

Midterm : 10/17, 5:10pm - 7:40pm
(2 hrs and 30 min)
in the classroom.

- 4 pages (double-sided) cheat sheet
- No electronics.
- Bring your student ID.
- Covers first 6 lectures
- Format. { 4 questions in total }
 - 1. True or false
 - 2. 3. 4 Calculations / proofs / descriptions.
(show your reasoning).
 - (multiple parts).

From last lecture: consistency of M-estimators (incl. MLE).

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} F_n(\theta) \quad \text{where } F_n(\theta) := \frac{1}{n} \sum_{i=1}^n f(\theta; x_i)$$

$$F(\theta) := \mathbb{E}[f(\theta; x)], \quad \theta^* = \arg \min_{\theta \in \Theta} F(\theta).$$

(for MLE, $f(\theta; x) = -\log p_\theta(x)$).

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

Holds true under the following conditions.

$$(i) \text{ (ULLN). } \sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)| \xrightarrow{P} 0.$$

$$(ii) \forall \varepsilon > 0, \exists \delta > 0 \\ \text{s.t. } F(\theta) \geq F(\theta^*) + \varepsilon$$

$$\text{when } \|\theta - \theta^*\| \geq \delta.$$

$$(\text{Equivalently. } \inf_{\theta: \|\theta - \theta^*\| \geq \delta} F(\theta) > F(\theta^*) \quad \forall \delta > 0)$$

(iii) Holds true when F is a continuous function, Θ is compact.
and θ^* is unique minimizer

(for finite-dim models, compactness \Leftrightarrow bounded and closed.)

Proof: F is cts function on the domain

$$\Theta \setminus \{\theta : \|\theta - \theta^*\| < \delta\}.$$

So $\inf_{\substack{\|\theta - \theta^*\| \geq \delta \\ \theta \in \Theta}} F(\theta)$ is attained by some θ_0

since θ^* is unique m.m., $F(\theta_0) > F(\theta^*)$.

When specialised to MLE:

$$F(\theta) = -\mathbb{E}_{\theta^*} [\log P_\theta(X)]$$

$$= \underbrace{\mathbb{E}_{\theta^*} [\log P_{\theta^*}(X)]}_{\text{const, indep of } \theta} + \underbrace{\mathbb{E}_{\theta^*} \left[\log \frac{P_{\theta^*}(X)}{P_\theta(X)} \right]}_{\text{Kullback-Leibler}}$$

$$D_{KL}(P||Q) = \mathbb{E}_P \left[\log \frac{P(X)}{Q(X)} \right] \quad \text{divergence.}$$

where p, q are densities of P, Q .

$$\cdot D_{KL}(P \parallel Q) \geq 0 \quad (\forall P, Q)$$

$$D_{KL}(P \parallel Q) \geq \frac{1}{2} \left(\int |P(x) - Q(x)| dx \right)^2$$

(Pinsker's inequality).

$$D_{KL}(P \parallel Q) = 0 \text{ if and only if } P = Q.$$

For MLE

$$F(\theta) = \text{constant} + D_{KL}(P_{\theta^*} \parallel P_0)$$

minimized at $\theta = \theta^*$

When the model is "identifiable", θ^* uniquely minimizes F .
 (If $\theta_1 \neq \theta_2$ then $P_{\theta_1} \neq P_{\theta_2}$)

On identifiability:

$$\text{e.g. } P_\theta = \frac{1}{2} N(\theta, 1) + \frac{1}{2} N(-\theta, 1).$$

$\theta \in \mathbb{R}$, then it's not identifiable

$$P_\theta = P_{-\theta}$$

④ $\{\theta : \theta \geq 0\}$: the model is identifiable.

Still need to show (i)

Thm (Wald) Suppose Θ is compact

(closed and bounded in \mathbb{R}^d)

Assuming:

- $E\left[\sup_{\theta \in \Theta} |f(\theta; x)|\right] < \infty$ (from DCT)
- $f(\theta; x)$ is a ces function of θ , $\forall x$.

then $\sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)| \xrightarrow{P} 0$.

Rmk: . Holds true under very general setting.

Θ doesn't even need to be finite-dim.

• Intuition :

If Θ is finite,

$$P\left(\sup_{\theta \in \Theta} |F_n(\theta) - F(\theta)| > \varepsilon\right)$$

$$\leq \sum_{\theta \in \Theta} P(|F_n(\theta) - F(\theta)| > \varepsilon) \rightarrow 0$$

Compactness + Continuity allows us
to use discrete approximations to \mathbb{H} .

Implementation for MLE:

Identifiability + log-likelihood + $\mathbb{E}[\sup_{\theta \in \Theta}] < \infty$

$$\Rightarrow \hat{\theta}_n \xrightarrow{P} \theta^*$$

furthermore, (from last lecture)
under additional smoothness assumptions $I(\theta^*) > 0$.

then $\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{d} N(0, I(\theta^*))$

(Since we can show MLE consistency,
we don't need to restrict in a local neighborhood
of θ^*).

Jackknife.

Motivation: in many cases, the estimator
satisfies asymptotic expansion

$$\widehat{\theta}_n - \theta^* = \frac{1}{n} \sum_{i=1}^n \beta_i + \frac{a}{n} + \frac{b}{n^2} + O_p\left(\frac{1}{n^3}\right)$$

$(\beta_i \stackrel{iid}{\sim}, \mathbb{E}[\beta] = 0, \mathbb{E}|\beta|^2 < \infty)$.

$$\left| \frac{1}{n} \sum_{i=1}^n \beta_i \right| \asymp O\left(\frac{1}{\sqrt{n}}\right)$$

Jackknife bias estimator

Let $\widehat{\theta}_{(-i)}$ be the estimator

computed from $(X_j)_{j \neq i}$. ($i = 1, 2, \dots, n$)

$$\widehat{\theta}_{(-i)} - \theta^* = \frac{1}{n-1} \sum_{j \neq i} \beta_j + \frac{a}{n-1} + \frac{b}{(n-1)^2} + O_p\left(\frac{1}{n^3}\right)$$

"leave-one-out" estimators.

$$\widehat{\delta}_{\text{jack}} = \frac{n-1}{n} \sum_{i=1}^n \widehat{\theta}_{(-i)} - (n-1) \widehat{\theta}_n$$

$$\begin{aligned} &= \left(\frac{n-1}{n} \sum_{i=1}^n \beta_i + \frac{(n-1)}{n} \cdot n \cdot \frac{a}{n-1} + O_p\left(\frac{1}{n^2}\right) \right) \quad \text{ignore} \\ &\quad - \left(\frac{n-1}{n} \sum_{i=1}^n \beta_i + \frac{n-1}{n} \cdot a + O_p\left(\frac{1}{n^2}\right) \right) \end{aligned}$$

$\widehat{\delta}_{\text{jack}}$ as estimator for $\frac{a}{n}$.

Jackknife debiasing.

$$\widehat{\theta}_{\text{debias}} = \widehat{\theta}_n - \widehat{b}_{\text{jack}}.$$

- useful in some high-dim applications.
- expansion holds true only w/ $n \rightarrow \infty$
for finite n , $\widehat{\theta}_{\text{debias}}$ may still be biased.

Computational algorithms.

- Newton's method.

$$\min_{\theta \in \mathbb{R}^d} F(\theta).$$

Start from $\theta^{(0)}$, iterate

$$\theta^{(t+1)} = \theta^{(t)} - \nabla^2 F(\theta^{(t)})^{-1} \nabla F(\theta^{(t)})$$

($t = 0, 1, 2, \dots$).

. Local super-exponential convergence

One-step Newton estimator:

Start from $\widetilde{\theta}_n$

$$\widehat{\theta}_n = \widetilde{\theta}_n - \nabla^2 F_n(\widetilde{\theta}_n)^{-1} \cdot \nabla F_n(\widetilde{\theta}_n).$$

Fact. If $\|\widetilde{\theta}_n - \theta^*\| = O_p\left(\frac{1}{\sqrt{n}}\right)$

(this holds true when $J_n(\widetilde{\theta}_n - \theta^*) \xrightarrow{d} N(0, \Sigma_0)$ for some potentially large Σ_0).

then $\sqrt{n}(\widehat{\theta}_n - \theta^*) \xrightarrow{d} N(0, H_{\theta^*}^{-1} \Sigma_{\theta^*} H_{\theta^*}^{-1})$

where $H_{\theta^*} = \nabla^2 F(\theta^*)$

$$\Sigma_{\theta^*} = \text{cov}(\nabla f(\theta^*, X))$$

(the same limiting distribution as M-estimator).

Proof. $\widetilde{\theta}_n - \nabla^2 F_n(\widetilde{\theta}_n)^{-1} \nabla F_n(\widetilde{\theta}_n)$

$$\approx \underbrace{\widetilde{\theta}_n}_{\text{cancellation}} - \underbrace{\nabla^2 F_n(\theta^*)^{-1} \cdot (\nabla F_n(\theta^*) + \nabla^2 F_n(\theta^*) \cdot (\widetilde{\theta}_n - \theta^*))}_{\text{cancellation}}$$

$$\approx \theta^* - \underbrace{\nabla^2 F(\theta^*)^{-1} \cdot \nabla F_n(\theta^*)}_{\text{cancellation}}$$

$$J_n(\dots) \xrightarrow{d} N(0, H_{\theta^*}^{-1} \Sigma_{\theta^*} H_{\theta^*}^{-1}).$$

• SGD.

$$\theta_{t+1} = \theta_t - \alpha_{t+1} \nabla f(\theta_t; X_{t+1})$$

(for $t=0, 1, 2, \dots, n-1$)

$(\alpha_t)_{t \geq 1}$ step sequence (tuning parameters)

(Polyak - Juditsky - Ruppert)

$$\hat{\theta}_n := \frac{1}{n} \sum_{t=1}^n \theta_t.$$

Fact. Under certain regularity conditions (incl convexity)

with suitable stepsize sequence,

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{d} N(0, H_*^{-1} \Sigma_* H_*^{-1})$$

(same as M-estimator).

Sufficient statistics.

Def. $(P_\theta : \theta \in \Theta)$ is a class of distributions

We call $T(X)$ sufficient if

H_θ and H_T conditional distribution of X

under P_θ conditioned on $T(X) = t$

is indep of θ .

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

$$T(X_1, X_2, \dots, X_n) = X_1 + X_2 + \dots + X_n \sim \text{Binom}(n, p)$$

Conditionally on $T(X_1, \dots, X_n) = t$

(X_1, \dots, X_n) are indicators of random subset
with cardinality t , indep of p .

Theorem (factorization)

T is sufficient $\Leftrightarrow \exists g_\theta, h$
s.t. $P_\theta(x) = g_\theta(T(x)) \cdot h(x)$
 $(\forall \theta, x)$

e.g. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, 1)$

$$P_\mu(x_1, \dots, x_n) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (\mu - x_i)^2\right)$$

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

So $T(x_i) = \sum_{i=1}^n x_i$ is sufficient.

Proof. " \Leftarrow "

$$P_\theta(x | T(X)=t) = \frac{g_\theta(t) \cdot h(x) \cdot \mathbb{1}_{\{T(x)=t\}}}{\int_{T(z)=t} g_\theta(z) \cdot h(z) dz}$$

indp of θ .

" \Rightarrow ". Construction: $g_\theta(t) = P_\theta(T(X)=t)$

(density of $T(X)$ when $X \sim P_\theta$)

$$h(x) = P_{\theta_0}(X=x | T(X)=t) (= P_\theta(X=x | T(X)=t))$$

(conditional density of X conditioned on $T(X)=t$
under model $X \sim P_{\theta_0}$).

Fact. For any estimator $\hat{\theta}_n$ of θ .

\exists another estimator $\tilde{\theta}_n$ which depends on X
only through $T(X)$

such that $\hat{\theta}_n \stackrel{d}{=} \tilde{\theta}_n$ under $P_\theta(\theta)$.

Proof: $X | T(X)$ is indep of θ .

So we can sample from this conditional distribution.

- Sample $\tilde{X} \sim p(x | T(X) = t)$
- Compute $\hat{\theta}_n$ on \tilde{X} , and let $\tilde{\theta}_n$ be the result.

Improving an estimator using sufficient stats.

Rao-Blackwellization.

Given estimator $\hat{\theta}_n$ and sufficient stats $T(X)$

$$\hat{\theta}_{RB} := \mathbb{E}[\hat{\theta}_n | T(X)].$$

$$\mathbb{E}[\hat{\theta}_{RB}] = \mathbb{E}_{\theta}[\hat{\theta}_n]$$

(If $\hat{\theta}_n$ is unbiased, so is $\hat{\theta}_{RB}$).

$$\text{Var}_{\theta}(\hat{\theta}_{RB}) \leq \text{Var}_{\theta}(\hat{\theta}_n).$$

Exponential family.

Suppose $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} P_\theta$

$$P_\theta(x) = \exp\left(\eta(\theta)^T T(x) - B(\theta)\right) \cdot h(x).$$

For joint distribution

$$P_\theta^{(n)}(x_1, \dots, x_n) = \exp\left(\eta(\theta)^T \cdot \sum_{i=1}^n T(x_i) - n B(\theta)\right) \cdot \prod_{i=1}^n h(x_i).$$

By factorization thm.

$\sum_{i=1}^n T(x_i)$ is sufficient for θ .

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

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

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \cdot \exp\left(\begin{bmatrix} \mu/\sigma^2 \\ -\frac{1}{2\sigma^2} \end{bmatrix}^T \begin{bmatrix} x \\ x^2 \end{bmatrix}\right).$$

$$T(X_1, X_2, \dots, X_n) = \begin{bmatrix} \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i^2 \end{bmatrix} \text{ is sufficient}$$