

Extremum Estimation and Numerical Derivatives

Han Hong^a, Aprajit Mahajan^b and Denis Nekipelov^{c*}

^{a,b} Department of Economics, Stanford University, Stanford, CA 94305.

^c Department of Economics, University of California at Berkeley, Berkeley, California, 94720

This Version: September 2010

Abstract

Many empirical researchers rely on the use of finite-difference approximation to evaluate derivatives of estimated functions. For instance commonly used optimization routines implicitly use finite-difference formulas for the gradients, which require the choice of step size parameters. This paper investigates the statistical properties of numerically evaluated gradients and of extremum estimators computed using numerical gradients. We find that first, one needs to adjust the step size or the tolerance as a function of the sample size. Second, higher-order finite difference formulas reduce the asymptotic bias similar to higher order kernels. Third, we provide weak sufficient conditions for uniform consistency of the finite-difference approximations for gradients and directional derivatives. Fourth, we analyze the numerical gradient-based extremum estimators and find that the asymptotic distribution of the resulting estimators can sometimes depend on the sequence of step sizes. We state conditions under which the numerical derivative estimator is consistent and asymptotically normal. Fifth, we generalize our results to semiparametric estimation problems. Finally, we show that the theory is also useful in a range of nonstandard estimation procedures.

JEL Classification: C14; C52

Keywords: Numerical derivative, entropy condition, stochastic equicontinuity

1 Introduction

The implementation of extremum estimators often involves the use of computational maximization routines. When the analytical gradient of the objective function is not available these routines use finite-difference approximations to the gradient which involve the use of step size parameters. The statistical noise in this approximation algorithm of the optimization routine is typically ignored in empirical work. In this paper we provide weak conditions for the consistency of numerical derivative estimates, and demonstrate that the use of finite approximation can affect both the rate of convergence and the asymptotic distribution of the resulting estimator. This result has important

*Corresponding Author. Email: nekipelov@econ.berkeley.edu. We would like to thank Tim Armstrong for insightful comments and crucial inputs in obtaining the most general version of the main theorem of the paper.

implications for the practical use of numerical optimization routines. First, the choice of numerical tolerance and the step size depend on the sample size. Second, the asymptotic distribution and, consequently, the shape of the confidence region depends on the particular sequence of step sizes. We focus on numerical gradient-based optimization routines that use finite-difference formulas to approximate the gradient of the objective function. We consider a general framework where the objective function is computed from an i.i.d. sample and can depend on finite or infinite-dimensional unknown parameters.

The importance of numerical derivatives has received some attention in the literature. Pakes and Pollard (1989), Newey and McFadden (1994) and Murphy and Van der Vaart (2000) provided sufficient conditions for using numerical derivatives to consistently estimate the asymptotic variance in a parametric model. The properties of numerical derivatives have predominantly been investigated for very smooth models. For instance, Anderssen and Bloomfield (1974) analyzed derivative computations for functions that are approximated using polynomial interpolation. L'Ecuyer and Perron (1994) considered asymptotic properties of numerical derivatives for the class of general smooth regression models. Andrews (1997) has considered the relationship between numerical tolerance for the computation of GMM-type objective functions and their sample variance. However, to the best of our knowledge there have been no studies of the impact of the numerical optimization on the statistical properties of extremum estimators.

Our results include fairly weak sufficient conditions for consistency, rates of convergence and the asymptotic distribution for several classes of numerically computed extremum estimators. Our analysis applies to M-estimators, generalized method of moment (GMM) estimators, and estimators that maximize a function involving second-order U-statistics. Numerical M-estimation is considered both for finite-dimensional and infinite-dimensional unknown parameters. We find that the choice of the step size for consistency and convergence to the asymptotic distribution depends on the interplay between the smoothness of the population objective function, the order of chosen approximation, and on the properties of the sample objective function. Specifically, we find that if the sample objective function is very smooth, then the step size for numerical differentiation can be chosen to approach zero at an arbitrarily fast rate. For a discontinuous objective function, the step size should not converge to zero too rapid as the sample increases.

We illustrate our findings with several empirical examples. In one example, we apply numerical gradient-based optimization to the maximum score (Manski (1975)), and find that for an appropriate step size sequence, the behavior of the resulting estimator is similar to the smoothed maximum score estimator of Horowitz (1992).

Section 2 analyzes uniformly consistent estimation of numerical derivatives for both parametric and semiparametric models. Section 3 and 4 study the impact of numerical derivative based optimization method on the asymptotic properties of the resulting extreme estimators. Section 5 extends these results to U-statistics based objective functions, and section 6 considers applications. Section 8 presents Monte Carlo simulation evidence. Finally section 9 concludes.

2 Estimation of derivatives from non-smooth sample functions

2.1 Derivatives of semiparametric moment functions

We cast our analysis in the context of a general conditional moment model of the form of

$$m(z, \theta, \eta(\cdot)) = E[\rho(Y, \theta, \eta(\cdot)) | Z = z] = 0, \quad \text{if and only if } (\theta, \eta(\cdot)) = (\theta_0, \eta_0(\cdot)).$$

The parameters include the finite dimensional $\theta \in \Theta \subset \mathbb{R}^d$ and the infinite dimensional $\eta(\cdot) \in \mathcal{H}$ parameters. This setup includes the unconditional moment as a special case when z is a constant. Because the moment condition $m(\cdot)$ can be multi-dimensional, this setup also includes two step and multi-step step estimators, when some of the moment conditions corresponding to initial stage estimators only depend on the infinite dimensional functions $\eta(\cdot)$. Semiparametric estimators for this general model and their asymptotic distributions are studied extensively in the literature. In some models, the moment conditions $\rho(y, \theta, \eta(\cdot))$ depend only on the value of the function $\eta(\cdot)$ evaluated at the argument y . In some other models, such as in dynamic discrete choice models and dynamic games, $\rho(y, \theta, \eta(\cdot))$ may depend on the entire function of $\eta(\cdot)$ in complex ways.

The sieve approach, studied in a sequence of papers by Newey and Powell (2003), Chen and Shen (1998), Ai and Chen (2003) and Chen and Pouzo (2009), approximates class of infinite dimensional functions \mathcal{H} using a parametric family of function \mathcal{H}_n whose dimension increases to infinity as the sample size n increases.

For any $w \in \mathcal{H}$ and $\alpha = (\theta, \eta)$, denote by $\frac{\partial m(Z, \alpha)}{\partial \eta}[w] = \left. \frac{dm(Z, \theta, \eta + \tau w)}{d\tau} \right|_{\tau=0}$ the directional derivative of $m(Z, \alpha)$ with respect to the η component in the w direction. The literature has demonstrate that $\hat{\theta}$ can be \sqrt{n} consistent and asymptotically normal while $\hat{\eta}$ can obtain the optimal nonparametric convergence rate for η , and has shown that consistent inference of θ depends on the ability to estimate the finite dimensional and infinite dimensional directional derivatives $D_{w_j}(z) \equiv \frac{\partial m(Z, \alpha)}{\partial \theta_j} - \frac{\partial m(Z, \alpha)}{\partial \eta}[w_j]$ uniformly consistently in various directions w_j , where $\alpha = (\theta, \eta)$. Akerberg, Chen, and Hahn (2009) further shows that treating the entire estimation procedure for α as parametric and reading off the variance of $\hat{\theta}$ from the upper-left block of an estimate of the asymptotic variance-covariance matrix of $\hat{\alpha} = (\hat{\theta}, \hat{\eta})$ will give consistent estimates of the asymptotic variance of the parametric component. However, in many practical estimation problems, the derivatives of $\frac{\partial \hat{m}(Z, \hat{\alpha})}{\partial \theta_j}$ and $\frac{\partial \hat{m}(Z, \hat{\alpha})}{\partial \eta}[w_j]$ do not have analytic solutions and have to be evaluated numerically. This might be the case even if $\rho(\cdot)$ appears to be linear, e.g. when $\rho(x; \theta_0, \eta_0(\cdot)) = \eta_0(z) - \alpha \eta_0(x) - f(x, z; \theta)$ for a known parametric function $f(x, z; \theta)$. The goal of this paper is to analyze the impact of numerical approximation on statistical properties of the estimator for $D_w(z)$ and the parameter of interest.

2.2 Numerical differentiation using finite differences

Finite difference methods (e.g. Judd (1998)) are often used for numerical approximation of derivatives. To illustrate, for a univariate function $g(x)$, we can use a step size ϵ to construct a one-sided

derivative estimate $\hat{g}'(x) = \frac{g(x+\epsilon) - g(x)}{\epsilon}$, or a two-sided derivative estimate $\hat{g}'(x) = \frac{g(x+\epsilon) - g(x-\epsilon)}{2\epsilon}$. More generally, the j th derivative of $g(x)$ can be estimated by a linear operator, denoted by $L_{k,p}^\epsilon g(\theta)$, that makes use of a p th order two-sided formula:

$$L_{k,p}^\epsilon g(x) = \frac{1}{\epsilon^j} \sum_{l=-p}^p c_l g(x + l\epsilon).$$

The usual two sided derivative refers to the case when $p = 1$. When $p \geq 1$, these are called higher order differentiation. For a given p , when the weights $c_l, l = 1, \dots, p$ are chosen appropriately, the error in approximating $g^{(k)}(x)$ with $L_{j,p}^\epsilon g(x)$ will be small:

$$L_{k,p}^\epsilon g(x) - g^{(k)}(x) = O(\epsilon^{2p+1-k}).$$

For $r = 2p + 1$, consider the following Taylor expansion:

$$L_{k,p}^\epsilon g(x) = \frac{1}{\epsilon^k} \sum_{l=-p}^p c_l \left[\sum_{i=0}^r \frac{g^{(i)}(x)}{i!} (l\epsilon)^i + O(\epsilon^{r+1}) \right] = \sum_{i=0}^r g^{(i)}(x) \frac{\epsilon^i}{\epsilon^j} \sum_{l=-p}^p \frac{c_l l^i}{i!} + O(\epsilon^{r+1-k}).$$

The coefficients c_l are therefore determined by a system of equations where $\delta_{i,k}$ is the Kronecker symbol that equals 1 if and only if $i = k$ and equals zero otherwise:

$$\sum_{l=-p}^p c_l l^i = i! \delta_{i,k}, \quad \text{for } i = 0, \dots, r.$$

We are mostly concerned with first derivatives where $k = 1$. The usual two sided formula corresponds to $p = 1$, $c_{-1} = -1/2$, $c_0 = 0$ and $c_1 = 1/2$. For second order first derivatives where $p = 2$ and $j = 1$, $c_1 = 1/12$, $c_{-1} = -1/12$, $c_2 = -2/3$, $c_{-2} = +2/3$, $c_0 = 0$. In addition to central numerical derivative, left and right numerical derivatives can also be defined analogously. Since they generally have larger approximation errors than central numerical derivatives, we will restrict most attention to central derivatives. In multivariate functions, the notation of p th order central derivatives can also be extended straightforwardly to partial derivatives. Since we are only concerned with $k = 1$, we only need to use $L_{1,p}^{\epsilon, x_j}$ to highlight the element of x for which the linear operator applies to. Also in multivariate functions, $g(x + \epsilon)$ means the vector of $[g(x + \epsilon e_k)]$, $k = 1, \dots, d$, where e_k is the vector with 1 in the k th position and 0 elsewhere, and d is the dimension of x .

Each j th component of $D_w(z)'$ in the asymptotic variance formula can then be estimated by

$$\tilde{D}_w^j(z) = L_{1,p}^{\epsilon_n, \theta_j} \hat{m}(z; \hat{\theta}, \hat{\eta}(\cdot)) - L_{1,p}^{\tau_n, w_j} \hat{m}(z; \hat{\theta}, \hat{\eta}(\cdot)),$$

where ϵ_n and τ_n are the relevant step sizes for the numerical derivatives with respect to the finite and infinite-dimensional parameters. In general the step sizes ϵ_n and τ_n can be chosen differently for different elements of the parametric and nonparametric components. It might also be possible to adapt the equal distance grid to a variable distance grid of the form $L_{k,p}^\epsilon g(x) = \frac{1}{\epsilon^j} \sum_{l=-p}^p c_l g(x + t_l \epsilon)$, where t_l can be different from 1. In addition both the step size and the grid distance can also be

made to be dependent on the observations and data driven. These possibilities are left for future research.

For most of the statistical analysis in the rest of the paper we assume away machine imprecision. Machine precisions also impose a lower bound on the step size in conjunction with the statistical lower bound (see, e.g. Press, Teukolsky, Vetterling, and Flannery (1992)). This and related implementation issues are discussed in section 7.

2.3 Sufficient conditions for consistency of finite-difference derivatives

Before we strive to obtain the weakest possible sufficient condition for consistency in next section, we first show that the existing sufficient conditions in the literature (e.g. Newey and McFadden (1994) and Powell (1984)) for parametric models can be straightforwardly generalized to semiparametric models.

ASSUMPTION 1. *For a linear operator $\Delta_{p,\theta^{p_1},\delta^{p_2}}[\delta]^{p_1}$ that is p^1 th linear in θ , that has a finite second moment and that is linear in each argument, e.g., $\Delta_{2,\theta,\eta}[t\delta](\theta - \theta_0) = t\Delta_{2,\theta,\eta}[\delta](\theta - \theta_0)$, the following approximation holds at (θ_0, η_0) :*

$$E \left[\left\| m(z; \theta, \eta(\cdot)) - \Delta_{1\theta}(\theta - \theta_0) - \Delta_{1\eta}[\delta] - \dots - \sum_{p_1+p_2=p} \Delta_{p,\theta^{p_1},\delta^{p_2}}[\delta]^{p_1}(\theta - \theta_0)^{p_2} \right\|^2 \right] = o \left(\|\delta\|_{\mathbf{L}^2}^{2p} + \|\theta - \theta_0\|^{2p} \right).$$

Assumption 1 requires that the conditional moment $m(z; \theta, \eta(\cdot))$ is mean square differentiable in L^2 norm with respect to the distribution of z . The next assumption relates to the rate of convergence of the nonparametric conditional moment estimate. Define U_γ as a neighborhood of θ_0, η_0 with radius γ : $U_\gamma = \{\theta, \eta(\cdot) : \|\theta - \theta_0\| < \gamma, |\eta(\cdot) - \eta_0(\cdot)| < \gamma\}$.

ASSUMPTION 2. *For some $k \in \mathbb{N}$, $k \leq 2$, uniformly in $z \in \mathcal{Z}$, as $\gamma \rightarrow 0$,*

$$\sup_{(\theta, \eta(\cdot)) \in U_\gamma} \frac{n^{1/k} \|\hat{m}(z; \theta, \eta(\cdot)) - m(z; \theta, \eta(\cdot)) - \hat{m}(z; \theta_0, \eta_0(\cdot))\|}{1 + n^{1/k} \|\hat{m}(z; \theta, \eta(\cdot))\| + n^{1/k} \|m(z; \theta, \eta(\cdot))\|} = o_p(1).$$

For unconditional moment models, typically $k = 2$. For conditional moment models, $k \geq 2$. The particular rate will depend on the method and the choice of the tuning parameters used in the estimation procedure.

In addition, the parametric component of the model is assumed to converge at the usual \sqrt{n} rate, while the functional component is assumed to converge at a slower nonparametric rate.

ASSUMPTION 3. *For $k_1 \geq 2$, $n^{1/k_1} \|\hat{\eta}(\cdot) - \eta_0(\cdot)\| = O_p(1)$, and $n^{1/2} \|\hat{\theta} - \theta_0\| = O_p(1)$.*

THEOREM 1. *Under assumptions 1, 2 and 3, if $\varepsilon_n n^{1/\max\{k, k_1\}} \rightarrow \infty$, $\varepsilon_n \rightarrow 0$, $\tau_n n^{1/\max\{k, k_1\}} \rightarrow \infty$, $\tau_n \rightarrow 0$, then $\sup_{z \in \mathcal{Z}} |\tilde{D}_w(z) - D_w(z)| \xrightarrow{p} 0$.*

The proof of the consistency theorem follows closely the arguments in the literature. The basic idea in the consistency argument is that while the step size should converge to zero to eliminate the bias, it should converge slowly so that the noise in the parameter estimation and in estimating the moment condition should not dominate the step size. In a parametric model, both the noise in the parameter estimation and in estimating the moment condition is of the order of $1/\sqrt{n}$. Therefore as shown in Newey and McFadden (1994) and Powell (1984), sufficiency will hold if $1/\sqrt{n} \ll \epsilon_n$. The extension of this argument to the semiparametric case is straightforward. The difference is that now the converge rates for both the (infinite-dimensional) parameters and the conditional moment equation are slower than $1/\sqrt{n}$ and therefore imposes a more stringent requirement on rate at which ϵ_n is allowed to converge to zero. However, as we will see in the next section, these sufficient conditions can be substantially weakened because of a local uniformity feature of the variance of the numerical derivatives.

2.4 Weak sufficient conditions for consistency for parametric models

In this section we show that the conditions on the step size for consistent derivative estimation is much weaker than previously known in the literature. In particular, as long as a local uniformity condition holds, there is no interaction between the step size choice and the statistical uncertainty in parameter estimation. We consider the unconditional parametric and conditional semiparametric cases separately to best convey intuitions.

Consider a parametric unconditional moment model defined by the sample and population moment conditions: $\hat{g}(\theta) = \frac{1}{n} \sum_{i=1}^n g(Z_i, \theta)$ and $g(\theta) = Eg(Z_i, \theta)$ where $g(\theta) = 0$ if and only if $\theta = \theta_0$, where θ_0 lies in the interior of the parameter space Θ . The goal is to estimate $G(\theta_0) = \frac{\partial g(\theta_0)}{\partial \theta}$ using $L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta}) = \left(L_{1,p}^{\epsilon_n, \hat{\theta}_j} \hat{g}(\hat{\theta}), j = 1, \dots, d \right)$.

In the following, we decompose the error of approximating $G(\theta_0)$ with $L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta})$ into three components: $L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta}) - G(\theta_0) = \hat{G}_1(\hat{\theta}) + G_2(\hat{\theta}) + G_3(\hat{\theta})$, where

$$\hat{G}_1(\hat{\theta}) = L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta}) - L_{1,p}^{\epsilon_n} g(\hat{\theta}), \quad (2.1)$$

and

$$G_2(\hat{\theta}) = L_{1,p}^{\epsilon_n} g(\hat{\theta}) - G(\hat{\theta}), \quad G_3(\hat{\theta}) = G(\hat{\theta}) - G(\theta_0).$$

We discuss how to control each of these three terms in turn. Notice first that the step size ϵ_n does not play a role in $G_3(\hat{\theta})$. The bias term $G_2(\hat{\theta})$ can be controlled if the bias reduction is uniformly small in a neighborhood of θ_0 .

The following assumption is a parametric version of Assumption 1.

ASSUMPTION 4. *A $2p + 1$ th order mean value expansion applies to the limiting function $g(\theta)$*

uniformly in a neighborhood of θ_0 . For all sufficiently small $|\epsilon|$ and $r = 2p + 1$,

$$\sup_{\theta \in \mathcal{N}(\theta_0)} \left| g(\theta + \epsilon) - \sum_{l=0}^r \frac{\epsilon^l}{l!} g^{(l)}(\theta) \right| = O(|\epsilon|^{r+1}).$$

An immediate consequence of this assumption is that $\hat{G}_2(\hat{\theta}) = O(\epsilon^{2p})$. We are left with $\hat{G}_1(\hat{\theta})$. The weakest possible condition to control $\hat{G}_1(\hat{\theta})$ that covers all the models that we are aware of seems to come from a convergence rate result in Pollard (1984).

ASSUMPTION 5. Consider functions $g(z, \theta)$ contained in class $\mathcal{F} = \{g(\cdot, \theta), \theta \in \Theta\}$. Then

- (i) All $g \in \mathcal{F}$ are globally bounded such that $\|F\| = \sup_{\theta \in \Theta} |g(Z_i, \theta)| < C \ll \infty$.
- (ii) The sample moment function is Lipschitz-continuous in **mean square** in some neighborhood of θ_0 . That is for sufficiently small $\epsilon > 0$

$$\sup_{\theta \in \mathcal{N}(\theta_0)} E \left[(g(Z_i, \theta + \epsilon) - g(Z_i, \theta - \epsilon))^2 \right] = O(\epsilon).$$

- (iii) The graphs of functions from \mathcal{F} form a polynomial class of sets.

Most of the functions in econometric applications fall in this category. By Lemmas 25 and 36 of Pollard (1984), assumption 5 implies that there exist universal constants $A > 0$ and $V > 0$ such that for any $\mathcal{F}_n \subset \mathcal{F}$ with envelope function $\|F_n\|$,

$$\sup_{\mathcal{Q}} N_1(\epsilon \mathcal{Q} F_n, \mathcal{Q}, \mathcal{F}_n) \leq A \epsilon^{-V}, \quad \sup_{\mathcal{Q}} N_2(\epsilon (\mathcal{Q} F_n^2)^{1/2}, \mathcal{Q}, \mathcal{F}_n) \leq A \epsilon^{-V}.$$

LEMMA 1. Under assumption 5, if $n\epsilon_n / \log n \rightarrow \infty$

$$\sup_{d(\theta, \theta_0) = o(1)} \|L_{1,p}^{\epsilon_n} \hat{g}(\theta) - L_{1,p}^{\epsilon_n} g(\theta)\| = o_p(1).$$

Consequently, Assumption 5 implies that $\hat{G}_1(\hat{\theta}) = o_p(1)$ if $d(\hat{\theta}, \theta_0) = o_p(1)$.

Proof: The argument follows directly from Theorem 2.37 in Pollard (1984) by verifying its conditions. For each n and each ϵ_n , consider the class of functions $\mathcal{F}_n = \{\epsilon_n L_{1,p}^{\epsilon_n} g(\cdot, \theta), \theta \in \mathcal{N}(\theta_0)\}$, with envelope function F , such that $PF \leq C$. Then we can write

$$\sup_{d(\theta, \theta_0) \leq o(1)} \epsilon_n \|L_{1,p}^{\epsilon_n} \hat{g}(\theta) - L_{1,p}^{\epsilon_n} g(\theta)\| \leq \sup_{f \in \mathcal{F}_n} |P_n f - P f|.$$

For each $f \in \mathcal{F}_n$, note that $E f^2 = E (\epsilon_n L_{1,p}^{\epsilon_n} g(\cdot, \theta))^2 = O(\epsilon_n)$ because of assumption 5.(ii). The lemma then follows immediately by taking $\alpha_n = 1$ and $\delta_n^2 = \epsilon_n$ in Theorem 2.37 in Pollard (1984). \square

THEOREM 2. *Under Assumptions 4 and 5, $L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta}) \xrightarrow{p} G(\theta_0)$ if $\epsilon_n \rightarrow 0$ and $n\epsilon_n/\log n \rightarrow \infty$, and if $d(\hat{\theta}, \theta_0) = o_p(1)$.*

In most situations $d(\hat{\theta}, \theta_0) = O_p(n^{-\eta})$ for some $\eta > 0$. Typically $\eta = 1/2$. One might hope to further weaken the requirement of the $\log n$ term when uniformity is only confined to a shrinking neighborhood of size $n^{-\eta}$. However, this turns out not possible unless the moment function exhibits an additional level of smoothness.

The result of Theorem 2 improved if we are willing to impose the following stronger assumption, which holds for smoother functions such as those that are Hölder-continuous.

ASSUMPTION 6. *In addition to assumption 5, for all sufficiently small ϵ and all $\theta \in \Theta$, if we define $\mathbb{G}_n(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (g(Z_i, \theta) - g(\theta))$, then*

$$E^* \sup_{\theta', \theta \in N(\theta_0)} |\mathbb{G}_n(\theta') - \mathbb{G}_n(\theta)| \lesssim \phi_n(\delta),$$

for functions $\phi_n(\cdot)$ such that $\delta \mapsto \phi_n(\delta)/\delta^\gamma$ is non-increasing for γ defined in part (i).

Assumption 6 is more stringent than Theorem 3.2.5 in Van der Vaart and Wellner (1996), and may fail in cases where Theorem 3.2.5 holds, for example with indicator functions as in assumption 5. Theorem 3.2.5 only requires that $E^* \sup_{d(\theta, \theta_0) < \delta} |\mathbb{G}(\theta) - \mathbb{G}(\theta_0)| \lesssim \phi_n(\delta)$. For i.i.d data, the tail bounds method used in Van der Vaart and Wellner (1996) can be modified to obtain assumption 6. In particular, define a class of functions $\mathcal{M}_\delta^\epsilon = \{g(Z_i, \theta_1) - g(Z_i, \theta_2), d(\theta_1, \theta_2) \leq \delta, d(\theta_1, \theta_0) < \epsilon, d(\theta_2, \theta_0) < \epsilon\}$. Then assumption 6, which requires bounding $E_P^* \|G_n\|_{\mathcal{M}_\delta^\epsilon}$, can be obtained by invoking the maximum inequalities in Theorems 2.14.1 and 2.14.2 in Van der Vaart and Wellner (1996). These inequalities provide that for M_δ^ϵ an envelope function of the class of functions $\mathcal{M}_\delta^\epsilon$,

$$\begin{aligned} E_P^* \|G_n\|_{\mathcal{M}_\delta^\epsilon} &\lesssim J(1, \mathcal{M}_\delta^\epsilon) \left(P^*(M_\delta^\epsilon)^2 \right)^{1/2}, \\ E_P^* \|G_n\|_{\mathcal{M}_\delta^\epsilon} &\lesssim J_{[]} (1, \mathcal{M}_\delta^\epsilon, L_2(P)) \left(P^*(M_\delta^\epsilon)^2 \right)^{1/2}, \end{aligned}$$

where $J(1, \mathcal{M}_\delta^\epsilon)$ and $J_{[]} (1, \mathcal{M}_\delta^\epsilon, L_2(P))$ are the uniform and bracketing entropy integrals defined in section 2.14.1 of Van der Vaart and Wellner (1996), and are generically finite for parametric functions. Therefore $\phi_n(\delta)$ depends mostly on the variance of the envelope functions $\left(P^*(M_\delta^\epsilon)^2 \right)^{1/2}$. For reasonably smooth functions that are Hölder-continuous, M_δ^ϵ depends only on δ as required by assumption 6.

THEOREM 3. *Under assumptions 4 and 6, $L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta}) \xrightarrow{p} G(\theta_0)$ if $\epsilon_n \rightarrow 0$ and $n\epsilon_n^{2-2\gamma} \rightarrow \infty$, and if $d(\hat{\theta}, \theta_0) = o_p(1)$.*

This result, which is an immediate consequence of Theorem 2.14.1 of Van der Vaart and Wellner (1996) and therefore stated without proof, shows that for continuous functions $g(Z_i, \theta)$ that are

Lipschitz in θ , the only condition needed for consistency is $\epsilon_n \rightarrow 0$. The result of Theorem 3 demonstrates that as long as the sample moment function does not have discontinuities, one can pick the step size to decrease at the polynomial rate with the sample size. If the moment function is discontinuous, Theorem 2 needs to be applied instead of Theorem 3 prescribing a slower logarithmic rate of decrease in the step size.

Example Consider the simple quantile case where the moment condition is defined by $g(z_i; \theta) = 1(z_i \leq \theta) - \tau$. In this case the numerical derivative estimate of the density of z_i at θ is given by

$$L_{1,2}^{\epsilon_n} \hat{g}(\hat{\theta}) = \frac{1}{n} \sum_{i=1}^n \frac{1(z_i \leq \hat{\theta} + \epsilon) - 1(z_i \leq \hat{\theta} - \epsilon)}{2\epsilon}.$$

This is basically the uniform kernel estimate of the density of z_i at θ :

$$\hat{f}(\hat{\theta} - \theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2\epsilon} 1\left(\frac{|z_i - \hat{\theta}|}{\epsilon} \leq 1\right) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2\epsilon} 1\left(\frac{|z_i - \theta_0 - (\hat{\theta} - \theta_0)|}{\epsilon} \leq 1\right).$$

The consistency conditions given in Powell (1984) and Newey and McFadden (1994), both of which require $\sqrt{n}\epsilon \rightarrow \infty$, are too strong. The intuition reason for this is because under this condition, the second part of the estimation noise due to $\hat{\theta} - \theta_0$, $\frac{\hat{\theta} - \theta_0}{\epsilon}$, will vanish. However, for the purpose of consistency this is not necessary. As long as $\hat{f}(x)$ is uniformly consistent for $f(x)$ for x in a shrinking neighborhood of 0 of size $n^{-\eta}$, it will follow that

$$\hat{f}(\hat{\theta} - \theta_0) \xrightarrow{p} f(0) = f_z(\theta_0).$$

2.5 Optimal rates for derivative estimation

The optimal choice of the step size depends typically on the smoothness of the empirical process indexed by γ and the smoothness of the population moment function indicated by the magnitude of the ‘‘Taylor residual’’ p . For the choice of the optimal rate for the step size of numerical differentiation, one can consider decomposition of the numerical derivatives into components \hat{G}_1 , G_2 and G_3 corresponding to the variance, deterministic and stochastic bias components. The optimal choice of the step size will provide the minimum mean-squared error for the estimated derivative by balancing the bias and the variance. When the sample moment function is discontinuous, conditions of Theorem 2 apply, delivering the logarithmic rate of decay for the variance component $\hat{G}_1(\hat{\theta}) = O_p\left(\sqrt{\frac{\log n}{n\epsilon}}\right)$. On the other hand, application of Assumption 4 to the population moment leads to $\hat{G}_2(\hat{\theta}) = O(\epsilon^{2p})$. Under conditions of Theorem 3, the variance term has a polynomial dependence from the step size with $\hat{G}_1(\hat{\theta}) = O_p\left(\frac{1}{\sqrt{n\epsilon^{1-\gamma}}}\right)$, while the bias term is still determined by Assumption 4. We note that for Lipschitz-continuous or differentiable models, in which generally $\gamma = 1$, there is no trade off between the variance and the bias, in which case the smaller the step size ϵ , the smaller the bias term. However, in this case the order of the root mean square is bounded from below by the variance term of $O(1/\sqrt{n})$ for sufficiently smaller ϵ_n . The next theorem formalizes.

THEOREM 4. *Under the conditions of Theorem 2, if $\hat{\theta} - \theta_0 = O_p(1/\sqrt{n})$, the optimal rate of ϵ satisfies $\epsilon = O\left((\log n/n)^{\frac{1}{4p+1}}\right)$, in which case the mean-squared error is $O_p\left((\log n/n)^{\frac{4p}{4p+1}}\right)$. When the conditions of theorem 3 hold instead, the optimal rate of ϵ is $O\left(n^{-\frac{1}{2(1-\gamma+2p)}}\right)$ if $\gamma < 1$, and $\epsilon \ll n^{-1/4p}$ if $\gamma = 1$. In both cases the error is $O_p\left(n^{-\frac{2p}{2(1-\gamma+2p)}}\right)$.*

2.6 Uniform consistency of directional derivatives for semiparametric models

This subsection extends the weak consistency condition to directional derivatives of semiparametric conditional moment models. As in Section 2.1, semiparametric conditional moment models are usually defined by conditional moment function $m(\theta, \eta; z) = E[\rho(Y_i, \theta, \eta) | Z_i = z]$. In this section we focus on two special cases where the conditional moment function is estimated nonparametrically using orthogonal series and when it is estimated using kernel smoothing. The infinite-dimensional parameter η is assumed to be estimated using sieves. The series estimator used to recover the conditional moment function is based on the vector of basis functions $p^N(z) = (p_{1N}(z), \dots, p_{NN}(z))'$,

$$\hat{m}(\theta, \eta, z) = p^{N'}(z) \left(\frac{1}{n} \sum_{i=1}^n p^N(z_i) p^{N'}(z_i) \right)^{-1} \frac{1}{n} \sum_{i=1}^n p^N(z_i) \rho(\theta, \eta; y_i). \quad (2.2)$$

The kernel estimator is defined using a multi-dimensional kernel function $K(\cdot)$ and a bandwidth sequence b_n as

$$\hat{m}(\theta, \eta, z) = \left(\frac{1}{nb_n^{d_z}} \sum_{i=1}^n K\left(\frac{z_i - z}{b_n^{d_z}}\right) \right)^{-1} \frac{1}{nb_n^{d_z}} \sum_{i=1}^n K\left(\frac{z_i - z}{b_n^{d_z}}\right) \rho(\theta, \eta; y_i). \quad (2.3)$$

In either case, we will denote the resulting estimate by $\hat{m}(\theta, \hat{\eta}; x)$. It turns out that the numerical derivative consistency results for η apply without any modification to the parametric component θ . Therefore with no loss of generality below we will focus on differentiating with respect to η .

The directional derivative of m in the direction $w \in \mathcal{H} - \eta_0$ with respect to η , $G_w = \left. \frac{d m(\theta_0, \eta_0 + \tau w, z)}{d\tau} \right|_{\tau=0}$, is estimated using $L_{1,p}^{\epsilon_n, w} \hat{m}(\hat{\theta}, \hat{\eta}, z)$, where an additional index is used to emphasize the direction for which the derivative is taken,

$$L_{1,p}^{\epsilon_n, w} \hat{m}(\hat{\theta}, \hat{\eta}, z) = \frac{1}{\epsilon_n} \sum_{l=-p}^p c_l \hat{m}(\hat{\theta}, \hat{\eta} + l w \epsilon_n, z).$$

Given that the direction w itself has to be estimated from the data as in section 2.1, we desire consistency results that hold uniformly both around the true parameter value and the directions of numerical differentiation. As in our analysis of parametric models, we focus on i.i.d data samples. We also impose standard assumptions on the basis functions as in Newey (1997).

ASSUMPTION 7. *For the basis functions $p^N(z)$ the following holds:*

- (i) The smallest eigenvalue of $E[p^N(Z_i)p^{N'}(Z_i)]$ is bounded away from zero uniformly in N^2
- (ii) For some $C > 0$, $\sup_{z \in \mathcal{Z}} \|p^N(z)\| \leq C < \infty$.
- (iii) The population conditional moment belongs to the completion of the sieve space and

$$\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}} \sup_{z \in \mathcal{Z}} \|m(\theta, \eta, z) - \text{proj}(m(\theta, \eta, z) | p^N(z))\| = O(N^{-\alpha}).$$

Assumption 7[ii] is convenient because $\rho(\cdot)$ is uniformly bounded. It can potentially be relaxed to allow for a sequence of constants $\zeta_0(N)$ with $\sup_{z \in \mathcal{Z}} \|p^N(z)\| \leq \zeta_0(N)$, where $\zeta_0(N)$ grows at appropriate rates as in Newey (1997) such as $\zeta_0(N)^2 N/n \rightarrow 0$ as $n \rightarrow \infty$.

The following assumption on the moment function $\rho(\cdot)$ will not require smoothness or continuity, and is related to Shen and Wong (1994) and Zhang and Gijbels (2003).

ASSUMPTION 8. (i) *Uniformly bounded moment functions:* $\sup_{\theta, \eta} \|\rho(\theta, \eta, \cdot)\| \leq C$. *The density of covariates Z is uniformly bounded from zero on its support.*

- (ii) Suppose that $0 \in \mathcal{H}_n$ and for $\epsilon_n \rightarrow 0$ and some $C > 0$,

$$\sup_{\substack{z \in \mathcal{Z}, \eta, w \in \mathcal{H}_n, |\eta|, |w| < C, \\ \theta \in \mathcal{N}(\theta_0)}} \text{Var}(\rho(\theta, \eta + \epsilon_n w; Y_i) - \rho(\theta, \eta - \epsilon_n w; Y_i) | z) = O(\epsilon_n),$$

- (iii) For each n , the class of functions $\mathcal{F}_n = \{\rho(\theta, \eta + \epsilon_n w; \cdot) - \rho(\theta, \eta - \epsilon_n w; \cdot), \theta \in \Theta, \eta, w \in \mathcal{H}_n\}$ is Euclidean whose coefficients depend on the number of sieve terms. In other words, there exist constants A , and $0^+ \leq r_0 < \frac{1}{2}$ such that the covering number satisfies

$$\log N(\delta, \mathcal{F}_n, \mathbf{L}_1) \leq A n^{2r_0} \log \left(\frac{1}{\delta} \right),$$

and $r_0 = 0^+$ corresponds to the case $\log n$.

The hardest condition to verify is 8 (iii). This assumption imposes a joint restriction both on the class of functions \mathcal{H}_n containing sieve estimators for η and the class of conditional moment functions parametrized both by θ and η . An example where this assumption holds is when $\rho(\cdot)$ is (weakly) monotone in η for each θ and \mathcal{H}_n is a orthogonal basis of dimensionality $K(n)$. For example, $\rho(\cdot)$ can be an indicator in nonparametric quantile regression. Lemma 5 in Shen and Wong (1994) suggests that the L_1 -metric entropy of the class of sieve \mathcal{F}_n has order $K(n) \log \frac{1}{\epsilon} \leq K(n) \epsilon^{-1}$ for sufficiently small $\epsilon > 0$ and $\|\eta_n - \eta_0\|_{\mathbf{L}_1} < \epsilon$. Then by Lemma 2.6.18 in Van der Vaart and Wellner (1996), if function $\rho(\cdot)$ is monotone, its application to η (for fixed θ) does not increase the metric entropy. In

²We note that the considered series basis may not be orthogonal with respect to the semi-metric defined by the distribution of Z_i .

addition, the proof of Theorem 3 in Chen, Linton, and Van Keilegom (2003) shows that the metric entropy for the entire class \mathcal{F}_n is a sum of metric entropies that are obtained by fixing η and θ . The choice $K(n) \sim n^{2r_0}$ delivers condition 8 (iii).

Denote $\pi_n \eta = \arg \inf_{\eta' \in \mathcal{H}_n} \|\eta' - \eta\|$. And let $d(\cdot)$ be the metric generated by the \mathbf{L}^1 norm. The following result is formulated in the spirit of Theorem 37 of Pollard (1984) and it requires its extension to the case of sieve estimators. A related idea for unconditional sieve estimation has been used in Zhang and Gijbels (2003).

LEMMA 2. *Suppose that $\rho(\pi_n \eta, \eta) = O_p(n^{-\phi})$. Under assumptions 7 and 8*

$$\sup_{d(\theta, \theta_0)=o(1), d(\eta, \eta_0)=o(1), \eta \in \mathcal{H}_n} |L_{1,p}^{\epsilon_n, w} \widehat{m}(\theta, \eta, z) - L_{1,p}^{\epsilon_n, w} m(\theta, \eta, z)| = o_p(1)$$

uniformly in z and w , provided that $\epsilon_n \rightarrow 0$ and $\min\{N^\alpha, n^\phi\}\epsilon_n \rightarrow \infty$, and $\frac{n\epsilon_n}{N^2 n^{2r_0} \log n} \rightarrow \infty$.

We note that provided result is uniform in z . In interesting feature of the series estimator for $m(\cdot)$ is that z is the argument of $p^N(\cdot)$ only and its boundedness is sufficient for the uniform result in Lemma 2. In some cases weaker conditions may be possible provided this feature. We show one of such cases in Section 4 for the case of the density-weighted sieve minimum distance estimators. We can provide a similar result for the case where the conditional moment function is estimated via kernel estimator. We begin with formulating the requirement on the kernel.

ASSUMPTION 9. *$K(\cdot)$ is the q -th order kernel function which is an element of the class of functions \mathcal{F} defined by Assumption 5. It integrates to 1, it is bounded and its square has a finite integral.*

Then we can formulate the following lemma that replicates the result of Lemma 2 for the case of the kernel estimator. We note that for uniformity we rely on Assumption 8(i) that requires the density of covariates to be uniformly bounded from zero.

LEMMA 3. *Under assumptions 8 and 9*

$$\sup_{d(\theta, \theta_0)=o(1), d(\eta, \eta_0)=o(1), \eta \in \mathcal{H}_n} |L_{1,p}^{\epsilon_n, w} \widehat{m}(\theta, \eta, z) - L_{1,p}^{\epsilon_n, w} m(\theta, \eta, z)| = o_p(1)$$

uniformly in w and z where $f(z)$ is strictly positive for the kernel estimator provided that $\epsilon_n \rightarrow 0$, $b_n \rightarrow 0$, $\epsilon_n \min\{b_n^{-q}, n^\phi\} \rightarrow \infty$ and $\frac{n\epsilon_n b_n^{d_z}}{n^{2r_0} \log n} \rightarrow \infty$.

Using Lemmas 2 and 3 we can formulate the consistency result for the directional derivative.

THEOREM 5. *Under assumptions 4, 8, and either 7 or 9, $L_{1,p}^{\epsilon_n, w} \widehat{m}(\hat{\theta}, \hat{\eta}, z) \xrightarrow{p} \frac{\partial m(\theta, \eta, z)}{\partial \eta}[w]$, uniformly in z and w , if $N \rightarrow \infty$, $\epsilon_n \min\{N^\alpha, n^\phi\} \rightarrow \infty$, and $\frac{n\epsilon_n}{N^2 n^{2r_0} \log n} \rightarrow \infty$ for series estimator, and $b_n \rightarrow 0$, $\epsilon_n \min\{b_n^{-q}, n^\phi\} \rightarrow \infty$, and $\frac{n\epsilon_n b_n^{d_z}}{n^{2r_0} \log n} \rightarrow \infty$ for kernel-based estimator, provided that $d(\hat{\theta}, \theta_0) = o_p(1)$ and $d(\hat{\eta}, \eta_0) = o_p(1)$.*

This theorem allows us to use finite-difference formulas for evaluation of directional derivatives. An interesting feature of this result is that it only indirectly depends on the rate of convergence of the infinite-dimensional parameter through our assumption 8[iii], which implicitly bounds the number of sieve terms that one can use by n^{2r_0} with $r_0 < \frac{1}{2}$, i.e. it has to increase slower than the sample size.

Remark: Our results in this section apply to the case where one is interested in obtaining a finite-difference based estimator for the directional derivative that is uniformly consistent over z . Such a need may arise where the direction of differentiation is also estimated, an example of which is the efficient sieve minimum distance estimator in Ai and Chen (2003). If one only needs to estimate the numerical derivative pointwise the conditions on the choice of the step size can be weakened. Such results may be relevant when one is interested in estimating the directional derivative at a point and a given direction.

2.7 Analysis with Hölder-continuous moment functions

In this section we consider a special case where finite differences of the moment function have a non-trivial envelope. Examples of such functions include Lipschitz and Hölder continuous functions. We introduce the following modification to Assumption 8(i):

ASSUMPTION 8.

(i') For any sufficiently small $\epsilon > 0$

$$\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}, w \in \mathcal{H}, |w| < C} \|\rho(\theta, \eta + w\epsilon, z) - \rho(\theta, \eta, z)\| \leq C(z)\epsilon^\gamma,$$

where $0 < \gamma \leq 1$ and $E[C(Z)^2] < \infty$.

This modification allows us to improve the rate results for the considered classes of functions. Specifically, we can adapt the exponential tail bounds provided in Alexander (1984) to establish a stronger result than the result that we had for discontinuous moment functions. This result of particular interest because it also implies appropriate properties of the numerical directional derivative in cases where the moment function is at least Lipschitz continuous.

LEMMA 4. Suppose that $\rho(\pi_n \eta, \eta) = O(n^{-\phi})$. Under either pair of assumptions 7 and 8(i'), (ii), (iii), (iv) or 9 and 8(i'), (ii), (iii), (iv)

$$\sup_{d(\hat{\theta}, \theta_0) = o_p(1), d(\hat{\eta}, \eta_0) = o_p(1), \eta \in \mathcal{H}_n} \left| L_{1,p}^{\epsilon_n, w} \hat{m}(\hat{\theta}, \hat{\eta}, z) - L_{1,p}^{\epsilon_n, w} m(\theta_0, \eta_0, z) \right| = o_p(1)$$

uniformly in z and w , provided that $\epsilon_n \rightarrow 0$, $\epsilon_n \min\{N^\alpha, n^\phi\} \rightarrow \infty$ and $\frac{\sqrt{n} \epsilon_n^{1-\gamma}}{N n^{r_0}} \rightarrow \infty$ for series estimator, and $b_n \rightarrow 0$, $\epsilon_n \min\{b_n^{-q}, n^\phi\} \rightarrow \infty$, $\frac{n^{1-2r_0} \epsilon_n^{1-\gamma} b_n^{d_z/2}}{\log n^{2r_0-1}} \rightarrow \infty$ for kernel estimator.

The consistency of the numerical derivative is a direct consequence of this lemma. We note that conditions of Lemma 4 are weaker than the conditions for the functions with trivial (constant)

envelopes. We can note that for functions that are Lipschitz-continuous the power $\gamma = 1$. This means that Lemma 4 will only set the requirement for the estimator of the conditional moment, but not the step size (except for the bias-correction term requiring that the bias from projection or kernel-smoothing should decay faster than the step size for the numerical derivative). This means that the use of the finite-difference formula will not affect the properties of the estimated conditional moment. One interesting observation from this case is when the moment function becomes substantially non-smooth (i.e. when γ is close to zero), the step size for numerical differentiation should be much larger than the bandwidth if one uses a first-order kernel to compute the conditional moment function.

3 Numerical optimization of non-smooth sample functions

3.1 Parametric extremum estimation: definitions

In this section we study the properties of estimators based on numerically solving the first-order conditions for likelihood-type objective functions. The original estimator of interest maximizes the sample objective function. However, either by explicit researcher choice or by the implicit choice of the maximization routine in the optimization software, the original problem is replaced by the search for the zero of the numerically computed gradient. Consider the problem of estimating parameter θ_0 in a metric space (Θ, d) with the metric d . The true parameter θ_0 is assumed to uniquely maximize the limiting objective function $Q(\theta) = Eg(Z_i; \theta)$. An M-estimator $\hat{\theta}$ of θ_0 is typically defined as

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \hat{Q}(\theta), \quad (3.4)$$

where $\hat{Q}(\theta) = \frac{1}{n} \sum_{i=1}^n g(Z_i; \theta)$. However, in practice, most sample objective functions $\hat{Q}(\theta)$ of interest cannot be optimized analytically and are optimized instead through numerical computation. The optimization routine often uses numerical derivatives either explicitly or implicitly. In this section we show that numerical differentiation can sometimes lead to a model that is different from the one usually studied under M-estimation. In particular, while numerical differentiation does not affect the asymptotic distribution for smooth models (under suitable conditions on the step size sequence), for nonsmooth models a numerical derivative based estimator can translate a nonstandard parametric model into a nonparametric one.

We focus on the class of optimization procedures that are based on numerical gradients, that are evaluated using the finite-difference formulas which we described in Section 2.2. We start by presenting a finite-difference numerical derivative version of the, M-estimator in (3.4). A numerical gradient-based optimization routine effectively substitutes (3.4) by a solution to the non-linear equation

$$\|L_{1,p}^{\varepsilon_n} \hat{Q}_n(\hat{\theta})\| = o_p\left(\frac{1}{\sqrt{n}}\right), \quad (3.5)$$

for some sequence of step sizes $\varepsilon_n \rightarrow 0$ and $\hat{Q}_n(\theta) = \frac{1}{n} \sum_{i=1}^n g(Z_i, \theta)$. We do not require the zeros of the first order condition to be exact in order to accommodate nonsmooth models. Many popular optimization packages use $p = 1$, corresponding to $\hat{D}_n^\varepsilon(\hat{\theta}) \equiv L_{1,1}^\varepsilon \hat{Q}_n(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n}}\right)$. The cases

with $p \geq 2$ correspond to a more general class of estimators that will have smaller asymptotic bias in nonsmooth models. As we will argue, the estimators (3.4) and (3.5) can have the same properties for models with continuous moment functions but for non-smooth models both their asymptotic distributions and the convergence rates can be substantially different.

3.2 Consistency of extremum estimators for parametric models

Our first step is to provide consistency of $\hat{\theta}$. The consistency analysis is based on the premise that the population problem has a unique maximum and the first-order condition has an isolated well-defined root corresponding to the global maximum. Many commonly used models have multiple local extrema, leading to multiple roots of the first-order condition. To facilitate our analysis we assume that the researcher is able to isolate a subset of the parameter space that contains the global maximum. For simplicity we will associate this subset with the entire parameter space Θ . The above discussion is formalized in the following identification assumption.

ASSUMPTION 10. *The map $\Theta \mapsto \mathbb{R}^k$ defined by $D(\theta) = \frac{\partial}{\partial \theta} E[g(Z_i, \theta)]$ is identified at $\theta_0 \in \Theta$. In other words from $\lim_{n \rightarrow \infty} \|D(\theta_n)\| = 0$ it follows that $\lim_{n \rightarrow \infty} \|\theta_n - \theta_0\| = 0$ for any sequence $\theta_n \in \Theta$. Moreover, $g(\theta) = E[g(Z_i, \theta)]$ is locally quadratic at θ_0 with $g(\theta) - g(\theta_0) \lesssim -d(\theta, \theta_0)^2$.*

The next assumption maintains suitable measurability requirements.

ASSUMPTION 11. *The parameter space Θ has a compact cover. For each n , there exists a countable subset $T_n \subset \Theta$ such that*

$$P^* \left(\sup_{\theta \in \Theta} \inf_{\theta' \in T_n} \|g(Z_i, \theta) - g(Z_i, \theta')\|^2 > 0 \right) = 0,$$

where P^* stands for the outer measure. In general, this condition states that the values of the moment function on the parameter space Θ can be approximated arbitrarily well (with probability one) by its values on a countable subset of Θ . If the moment function is continuous, it trivially satisfies this condition, but it also allows us to consider the moments defined by discontinuous functions. More precisely, Assumption 11 is a sufficient condition for the moment function to be *image admissible Suslin*. As it is discussed in Dudley (1999) and Kosorok (2008) this property will be required to establish the functional uniform law of large numbers needed for consistency.

For global consistency we require the population objective function to be sufficiently smooth not only at the true parameter, but also uniformly in the entire parameter space Θ for which we can rely on Assumption 4 that we previously used to establish uniform consistency for the estimate of the derivative of the sample moment function.

We will organize the discussion below by different classes of functions that can be used in practice. We start with functions that have absolutely locally bounded finite differences. Indicator functions and other functions with finite jumps fall into this category. Then we consider a class of functions that

have polynomial bounds on finite differences. This class includes Lipschitz and Hölder-continuous functions which can have “mildly explosive” finite differences (those which have infinite jumps but approach infinity slower than some power of the distance to the point of discontinuity).

3.2.1 Functions with absolutely bounded finite differences

For the proof of consistency of the extremum estimator we need to provide primitive conditions for the uniform convergence in probability of the numerical derivative of the sample objective function to the derivative of the population objective function. This proof usually invokes the use of maximum inequalities that can bound the expectation of extreme deviations of the sample objective function in small neighborhoods of the parameter space. It is hard to work with the maximum inequality directly when the sample objective function experiences finite jumps: in this case small deviations of the parameter may lead to finitely large changes in the objective function. However, establishing the uniform convergence in probability still remains possible if we are willing to analyze more delicate properties of the function class under consideration. In our analysis we focus on the class of functions outlined by Assumption 5.

The following theorem establishes the consistency of numerical gradient-based extremum estimators for the class of possibly discontinuous functions described by Assumption 5. It is a corollary of Theorem 2 following directly from the uniform convergence in probability as in Amemiya (1985) and, therefore, we omit the proof.

THEOREM 6. *Under assumptions 10, 11, 4, and 5, as long as $\varepsilon_n \rightarrow 0$ and $\frac{n\varepsilon_n}{\log n} \rightarrow \infty$,*

$$\sup_{\theta \in \Theta} \|L_{1,p}^{\varepsilon_n} \hat{Q}(\theta) - G(\theta)\| = o_p(1).$$

Consequently, $\hat{\theta} \xrightarrow{p} \theta_0$ if $\|L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta})\| = o_p(1)$.

This Theorem suggests that even though the sample objective function can be discontinuous with finite jumps, as long as it is Lipschitz-continuous in the mean square, one can use the numerical gradient-based routine for its optimization as long as the step size decays logarithmically with the sample size.

3.2.2 Functions with polynomial envelopes for finite differences

Our results in the previous subsection refer to the classes of objective functions for which the changes in the values of the objective function may not be directly related to the magnitude of the parameter changes. In this subsection we consider the case where such connection can be established. Surprisingly, our results are also valid for the cases of substantially irregular behavior of the objective function when its first derivative approach infinity in the vicinity of the maximum or minimum. A case in point is the objective function defined by $g(Z_i, \theta) = \sqrt{|Z_i - \theta|}$ for which local changes in θ around the origin lead to inversely proportional changes in the moment function. It turns out that

this explosive behavior can be compensated by an appropriate choice of the sequence of step sizes for the numerical derivative. It turns out that the objective functions of this type belong to a class of functions outlined in Assumption 6 which includes the Hölder-continuous functions.

We note that Assumption 6 (i) restricts our analysis to the functions that have a polynomial envelope on their finite differences. On the other hand, provided that γ can be very close to zero, it allows the finite differences of functions to be locally “explosive”. For instance, if we consider finite differences of the objective function $g(Z_i, \theta) = \sqrt{|Z_i - \theta|}$ around the origin, we note that they will be proportional to $1/\sqrt{\epsilon}$ and will not shrink with the decrease in ϵ . It turns out that this still allows us to provide consistency for the numerical gradient-based estimators.

The following theorem establishes consistency under a condition on the step size sequence that is a function of the sample size and the modulus of continuity of the empirical process. It is essentially a replica of theorem 3 and hence stated without proof.

THEOREM 7. *Under Assumptions 4, 6, 10, and 11, as long as $\epsilon_n \rightarrow 0$ and $n\epsilon_n^{2-2\gamma} \rightarrow \infty$,*

$$\sup_{\theta \in \Theta} \|L_{1,p}^{\epsilon_n} \hat{Q}(\theta) - G(\theta)\| = o_p(1).$$

Consequently, $\hat{\theta} \xrightarrow{p} \theta_0$ if $\|L_{1,p}^{\epsilon_n} \hat{Q}(\hat{\theta})\| = o_p(1)$.

For models that have Lipschitz-continuous sample objective functions (which include models with smooth sample objective functions) $\gamma = 1$. In this case the restriction $n\epsilon_n^{2-2\gamma} \rightarrow \infty$ holds trivially. This implies that for smooth models the sequence of step sizes can approach zero arbitrarily fast.³

3.3 Rate of convergence and asymptotic distribution in parametric case

3.3.1 Functions with absolutely bounded finite differences

In the previous section we provided sufficient conditions that determine consistency of the estimator that equates the finite-difference approximation of the gradient of objective function to zero. For the classes where local parameter changes do not lead to proportional changes in the sample objective function we restricted our attention to the functions with absolutely bounded finite differences forming Euclidean classes. Our condition provided the result that the numerical derivative of the sample objective function converges to the derivative of the population objective function uniformly in probability. In addition, making use of the result of Lemma 1 we can establish the precise rate of convergence for the objective function. For a given neighborhood $N(\theta_0)$,

$$\sup_{\theta \in N(\theta_0)} \sqrt{\frac{n\epsilon_n}{\log n}} \|L_{1,p}^{\epsilon_n} \hat{Q}(\theta) - L_{1,p}^{\epsilon_n} Q(\theta)\| = O_p(1),$$

under Assumption 5 and $n\epsilon_n/\log n \rightarrow \infty$.

³In Section 7 we point at some problems that are associated with “too fast” convergence of the step size sequence to zero. These problems, however, are not statistical and are connected with the machine computing precision.

Once we have “located” the parameter, we can investigate the behavior of the sample objective function in shrinking neighborhoods of size $\left(\frac{\log n}{n\varepsilon_n}\right)$. It turns out, that in such neighborhoods we can improve the rate result by choosing the step size in accordance with the radius of the neighborhoods containing the true parameter.

LEMMA 5. Suppose $\hat{\theta} \xrightarrow{p} \theta_0$, $L_{1,p}^{\varepsilon} \hat{Q}(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)$, and the assumptions of Theorem 6 hold.

(i) If $n\varepsilon_n/\log n \rightarrow \infty$, and $n\varepsilon^{1+4p} = O(1)$, then $\sqrt{\frac{n\varepsilon_n}{\log n}} d(\hat{\theta}, \theta_0) = o_{P^*}(1)$.

(ii) If $n\varepsilon_n^{1+4p} = o(1)$, and $\frac{n\varepsilon_n^3}{\log n} \rightarrow \infty$ we have

$$\sup_{d(\hat{\theta}, \theta_0) = O\left(\sqrt{\frac{\log n}{n\varepsilon_n}}\right)} \left(L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0) - L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) + L_{1,p}^{\varepsilon_n} Q(\theta_0) \right) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right).$$

We can use this result to establish the rate of convergence of the resulting estimator.

THEOREM 8. Suppose $\hat{\theta} \xrightarrow{p} \theta_0$ and $L_{1,p}^{\varepsilon} \hat{Q}(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)$. Under Assumptions of Theorem 6, if $n\varepsilon_n/\log n \rightarrow \infty$, $\varepsilon_n = o\left(\sqrt{\frac{\log n}{n\varepsilon_n}}\right)$, and $\sqrt{n\varepsilon^{1+p}} = O(1)$, then $\sqrt{n\varepsilon_n} d(\hat{\theta}, \theta_0) = O_P^*(1)$.

We have established the convergence rate for the finite-difference based extremum estimators with bounded finite differences. Under additional assumptions we can also establish the normality of the asymptotic distribution, which requires showing stochastic equicontinuity as a corollary of the rate of convergence.

COROLLARY 1. Under conditions of theorem 6 with $\sqrt{n\varepsilon_n^{1+p}} = o(1)$, and $n\varepsilon_n^3 \rightarrow \infty$ we have

$$\sup_{d(\hat{\theta}, \theta_0) = O\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)} \left(L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0) - L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) + L_{1,p}^{\varepsilon_n} Q(\theta_0) \right) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right).$$

Proof. The result in the corollary follows directly from the result in Lemma 5 if one notices that $n\varepsilon_n^3 \rightarrow \infty$ guarantees that $\varepsilon_n = o\left(\frac{\log n}{n\varepsilon_n}\right)$. Then by using parametrization $\theta_n = \theta_0 + \frac{t_n}{\sqrt{n\varepsilon_n}}$ we replicate the argument of Lemma 5. \square

In the proof of Lemma 5 we found that a convenient normalization for finite differences of the sample objective function is $(g(\cdot, \theta + \varepsilon_n) - g(\cdot, \theta - \varepsilon_n))/\sqrt{\varepsilon_n}$. We can impose an assumption on this object that assures its normality at θ_0 such as the standard Lindeberg condition which is often used to show pointwise asymptotic normality of kernel smoothers.

ASSUMPTION 12. A CLT holds: As $n \rightarrow \infty$ and $\varepsilon_n \rightarrow 0$,

$$\frac{\mathbb{G}_n(\theta_0 + \varepsilon_n) - \mathbb{G}_n(\theta_0 - \varepsilon_n)}{\sqrt{\varepsilon_n}} \xrightarrow{d} N(0, \Omega).$$

Based on this intuition we can provide the following Theorem.

THEOREM 9. *Assume that the conditions of theorem 6 hold but with $\sqrt{n\varepsilon_n^{1+p}} = o(1)$. In addition, suppose that the Hessian matrix $H(\theta)$ of $g(\theta)$ is continuous, nonsingular and finite at θ_0 . Then if Assumption 12 holds with $n\varepsilon_n^3 \rightarrow \infty$*

$$\sqrt{n\varepsilon_n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, H(\theta_0)^{-1} \Omega H(\theta_0)^{-1}\right).$$

3.3.2 Functions with polynomial envelopes for finite differences

In case where functions of interest permit power envelopes on the finite differences, we can establish the rate of convergence and describe the asymptotic distribution of the resulting estimator. Next, we establish the rate of convergence and the asymptotic distribution of the numerical derivative based M-estimator for the functions that admits polynomial envelopes on finite differences. We provide the following general result.

THEOREM 10. *Suppose $\hat{\theta} \xrightarrow{p} \theta_0$ and $L_{1,p}^\varepsilon \hat{Q}(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n^{1-\gamma}}}\right)$. Under Assumptions 11 and 6, if $n\varepsilon_n^{2-2\gamma}/\log n \rightarrow \infty$ and $\sqrt{n\varepsilon_n^{1-\gamma+2p}} = O(1)$, then $\sqrt{n\varepsilon_n^{1-\gamma}}d(\hat{\theta}, \theta_0) = O_P^*(1)$.*

This result is a Z-estimator version of Theorem 3.2.5 in Van der Vaart and Wellner (1996). Note that given the consistency assumption, the conditions required for obtaining the rate of convergence are weaker. For a typical two sided derivative $\nu = 2$. In this case, for a regular parametric model where $\gamma = 1$, the condition $\sqrt{n\varepsilon^2} \rightarrow 0$ is needed to obtain the usual asymptotic distribution without bias. This rate is compatible with $\sqrt{n\varepsilon} \rightarrow \infty$ and also allows for an ε sequence that delivers a consistent, albeit non-optimal, estimator of the asymptotic variance.

We now proceed with the analysis of the asymptotic distribution of the estimator that we obtain by equating the numerical derivative to zero. The following simplifying assumption suggests the asymptotic normality of the numerical derivative of the sample objective function at the true parameter value. It can be established for a sufficiently slow step size convergence rate by verifying the Lindeberg conditions and the finiteness of the covariance function. Section 3.4 establishes a more general form of the distribution of the estimator for the cases where the step sizes approaches zero at different rates. Asymptotic normality turns out to be a special case for the more general distribution.

ASSUMPTION 13. *A CLT holds: As $n \rightarrow \infty$ and $\varepsilon_n \rightarrow 0$,*

$$\frac{\mathbb{G}_n(\theta_0 + \varepsilon_n) - \mathbb{G}_n(\theta_0 - \varepsilon_n)}{\varepsilon_n^\gamma} \xrightarrow{d} N(0, \Omega).$$

The following theorem establishes the asymptotic normality of the numerical derivative-based estimator with an additional assumption regarding the derivatives of the population objective function.

THEOREM 11. *Assume that the conditions of theorem 10 hold but with $\sqrt{n\varepsilon_n^{1+2p-\gamma}} = o(1)$. In addition, suppose that the Hessian matrix $H(\theta)$ of $g(\theta)$ is continuous, nonsingular and finite at θ_0*

with Assumption 13 valid. Furthermore, $\sqrt{n}\varepsilon_n^{2-\gamma} \rightarrow \infty$. Then

$$\sqrt{n}\varepsilon_n^{1-\gamma} (\hat{\theta} - \theta_0) \xrightarrow{d} N(0, H(\theta_0)^{-1} \Omega H(\theta_0)^{-1}).$$

The additional assumption $\sqrt{n}\varepsilon_n^{2-\gamma} \rightarrow \infty$ turns out to be stronger for smooth models than for nonsmooth models. This is an artifact that we are relying on Assumption 6 and the convergence rate result in theorem 10 to obtain stochastic equicontinuity. When $\gamma = 1$, the conditions are consistent as long as $p \geq 1$, or as long as a two sided central derivative is used.

However, for smooth models when $\gamma = 1$, we might be willing to impose stronger assumptions on the sample objective function (e.g. Lemma 3.2.19 in Van der Vaart and Wellner (1996)) to weaken this requirement. The next theorem states such an alternative result.

THEOREM 12. *Suppose the conditions of theorem 11 hold except $\sqrt{n}\varepsilon_n^{2-\gamma} \rightarrow \infty$. Suppose $n\varepsilon_n \rightarrow \infty$. Suppose further that $g(z_i, \theta)$ is mean square differentiable in a neighborhood of θ_0 : for measurable functions $D(\cdot, \cdot) : \mathcal{Z} \times \Theta \rightarrow \mathbb{R}^p$ such that*

$$E[g(Z, \theta_1) - g(Z, \theta_2) - (\theta_2 - \theta_1)' D(Z, \theta_1)]^2 = o(\|\theta_1 - \theta_2\|^2),$$

$E\|D(Z, \theta_1)\|^2 < \infty$ for all θ_1 , and $\theta_2 \in \mathcal{N}_{\theta_0}$. Define $q_\varepsilon(z_i, \theta) = L_{1,p}^\varepsilon g(z_i; \theta) - D(z, \theta)$, Assume that

$$\sup_{d(\theta, \theta_0)=o(1), \varepsilon=o(1)} [\mathbb{G}q_\varepsilon(z_i, \theta_1) - \mathbb{G}q_\varepsilon(z_i, \theta_0)] = o_p(1),$$

and $D(z_i, \theta)$ is Donsker in $d(\theta, \theta_0) \leq \delta$, then the conclusion of theorem 11 holds.

Note that we still require $\sqrt{n}\varepsilon_n^2 \rightarrow 0$ to remove the asymptotic bias, and $n\varepsilon_n \rightarrow \infty$ to eliminate the variance in assumption 13, but we no longer require $\sqrt{n}\varepsilon \rightarrow \infty$. The conditions of this theorem are best understood in the context of a quantile regression estimator. Consider $g(z_i, \theta) = |z_i - \theta|$, and $p = 2$, so that $D(z, \theta) = \text{sgn}(z_i - \theta)$ and

$$q_\varepsilon(z_i, \theta) = \frac{(z_i - \theta)}{\varepsilon} 1(|z_i - \theta| \leq \varepsilon).$$

Then we can bound $q_\varepsilon(z_i, \theta_1) - q_\varepsilon(z_i, \theta_0)$ by, depending on which of $d(\theta, \theta_0)$ and ε is larger, the product between $\frac{1}{\varepsilon} \max(|z_i - \theta|, |z_i - \theta_0|)$ and the maximum of $1(|z_i - \theta| \leq \varepsilon + 1(|z_i - \theta_0| \leq \varepsilon))$, and $[1(\theta - \varepsilon \leq z_i \leq \theta + \varepsilon_0) + 1(\theta_0 - \varepsilon \leq z_i \leq \theta_0 + \varepsilon_0)]$. Since $\max(|z_i - \theta|, |z_i - \theta_0|) \leq \varepsilon$ when $q_\varepsilon(z_i, \theta) - q_\varepsilon(z_i, \theta_0)$ is nonzero, the last condition in theorem 12 is clearly satisfied by the euclidean property of the indicator functions. Alternatively, the $q_\varepsilon(z_i, \theta)$ function in the last condition can also be replaced directly by $L_{1,p}^\varepsilon g(z_i, \theta)$.

3.4 General distribution results

The previous asymptotic normality result turns out to be an artifact of an excessively slow rate of approach of the sequence of step sizes ε_n to zero. Our previous assumption required that we choose

$n\varepsilon_n^3 \rightarrow \infty$ for the discontinuous case and $\frac{\varepsilon_n^{\gamma-2}}{\sqrt{n}} = o(1)$ for continuous case. This assumption can be relaxed, at a cost of making the asymptotic distribution non-standard. However, this weakening also demonstrates that the numerical derivative-based estimators for non-smooth sample objective functions have interesting parallels with the cube-root asymptotics of Kim and Pollard (1990).

The following assumption has a simple implication for the covariance function of the sample objective function. It requires that the pairwise products of the sample objective function computed at different points in a vanishing neighborhood of the true parameter value have continuous expectations. Moreover, the variance of the numerical derivative of the sample objective function is infinitesimal at the points where pointwise derivative of the sample objective function may not exist. We first present the assumption for the case of discontinuous moment function.

ASSUMPTION 14. Suppose that for each sequence $\varepsilon_n \rightarrow 0$ with $\sqrt{n\varepsilon_n^{1+p}} = o(1)$

(i) The covariance function

$$H_{n,\varepsilon_n}(s, t) = \lim_{\alpha \rightarrow \infty} E \left[\frac{\alpha}{\varepsilon_n} \left(g \left(Z_i, \theta_0 + \varepsilon_n + \frac{s}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{s}{\alpha} \right) \right) \right. \\ \left. \times \left(g \left(Z_i, \theta_0 + \varepsilon_n + \frac{t}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{t}{\alpha} \right) \right) \right]$$

is finite and has a finite limit as $n \rightarrow \infty$ for $s, t \in \mathbb{R}$ and $\varepsilon_n = O\left(\frac{1}{\sqrt[3]{n}}\right)$.

(ii) For each t and each $\delta > 0$

$$\lim_{\alpha \rightarrow \infty} E \left[\frac{\alpha}{\varepsilon_n} \left(g \left(Z_i, \theta_0 + \varepsilon_n + \frac{t}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{t}{\alpha} \right) \right)^2 \right. \\ \left. \times \mathbf{1} \left\{ \left\| g \left(Z_i, \theta_0 + \varepsilon_n + \frac{t}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{t}{\alpha} \right) \right\| > \alpha \delta \right\} \right] = 0.$$

Similarly, we can impose an assumption for the case of Hölder-continuous objective functions.

ASSUMPTION 15. Suppose that for each sequence $\varepsilon_n \rightarrow 0$ with $\sqrt{n\varepsilon_n^{1+2p-\gamma}} = o(1)$

(i) The covariance function

$$H_{n,\varepsilon_n}(s, t) = \lim_{\alpha \rightarrow \infty} E \left[\frac{\alpha}{\varepsilon_n^{2\gamma}} \left(g \left(Z_i, \theta_0 + \varepsilon_n + \frac{s}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{s}{\alpha} \right) \right) \right. \\ \left. \times \left(g \left(Z_i, \theta_0 + \varepsilon_n + \frac{t}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{t}{\alpha} \right) \right) \right] < \infty$$

is finite and has a finite limit as $n \rightarrow \infty$ for $s, t \in \mathbb{R}$ and $\varepsilon_n = O\left(n^{-1/2(2-\gamma)}\right)$.

(ii) For each t and each $\delta > 0$

$$\lim_{\alpha \rightarrow \infty} E \left[\frac{\alpha}{\varepsilon_n^{2\gamma}} \left(g \left(Z_i, \theta_0 + \varepsilon_n + \frac{t}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{t}{\alpha} \right) \right)^2 \right. \\ \left. \times \mathbf{1} \left\{ \left\| g \left(Z_i, \theta_0 + \varepsilon_n + \frac{t}{\alpha} \right) - g \left(Z_i, \theta_0 - \varepsilon_n + \frac{t}{\alpha} \right) \right\| > \alpha \delta \right\} \right] = 0.$$

We combine Assumptions 14 and 15 with Assumptions 5 and 6 that restrict the attention to particular (large) parametric classes of functions.

THEOREM 13. 1. Suppose that assumptions 4, 5 and 15 hold. The population objective has a finite Hessian $H(\theta_0)$ at θ_0 and $\frac{n^3}{\varepsilon_n} = O(1)$. Then

$$\sqrt{\frac{\varepsilon_n}{n}} \sum_{i=1}^n L_{1,p}^{\varepsilon_n} g \left(Z_i, \theta_0 + \frac{t}{\sqrt{n\varepsilon_n}} \right) \rightsquigarrow Z(t).$$

2. Suppose that assumptions 4, 6 and 15 hold. The population objective has a finite Hessian $H(\theta_0)$ at θ_0 and $\frac{\varepsilon_n^{\gamma-2}}{\sqrt{n}} = O(1)$. Then

$$\frac{\varepsilon_n^{1-\gamma}}{\sqrt{n}} \sum_{i=1}^n L_{1,p}^{\varepsilon_n} g \left(Z_i, \theta_0 + \frac{t}{\varepsilon_n^{1-\gamma} \sqrt{n}} \right) \rightsquigarrow Z(t).$$

In these expressions $Z(t)$ is a mean-zero Gaussian process with covariance function $H(s, t) = \sum_l = -p^l l^2 c_l^2 \lim_{n \rightarrow \infty} H_{n, l\varepsilon_n}(s, t)$. Then $\sqrt{n}\varepsilon_n^{1-\gamma} (\hat{\theta} - \theta_0) \rightsquigarrow \hat{t}$, where \hat{t} is defined by $Z(\hat{t}) = H(\theta_0) \hat{t}$. In one dimension, \hat{t} can be interpreted as a boundary-crossing distribution.

In the special case where $\sqrt{n}\varepsilon_n^{2-\gamma} \rightarrow \infty$, $Z(t)$ is normal and does not depend on t .

Remark: Theorem 13 establishes that the lowest rate at which ε_n approaches zero is $n^{-1/2(2-\gamma)}$ for the case of Hölder-continuous moment functions and $n^{-1/3}$ for the case of discontinuous objective functions. A faster approach of the step size to zero leads to a loss of stochastic equicontinuity. This condition also provides the slowest convergence rate for the estimator of $n^{1/(2(2-\gamma))}$ for the Hölder case and $n^{1/3}$ for the discontinuous case.

Theorem 13 shows that depending on the step size sequence chosen, the asymptotic distribution is a “hybrid” between the distribution of the original extremum estimators and the distribution of smoothed estimators. The asymptotics in case of the “under-smoothed” numerical derivative is non standard and is driven by the boundary-crossing distribution of the limiting process (a discussion of this distribution can be found, for instance in Durbin (1971, 1985)).

4 Consistency of semiparametric extremum estimators

Our approach of transforming the problem of extremum estimation to the problem of solving a numerical first-order condition can be extended to the case where the parameter space is either entirely infinite-dimensional or contains an infinite-dimensional component. We consider a metric product space $\Theta \times \mathcal{H}$ where Θ is a compact subset of a Euclidean space \mathbb{R}^p and \mathcal{H} is a functional Banach space. Semiparametric extremum estimators typically arise from conditional moment equations. We consider a population moment equation

$$m(\eta, \theta, z) = E[\rho(\theta, \eta, Y_i) | Z_i = z] = 0$$

with $\rho : \Theta \times \mathcal{H} \times \mathcal{Y} \mapsto \mathcal{M} \subset \mathbb{R}^k$. The estimation problem is re-casted into an optimization problem by defining the objective $Q(\theta, \eta) = E[m(\theta, \eta, Z_i)' W(Z_i) m(\theta, \eta, Z_i)]$ using a $k \times k$ positive semi-definite (almost everywhere in \mathcal{Z}) weighting matrix $W(\cdot)$. The estimator minimizes the sample objective function with the infinite-dimensional component η over sieve space \mathcal{H}_n

$$(\hat{\theta}, \hat{\eta}) = \arg \min_{(\theta, \eta) \in \Theta \times \mathcal{H}_n} \hat{Q}_n(\theta, \eta) = \frac{1}{n} \sum_{i=1}^n \hat{m}(\theta, \eta, z_i)' \widehat{W}(z_i) \hat{m}(\theta, \eta, z_i).$$

A typical set of the necessary conditions for the optimum of $Q(\theta, \eta)$ can be found, for example, in the general class of mathematical programming problems in Pshenichnyi (1971). Consider a cone K in $\Theta \times \mathcal{H}$. If $(\theta_0, \eta_0) \in \Theta \times \mathcal{H}$ optimizes $Q(\theta, \eta)$, then there exists a number $\lambda_0 \geq 0$ such that $\theta(\lambda) = \theta_0 + \lambda \delta \in \Theta$ and $\eta(\lambda) = \eta_0 + \lambda w \in \mathcal{H}$ for all $(\delta, w) \in K$ and $\lambda \in \mathbb{R}$. Moreover, $\lambda_0 \zeta(\delta, w) = 0$, and

$$\lim_{\lambda \rightarrow +0} \frac{Q(\theta(\lambda), \eta(\lambda)) - Q(\theta_0, \eta_0)}{\lambda} \leq \zeta(\delta, w),$$

where $\zeta(\delta, w)$ is a functional which is convex with respect to (δ, w) . If we assume that the objective functional is strictly concave at (θ_0, η_0) then $\lambda_0 > 0$. This transforms the necessary conditions to

$$\lim_{\lambda \rightarrow +0} \frac{Q(\theta_0 + \lambda \delta, \eta_0 + \lambda w) - Q(\theta_0, \eta_0, x)}{\lambda} = 0.$$

In particular, this directional derivative should be equal to zero for all directions within the cone K . Specifically, if $\Theta \times \mathcal{H}$ is a linear space, then this should be valid for all directions in $(\Theta - \theta_0) \times (\mathcal{H} - \eta_0)$. If the functional is Frechét differentiable in (θ, η) , then the directional derivative exists in all directions in K and we can write the necessary condition for the extremum in the simple form:

$$\frac{d}{d\tau} Q(\theta_0 + \tau \delta, \eta_0 + \tau w) |_{\tau=0} = 0,$$

in all directions $w \in \mathcal{H} - \eta_0$. In particular if Θ is a finite-dimensional vector space and \mathcal{H} is a finite-dimensional functional space, then we can construct a system of first-order condition that exactly identifies a parameter pair (θ, η) as

$$\begin{aligned} \frac{\partial Q(\theta, \eta)}{\partial \theta_k} &= 0, \quad \text{for } k = 1, \dots, p, \\ \frac{\partial Q(\theta, \eta)}{\partial \eta} [\psi_j] &= 0, \quad \text{for } j = 1, \dots, G, \end{aligned} \tag{4.6}$$

where $\psi_j(\cdot)$ is a system of distinct elements of \mathcal{H} . In cases where the functional space \mathcal{H} is infinite-dimensional, then we define the population solution to the system of first-order condition as a limit of sequence of solutions in the finite-dimensional sieve spaces \mathcal{H}_n such that $\mathcal{H}_n \subseteq \mathcal{H}_{n+1} \subseteq \mathcal{H}$ for all n .

We consider the class of models where such substitution of maximization of the functional to finding a solution to the system of functional equations is possible. We formalize this concept by the following assumption.

ASSUMPTION 16. Suppose that (θ_0, η_0) is the maximizer of the functional $Q(\theta, \eta)$ and \mathcal{H}_n is the sieve space such that $\mathcal{H}_n \subset \mathcal{H}_{n+1} \subset \mathcal{H}$. The set \mathcal{H}_∞ is complete in \mathcal{H} . The sets \mathcal{H}_n share the same basis $\{\psi_j\}_{j=0}^\infty$ and $\langle \eta_0, \psi_j \rangle \rightarrow 0$ as $j \rightarrow \infty$. We assume that the left-hand side of (4.6) is continuous with respect to the strong product norm in $\Theta \times \mathcal{H}$. Suppose that (θ_n, η_n) solves (4.6). Then for any sequence of the sieve spaces \mathcal{H}_n satisfying the above conditions the corresponding system of solutions (θ_n, η_n) converges to (θ_0, η_0) in the strong norm.

This identification condition establishes the properties of the population objective function. We will now consider the properties of the sample objective function which can be expressed as

$$\widehat{Q}(\theta, \eta) = \frac{1}{n} \sum_{i=1}^n \widehat{m}(\theta, \eta, z_i)' \widehat{W}(z_i) \widehat{m}(\theta, \eta, z_i).$$

Following our analysis in Section 2.6 where we focused on computing directional derivatives of semiparametric models via finite-difference methods, we consider two cases where the estimator for the conditional moment function is obtained via series approximation (2.2) or kernel smoothing (2.3). The transformed system of equations that needs to be solved to obtain the estimator $\hat{\eta}$ and $\hat{\theta}$ can be obtained by directly applying the directional derivative to the objective function. Without loss of generality to simplify the algebra we focus on the case where the finite-dimensional parameter θ is scalar, the conditional moment function is one-dimensional, and the weighting matrix is a non-negative weighting function. We will use the notation for the step size ϵ_n for differentiation with respect to the finite-dimensional parameter and the notation τ_n for the step size of differentiation with respect to the infinite-dimensional parameter. This leads us to the expression for the numerical first-order condition in the form

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n L_{1,p}^{\epsilon_n} \widehat{m}(\theta, \eta, z_i)' \widehat{W}(z_i) \widehat{m}(\theta, \eta, z_i) &= o_p(1), \\ \frac{1}{n} \sum_{i=1}^n L_{1,p}^{\tau_n, \psi_j} \widehat{m}(\theta, \eta, z_i)' \widehat{W}(z_i) \widehat{m}(\theta, \eta, z_i) &= o_p(1), \end{aligned} \tag{4.7}$$

for $j = 1, \dots, G_n$.

In the following Theorem we provide the uniform convergence result. As we notice, under standard assumptions regarding the series expansion and the kernel estimator, there is an interference between the step size for numerical differentiation and the choice of the tuning parameter (number of terms in the expansion or the bandwidth).

THEOREM 14. Suppose that matrix $W(\cdot)$ is a.s. positive-definite and $\widehat{W}(\cdot) = W(\cdot) + o_p(1)$ uniformly in z . In addition, $\rho(\pi_n \eta, \eta) = O(n^{-\phi})$.

Under Assumption 7 and 8 for the series estimator, provided that $\epsilon_n \rightarrow 0$, $\min\{N^\alpha, n^\phi\} \min\{\epsilon_n, \tau_n\} \rightarrow \infty$, $\frac{n^{1-2r_0} \epsilon_n}{N^2 \log n^{1-2r_0}} \rightarrow \infty$, and $\frac{n^{1-2r_0} \tau_n}{N^2 \log n^{1-2r_0}} \rightarrow \infty$; under assumptions 8 and 9 for the kernel estimator, provided that $\epsilon_n \rightarrow 0$, $\min\{b_n^{-N}, n^\phi\} \min\{\epsilon_n, \tau_n\} \rightarrow \infty$, $b_n \rightarrow 0$, $\frac{n^{1-2r_0} \epsilon_n b_n^{dz}}{\log n^{1-2r_0}} \rightarrow \infty$, and $\frac{n^{1-2r_0} \epsilon_n b_n^{dz}}{\log n^{1-2r_0}} \rightarrow \infty$, $(\hat{\theta}, \hat{\eta}) \xrightarrow{p} (\theta_0, \eta_0)$ provided that system (4.7) is satisfied.

The result of this theorem is based on the uniform consistency result for the directional derivatives in Section 2.6. We note that those results provided consistency for the directional derivatives uniformly in θ , η , w and z . However, uniformity in z seems excessive for the estimator delivered by (4.7) provided that the estimator requires summation over the sample $\{z_i\}_{i=1}^n$. It turns out that in some special cases we can use this feature to improve the convergence result and deliver fast convergence rate for the estimator of the finite-dimensional parameter. Our result will rely on the properties of U-statistics covered in Section 5.

5 Numerical derivatives of U-statistics

5.1 Consistency of numerical derivatives

Numerical derivatives can also be used for objective functions that are based on U-statistics, an example of which is the maximum rank correlation estimator of Sherman (1993). The model considered in this section is parametric. Second order U-statistics, which we focus on, are the most commonly used in applications. A U-statistic objective function is defined from an i.i.d. sample $\{Z_i\}_{i=1}^n$ by a symmetric function $g(Z_i, Z_j, \theta)$ as as

$$\hat{g}(\theta) = \frac{1}{n(n-1)} S_n(f) \quad \text{where} \quad S_n(f) = \sum_{i \neq j} f(Z_i, Z_j, \theta). \quad (5.8)$$

We denote the expectation with respect to the independent product measure on $\mathcal{Z} \times \mathcal{Z}$ by E_{zz} and the expectation with respect to a single measure by E_z . The population value can then be written as $g(\theta) = E_{zz} f(Z_i, Z_j, \theta)$. This population objective function satisfies Assumption 4.

Following Serfling (1980), the following decomposition of the objective function into an empirical process and a degenerate U-process component can be used to establish the statistical properties of approximating $G(\theta_0) = \frac{\partial}{\partial \theta} g(\theta)$ by $L_{1,p}^{\varepsilon_n} \hat{g}(\hat{\theta})$,

$$\hat{g}(\theta) = g(\theta) + \hat{\mu}_n(\theta) + \frac{1}{n(n-1)} S_n(u), \quad (5.9)$$

where

$$\hat{\mu}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mu(Z_i, \theta), \quad \mu(z, \theta) = E_z f(Z_i, z, \theta) + E_z f(z, Z_i, \theta) - 2g(\theta),$$

and

$$u(z, z', \theta) = f(z, z', \theta) - E_z f(Z_i, z, \theta) - E_z f(z', Z_i, \theta) + g(\theta).$$

The condition for controlling the degenerate U-process will be weaker than that needed for the empirical process term because of its faster convergence rate. We maintain Assumption 5, with the new interpretation of functions $g(\cdot, \cdot, \theta)$. We also make the following additional assumptions.

ASSUMPTION 17. *The projections $\mu(z, \theta)$ are Lipschitz-continuous in θ uniformly over z .*

This assumption depends on the distribution of Z_i . For instance, when the kernel is defined by indicator functions, the expectation will be continuous in parameter for sufficiently smooth distribution of Z_i . It controls the impact of numerical differentiation on the projection term by the maximum inequality for Lipschitz-continuous functions:

$$E^* \sup_{d(\theta, \theta_0) = o(1)} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mu(Z_i, \theta + \epsilon_n) - \mu(Z_i, \theta - \epsilon_n) - g(\theta + \epsilon_n) + g(\theta - \epsilon_n)) \right| \leq C\epsilon_n,$$

for some $C > 0$.

ASSUMPTION 18. *For a neighborhood $N(\theta_0)$ around θ_0 ,*

$$\sup_{\theta \in N(\theta_0), z \in \mathcal{Z}} E |g(Z_i, z, \theta + \epsilon) - g(Z_i, z, \theta - \epsilon)|^2 = O(\epsilon).$$

This assumption allows us to establish an analog of Lemma 1 for the case of the U-processes, which is presented below.

LEMMA 6. *Suppose $\|F\| = \sup_{\theta \in N(\theta_0)} |g(Z_i, Z_j, \theta)| \ll C < \infty$. Under Assumptions 5 and 18, if $n^2 \epsilon_n / \log^2 n \rightarrow \infty$,*

$$\sup_{d(\theta, \theta_0) = o(1)} \|L_{1,p}^{\epsilon_n} \hat{g}(\theta) - L_{1,p}^{\epsilon_n} g(\theta)\| = o_p(1).$$

Consequently, assumption 4 implies that $\hat{G}_1(\hat{\theta}) = o_p(1)$ if $d(\hat{\theta}, \theta_0) = o_p(1)$, as defined in (2.1).

The consistency of the numerical derivatives of U-statistics follows directly from Lemma 6.

THEOREM 15. *Under assumptions, 4, 5 and the conditions of lemma 6, $L_{1,p}^{\epsilon_n} \hat{g}(\hat{\theta}) \xrightarrow{p} G(\theta_0)$ if $\epsilon_n \rightarrow 0$ and $n\epsilon_n^2 / \log^2 n \rightarrow \infty$, and if $d(\hat{\theta}, \theta_0) = o_p(1)$.*

As in the case of the empirical process, this theorem establishes the weakest possible condition on the step size of numerical differentiation when the envelope function of the differenced moment function does not necessarily decrease with the shrinking step size. We note that the resulting condition for the step size is weaker in the case of the U-statistics versus the case of the empirical sums. This is an artifact of the property that the projection of U-statistics tends to be smoother than the kernel function itself leading to a smaller scale of the U-statistic.

5.2 Numerical gradient-based estimation with U-statistics

We consider the solution of the empirical first-order condition $\hat{\theta}$ defined by $\|L_{1,p}^{\epsilon_n} \hat{Q}_n(\hat{\theta})\| = o_p\left(\frac{1}{\sqrt{n}}\right)$. In some cases when the objective function is not continuous, the value that sets the first-order

condition to zero might not exist, so we propose to choose the point that will set the first-order condition very close to zero. In this section we will only consider the distribution results regarding the first numerical derivative.

The structure of the consistency argument for the U-statistic defined estimation problem is similar to that for the standard sample means. In particular, when the kernel function is “sufficiently” smooth, the behavior of the objective function will be dominated by the empirical process component. In that case the analysis from our previous discussion will be valid for the objective function defined by the sample mean of $E_z f(z, Z_i, \theta)$. We maintain Assumption 10 applied to the map $D(\theta) = \frac{\partial}{\partial \theta} E_{zz} [f(Z_i, Z_j, \theta)]$. We also keep Assumption 11 following Arcones and Gine (1993) where the authors state that this assumption, along with the finiteness of the absolute moment of the U-statistic, constitute a sufficient measurability requirement. A deeper discussion of applicability of these conditions can be found in Section 10 of Dudley (1999).

5.3 U-statistics with kernels with absolutely bounded finite differences

We maintain Assumption 5, 17 and 18 for the class of kernels of the U-statistic.

Lemma 6 establishes the consistency result for this objective function which implies that under $n^2 \varepsilon_n / \log^2 n \rightarrow \infty$,

$$\sup_{d(\theta, \theta_0) = o(1)} \|L_{1,p}^{\varepsilon_n} \hat{g}(\theta) - L_{1,p}^{\varepsilon_n} g(\theta)\| = o_p(1).$$

Moreover, we can apply Lemma 10 in Nolan and Pollard (1987). This lemma states that for $t_n \geq \max\{\varepsilon_n^{1/2}, \frac{\log n}{n}\}$ we have for some constant $\beta > 0$

$$P\left(\sup_{\mathcal{F}_n} |S_n(f)| > \beta^2 n^2 t_n^2\right) \leq 2A \exp(-nt_n)$$

However, we note that provided that $\log n \sqrt{\varepsilon_n} / n \rightarrow \infty$, we can strengthen this result. In fact, provided that for sufficiently large n $t_n = \sqrt{\varepsilon_n}$, we note that we can provide condition

$$\sup_{d(\theta, \theta_0) = o(1)} \frac{n^2 \varepsilon_n}{\log^2 n} \|L_{1,p}^{\varepsilon_n} \hat{g}(\theta) - L_{1,p}^{\varepsilon_n} g(\theta)\| = O_p(1).$$

We next repeat the steps that we followed to determine the rate of convergence of the estimators given by sample means.

LEMMA 7. Suppose $\hat{\theta} \xrightarrow{p} \theta_0$ and $L_{1,p}^{\varepsilon} \hat{Q}(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)$. Suppose that Assumptions 5 (i) and (ii), 17 and 18 hold

(i) If $n\sqrt{\varepsilon_n}/\log n \rightarrow \infty$, and $\sqrt{n\varepsilon^{1+p}} = O(1)$, then $\frac{n^2 \varepsilon_n}{\log^2 n} d(\hat{\theta}, \theta_0) = o_{P^*}(1)$.

(ii) If $\sqrt{n\varepsilon_n^{1+p}} = o(1)$, and $\frac{n\varepsilon_n}{\log n} \rightarrow \infty$ we have

$$\sup_{d(\hat{\theta}, \theta_0) = O\left(\frac{\log^2 n}{n^2 \varepsilon_n}\right)} \left(L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0) - L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) + L_{1,p}^{\varepsilon_n} Q(\theta_0)\right) = o_p\left(\frac{1}{n}\right).$$

In this Lemma we, first established the “maximum” radius of the shrinking neighborhood containing the parameter. In the next step we consider the behavior of the objective function in the small neighborhood of order $O\left(\frac{\log^2 n}{n^2 \varepsilon_n}\right)$ of the true parameter. As we show, we can improve upon the rate of the objective function using the envelope property.

We can use this result to establish the rate of convergence of the resulting estimator.

THEOREM 16. *Suppose $\hat{\theta} \xrightarrow{p} \theta_0$ and $L_{1,p}^\varepsilon \hat{Q}(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n}}\right)$. Under Assumptions of Lemma 7, if $n\varepsilon_n/\log n \rightarrow \infty$, and $\sqrt{n\varepsilon_n^{1+p}} = O(1)$, then $\sqrt{nd}(\hat{\theta}, \theta_0) = O_P^*(1)$.*

Proof. We note that by Lemma 7 in the small neighborhoods of the true parameter the U-statistic part has a stochastic order $o_p\left(\frac{1}{n}\right)$. As a result, the sum will be dominated by the projection term. Provided that the projection is Lipschitz-continuous, we can apply the standard rate result in Newey and McFadden (1994) which gives the stochastic order for the first term $O_p\left(\frac{1}{\sqrt{n}}\right)$ and gives the corresponding parametric convergence. \square

The last relevant result concerns the asymptotic distribution of the estimator obtained from the maximization of the U-statistic.

ASSUMPTION 19. *The empirical process corresponding to the projection of the U-statistic $\mathbb{G}_{\mu,n}(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mu(z_i, \theta)$ satisfies*

$$\frac{\mathbb{G}_{\mu,n}(\theta_0 + \varepsilon_n) - \mathbb{G}_{\mu,n}(\theta_0)}{\varepsilon_n} \xrightarrow{d} N(0, \Omega),$$

as $n \rightarrow \infty$ and $\varepsilon \rightarrow 0$.

This assumption allows us to apply the result in Theorem 11 and obtain the following characterization of the asymptotic distribution of the estimator corresponding to the zero of the numerical gradient of the U-statistic.

THEOREM 17. *Assume that the conditions of theorem 16 hold and $\sqrt{n\varepsilon_n^{2p}} = o(1)$. In addition, suppose that the Hessian matrix $H(\theta)$ of $g(\theta)$ is continuous, nonsingular and finite at θ_0 and Assumption 19 holds. Then*

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, H(\theta_0)^{-1} \Omega H(\theta_0)^{-1}\right).$$

6 Applications

7 The choice of the magnitude of the step size

Our asymptotic results are concerned with the optimal choice of the rate of the step size for numerical differentiation. An important practical question will be the choice of the magnitude of the step size

for a particular data sample. In the non-parametric estimation literature there are approaches to the choice of the bandwidth for kernel smoothing. Survey of work on the choice of bandwidth for density estimation can be found in Jones, Marron, and Sheather (1996) with related results for non-parametric regression estimation and estimation of average derivatives in Hardle and Marron (1985) and Hart, Marron, and Tsybakov (1992) among others.

To a large extent, we can obtain the results for the optimal choice of constants simpler than in the case of non-parametric estimation because we will not be interested in the “uniform” step size. Previously we considered the decomposition: $L_{1,p}^{\varepsilon_n} \hat{g}(\hat{\theta}) - G = \hat{G}_1 + \hat{G}_2 + G_3 + G_4$, where

$$G_1 = \left[L_{1,p}^{\varepsilon_n} \hat{g}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} \hat{g}(\theta) \right] - \left[L_{1,p}^{\varepsilon_n} \hat{g}(\theta_0) - L_{1,p}^{\varepsilon_n} g(\theta_0) \right]$$

and

$$G_2 = L_{1,p}^{\varepsilon_n} \hat{g}(\theta_0) - L_{1,p}^{\varepsilon_n} g(\theta_0)$$

and

$$G_3 = L_{1,p}^{\varepsilon_n} g(\hat{\theta}) - L_{1,p}^{\varepsilon_n} g(\theta_0), \quad G_4 = L_{1,p}^{\varepsilon_n} g(\theta_0) - G.$$

We proved that $L_{1,p}^{\varepsilon_n} \hat{g}(\hat{\theta}) - G = O_p(\hat{G}_2 + G_4)$. We can now consider the problem of the optimal constant choice. We consider the mean-squared error as the criterion for the choice of the step size, i.e. the function of interest is

$$\text{MSE}(\varepsilon) = E \| L_{1,p}^{\varepsilon_n} \hat{g}(\hat{\theta}) - G \|^2,$$

which we approximate by the leading terms G_2 and G_4 . We note that

$$L_{1,p}^{\varepsilon_n} g(\theta) = \frac{1}{\varepsilon} \sum_{k=1}^p a_k g(\theta + t_k \varepsilon).$$

Assuming that function $g(\cdot)$ has at least $p+1$ derivatives, we can evaluate the result of application of the numerical derivative as

$$L_{1,p}^{\varepsilon_n} g(\theta) = g'(\theta) + \varepsilon_n^p g^{(p+1)}(\theta) \sum_{k=1}^p \frac{a_k t_k^p}{(p+1)!} + o(\varepsilon_n^p).$$

Thus $G_4 = \varepsilon_n^p g^{(p+1)}(\theta) \sum_{k=1}^p \frac{a_k t_k^p}{(p+1)!} + o(\varepsilon_n^p)$. We can evaluate the variance of G_2 as

$$\begin{aligned} \text{Var}(G_2) &= \frac{1}{n} E \left[\frac{1}{\varepsilon_n} \sum_{k=1}^p a_k (g(\theta + t_k \varepsilon_n, Z_i) - g(\theta + t_k \varepsilon_n)) \right]^2 \\ &= \varepsilon_n^{-(2-2\gamma)} n^{-1} \left[\sum_{k=1}^p a_k^2 \text{Var}(\varepsilon_n^{-\gamma} g(\theta + t_k \varepsilon_n, Z_i)) + \sum_{k,m=1}^p a_k a_m \text{Cov}(\varepsilon_n^{-\gamma} g(\theta + t_k \varepsilon_n, Z_i), \varepsilon_n^{-\gamma} g(\theta + t_m \varepsilon_n, Z_i)) \right] \\ &= \varepsilon_n^{-(2-2\gamma)} n^{-1} V_g(\varepsilon_n). \end{aligned}$$

In case where $\gamma = 1$, the variance of G_2 will not affect the estimation. However, there will still be the numerical error corresponding to the operating precision of computer operations. This error is known and fixed. We denote it $[\delta g]$. Then the total error can be evaluate as

$$\text{MSE}_1(\varepsilon) \approx \varepsilon_n^{2p} \left(g^{(p+1)}(\theta) \sum_{k=1}^p \frac{a_k t_k^p}{(p+1)!} \right)^2 + \frac{[\delta g]}{\varepsilon_n} \sum_{k=1}^p a_k.$$

In case where $\gamma < 1$ the numerical error will be exceeded by the sampling error. As a result, we can compute

$$\text{MSE}_{<1}(\varepsilon) \approx \varepsilon_n^{2p} \left(g^{(p+1)}(\theta) \sum_{k=1}^p \frac{a_k t_k^p}{(p+1)!} \right)^2 + \varepsilon_n^{-(2-2\gamma)} n^{-1} V_g(\varepsilon_n).$$

Then we can choose $\varepsilon_n = \frac{C}{n^r}$, where r is the optimal rate for ε_n if $\gamma < 1$ and $r = 0$ otherwise. The problem is to choose C . In most applications, however, the derivative $g^{(p+1)}$ is unknoww. One simple way of choosing C is the analog of biased cross-validation. We can choose a simple first-order formula for $g^{(p+1)}$ and pick a preliminary (over-smoothed) step size $\varepsilon_n^{**} = \frac{(p+1)\text{Var}(\theta, Z_i)}{n^{1/2(p+1)}}$ then evaluate

$$\widehat{g^{(p+1)}}(\theta) = \frac{1}{\varepsilon_n^{**}} \sum_{k=0}^{[p/2]} g\left(\theta + (-1)^k \varepsilon_n^{**}\right).$$

Plugging this expression into the expression for the mean-squared error, we can obtain the optimal step sizes. Then for $\gamma = 1$ we find that

$$C^{**} = \left(\frac{p! (p+1)! [\delta g] \sum_{k=1}^p a_k}{\left(\widehat{g^{(p+1)}}(\theta) \sum_{k=1}^p a_k t_k^p \right)^2} \right)^{1/(2p+1)}.$$

For $\gamma < 1$ we find

$$C^{**} = \left(\frac{(2-2\gamma)p! (p+1)! V_g(\varepsilon_n^{**})}{\left(\widehat{g^{(p+1)}}(\theta) \sum_{k=1}^p a_k t_k^p \right)^2} \right)^{1/(2p+2-2\gamma)}.$$

Note that if the function $g(\cdot)$ is intensive to compute, the choice of these constants allows one to use a relatively small subsample to calibrate the step sizes. Then one can use these constants to initialize the step sizes on a large scale using the entire sample.

In case where one can compute the function in a relatively straightforward way, calibration of the constants of interest can be performed by minimizing the approximate expression for the mean-squared error with respect to C , taking into account that the step size will enter both in the expression for the derivative $g^{(p+1)}$ and $V_g(\varepsilon_n)$. This approach is equivalent to the solve-the-equation plug-in approach in the bandwidth selection literature.

8 Monte-carlo evidence

We analyze the properties of the rules for selection of the step sizes using the so-called Chernoff's example (see Van der Vaart and Wellner (1996)). The example is based on the population objective function

$$Q(\theta) = E[\mathbf{1}\{x \in [\theta - 1, \theta + 1]\}],$$

with a continuously distributed scalar regressor x . Note that if X has an unbounded support and a continuous density $f_X(\cdot)$ then the maximizer θ_0 of the population objective satisfies the necessary first-order condition

$$f_X(\theta_0 + 1) = f_X(\theta_0 - 1).$$

In particular, if the distribution of X is symmetric and unimodal, then the unique solution is $\theta_0 = 0$. In our analysis we consider the case where X is standard normal, i.e. the population objective function indeed has a unique maximizer. The simulated sample $\{X_i\}_{i=1}^n$ comes from a standard normal distribution and the sample objective is computed as

$$Q_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \in [\theta - 1, \theta + 1]\} = \frac{1}{n} \sum_{i=1}^n U(X_i - \theta),$$

where $U(\cdot)$ is a standard uniform kernel. We consider using numerical gradient to compute the extremum for this objective function. To construct the numerical gradient, we use finite difference formulas of different orders. We use the step of numerical differentiation ε_n which depends on the sample size. In particular, the first-order right derivative formula is

$$\hat{D}_1(\hat{\theta}) = L_{1,1}^{\varepsilon_n} = \frac{\hat{Q}_n(\hat{\theta} + \varepsilon_n) - \hat{Q}_n(\hat{\theta})}{\varepsilon_n},$$

and the left derivative formula is

$$\hat{D}_1(\hat{\theta}) = L_{1,1}^{\varepsilon_n} = \frac{\hat{Q}_n(\hat{\theta}) - \hat{Q}_n(\hat{\theta} - \varepsilon_n)}{\varepsilon_n}.$$

The second-order formula is

$$\hat{D}_2(\hat{\theta}) = L_{1,2}^{\varepsilon_n} = \frac{\hat{Q}_n(\hat{\theta} + \varepsilon_n) - \hat{Q}_n(\hat{\theta} - \varepsilon_n)}{2\varepsilon_n},$$

and the third-order formula is

$$\hat{D}_3(\hat{\theta}) = L_{1,3}^{\varepsilon_n} = \frac{-\hat{Q}_n(\hat{\theta} - 2\varepsilon_n) + 8\hat{Q}_n(\hat{\theta} - \varepsilon_n) - 8\hat{Q}_n(\hat{\theta} + \varepsilon_n) + \hat{Q}_n(\hat{\theta} + 2\varepsilon_n)}{12\varepsilon_n}.$$

The estimator is then re-defined as a solution to the numerical first-order condition

$$\hat{D}_k(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right), \quad (8.10)$$

which is equivalent to the maximum of the empirical objective function that is achieved using a numerical gradient-based maximization routine. We can anticipate the properties of the analyzed estimator by analyzing its behavior analytically. For illustration we can use the numerical derivative formula $\hat{D}_2(\hat{\theta})$. Application of this formula to the sample objective function leads to the expression

$$\hat{D}_2(\hat{\theta}) = \frac{1}{n\varepsilon_n} \sum_{i=1}^n U\left(\frac{X_i - \theta - 1}{\varepsilon_n}\right) - \frac{1}{n\varepsilon_n} \sum_{i=1}^n U\left(\frac{X_i - \theta + 1}{\varepsilon_n}\right).$$

Using the fact that regular distributions of X the uniform kernel will be mean-square differentiable, we can use the following representation:

$$\frac{1}{\sqrt{n\varepsilon_n}} \sum_{i=1}^n U\left(\frac{X_i - \theta - 1}{\varepsilon_n}\right) = \frac{1}{\sqrt{n\varepsilon_n}} \sum_{i=1}^n U\left(\frac{X_i - \theta_0 - 1}{\varepsilon_n}\right) + f'_X(\theta_0 + 1) \sqrt{n\varepsilon_n} (\theta - \theta_0) + o_p(1).$$

Therefore, we can express the estimator, which solves (8.10) as

$$\sqrt{n\varepsilon_n}(\hat{\theta} - \theta_0) = (f'_X(\theta_0 + 1) - f'_X(\theta_0 - 1))^{-1} \frac{1}{\sqrt{n\varepsilon_n}} \sum_{i=1}^n \left[U\left(\frac{X_i - \theta_0 - 1}{\varepsilon_n}\right) - U\left(\frac{X_i - \theta_0 + 1}{\varepsilon_n}\right) \right] + o_p(1).$$

This demonstrates that in case of a “relatively slow” approach of the step size to zero, the properties of this estimator will be described by the case where $\sqrt[3]{n\varepsilon_n} \rightarrow \infty$ ($\gamma = \frac{1}{2}$). Once these conditions are satisfied, then we can use the Lindeberg-Levy CLT to establish that

$$\frac{1}{\sqrt{n\varepsilon_n}} \sum_{i=1}^n \left[U\left(\frac{X_i - \theta_0 - 1}{\varepsilon_n}\right) - U\left(\frac{X_i - \theta_0 + 1}{\varepsilon_n}\right) \right] \xrightarrow{d} \mathcal{N}(0, \Omega).$$

In this we can also use our result regarding the consistency of this estimator. The variance of the estimator can be evaluated similarly to the variance of kernel smoothers. In fact, we can evaluate

$$\begin{aligned} E \left\{ \left(\frac{1}{\sqrt{\varepsilon_n}} \left[U\left(\frac{X_i - \theta_0 - 1}{\varepsilon_n}\right) - U\left(\frac{X_i - \theta_0 + 1}{\varepsilon_n}\right) \right] \right)^2 \right\} \\ = \int U(u) (f_X(\theta_0 + 1 + u\varepsilon_n) + f_X(\theta_0 - 1 + u\varepsilon_n)) du = f_X(\theta_0 + 1) + f_X(\theta_0 - 1) + O(\varepsilon_n). \end{aligned}$$

Then the expression for the variance can be written as

$$V = (f'_X(\theta_0 + 1) - f'_X(\theta_0 - 1))^{-2} (f_X(\theta_0 + 1) + f_X(\theta_0 - 1))$$

If $x \sim \mathcal{N}(0, 1)$ then $V = \sqrt{\frac{\pi\varepsilon}{2}}$. For “fast” approach to zero, i.e. when $(\sqrt[3]{n\varepsilon_n})^{-1} = O(1)$ we should observe the non-standard “cube root” asymptotics.

It is clear that in relatively small samples when the step size of numerical differentiation is “small” the sample first-order condition will have multiple roots. Given the structure of the objective functions the roots will either be contained in the disjoint convex compact sets or will be singletons. To facilitate root finding, we use a dense grid over the state space of the model. For the step size ε_n we

choose the size of the grid cell to be $O(\varepsilon_n / \log n)$. This will assure that the error (measured as the Hausdorff distance between the true set of roots and the set of roots on the grid) will vanish at a faster rate than the numerical error from approximating the gradient using a finite-difference formula. For simplicity we use a uniform grid on $[-1, 1]$ such that the cell size is $\Delta_n = C \frac{\varepsilon_n}{\log n}$, the number of grid points is $N_{\Delta_n} = \left\lceil \frac{2 \log n}{C \varepsilon_n} \right\rceil + 1$ and the grid points can be obtained as $\theta_g = -1 + \Delta(g - 1)$ forming the set $G_{\Delta_n} = \{\theta_g\}_{g=1}^{N_{\Delta_n}}$. The grid search algorithm will identify the set of points

$$Z_n = \left\{ \theta \in G_{\Delta_n} : \left| \widehat{D}_k(\theta) \right| \leq \sqrt{\frac{\log n}{\varepsilon_n n}} \right\}.$$

We call this set the set of roots of the numerical first-order condition on a selected grid. Our Monte-Carlo study will analyze the structure of the set of roots on the grid to evaluate the performance of the numerical gradient-based estimator. The Monte-Carlo study proceeds in the following steps.

1. We generate 1000 Monte-Carlo samples with the number of observations from 500 to 4000. Each simulation sample is indexed by s and the sample size is denoted n_s .
2. We choose sample-adaptive step of numerical differentiation as $\varepsilon = C (n_s)^q$. We choose $C = 2$ and q from 0.2 to 2.
3. Using this step size, we set up the function that we associate with the empirical first-order condition with $\widehat{D}_k(\widehat{\theta}^s)$ for different orders of numerical derivatives.
4. Using the grid over the support $[-1, 1]$ (which we described above) we find all solutions satisfying (8.10). This will form the set of roots on the grid Z_{n_s} .
5. We store all roots on the grid and report the statistics averaged across the roots.
6. If $\#Z_{n_s}$ is the number of roots found in simulation s , we evaluate the mean-squared errors of estimation as:

$$\text{MSE}(\widehat{\theta}) = \sqrt{\frac{1}{S} \sum_{s=1}^S \frac{1}{\#Z_{n_s}} \sum_{r=1}^{\#Z_{n_s}} (\widehat{\theta}_{rs} - \theta_0)^2}$$

We illustrate our results with the set of graphs that show the dependence of the statistics for the set of roots on the sample size.

We emphasize here that the procedure that for correct characterization of the problem (especially in “under-smoothed” cases where the step size approaches to zero too fast) requires a root-finding routine that can correctly characterize the entire set of roots. We argue that an appropriate grid-based procedure will be appropriate because it provides a countable approximation to an uncountable set of roots. The profile of the first-order condition for different orders of numerical derivative formulas are presented in Figures 1-3. The panels from top to bottom show the numerical derivative

evaluated on the grid with a decreasing step size of numerical differentiation. The top panels show the smoothest case where the first-order condition has a pronounced single solution. However, one can see that this solution is biased in all cases. The pictures at the bottom show the profile of the numerical derivative where the step size is selected too small. One can see that there are large regions where the numerical derivative is identically equal to zero. This illustrates the bias-variance trade off in this model: once the step size for numerical differentiation is chosen to be large, then the first-order condition will yield the solution with small variance but potentially large bias. However, when the step size is small, then the solution will have a large variance (over the set of potential roots) but possibly small bias.

We show this intuition formally by analyzing the bias and the variance of the estimates. Figure 4 shows the decay of the root-mean squared error with the sample size over different selection of the sequence of step sizes. The blue line with the asterisk markers represents the over-smoothed case with the slowest approach of the step size of numerical differentiation to zero. With the increase of the rate of this approach, the mean-squared error decreases and reaches its optimum indicated by the line marked with “+”. Then it drops and demonstrates fluctuations indicating the case where to small step size leads to a large set of extraneous roots. Figure 5 shows the same dependence for the bias. One can see that bias tends to decrease with the sample size. Also an interesting observation is that the step size that corresponded to the optimal mean-squared error does not correspond to the minimal bias in the estimation. Moreover, the magnitude of the bias is much smaller than the magnitude of the mean-squared error. This implies that in all cases variance tends to dominate the bias in the absolute value. This shows that even in the non-smooth mods such as the one that we considered, one can use the numerical gradient-based procedure to search for the maximum, once the step size was chosen correctly.

[Figure 1 about here.]

[Figure 2 about here.]

[Figure 3 about here.]

[Figure 4 about here.]

[Figure 5 about here.]

9 Conclusion

In this paper we study the impact that the use of numerical finite-point approximation can make on the estimate of the gradient of a function. We focus on the case where the function is computed from a cross-sectional data sample. We find weak sufficient conditions that allow us to provide uniformly

consistent estimates for the gradients of functions and the directional derivatives of semiparametric moments. Such results may be used to compute the Hessians of sample function to estimate the asymptotic variance or they can be used as inputs in efficient two-step estimation procedures. We further investigate the role of finite-point approximation in computation of extremum estimators. Finite-point approximation formulas use tuning parameters such as step size and in this paper we find that the presence of such parameters may affect statistical properties of the original extremum estimation. We study extremum estimation for classical M-estimators, U-statistics, and semiparametric generalized minimum distance estimators where the optimization routine uses a finite-difference approximation to the numerical gradient. We find that the properties of the estimator obtained from the numerical optimization routine depends on the interference between the smoothness of the population objective function, the precision of approximation, and the smoothness of the sample objective function. While in smooth models the choice of the sequence of step sizes may not affect the properties of the estimator, in non-smooth models the presence of numerical optimization routine can alter both the rate of convergence of the estimator and its asymptotic distribution.

A Appendix

A.1 Proof of Theorem 1

Proof. We need to verify that uniformly in $z \in \mathcal{Z}$, $(\theta, \eta(\cdot)) \in \Theta \times \mathcal{H}$ and $w_j \in \mathcal{H}_n$ the numerical derivative will be converging to the population derivative at $(\theta_0, \eta_0(\cdot))$. We begin with noting that

$$\begin{aligned} \hat{m}(z; \hat{\theta} + e_j \varepsilon_n, \hat{\eta}(\cdot)) &= \hat{m}(z; \hat{\theta} + e_j \varepsilon_n, \hat{\eta}(\cdot)) - m(z; \hat{\theta} + e_j \varepsilon_n, \hat{\eta}(\cdot)) \\ &\quad - \hat{m}(z; \theta_0, \eta_0(\cdot)) + \hat{m}(z; \theta_0, \eta_0(\cdot)) - m(z; \theta_0, \eta_0(\cdot)) \\ &\quad + m(z; \hat{\theta} + e_j \varepsilon_n, \hat{\eta}(\cdot)) - m(z; \theta_0 + e_j \varepsilon_n, \hat{\eta}(\cdot)) \\ &\quad + m(z; \theta_0 + e_j \varepsilon_n, \hat{\eta}(\cdot)) - m(z; \theta_0 + e_j \varepsilon_n, \eta_0(\cdot)) \\ &\quad + m(z; \theta_0 + e_j \varepsilon_n, \eta_0(\cdot)) - m(z; \theta_0, \eta_0(\cdot)) \end{aligned}$$

Using the expansion representation above, we conclude that

$$\begin{aligned} \hat{m}(z; \hat{\theta} + e_j \varepsilon_n, \hat{\eta}(\cdot)) &\stackrel{\text{L}^2}{=} O_p(n^{-1/k}) + \Delta_{1\theta}(\hat{\theta} - \theta_0) + \Delta_{1\eta}[\hat{\eta} - \eta_0] \\ &\quad + (\hat{\theta} - \theta_0)' \Delta_{2\theta^2}(\hat{\theta} - \theta_0) + \Delta_{2\eta^2}[\hat{\eta} - \eta_0]^2 + \Delta_{2\theta\eta}[\hat{\eta} - \eta_0](\hat{\theta} - \theta_0) \\ &\quad + \Delta_{1\theta}^j \varepsilon_n + \Delta_{2\theta^2}^{jj} \varepsilon_n^2 + o_p(\|\hat{\eta} - \eta_0\|_{\text{L}^2}^2) + o_p(\|\hat{\theta} - \theta_0\|^2). \end{aligned}$$

Using a similar technique we can represent

$$\begin{aligned} \hat{m}\left(z; \hat{\theta} - e_j \varepsilon_n, \hat{\eta}(\cdot)\right) &\stackrel{\mathbf{L}^2}{=} O_p\left(n^{-1/k}\right) + \Delta_{1\theta}\left(\hat{\theta} - \theta_0\right) + \Delta_{1\eta}[\hat{\eta} - \eta_0] \\ &+ \left(\hat{\theta} - \theta_0\right)' \Delta_{2\theta^2}\left(\hat{\theta} - \theta_0\right) + \Delta_{2\eta^2}[\hat{\eta} - \eta_0]^2 + \Delta_{2\theta\eta}[\hat{\eta} - \eta_0]\left(\hat{\theta} - \theta_0\right) \\ &- \Delta_{1\theta}^j \varepsilon_n + \Delta_{2\theta^2}^{jj} \varepsilon_n^2 + o_p\left(\varepsilon_n^2\right). \end{aligned}$$

As a result, we can evaluate

$$\begin{aligned} \frac{\hat{m}\left(z; \hat{\theta} + e_j \varepsilon_n, \hat{\eta}(\cdot)\right) - \hat{m}\left(z; \hat{\theta} - e_j \varepsilon_n, \hat{\eta}(\cdot)\right)}{2\varepsilon_n} &\stackrel{\mathbf{L}^2}{=} \Delta_{1\theta}^j + O_p\left(\varepsilon_n^{-1}\left(n^{-1/k} + n^{-1/2} + n^{-1/k_1}\right)\right) \\ &+ o_p\left(\varepsilon_n\right). \end{aligned}$$

Note that $\varepsilon_n n^{1/\max\{k, k_1\}} \rightarrow \infty$ and $\varepsilon_n \rightarrow 0$ will assure uniform convergence of the moment function to its derivative. Next, we provide the result for the uniform convergence of the directional derivative with respect to the infinite-dimensional parameter. Consider a particular direction w_j (which in practice will be an element of the sieve space containing the approximation $\hat{\eta}(\cdot)$), then:

$$\begin{aligned} \hat{m}\left(z; \hat{\theta}, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) &= \hat{m}\left(z; \hat{\theta}, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) - m\left(z; \hat{\theta}, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) \\ &- \hat{m}\left(z; \theta_0, \eta_0(\cdot)\right) + \hat{m}\left(z; \theta_0, h_0(\cdot)\right) - m\left(z; \theta_0, \eta_0(\cdot)\right) \\ &+ m\left(z; \hat{\theta}, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) - m\left(z; \theta_0, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) \\ &+ m\left(z; \theta_0, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) - m\left(z; \theta_0, \eta_0(\cdot) + \tau_n w_j(\cdot)\right) \\ &+ m\left(z; \theta_0, \eta_0(\cdot) + \tau_n w_j(\cdot)\right) - m\left(z; \theta_0, \eta_0(\cdot)\right). \end{aligned}$$

Using the local \mathbf{L}^2 -representation, we can approximate the expansion above as

$$\begin{aligned} \hat{m}\left(z; \hat{\theta}, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) &\stackrel{\mathbf{L}^2}{=} O_p\left(n^{-1/k}\right) + \Delta_{1\theta}\left(\hat{\theta} - \theta_0\right) + \Delta_{1\eta}[\hat{\eta} - \eta_0] \\ &+ \left(\hat{\theta} - \theta_0\right)' \Delta_{2\theta^2}\left(\hat{\theta} - \theta_0\right) + \Delta_{2\eta^2}[\hat{\eta} - \eta_0]^2 + \tau_n \Delta_{2h^2}[\hat{\eta} - \eta_0, w_j] \\ &+ \Delta_{2\theta\eta}[\hat{\eta} - \eta_0]\left(\hat{\theta} - \theta_0\right) + \tau_n \Delta_{2\theta\eta}[w_j]\left(\hat{\theta} - \theta_0\right) \\ &+ \tau_n \Delta_{1\eta}[w_j] + \tau_n^2 \Delta_{2\eta^2}[w_j]^2 + o_p\left(\|\hat{\eta} - \eta_0\|_{\mathbf{L}^2}^2\right) + o_p\left(\|\hat{\theta} - \theta_0\|^2\right). \end{aligned}$$

We can write a similar expression for $\hat{m}\left(z; \hat{\theta}, \hat{\eta}(\cdot) - \tau_n w_j(\cdot)\right)$. As a result, the symmetrized numerical directional derivative will be approximated locally by

$$\begin{aligned} \frac{\hat{m}\left(z; \hat{\theta}, \hat{\eta}(\cdot) + \tau_n w_j(\cdot)\right) - \hat{m}\left(z; \hat{\theta}, \hat{\eta}(\cdot) - \tau_n w_j(\cdot)\right)}{2\tau_n} &\stackrel{\mathbf{L}^2}{=} \Delta_{1\eta}[w_j] + \Delta_{2\eta^2}[\hat{\eta} - \eta_0, w_j] \\ &+ \Delta_{2\theta\eta}[w_j]\left(\hat{\theta} - \theta_0\right) + O_p\left(\tau_n^{-1}\left(n^{-1/k} + n^{-1/2} + n^{-1/k_1}\right)\right) + o_p\left(\tau_n\right). \end{aligned}$$

Note that $\|\Delta_{2\eta^2}[\hat{\eta} - \eta_0, w_j]\|_{\mathbf{L}^2} = O_p\left(n^{-1/k_1}\right)$. For $k_1 > 2$ this term will dominate and determine the lower bound on the sub-parametric convergence rate for the numerical derivative. The conditions for τ_n will

be similar to the conditions for ε_n , that is $\tau_n n^{1/\max\{k, k_1\}} \rightarrow \infty$ and $\tau_n \rightarrow 0$. Moreover, it is clear that the convergence rate for τ_n is slower than n^{-1/k_1} . This result assures that for $z \in \mathcal{Z}$, $(\theta, \eta(\cdot)) \in \Theta \times \mathcal{H}$, and $w_j \in \mathcal{W}$, $\tilde{D}_{w_j}(z) \xrightarrow{p} D_{w_j}(z)$. \square

A.2 Proof of Lemma 2

It follows directly from Assumption 7.[iii] that for $\eta \in \mathcal{H}_n$

$$\left| L_{1,p}^{\varepsilon_n, w} m(\theta, \eta, z) - \text{proj} \left(L_{1,p}^{\varepsilon_n, w} m(\theta, \eta, z) | p^N(z) \right) \right| = O \left(\frac{1}{N^\alpha \varepsilon_n} + \frac{1}{n^\phi \varepsilon_n} \right),$$

that will converge to zero if $\min N^\alpha, n^\phi \varepsilon_n \rightarrow \infty$.

Therefore it suffices to prove Lemma 2 for

$$(*) = \left| L_{1,p}^{\varepsilon_n, w} \hat{m}(\theta, \eta, z) - \text{proj} \left(L_{1,p}^{\varepsilon_n, w} m(\theta, \eta, z) | p^N(z) \right) \right|.$$

As demonstrated in Newey (1997), for $P = (p^N(z_1), \dots, p^N(z_n))'$ and $\hat{Q} = P'P/n$

$$\|\hat{Q} - Q\| = O_p \left(\sqrt{\frac{N}{n}} \right), \quad \text{where his } \zeta_0(N) = C,$$

and Q is non-singular by Assumption 7.[i] with the smallest eigenvalue bounded from below by some constant $\underline{\lambda} > 0$. Hence the smallest eigenvalue of \hat{Q} will converge to $\underline{\lambda} > 0$. Following Newey (1997) we use the indicator 1_n to indicate the cases where the smallest eigenvalue of \hat{Q} is above $\frac{1}{2}$ to avoid singularities. Introduce the vector $\Delta(\theta, \eta, w, Y_i; \varepsilon_n) = (\rho(\theta, \eta + \varepsilon_n w; Y_i) - \rho(\theta, \eta - \varepsilon_n w; Y_i))_{i=1}^n$. We consider conditional expectation $E[\Delta(\theta, \eta, w, Y_i; \varepsilon_n) | Z_i = z]$ as a function of z (given θ, η, w , and ε_n). We can project this function of z on N basis vectors of the sieve space. Let β be the vector of coefficients of this projection. Also define $G(\theta, \eta, w, \varepsilon_n) = (E[\Delta(\theta, \eta, w, Y_i; \varepsilon_n) | Z_i])_{i=1}^n$. Then $(*)$ equals to a linear combination of $1_n |p^{N'}(z) (\hat{\beta} - \beta)| / \varepsilon_n$. Note that

$$p^{N'}(z) (\hat{\beta} - \beta) = p^{N'}(z) (\hat{Q}^{-1} P' (\Delta - G) / n + \hat{Q}^{-1} P' (G - P\beta) / n). \quad (\text{A.11})$$

For the first term in (A.11), we can use the result that smallest eigenvalue of \hat{Q} is converging to $\underline{\lambda} > 0$. Then application of the Cauchy-Schwartz inequality leads to

$$\left| p^{N'}(z) \hat{Q}^{-1} P' (\Delta - G) \right| \leq \|Q^{-1} p^N(z)\| \|P' (\Delta - G)\|.$$

Then $\|\hat{Q}^{-1} p^N(z)\| \leq \frac{C}{\underline{\lambda}} \sqrt{N}$, and

$$\begin{aligned} \|P' (\Delta - G)\| &= \sqrt{\sum_{k=1}^N \left(\sum_{i=1}^n p_{Nk}(z_i) (\Delta(\theta, \eta, w, Y_i; \varepsilon_n) - G(\theta, \eta, w, Z_i; \varepsilon_n)) \right)^2} \\ &\leq \sqrt{N} \max_k \left| \sum_{i=1}^n p_{Nk}(z_i) (\Delta(\theta, \eta, w, Y_i; \varepsilon_n) - G(\theta, \eta, w, Z_i; \varepsilon_n)) \right| \end{aligned}$$

Thus,

$$\left| p^{N'}(z) \hat{Q}^{-1} P' (\Delta - G) \right| \leq \frac{CN}{\underline{\lambda}} \max_k \left| \sum_{i=1}^n p_{Nk}(z_i) (\Delta(\theta, \eta, w, Y_i; \varepsilon_n) - G(\theta, \eta, w, Z_i; \varepsilon_n)) \right|.$$

Denote $\mu_n = \mu \frac{\epsilon_n}{N}$. Next we adapt the arguments for proving Theorem 37 in Pollard (1984) to provide the bound for $P\left(\sup_{\mathcal{F}_n} \frac{1}{n} \|p^{N'}(z) \hat{Q}^{-1} P'(\Delta - G)\| > N\mu_n\right)$. For N non-negative random variables Y_i we note that

$$P\left(\max_i Y_i > c\right) \leq \sum_{i=1}^N P(Y_i > c).$$

Using this observation, we can find that

$$P\left(\sup_{\mathcal{F}_n} \frac{1}{n} \|p^{N'}(z) \hat{Q}^{-1} P'(\Delta - G)\| > N\mu_n\right) \leq \sum_{k=1}^N P\left(\sup_{\mathcal{F}_n} \left\| \frac{1}{n} \sum_{i=1}^n p_{Nk}(z_i) (\Delta_i - G_i) \right\| > \mu_n\right)$$

This inequality allows us to substitute the tail bound for the class of functions $L_{1,p}^{\epsilon_n, w} \hat{m}(\theta, \eta; z)$ that is indexed by θ, η, w and z by a tail bound for a much simpler class

$$\mathcal{P}_n = \{p_{Nk}(\cdot) (\Delta(\theta, \eta, w, \cdot; \epsilon_n) - G(\theta, \eta, w; \epsilon_n)) : \theta \in N(\theta_0), \eta, w \in \mathcal{H}_n\}.$$

We note that, according to Lemma 2.6.18 in Van der Vaart and Wellner (1996), provided that each $p_{Nk}(\cdot)$ is a fixed function, the covering number for \mathcal{P}_n has the same order as the covering number for \mathcal{F}_n . Then we pick A to be the largest constant for the covering numbers $A_k n^{2r_0} \log\left(\frac{1}{\delta}\right)$ over classes \mathcal{P}_n . By Assumption 7.[i] and 8.[i] any $f \in \mathcal{P}_n$ is bounded $|f| < C < \infty$. Next we note that $\text{Var}(f) = O(\epsilon_n)$ for $f \in \mathcal{P}_n$ by Assumption 8.[ii]. The symmetrization inequality (30) in Pollard (1984) holds if $\epsilon_n / (16n\mu_n^2) \leq \frac{1}{2}$. This will occur if $\frac{n\epsilon_n}{N^2} \rightarrow \infty$. Provided that the symmetrization inequality holds, we can follow the steps of Theorem 37 in Pollard (1984) to establish the tail bound on the sample sum via a combination of the Hoeffding inequality and the covering number for the class \mathcal{P}_n . As a result, we obtain that

$$\begin{aligned} & P\left(\sup_{\mathcal{F}_n} \frac{1}{n} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta_i - G_i) \right\| > 8\mu_n\right) \\ & \leq 2 \exp\left(An^{2r_0} \log \frac{1}{\mu_n}\right) \exp\left(-\frac{1}{128} \frac{n\mu_n^2}{\epsilon_n}\right) + P\left(\sup_{\mathcal{F}_n} \frac{1}{n} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta_i - G_i) \right\|^2 > 64\epsilon_n\right). \end{aligned}$$

The second term can be evaluated with the aid of Lemma 33 in Pollard (1984):

$$P\left(\sup_{\mathcal{F}_n} \frac{1}{n} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta_i - G_i) \right\|^2 > 64\epsilon_n\right) \leq 4 \exp\left(An^{2r_0} \log \frac{1}{\epsilon_n}\right) \exp(-n\epsilon_n).$$

As a result, we find that

$$\begin{aligned} P\left(\sup_{\mathcal{F}_n} \frac{1}{n} \|p^{N'}(z) \hat{Q}^{-1} P'(\Delta - G)\| > N\mu_n\right) & \leq 2N \exp\left(An^{2r_0} \log \frac{1}{\mu_n} - \frac{1}{128} \frac{n\mu_n^2}{\epsilon_n}\right) \\ & \quad + 4N \exp\left(An^{2r_0} \log \frac{1}{\epsilon_n} - n\epsilon_n\right) \end{aligned}$$

We start the analysis with the first term. Consider the case with and $r_0 > 0$. Then the log of the first term takes the form

$$\begin{aligned} & An^{2r_0} \log(N/(\mu\epsilon_n)) - \frac{1}{128} \frac{n}{N^2} \mu^2 \epsilon_n + \log N \\ & = An^{2r_0} \log\left(\frac{N^2 n^{2r_0}}{\mu n \epsilon_n}\right) - \frac{1}{128} \frac{\mu^2 \epsilon_n n}{N^2} - An^{2r_0} \log \frac{N n^{2r_0}}{\mu n} + \log N. \end{aligned}$$

If $N \log n/n \rightarrow 0$, then one needs that $\frac{n\epsilon_n}{N^2 n^{2r_0} \log n} \rightarrow \infty$ if $r_0 > 0$ and $\frac{n\epsilon_n}{N^2 \log^2 n} \rightarrow \infty$ if $r_0 = 0^+$. Hence the first term is of $o(1)$. Now consider the second term. The exponent can be represented as

$$-n\epsilon_n + n^{2r_0} \log \frac{1}{\epsilon_n} + \log N,$$

which is guaranteed to converge to $-\infty$ if $\frac{n\epsilon_n}{N^2 n^{2r_0} \log n} \rightarrow \infty$.

We notice that in this proof uniformity of the directional derivative of the conditional moment in z follows directly from the boundedness of the sieve functions. This implies that in some cases this result can be weakened. \square

A.3 Proof of Lemma 3

Recall the definition of the kernel estimator

$$\hat{m}(\theta, \eta, z) = \left(\frac{1}{nb_n^{d_z}} \sum_{i=1}^n K\left(\frac{z - z_i}{b_n}\right) \right)^{-1} \frac{1}{nb_n^{d_z}} \sum_{i=1}^n \rho(\theta, \eta, z_i) K\left(\frac{z - z_i}{b_n}\right)$$

For the expression of interest, we can consider

$$\begin{aligned} \frac{\hat{m}(\theta, \eta + \epsilon_n w, z) - \hat{m}(\theta, \eta - \epsilon_n w, z)}{\epsilon_n} &= \left(\frac{1}{nb_n^{d_z}} \sum_{i=1}^n K\left(\frac{z - z_i}{b_n}\right) \right)^{-1} \\ &\times \frac{1}{nb_n^{d_z} \epsilon_n} \sum_{i=1}^n [\rho(\theta, \eta + \epsilon_n w, y_i) - \rho(\theta, \eta - \epsilon_n w, y_i)] K\left(\frac{z - z_i}{b_n}\right). \end{aligned}$$

Then we can consider a class of functions

$$\mathcal{G}_n = \{[\rho(\theta, \eta + \epsilon_n w, \cdot) - \rho(\theta, \eta - \epsilon_n w, \cdot)] K\left(\frac{z - \cdot}{b_n}\right), \theta \in \Theta, w, \eta \in \mathcal{H}_n, z \in \mathcal{Z}\}.$$

Consider the classe \mathcal{G}_n . We can represent it as

$$\mathcal{G}_n = \{g = f\kappa : f \in \mathcal{F}_n, \kappa \in \mathcal{F}\}.$$

$N_1(\cdot)$ and $N_2(\cdot)$ correspond to the L_1 and L_2 covering numbers. Consider the covering numbers for classes \mathcal{F}_n and \mathcal{F} . We select $\epsilon > 0$, then there exist $m_1 < N_1(\epsilon, \mathcal{F}_n, L_1(Q))$ and $m_2 < N_1(\epsilon, \mathcal{F}, L_1(Q))$ and covers $\{f_j\}_{j=1}^{m_1}$ and $\{\kappa_i\}_{i=1}^{m_2}$ such that for $f \in \mathcal{F}_n$ and $\kappa \in \mathcal{F}$ $\min_j Q|f - f_j| < \epsilon$ and $\min_i Q|\kappa - \kappa_i| < \epsilon$. We note that $|f| \leq C$ and $|g| \leq C$. Consider the cover $\{f_j \kappa_i\}_{j,i=1}^{j=m_1, i=m_2}$ noting that $f_j \kappa_i - f\kappa = (f_j - f)(\kappa_i - \kappa) + f(\kappa_i - \kappa) + \kappa(f_j - f)$. Then, in combination with Cauchy-Schwartz we have that

$$\min_{i,j} Q|\kappa_i f_j - \kappa f| \leq \min_j (Q|f_j - f|^2)^{1/2} \min_i (Q|\kappa_i - \kappa|^2)^{1/2} + C \min_j Q|f_j - f| + C \min_i Q|\kappa_i - \kappa|$$

Given the relationship between L_1 and L_2 covering numbers covers $\{f_j\}_{j=1}^{m_1}$ and $\{\kappa_i\}_{i=1}^{m_2}$ are sufficient to achieve $\min_j (Q|f_j - f|^2)^{1/2} < \epsilon$ and $\min_i (Q|\kappa_i - \kappa|^2)^{1/2} < \epsilon$. This means that $\min_{i,j} Q|\kappa_i f_j - \kappa f| < 3C\epsilon$. Thus, the L_1 covering number for \mathcal{G}_n is bounded by a product of L_2 covering numbers for \mathcal{F} and \mathcal{F}_n (which corresponds to the number of elements in the cover $\{f_j \kappa_i\}_{j,i=1}^{j=m_1, i=m_2}$).

Provided that classes \mathcal{F}_n and \mathcal{F} satisfy Euclidean property, we can apply Lemma 2.6.20 from Van der Vaart and Wellner (1996). This means that the class \mathcal{G}_n is Euclidean with parameters depending on n . Provided

that $\text{Var}(g) = O(\epsilon_n b_n)$ for $g \in \mathcal{G}_n$, we can use a similar logic as in the proof of Theorem 37 in Pollard (1984) with the results similar to those in the proof of Lemma 2. This leads to condition $\frac{n\epsilon_n b_n^{dz}}{n^{2r_0} \log n} \rightarrow \infty$. We note that the bias due to kernel smoothing $E[\hat{m}(\theta, \eta, Z_i) | Z_i = z] = O(b_n^m)$, where m is the order of the kernel, and the bias due to the sieve approximation is $n^{-\phi}$. Then

$$\|L_{1,p}^{\epsilon_n, w} E[\hat{m}(\theta, \eta, Z_i) | Z_i = z] - L_{1,p}^{\epsilon_n, w} m(\theta, \eta, z)\| = O\left(b_n^m \epsilon_n^{-1} + n^\phi \epsilon_n^{-1}\right),$$

which converges to zero if $\epsilon_n b_n^{-m} \rightarrow \infty$ and $\epsilon_n n^\phi \rightarrow \infty$. \square

A.4 Proof of Lemma 4

As before, we use $\Delta = (\rho(\theta, \eta + \epsilon_n h; Y_i) - \rho(\theta, \eta - \epsilon_n h; Y_i))_{i=1}^n$ and $G = (E[\Delta | Z_i])_{i=1}^n$. Then we note that $\|p_{Nk}(Z_i)(\Delta_i - G_i)\| \leq C\epsilon_n^\gamma$ due to Assumptions 7 and 8(i). Consider the class of functions

$$\mathcal{P}_n = \{p_{Nk}(\cdot)(\Delta(\theta, \eta, w, z, \cdot; \epsilon_n) - G(\theta, \eta, w, \cdot; \epsilon_n)), \theta \in N(\theta_0), \eta, w \in \mathcal{H}_n, z \in \mathcal{Z}\}.$$

We note that for this class we can find an envelope proportional to ϵ_n^γ . Then for $\theta_n = \epsilon_n^{-\gamma} \sup_{f \in \mathcal{P}_n} |P_n f^2|^{1/2}$ we can use the tail bound in Alexander (1984) which is analogous to Theorem 2.14.1 in Van der Vaart and Wellner (1996):

$$P \sup_{f \in \mathcal{P}_n} |\sqrt{n}(P_n f - P f)| \leq C\epsilon_n^\gamma P[J(\theta_n, \mathcal{P}_n)],$$

where $J(\theta, \mathcal{P}_n) = \sup_Q \int_0^\theta (\log N(\delta\epsilon_n^\gamma, \mathcal{P}_n, L_1(Q)))^{1/2} d\delta$. In our case $J(\theta, \mathcal{P}_n) = O\left(n^{r_0} \theta \sqrt{\log(1/\theta)}\right)$. This function achieves its maximum at $\theta = e^{-1/2}$. We can analyze the probability of large deviations for θ_n . From Lemma 33 in Pollard (1984) combined with Lemma 10 in Nolan and Pollard (1987), it follows that there exists β such that

$$P\left(\epsilon_n^{-\gamma} \sup_{f \in \mathcal{P}_n} |P_n f^2|^{1/2} > \beta t\right) \leq 4 \exp\left(An^{2r_0} \log\left(\frac{1}{t^2}\right) - n\beta^2 t^2\right).$$

Consider the sequence $t_n = o(1)$ such that $t_n \gg \sqrt{\frac{\log n^{1-2r_0}}{n^{1-2r_0}}}$. We note that for sufficiently large n $\beta t_n \epsilon_n^\gamma \ll \frac{1}{e}$. This means that

$$\begin{aligned} P[J(\theta_n, \mathcal{P}_n)] &= P[J(\theta_n, \mathcal{P}_n) \mathbf{1}(\theta_n > \beta t_n) + J(\theta_n, \mathcal{P}_n) \mathbf{1}(\theta_n \leq \beta t_n)] \\ &= n^{r_0} \frac{1}{\sqrt{2e}} P\left(\sup_{f \in \mathcal{P}_n} |P_n f^2|^{1/2} > \beta t_n \epsilon_n^\gamma\right) + n^{2r_0} \beta t_n \sqrt{\log \frac{1}{\beta t_n}}. \end{aligned}$$

Thus, we have established that

$$P \sup_{f \in \mathcal{P}_n} |\sqrt{n}(P_n f - P f)| \leq C_1 \epsilon_n^\gamma n^{r_0} \exp\left(An^{2r_0} \log \frac{1}{t_n} - n\beta^2 t_n^2 \epsilon_n^{2\gamma}\right) + C_2 \epsilon_n^\gamma n^{r_0} \beta t_n \sqrt{\log \frac{1}{\beta t_n}}$$

Then, returning to our previous argument, we recognize that

$$\begin{aligned} &P \frac{1}{n\epsilon_n} \sup_{\mathcal{F}_n} \left\| \sum_{i=1}^n p_{Nk}(x_i) (\Delta_i - G_i) \right\| \\ &= O\left(n^{r_0 - \frac{1}{2}} \epsilon_n^{\gamma-1} \exp\left(An^{2r_0} \log \frac{1}{t_n \epsilon_n^\gamma} - n\beta^2 t_n^2 \epsilon_n^{2\gamma}\right) + n^{r_0 - \frac{1}{2}} \epsilon_n^{\gamma-1} \beta t_n \sqrt{\log \frac{1}{\beta t_n}}\right). \end{aligned}$$

As a result, we established that

$$P \sup_{\mathcal{F}_n} \frac{1}{n\epsilon_n} \|p^{N'}(z)\hat{Q}^{-1}P'(\Delta - G)\| = O\left(Nn^{r_0-\frac{1}{2}}\epsilon_n^{\gamma-1} \exp\left(An^{2r_0} \log \frac{1}{t_n\epsilon_n^\gamma} - n\beta^2 t_n^2 \epsilon_n^{2\gamma}\right) + Nn^{r_0-\frac{1}{2}}\epsilon_n^{\gamma-1}\beta t_n \sqrt{\log \frac{1}{\beta t_n}}\right).$$

Next, we note that provided that $\frac{\sqrt{n}\epsilon_n^{1-\gamma}}{Nn^{2r_0}} \rightarrow \infty$, the expression on the right-hand side converges to zero. In fact, the first term is exponential. The second term has a factor that can be transformed into $\sqrt{2} \sqrt{\frac{\log \frac{1}{\beta^2 t_n^2}}{\frac{1}{\beta^2 t_n^2}}} = o(1)$, if $t_n = o(1)$.

Consider now similar conditions for kernel-based estimators. As before, we can represent the finite-difference formula for the directional derivative of interest as

$$\begin{aligned} \hat{m}(\theta, \eta + \epsilon_n w, z) - \hat{m}(\theta, \eta - \epsilon_n w, z) &= \left(\frac{1}{nb_n^{d_z}} \sum_{i=1}^n K\left(\frac{z - z_i}{b_n}\right) \right)^{-1} \\ &\times \frac{1}{nb_n^{d_z}} \sum_{i=1}^n [\rho(\theta, \eta + \epsilon_n w, z_i) - \rho(\theta, \eta - \epsilon_n w, z_i)] K\left(\frac{z - z_i}{b_n}\right). \end{aligned}$$

We consider the class

$$\mathcal{G}_n = \{[\rho(\theta, \eta + \epsilon_n w, \cdot) - \rho(\theta, \eta - \epsilon_n w, \cdot)] K\left(\frac{z - \cdot}{b_n}\right), \theta \in \Theta, w, \eta \in \mathcal{H}_n, z \in \mathcal{Z}\}.$$

We noted before that the L_1 covering number for this class can be constructed as a product of covering numbers for the classes forming the product. Provided that due to Hölder-continuity there exists an envelope ϵ_n for this class, we can represent

$$\log N_1(\epsilon \|F\|, \mathcal{G}_n, L_1) \leq An^{2r_0} \log\left(\frac{1}{\epsilon}\right),$$

for sufficiently large A . This makes our entire previous discussion to be analogous to the kernel case. The only difference arises in the exponential inequality. We notice that $\text{Var}(f) = O(\epsilon_n^{2\gamma} b_n^{d_z})$ for $f \in \mathcal{G}_n$. This means that the exponential inequality can be re-written as

$$P\left(\epsilon_n^{-\gamma} \sup_{f \in \mathcal{P}_n} |P_n f^2|^{1/2} > \beta t\right) \leq C \exp\left(An^{2r_0} \log \frac{1}{t} - n\beta^2 t^2\right).$$

We choose the sequence $t_n = O(c_n b_n^{d_z/2})$ for some arbitrary $c_n \rightarrow 0$, meaning that the decomposition into “large” and “small” θ_n will rely on the threshold $\beta c_n b_n^{d_z/2}$ as opposed to βt_n . Thus

$$\begin{aligned} P[J(\theta_n, \mathcal{P}_n)] &= P\left[J(\theta_n, \mathcal{P}_n) \mathbf{1}(\theta_n > \beta c_n b_n^{d_z/2}) + J(\theta_n, \mathcal{P}_n) \mathbf{1}(\theta_n \leq \beta c_n b_n^{d_z/2})\right] \\ &= n^{r_0} \frac{1}{\sqrt{2e}} P\left(\sup_{f \in \mathcal{P}_n} |P_n f^2|^{1/2} > \beta c_n b_n^{d_z/2} \epsilon_n^\gamma\right) + n^{r_0} \beta c_n b_n^{d_z/2} \sqrt{\log \frac{1}{\beta c_n b_n^{d_z/2}}}. \end{aligned}$$

Then decomposing the tail bound as before, we obtain

$$\begin{aligned} P \sup_{\mathcal{F}_n} \left\| \frac{1}{n\epsilon_n b_n^{d_z}} \sum_{i=1}^n [\rho(\theta, \eta + \epsilon_n h, z_i) - \rho(\theta, \eta - \epsilon_n h, z_i)] K\left(\frac{z - z_i}{b_n}\right) \right\| \\ = O\left(n^{r_0-\frac{1}{2}} b_n^{-d_z} \epsilon_n^{\gamma-1} \exp\left(An^{2r_0} \log \frac{1}{c_n b_n^{d_z/2}} - n\beta^2 c_n^2 b_n^{d_z} \epsilon_n^{2\gamma}\right) + n^{r_0-\frac{1}{2}} \epsilon_n^{\gamma-1} \beta c_n b_n^{-d_z/2} \sqrt{\log \frac{1}{\beta c_n b_n^{d_z/2}}}\right). \end{aligned}$$

We notice that $c_n = o(1)$ and $n^{2r_0-1} \epsilon_n^{1-\gamma} b_n^{d_z/2} / \log n^{2r_0-1} \rightarrow \infty$ assures the convergence to zero. \square

A.5 Proof of Lemma 5

Proof. (i)

To find the convergence rate for the estimator $\hat{\theta}$ we need to find the “balancing” sequence ρ_n that assures that $\rho_n d(\hat{\theta}, \theta_0) = O_P^*(1)$. Using the assumption of the compactness of the parameter space, we cover it using a grid with cells $S_{j,n} = \{\theta : 2^{j-1} < \rho_n d(\theta, \theta_0) < 2^j\}$. The idea of the proof is to show that for any given $\kappa > 0$ we can find finite β such that the probability of event $\rho_n d(\hat{\theta}, \theta_0) > \beta$ is below κ . Using a partition of the parameter space, we can pick a finite integer M such that $2^M < \beta$. Then we evaluate the probability of a large deviation $\rho_n d(\hat{\theta}, \theta_0) > 2^M$. We know that the estimator solves

$$\sqrt{n\varepsilon_n} L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) = o_p(1).$$

If $\rho_n d(\hat{\theta}, \theta_0)$ is larger than 2^M for a given M , then over the θ in one of the cells $S_{j,n}$, $\sqrt{n\varepsilon_n} L_{1,p}^{\varepsilon_n} Q_n(\theta)$ achieves a distance as close as desired to zero. Hence, for every $\delta > 0$,

$$P(\rho_n d(\hat{\theta}, \theta_0) > 2^M) \leq \sum_{\substack{j \geq M \\ 2^j < \delta \rho_n}} P\left(\sup_{\theta \in S_{j,n}} \left(-\|L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\|\right) \geq -o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)\right) + P(2d(\hat{\theta}, \theta_0) \geq \delta)$$

Then we evaluate the population objective, using the fact that it has p mean-square derivatives with Taylor residual of order ν :

$$\|L_{1,p}^{\varepsilon_n} Q(\theta)\| \geq C d(\theta, \theta_0) + C' \varepsilon_n^{\nu-1},$$

where θ_0 is the zero of the population first-order condition and the approximated derivative has a known order of approximation $\|L_{1,p}^{\varepsilon_n} Q(\theta_0)\| = C' \varepsilon_n^{\nu-1}$ for some constant C' . Substitution of this expression into the argument of interest leads to

$$\|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\| \geq \|L_{1,p}^{\varepsilon_n} Q(\theta)\| - \|L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\|.$$

Therefore

$$\sup_{\theta \in S_{j,n}} \left(-\|L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\|\right) \geq \sup_{\theta \in S_{j,n}} \|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\| - \sup_{\theta \in S_{j,n}} \|L_{1,p}^{\varepsilon_n} Q(\theta)\|$$

Then applying the Markov inequality to the re-centered process for $\theta \in S_{j,n}$

$$\begin{aligned} P\left(\sup_{\theta \in S_{j,n}} \|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\| \geq C d(\theta, \theta_0) + C' \varepsilon_n^{\nu-1} + o\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)\right) \\ \leq \frac{E^* \left[\sup_{\theta \in S_{j,n}} \|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\| \right]}{C d(\theta, \theta_0) + C' \varepsilon_n^{\nu-1} + o\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)} = O\left(\rho_n \sqrt{\frac{\log n}{n\varepsilon_n}}\right) 2^{-(j-1)}. \end{aligned}$$

Thus if $\rho_n = o\left(\sqrt{\frac{n\varepsilon_n}{\log n}}\right)$ then the probability of interest is $o(1)$.

(ii)

Consider a class of functions

$$\mathcal{G}_n = \left\{ g(\cdot, \theta_n + \varepsilon_n) - g(\cdot, \theta_n - \varepsilon_n) - g(\cdot, \theta_0 + \varepsilon_n) + g(\cdot, \theta_0 - \varepsilon_n), \theta_n = \theta_0 + t_n \sqrt{\frac{\log n}{n\varepsilon_n}} \right\},$$

with $\varepsilon_n \rightarrow 0$ and $t_n = O(1)$. We can evaluate the L^2 norm of the functions from class \mathcal{G}_n using Assumption 5 (ii). Note that

$$E \left[(g(Z_i, \theta_n + \varepsilon_n) - g(Z_i, \theta_n - \varepsilon_n))^2 \right] = O(\varepsilon_n),$$

with the same evaluation for the second term. On the other hand, we can change the notation to $\theta_{1n} = \theta_0 + \varepsilon_n + \frac{t_n}{2} \sqrt{\frac{\log n}{n\varepsilon_n}}$ and $\theta_{1n} = \theta_0 - \varepsilon_n + \frac{t_n}{2} \sqrt{\frac{\log n}{n\varepsilon_n}}$. Then we can group the first term with the third and the second one with the fourth. For the first group this leads to

$$E \left[\left(g \left(Z_i, \theta_{1n} + \frac{t_n}{2} \sqrt{\frac{\log n}{n\varepsilon_n}} \right) - g \left(Z_i, \theta_{1n} - \frac{t_n}{2} \sqrt{\frac{\log n}{n\varepsilon_n}} \right) \right)^2 \right] = O \left(\sqrt{\frac{\log n}{n\varepsilon_n}} \right),$$

and for the second group

$$E \left[\left(g \left(Z_i, \theta_{2n} + \frac{t_n}{2} \sqrt{\frac{\log n}{n\varepsilon_n}} \right) - g \left(Z_i, \theta_{2n} - \frac{t_n}{2} \sqrt{\frac{\log n}{n\varepsilon_n}} \right) \right)^2 \right] = O \left(\sqrt{\frac{\log n}{n\varepsilon_n}} \right).$$

Thus, two different ways of grouping the terms allow us to obtain two possible bounds on the norm of the entire term. As a result, we find that

$$P f^2 = O \left(\min \left\{ \varepsilon_n, \sqrt{\frac{\log n}{n\varepsilon_n}} \right\} \right), \quad f \in \mathcal{G}_n.$$

Next we denote $\delta_n = \min \left\{ \varepsilon_n, \sqrt{\frac{\log n}{n\varepsilon_n}} \right\}$. Invoking Lemma 33 in Pollard (1984) and using Assumption 5 (iii) we obtain that

$$P \left(\sup_{\mathcal{G}_n} P_n f^2 > 64\delta_n \right) \leq 4A (\delta_n)^{-V} \exp(-n\delta_n) = 4A \exp \left(-n\delta_n + V \log \left(\frac{1}{\delta_n} \right) \right).$$

If $n\varepsilon_n^3 / \log n \rightarrow \infty$, then this allows us to conclude that $\sup_{\mathcal{G}_n} P_n f^2 = o_p(\delta_n)$. Next we can apply the maximum inequality from Theorem 2.14.1 in Van der Vaart and Wellner (1996) which implies that for the functions with constant envelopes

$$\sqrt{\frac{n}{\varepsilon_n}} \sup_{\mathcal{G}_n} |P_n f - P f| \lesssim \frac{1}{\sqrt{\varepsilon_n}} J \left(\sup_{\mathcal{G}_n} P_n f^2, \mathcal{G}_n \right),$$

where $J(\cdot)$ is a covering integral:

$$J(\delta, \mathcal{F}) = \sup_Q \int_0^\delta \sqrt{1 + N(\epsilon, \mathcal{F}, \mathbf{L}^2(Q))} d\epsilon.$$

For Euclidean class as in Assumption 5 (iii) we can evaluate $J(\delta, \mathcal{G}_n) = O \left(\delta \sqrt{\log \left(\frac{1}{\delta} \right)} \right)$. Using the expression for $\sup_{\mathcal{G}_n} P f^2$, we can evaluate

$$\sqrt{\frac{n}{\varepsilon_n}} \sup_{\mathcal{G}_n} |P_n f - P f| = o_p \left(\frac{\delta_n}{\sqrt{\varepsilon_n}} \sqrt{\log \left(\frac{1}{\delta_n} \right)} \right).$$

Then, provided that $\frac{n\varepsilon_n^3}{\log n} \rightarrow \infty$, we see that $\delta_n = \sqrt{\log n / (n\varepsilon_n)}$ and

$$\sqrt{\frac{n}{\varepsilon_n}} \sup_{\mathcal{G}_n} |P_n f - P f| = o_p \left(\sqrt{\frac{\log \left(\frac{n\varepsilon_n}{\log n} \right)}{\frac{n}{\log n}}} \right) = o_p(1).$$

The statement of the Lemma follows directly from this result. \square

A.6 Proof of theorem 8

Proof. This proof will replicate the steps of proof of Lemma 5. We perform the triangulation of the parameter space according to the balancing rate ρ_n into segments $S_{j,n} = \{\theta : 2^{j-1} < \rho_n d(\theta, \theta_0) < 2^j\}$. For every $\delta > 0$,

$$P\left(\rho_n d(\hat{\theta}, \theta_0) > 2^M\right) \leq \sum_{\substack{j \geq M \\ 2^j < \delta \rho_n}} P\left(\sup_{\theta \in S_{j,n}} \left(-\left\|L_{1,p}^{\varepsilon_n} \hat{Q}(\theta)\right\|\right) \geq -o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right)\right) + P\left(2d(\hat{\theta}, \theta_0) \geq \delta\right)$$

We can then follow the steps of Lemma 5 to evaluate the upper bound for the elements in the sum on the right-hand side as

$$\begin{aligned} & P\left(\sup_{\theta \in S_{j,n}} \left\|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta) - L_{1,p}^{\varepsilon_n} Q(\theta_0) + L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0)\right\| \geq O\left(\frac{2^j}{\rho_n}\right)\right) \\ & \leq \frac{E^*\left[\sup_{\theta \in S_{j,n}} \left\|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta) - L_{1,p}^{\varepsilon_n} Q(\theta_0) + L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0)\right\|\right]}{\frac{2^j}{\rho_n}} = O\left(\rho_n \frac{1}{\sqrt{n\varepsilon_n}}\right) 2^{-(j-1)}. \end{aligned}$$

Thus if $\rho_n = O(\sqrt{n\varepsilon_n})$ then the probability of interest is $O(1)$. \square

A.7 Proof of theorem 9

Proof. We can note that the scaled bias due to numerical approximation can be evaluated as

$$\sqrt{n\varepsilon_n} (L_{1,p}^{\varepsilon_n} Q(\theta) - G(\theta)) = O_p\left(\sqrt{n\varepsilon_n^{1+\nu/2}}\right) = o_p(1).$$

We know that the estimator solves

$$L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) = o_p\left(\frac{1}{\sqrt{n\varepsilon_n}}\right).$$

Then

$$\sqrt{n\varepsilon_n} (L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) + G(\hat{\theta}) + L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) - G(\hat{\theta})) = o_p(1).$$

This means that locally

$$\begin{aligned} & \sqrt{n\varepsilon_n} (L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0) - L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) - L_{1,p}^{\varepsilon_n} Q(\theta_0)) \\ & + \sqrt{n\varepsilon_n} L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0) + \sqrt{n\varepsilon_n} H(\theta_0) (\hat{\theta} - \theta_0) = o_p(1). \end{aligned}$$

From Corollary 1 it follows that

$$\sqrt{n\varepsilon_n} (L_{1,p}^{\varepsilon_n} \hat{Q}(\hat{\theta}) - L_{1,p}^{\varepsilon_n} \hat{Q}(\theta_0) - L_{1,p}^{\varepsilon_n} Q(\hat{\theta}) - L_{1,p}^{\varepsilon_n} Q(\theta_0)) = o_p(1).$$

Then we can apply the CLT to obtain the desired result. \square

A.8 Proof of theorem 10

Proof. The rate of convergence adapts the proof of Theorem 3.2.5 of Van der Vaart and Wellner (1996) to our case. Denote the rate of convergence for the estimator $\hat{\theta}$ by ρ_n . Then we can partition the parameters space into sets $S_{j,n} = \{\theta : 2^{j-1} < \rho_n d(\theta, \theta_0) < 2^j\}$. Then we evaluate the probability of a large deviation $\rho_n d(\hat{\theta}, \theta_0) > 2^M$ for some integer M , where $\rho_n = \sqrt{n}\varepsilon_n^{1-\gamma}$. We know that the estimator solves

$$\sqrt{nr_n} L_{1,p}^{\varepsilon_n} Q_n(\hat{\theta}) = o_p(1).$$

If $\rho_n d(\hat{\theta}, \theta_0)$ is larger than 2^M for a given M , then over the θ in one of the shells $S_{j,n}$, $\sqrt{nr_n} L_{1,p}^{\varepsilon_n} Q_n(\theta)$ achieves a distance as close as desired to zero. Hence, for every $\delta > 0$,

$$P\left(\rho_n d(\hat{\theta}, \theta_0) > 2^M\right) \leq \sum_{\substack{j \geq M \\ 2^j < \delta \rho_n}} P\left(\sup_{\theta \in S_{j,n}} (-\|L_{1,p}^{\varepsilon_n} Q_n(\theta)\|) \geq -o_p\left(\frac{1}{\sqrt{nr_n}}\right)\right) + P\left(2d(\hat{\theta}, \theta_0) \geq \delta\right)$$

Note that mean square differentiability implies that for every θ in a neighborhood of θ_0 , $g(\theta) - g(\theta_0) \lesssim -d^2(\theta, \theta_0)$. Then we evaluate the population objective, using the fact that it has p mean-square derivatives:

$$\|L_{1,p}^{\varepsilon_n} Q(\theta)\| \geq C d(\theta, \theta_0) + C' \varepsilon_n^{\nu-1},$$

where θ_0 is the zero of the population first-order condition and the approximated derivative has a known order of approximation $\|L_{1,p}^{\varepsilon_n} Q(\theta_0)\| = C' \varepsilon_n^{\nu-1}$ for some constant C' . Substitution of this expression into the argument of interest leads to

$$\|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} Q_n(\theta)\| \geq \|L_{1,p}^{\varepsilon_n} Q(\theta)\| - \|L_{1,p}^{\varepsilon_n} Q_n(\theta)\| \geq C d(\theta, \theta_0) + C' \varepsilon_n^{\nu-1} + o_p\left(\frac{1}{\sqrt{nr_n}}\right).$$

Then applying the Markov inequality to the re-centered process for $\theta \in S_{j,n}$

$$P\left(r_n \sqrt{n} \|L_{1,p}^{\varepsilon_n} Q(\theta) - L_{1,p}^{\varepsilon_n} Q_n(\theta)\| \geq C r_n \sqrt{n} d(\theta, \theta_0) + C' r_n \sqrt{n} \varepsilon_n^{\nu-1} + o(1)\right) \leq C' r_n^{-1/2} n^{-1/2} \left(\frac{2^j}{\rho_n}\right)^{-1}.$$

Then $\rho_n = \sqrt{n}$ in the regular case and $\rho_n = r_n \sqrt{n}$ in cases where $\gamma \neq 1$.

Finally also note that the evaluation for the expectation holds for $\theta = \theta_0 \pm t_k \varepsilon_n$, as shown above. By Markov inequality according to Theorem 2.5.2 from van der Vaart and Wellner (1998) it follows that the process $r_n \sqrt{n} L_{1,p}^{\varepsilon_n} Q_n(\theta_0)$ indexed by ε_n is P-Donsker. \square

A.9 Proof of theorem 11

Proof. The result will follow if we can demonstrate that

$$\sqrt{nr_n} \left(L_{1,p}^{\varepsilon} \hat{g}(\hat{\theta}) - L_{1,p}^{\varepsilon} \hat{g}(\theta_0) - G(\hat{\theta}) + G(\theta_0) \right) = o_p(1). \quad (\text{A.12})$$

Because of the assumption that $\sqrt{n}\varepsilon^{\nu-\gamma} \rightarrow \infty$, the bias is sufficiently small. Therefore this is equivalent to showing that

$$\sqrt{nr_n} \left(L_{1,p}^{\varepsilon} \hat{g}(\hat{\theta}) - L_{1,p}^{\varepsilon} \hat{g}(\theta_0) - EL_{1,p}^{\varepsilon} \hat{g}(\hat{\theta}) + EL_{1,p}^{\varepsilon} \hat{g}(\theta_0) \right) = o_p(1).$$

Because of the convergence rate established in theorem 10, this will implied by, with $r_n = \varepsilon^{1-\gamma}$:

$$\sup_{d(\theta, \theta_0) \lesssim O\left(\frac{1}{\sqrt{n\varepsilon^{1-\gamma}}}\right)} \sqrt{nr_n} (L_{1,p}^\varepsilon \hat{g}(\theta) - L_{1,p}^\varepsilon \hat{g}(\theta_0) - EL_{1,p}^\varepsilon \hat{g}(\theta) + EL_{1,p}^\varepsilon \hat{g}(\theta_0)) = o_p(1).$$

The left hand side can be written as a linear combination of the empirical processes:

$$\sup_{d(\theta, \theta_0) \lesssim O\left(\frac{1}{\sqrt{n\varepsilon^{1-\gamma}}}\right)} \sqrt{n} \frac{r_n}{\varepsilon} [\mathbb{G}(\theta + t\varepsilon) - \mathbb{G}(\theta - t\varepsilon) - \mathbb{G}(\theta_0 + t\varepsilon) - \mathbb{G}(\theta_0 - t\varepsilon)].$$

Because of assumption 6, it is bounded stochastically by

$$O_p\left(\frac{r_n}{\varepsilon} \min(d(\theta, \theta_0), \varepsilon)^\gamma\right).$$

When $\sqrt{n}\varepsilon^{2-\gamma} \rightarrow \infty$, $d(\theta, \theta_0) \lesssim O\left(\frac{1}{\sqrt{n\varepsilon^{1-\gamma}}}\right) = o(\varepsilon)$. Hence the above display is $o_p(1)$. Therefore (A.12) holds.

Recall that $\hat{\theta}$ is defined by $\sqrt{nr_n}(L_{1,p}^\varepsilon \hat{g}(\hat{\theta}) - EL_{1,p}^\varepsilon \hat{g}(\hat{\theta})) = o_p(1)$. Then (A.12) implies that, using a first order taylor expansion of $G(\theta)$:

$$\sqrt{nr_n} (L_{1,p}^\varepsilon \hat{g}(\theta_0) - EL_{1,p}^\varepsilon \hat{g}(\theta_0)) + H(\theta_0) \sqrt{nr_n} (\hat{\theta} - \theta_0) = o_p(1).$$

□

A.10 Proof of Theorem 13

The result of the theorem follows directly from the results in Kim and Pollard (1990). First, we find that due to Assumption 4 and the result in Theorem 8, we can apply Lemma 4.1. Then from Assumptions 14 and 15 it follows that by Lemma 4.5 in Kim and Pollard (1990) applies. The result of Lemma 4.6. follows from Lemma 5 and Assumption 6 (for the Hölder-continuous case). This leads to the validity of Theorem 4.7. in Kim and Pollard (1990) which leads to the statement in our Theorem 13.

A.11 Proof of Theorem 14

Our proof will rely to a large extent to on the proof of Lemma 2. We first deliver the result regarding the consistency of $\hat{m}(\cdot)$. We consider the cases of the kernel and the sieve estimator for $m(\cdot)$ separately. Denote $\Delta(\theta, \eta, y_i) = (\rho(\theta, \eta; y_i) - m(\theta, \eta; z_i))_{i=1}^n$ and $G(\theta, \eta) = (E[\Delta(\theta, \eta, Y_i) | z_i])_{i=1}^n$. We note that from Assumption 8[iii] it follows that $\text{Var}(\rho(\theta, \eta; y_i) - m(\theta, \eta; z_i) | z) \leq \nu$. Select $\delta_n \rightarrow 0$ such that $\delta_n \gg N \sqrt{\frac{\log n^{1-2r_0}}{n^{1-2r_0}}}$. Denoting $\mu_n = \epsilon \delta_n^2 / N$ we find that $\text{Var}(P_n \rho(\theta, \eta; y_i)) / (4\mu_n)^2 \ll (\log n)^{-1}$, which implies that the symetrization inequality for empirical processes holds.

Utilizing the proof of Lemma 2 we can write that

$$\|\hat{m}(\theta, \eta, z) - m(\theta, \eta, z)\| \leq \frac{CN}{\Delta} \max_{k=1, \dots, N} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta(\theta, \eta, y_i) - G(\theta, \eta, z_i)) \right\|.$$

We note that each element in the latter sum belongs to the class described by Assumption 8[iv]. Considering

the individual elements in the sum we can apply the symmetrization argument in Pollard (1984) and write

$$\begin{aligned} & P \left(\sup_{(\theta, \eta, z) \in \Theta \times \mathcal{H}_n \times \mathcal{Z}} \frac{1}{n} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta - G) \right\| > 8\mu_n N \right) \\ & \leq 2 \exp \left(An^{2r_0} \log \frac{1}{\mu_n} \right) \exp \left(-\frac{1}{128} \frac{n\mu_n^2}{\nu} \right) + P \left(\sup_{(\theta, \eta, z) \in \Theta \times \mathcal{H}_n \times \mathcal{Z}} \frac{1}{n} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta_i - G_i) \right\|^2 > 64\nu \right). \end{aligned}$$

The second term can be evaluated with the aid of Lemma 33 in Pollard (1984):

$$P \left(\sup_{(\theta, \eta, z)} \frac{1}{n} \left\| \sum_{i=1}^n p_{Nk}(z_i) (\Delta_i - G_i) \right\|^2 > 64\nu \right) \leq 4 \exp \left(An^{2r_0} \log \frac{1}{\nu} \right) \exp(-n\nu).$$

As a result, we find that

$$\begin{aligned} P \left(\sup_{(\theta, \eta, z) \in \Theta \times \mathcal{H}_n \times \mathcal{Z}} \frac{1}{n} \|p^{N'}(z)\hat{Q}^{-1}P'(\Delta - G)\| > \mu_n \right) & \leq 2N \exp \left(An^{2r_0} \log \frac{N}{\epsilon\delta_n} - \frac{1}{128} \frac{n\epsilon^2\delta_n^2}{N^2\nu} \right) \\ & + 4N \exp \left(An^{2r_0} \log \frac{1}{\nu} - n\nu \right). \end{aligned}$$

A similar analysis can be conducted for the kernel-based estimator for the conditional moment function. We note that

$$\begin{aligned} & \|\hat{m}(\theta, \eta, z) - m(\theta, \eta, z)\| \\ & = \left(\frac{1}{nb_n^{d_z}} \sum_{i=1}^n K \left(\frac{z - z_i}{b_n} \right) \right)^{-1} \frac{1}{nb_n^{d_z}} \sum_{i=1}^n K \left(\frac{z - z_i}{b_n} \right) [\rho(\theta, \eta, y_i) - m(\theta, \eta, z)]. \end{aligned}$$

We note that the variance of each term in the summation has order $O(b_n^{d_z})$. Using the proof of Lemma 3, we find that the choice $\delta_n \gg \sqrt{\frac{\log n^{1-2r_0}}{b_n^{d_z} n^{1-2r_0}}}$ guarantees that the stochastic order $\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}_n, z \in \mathcal{Z}} \|\hat{m}(\theta, \eta, z) - m(\theta, \eta, z)\| = o_p(\delta_n)$.

Denote $\hat{G}(\theta, \eta, z_i) = L_{1,p}^{\varepsilon_n} \hat{m}(\theta, \eta, z_i)$, and $\hat{G}_j(\theta, \eta, z_i) = L_{1,p}^{\tau_n, \psi_j} \hat{m}(\theta, \eta, z_i)$ and $G(\theta, \eta, z_i)$ and $G_j(\theta, \eta, z_i)$ their population analogs. Then we can decompose

$$\begin{aligned} & \hat{G}(z_i, \theta, \eta)' \hat{W}(z_i) \hat{m}(\theta, \eta, z_i) - G(z_i, \theta, \eta) W(z_i) m(\theta, \eta, z_i) = \left(\hat{G}(z_i, \theta, \eta) \hat{W}(z_i) - G(z_i, \theta, \eta) W(z_i) \right) m(\theta, \eta, z_i) \\ & + G(z_i, \theta, \eta)' W(z_i) (\hat{m}(\theta, \eta, z_i) - m(\theta, \eta, z_i)) + \left(\hat{G}(z_i, \theta, \eta) \hat{W}(z_i) - G(z_i, \theta, \eta)' W(z_i) \right) (\hat{m}(\theta, \eta, z_i) - m(\theta, \eta, z_i)). \end{aligned}$$

We can provide the same expansion for $\hat{B}(\cdot)$.

Consider the difference $\hat{G}(z_i, \theta, \eta) \hat{W}(z_i) - G(z_i, \theta, \eta) W(z_i)$. From the proof of Lemma 2 it follows that $\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}_n, z \in \mathcal{Z}} \left\| \hat{G}(z_i, \theta, \eta) - G(z_i, \theta, \eta) \right\| = o_p(1)$. This means that

$$\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}_n, z \in \mathcal{Z}} \left\| \hat{G}(z_i, \theta, \eta) \hat{W}(z_i) \hat{m}(\theta, \eta, z_i) - G(z_i, \theta, \eta) W(z_i) m(\theta, \eta, z_i) \right\| = o_p(1).$$

We have also provided the proof that $\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}_n, z \in \mathcal{Z}} \|\hat{m}(\theta, \eta, z) - m(\theta, \eta, z)\| = o_p(1)$. This would imply that

$$\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}_n, z \in \mathcal{Z}} \left\| \left(\hat{G}(z_i, \theta, \eta) \hat{W}(z_i) - G(z_i, \theta, \eta)' W(z_i) \right) (\hat{m}(\theta, \eta, z_i) - m(\theta, \eta, z_i)) \right\| = o_p(1).$$

Then this means that

$$\frac{1}{n} \sum_{i=1}^n \left(\hat{G}(z_i, \theta, \eta)' \hat{W}(z_i) \hat{m}(\theta, \eta, z_i) - G(z_i, \theta, \eta) W(z_i) m(\theta, \eta, z_i) \right) = o_p(1),$$

provided that for n non-negative functions $\sup_{s \in S} \sum_{i=1}^n f(s) \leq \sum_{i=1}^n \sup_{s \in S} f(s)$. This proves that

$$\sup_{(\theta, \eta) \in \Theta \times \mathcal{H}_n} \left\| L_{1,p}^{\tau_n, \psi_j} \hat{Q}(\theta, \eta) - \frac{\partial Q(\theta, \eta)}{\partial \eta} [\psi_j] \right\| = o_p(1).$$

□

A.12 Proof of Lemma 6

The result of the theorem can be obtained by adapting the argument in Lemma 1 to the argument in the proof of Theorem 9 of Nolan and Pollard (1987). We define the class of functions $\mathcal{F}_n = \{\epsilon_n L_{1,p}^{\epsilon_n} g(\cdot, \cdot, \theta), \theta \in N(\theta_0)\}$, with envelope function F , such that $PF \leq C$. Then we can write

$$\sup_{d(\theta, \theta_0) \leq o(1)} \epsilon_n \|L_{1,p}^{\epsilon_n} \hat{g}(\theta) - L_{1,p}^{\epsilon_n} g(\theta)\| \leq \frac{1}{n(n-1)} \sup_{f \in \mathcal{F}_n} |S_n(f)|.$$

Noting (5.9), lemma 1 can be shown separately for the $\hat{\mu}_n(\theta)$ and $S_n(u)/n(n-1)$ components of the decomposition. Because assumption 17 is a special case of assumption 6, Theorem 3 applies with $\gamma = 1$. Therefore the result of lemma 6 holds for the $\hat{\mu}_n(\cdot)$ component as long as ϵ_n . We will hence with no loss of generality focus on $S_n(u)$ and assume that $g(\cdot, \cdot, \theta)$ is degenerate.

Due to Assumption 18, for each $f \in \mathcal{F}_n$, $E|f|^2 = E|\epsilon_n L_{1,p}^{\epsilon_n} g(\cdot, \cdot, \theta)|^2 = O(\epsilon_n)$. Define $t_n \geq \max\{\epsilon_n^{1/2}, \frac{\log n}{n}\}$ as in Lemma 10 of Nolan and Pollard (1987). Under the condition $n\sqrt{\epsilon_n}/\log n \rightarrow \infty$ in lemma 6, for large enough n , $t_n = \epsilon_n$. Denote $\delta_n = \mu t_n^2 n^2$. By the Markov inequality,

$$P\left(\sup_{f \in \mathcal{F}_n} |S_n(f)| > \delta_n\right) \leq \delta_n^{-1} P \sup_{f \in \mathcal{F}_n} |S_n(f)|.$$

By assumption 5, the covering integral of \mathcal{F}_n is bounded by a constant multiple of $H(s) = s[1 + \log(1/s)]$. the maximum inequality in Theorem 6 of Nolan and Pollard (1987) implies that

$$P \sup_{f \in \mathcal{F}_n} |S_n(f)|/n \leq C P H \left[\sup_{f \in \mathcal{F}_n} |T_n f^2|^{1/2}/n \right].$$

where T_n is the symmetrized measure defined in Nolan and Pollard (1987). The right-hand side can be further bounded by Lemma 10 in Nolan and Pollard (1987). This lemma states that there exists a constant β such that

$$P\left(\sup_{f \in \mathcal{F}_n} |S_{2n}(f)| > \beta^2 4n^2 t_n^2\right) \leq 2A \exp(-2n t_n),$$

where A is the Euclidean constant in assumption 5. Since $f(\cdot)$ is globally bounded, $|f(\cdot)|^2 \leq B|f(\cdot)|$ for a constant B . In addition, note that $|S_{2n}(f)| \geq |T_n f|$. Therefore, we find that $|T_n f^2| \leq B|S_{2n}(f)|$, which implies

$$P\left(\sup_{f \in \mathcal{F}_n} |T_n f^2| > \frac{4\beta^2}{B} n^2 t_n^2\right) \leq 2A \exp(-2n t_n).$$

Also note that $H[\cdot]$ achieves its maximum at 1 and is increasing for its argument less than 1. For sufficiently large n the term $\frac{4\beta^2}{B} t_n^2 \ll 1$. Then

$$\begin{aligned} P H \left[\sup_{f \in \mathcal{F}_n} |T_n f^2|^{1/2} / n \right] &= P \left(H \left[\frac{1}{n} \sup_{f \in \mathcal{F}_n} |T_n f^2|^{1/2} \right] \mathbf{1} \left\{ \sup_{f \in \mathcal{F}_n} |T_n f^2| > \frac{4\beta^2}{B} n^2 t_n^2 \right\} \right. \\ &\quad \left. + H \left[\frac{1}{n} \sup_{f \in \mathcal{F}_n} |T_n f^2|^{1/2} \right] \mathbf{1} \left\{ \sup_{f \in \mathcal{F}_n} |T_n f^2| < \frac{4\beta^2}{B} n^2 t_n^2 \right\} \right) \\ &\leq 1 \cdot P \left(\sup_{f \in \mathcal{F}_n} |T_n f^2| > \frac{4\beta^2}{B} n^2 t_n^2 \right) + H \left[\frac{2\beta}{\sqrt{B}} t_n \right] \cdot 1 \\ &\leq 2 A \exp(-2n t_n) + H \left(\frac{2\beta}{\sqrt{B}} t_n \right). \end{aligned}$$

Substituting this result into the maximum inequality one can obtain

$$\begin{aligned} P \left(\sup_{f \in \mathcal{F}_n} |S_n(f)| > \delta_n \right) &\leq n \delta_n^{-1} \left(H \left(\frac{2\beta}{\sqrt{B}} t_n \right) + 2 A \exp(-2n t_n) \right) \\ &= (t_n n)^{-1} + (n t_n)^{-2} \exp(-2n t_n) - (t_n n)^{-1} \log t_n. \end{aligned}$$

By assumption $t_n n \gg \log n \rightarrow \infty$, the first term vanishes. The second term also vanishes by showing that $n^{-1} t_n^{-2} \exp(-2n t_n) \rightarrow 0$, because it is bounded by, for some $C_n \rightarrow \infty$, $1/(\log n n^{C_n} t_n)$. Finally, considering the term $t_n^{-1} n^{-1} \log t_n$, we note that it can be decomposed into $t_n^{-1} n^{-1} \log(n t_n) - t_n^{-1} n^{-1} \log n$. Both terms converge to zero because both $t_n n \rightarrow \infty$ and $\frac{t_n n}{\log n} \rightarrow \infty$. We have thus shown that for any $\mu > 0$

$$P \left(\sup_{f \in \mathcal{F}_n} \left| \frac{1}{n(n-1)} S_n(f) \right| > \mu \varepsilon_n \right) = o(1).$$

This proves the statement of the theorem. \square

A.13 Proof of Lemma 7

Proof. (i)

We note that for the projection part

$$\sup_{d(\theta, \theta_0) = o(1)} \frac{1}{\sqrt{n}} \|L_{1,p}^{\varepsilon_n} \hat{\mu}(\theta) - L_{1,p}^{\varepsilon_n} \mu(\theta)\| = o_p(1).$$

As a result the U-process part will dominate and the convergence rate will be determined by its order $\frac{\log^2 n}{n^2 \varepsilon_n}$. The rest follows from the proof of Lemma 5.

(ii)

The proof of this lemma largely relies on the proof of Lemma 5 Consider a class of functions

$$\mathcal{G}_n = \left\{ g(\cdot, \theta_n + \varepsilon_n) - g(\cdot, \theta_n - \varepsilon_n) - g(\cdot, \theta_0 + \varepsilon_n) + g(\cdot, \theta_0 - \varepsilon_n), \theta_n = \theta_0 + t_n \frac{\log^2 n}{n^2 \varepsilon_n} \right\},$$

with $\varepsilon_n \rightarrow 0$ and $t_n = O(1)$. We can evaluate the L^2 norm of the functions from class \mathcal{G}_n using Assumption 5 (ii). Note that

$$E \left[(g(Z_i, z, \theta_n + \varepsilon_n) - g(Z_i, z, \theta_n - \varepsilon_n))^2 \right] = O(\varepsilon_n),$$

with the same evaluation for the second term. On the other hand, we can change the notation to $\theta_{1n} = \theta_0 + \varepsilon_n + \frac{t_n}{2} \frac{\log^2 n}{n^2 \varepsilon_n}$ and $\theta_{1n} = \theta_0 + \frac{\varepsilon_n}{2} + t_n \frac{\log^2 n}{n^2 \varepsilon_n}$. Then we can group the first term with the third and the second one with the fourth. For the first group this leads to

$$E \left[\left(g \left(Z_i, z, \theta_{1n} + \frac{t_n}{2} \frac{\log^2 n}{n^2 \varepsilon_n} \right) - g \left(Z_i, z, \theta_{1n} - \frac{t_n}{2} \frac{\log^2 n}{n^2 \varepsilon_n} \right) \right)^2 \right] = O \left(\frac{\log^2 n}{n^2 \varepsilon_n} \right),$$

and for the second group

$$E \left[\left(g \left(Z_i, z, \theta_{2n} + \frac{\varepsilon_n}{2} \right) - g \left(Z_i, z, \theta_{2n} - \frac{\varepsilon_n}{2} \right) \right)^2 \right] = O(\varepsilon_n).$$

Thus, two different ways of grouping the terms allow us to obtain two possible bounds on the norm of the entire term. As a result, we find that

$$P f^2 = O \left(\min \left\{ \varepsilon_n, \frac{\log^2 n}{n^2 \varepsilon_n} \right\} \right), \quad f \in \mathcal{G}_n.$$

Next we denote $\delta_n = \min \left\{ \varepsilon_n, \frac{\log^2 n}{n^2 \varepsilon_n} \right\}$.

Due to Assumption 18, for each $f \in \mathcal{F}_n$, $E|f|^2 = E|\epsilon_n L_{1,p}^{\epsilon_n} g(\cdot, \theta)|^2 = O(\epsilon_n)$. Define $t_n \geq \max\{\delta_n^{1/2}, \frac{\log n}{n}\}$ as in Lemma 10 of Nolan and Pollard (1987) then for $n\sqrt{\delta_n}/\log n \rightarrow \infty$

$$\sup_{\mathcal{F}_n} \left\| \frac{1}{n(n-1)} T_n(f^2) \right\| = o_p(\delta_n^2),$$

where T_n is the symmetrized measure defined in Nolan and Pollard (1987). By Assumption 5 (iii), the covering integral of \mathcal{F}_n is bounded by a constant multiple of $H(s) = s[1 + \log(1/s)]$. The maximum inequality in Theorem 6 of Nolan and Pollard (1987) implies that

$$P \sup_{f \in \mathcal{F}_n} |S_n(f)|/n \leq C P H \left[\sup_{f \in \mathcal{F}_n} |T_n f^2|^{1/2}/n \right].$$

Then the stochastic order of $\frac{1}{n\varepsilon_n} \sup_{f \in \mathcal{F}_n} |S_n(f)|$ can be evaluated as

$$\frac{\sqrt{n}}{\varepsilon_n} \frac{1}{n\varepsilon_n} \sup_{f \in \mathcal{F}_n} |S_n(f)| = O_p \left(\frac{\delta_n}{\varepsilon_n} \log \delta_n \right) = O_p \left(\frac{\log \left(\frac{n^2 \varepsilon_n}{\log n} \right)}{\frac{n^2 \varepsilon_n^2}{\log n}} \right) = o_p(1).$$

This delivers the result in the Lemma. □

Acknowledgments: The authors acknowledge generous research supports from the National Science Foundation, the University of California at Berkeley and Stanford University. We thank Tim Armstrong and numerous conference and seminar participants for insightful comments. The usual disclaimer applies.

References

ACKERBERG, D., X. CHEN, AND J. HAHN (2009): “A Practical Asymptotic Variance Estimator for Two-Step Semiparametric Estimators,” unpublished manuscript, UCLA and Yale University.

- AI, C., AND X. CHEN (2003): "Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions," *Econometrica*, 71, 1795–1843.
- ALEXANDER, K. (1984): "Probability inequalities for empirical processes and a law of the iterated logarithm," *The Annals of Probability*, 12(4), 1041–1067.
- AMEMIYA, T. (1985): *Advanced Econometrics*. Harvard University Press.
- ANDERSSEN, R., AND P. BLOOMFIELD (1974): "Numerical differentiation procedures for non-exact data," *Numerische Mathematik*, 22, 157–182.
- ANDREWS, D. (1997): "A stopping rule for the computation of generalized method of moments estimators," *Econometrica*, 65(4), 913–931.
- ARCONES, M., AND E. GINE (1993): "Limit theorems for U-processes," *The Annals of Probability*, 21, 1494–1542.
- CHEN, X., O. LINTON, AND I. VAN KEILEGOM (2003): "Estimation of Semiparametric Models when the Criterion Function Is Not Smooth," *Econometrica*, 71, 1591–1608.
- CHEN, X., AND D. POUZO (2009): "Estimation of nonparametric conditional moment models with possibly nonsmooth moments," SSRN working paper.
- CHEN, X., AND X. SHEN (1998): "Sieve Extremum Estimates for Weakly Dependent Data," *Econometrica*, 66, 289–314.
- DUDLEY, R. (1999): *Uniform central limit theorems*. Cambridge university press.
- DURBIN, J. (1971): "Boundary-crossing probabilities for the Brownian motion and Poisson processes and techniques for computing the power of the Kolmogorov-Smirnov test," *Journal of Applied Probability*, pp. 431–453.
- (1985): "The first-passage density of a continuous Gaussian process to a general boundary," *Journal of Applied Probability*, pp. 99–122.
- HARDLE, W., AND J. MARRON (1985): "Optimal bandwidth selection in nonparametric regression function estimation," *The Annals of Statistics*, pp. 1465–1481.
- HART, J., J. MARRON, AND A. TSYBAKOV (1992): "Bandwidth choice for average derivative estimation," *Journal of the American Statistical Association*, 87, 218–226.
- HOROWITZ, J. (1992): "A Smoothed Maximum Score Estimator for the Binary Response Models," *Econometrica*, 60.
- JONES, M., J. MARRON, AND S. SHEATHER (1996): "A brief survey of bandwidth selection for density estimation," *Journal of the American Statistical Association*, 91.
- JUDD, K. (1998): *Numerical Methods in Economics*. MIT Press.
- KIM, J., AND D. POLLARD (1990): "Cube root asymptotics," *Ann. Statist.*, 18, 191–219.

- KOSOROK, M. (2008): *Introduction to empirical processes and semiparametric inference*. Springer Verlag.
- L'ECUYER, P., AND G. PERRON (1994): "On the Convergence Rates of IPA and FDC Derivative Estimators," *Operations Research*, 42, 643–656.
- MANSKI, C. F. (1975): "Maximum Score Estimation of the Stochastic Utility Model of Choice," *Journal of Econometrics*, 3, 205–228.
- MURPHY, S., AND A. VAN DER VAART (2000): "On Profile Likelihood.," *Journal of the American Statistical Association*, 95.
- NEWHEY, W. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators," *Journal of Econometrics*, 79, 147–168.
- NEWHEY, W., AND D. MCFADDEN (1994): "Large Sample Estimation and Hypothesis Testing," in *Handbook of Econometrics*, Vol. 4, ed. by R. Engle, and D. McFadden, pp. 2113–2241. North Holland.
- NEWHEY, W., AND J. POWELL (2003): "Instrumental variable estimation of nonparametric models," *Econometrica*, 71, 1565–1578.
- NOLAN, D., AND D. POLLARD (1987): "U-processes: rates of convergence," *The Annals of Statistics*, pp. 780–799.
- PAKES, A., AND D. POLLARD (1989): "Simulation and the Asymptotics of Optimization Estimators," *Econometrica*, 57(5), 1027–1057.
- POLLARD, D. (1984): *Convergence of Stochastic Processes*. Springer Verlag.
- POWELL, J. L. (1984): "Least Absolute Deviations Estimation For The Censored Regression Model," *Journal of Econometrics*, 25, 303–325.
- PRESS, W., S. A. TEUKOLSKY, W. VETTERING, AND B. FLANNERY (1992): *Numerical Recipes in C, The Art of Scientific Computing*. Cambridge.
- PSHENICHNYI, B. (1971): *Necessary conditions for an extremum*. CRC.
- SERFLING, R. (1980): *Approximation Theorems in Mathematical Statistics*. John Wiley and Sons.
- SHEN, X., AND W. WONG (1994): "Convergence rate of sieve estimates," *The Annals of Statistics*, 22, 580–615.
- SHERMAN, R. P. (1993): "The limiting distribution of the maximum rank correlation estimator," *Econometrica*, 61, 123–137.
- VAN DER VAART, A. W., AND J. A. WELLNER (1996): *Weak convergence and empirical processes*. Springer-Verlag, New York.
- ZHANG, J., AND I. GIJBELS (2003): "Sieve empirical likelihood and extensions of the generalized least squares," *Scandinavian Journal of Statistics*, pp. 1–24.

Figure 1: Sample numerical first derivative (1st-order) for decreasing step sizes (from top to bottom)

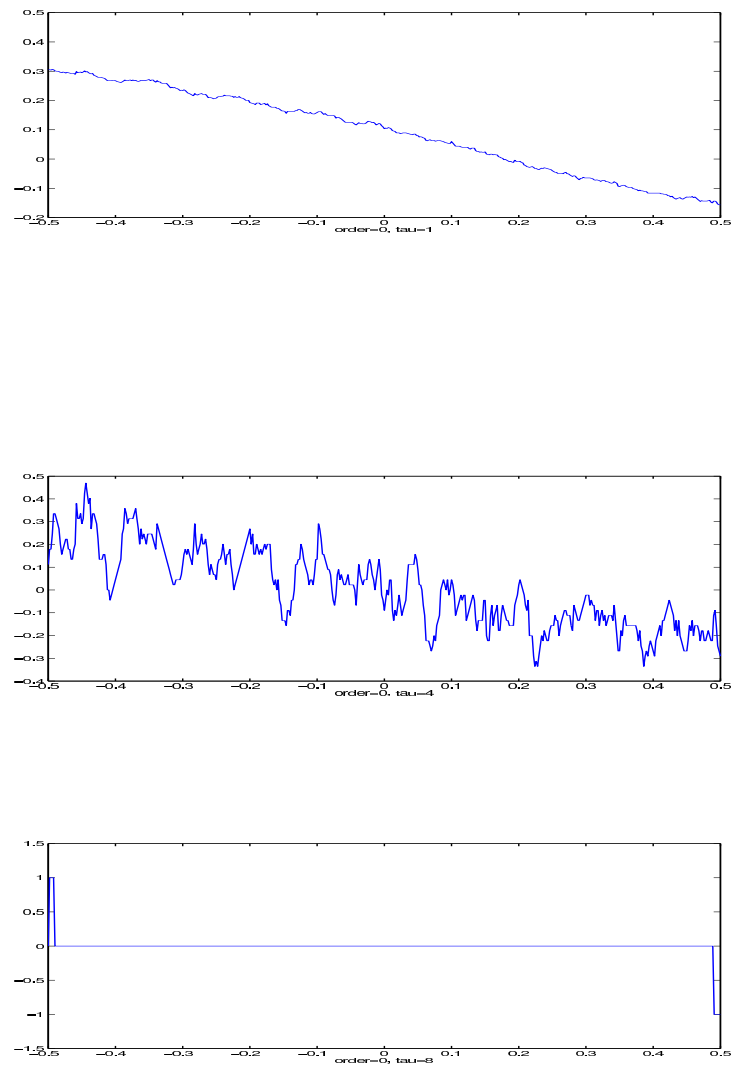


Figure 2: Sample numerical first derivative (2nd-order) for decreasing step sizes (from top to bottom)

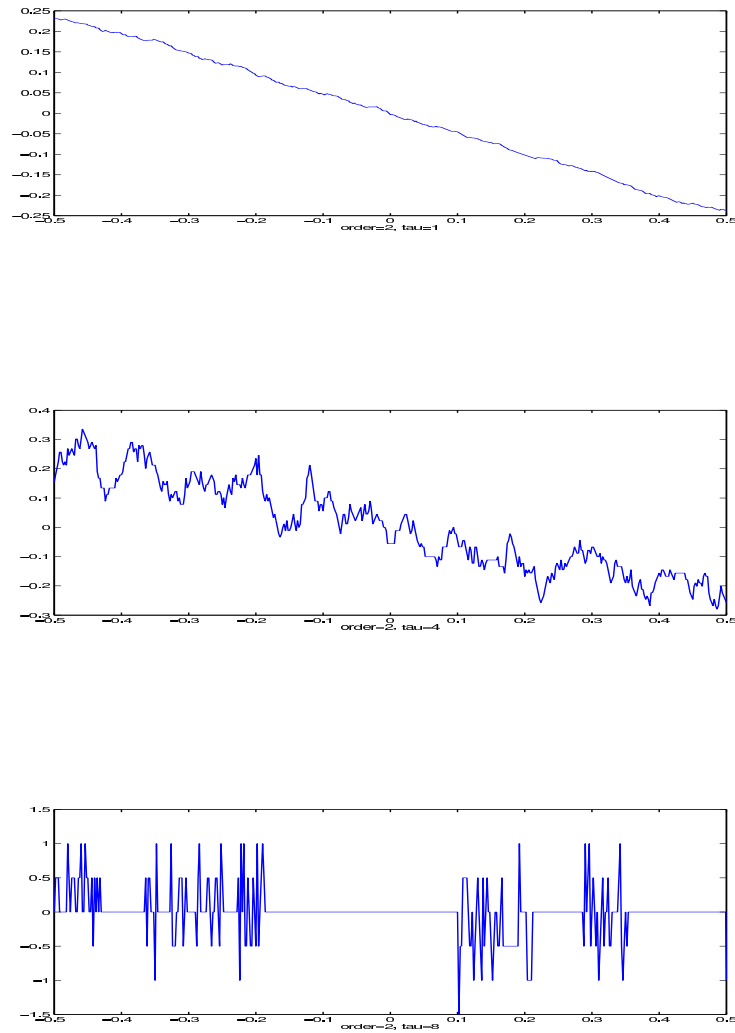


Figure 3: Sample numerical first derivative (3rd-order) for decreasing step sizes (from top to bottom)

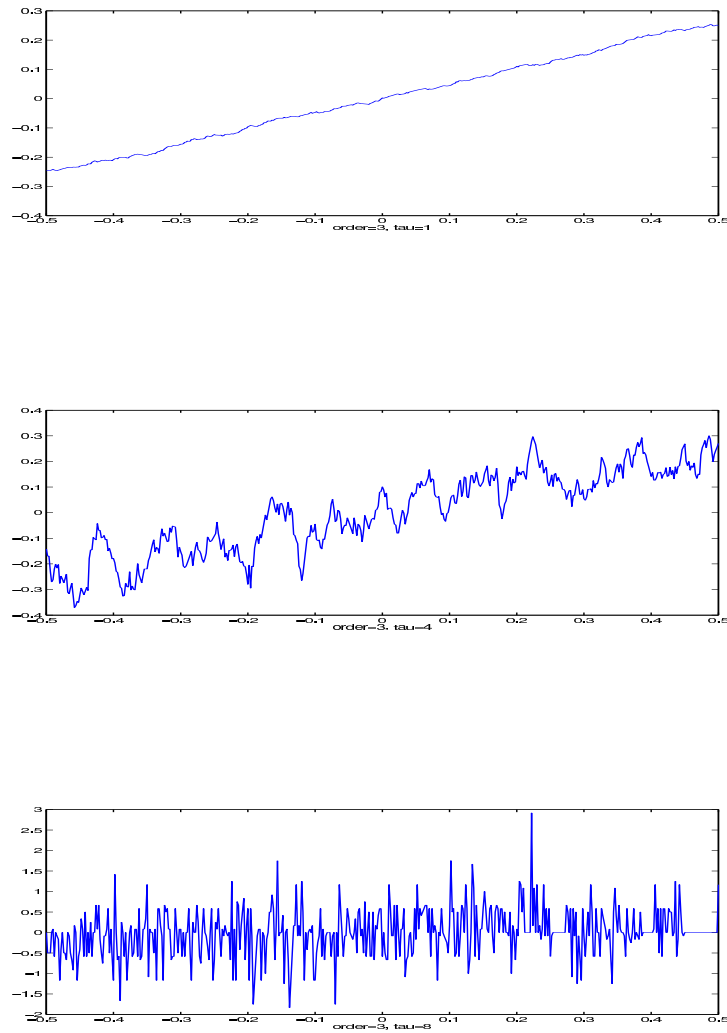


Figure 4: Mean-squared error of the estimated parameter

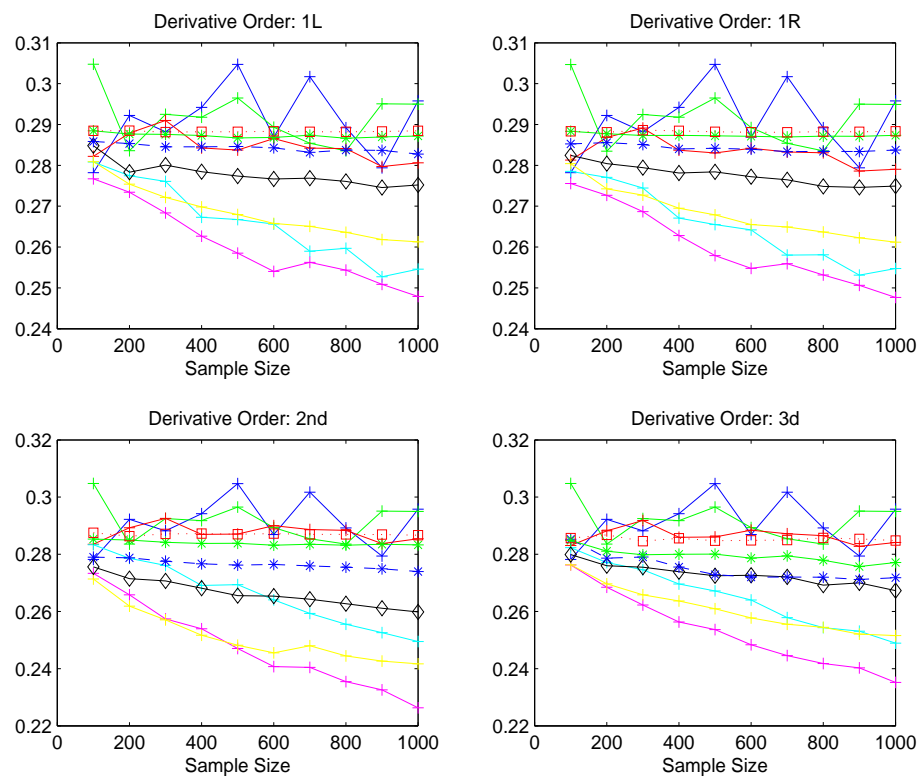


Figure 5: Bias of the estimated parameter

