

## Statistics 581

### Revision of Section 4.4: Consistency of Maximum Likelihood Estimates

Wellner; 11/30/2001

#### Some Uniform Strong Laws of Large Numbers

Suppose that:

- A.  $X, X_1, \dots, X_n$  are i.i.d.  $P$  on the measurable space  $(\mathcal{X}, \mathcal{A})$ .
- B. For each  $\theta \in \Theta$ ,  $f(x, \theta)$  is a measurable, real-valued function of  $x$ ,  $f(\cdot, \theta) \in L_1(P)$ .

Let  $\mathcal{F} = \{f(\cdot, \theta) : \theta \in \Theta\}$ . Since  $f(\cdot, \theta) \in L_1(P)$  for each  $\theta$ ,

$$g(\theta) \equiv Ef(X, \theta) = \int f(x, \theta) dP(x) \equiv Pf(\cdot, \theta)$$

exists and is finite. Moreover, by the strong law of large numbers,

$$\begin{aligned} \mathbb{P}_n f(\cdot, \theta) &\equiv \int f(x, \theta) d\mathbb{P}_n(x) = \frac{1}{n} \sum_{i=1}^n f(X_i, \theta) \\ (0.1) \quad &\rightarrow_{a.s.} Ef(X, \theta) = Pf(\cdot, \theta) = g(\theta). \end{aligned}$$

It is often useful and important to strengthen (0.1) to hold uniformly in  $\theta \in \Theta$ :

$$(0.2) \quad \sup_{\theta \in \Theta} |\mathbb{P}_n f(\cdot, \theta) - Pf(\cdot, \theta)| \rightarrow_{a.s.} 0.$$

Note that the left side in (0.2) is equal to

$$\|\mathbb{P}_n - P\|_{\mathcal{F}} \equiv \sup_{f \in \mathcal{F}} |\mathbb{P}_n f - Pf|.$$

Here is how (0.2) can be used: suppose that we have a sequence  $\widehat{\theta}_n$  of estimators, possibly dependent on  $X_1, \dots, X_n$ , such that  $\widehat{\theta}_n \rightarrow_{a.s.} \theta_0$ . Suppose that  $g(\theta)$  is continuous at  $\theta_0$ . We would like to conclude that

$$(0.3) \quad \mathbb{P}_n f(\cdot, \widehat{\theta}_n) = \frac{1}{n} \sum_{i=1}^n f(X_i, \widehat{\theta}_n) \rightarrow_{a.s.} g(\theta_0).$$

The convergence (0.3) does not follow from (0.1); but (0.3) does follow from (0.2):

$$\begin{aligned} |\mathbb{P}_n f(\cdot, \widehat{\theta}_n) - g(\theta_0)| &\leq |\mathbb{P}_n f(\cdot, \widehat{\theta}_n) - g(\widehat{\theta}_n)| + |g(\widehat{\theta}_n) - g(\theta_0)| \\ &\leq \sup_{\theta \in \Theta} |\mathbb{P}_n f(\cdot, \theta) - g(\theta)| + |g(\widehat{\theta}_n) - g(\theta_0)| \\ &= \|\mathbb{P}_n - P\|_{\mathcal{F}} + |g(\widehat{\theta}_n) - g(\theta_0)| \\ &\rightarrow_{a.s.} 0 + 0 = 0. \end{aligned}$$

The following theorems, due to Le Cam, give conditions on  $f$  and  $P$  under which (2) holds. The first theorem is a prototype for what are now known in empirical process theory as “Glivenko-Cantelli theorems”.

**Theorem 1.** Suppose that:

- (a)  $\Theta$  is compact.
- (b)  $f(x, \cdot)$  is continuous in  $\theta$  for all  $x$ .
- (c) There exists a function  $F(x)$  such that  $EF(X) < \infty$  and  $|f(x, \theta)| \leq F(x)$  for all  $x \in \mathcal{X}, \theta \in \Theta$ .

Then (0.2) holds; i.e.

$$\sup_{\theta \in \Theta} |\mathbb{P}_n f(\cdot, \theta) - Pf(\cdot, \theta)| \rightarrow_{a.s.} 0.$$

The second theorem is a “one-sided” version of theorem 1 which is useful for the theory of maximum likelihood estimation.

**Theorem 2.** Suppose that:

- (a)  $\Theta$  is compact.
- (b)  $f(x, \cdot)$  is upper 1 in  $\theta$  for all  $x$ .
- (c) There exists a function  $F(x)$  such that  $EF(X) < \infty$  and  $f(x, \theta) \leq F(x)$  for all  $x \in \mathcal{X}, \theta \in \Theta$ .
- (d) For all  $\theta$  and all sufficiently small  $\rho > 0$

$$\sup_{|\theta' - \theta| < \rho} f(x, \theta')$$

is measurable in  $x$ .

Then

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} \mathbb{P}_n f(\cdot, \theta) \leq_{a.s.} \sup_{\theta \in \Theta} Pf(\cdot, \theta) = \sup_{\theta \in \Theta} g(\theta).$$

We proceed by first proving Theorem 2. Then Theorem 1 will follow as a consequence of Theorem 2.

**Proof of Theorem 2.** Let

$$\psi(x, \theta, \rho) \equiv \sup_{|\theta' - \theta| < \rho} f(x, \theta').$$

Then  $\psi$  is measurable (for  $\rho$  sufficiently small), bounded by an integrable function  $F$ , and

$$\psi(x, \theta, \rho) \searrow f(x, \theta) \quad \text{as} \quad \rho \searrow 0 \quad \text{by (b).}$$

Thus by the monotone convergence theorem

$$\int \psi(x, \theta, \rho) dP(x) \searrow \int f(x, \theta) dP(x) = g(\theta).$$

Let  $\epsilon > 0$ . For each  $\theta$ , find  $\rho_\theta$  so that

$$\int \psi(x, \theta, \rho) dP(x) < g(\theta) + \epsilon.$$

The spheres

$$S(\theta, \rho_\theta) = \{\theta' : |\theta' - \theta| < \rho_\theta\}$$

cover  $\Theta$ , so by (a) there exists a finite sub cover:  $\Theta \subset \bigcup_{j=1}^m S(\theta_j, \rho_{\theta_j})$ . for each  $\theta \in \Theta$  there is some  $j$ ,  $1 \leq j \leq m$ , such that  $\theta \in S(\theta_j, \rho_{\theta_j})$ ; hence from the definition of  $\psi$  it follows that

$$f(x, \theta) \leq \psi(x, \theta_j, \rho_{\theta_j})$$

for all  $x$ . Therefore

$$\mathbb{P}_n f(\cdot, \theta) \leq \mathbb{P}_n \psi(\cdot, \theta_j, \rho_{\theta_j}),$$

and hence

$$\begin{aligned} \sup_{\theta \in \Theta} \mathbb{P}_n f(\cdot, \theta) &\leq \sup_{1 \leq j \leq m} \mathbb{P}_n \psi(\cdot, \theta_j, \rho_{\theta_j}) \\ &\xrightarrow{a.s.} \sup_{1 \leq j \leq m} P \psi(\cdot, \theta_j, \rho_{\theta_j}) \\ &\leq \sup_{1 \leq j \leq m} g(\theta_j) + \epsilon \\ &\leq \sup_{\theta \in \Theta} g(\theta) + \epsilon. \end{aligned}$$

Hence

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} \mathbb{P}_n f(\cdot, \theta) \leq_{a.s.} \sup_{\theta \in \Theta} g(\theta) + \epsilon.$$

Letting  $\epsilon \downarrow 0$  completes the proof. □

**Proof of Theorem 1.** Since  $f$  is continuous in  $\theta$ , condition (d) of Theorem 2 is satisfied: for any countable set  $D$  dense in  $\{\theta' : |\theta' - \theta| < \rho\}$ ,

$$\sup_{|\theta' - \theta| < \rho} f(x, \theta') = \sup_{\theta' \in D} f(x, \theta')$$

where the right side is measurable since it is a countable supremum of measurable functions. Furthermore,  $g(\theta)$  is continuous in  $\theta$ :

$$\lim_{\theta' \rightarrow \theta} g(\theta) = \lim_{\theta' \rightarrow \theta} \int f(x, \theta') dP(x) = \int f(x, \theta) dP(x)$$

by the dominated convergence theorem. Now Theorem 1 follows from Theorem 2 applied to the functions  $h(x, \theta) \equiv f(x, \theta) - g(\theta)$  and  $-h(x, \theta)$ : by Theorem 2 applied to  $\{h(x, \theta) : \theta \in \Theta\}$ ,

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} (\mathbb{P}_n f(\cdot, \theta) - g(\theta)) \leq 0 \quad \text{a.s.}$$

By Theorem 2 applied to  $\{-h(x, \theta) : \theta \in \Theta\}$ ,

$$\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} (g(\theta) - \mathbb{P}_n f(\cdot, \theta)) \leq 0 \quad \text{a.s.}$$

The conclusion of Theorem 1 follows since

$$\begin{aligned} 0 &\leq \sup_{\theta \in \Theta} |\mathbb{P}_n f(\cdot, \theta) - g(\theta)| \\ &= \sup_{\theta \in \Theta} (\mathbb{P}_n f(\cdot, \theta) - g(\theta)) \vee \sup_{\theta \in \Theta} (g(\theta) - \mathbb{P}_n f(\cdot, \theta)). \end{aligned}$$

□

For our application of Theorem 2 to consistency of maximum likelihood, the following Lemma will be useful.

**Lemma 1.** If the conditions of Theorem 2 hold, then  $g(\theta)$  is upper-semicontinuous: i.e.

$$\limsup_{\theta' \rightarrow \theta} g(\theta') \leq g(\theta).$$

**Proof.** Since  $f(x, \theta)$  is upper semicontinuous,

$$\limsup_{\theta' \rightarrow \theta} f(x, \theta') \leq f(x, \theta) \quad \text{for all } x;$$

i.e.

$$\liminf_{\theta' \rightarrow \theta} \{f(x, \theta) - f(x, \theta')\} \geq 0 \quad \text{for all } x.$$

Hence it follows by Fatou's lemma that

$$\begin{aligned} 0 &\leq \text{Eliminf}_{\theta' \rightarrow \theta} \{f(X, \theta) - f(X, \theta')\} \\ &\leq \liminf_{\theta' \rightarrow \theta} E \{f(X, \theta) - f(X, \theta')\} \\ &= Ef(X, \theta) - \limsup_{\theta' \rightarrow \theta} Ef(X, \theta'); \end{aligned}$$

i.e.

$$\limsup_{\theta' \rightarrow \theta} Ef(X, \theta') \leq Ef(X, \theta) = g(\theta).$$

□

Now we are prepared to tackle consistency of maximum likelihood estimates.

**Theorem 3. (Wald, 1949).** Suppose that  $X, X_1, \dots, X_n$  are i.i.d.  $P_{\theta_0}$ ,  $\theta_0 \in \Theta$  with density  $p(x, \theta_0)$  with respect to the dominating measure  $\nu$ , and that:

- (a)  $\Theta$  is compact.
- (b)  $p(x, \cdot)$  is upper semi-continuous in  $\theta$  for all  $x$ .
- (c) There exists a function  $F(x)$  such that  $EF(X) < \infty$  and

$$f(x, \theta) \equiv \log p(x, \theta) - \log p(x, \theta_0) \leq F(x)$$

for all  $x \in \mathcal{X}$ ,  $\theta \in \Theta$ .

- (d) For all  $\theta$  and all sufficiently small  $\rho > 0$

$$\sup_{|\theta' - \theta| < \rho} p(x, \theta')$$

is measurable in  $x$ .

- (e)  $p(x, \theta) = p(x, \theta_0)$  a.e.  $\nu$  implies that  $\theta = \theta_0$ .

Then for any sequence of maximum likelihood estimates  $\hat{\theta}_n$  of  $\theta_0$ ,

$$\hat{\theta}_n \rightarrow_{a.s.} \theta_0.$$

**Proof.** Let  $\rho > 0$ . The functions  $\{f(x, \theta) : \theta \in \Theta\}$  satisfy the conditions of theorem 2. But we will apply Theorem 2 with  $\Theta$  replaced by the subset

$$S \equiv \{\theta : |\theta - \theta_0| \geq \rho\} \subset \Theta.$$

Then  $S$  is compact, and by Theorem 2

$$P_{\theta_0} \left( \limsup_{n \rightarrow \infty} \sup_{\theta \in S} \mathbb{P}_n f(\cdot, \theta) \leq \sup_{\theta \in S} g(\theta) \right) = 1$$

where

$$\begin{aligned} g(\theta) &= E_{\theta_0} f(X, \theta) = E_{\theta_0} \left\{ \log \frac{p(X, \theta)}{p(X, \theta_0)} \right\} \\ &= -K(P_{\theta_0}, P_{\theta}) < 0 \quad \text{for } \theta \in S. \end{aligned}$$

Furthermore by the Lemma,  $g(\theta)$  is upper semicontinuous and hence achieves its supremum on the compact set  $S$ . Let  $\delta = \sup_{\theta \in S} g(\theta)$ . Then by Lemma 4.1.2 it follows that  $\delta < 0$  and we have

$$P_{\theta_0} \left( \limsup_{n \rightarrow \infty} \sup_{\theta \in S} \mathbb{P}_n f(\cdot, \theta) \leq \delta \right) = 1.$$

Thus with probability 1 there exists an  $N$  such that for all  $n > N$

$$\sup_{\theta \in S} \mathbb{P}_n f(\cdot, \theta) \leq \delta/2 < 0.$$

But

$$\begin{aligned} \mathbb{P}_n f(\cdot, \hat{\theta}_n) &= \sup_{\theta \in \Theta} \mathbb{P}_n f(\cdot, \theta) \\ &= \sup_{\theta \in \Theta} \frac{1}{n} \{l_n(\theta) - l_n(\theta_0)\} \geq 0. \end{aligned}$$

Hence  $\hat{\theta}_n \notin S$  for  $n > N$ ; that is,  $|\hat{\theta}_n - \theta_0| < \rho$  with probability 1. Since  $\rho$  was arbitrary,  $\hat{\theta}_n$  is a.s. consistent.

**Remark 3.** Theorem 3 is due to Wald (1949). The present writeup is an adaptation of Chapters 16 and 17 of Ferguson (1996). For further Glivenko - Cantelli theorems, see chapter 2.4 of Van der Vaart and Wellner (1996).

### References:

- Le Cam, L. (1953). On some asymptotic properties of maximum likelihood estimates and related estimates. *Univ. Calif. Publ. in Statist.* **1**, 277 - 330.
- Ferguson, T. (1996). *A Course in Large Sample Theory*. Chapman and Hall, London.
- Van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer-Verlag, New York.
- Wald, A. (1949). Note on the consistency of the maximum likelihood estimate. *Ann. Math. Statist.* **20**, 595 - 601.