] = \mathbb{E}_{ heta^*}[rac{\partial^2}{\partial heta^2}\log f(X \mid heta)]$$

Then in particular, at the true parameter  $\theta^*$ ,

$$-L''( heta^*) = -\mathbb{E}_{ heta^*} \Big[ rac{\partial^2}{\partial heta^2} \log f(X \mid heta) \Big|_{ heta = heta_*} \Big] = I( heta^*)$$

So  $I(\theta^*)$  measures the curvature of  $L(\theta)$  around its maximizer  $\theta^*$ .

## Example of large Fisher information



When Fisher information  $I(\theta^*) = -L''(\theta^*)$  is large, moving  $\theta$  slightly away from  $\theta^*$  leads to a large decrease in the log-likelihood  $L(\theta)$ , i.e. a small change in  $\theta$  makes the observed data much less probable. In this sense, the data contains more "information" that the parameter must be close to the true parameter  $\theta^*$ .

## Example of small Fisher information



Conversely, when Fisher information  $I(\theta^*) = -L''(\theta^*)$  is small, large changes in  $\theta$  only decrease the log-likelihood value  $L(\theta)$  by a small amount, i.e. the observed data remains highly probable over large ranges of  $\theta$ . In this sense, our data contains less "information" that the parameter is close to the true parameter  $\theta^*$ . Fisher information as an intrinsic property of the model

To summarize:

- *I*(θ) represents the curvature of the expected log-likelihood function around its maximum, when θ is the true parameter.
- *I*(θ) quantifies the rate at which the expected log-likelihood decreases in value around the true parameter, and thus may be interpreted as a measure of the expected amount of information the data contains about θ.
- We have shown that I(θ) is related to the asymptotic variance of the MLE, but it may also be understood as a fundamental quantity about the model f(x | θ) irrespective of any specific estimation procedure.

#### Theorem (Cramer-Rao bound)

Consider a parametric model  $f(x \mid \theta)$  where  $\theta \in \mathbb{R}$  is a single parameter. Let  $T(X_1, \ldots, X_n)$  be any unbiased estimator of  $\theta$  based on data  $X_1, \ldots, X_n \stackrel{\text{IID}}{\sim} f(x \mid \theta)$ . Then

$$\mathsf{Var}_{ heta}[\mathcal{T}] \geq rac{1}{n I( heta)}$$

This is an example of an information-theoretic lower bound, stating that  $\frac{1}{nI(\theta)}$  is a fundamental limit to how accurate *any* unbiased estimator of  $\theta$  can be.<sup>1</sup>

<sup>1</sup>Related results show that no estimator — unbiased or not — can achieve a lower asymptotic mean-squared-error than  $\frac{1}{nI(\theta)}$ , except possibly at a measure-zero set of special values for  $\theta$ . We will not discuss these results in our course.

Proof: Define Z= 30 l, (0) = 2 30 log f(X. (0) score at X: Recall, at the tru punch O, IE [ ] ly f(x10)]= 0, Varo [ ] ly f(x10)]= I(0) ⇒ E\_0[2]=0, V~0[2]=n·I(0). Correlation between Z and my unbiased estimate T is Cong [Z,T] JVng [2] JVng [7] e[-1,1]



T is an unbiased estimater of O: For all O.  $\theta = \mathbb{E}_{\theta}[T] = \int T(x_{y,y}x_{y}) \cdot f(x_{y,y}x_{y}|\theta) dx_{y,y}dx_{y}$ Differentiale in O:

 $I = \int T(x_{v_1}, x_n) \cdot \frac{2}{30} \cdot \frac{f(x_{v_1}, x_n \mid 0)}{20} dx_{v_1} - dx_n$ =  $\frac{2}{30} \log \cdot \frac{f(x_{v_1}, x_n \mid 0)}{20} \times \frac{f(x_{v_2}, x_n \mid 0)}{20}_{10}$ 

 $= IE_{\theta} \left[ T(X_{i}, X_{n}) \cdot \frac{2}{5\theta} l_{y} f(X_{i}, X_{n}(\theta)) \right]$ 

=> [= E, [TZ] So Covo[T, Z]= Eo[TZ]-Eo[T]·Eo[Z]

= EA [TZ]=1

## Asymptotic efficiency

An estimator  $\hat{\theta}$  is **asymptotically efficient** if, in distribution as  $n \to \infty$ ,

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}\left(0, \frac{1}{I(\theta)}\right)$$

I.e., its variance for large *n* is approximately equal to the Cramer-Rao lower bound of  $\frac{1}{nl(\theta)}$ . Thus, our main theorem from last lecture shows that the MLE is asymptotically efficient.

If two estimators  $\hat{ heta}_1$  and  $\hat{ heta}_2$  based on  $X_1, \ldots, X_n \overset{IID}{\sim} f(x \mid heta)$  satisfy

$$\sqrt{n}(\hat{\theta}_1 - \theta) \rightarrow \mathcal{N}(0, v_1(\theta))$$
  
 $\sqrt{n}(\hat{\theta}_2 - \theta) \rightarrow \mathcal{N}(0, v_2(\theta))$ 

as  $n \to \infty$  where  $v_2(\theta) < v_1(\theta)$ , then  $v_2(\theta)/v_1(\theta)$  is the asymptotic relative efficiency of  $\hat{\theta}_1$  relative to  $\hat{\theta}_2$ .

# Asymptotic relative efficiency

Since  $\operatorname{Var}[\hat{\theta}_1] \approx \frac{v_1(\theta)}{n}$  and  $\operatorname{Var}[\hat{\theta}_2] \approx \frac{v_2(\theta)}{n}$  for large *n*, relative efficiency may also be interpreted as the ratio of sample sizes needed for  $\hat{\theta}_1$  and  $\hat{\theta}_2$  to achieve the same variance: E.g. if  $v_2(\theta) = 80\% \times v_1(\theta)$ , then  $\hat{\theta}_1$  using *n* samples has roughly the same variance as  $\hat{\theta}_2$  using  $80\% \times n$  samples.

Example: For  $X_1, \ldots, X_n \stackrel{ID}{\sim} \text{Pareto}(\theta, 1)$ , we have seen that the method-of-moments estimator  $\hat{\theta}_{MoM}$  and MLE  $\hat{\theta}_{MLE}$  satisfy

$$egin{aligned} &\sqrt{n}(\hat{ heta}_{\mathsf{MoM}}- heta) o \mathcal{N}(0,rac{ heta( heta-1)^2}{ heta-2}) \ &\sqrt{n}(\hat{ heta}_{\mathsf{MLE}}- heta) o \mathcal{N}(0, heta^2) \end{aligned}$$

so the asymptotic relative efficiency of method-of-moments to maximum likelihood is  $\frac{\theta(\theta-2)}{(\theta-1)^2}$  (which is less than 1).

## Cramer-Rao bound for plug-in estimators

The Cramer-Rao lower bound holds also for plug-in estimators: Letting  $\hat{\theta}$  be the MLE, the estimate  $g(\hat{\theta})$  for  $g(\theta)$  satisfies

$$\sqrt{n}(g(\hat{\theta}) - g(\theta)) \rightarrow \mathcal{N}\left(0, \frac{g'(\theta)^2}{I(\theta)}\right)$$

by the delta method. Thus  $\operatorname{Var}[g(\hat{\theta})] \approx \frac{g'(\theta)^2}{nI(\theta)}$  for large n.

#### Theorem

For a parametric model  $f(x \mid \theta)$  with a single parameter  $\theta \in \mathbb{R}$ , if  $T(X_1, \ldots, X_n)$  is any unbiased estimator of  $g(\theta)$  based on data  $X_1, \ldots, X_n \stackrel{IID}{\sim} f(x \mid \theta)$ , then

$$\operatorname{Var}_{\theta}[T] \geq rac{g'( heta)^2}{nI( heta)}.$$

## Efficiency of plug-in estimators

An estimator T for  $g(\theta)$  based on  $X_1, \ldots, X_n \stackrel{IID}{\sim} f(x \mid \theta)$  is asymptotically efficient if, in distribution as  $n \to \infty$ ,

$$\sqrt{n}(T-g(\theta)) \rightarrow \mathcal{N}\left(0, \frac{g'(\theta)^2}{I(\theta)}\right)$$

Thus, the plug-in MLE  $g(\hat{\theta})$  is asymptotically efficient.

Example: For  $X_1, \ldots, X_n \stackrel{ID}{\sim} \text{Pareto}(\theta, 1)$ , the mean is  $g(\theta) = \frac{\theta}{\theta - 1}$ . We have seen that

$$egin{aligned} &\sqrt{n}(ar{X}-g( heta)) 
ightarrow \mathcal{N}(0,rac{ heta}{( heta-1)^2( heta-2)}) \ &\sqrt{n}(g(\hat{ heta}_{\mathsf{MLE}})-g( heta)) 
ightarrow \mathcal{N}(0,rac{ heta^2}{( heta-1)^4}) \end{aligned}$$

so the asymptotic relative efficiency of  $\bar{X}$  to the plug-in MLE  $g(\hat{\theta}_{\text{MLE}})$  is also  $\frac{\theta(\theta-2)}{(\theta-1)^2}$ .

#### Fisher information for multiple parameters

For a parametric model with k parameters  $\theta \in \mathbb{R}^k$ , the **Fisher** information matrix  $I(\theta) \in \mathbb{R}^{k \times k}$  is the matrix whose (i, j) entry is defined by the equivalent expressions

$$I(\theta)_{ij} = \operatorname{Cov}_{\theta} \left[ \frac{\partial}{\partial \theta_i} \log f(X|\theta), \ \frac{\partial}{\partial \theta_j} \log f(X|\theta) \right]$$
$$= -\mathbb{E}_{\theta} \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(X|\theta) \right]$$

For k = 1, this " $1 \times 1$  matrix"  $I(\theta)$  is the same as our previous definition of the Fisher information for a single parameter.

## Asymptotic normality of the MLE

The inverse of this Fisher information matrix describes the asymptotic covariance matrix of the MLE  $\hat{\theta} \in \mathbb{R}^k$ .

#### Theorem

Let  $f(x \mid \theta)$  be a parametric model where  $\theta \in \mathbb{R}^k$ . Let  $\hat{\theta} \in \mathbb{R}^k$  be the MLE based on  $X_1, \ldots, X_n \stackrel{ID}{\sim} f(x \mid \theta)$ . Then, under regularity assumptions,

$$\sqrt{n}(\hat{\theta}-\theta) \rightarrow \mathcal{N}(0, I(\theta)^{-1}).$$

The right side is a k-dimensional multivariate normal distribution, whose covariance  $I(\theta)^{-1}$  is the  $k \times k$  matrix inverse of  $I(\theta)$ .

#### MLE in the Gamma model

Example: Let  $X_1, \ldots, X_n \stackrel{ID}{\sim}$  Gamma $(\alpha, \beta)$ . Recall that the MLEs  $(\hat{\alpha}, \hat{\beta})$  do not have closed-form expressions and are typically computed numerically.

$$f(x \mid \alpha, \beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} x^{\alpha-i} e^{-\beta x}$$

$$\Rightarrow \log f(x \mid \alpha, \beta) = \alpha \log \beta - \log \Gamma(\alpha) + (\alpha - i) \log x - \beta x$$

$$\frac{2}{2\alpha} \log f(x \mid \alpha, \beta) = \log \beta - \frac{2}{2\alpha} \log \Gamma(\alpha) + \log x$$

$$= \psi(\alpha) + \log \frac{1}{2\alpha} + \log \frac{1}{2\alpha} + \log \frac{1}{2\alpha}$$

$$\frac{\partial}{\partial \beta} \log - (x | x, \beta) = \frac{\alpha}{\beta} - x$$

MLE in the Gamma model

32 log + (x(α,β)= - 2 lg P(a) = 1/ (a) trigamon Comprise

The log f(x(x, B) = 1 , De log f(x(x, B) = - de log f(x, B) = - de log 5.  $I(x,\beta) = -E_{x,\beta} \begin{bmatrix} \frac{\partial^2}{\partial x^2} l_3 f(X|x,\beta) & \frac{\partial^2}{\partial x^2} l_3 f(X|x,\beta) \\ \frac{\partial^2}{\partial x^2} l_3 f(X|x,\beta) & \frac{\partial^2}{\partial p^2} l_3 f(X|x,\beta) \end{bmatrix}$ 

$$= \begin{bmatrix} \gamma_{i}(\alpha) & -1/\beta \\ -1/\beta & \alpha/\beta^{2} \end{bmatrix}$$
  
$$\Rightarrow I(\alpha, \beta)^{-1} = \frac{1}{\gamma_{i}(\alpha)\alpha'\beta^{2} - 1/\beta^{2}} \begin{bmatrix} \alpha/\beta^{2} & 1/\beta \\ 1/\beta & \gamma_{i}(\alpha) \end{bmatrix}$$

### MLE in the Gamma model

Informally, for large n, the MLEs  $\hat{\alpha}$  and  $\hat{\beta}$  have an approximate bivariate normal distribution

$$\mathcal{N}\left(\begin{pmatrix}\alpha\\\beta\end{pmatrix}, \frac{1}{n}I(\alpha,\beta)^{-1}\right)$$

The approximate variance of  $\hat{\alpha}$  is the upper-left entry of  $\frac{1}{n}I(\alpha,\beta)^{-1}$ , which is  $\frac{1}{n} \cdot \frac{\alpha}{\psi_1(\alpha)\alpha-1}$ .

The approximate covariance of  $\hat{\alpha}$  and  $\hat{\beta}$  is the off-diagonal entry of  $\frac{1}{n}I(\alpha,\beta)^{-1}$ , which is  $\frac{1}{n} \cdot \frac{\beta}{\psi_1(\alpha)\alpha-1}$ . This is always positive, implying that the errors  $\hat{\alpha} - \alpha$  and  $\hat{\beta} - \beta$  are positively correlated for large *n*.

## Comparison of variances when a parameter is known

Consider any model with two parameters  $(\alpha, \beta)$  and Fisher information matrix

$$I(\alpha,\beta) = \begin{pmatrix} \mathsf{a} & \mathsf{c} \\ \mathsf{c} & \mathsf{b} \end{pmatrix}$$

In the joint MLE  $(\hat{\alpha}, \hat{\beta})$ , the asymptotic variance of  $\hat{\alpha}$  is the upper-left entry of  $\frac{1}{n}I(\alpha, \beta)^{-1}$ , which is

$$\frac{1}{n} \cdot \frac{1}{a - c^2/b}$$

Suppose instead that  $\beta$  is known, and let  $\hat{\alpha}'$  be the MLE in model with a single parameter  $\alpha \in \mathbb{R}$ . The Fisher information in this one-parameter model is the upper-left entry of  $I(\alpha, \beta)$ , so the asymptotic variance of  $\hat{\alpha}'$  is

$$\frac{1}{n} \cdot \frac{1}{a}$$

# Comparison of variances when a parameter is known

Here 
$$a - c^2/b > 0$$
 and  $a, b > 0$ , so  $\label{eq:absolution} \frac{1}{a - c^2/b} > \frac{1}{a}$ 

always. The difference between  $\frac{1}{a-c^2/b}$  and  $\frac{1}{a}$  represents the difference in asymptotic variance of the MLE for estimating  $\alpha$  in settings when  $\beta$  is unknown vs. when  $\beta$  is known.

This is an example of a trade-off between model complexity and accuracy of estimation: A complex model with more parameters might better capture the true distribution of data, but each individual parameter may be more difficult to estimate than in a simpler model with fewer parameters.