

Parametric inference

- Z-estimators

$h(\cdot, \cdot) \in \mathbb{R}^d$
↑
param data.

Suppose $H(\theta) = \mathbb{E}[h(\theta; X)]$
 satisfies $H(\theta^*) = 0 \quad (\theta \in \mathbb{R}^d)$

Goal: estimate θ^* (point estimation, CI)

$$\hat{H}_n(\theta) = \frac{1}{n} \sum_{i=1}^n h(\theta; x_i).$$

Try to find θ^* by solving $\hat{H}_n(\theta) = 0$.

e.g. "method of moments".

Given $(P_\theta : \theta \in \Theta)$, for each $\theta \in \Theta, j \in \mathbb{N}$

we have moment

$$\alpha_j(\theta) = \mathbb{E}_\theta[X^j].$$

Suppose we have data $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} P_{\theta^*}$

$$\forall j, \quad \mathbb{E}_{\theta^*}[X^j - \alpha_j(\theta^*)] = 0$$

Let $h(\theta; x) = \begin{bmatrix} x - \alpha_1(\theta) \\ x^2 - \alpha_2(\theta) \\ \vdots \\ x^d - \alpha_d(\theta) \end{bmatrix} \in \mathbb{R}^d$.

$$H(\theta^*) = \mathbb{E}[h(\theta^*; x)] = 0$$

We can estimate θ^* by $\hat{H}_n(\hat{\theta}_n) = 0$

$$\forall j \in \{1, 2, \dots, d\}, \quad \frac{1}{n} \sum_{i=1}^n x_i^j = \alpha_j(\hat{\theta}_n).$$

Rank: does not have to be first d moments.

any set of ^v_{positive} integers could possibly work.

e.g. "M-estimation":

$$\hat{\theta}_n := \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n f(\theta; x_i)$$

$$\text{where we also define } F(\theta) := \mathbb{E}[f(\theta; x)]$$

$$\text{and } \theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} F(\theta).$$

Generally, M-estimators and Z-estimators
are not comparable.

When $f(\cdot; x)$ is differentiable $\forall x$.

and $\mathbb{H} = \mathbb{R}^d$, then M-estimator
is also a Z-estimator.

$$h(\theta; x) = \nabla_{\theta} f(\theta; x).$$

(in that case, M-estimator is a Z-estimator,

while Z-estimators may not be M-estimators
e.g. first-order saddle point, local minima)

In general, Z-estimators may not exist / satisfy good properties.

But it works "locally".

Thm: Assume $\exists \theta^* \text{ s.t. } H(\theta^*) = 0$.

Assume $h(\theta; x)$ is cts differentiable
at a local neighborhood of θ^* , $\forall x$.

Assume $\nabla_{\theta} H(\theta^*) \in \mathbb{R}^{d \times d}$ is non-degenerate.

then $\mathbb{P}(\exists \hat{\theta}_n, \hat{H}_n(\hat{\theta}_n) = 0) \rightarrow 1$, and
(Assuming that we search $\hat{\theta}_n$
at sufficiently small neighborhood
of θ^*)

$$\sqrt{n} (\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \nabla H(\theta^*)^{-1} \Sigma^* \nabla H(\theta^*)^{-1})$$

where $\Sigma^* = \text{Cov}_{\theta^*}(h(\theta^*, x))$

$$= \mathbb{E}[h(\theta^*, x) \cdot h(\theta^*, x)^T].$$

$$A^{-!} := (A^{-1})^T.$$

(Assuming all the expectations in the formula exist).

Remark:
• Σ^* , $\nabla H(\theta^*)$ can be estimated through empirical data

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n h(\hat{\theta}_n; x_i) h(\hat{\theta}_n; x_i)^T$$

$$\hat{J}_n = \frac{1}{n} \sum_{i=1}^n \nabla h(\hat{\theta}_n; x_i)$$

Easy to show $\hat{\Sigma}_n \xrightarrow{P} \Sigma^*$

$$\hat{J}_n \xrightarrow{P} \nabla H(\theta^*)$$

So we also have

$$\sqrt{n} \left(\hat{J}_n^{-1} \hat{\Sigma}_n \hat{J}_n^{-1} \right)^{-\frac{1}{2}} (\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \text{Id})$$

(in 1-d, $\frac{\sqrt{n} (\hat{\theta}_n - \theta^*) \cdot \hat{J}_n}{\sqrt{\hat{\Sigma}_n}} \xrightarrow{d} \mathcal{N}(0, 1)$)

- Alternatively, CI can be built using bootstrap.

Proof of Z-estimator result:

$$\begin{aligned} \theta &= \hat{H}_n(\hat{\theta}_n) \\ &= \underbrace{\hat{H}_n(\theta^*)}_{\text{u}} + \underbrace{\nabla \hat{H}_n(\theta^*) \cdot (\hat{\theta}_n - \theta^*)}_{\text{ignore}} + o_p\left(\hat{\theta}_n - \theta^*\right) \\ &\quad = \frac{1}{n} \sum_{i=1}^n \nabla h(\theta^*; x_i) \\ \frac{1}{n} \sum_{i=1}^n h(\theta^*; x_i) &\xrightarrow{P} 0 \qquad \xrightarrow{P} \nabla H(\theta^*) \\ \sqrt{n} \hat{H}_n(\theta^*) &\xrightarrow{d} \mathcal{N}(0, \Sigma^*) \qquad \begin{matrix} \text{invertible} \\ \text{matrix.} \end{matrix} \end{aligned}$$

Putting them together

$$\nabla H(\theta^*) \cdot (\hat{\theta}_n - \theta^*) = -\frac{1}{n} \sum_{i=1}^n h(\theta^*; x_i) + o_p(\hat{\theta}_n - \theta^*)$$

$$\therefore \sqrt{n} \cdot \nabla H(\theta^*) (\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \Sigma^*)$$

Applying $\nabla H(\theta^*)^{-1}$, we get the result.

(Existence of $\hat{\theta}_n$ can be shown
e.g. by first of an iterative procedure.)

Maximal Likelihood.

General principle:

$$X \sim P_{\theta^*} \quad (\theta^* \text{ does not have to be parametric})$$

Suppose $\forall \theta \in \Theta$, P_θ has a density $p_\theta(x)$

- obs r.v. pdf
- discrete r.v. pmf.

Suppose that $p_\theta(x)$ has common support $\forall \theta \in \Theta$.

then we can estimate θ^* by solving

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} p_\theta(x)$$

For this lecture: $\Theta \subseteq \mathbb{R}^d$

$$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} P_{\theta^*}.$$

In such a case

$$p_\theta^{(n)}(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p_\theta(x_i).$$

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p_\theta(x_i)$$

e.g. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(p)$.

$$P_p(x) = \begin{cases} p & x=1 \\ 1-p & x=0. \end{cases}$$

$$\hat{p}_n = \underset{p \in [0,1]}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \left(1_{X_i=1} \cdot \log p + 1_{X_i=0} \cdot \log (1-p) \right)$$

$$= \underset{p \in [0,1]}{\operatorname{argmax}} \left\{ \bar{X}_n \log p + (1-\bar{X}_n) \cdot \log (1-p) \right\}$$

Solve F.O.C. for p , get $\hat{p}_n = \hat{X}_n$.

e.g. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$P_{\mu, \sigma^2}(x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$L(X_1^n) = \sum_{i=1}^n \left(-\frac{(x_i - \mu)^2}{2\sigma^2} - \log(2\pi\sigma^2) \right).$$

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{easily from quadratic function})$$

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_n)^2 \quad (\text{by F.O.C. after plugging in } \hat{\mu}_n)$$

Applying 8-estimator results, we get

consistency & asymptotic normality

e.g. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Unif}([0, \theta]).$

$$P_\theta(x_i) = \begin{cases} \frac{1}{\theta} & x_i \in [0, \theta] \\ 0 & x_i \notin [0, \theta]. \end{cases}$$

$$P_\theta^{(n)}(x_1, x_2, \dots, x_n) = \begin{cases} \frac{1}{\theta^n} & \theta \geq \max_{1 \leq i \leq n} x_i \\ 0 & \theta < \max_{1 \leq i \leq n} x_i. \end{cases}$$

MLE estns. $\hat{\theta}_n = \max_{1 \leq i \leq n} x_i.$

$$\text{PL}(\hat{\theta}_n \leq t) = \text{PL}(X_1 \leq t)^n = \left(\frac{t}{\theta}\right)^n.$$

$$\mathbb{E}_\theta[(\hat{\theta}_n - \theta)^2] = O\left(\frac{1}{n^2}\right)$$

faster than MLE convergence rate in the smooth case.

Properties of MLE.

Suppose $\log p_\theta(x)$ twice cts differentiable.

and assuming $\mathbb{E} |\nabla \log p_\theta(x)|^2 < \infty$
 $\mathbb{E} |\nabla^2 \log p_\theta(x)| < \infty.$

we can then apply Z-estimator results

$$\sqrt{n} (\hat{\theta}_n - \theta^*) \xrightarrow{d} N(0, (\mathcal{J}^*)^{-1} \Sigma^* (\mathcal{J}^*)^{-T})$$

where $\Sigma^* := E_{\theta} \left[\nabla_{\theta} \log P_{\theta}(X) \nabla_{\theta} \log P_{\theta}(X)^T \right]$

$$\mathcal{J}^* = E_{\theta} [\nabla^2 \log P_{\theta}(X)]. \quad (\text{symmetric})$$

score function: $\nabla \log P_{\theta}(X)$

• $E_{\theta} [\nabla \log P_{\theta}(X)] = 0$

(Proof: $\int \nabla_{\theta} \log P_{\theta}(x) \cdot P_{\theta}(x) dx$
 $= \int \nabla_{\theta} (P_{\theta}(x)) dx = \nabla_{\theta} (1) = 0.$)

• (Fisher's identity) $\mathcal{J}^* = -\Sigma^*$.

Proof: note that

$\forall \theta \in \mathbb{R}, \int \partial_{\theta_j} \log P_{\theta}(x) \cdot P_{\theta}(x) dx = 0$

Taking ∂_{θ_k} on both sides, we have

$$\int \partial_{\theta_j} \partial_{\theta_k} \log p_{\theta}(x) \cdot p_{\theta}(x) dx$$

$$+ \int (\partial_{\theta_j} \log p_{\theta}(x) \cdot \partial_{\theta_k} \log p_{\theta}(x)) \cdot p_{\theta}(x) dx = 0$$

holding true for any $j, k \in \{1, 2, \dots, d\}$.

Putting together in matrix form, we have

$$-\mathcal{J}^* = \left[-\mathbb{E}_{\theta} [\partial_{\theta_j} \partial_{\theta_k} \log p_{\theta}(x)] \right]_{j, k \in \{d\}} = \left[\mathbb{E}_{\theta} [\partial_{\theta_j} \log p_{\theta}(x) \cdot \partial_{\theta_k} \log p_{\theta}(x)] \right]_{j, k} \\ = \Sigma^*.$$

"Fisher information matrix"

$$\mathcal{I}(\theta) = \mathbb{E}_{\theta} [\nabla \log p_{\theta}(x) \nabla \log p_{\theta}(x)^T] \\ = -\mathbb{E}_{\theta} [\nabla_{\theta}^2 \log p_{\theta}(x)]$$

$$\mathcal{J}_n (\hat{\theta}_n - \theta^*) \xrightarrow{d} \mathcal{N}(0, \mathcal{I}(\theta)^{-1})$$

Remark.

- $\mathcal{I}(\theta)$ is a measure of information about θ contained in X
- If $\mathcal{I}(\theta)$ is not invertible, this theory

does not apply, but MLE may still exist
(and it may still be consistent)

- For CI, we just need to estimate $I(\theta)$
(using either one of the representations).

Asymptotic optimality.

Under certain regularity conditions,
among all the possible estimators for θ ,
MLE has the smallest possible MSE, asymptotically

More formally, if $\tilde{\theta}_n$ is another estimator,

$$\text{st. } E_{\theta} [(\tilde{\theta}_n - \theta)^2] \xrightarrow{n \rightarrow \infty} a(\theta) \quad \forall \theta \in \Theta.$$

then for almost all $\theta \in \Theta$ (up to a measure-0 set)

$$a(\theta) \geq I(\theta)^{-1}, \text{ or } \text{tr}(I(\theta)^{-1}) \text{ in the multivariate case}$$

"partial evidence of this optimality"

Cramér-Rao lower bound.

Thm (CRLB) If $\hat{\theta}_n$ is unbiased,

then $\forall \theta \in \mathbb{R}$,

$$\text{var}_{\theta}(\hat{\theta}_n) \geq \frac{1}{n} \text{tr}(I(\theta)^{-1})$$

(only for unbiased estimators,

for biased estimator, asymptotic LB above)

Proof (l-D for simplicity).

$$\forall \theta \in \mathbb{R} \quad \int \nabla_{\theta} \log p_{\theta}(x) \cdot p_{\theta}(x) dx = 0. \quad (\text{by previous derivation}) \quad (1)$$

$$\int \hat{\theta}_n(x) \cdot p_{\theta}(x) dx = \theta \quad (\text{by unbiasedness}) \quad (2)$$

$\frac{d}{d\theta}(z)$, we get

$$\int \hat{\theta}_n(x) \cdot \nabla_{\theta} \log p_{\theta}(x) \cdot p_{\theta}(x) dx = 0 \quad (3)$$

(3) - $\theta \times (1)$.

$$\begin{aligned} 1 &= \int (\hat{\theta}_n(x) - \theta) \cdot \nabla_{\theta} \log p_{\theta}(x) \cdot p_{\theta}(x) dx \\ &= \mathbb{E}_{\theta}[(\hat{\theta}_n - \theta) \cdot (\nabla_{\theta} \log p_{\theta}(x))]. \end{aligned}$$

$$\leq \sqrt{\mathbb{E}_\theta[(\hat{\theta}_n - \theta)^2]} \cdot \sqrt{\mathbb{E}_\theta[\text{Var}_{\theta}(\log p_\theta(x))]}$$

(Cauchy-Schwarz)

Re-arranging, we get $\mathbb{E}[(\hat{\theta}_n - \theta)^2] \geq \frac{1}{\mathbb{E}_\theta[\log p_\theta(x)]^2}$

Above derivations are for joint distribution

of $X = (X_1, X_2, \dots, X_n)$.

$$\begin{aligned} I^{(n)}(\theta) &= \mathbb{E}_\theta[\nabla_{\theta} \log P_\theta^{(n)}(X_1, \dots, X_n) \cdot \nabla_{\theta} \log P_\theta^{(n)}(X_1, \dots, X_n)^T] \\ &= \sum_{i=1}^n \mathbb{E}_\theta[\nabla_{\theta} \log p_\theta(X_i) \cdot \nabla_{\theta} \log p_\theta(X_i)^T] \\ &= n \cdot I(\theta). \end{aligned}$$

So $\text{Var}_{\theta}(\hat{\theta}_n) \geq \frac{1}{n I(\theta)}$.

Asymptotic normality of MLE requires

strong assumptions (e.g. differentiability, $I(\theta)$ invertible, etc).

We can still have consistency w/o strong conditions.

- Consistency arguments extend beyond parametric models.
- Consistency is "global": do not need to restrict $\hat{\theta}_n$ at a local neighborhood of θ^* .

Thm. Assume that

"ULLN"

$$(i) \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i) - \mathbb{E}_{\theta^*} [\log p_\theta(X_i)] \right| \xrightarrow{P} 0 \quad (\text{under } P_{\theta^*})$$

$$(ii) \forall \varepsilon > 0, \sup_{\theta: |\theta - \theta^*| \geq \varepsilon} \mathbb{E}_{\theta^*} [\log p_\theta(X)] < \mathbb{E}_{\theta^*} [\log p_{\theta^*}(X)]$$

then MLE is consistent, $\hat{\theta}_n \xrightarrow{P} \theta^*$.

Proof: Notations: $f(\theta; x) = \log p_\theta(x)$.

$$F(\theta) := \mathbb{E}_{\theta^*} [f(\theta; X)].$$

$$\hat{F}_n(\theta) := \frac{1}{n} \sum_{i=1}^n f(\theta; X_i).$$

Condition (ii) becomes

$$\sup_{\theta: |\theta - \theta^*| \geq \varepsilon} F(\theta) < F(\theta^*). \quad (*)$$

We only need to study $F(\theta^*) - F(\hat{\theta}_n)$.

$$0 \leq F(\theta^*) - F(\hat{\theta}_n)$$

$$= \underbrace{(F(\theta^*) - \hat{F}_n(\theta^*))}_{\leq \sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)|} + \underbrace{(\hat{F}_n(\theta^*) - \hat{F}_n(\hat{\theta}_n))}_{\leq 0} + \underbrace{(\hat{F}_n(\hat{\theta}_n) - F(\hat{\theta}_n))}_{\leq \sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)|}$$

$$\xrightarrow{P} 0$$

$$\leq 0$$

$$\xrightarrow{P} 0$$

So we have

$$F(\hat{\theta}_n) \xrightarrow{P} F(\theta^*).$$

by (*), we have $\hat{\theta}_n \xrightarrow{P} \theta^*$.

Now about the two conditions:

Condition (ii):

If $F(\theta) := E_{\theta^*}[\log P_\theta(X)]$ is CTS in θ

(ii) $\Leftrightarrow \theta^*$ uniquely maximizes F .

$$F(\theta) = \underbrace{E_{\theta^*}[\log P_{\theta^*}(X)]}_{\text{indp of } \theta} - \underbrace{E_{\theta^*}[\log \frac{P_{\theta^*}}{P_\theta}(X)]}_{D_{KL}(P_{\theta^*} || P_\theta)}$$

Properties of KL divergence:

$D_{KL}(P || Q) \geq 0$, with θ attained
by $P = Q$.