

Theorem: Let A1-2-3 hold, and assume that

$$\frac{\partial^2}{\partial \vartheta_j \partial \vartheta_i} \int \log \frac{P_\vartheta(x)}{P_{\vartheta'}(x)} P_\vartheta(x) \mu(dx) =$$

$$= \int \frac{\partial^2}{\partial \vartheta_j \partial \vartheta_i} \left(\log \frac{P_\vartheta(x)}{P_{\vartheta'}(x)} P_\vartheta(x) \right) \mu(dx).$$

Then,

$$I(\vartheta)_{ij} = \frac{\partial^2}{\partial \vartheta_i' \partial \vartheta_j'} D_{KL}(\vartheta \parallel \vartheta') \Big|_{\vartheta=\vartheta'}.$$

Note: ① At $\vartheta = \vartheta'$ KL is 0

② The second derivative tells us how fast the KL divergence grows as we move away from ϑ .

(3) Fisher information is the local curvature of KL divergence at θ .

Proof:

$$D_{KL}(\theta | \theta') = \int [\log p_{\theta}(x) - \underbrace{\log p_{\theta'}(x)}_{\text{depends on } \theta'}] p_{\theta}(x) \mu(dx)$$

$$\frac{\partial}{\partial \theta'_i} D_{KL}(\theta | \theta') = - \int \frac{\partial}{\partial \theta'_i} \log p_{\theta'}(x) \cdot p_{\theta}(x) \mu(dx)$$

$$\frac{\partial}{\partial \theta'_i \partial \theta'_j} D_{KL}(\theta | \theta') = - \int \frac{\partial^2}{\partial \theta'_i \partial \theta'_j} \log p_{\theta'}(x) \cdot p_{\theta}(x) \mu(dx)$$

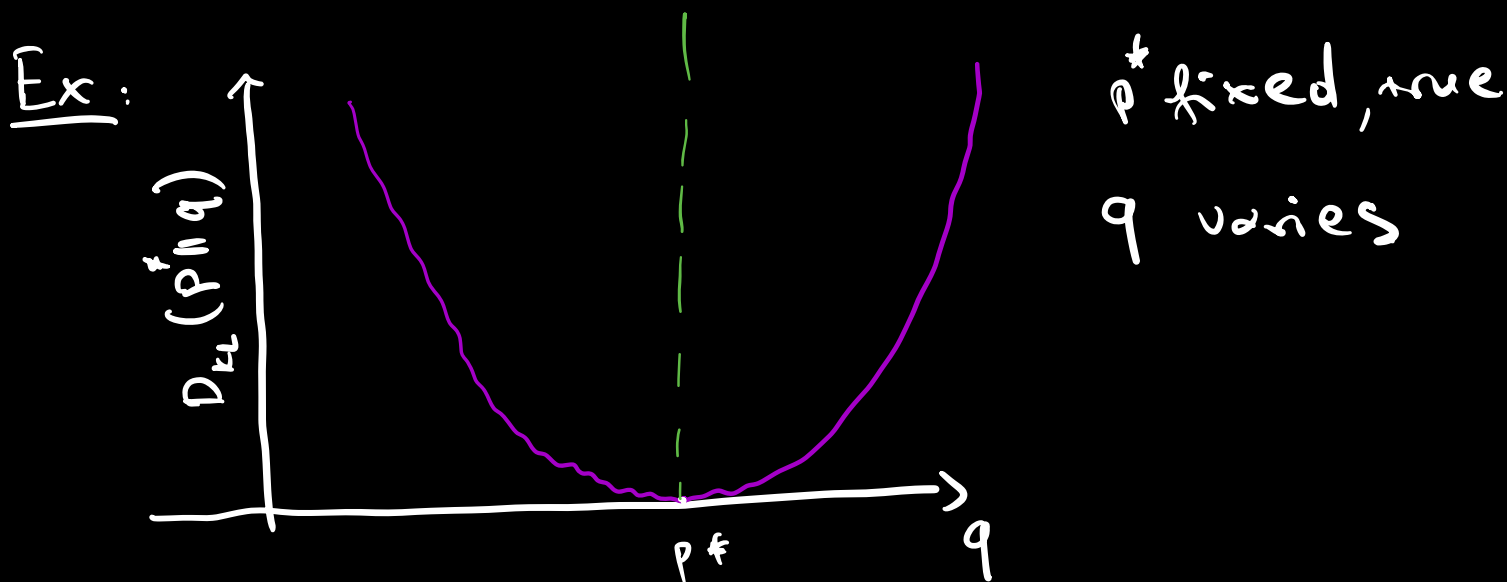
$$\left. \frac{\partial}{\partial \theta'_i \partial \theta'_j} D_{KL}(\theta | \theta') \right|_{\theta = \theta'} = - \left. \int \frac{\partial^2}{\partial \theta'_i \partial \theta'_j} \log p_{\theta'}(x) \right|_{\theta = \theta'} p_{\theta}(x) \mu(dx)$$

By Fisher information matrix.

$$I(\theta)_{ij} = E_{\theta} \left[\frac{\partial}{\partial \theta'_i} \log p_{\theta}(x) \cdot \frac{\partial}{\partial \theta'_j} \log p_{\theta}(x) \right] = - E_{\theta} \left(\frac{\partial^2}{\partial \theta'_i \partial \theta'_j} \log p_{\theta}(x) \right)$$

So we conclude that:

$$I(\theta)_{ij} = \frac{\partial^2}{\partial \theta'_i \partial \theta_j} D_{KL}(\theta | \theta') \Big|_{\theta' = \theta} \quad \square$$



Asymptotic Statistics:

Q: What happens when $n \rightarrow \infty$?

Consider a family $\{P_\theta\}_{\theta \in \Theta}$ of probability measures on (X, \mathcal{G}) and the statistical model

$$(X, \mathcal{F}, \{P_\theta\}_{\theta \in \Theta}) = (X^N, \mathcal{G}^{\otimes N}, \{P_\theta^{\otimes N}\}_{\theta \in \Theta})$$

allowing us to define the sequence:

$$\{\Sigma_n\}_{n=1}^{\infty} = \{(x_1, \dots, x_n)\}_{n=1}^{\infty}$$

Theorem (Ionescu-Tulcea): for $X = \mathbb{R}$

$\mathcal{F} = \mathcal{B}(\mathbb{R})$, fixed P_θ , there exists
unique $P_\theta = P_\theta^{\otimes \infty}$ with

$$\forall \{A_i \in \mathcal{F}\}_{i=1}^n \quad P_\theta(A_1 \times \dots \times A_n \times \mathbb{R} \times \dots) \\ = P_\theta(A_1) \dots P_\theta(A_n)$$

Notation: let $g: \Theta \rightarrow \mathbb{R}^p$,

$$\Gamma = \{g(\vartheta) : \vartheta \in \Theta\}.$$

Definition: we call $\{T_n\}_{n=1}^{\infty}$ for

$T_n: X^{\otimes n} \rightarrow \Gamma$ a sequence of estimators of $g(\vartheta)$ if T_n is an estimator

of $g(\theta)$ $\forall n \in \mathbb{N}$.

Example: $x_j \sim \text{Pois}(\theta)$ $g(\theta) = \theta$

$$T_n(x_n) = \frac{1}{n} \sum_{j=1}^n x_j$$

Definition: An estimator T_n of $g(\theta) \in \mathbb{R}^p$ is an μ -estimator if there exists $\rho: \mathcal{X} \times \mathbb{R}^p \rightarrow \mathbb{R}$ such that

$$T_n(\mathbf{x}) = \arg \min_{\gamma} \sum_{j=1}^n \rho(x_j, \gamma)$$

Ex(MLE): if $g(\theta) = \theta$ $\rho(x, \theta) = -\log p_\theta(x)$

we use that MLE is a special case of μ -estimator.

Ex(Huber) consider the loss

$$\rho_c(x, \gamma) = \begin{cases} \frac{1}{2} \|x_j - \gamma\|^2 & \|x - \gamma\| < c \\ c \|x_j - \gamma\|^2 & \|x - \gamma\| \geq c \end{cases}$$

Then the huber estimator

$$T_H(x) = \operatorname{argmin}_{\theta} \sum_{j=1}^n \rho(x_j, \theta)$$

Note that if $\rho(x, \cdot) \in C^1$ for all $x \in X$ we find the minimum by differentiation. $T_n(x)$ solves:

$$\sum_{j=1}^n \underbrace{\nabla_{\theta} \rho(x_j, \theta)}_{\psi(x, \theta)} = 0$$

Definition: An estimator T_n of $g(\theta) \in \mathbb{R}^p$ is a z -estimator if there exists $\psi: X = \mathbb{R}^p \rightarrow \mathbb{R}$ such that $T_n(x)$ solves:

$$\sum_{j=1}^n \psi(x_j, T_n(x)) = 0$$

we would like to have

$$T_n \longrightarrow g(\theta) \text{ as } n \rightarrow \infty.$$

Definition: a sequence of r.v. $\{T_n\}$ converges:

- in probability towards the r.v. T_*
 $(T_n \xrightarrow[n \rightarrow \infty]{P} T_*)$ if

$$\lim_{n \rightarrow \infty} \underline{P}(\|T_n - T\| > \varepsilon) = 0 \quad \forall \varepsilon > 0$$

- almost surely towards the r.v. T_*
 $(T_n \xrightarrow[n \rightarrow \infty]{a.s.} T_*)$ if

$$\underline{P}\left(\lim_{n \rightarrow \infty} T_n = T_*\right) = 1$$

Note: $T_n \xrightarrow{a.s.} T_* \xRightarrow{\text{Fatou lemma}} T_n \rightarrow T_*$

Definition: A sequence T_n of estimators of $g(\theta)$ is called consistent / strongly consistent if

$$\underline{T_n \xrightarrow[n \rightarrow \infty]{P} g(\theta)} \quad / \quad \underline{T_n \xrightarrow[n \rightarrow \infty]{a.s.} g(\theta) \quad \forall \theta \in \Theta}$$

Ex (Poisson) $x_j \stackrel{i.i.d}{\sim} \text{Pois}(\theta)$

$$g(\theta) = \theta \quad T_n(\mathbf{x}) = \frac{1}{n} \sum_j x_j$$

Prove that T_n is consistent.

$$\mathbb{P}[\|T_n(\mathbf{x}) - \theta\| > \varepsilon] \stackrel{\text{Chebyshev}}{\leq} \frac{\text{Var}(T_n)}{\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0$$

\downarrow
 $\frac{\theta}{n}$

$\Rightarrow T_n$ is consistent.

By the law of large numbers
we could expect that

$$\frac{1}{n} \sum p(x_j, \theta) \xrightarrow{\text{P.L.N.}} E_{\theta}(p(x, \theta))$$

and therefore $T_n \xrightarrow{n \rightarrow \infty} \underset{\theta}{\text{argmin}} E_{\theta}(p(x, \theta))$
 $\forall \theta \in \Theta$

Consistency of MLE estimators:

Theorem: let $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} p_\theta(x)$
for $\theta \in \Theta \subseteq \mathbb{R}^k$ and let:

(A1) identifiability

$\forall \theta \neq \theta_0, p_\theta(x) \neq p_{\theta_0}(x)$ on a set of positive measures

(A2) Parameter space $\Theta \subseteq \mathbb{R}^k$ compact.

(A3) The map $\theta \rightarrow \log p_\theta(x)$ is continuous for all x .

(A4) There exists a function $M(x) \geq 0$
s.t. $|\log p_\theta(x)| \leq M(x) \quad \forall \theta \in \Theta$.

and $\mathbb{E}_{\theta_0}[M(x)] < \infty$

(A5) for each $\theta \in \Theta$, x_1, \dots, x_n are i.i.d. for p_θ .

where θ_0 is the true parameter

Then the MLE estimator T_n is consistent i.e.:

$$T_n \xrightarrow[n \rightarrow \infty]{P} \vartheta_0$$

Proof: we want to show that

$$T_n = \arg \max_{\vartheta \in \Theta} l_n(\vartheta) \xrightarrow{P} \vartheta_0$$

$$\text{Let } l(\vartheta) = \mathbb{E}_{\vartheta_0} [\log p_{\vartheta}(x)] = \int \log p_{\vartheta}(x) p_{\vartheta_0}(dx)$$

By (A1) $l(\vartheta)$ is uniquely maximized at $\vartheta = \vartheta_0$.

By LLN (pointwise convergence) for fixed $\vartheta \in \Theta$ the log likelihood satisfies:

$$\frac{1}{n} l_n(\vartheta) = \frac{1}{n} \sum \log p_{\vartheta}(x_i) \xrightarrow{\text{a.s.}} l(\vartheta)$$

$\log p_{\vartheta}(x)$ is integrable under ϑ_0 (A4)

uniform convergence over Θ .

By A2 - A3 - A4 and the uniform law of large numbers ULLN

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \ln(\theta) - \ell(\theta) \right| \xrightarrow{P} 0$$

From arg max consistency

$$T_n = \arg \max_{\theta} \frac{1}{n} \ln(\theta)$$

$$\theta_0 = \arg \max_{\theta} \ell(\theta)$$

Then if Θ is compact

$$\frac{\ln(\theta)}{n} \rightarrow \ell(\theta)$$

$\ell(\theta)$ has a unique maximizer
at θ_0 .

$$\Rightarrow T_n \xrightarrow{P} \theta_0.$$



