

LOCAL M-ESTIMATION WITH DISCONTINUOUS CRITERION FOR DEPENDENT AND LIMITED OBSERVATIONS

BY MYUNG HWAN SEO¹ AND TAISUKE OTSU²

Seoul National University and London School of Economics

We examine the asymptotic properties of local M-estimators under three sets of high-level conditions. These conditions are sufficiently general to cover the minimum volume predictive region, the conditional maximum score estimator for a panel data discrete choice model and many other widely used estimators in statistics and econometrics. Specifically, they allow for discontinuous criterion functions of weakly dependent observations which may be localized by kernel smoothing and contain nuisance parameters with growing dimension. Furthermore, the localization can occur around parameter values rather than around a fixed point and the observations may take limited values which lead to set estimators. Our theory produces three different nonparametric cube root rates for local M-estimators and enables valid inference building on novel maximal inequalities for weakly dependent observations. The standard cube root asymptotics is included as a special case. The results are illustrated by various examples such as the Hough transform estimator with diminishing bandwidth, the maximum score-type set estimator and many others.

1. Introduction. There is a class of estimation problems in statistics where a point (or set-valued) estimator is obtained by maximizing a discontinuous and possibly localized criterion function. As a prototype, consider estimation of a simplified version of the minimum volume predictive region for y at $x = c$ [Polonik and Yao (2000)]. Let $\mathbb{I}\{\cdot\}$ be the indicator function, $K(\cdot)$ be a kernel function, and h_n be a bandwidth. At a significance level α , the estimator $[\hat{\theta} \pm \hat{v}]$ is obtained by the M-estimation

$$(1) \quad \max_{\theta \in \Theta} \sum_{t=1}^n \mathbb{I}\{|y_t - \theta| \leq \hat{v}\} K\left(\frac{x_t - c}{h_n}\right),$$

where Θ is some parameter space and

$$\hat{v} = \inf \left\{ v \in \mathbb{R} : \max_{\theta \in \Theta} \frac{\sum_{t=1}^n \mathbb{I}\{|y_t - \theta| \leq v\} K\left(\frac{x_t - c}{h_n}\right)}{\sum_{t=1}^n K\left(\frac{x_t - c}{h_n}\right)} \geq \alpha \right\}.$$

Received March 2016; revised October 2016.

¹Supported by Promising-Pioneering Researcher Program through Seoul National University.

²Supported by ERC Consolidator Grant SNP 615882.

MSC2010 subject classifications. Primary 62F12; secondary 60F17, 60G15, 62G20.

Key words and phrases. Cube root asymptotics, maximal inequality, mixing process, partial identification, parameter-dependent localization.

This problem exhibits several distinguishing features such as discontinuity of the criterion function, localization by kernel smoothing and serial dependence in time series data which have prevented a full-blown asymptotic analysis of the M-estimator $\hat{\theta}$. Only consistency is reported in the literature.

This type of M-estimation has numerous applications. Since Chernoff's (1964) study on estimation of the mode, many papers have raised such estimation problems, for example, the shorth [Andrews et al. (1972)], least median of squares [Rousseeuw (1984)], nonparametric monotone density estimation [Prakasa Rao (1969)] and maximum score estimation [Manski (1975)]. These classical examples are studied in a seminal work by Kim and Pollard (1990), which explained elegantly how this type of estimation problem induces so-called cube root asymptotics in a unified framework by means of empirical process theory. See also van der Vaart and Wellner (1996) and Kosorok (2008) for a general theory of M-estimation via empirical processes. However, these works do not cover the estimation problem in (1) due to their focus on cross-sectional data among other things. It should be emphasized that this is not a pathological example. We provide various relevant examples in Section 3 and the supplementary material [Seo and Otsu (2018), Section B] including the well-known Honoré and Kyriazidou's (2000) estimator for a dynamic panel discrete choice model and a localized maximum score estimator for a new binary choice model with random coefficients.

This paper covers a broader class of M-estimators than the above examples suggest. The baseline scenario (called *local M-estimation*) is generalized in two directions. First, we accommodate not only variables taking limited values (e.g., interval-valued data) which typically lead to estimation of a set rather than a point, but also nuisance parameters with growing dimension. Set estimation problems due to limited observations are also known as partial identification problems in econometrics [e.g., Manski and Tamer (2002)]. It is also novel to accommodate high-dimensional nuisance parameters in M-estimation with discontinuous criterion functions. Second, we allow for localization to be dependent on parameter values instead of prespecified values. For instance, the criterion function may take the form of $\sum_{t=1}^n \mathbb{I}\{|y_t - \theta| \leq h_n\}$ with $h_n \rightarrow 0$. Relevant examples include mode estimation [Chernoff (1964) and Lee (1989)] and the Hough transform estimator in image analysis [Goldenshluger and Zeevi (2004)]. Henceforth, we call this case *parameter-dependent local M-estimation*. Parameter-dependence brings some new features into our asymptotic analysis but in a different way from a classical example of parameter-dependency on the support such as the maximum likelihood estimator for Uniform $[0, \theta]$.

The main contribution of this paper is to develop a general asymptotic theory for such M-estimation problems. Our theoretical results cover all the examples above and can be used to establish limit laws for point estimators and convergence rates for set estimators. To this end, we develop suitable maximal inequalities which enable us to obtain nonparametric cube root rates of $(nh_n)^{1/3}$,

$\{nh_n/\log(nh_n)\}^{1/3}$ and $(nh_n^2)^{1/3}$ for the cases of local M-estimation, limited observations and parameter-dependent localization, respectively. These inequalities are extended to establish stochastic asymptotic equicontinuity of normalized processes of the criterion functions so that an argmax theorem delivers limit laws of the M-estimators. It is worth noting that all the conditions are characterized through moment conditions and can be easily verified as illustrated in the examples. Thus, our results can be applied without prior knowledge of empirical process theory. It is often not trivial to verify entropy conditions such as uniform manageability in [Kim and Pollard \(1990\)](#). Particularly for dependent data, the covering or bracketing numbers often need to be calculated using a norm that hinges on the mixing coefficients and distribution of the data [e.g., the $L_{2,\beta}$ -norm in [Doukhan, Massart and Rio \(1995\)](#)].

Another contribution is that we allow for weakly dependent data which are associated with absolutely regular processes with exponentially decaying mixing coefficients. In some applications, the cube root asymptotic theory has been extended to time series data, for example, [Anevski and Hössjer \(2006\)](#) for monotone density estimation, [Zinde-Walsh \(2002\)](#) for least median of squares, [de Jong and Woutersen \(2011\)](#) for maximum score and [Koo and Seo \(2015\)](#) for break estimation under misspecification. However, it is not clear whether they are able to handle the general class of estimation problems in this paper.

The paper is organized as follows. Section 2 develops an asymptotic theory for local M-estimation and Section 3 provides several examples. In Section 4, we generalize the asymptotic theory to the cases of limited observations (Section 4.1) and parameter-dependent localization (Section 4.2). Section 5 concludes. All proofs, details for illustrations and additional examples are contained in the supplementary material [[Seo and Otsu \(2018\)](#)].

2. Local M-estimation. This section studies the M-estimator $\hat{\theta}$ that maximizes

$$\mathbb{P}_n f_{n,\theta} = \frac{1}{n} \sum_{t=1}^n f_{n,\theta}(z_t),$$

where $\{f_{n,\theta}\}$ is a sequence of criterion functions indexed by the parameters $\theta \in \Theta \subseteq \mathbb{R}^d$ and $\{z_t\}$ is a strictly stationary sequence of random variables with marginal P . We introduce a set of conditions for $f_{n,\theta}$ that induces a possibly localized counterpart of [Kim and Pollard's \(1990\)](#) cube root asymptotics. Their cube root asymptotics can be viewed as a special case of ours, where $f_{n,\theta}$ does not vary with n . Let $Pf = \int f dP$ for a function f , $|\cdot|$ be the Euclidean norm, and $\|\cdot\|_2$ be the $L_2(P)$ -norm of a random variable. The class of criterion functions of interest is characterized as follows.

ASSUMPTION M. For a sequence $\{h_n\}$ of positive numbers with $nh_n \rightarrow \infty$, $f_{n,\theta}$ satisfies the following conditions:

(i) $h_n f_{n,\theta}$ is uniformly bounded, $\lim_{n \rightarrow \infty} P f_{n,\theta}$ is uniquely maximized at θ_0 , $P f_{n,\theta}$ is twice continuously differentiable at θ_0 for all n large enough, and satisfies

$$(2) \quad P(f_{n,\theta} - f_{n,\theta_0}) = \frac{1}{2}(\theta - \theta_0)' V(\theta - \theta_0) + o(|\theta - \theta_0|^2) + o((nh_n)^{-2/3}),$$

for a negative definite matrix V .

(ii) There exist positive constants C and C' such that

$$|\theta_1 - \theta_2| \leq Ch_n^{1/2} \|f_{n,\theta_1} - f_{n,\theta_2}\|_2,$$

for all n large enough and $\theta_1, \theta_2 \in \{\Theta : |\theta - \theta_0| \leq C'\}$.

(iii) There exists a positive constant C'' such that

$$(3) \quad P \sup_{\theta \in \Theta: |\theta - \theta'| < \varepsilon} h_n |f_{n,\theta} - f_{n,\theta'}|^2 \leq C'' \varepsilon,$$

for all n large enough, $\varepsilon > 0$ small enough and θ' in a neighborhood of θ_0 .

$\{h_n\}$ is usually a sequence of bandwidths for localization. Although we are primarily interested in the case of $h_n \rightarrow 0$, we do not exclude the case of $h_n = 1$ which corresponds to the conventional cube root asymptotics in [Kim and Pollard \(1990\)](#). Also, we note that our conditions appear somewhat different from [Kim and Pollard \(1990\)](#). In fact, our conditions consist of directly verifiable moment conditions without resorting to the notion of empirical process theory such as uniform manageability.

Assumption **M(i)** requires boundedness, point identification and quadratic approximation of $P f_{n,\theta}$. Boundedness of $h_n f_{n,\theta}$ is a major requirement but is satisfied for all examples in this paper and for [Kim and Pollard \(1990\)](#). In Section 4, we relax the assumption of point identification. When the criterion function involves kernel smoothing for localization, it typically takes the form of a product of a bounded function and smoothing term $\frac{1}{h_n} K(\frac{x-c}{h_n})$ [see (1) and examples in Section 3].

Despite discontinuity of $f_{n,\theta}$, its population counterpart $P f_{n,\theta}$ is smooth and approximated by a quadratic function as in (2). This distinguishes our estimation problem from that of a change-point in a regression model, which also involves a discontinuous criterion function but the change-point estimator is super-consistent [e.g., [Chan \(1993\)](#)] unless the estimating equation is misspecified as in the split point estimator for decision trees [[Bühlmann and Yu \(2002\)](#) and [Banerjee and McKeague \(2007\)](#)].

Assumption **M(ii)** is used to relate the $L_2(P)$ -norm for the criterion functions to the Euclidean norm for the parameters. This condition is implicit in [Kim and Pollard \[\(1990\), Condition \(v\)\]](#) under independent observations and is often verified in the course of checking the expansion in (2).

Assumption **M(iii)**, an envelope condition for the class $\mathcal{F}_n = \{f_{n,\theta} - f_{n,\theta'} : |\theta - \theta'| \leq \varepsilon\}$, plays a key role for cube root asymptotics. It should be noted that for

the familiar square root asymptotics, the upper bound in (3) is of order ε^2 instead of ε . It is often the case that verifying the envelope condition for arbitrary θ' in a neighborhood of θ_0 is no more demanding than that for θ_0 .

In particular, Assumption M(iii) is used to guarantee an integrability condition on the metric entropy with bracketing for \mathcal{F}_n in the $L_{2,\beta}$ -norm so that the maximal inequality in Doukhan, Massart and Rio [(1995), Theorem 3] can be applied to establish Lemma M below. On the other hand, Kim and Pollard (1990) used the concept of uniform manageability [Pollard (1989)] to control the size of \mathcal{F}_n which is defined by the ϵ -capacity, or metric entropy, by covering numbers. Generally the bracketing and covering numbers approaches are not directly comparable [see, e.g., Section 2.5 of van der Vaart and Wellner (1996)]. It would be interesting to explore how the symmetrization argument combined with the suitable manageability concept can be applied in our setup.

We now study the asymptotic properties of the M-estimator which is precisely defined as a random variable $\hat{\theta}$ satisfying

$$(4) \quad \mathbb{P}_n f_{n,\hat{\theta}} \geq \sup_{\theta \in \Theta} \mathbb{P}_n f_{n,\theta} - o_p((nh_n)^{-2/3}).$$

The first step is to establish weak consistency $\hat{\theta} \xrightarrow{P} \theta_0$, which is rather standard and usually shown by establishing the uniform convergence of $\mathbb{P}_n f_{n,\theta}$. In this section, we simply assume the consistency of $\hat{\theta}$. See the supplementary material [Seo and Otsu (2018)] for some illustrations to show consistency.

The next step is to derive the convergence rate of $\hat{\theta}$. A key ingredient for this is to obtain the modulus of continuity of the empirical process $\{\mathbb{G}_n h_n^{1/2}(f_{n,\theta} - f_{n,\theta_0})\}$ by some maximum inequality, where $\mathbb{G}_n f = \sqrt{n}(\mathbb{P}_n f - Pf)$ for a function f . If $f_{n,\theta}$ does not vary with n and $\{z_t\}$ is independent, several maximal inequalities are available in the literature [e.g., page 199 of Kim and Pollard (1990)]. If $f_{n,\theta}$ varies with n and $\{z_t\}$ is dependent, to the best of our knowledge, there is no maximal inequality which can be applied to the class of functions satisfying Assumption M. Our first task is to establish such an inequality.

To proceed, we now characterize the dependence structure of the data. Among several notions of dependence, this paper focuses on an absolutely regular process. See Doukhan, Massart and Rio (1995) for a discussion on empirical process theory of absolutely regular processes. Let $\mathcal{F}_{-\infty}^0$ and \mathcal{F}_m^∞ be σ -fields of $\{\dots, z_{t-1}, z_0\}$ and $\{z_m, z_{m+1}, \dots\}$, respectively. Define the β -mixing coefficient as $\beta_m = \frac{1}{2} \sup \sum_{(i,j) \in I \times J} |P\{A_i \cap B_j\} - P\{A_i\}P\{B_j\}|$, where the supremum is taken over all finite partitions $\{A_i\}_{i \in I}$ and $\{B_j\}_{j \in J}$, respectively, $\mathcal{F}_{-\infty}^0$ and \mathcal{F}_m^∞ measurable. Throughout the paper, we maintain the following assumption on $\{z_t\}$.

ASSUMPTION D. $\{z_t\}$ is a strictly stationary and absolutely regular process with β -mixing coefficients $\{\beta_m\}$ such that $\beta_m = O(\rho^m)$ for some $0 < \rho < 1$.

This assumption obviously covers the case of independent observations. It also says the mixing coefficient β_m should decay at an exponential rate.³ For example, various Markov, GARCH and stochastic volatility models satisfy this assumption [Carrasco and Chen (2002)]. See Section 2.2.2 below for further discussions.

Under this assumption, we obtain the following maximal inequality.

LEMMA M. *Under Assumptions M and D, there exist positive constants C and C' such that*

$$P \sup_{\theta \in \Theta: |\theta - \theta_0| < \delta} |\mathbb{G}_n h_n^{1/2}(f_{n,\theta} - f_{n,\theta_0})| \leq C \delta^{1/2},$$

for all n large enough and $\delta \in [(nh_n)^{-1/2}, C']$.

This lemma provides a preliminary lemma to derive the convergence rate.

LEMMA 1. *Under Assumptions M and D, for each $\varepsilon > 0$, there exist random variables $\{R_n\}$ of order $O_p(1)$ and a positive constant C such that*

$$|\mathbb{P}_n(f_{n,\theta} - f_{n,\theta_0}) - P(f_{n,\theta} - f_{n,\theta_0})| \leq \varepsilon |\theta - \theta_0|^2 + (nh_n)^{-2/3} R_n^2,$$

for all $\theta \in \{\Theta : (nh_n)^{-1/3} \leq |\theta - \theta_0| \leq C\}$.

We now derive the convergence rate of $\hat{\theta}$. Suppose $|\hat{\theta} - \theta_0| \geq (nh_n)^{-1/3}$. Then by (4), Lemma 1 and Assumption M(i), we can take a positive constant c such that

$$\begin{aligned} o_p((nh_n)^{-2/3}) &\leq \mathbb{P}_n(f_{n,\hat{\theta}} - f_{n,\theta_0}) \\ &\leq P(f_{n,\hat{\theta}} - f_{n,\theta_0}) + \varepsilon |\hat{\theta} - \theta_0|^2 + (nh_n)^{-2/3} R_n^2 \\ &\leq (-c + \varepsilon) |\hat{\theta} - \theta_0|^2 + o(|\hat{\theta} - \theta_0|^2) + O_p((nh_n)^{-2/3}), \end{aligned}$$

for each $\varepsilon > 0$. Taking ε small enough to satisfy $\varepsilon < c$ yields the convergence rate $\hat{\theta} - \theta_0 = O_p((nh_n)^{-1/3})$.

Given this, the final step is to establish the limiting distribution of $\hat{\theta}$. To this end, we apply a continuous mapping theorem of an argmax element [e.g., Theorem 2.7 of Kim and Pollard (1990)] and it is enough to show weak convergence of the normalized empirical process

$$Z_n(s) = n^{1/6} h_n^{2/3} \mathbb{G}_n(f_{n,\theta_0+s(nh_n)^{-1/3}} - f_{n,\theta_0}),$$

for $|s| \leq K$ with any $K > 0$. Weak convergence of Z_n may be characterized by its finite dimensional convergence and stochastic asymptotic equicontinuity (or tightness). If $f_{n,\theta}$ does not vary with n and $\{z_t\}$ is independent as in Kim and Pollard

³Polynomial decays of β_m are often associated with strong dependence and long memory type behaviors in sample statistics. See Chen, Hansen and Carrasco (2010) and references therein. In this case, asymptotic analysis for the M-estimator will be very different.

(1990), a classical central limit theorem combined with the Cramér–Wold device implies finite dimensional convergence, and a maximal inequality on a suitable class of functions guarantees stochastic asymptotic equicontinuity of the normalized empirical process. We adapt this approach to our local M-estimation problem with dependent observations.

Consider a function $\beta(\cdot)$ such that $\beta(t) = \beta_{[t]}$ if $t \geq 1$ and $\beta(t) = 1$ otherwise and denote its càdlàg inverse by $\beta^{-1}(\cdot)$. Let $Q_g(u)$ be the inverse function of the tail probability function $x \mapsto P\{|g(z_t)| > x\}$. For finite dimensional convergence, we employ Rio’s (1997), Corollary 1, central limit theorem for α -mixing arrays to our setup.

LEMMA C. *Suppose Assumption D holds true, $Pg_n = 0$ and*

$$(5) \quad \sup_{n \in \mathbb{N}} \int_0^1 \beta^{-1}(u) Q_{g_n}(u)^2 du < \infty.$$

Then $\Sigma = \lim_{n \rightarrow \infty} \text{Var}(\mathbb{G}_n g_n)$ exists and $\mathbb{G}_n g_n \xrightarrow{d} N(0, \Sigma)$.

The finite dimensional convergence of Z_n follows from Lemma C by setting g_n as any finite dimensional projection of the process $\{g_{n,s} - Pg_{n,s}\}$ with

$$(6) \quad g_{n,s} = n^{1/6} h_n^{2/3} (f_{n,\theta_0+s(nh_n)^{-1/3}} - f_{n,\theta_0}).$$

The requirement in (5) is the Lindeberg-type condition in Rio [(1997), Corollary 1] and excludes polynomial decay of β_m . Note that for criterion functions satisfying Assumption M, the $(2 + \delta)$ th moments $P|g_{n,s}|^{2+\delta}$ typically diverge because $g_{n,s}$ usually involves indicator functions. To verify (5), the following lemma is often useful.

LEMMA 2. *Suppose Assumptions M and D hold true and there is a positive constant c such that*

$$(7) \quad P\{|g_{n,s}| \geq c\} \leq c(nh_n^{-2})^{-1/3},$$

for all s and n large enough. Then (5) holds true.

In our examples, $g_{n,s}$ is zero or close to zero with high probability so that (7) is easily satisfied. See Section 3 for illustrations.

We provide another maximal inequality that is useful to establish stochastic asymptotic equicontinuity of the process Z_n .

LEMMA M'. *Suppose Assumption D holds true. Consider a class of functions $\mathcal{G}_n = \{g_{n,s} : |s| \leq K\}$ for some $K > 0$ with envelope G_n . Suppose there is a positive constant C such that*

$$(8) \quad P \sup_{s:|s-s'|<\varepsilon} |g_{n,s} - g_{n,s'}|^2 \leq C\varepsilon,$$

for all n large enough, $|s'| \leq K$, and $\varepsilon > 0$ small enough. Also, assume that there exist $0 \leq \kappa < 1/2$ and $C' > 0$ such that $G_n \leq C'n^\kappa$ and $\|G_n\|_2 \leq C'$ for all n large enough. Then for any $\sigma > 0$, there exist $\delta > 0$ and a positive integer N_δ such that

$$P \sup_{(s,s'):|s-s'|<\delta} |\mathbb{G}_n(g_{n,s} - g_{n,s'})| \leq \sigma,$$

for all $n \geq N_\delta$.

Stochastic asymptotic equicontinuity of Z_n is implied from this lemma by setting $g_{n,s}$ as in (6). Note that (8) is satisfied by Assumption **M**(iii).⁴ Compared to Lemma **M** used to derive the convergence rate of $\hat{\theta}$, Lemma **M'** is applied only to establish stochastic asymptotic equicontinuity of Z_n . Therefore, we do not need an exact decay rate on the right-hand side of the maximal inequality.⁵

By finite dimensional convergence and stochastic asymptotic equicontinuity of Z_n , its weak convergence is implied and the continuous mapping theorem for an argmax element [Theorem 2.7 of **Kim and Pollard (1990)**] yields the limiting distribution of $\hat{\theta}$. Define the covariance kernel

$$H(s_1, s_2) = \lim_{n \rightarrow \infty} \sum_{t=-n}^n \text{Cov}(g_{n,s_1}(z_t), g_{n,s_2}(z_t)),$$

if it exists. Throughout the paper, we use this notation for different choices of $g_{n,s}$. The main theorem of this section is presented as follows.

THEOREM 1. *Suppose that Assumptions **M** and **D** hold, $\hat{\theta}$ defined in (4) converges in probability to $\theta_0 \in \text{int } \Theta$, and (5) holds with $g_{n,s} - P g_{n,s}$ defined in (6) for each s . Then*

$$(9) \quad (nh_n)^{1/3}(\hat{\theta} - \theta_0) \xrightarrow{d} \arg \max_{s \in \mathbb{R}^d} Z(s),$$

where $Z(s)$ is a Gaussian process with continuous sample paths, expected value $s'Vs/2$, and covariance kernel $H(s_1, s_2)$.

This theorem can be considered as an extension of the main theorem of **Kim and Pollard (1990)** to the cases where the criterion function can vary with the sample size and the observations can obey a dependent process. To the best of our knowledge, the (nonparametric) cube root convergence rate $(nh_n)^{1/3}$ is new in the literature. It is interesting to note that similar to standard nonparametric estimation, nh_n still plays the role of the “effective sample size.”

⁴The upper bound in (8) can be relaxed to $\varepsilon^{1/p}$ for $1 \leq p < \infty$. However, it is typically satisfied with $p = 1$ for the examples we consider.

⁵In particular, Z_n itself does not satisfy Assumption **M**(ii).

2.1. *Nuisance parameters.* It is often the case that the criterion function contains some nuisance parameters, which can be estimated with rates faster than $(nh_n)^{1/3}$. For the rest of this section, let $\hat{\theta}$ and $\tilde{\theta}$ satisfy

$$\begin{aligned} \mathbb{P}_n f_{n,\hat{\theta},\hat{v}} &\geq \sup_{\theta \in \Theta} \mathbb{P}_n f_{n,\theta,\hat{v}} + o_p((nh_n)^{-2/3}), \\ \mathbb{P}_n f_{n,\tilde{\theta},v_0} &\geq \sup_{\theta \in \Theta} \mathbb{P}_n f_{n,\theta,v_0} + o_p((nh_n)^{-2/3}), \end{aligned}$$

respectively, where v_0 is a vector of nuisance parameters and \hat{v} is its estimator satisfying $\hat{v} - v_0 = o_p((nh_n)^{-1/3})$. Theorem 1 is extended as follows.

THEOREM 2. *Suppose Assumption D holds true. Let $\{f_{n,\theta,v_0} : \theta \in \Theta\}$ satisfy Assumption M and $\{f_{n,\theta,v} : \theta \in \Theta, v \in \Lambda\}$ satisfy Assumption M(iii). Also assume that there exists a negative definite matrix V_1 such that*

$$\begin{aligned} &P(f_{n,\theta,v} - f_{n,\theta_0,v_0}) \\ (10) \quad &= \frac{1}{2}(\theta - \theta_0)' V_1(\theta - \theta_0) + o(|\theta - \theta_0|^2) \\ &+ O(|v - v_0|^2) + o((nh_n)^{-2/3}), \end{aligned}$$

for all θ and v in neighborhoods of θ_0 and v_0 , respectively. Then $\hat{\theta} = \tilde{\theta} + o_p((nh_n)^{-1/3})$. Additionally, if (5) holds with $(g_{n,s} - P g_{n,s})$ for each s with $g_{n,s}$ being $n^{1/6} h_n^{2/3} (f_{n,\theta_0+s(nh_n)^{-1/3},v_0} - f_{n,\theta_0,v_0})$, then

$$(nh_n)^{1/3}(\hat{\theta} - \theta_0) \xrightarrow{d} \arg \max_{s \in \mathbb{R}^d} Z(s),$$

where $Z(s)$ is a Gaussian process with continuous sample paths, expected value $s' V_1 s/2$ and covariance kernel $H(s_1, s_2)$.

A key step for the proof of this theorem is to confirm that the empirical process $\mathbb{G}_n f_{n,\theta,v_0+c(nh_n)^{-1/3}}$ is well approximated by $\mathbb{G}_n f_{n,\theta,v_0}$ over $|\theta - \theta_0| \leq \epsilon$ and $|c| \leq \epsilon$ [see (A.10) in the supplementary material, Seo and Otsu (2018)]. This is shown by applying Lemma M' with $g_{n,s} = n^{1/6} h_n^{2/3} (f_{n,\theta,v_0+c(nh_n)^{-1/3}} - f_{n,\theta,v_0})$. Condition (8) in Lemma M' demands more precise control on the size of the envelope for the class of $g_{n,s}$ than the comparable condition in Z-estimation with nuisance parameters [e.g., equation (3) of van der Vaart and Wellner (2007)].

2.2. *Discussions.*

2.2.1. *Inference.* Once we show that the M-estimator has a proper limiting distribution, Politis, Romano and Wolf [(1999), Theorem 3.3.1] justify the use of subsampling to construct confidence intervals. Since Assumption D satisfies the

requirement of their theorem, subsampling inference based on s consecutive observations with $s/n \rightarrow \infty$ is asymptotically valid (in a pointwise sense explained below). See Politis, Romano and Wolf [(1999), Section 3.6] for a discussion on data-dependent choices of s .

We note that this asymptotic validity of subsampling inference is in a pointwise sense rather than uniform. To be specific, suppose $\{z_t\}$ is an independent and identically distributed (i.i.d.) sample from the probability measure P that belongs to a class of probability measures \mathcal{P} . Also denote the true parameters by $\theta_0(P)$ to make explicit the dependence on P . Based on Romano and Shaikh (2008), a confidence set \mathcal{C}_n for $\theta_0(P)$ is called *pointwise* valid in $(1 - \alpha)$ level if

$$\liminf_{n \rightarrow \infty} P\{\theta_0(P) \in \mathcal{C}_n\} \geq 1 - \alpha,$$

for each $P \in \mathcal{P}$ and is called *uniformly* valid in $(1 - \alpha)$ level if

$$\liminf_{n \rightarrow \infty} \inf_{P \in \mathcal{P}} P\{\theta_0(P) \in \mathcal{C}_n\} \geq 1 - \alpha.$$

Our Theorems 1 and 2 combined with Politis, Romano and Wolf [(1999), Theorem 3.3.1] guarantee the pointwise validity of the subsampling confidence set based on quantiles of the subsample statistic $(sh_s)^{1/3}(\hat{\theta}_s - \hat{\theta})$, where $\hat{\theta}_s$ and $\hat{\theta}$ are the M-estimators based on the subsample and full sample, respectively. Also a pointwise valid confidence interval for each element of $\theta_0(P)$ can be obtained in a similar manner.

To investigate whether we can construct a uniformly valid confidence set in our setup, we assume that $\{z_t\}$ is i.i.d. and the distribution $J_n(\cdot, \theta, P)$ of $Q_n(\theta) = (nh_n)^{2/3}\{\max_{\vartheta \in \Theta} \mathbb{P}_n f_{n,\vartheta} - \mathbb{P}_n f_{n,\theta}\}$ satisfies

$$(11) \quad \limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta} \sup_{P \in \mathcal{P}: \theta = \theta_0(P)} \sup_{x \in \mathbb{R}} \{J_s(x, \theta, P) - J_n(x, \theta, P)\} \leq 0,$$

Then Romano and Shaikh [(2008), Theorems 3.1 and 3.3] imply the uniform validity of the confidence set

$$\mathcal{C}_n = \{\theta \in \Theta : Q_n(\theta) \leq q_s(\theta, 1 - \alpha)\},$$

over \mathcal{P} , where $q_s(\theta, 1 - \alpha)$ is the $(1 - \alpha)$ th quantile of the distribution of the subsample statistic $Q_s(\theta)$. By inspection of Romano and Shaikh (2008), we can see that (11) is satisfied if $Q_n(\theta_0(P_n))$ converges in law to a unique continuous distribution for any sequence of $P_n \in \mathcal{P}$ yielding a row-wise i.i.d. triangular array. Our lemmas to obtain Theorem 1 can be readily extended to the array setting by restating Assumptions M and D and the additional conditions for Theorem 1 in the array setup. We note that computation of \mathcal{C}_n may require an extensive numerical search over Θ , where the quantile $q_s(\theta, 1 - \alpha)$ needs to be computed for each θ .

The above uniformity result relies upon the general results in Romano and Shaikh [(2008), Theorems 3.1 and 3.3] and there are at least three issues to be further considered. First, the i.i.d. assumption for the sample does not allow serial

dependence as in Assumption **D**. To accommodate dependent data, the high-level assumptions provided by Romano and Shaikh [(2008), Theorems 3.1] for uniform validity should be modified. Second, it is not a trivial task to extend the results in Romano and Shaikh (2008) to inference on subvectors (or functions) of θ except for a conservative projection of \mathcal{C}_n to a lower dimension. Third, a key result in Romano and Shaikh [(2008), Theorems 3.1] holds for objects centered at the true parameter $\theta_0(P)$ instead of the estimator $\hat{\theta}$. Therefore, their result does not apply to the subsample statistic $(sh_s)^{1/3}(\hat{\theta}_s - \hat{\theta})$. All of these issues require full length papers and are beyond the scope of this paper.

Another candidate to conduct inference based on the M-estimator is the bootstrap. However, even for independent observations, it is known that the naive non-parametric bootstrap is typically invalid under cube root asymptotics [Abrevaya and Huang (2005) and Sen, Banerjee and Woodroffe (2010)].

2.2.2. *Generalization of Assumption D.* All the results in this section build upon Assumption **D** which requires $\{z_t\}$ to be strictly stationary and absolutely regular (or β -mixing) with exponentially decaying mixing coefficients. Assumption **D** is used for both the maximal inequality (Lemma **M**) and central limit theorem (Lemma **C**) which are building blocks to derive the asymptotic distribution of $\hat{\theta}$. It is of interest whether we can establish analogous results under more general setups, such as α -mixing, by utilizing some recent developments in the empirical process theory for dependent data. For instance, Merlevède, Peligrad and Rio (2009, 2011) obtained Bernstein-type inequalities for α -mixing processes and Baraud (2010) and Nickl and Söhl [(2016), Section 3] explored the generic chaining argument by Talagrand (2005) for Markov chains.

Since the central limit theorem in Rio [(1997), Corollary 1] holds for α -mixing arrays, we can modify Lemma **C** to accommodate α -mixing processes. Thus, we focus on extending Lemma **M**, the maximal inequality. A crucial step for this extension is whether we can replace the key lemma in Doukhan, Massart and Rio [(1995), Lemma 3], which leads to the maximal inequality for β -mixing processes [in equation (A.6) of the supplementary material, Seo and Otsu (2018)] through a chaining argument. Specifically, consider a finite subclass \mathcal{F} of bounded functions with cardinality $p \geq \exp(1)$. By a decoupling technique for β -mixing processes, Doukhan, Massart and Rio [(1995), Lemma 3] showed that for positive constants c and c_1 , there exists a universal positive constant C such that

$$P \max_{f \in \mathcal{F}} |\mathbb{G}_n f| \leq C \left(c \sqrt{\log p} + c_1 q \frac{\log p}{\sqrt{n}} + c_1 \beta_q \sqrt{n} \right),$$

for all $q = 1, \dots, n$. Note that the above upper bound reduces to the first term $Cc\sqrt{\log p}$ for the i.i.d. case. By properly choosing q , the first term still dominates in the β -mixing case even if $\log p$ is close to n so that Lemma **M** can be established. In contrast, the maximal inequality implied by Merlevède, Peligrad and Rio [(2009), (2.1) in Theorem 1] for α -mixing would be written in the form of

$C(c\sqrt{\log p} + c_1 \log n \log \log n \frac{\log p}{\sqrt{n}})$. Therefore, as $\log p$ becomes close to n , the second term will dominate. Since this order of cardinality p (i.e., $\log p$ close to n) is required in the proof of Doukhan, Massart and Rio [(1995), Theorem 2], the upper bound in Lemma M for α -mixing processes would become larger.⁶

Another direction to extend our result is to accommodate general Markov chains that may not be covered by Assumption D. To this end, a chaining argument [see Baraud (2010) and Nickl and Söhl (2016)] based on Bernstein-type inequalities for Markov chains [e.g., Adamczak (2008) and Paulin (2015)] may yield an analog of Lemma M. Although this is an intriguing question, existing time series examples on cube root asymptotics mostly focus on mixing data [e.g., Polonik and Yao (2000) and de Jong and Woutersen (2011)] and also typically involve additional conditioning or exogenous variables. Thus, we leave this extension for future work.

3. Examples. We provide several examples to demonstrate the usefulness of the asymptotic theory in the last section. For the sake of space, we only sketch the arguments to verify the conditions to apply the theorems in Section 2. Detailed verifications under primitive conditions are relegated to the supplementary material [Seo and Otsu (2018)].

3.1. *Dynamic panel discrete choice.* For a binary response y_{it} and k -dimensional covariates x_{it} , consider a dynamic panel data model

$$P\{y_{i0} = 1|x_i, \alpha_i\} = F_0(x_i, \alpha_i),$$

$$P\{y_{it} = 1|x_i, \alpha_i, y_{i0}, \dots, y_{it-1}\} = F(x'_{it}\beta_0 + \gamma_0 y_{it-1} + \alpha_i),$$

for $i = 1, \dots, n$ and $t = 1, 2, 3$, where α_i is unobservable and both F_0 and F are unknown. Honoré and Kyriazidou (2000) proposed the conditional maximum score estimator $(\hat{\beta}, \hat{\gamma})$ that maximizes

$$\sum_{i=1}^n K\left(\frac{x_{i2} - x_{i3}}{b_n}\right)(y_{i2} - y_{i1}) \operatorname{sgn}\{(x_{i2} - x_{i1})'\beta + (y_{i3} - y_{i0})\gamma\},$$

where K is a kernel function and b_n is a bandwidth. Kernel smoothing is introduced to deal with the unknown link function F . Honoré and Kyriazidou (2000) obtained consistency of this estimator but the convergence rate and limiting distribution are unknown. Since the criterion function varies with the sample size through the bandwidth b_n , the cube root asymptotic theory of Kim and Pollard (1990) is not applicable here.

⁶Although a full investigation is beyond the scope of this paper, we conjecture that it is also the case for the generic chaining argument by Talagrand (2005). Indeed, equation (1.9) on page 10 of Talagrand (2005) explains that generic chaining needs partitions of cardinality up to 2^{2^n} .

This open question can be addressed by Theorem 1. Let $z = (z'_1, z_2, z'_3)'$ with $z_1 = x_2 - x_3$, $z_2 = y_2 - y_1$, and $z_3 = ((x_2 - x_1)', y_3 - y_0)'$. The above estimator for $\theta_0 = (\beta'_0, \gamma_0)'$ can be written as an M-estimator using the criterion function

$$(12) \quad f_{n,\theta}(z) = e_n(z)(\mathbb{I}\{z'_3\theta \geq 0\} - \mathbb{I}\{z'_3\theta_0 \geq 0\}),$$

where $e_n(z) = b_n^{-k}K(b_n^{-1}z_1)z_2$. To apply Theorem 1, it is enough to show that $f_{n,\theta}$ in (12) satisfies Assumption M with $h_n = b_n^k$ and the condition in (7). Then the limiting distribution of Honoré and Kyriazidou's (2000) estimator is obtained as in (9).

Here, we sketch the verification. See Section B.1 of the supplementary material [Seo and Otsu (2018)] for detailed verifications and primitive conditions. For Assumption M(i), $\{h_n f_{n,\theta}\}$ is bounded for the bounded kernel K and (2) is obtained by a Taylor expansion combined with the argument in Kim and Pollard [(1990), pages 214–215]. For Assumption M(ii), take any θ_1 and θ_2 and note that

$$\begin{aligned} h_n^{1/2} \|f_{n,\theta_1} - f_{n,\theta_2}\|_2 &= \sqrt{P\{h_n E[e_n(z)^2|z_3]|\mathbb{I}\{z'_3\theta_1 \geq 0\} - \mathbb{I}\{z'_3\theta_2 \geq 0\}|\}} \\ &\geq c^{1/2} P\{z'_3\theta_1 \geq 0 > z'_3\theta_2 \text{ or } z'_3\theta_2 \geq 0 > z'_3\theta_1\}, \end{aligned}$$

for some $c > 0$, where the inequality follows from $h_n E[e_n(z)^2|z_3] > c$ [by a change of variables and the condition on the density $z_1|(z_2 \neq 0, z_3)$ being bounded away from zero] and Jensen's inequality. The right-hand side is the probability for a pair of wedge shaped regions with an angle of order $|\theta_1 - \theta_2|$. Thus, Assumption M(ii) is satisfied if the density of z_3 is bounded away from zero in a neighborhood of the origin. Assumption M(iii) can be verified in a similar way (by considering the upper bound instead). The Markov inequality and boundedness of the density imply (7).

3.2. *Random coefficient binary choice.* As a new statistical model which can be covered by our asymptotic theory, consider a regression model $y_t = x'_t\theta(w_t) + u_t$ with random coefficients. Suppose we observe $\{\text{sgn}(y_t), x_t, w_t\}$ and wish to estimate $\theta_0 = \theta(c)$ at some given c .⁷ We propose a localized version of the maximum score estimator

$$(13) \quad \hat{\theta} = \arg \max_{\theta \in S} \sum_{t=1}^n K\left(\frac{w_t - c}{b_n}\right) [\mathbb{I}\{y_t \geq 0, x'_t\theta \geq 0\} + \mathbb{I}\{y_t < 0, x'_t\theta < 0\}],$$

where S is the surface of the unit sphere. Again, the cube root asymptotic theory of Kim and Pollard (1990) is not applicable due to the bandwidth.

⁷Gautier and Kitamura (2013) studied identification and estimation of the random coefficient binary choice model, where $\theta_t = \theta(w_t)$ is unobservable. Here, we study the model where heterogeneity in the slope is caused by the observables w_t .

Theorem 1 can be applied to obtain the limiting distribution of this estimator. Note that $\hat{\theta}$ in (13) can be written as an M-estimator using the criterion function

$$(14) \quad f_{n,\theta}(x, w, u) = \frac{1}{h_n} K\left(\frac{w - c}{h_n^{1/k}}\right) h(x, u) [\mathbb{I}\{x'\theta \geq 0\} - \mathbb{I}\{x'\theta_0 \geq 0\}],$$

for $h_n = b_n^k$ and $h(x, u) = \mathbb{I}\{x'\theta_0 + u \geq 0\} - \mathbb{I}\{x'\theta_0 + u < 0\}$. Once we check Assumption M and (7), Theorem 1 implies the limiting distribution.

The verification is sketched as follows. See Section B.2 of the supplementary material [Seo and Otsu (2018)] for detailed verifications and primitive conditions. Assumption M(i)–(ii) and (7) can be checked similarly as in Section 3.1. Here, we verify Assumption M(iii). By a change of variables and $h(x, u)^2 = 1$, there exists a positive constant C' such that

$$\begin{aligned} P \sup_{\theta \in \Theta: |\theta - \vartheta| < \varepsilon} h_n |f_{n,\theta} - f_{n,\vartheta}|^2 &= \iint K(s)^2 \sup_{\theta \in \Theta: |\theta - \vartheta| < \varepsilon} |[\mathbb{I}\{x'\theta \geq 0\} - \mathbb{I}\{x'\vartheta \geq 0\}]|^2 p(x, c + sb_n) dx ds \\ &\leq C' E \left[\sup_{\theta \in \Theta: |\theta - \vartheta| < \varepsilon} |[\mathbb{I}\{x'\theta \geq 0\} - \mathbb{I}\{x'\vartheta \geq 0\}]|^2 |w = c \right], \end{aligned}$$

for all $\varepsilon > 0$, ϑ in a neighborhood of θ_0 , and n large enough, where p is the joint density of (x_t, w_t) . Since the right-hand side is the conditional probability for a pair of wedge shaped regions with an angle of order ε , Assumption M(iii) is guaranteed by some boundedness condition on the conditional density of x_t given $w_t = c$.

3.3. *Minimum volume predictive region.* As an illustration of Theorem 2, we now consider the example in (1), a simplified version of Polonik and Yao’s (2000) minimum volume predictor. For notational convenience, assume $\theta_0 = 0$ and $v_0 = 1$. By applying Lemma M’, the convergence rate of the nuisance parameter estimator is obtained as $\hat{v} - 1 = O_p((nh_n)^{-1/2} + h_n^2)$ [see Section B.3 in the supplementary material, Seo and Otsu (2018)].

The criterion function for the maximization in (1) can be written as

$$f_{n,\theta,\hat{v}}(y, x) = \frac{1}{h_n} K\left(\frac{x - c}{h_n}\right) [\mathbb{I}\{y \in [\theta - \hat{v}, \theta + \hat{v}]\} - \mathbb{I}\{y \in [-\hat{v}, \hat{v}]\}].$$

We apply Theorem 2 to obtain the convergence rate of $\hat{\theta}$. Details are provided in Section B.3 of the supplementary material. Assumptions M for $f_{n,\theta,1}$ and M(iii) for $f_{n,\theta,v}$ are verified similarly as in Sections 3.1 and 3.2. To check (10), a Taylor expansion yields

$$\begin{aligned} P(f_{n,\theta,v} - f_{n,0,1}) &= \frac{1}{2} V_1 \theta^2 + \{\dot{\gamma}_{y|x}(1|c) + \dot{\gamma}_{y|x}(-1|c)\} \gamma_x(c) \theta v + o(\theta^2 + |v - 1|^2) + O(h_n^2), \end{aligned}$$

for $V_1 = \{\dot{\gamma}_{y|x}(1|c) - \dot{\gamma}_{y|x}(-1|c)\}\gamma_x(c)$, where γ and $\dot{\gamma}$ mean the density and its derivative, respectively.

Therefore, Theorem 2 implies $\hat{\theta} - \theta_0 = O_p((nh_n)^{-1/3} + h_n)$, which confirms positively the conjecture of Polonik and Yao [(2000), Remark 3] on the exact convergence rate of $[\hat{\theta} \pm \hat{\nu}]$.

3.4. *Dynamic maximum score.* To illustrate the derivation of the covariance kernel H in Theorem 1 for dependent data, we consider the maximum score estimator [Manski (1975)] for a regression model $y_t = x_t'\theta_0 + u_t$, that is,

$$\hat{\theta} = \arg \max_{\theta \in S} \sum_{t=1}^n [\mathbb{I}\{y_t \geq 0, x_t'\theta \geq 0\} + \mathbb{I}\{y_t < 0, x_t'\theta < 0\}],$$

where S is the surface of the unit sphere. This estimator can be written as an M-estimator using the criterion function

$$f_{\theta}(x, u) = h(x, u)[\mathbb{I}\{x'\theta \geq 0\} - \mathbb{I}\{x'\theta_0 \geq 0\}],$$

where $h(x, u) = \mathbb{I}\{x'\theta_0 + u \geq 0\} - \mathbb{I}\{x'\theta_0 + u < 0\}$. The conditions to apply Theorem 1 can be verified similarly as in the above examples [see Section B.4 of the supplementary material, Seo and Otsu (2018)]. Here we focus on the derivation of the covariance kernel for the limiting distribution under Assumption D.

Let $q_{n,t} = f_{\theta_0+n^{-1/3}s_1}(x_t, u_t) - f_{\theta_0+n^{-1/3}s_2}(x_t, u_t)$. The covariance kernel is written as $H(s_1, s_2) = \frac{1}{2}\{L(s_1, 0) + L(0, s_2) - L(s_1, s_2)\}$, where

$$L(s_1, s_2) = \lim_{n \rightarrow \infty} n^{1/3} \left\{ \text{Var}(q_{n,t}) + \sum_{m=1}^{\infty} \text{Cov}(q_{n,t}, q_{n,t+m}) \right\}.$$

The limit of $n^{1/3} \text{Var}(q_{n,t})$ is given in Kim and Pollard [(1990), page 215]. For the covariance, we note that $q_{n,t}$ takes only three values: $-1, 0$ or 1 . The definition of β_m and Assumption D imply

$$|P\{q_{n,t} = j, q_{n,t+m} = k\} - P\{q_{n,t} = j\}P\{q_{n,t+m} = k\}| \leq n^{-2/3} \beta_m,$$

for all $n, m \geq 1$ and $j, k = -1, 0, 1$. Thus, $\{q_{n,t}\}$ is a β -mixing array with mixing coefficients bounded by $n^{-2/3} \beta_m$. This in turn implies that $\{q_{n,t}\}$ is an α -mixing array with mixing coefficients bounded by $2n^{-2/3} \beta_m$. By applying the α -mixing inequality, the covariance is bounded as

$$\text{Cov}(q_{n,t}, q_{n,t+m}) \leq Cn^{-2/3} \beta_m \|q_{n,t}\|_p^2,$$

for some $C > 0$ and $p > 2$. Note that

$$\begin{aligned} \|q_{n,t}\|_p^2 &\leq [P\{\mathbb{I}\{x'(\theta_0 + s_1 n^{-1/3}) > 0\} - \mathbb{I}\{x'(\theta_0 + s_2 n^{-1/3}) > 0\}\}]^{2/p} \\ &= O(n^{-2/(3p)}). \end{aligned}$$

Combining these results, we get $n^{1/3} \sum_{m=1}^{\infty} \text{Cov}(q_{n,t}, q_{n,t+m}) \rightarrow 0$ as $n \rightarrow \infty$. Therefore, the covariance kernel H is the same as the independent case in Kim and Pollard [(1990), page 215].

3.5. *Other examples.* In the supplementary material [Seo and Otsu (2018)], we present additional examples on the dynamic least median of squares estimator (Section B.5) and the monotone density estimator (Section B.6).

4. Generalizations. In this section, we consider two generalizations of the asymptotic theory in Section 2. The first concerns data taking limited values such as interval-valued regressors and the second is to allow for localization to depend on the parameter values.

4.1. *Limited observations.* We consider the case where some of the variables take limited values. In particular, we relax the assumption of point identification of θ_0 and study the case where the limiting criterion function is maximized at any element of a set $\Theta_I \subset \Theta$. The set Θ_I is called the identified set. In order to estimate Θ_I , we consider a collection of approximate maximizers of the sample criterion function

$$\hat{\Theta} = \left\{ \theta \in \Theta : \max_{\theta \in \Theta} \mathbb{P}_n f_{n,\theta} - \mathbb{P}_n f_{n,\theta} \leq \hat{c}(nh_n)^{-1/2} \right\},$$

that is, the level set based on $\mathbb{P}_n f_{n,\theta}$ from its maximum by a cutoff value $\hat{c}(nh_n)^{-1/2}$. This section studies the convergence rate of $\hat{\Theta}$ to Θ_I under the Hausdorff distance defined below. We assume that Θ_I is convex. Then the projection $\pi_\theta = \arg \min_{\theta' \in \Theta_I} |\theta' - \theta|$ of θ on Θ_I is uniquely defined. To deal with the partially identified case, we modify Assumption M as follows.

ASSUMPTION S. For a sequence $\{h_n\}$ of positive numbers satisfying $nh_n \rightarrow \infty$, $f_{n,\theta}$ satisfies the following conditions:

(i) $h_n f_{n,\theta}$ is uniformly bounded, $\lim_{n \rightarrow \infty} P f_{n,\theta}$ is maximized at any θ in a bounded convex set Θ_I , and there exist positive constants c and c' such that

$$(15) \quad P(f_{n,\pi_\theta} - f_{n,\theta}) \geq c|\theta - \pi_\theta|^2 + o(|\theta - \pi_\theta|^2) + o((nh_n)^{-2/3}),$$

for all n large enough and all $\theta \in \{\Theta : 0 < |\theta - \pi_\theta| \leq c'\}$.

(ii) There exist positive constants C and C' such that

$$|\theta - \pi_\theta| \leq Ch_n^{1/2} \|f_{n,\theta} - f_{n,\pi_\theta}\|_2,$$

for all n large enough and all $\theta \in \{\Theta : 0 < |\theta - \pi_\theta| \leq C'\}$.

(iii) There exists a positive constant C'' such that

$$P \sup_{\theta \in \Theta: 0 < |\theta - \pi_\theta| < \varepsilon} h_n |f_{n,\theta} - f_{n,\pi_\theta}|^2 \leq C'' \varepsilon,$$

for all n large enough and all $\varepsilon > 0$ small enough.

We allow $h_n = 1$ for the case without a bandwidth in the criterion function. Similar comments to Assumption **M** apply. The main difference is that the conditions are imposed on the contrast $f_{n,\theta} - f_{n,\pi_\theta}$ using the projection π_θ . Assumption **S(i)** contains boundedness and expansion conditions. The inequality in (15) can be checked by a one-sided Taylor expansion using the directional derivative. Assumption **S(ii)** and (iii) play similar roles as Assumption **M(ii)** and (iii) and can be verified in a similar way.

We first establish the maximal inequality for the criterion functions satisfying Assumption **S**. Let $r_n = nh_n / \log(nh_n)$.

LEMMA MS. *Under Assumptions **D** and **S**, there exist positive constants C and $C' < 1$ such that*

$$P \sup_{\theta \in \Theta: 0 < |\theta - \pi_\theta| < \delta} |\mathbb{G}_n h_n^{1/2} (f_{n,\theta} - f_{n,\pi_\theta})| \leq C(\delta \log(1/\delta))^{1/2},$$

for all n large enough and $\delta \in [r_n^{-1/2}, C']$.

Compared to Lemma **M**, the additional log term on the right-hand side is due to the fact that the supremum is taken over the δ -tube (or manifold) instead of the δ -ball, which increases the entropy. This maximal inequality is applied to obtain an analog of Lemma 1.

LEMMA 3. *Under Assumptions **D** and **S**, for each $\varepsilon > 0$, there exist random variables $\{R_n\}$ of order $O_p(1)$ and a positive constant C such that*

$$|\mathbb{P}_n(f_\theta - f_{\pi_\theta}) - P(f_\theta - f_{\pi_\theta})| \leq \varepsilon |\theta - \pi_\theta|^2 + r_n^{-2/3} R_n^2,$$

for all $\theta \in \{\Theta : r_n^{-1/3} \leq |\theta - \pi_\theta| \leq C\}$.

Let $\rho(A, B) = \sup_{a \in A} \inf_{b \in B} |a - b|$ and $H(A, B) = \max\{\rho(A, B), \rho(B, A)\}$ be the Hausdorff distance of sets $A, B \subset \mathbb{R}^d$. Based on these lemmas, the convergence rate of the set estimator $\hat{\Theta}$ is obtained as follows.

THEOREM 3. *Suppose that Assumptions **D** and **S** hold, $H(\hat{\Theta}, \Theta_I) \xrightarrow{p} 0$, $\{h_n^{1/2} f_{n,\theta} : \theta \in \Theta_I\}$ is P -Donsker and $\hat{c} = o_p((nh_n)^{1/2})$. Then*

$$\rho(\hat{\Theta}, \Theta_I) = O_p(\hat{c}^{1/2} (nh_n)^{-1/4} + r_n^{-1/3}).$$

Furthermore, if $\hat{c} \rightarrow \infty$, then $P\{\Theta_I \subset \hat{\Theta}\} \rightarrow 1$ and

$$H(\hat{\Theta}, \Theta_I) = O_p(\hat{c}^{1/2} (nh_n)^{-1/4}).$$

Note that ρ is asymmetric in its arguments. In contrast to the convergence rate of $\rho(\hat{\Theta}, \Theta_I)$ obtained in the first part of this theorem, the second part says $P\{\Theta_I \subset \hat{\Theta}\} \rightarrow 1$ [i.e., $\rho(\Theta_I, \hat{\Theta})$ can converge to zero at an arbitrary rate] as far as $\hat{c} \rightarrow \infty$. For example, we may set $\hat{c} = \log(nh_n)$. These results are combined to imply the convergence rate $H(\hat{\Theta}, \Theta_I) = O_p(\hat{c}^{1/2}(nh_n)^{-1/4})$ under the Hausdorff distance. When $\hat{c} \rightarrow \infty$, the cube root term of order $r_n^{-1/3}$ in the rate of $\rho(\hat{\Theta}, \Theta_I)$ is dominated by the term of order $\hat{c}^{1/2}(nh_n)^{-1/4}$.

We next consider the case where the criterion function contains nuisance parameters. In particular, we allow the dimension k_n of the nuisance parameters ν to grow as the sample size increases. For instance, the nuisance parameters might be coefficients in a sieve estimation procedure. It is important to allow the growing dimension of ν to cover [Manski and Tamer's \(2002\)](#) set estimator, where the criterion function contains some nonparametric estimate and its transform by the indicator. The rest of this subsection considers the set estimator

$$\hat{\Theta} = \left\{ \theta \in \Theta : \max_{\theta \in \Theta} \mathbb{P}_n f_{n,\theta,\hat{\nu}} - \mathbb{P}_n f_{n,\theta,\hat{\nu}} \leq \hat{c}(nh_n)^{-1/2} \right\},$$

with some preliminary estimator $\hat{\nu}$ and cutoff value \hat{c} .

Let $g_{n,s} = h_n^{1/2}(f_{n,\theta,\nu} - f_{n,\theta,\nu_0})$ with $s = (\theta', \nu')$ and consider $\mathcal{G}_n = \{g_{n,s} : |\theta - \pi_\theta| \leq K_1, |\nu - \nu_0| \leq a_n K_2\}$ for some $K_1, K_2 > 0$ with the envelope function $G_n = \sup_{\mathcal{G}_n} |g_{n,s}|$. The maximal inequality in [Lemma MS](#) is modified as follows.

LEMMA MS'. *Suppose Assumption D holds true and there exists a positive constant C such that*

$$(16) \quad P \sup_{s:\theta \in \Theta, |\nu - \nu_0| \leq \varepsilon} |g_{n,s}|^2 \leq C\sqrt{k_n}\varepsilon,$$

$$(17) \quad \sup_{s:\theta \in \Theta, |\nu - \nu_0| \leq \varepsilon} \{|\nu - \nu_0| - C\|g_{n,s}\|_2\} \leq 0,$$

for all n large enough and ε small enough. Also assume that there exist $0 \leq \kappa < 1/4$ and $C' > 0$ such that $G_n \leq C'n^\kappa$ and $\|G_n\|_2 \leq C'$ for all n large enough. Then there exists $K_3 > 0$ such that

$$P \sup_{g_{n,s} \in \mathcal{G}_n} |\mathbb{G}_n g_{n,s}| \leq K_3 a_n^{1/2} k_n^{3/4} \sqrt{\log k_n a_n^{-1}},$$

for all n large enough.

The increasing dimension k_n of ν affects the upper bound via two routes. First, it increases the size of envelope by a factor of $\sqrt{k_n}$, which in turn increases the entropy of the space. Second, it also demands us to consider an inflated class of functions to apply the more fundamental maximal inequality by [Doukhan, Massart and Rio \(1995\)](#), which relies on the $\|\cdot\|_{2,\beta}$ norm. Note that the envelope condition in (16) allows for step functions containing some nonparametric estimates.

Based on this lemma, the convergence rate of the set estimator $\hat{\Theta}$ is characterized as follows.

THEOREM 4. *Suppose Assumption D holds true. Let $\{f_{n,\theta,v_0} : \theta \in \Theta\}$ satisfy Assumption S and $\{h_n^{1/2} f_{n,\theta,v_0} : \theta \in \Theta_I\}$ be a P -Donsker class. Assume $\rho(\hat{\Theta}, \Theta_I) \xrightarrow{P} 0$, $\hat{c} = o_p((nh_n)^{1/2})$, $k_n \rightarrow \infty$, and $|\hat{v} - v_0| = o_p(a_n)$ for some $\{a_n\}$ such that $h_n/a_n \rightarrow \infty$. Furthermore, there exist some $\varepsilon > 0$ and neighborhoods $\{\theta \in \Theta : |\theta - \pi_\theta| < \varepsilon\}$ and $\{v : |v - v_0| \leq \varepsilon\}$, where $h_n^{1/2}(f_{n,\theta,v} - f_{n,\theta,v_0})$ satisfies (16) and (17) and*

$$(18) \quad \begin{aligned} &P(f_{n,\theta,v} - f_{n,\pi_\theta,v} - f_{n,\theta,v_0} + f_{n,\pi_\theta,v_0}) \\ &= o(|\theta - \pi_\theta|^2) + O(|v - v_0|^2 + r_n^{-2/3}). \end{aligned}$$

Then

$$(19) \quad \begin{aligned} \rho(\hat{\Theta}, \Theta_I) &= O_p(\hat{c}^{1/2}(nh_n)^{-1/4} + r_n^{-1/3}) \\ &+ O_p((nh_n a_n^{-1})^{-1/4} (\log k_n)^{1/2}) + o(a_n). \end{aligned}$$

Furthermore, if $\hat{c} \rightarrow \infty$, then $P\{\Theta_I \subset \hat{\Theta}\} \rightarrow 1$ and

$$(20) \quad H(\hat{\Theta}, \Theta_I) = O_p(\hat{c}^{1/2}(nh_n)^{-1/4} + (nh_n)^{-1/4} a_n^{1/4} k_n^{3/8} \log^{1/4} n) + o(a_n).$$

Compared to Theorem 3, we have two extra terms in (20) due to the (non-parametric) estimation of v_0 . However, they can be shown to be dominated by the first term under standard conditions. Suppose that $k_n^4 \log k_n/n \rightarrow 0$ and the preliminary estimator \hat{v} satisfies $|\hat{v} - v_0| = O_p(n^{-1/2}(k_n \log k_n)^{1/2})$, which is often the case as in sieve estimation [see, e.g., [Chen \(2007\)](#)].⁸ Then we can set $a_n = n^{-1/2}(k_n \log k_n)^{1/2}$ so that $a_n^{1/4} k_n^{3/8} \rightarrow 0$. Now by choosing $\hat{c} = \log n$, the first term in (20) dominates the other terms.

4.1.1. Example: Binary choice with interval regressor. As an illustration of partially identified models, we consider a binary choice model with an interval-valued regressor studied by [Manski and Tamer \(2002\)](#). Let $y = \mathbb{I}\{x'\theta_0 + w + u \geq 0\}$ where x is a vector of observable regressors, w is an unobservable regressor, and u is an unobservable error term satisfying $P\{u \leq 0|x, w\} = \alpha$ (we set $\alpha = 0.5$ to simplify the notation). Instead of w , we observe an interval $[w_l, w_u]$ such that $P\{w_l \leq w \leq w_u\} = 1$. Here, we normalize the coefficient of w to be one. In this setup, the parameter θ_0 is partially identified and its identified set is written as [[Manski and Tamer \(2002\)](#), Proposition 2]

$$\Theta_I = \{\theta \in \Theta : P\{x'\theta + w_u \leq 0 < x'\theta_0 + w_l \text{ or } x'\theta_0 + w_u \leq 0 < x'\theta + w_l\} = 0\}.$$

⁸Alternatively, v_0 can be estimated by some high-dimensional method [e.g., [Belloni, Chen, Chernozhukov and Hansen \(2012\)](#)] which also typically guarantees $a_n = o(n^{-1/4})$.

Let $\tilde{x} = (x', w_l, w_u)'$ and $q_{\hat{v}}(\tilde{x})$ be an estimator of $q_{v_0}(\tilde{x}) = P\{y = 1|\tilde{x}\}$ with the estimated parameters \hat{v} . By exploring the maximum score approach, **Manski and Tamer (2002)** developed the set estimator for Θ_I

$$(21) \quad \hat{\Theta} = \left\{ \theta \in \Theta : \max_{\theta \in \Theta} S_n(\theta) - S_n(\theta) \leq \epsilon_n \right\},$$

where

$$S_n(\theta) = \mathbb{P}_n(y - 0.5) [\mathbb{I}\{q_{\hat{v}}(\tilde{x}) > 0.5\} \operatorname{sgn}(x'\theta + w_u) + \mathbb{I}\{q_{\hat{v}}(\tilde{x}) \leq 0.5\} \operatorname{sgn}(x'\theta + w_l)].$$

Manski and Tamer (2002) established $H(\hat{\Theta}, \Theta_I) \xrightarrow{P} 0$ by assuming that the cutoff value ϵ_n is bounded from below by the (almost sure) uniform convergence rate of $S_n(\theta)$ to the limiting object. As **Manski and Tamer [(2002), Footnote 3]** argued characterization of this rate is a complex task because $S_n(\theta)$ is a step function and $\mathbb{I}\{q_{\hat{v}}(\tilde{x}) > 0.5\}$ is a step function transform of the nonparametric estimate of $P\{y = 1|\tilde{x}\}$. As such, it has been an open question. Obtaining the lower bound of ϵ_n is important because we wish to minimize the volume of the estimator $\hat{\Theta}$ without losing the asymptotic validity. By applying **Theorem 4**, we can explicitly characterize the lower bound of ϵ_n and establish the convergence rate of $\hat{\Theta}$.

A little algebra shows that the set estimator in (21) is written as

$$\hat{\Theta} = \left\{ \theta \in \Theta : \max_{\theta \in \Theta} \mathbb{P}_n f_{\theta, \hat{v}} - \mathbb{P}_n f_{\theta, \hat{v}} \leq \hat{c}n^{-1/2} \right\},$$

where $z = (x', w, w_l, w_u, u)'$, $h(x, w, u) = \mathbb{I}\{x'\theta_0 + w + u \geq 0\} - \mathbb{I}\{x'\theta_0 + w + u < 0\}$, and

$$f_{\theta, v}(z) = h(x, w, u) [\mathbb{I}\{x'\theta + w_u \geq 0, q_v(\tilde{x}) > 0.5\} - \mathbb{I}\{x'\theta + w_l < 0, q_v(\tilde{x}) \leq 0.5\}].$$

To apply **Theorem 4**, we check **Assumption S** with $h_n = 1$. See **Section B.7** of the supplementary material [**Seo and Otsu (2018)**] for details. Here, we illustrate the verifications of (16) and (18). Let $I_v(\tilde{x}) = \mathbb{I}\{q_v(\tilde{x}) > 0.5 \geq q_{v_0}(\tilde{x}) \text{ or } q_v(\tilde{x}) \leq 0.5 < q_{v_0}(\tilde{x})\}$. Note that

$$|f_{\theta, v} - f_{\theta, v_0}|^2 \leq \mathbb{I}\{x'\theta \geq -w_u \geq x'\pi_\theta \text{ or } x'\theta < -w_u < x'\pi_\theta\} I_v(\tilde{x}) \leq I_v(\tilde{x}).$$

Also, we have

$$P \sup_{v \in \Lambda: |v - v_0| < \epsilon} I_v(\tilde{x}) \leq CP \sup_{v \in \Lambda: |v - v_0| < \epsilon} |q_v(\tilde{x}) - q_{v_0}(\tilde{x})| \leq C\sqrt{k_n}\epsilon,$$

for some $C > 0$, where the first inequality holds under boundedness of the conditional density of $q_{v_0}(\tilde{x})$ and the second under smoothness of q_v . This verifies (16). Also, (17) is verified in the same manner as **Assumption S(ii)**.

For (18), note that

$$\begin{aligned}
 & |P(f_{\theta,v} - f_{\theta,v_0}) - P(f_{\pi_\theta,v} - f_{\pi_\theta,v_0})| \\
 (22) \quad & \leq P\mathbb{I}\{x'\theta \geq -w_u \geq x'\pi_\theta \text{ or } x'\theta < -w_u < x'\pi_\theta\}I_v(\tilde{x}) \\
 & \quad + P\mathbb{I}\{x'\theta \geq -w_l \geq x'\pi_\theta \text{ or } x'\theta < -w_l < x'\pi_\theta\}I_v(\tilde{x}),
 \end{aligned}$$

for each $\theta \in \{\Theta : |\theta - \pi_\theta| < \varepsilon\}$ and v in a neighborhood of v_0 . For the first term of (22), the law of iterated expectation and an expansion of $q_v(\tilde{x})$ around v_0 imply

$$\begin{aligned}
 & P\mathbb{I}\{x'\theta \geq -w_u \geq x'\pi_\theta \text{ or } x'\theta < -w_u < x'\pi_\theta\}I_v(\tilde{x}) \\
 & \leq P\mathbb{I}\{x'\theta \geq -w_u \geq x'\pi_\theta \text{ or } x'\theta < -w_u < x'\pi_\theta\}A(w_u, x)|v - v_0|,
 \end{aligned}$$

for some bounded function A . The second term of (22) is bounded in the same manner. Therefore, $|P(f_{\theta,v} - f_{\theta,v_0}) - P(f_{\pi_\theta,v} - f_{\pi_\theta,v_0})| = O(|\theta - \pi_\theta||v - v_0|)$ and (18) is verified. By applying Theorem 4, we can conclude that the convergence rate of Manski and Tamer’s (2002) set estimator $\hat{\Theta}$ in (21) is characterized by (19) and (20).

Compared to Manski and Tamer (2002), we provide a sharper lower bound on the the tuning parameter ϵ_n , which is $\hat{c}n^{-1/2}$ with $\hat{c} \rightarrow \infty$. For example, if we set $\hat{c} = \log n$, the convergence rate becomes $H(\hat{\Theta}, \Theta_I) = O_p(n^{-1/4}(\log n)^{1/2})$. We basically verify the high-level assumption of Chernozhukov, Hong and Tamer [(2007), Condition C.2] in the cube root context. However, we mention that in the above setup, the criterion function contains nuisance parameters with increasing dimension and the result in Chernozhukov, Hong and Tamer (2007) does not apply directly.

Furthermore, our result enables us to construct a confidence set by subsampling as described in Chernozhukov, Hong and Tamer (2007). Specifically, the maximal inequality in Lemma MS’ and the assumption that $\{h_n^{1/2} f_{n,\theta,v_0} : \theta \in \Theta_I\}$ is P -Donsker are sufficient to satisfy their Conditions C.4 and C.5.

4.2. *Parameter-dependent localization.* We now consider a setup where localization of the criterion function depends on the parameter values. A leading example is mode estimation. Chernoff (1964) studied asymptotic properties of the mode estimator that maximizes $(nh)^{-1} \sum_{t=1}^n \mathbb{I}\{|y_t - \beta| \leq h\}$ with respect to β for some fixed h . Lee (1989) extended this estimator to regression models, established its consistency, and conjectured the cube root convergence rate. To estimate β consistently for a broader family of distributions, however, we need to treat h as a bandwidth parameter and let $h \rightarrow 0$ as in Yao, Lindsay and Li (2012), for example.

This parameter-dependent localization alters Assumption M(iii) because it increases the size [in terms of the $L_2(P)$ -norm] of the envelope of the class $\{h^{-1}(\mathbb{I}\{|y_t - \beta| \leq h\} - \mathbb{I}\{|y_t - \beta_0| \leq h\}) : |\beta - \beta_0| \leq \varepsilon\}$. More precisely, we replace Assumption M(iii) with the following.

ASSUMPTION M. (iii') *There exists a positive constant C'' such that*

$$P \sup_{\theta \in \Theta: |\theta - \theta'| < \varepsilon} h_n^2 |f_{n,\theta} - f_{n,\theta'}|^2 \leq C'' \varepsilon,$$

for all n large enough, $\varepsilon > 0$ small enough, and θ' in a neighborhood of θ_0 .

Under this assumption, Lemma M in Section 2 is modified as follows.

LEMMA M1. *Under Assumption M(i), (ii) and (iii'), there exist positive constants C and C' such that*

$$P \sup_{\theta \in \Theta: |\theta - \theta_0| < \delta} |\mathbb{G}_n h_n^{1/2} (f_{n,\theta} - f_{n,\theta_0})| \leq C h_n^{-1/2} \delta^{1/2},$$

for all n large enough and $\delta \in [(nh_n^2)^{-1/2}, C']$.

Parameter dependency arises in different contexts and may lead to different types of nonstandard distributions. For instance, the maximum likelihood estimator for Uniform $[0, \theta]$ yields super consistency [see Hirano and Porter (2003) for a general discussion]. This contrast is similar to the difference between estimation of a change point in regression analysis and mode regression.

Once we have obtained this lemma, the remaining steps are similar to those in Section 2 by replacing “ h_n ” with “ h_n^2 .” Here, we present the result without nuisance parameters ν for the sake of expositional simplicity.

THEOREM 5. *Suppose that Assumptions D and M(i), (ii) and (iii') hold. Also suppose (5) holds with $(g_{n,s} - P g_{n,s})$ for each s , where $g_{n,s} = n^{1/6} h_n^{4/3} \times (f_{n,\theta_0+s(nh_n^2)^{-1/3}} - f_{n,\theta_0})$. Then*

$$(23) \quad (nh_n^2)^{1/3} (\hat{\theta} - \theta_0) \xrightarrow{d} \arg \max_{s \in \mathbb{R}^d} Z(s),$$

where $Z(s)$ is a Gaussian process with continuous sample paths, expected value $s'Vs/2$ and covariance kernel $H(s_1, s_2)$.

4.2.1. *Example: Hough transform estimator.* In the statistics literature on the computer vision algorithm, Goldenshluger and Zeevi (2004) investigated the so-called Hough transform estimator for regression models

$$(24) \quad \hat{\beta} = \arg \max_{\beta \in B} \sum_{t=1}^n \mathbb{I}\{|y_t - x_t' \beta| \leq h|x_t|\},$$

where B is some parameter space, $x_t = (1, \tilde{x}_t)'$ for a scalar \tilde{x}_t and h is a fixed tuning constant. Goldenshluger and Zeevi (2004) derived the cube root asymptotics for $\hat{\beta}$ with fixed h and discussed carefully the practical choice of h . However, for

this estimator, h plays the role of the bandwidth and the analysis for the case of $h_n \rightarrow 0$ is a substantial open question [see pages 1915–1916 of [Goldenshluger and Zeevi \(2004\)](#)]. Here we study the asymptotic property of $\hat{\beta}$ in (24) with $h = h_n \rightarrow 0$. The estimators by [Chernoff \(1964\)](#) and [Lee \(1989\)](#) with varying h can be analyzed in the same manner.

Let $z = (x', u)'$. Note that $\hat{\theta} = \hat{\beta} - \beta_0$ is written as an M-estimator using the criterion function

$$f_{n,\theta}(z) = h_n^{-1} \mathbb{I}\{|u - x'\theta| \leq h_n|x|\}.$$

To apply Theorem 5, we need to verify that $f_{n,\theta}$ satisfies Assumption **M**(i), (ii) and (iii'). Here, we focus on showing **M**(iii') while other details are found in Section B.8 of the supplementary material [[Seo and Otsu \(2018\)](#)]. Observe that

$$\begin{aligned} & P \sup_{\theta \in \Theta: |\theta - \vartheta| < \varepsilon} h_n^2 |f_{n,\theta} - f_{n,\vartheta}|^2 \\ & \leq P \sup_{\theta \in \Theta: |\theta - \vartheta| < \varepsilon} \mathbb{I}\{|u - x'\vartheta| \leq h_n|x|, |u - x'\theta| > h_n|x|\} \\ & \quad + P \sup_{\theta \in \Theta: |\theta - \vartheta| < \varepsilon} \mathbb{I}\{|u - x'\theta| \leq h_n|x|, |u - x'\vartheta| > h_n|x|\}, \end{aligned}$$

for all ϑ in a neighborhood of 0. Since the same argument applies to the second term, we focus on the first term (say, T). If $\varepsilon \leq 2h_n$, an expansion around $\varepsilon = 0$ implies

$$T \leq P\{(h_n - \varepsilon)|x| \leq u \leq h_n|x|\} = P\gamma(h_n|x|)|x|\varepsilon + o(\varepsilon),$$

assuming independence between u and x . Also, if $\varepsilon > 2h_n$, an expansion around $h_n = 0$ implies

$$T \leq P\{-h_n|x| \leq u \leq h_n|x|\} \leq P\gamma(0)|x|\varepsilon + o(h_n).$$

Therefore, Assumption **M**(iii') is satisfied.

Finally, the covariance kernel is obtained in a similar way as in Section 3.4. Let $r_n = (nh_n^2)^{1/3}$. The covariance kernel is written by $H(s_1, s_2) = \frac{1}{2}\{L(s_1, 0) + L(0, s_2) - L(s_1, s_2)\}$, where $L(s_1, s_2) = \lim_{n \rightarrow \infty} \text{Var}(r_n^2 \mathbb{P}_n g_{n,t})$ with $g_{n,t} = f_{n,s_1/r_n} - f_{n,s_2/r_n}$. An expansion implies $n^{-1} \text{Var}(r_n^2 g_{n,t}) \rightarrow 2\gamma(0)P|x'(s_1 - s_2)|$, where γ is the density of u . We can also see that the covariance term is negligible. Therefore, by Theorem 5, the limiting distribution of the Hough transform estimator with diminishing bandwidth is obtained as in (23) with $V = \ddot{\gamma}(0)P(|x|xx')$ and $H(s_1, s_2) = 2\gamma(0)P|x'(s_1 - s_2)|$.

5. Conclusion. This paper develops general asymptotic theory, which encompasses a wide class of nonregular M-estimation problems. Many of these problems have been left without a proper inference method for a long time. It is worthwhile to emphasize that our theory validates inference based on subsampling for this important class of estimators, including construction of confidence sets for set-valued

parameters in Manski and Tamer's (2002) binary choice model with an interval regressor. An interesting line of future research is to develop valid bootstrap methods for these estimators. Naive applications of standard bootstrap resampling lead to inconsistent inference as shown by Abrevaya and Huang (2005) and Sen, Banerjee and Woodroffe (2010) among others.

Acknowledgments. The authors would like to thank Marine Carrasco, Aureo de Paula, Javier Hidalgo, Dennis Kristensen, Benedikt Pötscher, Peter Robinson, Kyungchul Song, Yoon-Jae Whang and seminar and conference participants at Cambridge, CIREQ in Montreal, CORE in Louvain, CREATES in Aarhus, IHS in Vienna, ISNPS in Cádiz, LSE, Surrey, UCL, Vienna and York for helpful comments. The authors also acknowledge helpful comments from an Associate Editor and anonymous referees.

SUPPLEMENTARY MATERIAL

Supplement to “Local M-estimation with discontinuous criterion for dependent and incomplete observation” (DOI: [10.1214/17-AOS1552SUPP](https://doi.org/10.1214/17-AOS1552SUPP); .pdf). The supplement contains all the proofs of the theorems and lemmas, details for illustrations and additional examples.

REFERENCES

- ABREVAYA, J. and HUANG, J. (2005). On the bootstrap of the maximum score estimator. *Econometrica* **73** 1175–1204. [MR2149245](#)
- ADAMCZAK, R. (2008). A tail inequality for suprema of unbounded empirical processes with applications to Markov chains. *Electron. J. Probab.* **13** 1000–1034. [MR2424985](#)
- ANDREWS, D. F., BICKEL, P. J., HAMPEL, F. R., HUBER, P. J., ROGERS, W. H. and TUKEY, J. W. (1972). *Robust Estimates of Location: Survey and Advances*. Princeton Univ. Press, Princeton, NJ. [MR0331595](#)
- ANEVSKI, D. and HÖSSJER, O. (2006). A general asymptotic scheme for inference under order restrictions. *Ann. Statist.* **34** 1874–1930. [MR2283721](#)
- BANERJEE, M. and MCKEAGUE, I. W. (2007). Confidence sets for split points in decision trees. *Ann. Statist.* **35** 543–574. [MR2336859](#)
- BARAUD, Y. (2010). A Bernstein-type inequality for suprema of random processes with applications to model selection in non-Gaussian regression. *Bernoulli* **16** 1064–1085. [MR2759169](#)
- BELLONI, A., CHEN, D., CHERNOZHUKOV, V. and HANSEN, C. (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* **80** 2369–2429. [MR3001131](#)
- BÜHLMANN, P. and YU, B. (2002). Analyzing bagging. *Ann. Statist.* **30** 927–961. [MR1926165](#)
- CARRASCO, M. and CHEN, X. (2002). Mixing and moment properties of various GARCH and stochastic volatility models. *Econometric Theory* **18** 17–39. [MR1885348](#)
- CHAN, K. S. (1993). Consistency and limiting distribution of the least squares estimator of a threshold autoregressive model. *Ann. Statist.* **21** 520–533. [MR1212191](#)
- CHEN, X. (2007). Chapter 76. Large sample sieve estimation of semi-nonparametric models. In *Handbook of Econometrics* **6B** 5549–5632. Elsevier, Amsterdam.
- CHEN, X., HANSEN, L. P. and CARRASCO, M. (2010). Nonlinearity and temporal dependence. *J. Econometrics* **155** 155–169. [MR2607192](#)

- CHERNOFF, H. (1964). Estimation of the mode. *Ann. Inst. Statist. Math.* **16** 31–41. [MR0172382](#)
- CHERNOZHUKOV, V., HONG, H. and TAMER, E. (2007). Estimation and confidence regions for parameter sets in econometric models. *Econometrica* **75** 1243–1284. [MR2347346](#)
- DE JONG, R. M. and WOUTERSEN, T. (2011). Dynamic time series binary choice. *Econometric Theory* **27** 673–702. [MR2822362](#)
- DOUKHAN, P., MASSART, P. and RIO, E. (1995). Invariance principles for absolutely regular empirical processes. *Ann. Inst. Henri Poincaré Probab. Stat.* **31** 393–427. [MR1324814](#)
- GAUTIER, E. and KITAMURA, Y. (2013). Nonparametric estimation in random coefficients binary choice models. *Econometrica* **81** 581–607. [MR3043343](#)
- GOLDENSHLUGER, A. and ZEEVI, A. (2004). The Hough transform estimator. *Ann. Statist.* **32** 1908–1932. [MR2102497](#)
- HIRANO, K. and PORTER, J. R. (2003). Asymptotic efficiency in parametric structural models with parameter-dependent support. *Econometrica* **71** 1307–1338. [MR2000249](#)
- HONORÉ, B. E. and KYRIAZIDOU, E. (2000). Panel data discrete choice models with lagged dependent variables. *Econometrica* **68** 839–874. [MR1771585](#)
- KIM, J. and POLLARD, D. (1990). Cube root asymptotics. *Ann. Statist.* **18** 191–219. [MR1041391](#)
- KOO, B. and SEO, M. H. (2015). Structural-break models under mis-specification: Implications for forecasting. *J. Econometrics* **188** 166–181. [MR3371666](#)
- KOSOROK, M. R. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. Springer, New York. [MR2724368](#)
- LEE, M.-J. (1989). Mode regression. *J. Econometrics* **42** 337–349. [MR1040748](#)
- MANSKI, C. F. (1975). Maximum score estimation of the stochastic utility model of choice. *J. Econometrics* **3** 205–228. [MR0436905](#)
- MANSKI, C. F. and TAMER, E. (2002). Inference on regressions with interval data on a regressor or outcome. *Econometrica* **70** 519–546. [MR1913822](#)
- MERLEVÈDE, F., PELIGRAD, M. and RIO, E. (2009). Bernstein inequality and moderate deviations under strong mixing conditions. In *High Dimensional Probability V: The Luminy Volume*. *Inst. Math. Stat. (IMS) Collect.* **5** 273–292. IMS, Beachwood, OH. [MR2797953](#)
- MERLEVÈDE, F., PELIGRAD, M. and RIO, E. (2011). A Bernstein type inequality and moderate deviations for weakly dependent sequences. *Probab. Theory Related Fields* **151** 435–474. [MR2851689](#)
- NICKL, R. and SÖHL, J. (2016). Nonparametric Bayesian posterior contraction rates for discretely observed scalar diffusions. Working paper, [arXiv:1510.05526v2](#).
- PAULIN, D. (2015). Concentration inequalities for Markov chains by Marton couplings and spectral methods. *Electron. J. Probab.* **20** 1–32. [MR3383563](#)
- POLITIS, D. N., ROMANO, J. P. and WOLF, M. (1999). *Subsampling*. Springer, New York. [MR1707286](#)
- POLLARD, D. (1989). Asymptotics via empirical processes. *Statist. Sci.* **4** 341–366. [MR1041762](#)
- POLONIK, W. and YAO, Q. (2000). Conditional minimum volume predictive regions for stochastic processes. *J. Amer. Statist. Assoc.* **95** 509–519. [MR1803169](#)
- PRAKASA RAO, B. L. S. (1969). Estimation of a unimodal density. *Sankhyā Ser. A* **31** 23–36. [MR0267677](#)
- RIO, E. (1997). About the Lindeberg method for strongly mixing sequences. *ESAIM Probab. Stat.* **1** 35–61. [MR1382517](#)
- ROMANO, J. P. and SHAIKH, A. M. (2008). Inference for identifiable parameters in partially identified econometric models. *J. Statist. Plann. Inference* **138** 2786–2807. [MR2422399](#)
- ROUSSEEUW, P. J. (1984). Least median of squares regression. *J. Amer. Statist. Assoc.* **79** 871–880. [MR0770281](#)
- SEN, B., BANERJEE, M. and WOODROOFE, M. (2010). Inconsistency of bootstrap: The Grenander estimator. *Ann. Statist.* **38** 1953–1977. [MR2676880](#)

- SEO, M. H and OTSU, T. (2018). Supplement to “Local M-estimation with discontinuous criterion for dependent and limited observations.” DOI:[10.1214/17-AOS1552SUPP](https://doi.org/10.1214/17-AOS1552SUPP).
- TALAGRAND, M. (2005). *The Generic Chaining*. Springer, Berlin. [MR2133757](#)
- VAN DER VAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer, New York. [MR1385671](#)
- VAN DER VAART, A. W. and WELLNER, J. A. (2007). Empirical processes indexed by estimated functions. In *Asymptotics: Particles, Processes and Inverse Problems. Institute of Mathematical Statistics Lecture Notes—Monograph Series* **55** 234–252. IMS, Beachwood, OH. [MR2459942](#)
- YAO, W., LINDSAY, B. G. and LI, R. (2012). Local modal regression. *J. Nonparametr. Stat.* **24** 647–663. [MR2968894](#)
- ZINDE-WALSH, V. (2002). Asymptotic theory for some high breakdown point estimators. *Econometric Theory* **18** 1172–1196. [MR1926018](#)

DEPARTMENT OF ECONOMICS
SEOUL NATIONAL UNIVERSITY
KWAN-AK RO 1
KWAN-AK GU
SEOUL
KOREA
E-MAIL: myunghseo@snu.ac.kr

DEPARTMENT OF ECONOMICS
LONDON SCHOOL OF ECONOMICS
HOUGHTON STREET
LONDON, WC2A 2AE
UNITED KINGDOM
E-MAIL: t.otsu@lse.ac.uk