

NONLINEAR LEAST-SQUARES ESTIMATION

DAVID POLLARD AND PETER RADCHENKO

ABSTRACT. The paper uses empirical process techniques to study the asymptotics of the least-squares estimator for the fitting of a nonlinear regression function. By combining and extending ideas of Wu and Van de Geer, it establishes new consistency and central limit theorems that hold under only second moment assumptions on the errors. An application to a delicate example of Wu's illustrates the use of the new theorems, leading to a normal approximation to the least-squares estimator with unusual logarithmic rescalings.

1. INTRODUCTION

Consider the model where we observe y_i for $i = 1, \dots, n$ with

$$(1) \quad y_i = f_i(\theta) + u_i \quad \text{where } \theta \in \Theta.$$

The unobserved f_i can be random or deterministic functions. The unobserved errors u_i are random with zero means and finite variances. The index set Θ might be infinite dimensional. Later in the paper it will prove convenient to also consider triangular arrays of observations.

Think of $f(\theta) = (f_1(\theta), \dots, f_n(\theta))'$ and $u = (u_1, \dots, u_n)'$ as points in \mathbb{R}^n . The model specifies a surface $M_\Theta = \{f(\theta) : \theta \in \Theta\}$ in \mathbb{R}^n . The vector of observations

Date: Original version April 1995; revised March 2002; revised 2 December 2003.

2000 Mathematics Subject Classification. Primary 62E20. Secondary: 60F05, 62G08, 62G20.

Key words and phrases. Nonlinear least squares, empirical processes, subgaussian, consistency, central limit theorem.

$y = (y_1, \dots, y_n)'$ is a random point in \mathbb{R}^n . The least squares estimator (LSE) $\hat{\theta}_n$ is defined to minimize the distance of y to M_Θ ,

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} |y - f(\theta)|^2,$$

where $|\cdot|$ denotes the usual Euclidean norm on \mathbb{R}^n . Many authors have considered the behavior of $\hat{\theta}_n$ as $n \rightarrow \infty$ when the y_i are generated by the model for a fixed θ_0 in Θ .

When the f_i are deterministic, it is natural to express assertions about convergence of $\hat{\theta}_n$ in terms of the n -dimensional Euclidean distance $\kappa_n(\theta_1, \theta_2) := |f(\theta_1) - f(\theta_2)|$. For example, Jennrich (1969) took Θ to be a compact subset of \mathbb{R}^p , the errors $\{u_i\}$ to be iid with zero mean and finite variance, and the f_i to be continuous functions in θ . He proved strong consistency of the least squares estimator under the assumption that $n^{-1}\kappa_n(\theta_1, \theta_2)^2$ converges uniformly to a continuous function that is zero if and only if $\theta_1 = \theta_2$. He also gave conditions for asymptotic normality.

Under similar assumptions Wu (1981, Theorem 1) proved that existence of a consistent estimator for θ_0 implies that

$$(2) \quad \kappa_n(\theta) := \kappa_n(\theta, \theta_0) \rightarrow \infty \quad \text{at each } \theta \neq \theta_0.$$

If Θ is finite, the divergence (2) is also a sufficient condition for the existence of a consistent estimator (Wu 1981, Theorem 2). His main consistency result (his Theorem 3) may be reexpressed as a general convergence assertion.

Theorem 1. *Suppose the $\{f_i\}$ are deterministic functions indexed by a subset Θ of \mathbb{R}^p . Suppose also that $\sup_i \operatorname{var}(u_i) < \infty$ and $\kappa_n(\theta) \rightarrow \infty$ at each $\theta \neq \theta_0$. Let S be a bounded subset of $\Theta \setminus \{\theta_0\}$ and let $R_n := \inf_{\theta \in S} \kappa_n(\theta)$. Suppose there exist constants $\{L_i\}$ such that*

- (i) $\sup_{\theta \in S} |f_i(\theta) - f_i(\theta_0)| \leq L_i$ for each i ;
- (ii) $|f_i(\theta_1) - f_i(\theta_2)| \leq L_i |\theta_1 - \theta_2|$ for all $\theta_1, \theta_2 \in S$;
- (iii) $\sum_{i \leq n} L_i^2 = O(R_n^\alpha)$ for some $\alpha < 4$.

Then $\mathbb{P}\{\hat{\theta}_n \notin S \text{ eventually}\} = 1$.

Remark. Assumption (i) implies $\sum_{i \leq n} L_i^2 \geq \kappa_n(\theta)^2 \rightarrow \infty$ for each θ in S , which forces $R_n \rightarrow \infty$.

If Θ is compact and if for each $\theta \neq \theta_0$ there is a neighborhood $S = S_\theta$ satisfying the conditions of the Lemma then $\hat{\theta}_n \rightarrow \theta_0$ almost surely.

Wu's paper was the starting point for several authors. For example, both Lai (1994) and Skouras (2000) generalized Wu's consistency results by taking the functions $f_i(\theta) = f_i(\theta, \omega)$ as random processes indexed by θ . They took the $\{u_i\}$ as a martingale difference sequence, with $\{f_i\}$ a predictable sequence of functions with respect to a filtration $\{\mathcal{F}_i\}$.

Another line of development is typified by the work of Van de Geer (1990) and Van de Geer and Wegkamp (1996). They took $f_i(\theta) = f(x_i, \theta)$, where $\mathcal{F} = \{f_\theta : \Theta\}$ is a set of deterministic functions and the (x_i, u_i) are iid pairs. (In fact they identified Θ with the index set \mathcal{F} .) Under a stronger assumption about the errors, they established rates of convergence of $\kappa_n(\hat{\theta}_n)$ in terms of \mathcal{L}^2 entropy conditions on \mathcal{F} , using empirical process methods that were developed after Wu's work.

The stronger assumption was that the errors are uniformly subgaussian. In general, we say that a random variable W has a subgaussian distribution if there exists some finite τ such that

$$\mathbb{P} \exp(tW) \leq \exp\left(\frac{1}{2}\tau^2 t^2\right) \quad \text{for all } t \in \mathbb{R}.$$

We write $\tau(W)$ for the smallest such τ . Van de Geer assumed that $\sup_i \tau(u_i) < \infty$.

Remark. Notice that we must have $\mathbb{P}W = 0$ when W is subgaussian because the linear term in the expansion of $\mathbb{P} \exp(tW)$ must vanish. When $\mathbb{P}W = 0$, subgaussianity is equivalent to existence of a finite constant β for which $\mathbb{P}\{|W| \geq x\} \leq 2 \exp(-x^2/\beta^2)$ for all $x \geq 0$.

In our paper we try to bring together the two lines of development. Our main motivation for working on nonlinear least squares was an example presented by Wu (1981, page 507). He noted that his consistency theorem has difficulties with a simple model,

$$(3) \quad f_i(\theta) = \lambda i^{-\mu} \quad \text{for } \theta = (\lambda, \mu) \in \Theta, \text{ a compact subset of } \mathbb{R} \times \mathbb{R}^+.$$

For example, condition (2) does not hold for $\theta_0 = (0, 0)$ at any θ with $\mu > 1/2$. When $\theta_0 = (\lambda_0, 1/2)$, Wu's method fails in a more subtle way, but Van de Geer (1990)'s method would work if the errors satisfied the subgaussian assumption. In Section 4, under only second moment assumptions on the errors, we establish weak consistency and a central limit theorem.

The main idea behind all the proofs—ours, as well as those of Wu and Van de Geer—is quite simple. The LSE also minimizes the random function

$$(4) \quad G_n(\theta) := |y - f(\theta)|^2 - |u|^2 = \kappa_n(\theta)^2 - 2Z_n(\theta)$$

where $Z_n(\theta) := u'f(\theta) - u'f(\theta_0)$.

In particular, $G_n(\hat{\theta}_n) \leq G_n(\theta_0) = 0$, that is, $\frac{1}{2}\kappa_n(\hat{\theta}_n)^2 \leq Z_n(\hat{\theta}_n)$. For every subset S of Θ ,

$$(5) \quad \mathbb{P}\{\hat{\theta}_n \in S\} \leq \mathbb{P}\{\exists \theta \in S : Z_n(\theta) \geq \frac{1}{2}\kappa_n(\theta)^2\} \leq 4\mathbb{P}\sup_{\theta \in S} |Z_n(\theta)|^2 / \inf_{\theta \in S} \kappa_n(\theta)^4.$$

The final bound calls for a maximal inequality for Z_n .

Our methods for controlling Z_n are similar in spirit to those of Van de Geer. Under her subgaussian assumption, for every class of real functions $\{g_\theta : \theta \in \Theta\}$, the process

$$(6) \quad X(\theta) = \sum_{i \leq n} u_i g_i(\theta)$$

has subgaussian increments. Indeed, if $\tau(u_i) \leq \tau$ for all i then

$$\tau \left(X(\theta_1) - X(\theta_2) \right)^2 \leq \sum_{i \leq n} \tau(u_i)^2 \left(g_i(\theta_1) - g_i(\theta_2) \right)^2 \leq \tau^2 |g(\theta_1) - g(\theta_2)|^2.$$

That is, the tails of $X(\theta_1) - X(\theta_2)$ are controlled by the n -dimensional Euclidean distance between the vectors $g(\theta_1)$ and $g(\theta_2)$. This property allowed her to invoke a chaining bound (similar to our Theorem 2) for the tail probabilities of $\sup_{\theta \in S} |Z_n(\theta)|$ for various annuli $S = \{\theta : R \leq \kappa_n(\theta) < 2R\}$.

Under the weaker second moment assumption on the errors, we apply symmetrization arguments to transform to a problem involving a new process $Z_n^\circ(\theta)$ with conditionally subgaussian increments. We avoid Van de Geer's subgaussianity assumption at the cost of extra Lipschitz conditions on the $f_i(\theta)$, analogous to Assumption (ii) of Theorem 1,

which lets us invoke chaining bounds for conditional second moments of $\sup_{\theta \in S} |Z_n^\circ(\theta)|$ for various S .

In Section 3 we prove a new consistency theorem (Theorem 3) and a new central limit theorem (Theorem 4, generalizing Wu's Theorem 5) for nonlinear LSEs. More precisely, our consistency theorem corresponds to an explicit bound for $\mathbb{P}\{\kappa_n(\hat{\theta}_n) \geq R\}$, but we state the result in a form that makes comparison with Theorem 1 easier. Our Theorem does not imply almost sure convergence, but our techniques could easily be adapted to that task. We regard the consistency as a preliminary to the next level of asymptotics and not as an end in itself. We describe the local asymptotic behavior with another approximation result, Theorem 4, which can easily be transformed into a central limit theorem under a variety of mild assumptions on the $\{u_i\}$ errors. For example, in Section 4 we apply the Theorem to the model (3), to sharpen the consistency result at $\theta_0 = (1, 1/2)$ into the approximation

$$(7) \quad \left(\ell_n^{1/2}(\hat{\lambda}_n - 1), \ell_n^{3/2}(1 - 2\hat{\mu}_n) \right) = \sum_{i \leq n} u_i \zeta'_{i,n} + o_p(1)$$

where $\ell_n := \log n$ and

$$\zeta_{i,n} = i^{-1/2} \ell_n^{-1/2} \begin{pmatrix} 2 & -6 \\ -6 & 24 \end{pmatrix} \begin{pmatrix} 2 \\ \ell_i/\ell_n \end{pmatrix}.$$

The sum on the right-hand side of (7) is of order $O_p(1)$ when $\sup_i \text{var}(u_i) < \infty$. If the $\{u_i\}$ are also identically distributed, the sum has a limiting multivariate normal distribution.

2. MAXIMAL INEQUALITIES

Assumption (ii) of Theorem 1 ensures that the increments $Z_n(\theta_1) - Z_n(\theta_2)$ are controlled by the ordinary Euclidean distance in Θ ; we allow for control by more general metrics. Wu invoked a maximal inequality for sums of random continuous processes, a result derived from a bound on the covering numbers for M_θ as a subset of \mathbb{R}^n under the usual Euclidean distance; we work with covering numbers for other metrics.

Definition 1. Let (T, d) be a metric space. The covering number $N(\delta, S, d)$ is defined as the size of the smallest δ -net for S , that is, the smallest N for which there are points t_1, \dots, t_N in T with $\min_i d(s, t_i) \leq \delta$ for every s in S .

Remark. The definition is the same for a pseudometric space, that is, a space where $d(\theta_1, \theta_2) = 0$ need not imply $\theta_1 = \theta_2$. In fact, all results in our paper that refer to metric spaces also apply to pseudometric spaces. The slight increase in generality is sometimes convenient when dealing with metrics defined by \mathcal{L}^p norms on functions.

Standard chaining arguments (see, for example, Pollard 1989), give maximal inequalities for processes with subgaussian increments controlled by a metric on the index set.

Theorem 2. Let $\{W_t : t \in T\}$ be a stochastic process, indexed by a metric space (T, d) , with subgaussian increments. Let T_δ be a δ -net for T . Suppose:

- (i) there is a constant K such that $\tau(W_s - W_t) \leq Kd(s, t)$ for all $s, t \in T$;
- (ii) $J_\delta := \int_0^\delta \rho(N(y, S, d)) dy < \infty$, where $\rho(N) := \sqrt{1 + \log N}$.

Then there is a universal constant c_1 such that

$$\frac{1}{c_1} \sqrt{\mathbb{P} \sup_t |W_t|^2} \leq K J_\delta + \rho(N(\delta, T, d)) \max_{s \in T_\delta} \tau(W_s).$$

Remark. We should perhaps work with outer expectations because, in general, there is no guarantee that a supremum of uncountably many random variables is measurable. For concrete examples, such as the one discussed in Section 4, measurability can usually be established by routine separability arguments. Accordingly, we will ignore the issue in this paper.

Under the assumption that $\text{var}(u_i) \leq \sigma^2$, the X process from (6) need not have subgaussian increments. However, it can be bounded in a stochastic sense by a symmetrized process $X^\circ(\theta) := \sum_{i \leq n} \epsilon_i u_i g_i(\theta)$, where the $2n$ random variables $\epsilon_1, \dots, \epsilon_n, u_1, \dots, u_n$ are mutually independent with $\mathbb{P}\{\epsilon_i = +1\} = 1/2 = \mathbb{P}\{\epsilon_i = -1\}$. In fact, for each subset S of the index set Θ ,

$$(8) \quad \mathbb{P} \sup_{\theta \in S} |X(\theta)|^2 \leq 4 \mathbb{P} \sup_{\theta \in S} |X^\circ(\theta)|^2.$$

For a proof see, for example, van der Vaart and Wellner (1996, Lemma 2.3.1). The process X° has conditionally subgaussian increments with

$$(9) \quad \tau_u \left(X_{\theta_1}^\circ - X_{\theta_2}^\circ \right)^2 \leq \sum_{i \leq n} u_i^2 \left(g_i(\theta_1) - g_i(\theta_2) \right)^2.$$

The subscript u indicates the conditioning on u .

Corollary 1. *Let S_δ be a δ -net for S and let X be as in (6). Suppose*

- (i) $\mathbb{P}u_i = 0$ and $\text{var}(u_i) \leq \sigma^2$ for $i = 1, \dots, n$
- (ii) *there is a metric d for which $J_\delta := \int_0^\delta \rho(N(y, S, d)) dy < \infty$*
- (iii) *there are constants L_1, \dots, L_n for which*

$$|g_i(\theta_1) - g_i(\theta_2)| \leq L_i d(\theta_1, \theta_2) \quad \text{for all } i \text{ and all } \theta_1, \theta_2 \in S$$

- (iv) *there are constants b_1, \dots, b_n for which $|g_i(\theta)| \leq b_i$ for all i and all θ in S .*

Then there is a universal constant c_2 such that

$$\mathbb{P} \sup_{\theta \in S} |X_\theta|^2 \leq c_2^2 \sigma^2 (L J_\delta + B \rho(N(\delta, S, d)))^2$$

where $L := \sqrt{\sum_i L_i^2}$ and $B := \sqrt{\sum_i b_i^2}$.

Proof. From (9),

$$\tau_u(X_{\theta_1}^\circ - X_{\theta_2}^\circ) \leq L_u d(\theta_1, \theta_2) \quad \text{where } L_u := \sqrt{\sum_{i \leq n} L_i^2 u_i^2}$$

and

$$\tau_u(X_\theta^\circ) \leq B_u := \sqrt{\sum_{i \leq n} b_i^2 u_i^2}$$

Apply Theorem 2 conditionally to the process X° to bound $\mathbb{P}_u \sup_{\theta \in S} |X_\theta^\circ|^2$. Then invoke inequality (8), using the fact that $\mathbb{P}L_u^2 \leq \sigma^2 L^2$ and $\mathbb{P}B_u^2 \leq \sigma^2 B^2$. \square

3. LIMIT THEOREMS

Inequality (5) and Corollary 1, with $g_i(\theta) = f_i(\theta) - f_i(\theta_0)$, give us some probabilistic control over $\widehat{\theta}_n$.

Theorem 3. *Let S be a subset of Θ equipped with a pseudometric d . Let $\{L_i : i = 1, \dots, n\}$, $\{b_i : i = 1, \dots, n\}$, and δ be positive constants such that*

- (i) $|f_i(\theta_1) - f_i(\theta_2)| \leq L_i d(\theta_1, \theta_2)$ for all $\theta_1, \theta_2 \in S$
- (ii) $|f_i(\theta) - f_i(\theta_0)| \leq b_i$ for all $\theta \in S$
- (iii) $J_\delta := \int_0^\delta \rho \left(N(y, S, d) \right) dy < \infty$

Then

$$\mathbb{P}\{\widehat{\theta}_n \in S\} \leq 4c_2^2 \sigma^2 \left(B \rho \left(N(\delta, S, d) \right) + L J_\delta \right)^2 / R^4,$$

where $R := \inf\{\kappa(\theta) : \theta \in S\}$, and $L^2 = \sum_i L_i^2$, and $B^2 := \sum_i b_i^2$.

The Theorem becomes more versatile in its application if we partition S into a countable union of subsets S_k , each equipped with its own pseudometric and Lipschitz constants. We then have $\mathbb{P}\{\widehat{\theta}_n \in \cup_k S_k\}$ smaller than a sum over k of bounds analogous to those in the Theorem. As shown in Section 4, this method works well for the Wu example if we take $S_k = \{\theta : R_k \leq \kappa_n(\theta) < R_{k+1}\}$, for an $\{R_k\}$ sequence increasing geometrically.

A similar appeal to Corollary 1, with the $g_i(\theta)$ as partial derivatives of $f_i(\theta)$ functions, gives us enough local control over Z_n to go beyond consistency. To accommodate the application in Section 4, we change notation slightly by working with a triangular array: for each n ,

$$y_{in} = f_{in}(\theta_0) + u_{in}, \quad \text{for } i = 1, 2, \dots, n,$$

where the $\{u_{in} : i = 1, \dots, n\}$ are unobserved independent random variables with mean zero and variance bounded by σ^2 .

Theorem 4. *Suppose $\widehat{\theta}_n \rightarrow \theta_0$ in probability, with θ_0 an interior point of Θ , a subset of \mathbb{R}^p . Suppose also:*

- (i) Each f_{in} is continuously differentiable in a neighborhood \mathcal{N} of θ_0 with derivatives $D_{in}(\theta) = \partial f_{in}(\theta) / \partial \theta$.
- (ii) $\gamma_n^2 := \sum_{i \leq n} |D_{in}(\theta_0)|^2 \rightarrow \infty$ as $n \rightarrow \infty$.
- (iii) There are constants $\{M_{in}\}$ with $\sum_{i \leq n} M_{in}^2 = O(\gamma_n^2)$ and a metric d on \mathcal{N} for which $|D_{in}(\theta_1) - D_{in}(\theta_2)| \leq M_{in}d(\theta_1, \theta_2)$ for $\theta_1, \theta_2 \in \mathcal{N}$.
- (iv) The smallest eigenvalue of the matrix $V_n = \gamma_n^{-2} \sum_{i \leq n} D_{in}(\theta_0) D_{in}(\theta_0)'$ is bounded away from zero for n large enough.
- (v) $\int_0^1 \rho(N(y, \mathcal{N}, d)) dy < \infty$
- (vi) $d(\theta, \theta_0) \rightarrow 0$ as $\theta \rightarrow \theta_0$.

Then $\kappa_n(\hat{\theta}_n) = O_p(1)$ and

$$\gamma_n(\hat{\theta}_n - \theta_0) = \sum_{i \leq n} \xi_{i,n} u_{in} + o_p(1) = O_p(1).$$

where $\xi_{i,n} = \gamma_n^{-1} V_n^{-1} D_{in}(\theta_0)$.

Proof. Let D be the $p \times n$ matrix with i th column $D_{in}(\theta_0)$, so that $\gamma_n^2 = \text{trace}(DD')$ and $V_n = \gamma_n^{-2} DD'$. The main idea of the proof is to replace $f(\theta)$ by $f(\theta_0) + D'(\theta - \theta_0)$, thereby approximating $\hat{\theta}_n$ by the least-squares solution

$$\bar{\theta}_n := \theta_0 + (DD')^{-1} Du = \underset{\theta \in \mathbb{R}^p}{\text{argmin}} |y - f(\theta_0) - D'(\theta - \theta_0)|.$$

To simplify notation, assume with no loss of generality, that $f(\theta_0) = 0$ and $\theta_0 = 0$. Also, drop extra n subscripts when the meaning is clear. The assertion of the Theorem is that $\hat{\theta}_n = \bar{\theta}_n + o_p(\gamma_n^{-1})$.

Without loss of generality, suppose the smallest eigenvalue of V_n is larger than a fixed constant $c_0^2 > 0$. Then

$$\gamma_n^2 = \text{trace}(DD') \geq \sup_{|t| \leq 1} |D't|^2 \geq \inf_{|t| \leq 1} |D't|^2 = c_0^2 \gamma_n^2,$$

from which it follows that

$$(10) \quad c_0 |t| \leq |D't| / \gamma_n \leq |t| \quad \text{for all } t \in \mathbb{R}^p.$$

Similarly, $\mathbb{P}|Du|^2 = \text{trace} \left(D\mathbb{P}(uu')D' \right) \leq \sigma^2\gamma_n^2$, implying that $|Du| = O_p(\gamma_n)$ and

$$\bar{\theta}_n = \gamma_n^{-2}V_n^{-1}Du = O_p(\gamma_n^{-1}).$$

In particular, $\mathbb{P}\{\bar{\theta}_n \in \mathcal{N}\} \rightarrow 1$, because θ_0 is an interior point of Θ . Note also that

$$\mathbb{P}|\sum_{i \leq n} \xi_i u_i|^2 \leq \sigma^2 \text{trace}(\sum_{i \leq n} \xi_i \xi_i') = \sigma^2 \text{trace}(V_n^{-1}) < \infty.$$

Consequently $\sum_{i \leq n} \xi_i u_i = O_p(1)$.

From the assumed consistency, we know that there is a sequence of balls $\mathcal{N}_n \subseteq \mathcal{N}$ that shrink to $\{0\}$ for which $\mathbb{P}\{\hat{\theta}_n \in \mathcal{N}_n\} \rightarrow 1$. From (vi) and (v), it follows that both $r_n := \sup\{d(\theta, 0) : \theta \in \mathcal{N}_n\}$ and $J_{r_n} = \int_0^{r_n} \rho \left(N(y, \mathcal{N}, d) \right) dy$ converge to zero as $n \rightarrow \infty$.

The $n \times 1$ remainder vector $R(\theta) := f(\theta) - D'\theta$ has i th component

$$(11) \quad R_i(\theta) = f_i(\theta) - D_i(0)'\theta = \theta' \int_0^1 D_i(t\theta) - D_i(0) dt.$$

Uniformly in the neighborhood \mathcal{N}_n we have

$$|R(\theta)| \leq |\theta| \left(\sum_{i \leq n} M_{in}^2 \right)^{1/2} r_n = o(|\theta|\gamma_n),$$

which, together with the upper bound from inequality (10), implies

$$(12) \quad |f(\theta)|^2 = |D'\theta|^2 + o(\gamma_n^2|\theta|^2) = O(\gamma_n^2|\theta|^2) \quad \text{as } |\theta| \rightarrow 0.$$

In the neighborhood \mathcal{N}_n , via (11) we also have,

$$|u'R(\theta)| \leq |\theta| \sup_{s \in \mathcal{N}_n} \left| \sum_i u_i \left(D_i(s) - D_i(0) \right) \right|.$$

From Corollary 1 with $g_i(\theta) = D_i(\theta) - D_i(0)$ deduce that

$$\mathbb{P} \sup_{s \in \mathcal{N}_n} \left| \sum_i u_i \left(D_i(s) - D_i(0) \right) \right|^2 \leq c_2^2 \sigma^2 J_{r_n}^2 \sum_i M_{in}^2 = o(\gamma_n^2),$$

which implies

$$(13) \quad |u'R(\theta)| = o_p(\gamma_n|\theta|) \quad \text{uniformly for } \theta \in \mathcal{N}_n.$$

Approximations (12) and (13) give us uniform approximations for the criterion functions in the shrinking neighborhoods \mathcal{N}_n :

$$\begin{aligned}
 G_n(\theta) &= |u - f(\theta)|^2 - |u|^2 \\
 &= -2u'f(\theta) + |f(\theta)|^2 \\
 (14) \quad &= -2u'D'\theta + |D'\theta|^2 + o_p(\gamma_n|\theta|) + o_p(\gamma_n^2|\theta|^2) \\
 &= |u - D'\bar{\theta}_n|^2 - |u|^2 + |D'(\theta - \bar{\theta}_n)|^2 + o_p(\gamma_n|\theta|) + o_p(\gamma_n^2|\theta|^2).
 \end{aligned}$$

The uniform smallness of the remainder terms lets us approximate G_n at random points that are known to lie in \mathcal{N}_n .

The rest of the argument is similar to that of Chernoff (1954). When $\hat{\theta}_n \in \mathcal{N}_n$ we have $G_n(\hat{\theta}_n) \leq G_n(0)$, implying

$$|D'(\hat{\theta}_n - \bar{\theta}_n)|^2 + o_p(\gamma_n|\hat{\theta}_n|) + o_p(\gamma_n^2|\hat{\theta}_n|^2) \leq |D'\bar{\theta}_n|^2.$$

Invoke (10) again, simplifying the last approximation to

$$c_0^2|\gamma_n\hat{\theta}_n - \gamma_n\bar{\theta}_n|^2 \leq O_p(1) + o_p(|\gamma_n\hat{\theta}_n| + |\gamma_n\bar{\theta}_n|^2).$$

It follows that $|\hat{\theta}_n| = O_p(\gamma_n^{-1})$ and, via (12),

$$\kappa_n(\hat{\theta}_n) = |f(\hat{\theta}_n)| = O_p(1).$$

We may also assume that \mathcal{N}_n shrinks slowly enough to ensure that $\mathbb{P}\{\bar{\theta}_n \in \mathcal{N}_n\} \rightarrow 1$. When both $\hat{\theta}_n$ and $\bar{\theta}_n$ lie in \mathcal{N}_n the inequality $G_n(\hat{\theta}_n) \leq G_n(\bar{\theta}_n)$ and approximation (14) give

$$|D'(\hat{\theta}_n - \bar{\theta}_n)|^2 + o_p(1) \leq o_p(1).$$

It follows that $\hat{\theta}_n = \bar{\theta}_n + o_p(\gamma_n^{-1})$. □

Remark. If the errors are iid and $\max |\xi_{i,n}| = o(1)$ then the distribution of $\sum_{i \leq n} \xi_{i,n} u_{in}$ is asymptotically $N(0, \sigma^2 V_n^{-1})$.

4. ANALYSIS OF MODEL (3): WU'S EXAMPLE

The results in this section illustrate the work of our limit theorems in a particular case where Wu's method fails. We prove both consistency and a central limit theorem for the model (3) with $\theta_0 = (\lambda_0, 1/2)$. In fact, without loss of generality, $\lambda_0 = 1$.

As before, let $\ell_n = \log n$. Remember $\theta = (\lambda, \mu)$ with $\lambda \in \mathbb{R}$ and $0 \leq \mu \leq C_\mu$ for a finite constant C_μ greater than $1/2$, which ensures that $\theta_0 = (1, 1/2)$ is an interior point of the parameter space. Taking $C_\mu = 1/2$ would complicate the central limit theorem only slightly. The behavior of $\hat{\theta}_n$ is determined by the behavior of the function

$$G_n(\gamma) := \sum_{i \leq n} i^{-1+\gamma} \quad \text{for } \gamma \leq 1,$$

or its standardized version

$$g_n(\beta) := G_n(\beta/\ell_n)/G_n(0) = \sum_{i \leq n} \left(i^{-1}/G_n(0) \right) \exp \left(\beta \ell_i/\ell_n \right),$$

which is the moment generating function of the probability distribution that puts mass $i^{-1}/G_n(0)$ at ℓ_i/ℓ_n , for $i = 1, \dots, n$. For large n , the function g_n is well approximated by the increasing, nonnegative function

$$g(\beta) = \begin{cases} (e^\beta - 1)/\beta & \text{for } \beta \neq 0 \\ 1 & \text{for } \beta = 0 \end{cases},$$

the moment generating function of the uniform distribution on $(0, 1)$. More precisely, comparison of the sum with the integral $\int_1^n x^{-1+\gamma} dx$ gives

$$(15) \quad G_n(\gamma) = \ell_n g(\gamma \ell_n) + r_n(\gamma) \quad \text{with } 0 \leq r_n(\gamma) \leq 1 \text{ for } \gamma \leq 1.$$

The distributions corresponding to both g_n and g are concentrated on $[0, 1]$. Both functions have the properties described in the following lemma.

Lemma 1. *Suppose $h(\gamma) = P \exp(\gamma x)$, the moment generating function of a probability distribution concentrated on $[0, 1]$. Then*

- (i) $\log h$ is convex

(ii) $h(\gamma)^2/h(2\gamma)$ is unimodal: increasing for $\gamma < 0$, decreasing for $\gamma > 0$, achieving its maximum value of 1 at $\gamma = 0$

(iii) $h'(\gamma) \leq h(\gamma)$

Proof. Assertion (i) is just the well known fact that the logarithm of a moment generating function is convex. Thus h'/h , the derivative of $\log h$, is an increasing function, which implies (ii) because

$$\frac{d}{d\gamma} \log \left(\frac{h(\gamma)^2}{h(2\gamma)} \right) = 2 \frac{h'(\gamma)}{h(\gamma)} - 2 \frac{h'(2\gamma)}{h(2\gamma)}.$$

Property (iii) comes from the representation $h'(\gamma) = P \left(x e^{\gamma x} \right)$. \square

Remark. Direct calculation shows that $g(\gamma)^2/g(2\gamma)$ is a symmetric function.

Reparametrize by putting $\beta = (1 - 2\mu)\ell_n$, with $(1 - 2C_\mu)\ell_n \leq \beta \leq \ell_n$, and $\alpha = \lambda \sqrt{G_n(\beta/\ell_n)}$. Notice that $|f(\theta)| = |\alpha|$ and that θ_0 corresponds to $\alpha_0 = \sqrt{G_n(0)} \approx \sqrt{\ell_n}$ and $\beta_0 = 0$. Also

$$f_i(\theta) = \alpha \nu_i(\beta/\ell_n) \quad \text{where} \quad \nu_i(\gamma) := i^{-1/2} \exp(\gamma \ell_i/2) / \sqrt{G_n(\gamma)},$$

and

$$(16) \quad \kappa_n(\theta)^2 = G_n(0) \left(\lambda^2 g_n(\beta) - 2\lambda g_n(\beta/2) + 1 \right).$$

We define $\nu_i := \sup_{\gamma \leq 1} \nu_i(\gamma)$.

Lemma 2. For all (α, β) corresponding to $\theta = (\lambda, \mu) \in \mathbb{R} \times [0, C_\mu]$:

- (i) $\kappa_n(\theta) - \sqrt{G_n(0)} \leq |\alpha| \leq \kappa_n(\theta) + \sqrt{G_n(0)}$
- (ii) $\sum_{i \leq n} \nu_i^2 = O \left(\log \log n \right)$
- (iii) $|d\nu_i(\beta/\ell_n)/d\beta| \leq \frac{1}{2} \nu_i(\beta/\ell_n)$
- (iv) $|f_i(\alpha_1, \beta_1) - f_i(\alpha_2, \beta_2)| \leq \left(|\alpha_1 - \alpha_2| + \frac{1}{2} |\alpha_2| |\beta_1 - \beta_2| \right) \nu_i$
- (v) $|f_i(\theta) - f_i(\theta_0)| \leq i^{-1/2} + |\alpha| \nu_i$

Proof. Inequalities (i) and (v) follow from the triangle inequality.

For inequality (ii), first note that $\nu_1^2 \leq 1$. For $i \geq 2$, separate out contributions from three ranges:

$$\nu_i^2 = \max \left(\sup_{1 \geq \gamma \geq 1/\ell_n} \nu_i(\gamma)^2, \sup_{|\gamma| < 1/\ell_n} \nu_i(\gamma)^2, \sup_{\gamma \leq -1/\ell_n} \nu_i(\gamma)^2 \right).$$

For $\gamma \geq 1/\ell_n$, invoke (15) to get a tractable upper bound:

$$\nu_i(\gamma)^2 \leq i^{-1} \frac{\exp(\gamma \ell_i)}{\ell_n g(\gamma \ell_n)} \leq i^{-1} \gamma \frac{\exp(\gamma \ell_i)}{\exp(\gamma \ell_n) - 1} \leq i^{-1} \frac{\exp(\log \gamma + \gamma \log(i/n))}{1 - e^{-1}}.$$

The last expression achieves its maximum over $[1/\ell_n, 1]$ at

$$\gamma_0 := \begin{cases} 1/\log(n/i) & \text{if } 1 \leq i \leq n/e \\ 1 & \text{if } n/e \leq i \leq n \end{cases},$$

which gives

(17)

$$\sup_{1 \geq \gamma \geq 1/\ell_n} \nu_i(\gamma)^2 \leq \frac{(e-1)^{-1}}{n} H\left(\frac{i \wedge (n/e)}{n}\right) \quad \text{where } H(x) := 1/(x \log(1/x)).$$

Similarly, if $-1 < \gamma \ell_n < 1$,

$$\nu_i(\gamma)^2 \leq \frac{\exp(\gamma \ell_i)}{i \ell_n g(\gamma \ell_n)} \leq \frac{\exp(\ell_i/\ell_n)}{i \ell_n g(-1)} \leq \frac{e/g(-1)}{i \ell_n}.$$

The last term is smaller than a constant multiple of the bound from (17). Finally, if $-\gamma = \delta \geq 1/\ell_n$ and $i \geq 2$ then

$$\nu_i(\gamma)^2 \leq i^{-1} \delta \frac{\exp(-\delta \ell_i)}{1 - \exp(-\delta \ell_n)} \leq i^{-1} \frac{\exp(\log \delta - \delta \ell_i)}{1 - e^{-1}} \leq \frac{e^{-1}/(1 - e^{-1})}{i \ell_i}.$$

In summary, for some universal constant C ,

$$\nu_i^2 \leq C \max \left(n^{-1} H\left(\frac{i \wedge (n/e)}{n}\right), \frac{1}{i \log i} \right) \quad \text{if } 2 \leq i \leq n.$$

Bounding sums by integrals we thus have

$$C^{-1} \sum_{i=2}^n \nu_i^2 \leq \int_{1/n}^{1/e} H(x) dx + H(1/e)/n + \int_2^n (x \log x)^{-1} dx = O(\log \log n).$$

For (iii) note that

$$2\frac{d}{d\beta}\nu_i(\beta/\ell_n) = 2\frac{d}{d\beta}\exp\left(\frac{1}{2}\beta\ell_i/\ell_n\right)\left(G_n(0)g_n(\beta)\right)^{-1/2} = \left(\frac{\ell_i}{\ell_n} - \frac{g'_n(\beta)}{g_n(\beta)}\right)\nu_i(\beta),$$

which is bounded in absolute value by $\nu_i(\beta)$ because $0 \leq g'_n(\beta) \leq g_n(\beta)$.

For (iv)

$$\begin{aligned} |f_i(\alpha_1, \beta_1) - f_i(\alpha_2, \beta_2)| &\leq |(\alpha_1 - \alpha_2)\nu_i(\beta_1/\ell_n)| + |\alpha_2||\nu_i(\beta_1/\ell_n) - \nu_i(\beta_2/\ell_n)| \\ &\leq |(\alpha_1 - \alpha_2)|\nu_i + |\alpha_2||(\beta_1 - \beta_2)|\frac{1}{2}\nu_i, \end{aligned}$$

the bound for the second term coming from the mean-value theorem and (iii). \square

Lemma 3. For $\epsilon > 0$, let $\mathcal{N}_\epsilon = \{\theta : \max(|\lambda - 1|, |\beta|) \geq \epsilon\}$. If ϵ is small enough, there exists a constant $C_\epsilon > 0$ such that $\inf\{\kappa_n(\theta) : \theta \notin \mathcal{N}_\epsilon\} \geq C_\epsilon\sqrt{\ell_n}$ when n is large enough.

Proof. Suppose $|\beta| \geq \epsilon$. Remember that $G_n(0) \geq \ell_n$. Minimize over λ the lower bound (16) for $\kappa_n(\theta)^2$ by choosing $\lambda = g_n(\beta/2)/g_n(\beta)$, then invoke Lemma 1(ii).

$$\frac{\kappa_n(\theta)^2}{\ell_n} \geq 1 - \frac{g_n(\beta/2)^2}{g_n(\beta)} \geq 1 - \max\left(\frac{g_n(\epsilon/2)^2}{g_n(\epsilon)}, \frac{g_n(-\epsilon/2)^2}{g_n(-\epsilon)}\right) \rightarrow 1 - \frac{g(\epsilon/2)^2}{g(\epsilon)} > 0.$$

If $|\beta| \leq \epsilon$ and ϵ is small enough to make $(1 - \epsilon)e^{\epsilon/2} < 1 < (1 + \epsilon)e^{-\epsilon/2}$, use

$$\kappa_n(\theta)^2 = \sum_{i \leq n} i^{-1} \left(\lambda \exp(\beta\ell_i/2\ell_n) - 1 \right)^2.$$

If $\lambda \geq 1 + \epsilon$ bound each summand from below by $i^{-1}((1 + \epsilon)e^{-\epsilon/2} - 1)^2$. If $\lambda \leq 1 - \epsilon$ bound each summand from below by $i^{-1}(1 - (1 - \epsilon)e^{\epsilon/2})^2$. \square

4.1. Consistency. On the annulus $S_R := \{R \leq \kappa_n(\theta) < 2R\}$ we have

$$\begin{aligned} |a| &\leq K_R := 2R + \sqrt{G_n(0)} \\ |f_i(\theta_1) - f_i(\theta_2)| &\leq K_R \nu_i d_R(\theta_1, \theta_2) \\ \text{where } d_R(\theta_1, \theta_2) &:= |\alpha_1 - \alpha_2|/K_R + \frac{1}{2}|\beta_1 - \beta_2| \\ |f_i(\theta) - f_i(\theta_0)| &\leq b_i := i^{-1/2} + K_R \nu_i. \end{aligned}$$

Note that

$$\sum_{i \leq n} \left(i^{-1/2} + K_R \nu_i \right)^2 = O(\ell_n + K_R^2 \log \ell_n) = O(K_R^2 \mathcal{L}_n) \quad \text{where } \mathcal{L}_n := \log \log n.$$

The rectangle $\{|\alpha| \leq K_R, |\beta| \leq c\ell_n\}$ can be partitioned into $O(y^{-1}\ell_n/y)$ subrectangles of d_R -diameter at most y . Thus $N(y, S_R, d_R) \leq C_0\ell_n/y^2$ for a constant C_0 that depends only on C_μ , which gives

$$\int_0^1 \rho\left(N(y, S_R, d_R)\right) dy = O\left(\sqrt{\mathcal{L}_n}\right).$$

Apply Theorem 3 with $\delta = 1$ to conclude that

$$\mathbb{P}\{\widehat{\theta}_n \in S_R\} \leq C_1 K_R^2 \mathcal{L}_n^2 / R^4 \leq C_2 (R^2 + \ell_n) \mathcal{L}_n^2 / R^4.$$

Put $R = C_3 2^k (\ell_n \mathcal{L}_n^2)^{1/4}$ then sum over k to deduce that

$$\mathbb{P}\{\kappa_n(\widehat{\theta}_n) \geq C_3 (\ell_n \mathcal{L}_n^2)^{1/4}\} \leq \epsilon \quad \text{eventually}$$

if the constant C_3 is large enough. That is $\kappa_n(\widehat{\theta}_n) = O_p\left((\ell_n \mathcal{L}_n^2)^{1/4}\right)$ and, via Lemma 3,

$$|\widehat{\lambda}_n - 1| = o_p(1) \quad \text{and} \quad 2\ell_n |\widehat{\mu}_n - \mu_0| = |\widehat{\beta}| = o_p(1).$$

4.2. Central limit theorem. This time work with the (λ, β) reparametrization, with

$$\begin{aligned} f_i(\lambda, \beta) &= \lambda i^{-1/2+\beta/2\ell_n} \\ D_i(\lambda, \beta)' &= \left(\frac{\partial f_i(\lambda, \beta)}{\partial \lambda}, \frac{\partial f_i(\lambda, \beta)}{\partial \beta} \right) = \left(1/\lambda, \ell_i/2\ell_n \right) f_i(\lambda, \beta) \end{aligned}$$

and $\theta_0 = (\lambda_0, \beta_0) = (1, 0)$. Take d as the usual two-dimensional Euclidean distance in the (λ, β) space. For simplicity of notation, we omit some n subscripts, even though the relationship between θ and (λ, β) changes with n .

We have just shown that the LSE $(\widehat{\lambda}_n, \widehat{\beta}_n)$ is consistent.

Comparison of sums with analogous integrals gives the approximations

$$(18) \quad \sum_{i \leq n} i^{-1} \ell_i^{p-1} = \ell_n^p / p + r_p \quad \text{with } |r_p| \leq 1 \text{ for } p = 0, 1, 2, \dots$$

In consequence,

$$\gamma_n^2 = \sum_{i \leq n} |D_i(\lambda_0, \beta_0)|^2 = \sum_{i \leq n} i^{-1} (1 + \ell_i^2 / 4\ell_n^2) = \frac{13}{12} \ell_n + O(1)$$

and

$$V_n = \gamma_n^{-2} \sum_{i \leq n} i^{-1} \begin{pmatrix} 1 & \ell_i/2\ell_n \\ \ell_i/2\ell_n & \ell_i^2/4\ell_n^2 \end{pmatrix} = V + O(1/\ell_n) \quad \text{where } V = \frac{1}{13} \begin{pmatrix} 12 & 3 \\ 3 & 1 \end{pmatrix}.$$

The smaller eigenvalue of V_n converges to the smaller eigenvalue of the positive definite matrix V , which is strictly positive.

Within the neighborhood $\mathcal{N}_\epsilon := \{\max(|\lambda - 1|, |\beta|) \leq \epsilon\}$, for a fixed $\epsilon \leq 1/2$, both $|f_i(\lambda, \beta)|$ and $|D_i(\lambda, \beta)|$ are bounded by a multiple of $i^{-1/2}$. Thus

$$|D_i(\theta_1) - D_i(\theta_2)| \leq \left| \lambda_1^{-1} - \lambda_2^{-1} \right| |f_i(\theta_1)| + 3|f_i(\theta_1) - f_i(\theta_2)| \leq C_\epsilon i^{-1/2} d(\theta_1, \theta_2).$$

That is, we may take M_i as a multiple of $i^{-1/2}$, which gives $\sum_{i \leq n} M_i^2 = O(\ell_n)$.

All the conditions of Theorem 4 are satisfied. We have

$$\sqrt{\ell_n}(\widehat{\lambda_n} - 1, \widehat{\beta_n}) = \frac{12}{13} \sum_{i \leq n} u_i i^{-1/2} \ell_n^{-1/2} (1, \ell_i/2\ell_n) V^{-1} + o_p(1).$$

REFERENCES

- Chernoff, H. (1954). On the distribution of the likelihood ratio. *Annals of Mathematical Statistics* 25, 573–578.
- Jennrich, R. I. (1969). Asymptotic properties of non-linear least squares estimators. *Annals of Mathematical Statistics* 40, 633–643.
- Lai, T. L. (1994). Asymptotic properties of nonlinear least-squares estimates in stochastic regression models. *Annals of Statistics* 22, 1917–1930.
- Pollard, D. (1989). Asymptotics via empirical processes (with discussion). *Statistical Science* 4, 341–366.
- Skouras, K. (2000). Strong consistency in nonlinear stochastic regression models. *Annals of Statistics* 28, 871–879.
- Van de Geer, S. (1990). Estimating a regression function. *Annals of Statistics* 18, 907–924.

- Van de Geer, S. and M. Wegkamp (1996). Consistency for the least squares estimator in nonparametric regression. *Annals of Statistics* 24, 2513–2523.
- van der Vaart, A. W. and J. A. Wellner (1996). *Weak Convergence and Empirical Process: With Applications to Statistics*. Springer-Verlag.
- Wu, C.-F. (1981). Asymptotic theory of nonlinear least squares estimation. *Annals of Statistics* 9, 501–513.

STATISTICS DEPARTMENT, YALE UNIVERSITY, BOX 208290 YALE STATION, NEW HAVEN, CT 06520-8290.

E-mail address: david.pollard@yale.edu; peter.radchenko@yale.edu

URL: <http://www.stat.yale.edu/~pollard/>; <http://pantheon.yale.edu/~pvr4>