

Maximum Likelihood Estimation for Parametric Models

↳ X has pmf/pdf $f(x|\theta)$, with known f ,
 θ is finite dimensional but unknown

DEF (Correct specification)

A model is correctly specified when there is a unique $\theta_0 \in \Theta$ s.t. $f(x|\theta_0)$ coincides with the true density of X .

DEF (misspecification)

A model is misspecified if no value of $\theta \in \Theta$

DEF (likelihood function (of iid data))

$$L_n(\theta) = f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

joint density / pmf of the data

DEF (MLE)

A MLE $\hat{\theta}_{ML}$ maximizes

• the likelihood function $L_n(\theta) = \prod_{i=1}^n f(x_i | \theta)$

• the log-likelihood function $l_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log f(x_i | \theta)$

a function of parameter
 = what value of θ maximizes the likelihood of observing sample $\{x_1, \dots, x_n\}$

DEF (Expected likelihood function)

$$l(\theta) = E[\log f(x|\theta)]$$

pdf of the RV X , not indexed by n

Theorem When a model is correctly specified, θ_0 maximizes $l(\theta)$.

$$\theta_0 = \operatorname{argmax}_{\theta \in \Theta} E[\log f(x|\theta)]$$

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta \in \Theta} \left(\frac{1}{n} \sum_{i=1}^n \log f(x_i | \theta) \right)$$

→ sample analog

DEF (Efficient score) $s = \frac{\partial}{\partial \theta} \log f(x|\theta_0)$

DEF (Fisher information) $\mathcal{F}_\theta = E[ss']$

Theorem When a model is correctly specified, support X doesn't depend on θ , $\theta_0 \in \text{int}(\Theta)$, then $E[S] = 0$, $\text{var}(S) = \mathcal{F}_\theta$.

Proof. $E[S] = E\left[\frac{\partial}{\partial \theta} \log f(x|\theta_0)\right]$
 $\stackrel{\text{Leibniz rule}}{=} \frac{\partial}{\partial \theta} E[\underbrace{\log f(x|\theta_0)}_{\ell(\theta_0)}]$

$$= \frac{\partial}{\partial \theta} \ell(\theta_0)$$

$$= 0 \quad \text{b/c } \theta_0 \in \Theta, \text{ interior Sol given by FOC.}$$

$$\text{var}(S) = E[SS'] - \underbrace{(E[S])^2}_0 = \mathcal{F}_\theta.$$

Theorem (Information matrix equality)

$$E\left[\underbrace{\begin{bmatrix} \frac{\partial \log f(x|\theta_0)}{\partial \theta} & \frac{\partial \log f(x|\theta_0)}{\partial \theta'} \end{bmatrix}}_{\text{Fisher information}}\right] = - \underbrace{E\left[\frac{\partial^2}{\partial \theta \partial \theta'} \log f(x|\theta_0)\right]}_{\mathcal{F}_\theta \text{ expected Hessian}}$$

Theorem Assume model is correctly specified, support X doesn't depend on θ , $\theta_0 \in \text{int}(\Theta)$. If $\tilde{\theta}$ is an unbiased estimator of θ , then

$$\text{var}(\tilde{\theta}) \geq \underbrace{(n\mathcal{F}_\theta)^{-1}}_{\text{CRLB}}$$

An estimator is Cramer-Rao efficient if it is unbiased and $\text{var}(\tilde{\theta}) = (n\mathcal{F}_\theta)^{-1}$.

Intuition: more curvature of log likelihood function wrt θ , more information data provides about θ b/c $\text{var}(\theta)$ is lower bounded by the inverse of \mathcal{F}_θ . more curvature \Rightarrow more precise estimate of $\hat{\theta}_{MLE}$

Asymptotic Properties of MLE

MLE is • consistent

• asymptotically normal $\sqrt{n}(\hat{\theta}_{MLE} - \theta) \xrightarrow{d} \mathcal{N}(0, \mathcal{F}_\theta^{-1})$

• asymptotically Cramer-Rao efficient.

$$\text{var}(\hat{\theta}_{MLE}) = (n\mathcal{F}_\theta)^{-1}$$

4. Based on the notation in the slides on *Estimation*, let us prove the Information Matrix Equality

$$\mathbb{E} \left[\frac{\partial^2 \log f(X|\theta_0)}{\partial \theta \partial \theta'} \right] = -\mathbb{E} \left[\frac{\partial \log f(X|\theta_0)}{\partial \theta} \frac{\partial \log f(X|\theta_0)}{\partial \theta'} \right].$$

Let $f = f(x|\theta_0)$, ∇_j means derivative with respect to the j -th element $\theta^{(j)}$, and ∇_{jk} mean 2nd-order derivative with respect to $\theta^{(j)}$ and $\theta^{(k)}$. Suppose we can exchange the integral “ \int ” and derivatives “ ∇_j ”.

(a) By differentiating $\int f dx = 1$ with respect to $\theta^{(j)}$, show that $\mathbb{E}[\nabla_j \log f] = 0$.

(b) By differentiating $\mathbb{E}[\nabla_j \log f] = 0$ with respect to $\theta^{(k)}$, show that

$$\mathbb{E}[\nabla_{jk} \log f] + \mathbb{E}[(\nabla_j \log f)(\nabla_k \log f)] = 0,$$

which yields the Information Matrix Equality.

$$(a) \quad \nabla_j \int f dx = \int \nabla_j f dx = 0$$

↑
exchange \int, ∇_j

$$0 = \int (\nabla_j f) dx$$

$$= \int (\nabla_j f) \underbrace{\frac{1}{f}}_1 \cdot f dx$$

Notice: $\nabla_j \log f = \frac{1}{f} \nabla_j f$
by chain rule

$$= \int (\nabla_j \log f) f dx$$

a function of x density $f(x|\theta_0)$

$$\stackrel{\text{def.}}{=} \mathbb{E}[\nabla_j \log f] \quad \rightarrow \text{"population" FOC of the } j\text{-th element}$$

↑
 $f(x|\theta_0)$

$$(b) \quad 0 = \nabla_k \mathbb{E}[\nabla_j \log f]$$

$$\stackrel{\text{plugin}}{=} \nabla_k \int (\nabla_j \log f) f dx$$

$$\stackrel{\text{exchange } \int, \nabla}{=} \int \nabla_k ((\nabla_j \log f) f) dx$$

$$\stackrel{\text{product rule}}{=} \int \{(\nabla_k \log f) f + (\nabla_j \log f)(\nabla_k f)\} dx$$

$$\stackrel{\int \text{ is linear}}{=} \int (\nabla_{jk} \log f) f dx + \int (\nabla_j \log f)(\nabla_k f) dx$$

$$\underbrace{\int (\nabla_{jk} \log f) f dx}_{\mathbb{E}[\nabla_{jk} \log f]} \quad \uparrow f(x|\theta_0) \quad = \int (\nabla_j \log f) \left\{ (\nabla_k f) \cdot \frac{1}{f} \cdot f \right\} dx$$

$$\nabla_k \log f = \frac{1}{f} \nabla_k f$$

$$= \int (\nabla_j \log f)(\nabla_k \log f) \cdot f dx$$

$$= \mathbb{E}[(\nabla_j \log f)(\nabla_k \log f)]$$

3. Suppose X follows a normal distribution with unknown mean μ and variance $\sigma^2 > 0$. The density of X is

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Given a random sample $\{X_i, i = 1 \dots n\}$ drawn from X , find the MLE estimator for (μ, σ^2) .

likelihood function: $L_n(\mu, \sigma^2) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \right)$ b/c random sampling

$$= \left(\frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} e^{-\sum_{i=1}^n \frac{(x_i-\mu)^2}{2\sigma^2}}$$

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

MLE $(\hat{\mu}, \hat{\sigma}^2)$ should satisfy FOC:

$$\left. \frac{\partial \ln(\mu, \sigma^2)}{\partial \mu} \right|_{\substack{\mu=\hat{\mu} \\ \sigma^2=\hat{\sigma}^2}} = 2 \left(\frac{1}{2\hat{\sigma}^2} \right) \sum_{i=1}^n (x_i - \hat{\mu}) = 0 \quad \dots \textcircled{1}$$

$$\left. \frac{\partial \ln(\mu, \sigma^2)}{\partial \sigma^2} \right|_{\substack{\mu=\hat{\mu} \\ \sigma^2=\hat{\sigma}^2}} = -\frac{n}{2} \frac{1}{2\pi\hat{\sigma}^2} 2\pi + \sum_{i=1}^n \frac{(x_i - \hat{\mu})^2}{2(\hat{\sigma}^2)^2} = 0 \quad \dots \textcircled{2}$$

take derivative with (σ^2) as a whole

$\textcircled{1}$: we get $\sum_{i=1}^n (x_i - \hat{\mu}) = 0$

$$\Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

plug $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$ into $\textcircled{2}$: $\frac{n}{2\hat{\sigma}^2} = \sum_{i=1}^n \frac{(x_i - \hat{\mu})^2}{2(\hat{\sigma}^2)^2}$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{1}{n} \sum_{i=1}^n \left(x_i - \frac{1}{n} \sum_{i=1}^n x_i \right)^2.$$

To show $\hat{\mu}, \hat{\sigma}^2$ are the maximizer of $\ln(\mu, \sigma^2)$, need to check Hessian

let $\theta = (\mu, \sigma^2)$, $\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2)$

$$\begin{aligned} \frac{\partial^2 \ln(\theta)}{\partial \theta \partial \theta'} &= \begin{pmatrix} \frac{\partial^2 \ln(\mu, \sigma^2)}{\partial \mu \partial \mu} & \frac{\partial^2 \ln(\mu, \sigma^2)}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2 \ln(\mu, \sigma^2)}{\partial \sigma^2 \partial \mu} & \frac{\partial^2 \ln(\mu, \sigma^2)}{\partial \sigma^2 \partial \sigma^2} \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial}{\partial \mu} \left(\frac{1}{\sigma^2} \sum_{i=1}^n x_i - \frac{1}{\sigma^2} n \mu \right) & \frac{\partial}{\partial \sigma^2} \left(\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \right) \\ \frac{\partial}{\partial \sigma^2} \left(\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \right) & \frac{\partial}{\partial \sigma^2} \left(-\frac{n}{2} \frac{1}{\sigma^2} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^4} \right) \end{pmatrix} \\ &= \begin{pmatrix} -\frac{n}{\sigma^2} & -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) \\ -\frac{1}{\sigma^4} \sum_{i=1}^n (x_i - \mu) & \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \mu)^2 \end{pmatrix} \end{aligned}$$

$$\frac{\partial^2 \ln(\theta)}{\partial \theta \partial \theta'} \bigg|_{\theta = \hat{\theta}} = \begin{pmatrix} -\frac{n}{\hat{\sigma}^2} & 0 \\ 0 & -\frac{n}{2\hat{\sigma}^4} \end{pmatrix} \text{ is ND.}$$

$$\begin{aligned} & \leftarrow \text{b/c } \frac{1}{\hat{\sigma}^4} \sum_{i=1}^n (x_i - \frac{1}{n} \sum_{i=1}^n x_i) \\ & = \frac{1}{\hat{\sigma}^4} \left(\underbrace{\sum_{i=1}^n x_i - n \left(\frac{1}{n} \sum_{i=1}^n x_i \right)}_{=0} \right) = 0 \end{aligned}$$

$$\begin{aligned} & \frac{n}{2\hat{\sigma}^4} - \frac{1}{\hat{\sigma}^6} \sum_{i=1}^n (x_i - \hat{\mu})^2 \\ & = \frac{n}{2\hat{\sigma}^4} - \frac{\cancel{n}}{\hat{\sigma}^6} \underbrace{\frac{1}{\cancel{n}} \sum_{i=1}^n (x_i - \hat{\mu})^2}_{=\hat{\sigma}^2} = \frac{n}{2\hat{\sigma}^4} - \frac{n}{\hat{\sigma}^4} = -\frac{n}{2\hat{\sigma}^4} \end{aligned}$$