

Chapter 5: Likelihood estimation JAW Ch.4, Sev. Ch. 4

5.1 Maximum likelihood estimation: basics

(i) Suppose:

(a) **A0:** X_1, \dots, X_n are i.i.d. $P_\theta: \theta \in \Theta \subset \mathfrak{R}^k$

(b) **A1: Identifiability:** $\theta \neq \theta^* \Rightarrow P_\theta \neq P_{\theta^*}$

(c) **A2:** P_θ has density $f(\cdot; \theta)$ w.r.t σ -finite μ .

(d) **A3:** $A = \{x : f(x; \theta) > 0\}$ does not depend on θ .

(ii) Given A0,A1,A2:

the *likelihood* $L_n(\theta) = L(\theta; X^{(n)}) = f_n(X^{(n)}; \theta) = \prod_{i=1}^n f(X_i; \theta)$

The *log-likelihood* $\ell_n(\theta) = \log_e L_n(\theta) = \sum_1^n \log f(X_i; \theta)$

For set $B \subset \Theta$, $\ell_n(B) \equiv \sup_{\theta \in B} \ell_n(\theta)$.

(iii) Given A0,A1,A2: the value $\widehat{\theta}_n$ of θ which maximises the likelihood $L_n(\theta)$, if it exists and is unique, is the maximum likelihood estimator (MLE) of θ . Note $\ell_n(\Theta) = \ell_n(\widehat{\theta}_n)$.

(iv) Given A0-A3, and differentiability of $L_n(\theta)$ the MLE may be found by solving the *likelihood equation* or *score equation*, $\nabla \ell_n(\theta) = 0$. (However, this equation may have no roots in Θ , or multiple roots.)

(v) Basic properties:

(a) the MLE depends only on the minimal sufficient statistic

(b) MLE's are NOT necessarily unbiased – if consistent, then unbiased in the limit ($E(T_n) \rightarrow q(\theta)$, provided $E(T_n) < M < \infty$) (cf TPE “asymptotically unbiased”).

(c) If an unbiased estimator attaining the CRLB exists, it is the MLE.

(d) If $q(\theta)$ is any 1-1 function of θ , $q(\widehat{\theta}) = q(\widehat{\theta})$.

5.2 Kullback-Leibler Information (JAW 4.3)

(i) **Defn:** Let P and Q be probability measures (Q may be sub-prob.meas.) with densities p and q . Then $K(P, Q) \equiv E_P(\log(p(X)/q(X)))$.

(ii) $K(P, Q)$ is well-defined, and ≥ 0 (possibly ∞), and $= 0$ iff $Q = P$. (Proof by Jensen's inequality, or by $\log x \leq (x - 1)$.)

(iii) If A0-A3, the SLLN gives, under P_{θ_0} , for $\theta \neq \theta_0$

$$\frac{1}{n} \log \frac{L(\theta_0 : X^{(n)})}{L(\theta; X^{(n)})} = \frac{1}{n} \sum_1^n \log \frac{P_{\theta_0}(X_i)}{P_{\theta}(X_i)} \rightarrow_{a.s.} K(P_{\theta_0}, P_{\theta}) > 0$$

(iv) Thus $P_{\theta_0}(L(\theta_0 : X^{(n)}) > L(\theta; X^{(n)})) \rightarrow 1$ as $n \rightarrow \infty$. This motivates the definition of $\widehat{\theta}_n$, but

“Likelihood is a pointwise function on Θ ”

To proceed we need some metric/uniformity/smoothness w.r.t θ .

(iv) **A4:** $\Theta \supset \Theta_0$, an open set in \mathfrak{R}^k , and for $\theta \in \Theta_0$:

(a) $\ell(\theta; x) \equiv \log p_{\theta}(x)$ is twice continuously diffble in θ (μ (a.e.) x)

(b) μ (a.e.) x , third order derivatives exist, with $\frac{\partial^3 \ell}{\partial \theta_j \partial \theta_l \partial \theta_u}$ bounded by $M_{jlu}(x)$ for all $\theta \in \Theta_0$, and $E_{\theta_0}(M_{jlu}(X)) < \infty$ for all $j, l, u = 1, \dots, k$

(v) **A5:** (a) $E_{\theta_0}(\nabla \ell(\theta; X)|_{\theta=\theta_0}) = 0$.

(b) $E_{\theta_0}((\nabla \ell(\theta; X))^t (\nabla \ell(\theta; X))|_{\theta=\theta_0}) = \sum_{j=1}^k \left(\frac{\partial \ell}{\partial \theta_j} \right)^2 < \infty$.

(c) $I(\theta_0) = - \left(E_{\theta_0} \left(\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_l} \right) |_{\theta=\theta_0} \right)$ is positive definite.

5.3 Consistency of MLE (JAW 4.7, Severini 4.2)

(i) **Theorem:** Suppose A0-A5. Then, with probability $\rightarrow 1$ as $n \rightarrow \infty$, \exists solution $\tilde{\theta}_n$ of the likelihood equations s.t. $\tilde{\theta}_n \rightarrow_p \theta_0$ when P_{θ_0} is true.

(ii) For $a > 0$ let $Q_a \equiv \{\theta \in \Theta : |\theta - \theta_0| = a\}$. We show below that, provided $Q_a \subset \Theta_0$, then $P_{\theta_0}(\sup_{\theta \in Q_a} \ell(\theta) < \ell(\theta_0)) \rightarrow 1$ as $n \rightarrow \infty$. Hence there is a local max, which must be root of the likelihood eqn, inside Q_a .

(Note $\sup_{Q_a} \ell(\theta)$ is attained on Q_a by some $\theta \in Q_a$.)

(iii) Define “observed information” $J(\theta_0) = -\left(\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \Big|_{\theta=\theta_0}\right)$.

(Note that $J(\theta)$, unlike $I(\theta)$, is a r.v.)

(iv) $n^{-1}(\ell(\theta; X^{(n)}) - \ell(\theta_0; X^{(n)}))$

$$\begin{aligned} &= n^{-1}(\theta - \theta_0)^t \nabla (\ell(\theta_0)) + \frac{1}{2}(\theta - \theta_0)^t (-n^{-1}J(\theta_0))(\theta - \theta_0) \\ &\quad + (6n)^{-1} \sum_{jlu} (\theta_j - \theta_{0,j})(\theta_l - \theta_{0,l})(\theta_u - \theta_{0,u}) \sum_i \gamma_{jlu}(X_i) M_{jlu}(X_i) \\ &= S_1 + S_2 + S_3 \end{aligned}$$

where $|\gamma_{jlu}| < 1$ (by A4 (b)).

(v) By A5 and WLLN: $S_1 \rightarrow_p 0$, $-2S_2 \rightarrow_p (\theta - \theta_0)^t I(\theta_0)(\theta - \theta_0) \geq \lambda_k a^2$ where λ_k is smallest eigenvalue of $I(\theta_0)$.

$S_3 \rightarrow_p (1/6) \sum_{jlu} (\theta_j - \theta_{0,j})(\theta_l - \theta_{0,l})(\theta_u - \theta_{0,u}) E_{\theta_0}(\gamma_{jlu}(X_1) M_{jlu}(X_1))$ and $|S_3| \leq (1/3)(ka)^3 \sum_{jlu} m_{jlu} \equiv Ba^3$ as $n \rightarrow \infty$.

(vi) For large enough n

$$\begin{aligned} \sup_{\theta \in Q_a} (S_1 + S_2 + S_3) &\leq \sup_{\theta \in Q_a} |S_1 + S_3| + \sup_{\theta \in Q_a} (S_2) \\ &< (k + B)a^3 - \lambda_k a^2 / 4 \\ &< 0 \text{ for small enough } a \end{aligned}$$

5.4 Aymptotic normality and efficiency of $\tilde{\theta}_n$. (Sev. 4.2.2)

(i) Notation

$$Z_n \equiv n^{-\frac{1}{2}} \sum_i \nabla(\ell(\theta_0; X_i)) = n^{-\frac{1}{2}} \nabla \ell_n(\theta_0; X^{(n)}),$$

$$\tilde{\ell}(\theta_0; X) \equiv I^{-1}(\theta_0) \nabla \ell(\theta_0; X) \text{ so } n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}(\theta_0; X_i) = I^{-1}(\theta_0) Z_n.$$

Define $G_n(\epsilon) \equiv \{\tilde{\theta}_n; \nabla \ell_n(\tilde{\theta}_n) = 0, |\tilde{\theta}_n - \theta_0| < \epsilon\}$ non-empty as $n \rightarrow \infty, \forall \epsilon > 0$. (Also note here $I(\theta) \equiv I_1(\theta)$.)

(ii) Theorem

$$(a) \quad (n^{\frac{1}{2}}(\tilde{\theta}_n - \theta_0) - n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}(\theta_0; X_i)) \rightarrow_p 0$$

$$(b) \quad n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}(\theta_0; X_i) \rightarrow_d I^{-1}(\theta_0) Z \equiv D \sim N_k(0, I^{-1}(\theta_0)).$$

(iii) First (b): by CLT $Z_n \rightarrow_d N(0, I(\theta_0))$,

so $n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}(\theta_0; X_i) = I^{-1}(\theta_0) Z_n \rightarrow_d N(0, I^{-1}(\theta_0))$

(iv) On G_n ,

$$0 = n^{-\frac{1}{2}} \nabla \ell_n(\tilde{\theta}_n) = n^{-\frac{1}{2}} \nabla \ell_n(\theta_0) - n^{-1} J_n(\theta_n^*) n^{\frac{1}{2}}(\tilde{\theta}_n - \theta_0)$$

where $|\theta_n^* - \theta_0| < |\tilde{\theta}_n - \theta_0|$.

Or $n^{\frac{1}{2}}(\tilde{\theta}_n - \theta_0) = (n^{-1} J_n(\theta_n^*))^{-1} Z_n$ if $J_n^{-1}(\theta_n^*) \exists$.

(v) $\tilde{\theta}_n \rightarrow_p \theta_0$, so using one-term expansion of 2 nd. deriv, and boundedness of 3 rd., $n^{-1}(J_n(\theta_n^*) - J_n(\theta_0)) \rightarrow_p 0$. By continuity, $(n^{-1} J_n(\theta_n^*))^{-1} \rightarrow_p (E(J_1(\theta_0)))^{-1} = I^{-1}(\theta_0)$. (and $J_n(\theta_n^*)$ is pos def with prob approaching 1).

Now $I^{-1}(\theta_0) Z_n = n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}(\theta_0; X_i)$ hence (a).

(vi) Transforming (ii) to $q(\theta)$: $\dim(q) = k^*, 1 \leq k^* \leq k$

$$(a) \quad (n^{\frac{1}{2}}(q(\tilde{\theta}_n) - q(\theta_0)) - n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}_q(\theta_0; X_i)) \rightarrow_p 0$$

$$(b) \quad n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}_q(\theta_0; X_i) \rightarrow_d N_{k^*}(0, (\nabla q(\theta_0))^t I^{-1}(\theta_0) (\nabla q(\theta_0))).$$

where $\tilde{\ell}_q(\theta_0; X_i) = (\nabla q(\theta_0))^t I^{-1}(\theta_0) (\nabla \ell(\theta_0, X_i))$

5.5 Bits and pieces

5.5.1 Estimation of $I(\theta)$: Suppose we need to estimate $I(\theta_0)$, and have A0-A5, as above, so ℓ_n is twice continuously diffble, and expectations \exists :

(a) $\tilde{\theta}_n \rightarrow_p \theta_0$, so $I(\tilde{\theta}_n) \rightarrow_p I(\theta_0)$, but $I(\cdot)$ can be hard to compute.

(b) $n^{-1} \sum_{i=1}^n (\nabla \ell(\tilde{\theta}_n; X_i)) (\nabla \ell(\tilde{\theta}_n; X_i))^t$ is also a consistent estimator of $I(\theta_0)$, since $\nabla \ell(\tilde{\theta}_n; X_i) \rightarrow_p \nabla \ell(\theta_0; X_i)$.

(c) Often easiest is to use the second derivatives:

$$\left(-n^{-1} \sum_{i=1}^n \frac{\partial^2 \ell(\theta; X_i)}{\partial \theta_j \partial \theta_l} \right) \Big|_{\theta=\tilde{\theta}_n} = \left(n^{-1} J_n(\tilde{\theta}_n) \right)$$

is also a consistent estimator of $I(\theta_0)$.

(d) If CRLB attained: $\nabla \ell_n(\theta) = nI(\theta)(\widehat{\theta}_n - \theta)$
Hence (differentiating), $n^{-1} J_n(\widehat{\theta}_n) = I(\widehat{\theta}_n)$.

5.5.2 The one-step estimator

We want to solve $\nabla \ell_n(\theta; X^{(n)}) = 0$. This can be hard. Suppose we have a preliminary estimator $\bar{\theta}_n$. Then we can do one-step Newton-Raphson:

$$0 = \nabla \ell_n(\theta; X^{(n)}) \approx \nabla \ell_n(\bar{\theta}_n; X^{(n)}) + \left(\frac{\partial^2 \ell_n}{\partial \theta_j \partial \theta_l} \right) \Big|_{\theta=\bar{\theta}_n} (\theta - \bar{\theta}_n)$$

Thus, replacing the second derivatives by some consistent estimator $-\hat{I}$ from (a),(b) or (c) above, new θ_n^* is

$$\theta_n^* = \bar{\theta}_n + (nI(\widehat{\theta}_n))^{-1} \nabla \ell_n(\bar{\theta}_n; X^{(n)})$$

If $n^{1/4}(\bar{\theta}_n - \theta_0) \rightarrow_p 0$ then θ_n^* satisfies same Theorem 5.4(ii) (a) and (b) as $\tilde{\theta}_n$ –see JAW 4.7.

5.6 Appendix: summary of notation

| Notation, definition | description, result, etc. |
|---|--|
| $X^{(n)} = (X_1, \dots, X_n)$ | sample of i.i.d. X_i from pdf f |
| $\ell_n(\theta; X^{(n)}) = \sum_1^n \log f(X_i; \theta)$ | log-likelihood function. |
| $\nabla \ell(\theta; X_i) = \left(\frac{\partial \log f(X_i; \theta)}{\partial \theta_j} \right)$ | contribution of X_i to score |
| $\nabla \ell_n(\theta) = \sum_1^n \nabla \ell(\theta; X_i)$ | the score function: deriv. of ℓ_n |
| $I(\theta) = \mathbb{E} \left(-\frac{\partial^2 \log f(X_i; \theta)}{\partial \theta_j \partial \theta_l} \right)$ | $I_n(\theta) = nI(\theta)$: Fisher Information |
| $J_n(\theta; X^{(n)}) = \left(-\frac{\partial^2 \ell_n}{\partial \theta_j \partial \theta_l} \right)$ | observed information |
| | $n^{-1} J_n(\theta) \rightarrow_p I(\theta)$ |
| $\theta_0 \in \Theta \subset \mathfrak{R}^k$ | true value of θ |
| $\nabla \ell_n(\theta) = 0$ | the likelihood equation |
| $\widehat{\theta}_n$; MLE | maximizes $\ell_n(\theta)$ in Θ |
| $\tilde{\theta}_n$ | root of the likelihood eqn |
| $Z_n \equiv n^{-\frac{1}{2}} \nabla \ell_n(\theta_0; X^{(n)})$ | $Z_n \rightarrow_d N(0, I(\theta_0))$ |
| $D_n = I^{-1}(\theta_0) Z_n$ | $D_n \rightarrow_d D \sim N_k(0, I^{-1}(\theta_0))$ |
| $\tilde{\ell}(\theta_0; X) \equiv I^{-1}(\theta_0) \nabla \ell(\theta_0; X)$ | $n^{-\frac{1}{2}} \sum_{i=1}^n \tilde{\ell}(\theta_0; X_i) = I^{-1}(\theta_0) Z_n$ |
| $\Delta_n \equiv n^{\frac{1}{2}}(\tilde{\theta}_n - \theta_0)$ | $(\Delta_n - D_n) \rightarrow_p 0$ |
| $\tilde{\ell}_n(\theta, X^{(n)}) = \sum_{i=1}^n \tilde{\ell}(\theta; X_i)$ | influence fn for θ |
| $\tilde{\ell}_q(\theta; X^{(n)}) = (\nabla q(\theta))^t \tilde{\ell}_n(\theta; X^{(n)})$ | influence fn for q . |