

Higher-order Properties of Approximate Estimators

Dennis Kristensen* Bernard Salanié†

JULY 1, 2015

Abstract

Many modern estimation methods in econometrics approximate an objective function, for instance, through simulation or discretization. These approximations typically affect both bias and variance of the resulting estimator. We first provide a higher-order expansion of such “approximate” estimators that takes into account the errors due to the use of approximations. We show how a Newton-Raphson adjustment can reduce the impact of approximations. Then we use our expansions to develop inferential tools that take into account approximation errors: we propose adjustments of the approximate estimator that remove its first-order bias and adjust its standard errors. These corrections apply to a class of approximate estimators that includes all known simulation-based procedures. A Monte Carlo simulation on the mixed logit model shows that our proposed adjustments can yield spectacular improvements at a low computational cost.

1 Introduction

The complexity of econometric models has grown steadily over the past three decades. The increase in computer power contributed to this development in various ways, and in particular by allowing econometricians to estimate more complicated models using methods that rely on approximations. A leading example is simulation-based inference, where a function of observables and parameters is approximated using simulations. In this case, the function is an integral such as a moment, as in the simulated method of moments (McFadden (1989), Duffie and Singleton (1993)) and in simulated pseudo-maximum likelihood (Laroque and Salanié (1989, 1993, 1994)). It may also be an integrated density/cdf, as in simulated maximum likelihood (Lee (1992, 1995)), Kolmogorov-Smirnov type statistics (Corradi and Swanson (2007)), or a value function (Rust (1987)).¹ Then the approximation technique often amounts to Monte Carlo integration. Other numerical integration techniques may be

*UCL, IFS and CREATES. E-mail: d.kristensen@ucl.ac.uk.

†Columbia University. E-mail: bsalanie@columbia.edu

¹Simulation-based inference is surveyed in Gouriéroux and Monfort (1996), van Dijk, Monfort and Brown (1995) and Mariano, Schuerman and Weeks (2001) among others.

preferred for low-dimensional integrals, e.g. Gaussian quadrature, or both techniques can be mixed (see for example Lee (2001)). Within the class of simulation-based methods, some nonparametric alternatives rely on kernel sums instead of integration (e.g. Fermanian and Salanié (2004); Creel and Kristensen (2012); Kristensen and Shin (2012)), or on sieve methods (Kristensen and Schjerning (2011); Norets (2012)). Other estimation methods involve numerical approximations, such as discretization of continuous processes, using a finite grid in the state space for dynamic programming models, and so on.

In all of these cases, we call the *approximator* the term that replaces the component of the objective function that we cannot evaluate exactly. Then the *exact estimator* is the infeasible estimator that reduces the approximation error to zero. In simulation-based inference, the exact estimator would be obtained with an infinite number of simulations; in dynamic programming models it would rely on an infinitely fine grid. We call the estimator that relies on a finite approximation an *approximate estimator*.

The use of approximations usually deteriorates the properties of the approximate estimator relative to those of the corresponding exact estimator: the former may suffer from additional biases and/or lose efficiency compared to the latter. When the approximation error is not stochastic, its main effect is to impart additional bias to the estimator. On the other hand, stochastic approximations not only create bias: they also reduce efficiency. These are generic statements, of course. In some important special cases, such as the simulated method of moments, using approximations does not create additional bias, although it does reduce efficiency.

The effect of the approximation on the estimator can usually be reduced by choosing a sufficiently fine approximation; but this comes at the cost of increased computation time. In many applications this may be a seriously limiting factor; increased computer power helps, but it also motivates researchers to work on more complex models. It is therefore important to quantify the additional estimation errors that approximators generate, and also to account for these additional errors when drawing inference. As we will show, standard confidence intervals on the estimated parameters can be quite misleading unless they are properly adjusted for the errors induced by the approximation. As a first step in this direction, we analyze higher-order properties of the approximate estimator relative to the exact one in a very general setting. These expansions apply to a very large class of models, and they can be used to develop a number of adjustments to estimators and/or standard errors that open the way to better inference. To show this, we develop analytical bias and variance adjustments for a large class of approximate estimators where the approximation is stochastic. In simulation-based inference for instance, these adjustments remove the leading terms due to simulations. We also propose a very generally applicable two-step method; it consists of updating the approximate estimator obtained by one or several Newton-Raphson iterations based on the same objective function, but with a much finer degree of approximation. These different

methods can of course be combined when they both apply.

Our theoretical results applies to generalized method of moment estimators as well as M-estimators, both when the approximation is stochastic and when it is not. The results encompass and extend results in the literature on simulation-based estimators, such as Laroque and Salanié (1989), MacFadden (1989), Pakes and Pollard (1989), Lee (1995, 1999) and Gouriéroux and Monfort (1996). Moreover, the expansion can be used to analyze the behavior of estimators that rely on numerical approximation. Many structural estimates rely on such approximations, in asset pricing models (Tauchen and Hussey, 1991), DSGE models (Fernández-Villaverde, Rubio-Ramirez and Santos, 2006), or dynamic discrete choice models (Rust, 1987). Our results also apply to many estimators used in empirical IO, which combine simulation and numerical approximation. And it also covers situations where numerical derivatives are used, either for computation of variance estimators or optimization algorithms based on Newton iterations². To the best of our knowledge, this is the first paper to provide results for such a general class of models.

To test the practical performance of our proposed adjustment methods, we run a simulation study on a mixed logit model. The mixed logit is one of the basic building blocks in much work in demand analysis, for instance; and it is simple enough that we can compute the true value of the biases and efficiency losses, as well as our estimated corrections. We show that uncorrected SML has non-negligible bias, even for large sample sizes; and that standard confidence intervals can be wildly off the mark. Our analytical adjustment removes most of the bias at almost no additional computational cost; and it yields very reliable confidence intervals. The Newton-Raphson correction also reduces the bias and improves confidence intervals, but it does so less effectively than the analytical adjustment.

In a recent paper, Freyberger (2015) derived analytical adjustments for the BLP model when the numbers of consumers and/or the number of simulation draws are finite. His results are complementary to ours: they are less general, but since he only deals with a specific model his assumptions are weaker and his formulæ are more explicit.

The paper is organized as follows. Section 2 presents our framework and some examples. In Section 3, we derive a higher-order expansion of the approximate estimator relative to the exact one. We describe our Newton-Raphson correction in section 4. Then in Section 5 we build on the expansion to propose adjusted estimators, standard errors, and confidence intervals. Section 6 applies the general theory to two specific approximate estimators, while Section 7 presents the results of a Monte Carlo simulation study using the simulated MLE of the mixed logit model as an example. We discuss possible extensions of our results in Section 8. Appendix A and B contain proofs of the main results and lemmas, respectively. Appendix C provides details for one example of our theory, and Appendix D outlines how the theory

²However, in most of our examples, we abstract away from issues with numerical maximization that sometimes arise when computing extremum estimators.

can be generalized to handle multiple approximators with different properties.

2 Framework

At the most general level, our framework can be described as follows. Given a sample $\mathcal{Z}_n = \{z_1, \dots, z_n\}$ of n observations, the econometrician proposes to estimate a parameter $\theta_0 \in \Theta \subseteq \mathbb{R}^k$ through an estimating equation, $G_n(\theta, \gamma) = \mathcal{G}_n(\mathcal{Z}_n, \theta, \gamma)$ that the estimator $\hat{\theta}_n$ is set to solve, $G_n(\hat{\theta}_n, \gamma) = 0$. The estimating equation depends on the data \mathcal{Z}_n both directly and, possibly, via a nuisance parameter γ_0 that may also depend on the unknown parameter. The parameter γ_0 could be finite-dimensional (e.g. it could be an estimated variance), but in most situations it is a parameter dependent function, $u \mapsto \gamma(u; \theta)$, and so is infinite-dimensional. The nature of the argument u of the function γ will vary a lot with the applications; it could be covariates relative to one observation, the value of a conditional moment, or more complex objects. This is irrelevant for our general theory.

Our paper considers situations where the object γ_0 is not known in closed form to the econometrician, and instead has to be approximated. In this case, a feasible estimator is obtained by solving the analog estimating equation $G_n(\theta, \hat{\gamma}_S) = 0$ w.r.t. θ , where $\hat{\gamma}_S$ is the chosen “approximator” that depends on some approximation scheme of order S (e.g. S simulations, or a discretization on a grid of size S). The resulting estimator, denoted $\hat{\theta}_{n,S}$, will be referred to as the “approximate” estimator. We throughout restrict attention to the case of smooth approximators where $\hat{\gamma}_S$ is, as a minimum, differentiable.

The object γ may be a vector-valued function. If it is, then its components may entail very different biases and variances to the approximate estimator. To save on notation, in the body of the paper we assume that the biases and variances due to approximations of the different components vanish at the same rate. Appendix D provides results for the case of multiple approximators with possibly different rates.

Note that our problem is similar to two-step semiparametric estimation, where in the first step a nuisance parameter (γ_0) is replaced by its estimator (the approximator $\hat{\gamma}_S$), which in turn is used to obtain an estimator ($\hat{\theta}_{n,S}$) of θ_0 . Some themes of that literature (see e.g. Andrews 1994 and Chen et al. 2003) recur in our analysis.

2.1 Examples

We now present a few examples that fall within the above setting.

Example 1: Simulated maximum likelihood (SML). Suppose we want to estimate a parameterized conditional distribution $p(y|x; \theta)$. The natural choice is the maximum-likelihood estimator, which maximizes $L_n(\theta; p) = \frac{1}{n} \sum_{i=1}^n \log p(y_i|x_i; \theta)$. Sometimes the density p cannot be written in closed form. For example, in models with unobserved hetero-

generity, $p(y|x; \theta) = \int w(y|x, \varepsilon; \theta) f(\varepsilon) d\varepsilon$ for some densities w and f . If the integral cannot be computed analytically, a simulated version can be obtained by drawing $(\varepsilon_s, s = 1, \dots, S)$ from the distribution of f to obtain a simulated version by $\hat{p}_S(y|x; \theta) = \sum_{s=1}^S w(y|x, \varepsilon_s; \theta) / S$. Maximizing $L_n(\theta; \hat{p}_S)$ gives a simulated maximum likelihood estimator (SMLE). If the density w is a twice differentiable function of θ , then this fits in our framework, with $z = (y, x)$, $\gamma(z; \theta) := p(y|x; \theta)$, and

$$G_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \log \gamma}{\partial \theta}(z_i, \theta).$$

More recently, Fermanian and Salanié (2004) proposed using kernel estimators as approximators. Suppose that $y = r(x, \varepsilon; \theta_0)$, with implied conditional density $\gamma(z; \theta) = p(y|x, \theta)$. Then generate samples, $y_s(x, \theta) = r(x, \varepsilon_s; \theta)$ for $s = 1, \dots, S$, and approximate the density γ with a kernel density estimator based on the y_s 's: $\hat{\gamma}_S(z; \theta) = \sum_{s=1}^S K_h(y - y_s(x, \theta)) / S$. Maximizing $L_n(\theta; \hat{\gamma}_S)$ defines a nonparametric simulated maximum likelihood estimator (NPSMLE). For a similar approach in time series models, see Altissimo and Mele (2009) and Kristensen and Shin (2012).

Example 2: Simulated pseudo-maximum likelihood (SPML). Suppose that we have the following conditional moment restriction, $E[y|x] = m(x; \theta)$, where, for some function w and some unobserved error ε , $m(x; \theta) = E[w(x, \varepsilon; \theta) | x]$. Nonlinear least squares would choose θ to minimize $\sum_{i=1}^n (y_i - m(x_i; \theta))^2 / (2n)$. If the conditional expectation m cannot be evaluated analytically, Laroque and Salanié (1989) proposed simulated pseudo-maximum likelihood (SPML) estimators: Draw i.i.d. random variables ε_s , $s = 1, \dots, S$, and define $\hat{m}_S(x; \theta) = S^{-1} \sum_{s=1}^S w(x, \varepsilon_s; \theta)$. Then an SNLS estimator is obtained by replacing m with \hat{m}_S . Again, this fits into our framework if m is a twice differentiable function of θ . Then we take $\gamma(x; \theta) = m(x; \theta)$ and $G_n(\theta, \gamma_0) = \partial Q_n(\theta, \gamma_0) / (\partial \theta)$, where $Q_n(\theta, \gamma_0) = \sum_{i=1}^n (y_i - \gamma_0(x_i; \theta))^2 / (2n)$. The above idea can be extended to incorporate information regarding the conditional variance of y , or more generally to a variety of PML and QML estimators.

Example 3: Simulated method of moments (SMM). The parameter of interest is identified through a set of moment conditions $M(\theta_0) := Em(z, \theta_0) = 0$. Given a weighting matrix W_n , the GMM estimator would minimize $Q_n(\theta, \gamma_0) = M_n(\theta)' W_n M_n(\theta)$, where $M_n(\theta) = \sum_{i=1}^n m(z_i, \theta) / n$. Here, $\gamma(\theta) = M(\theta)$, which may be hard to evaluate, as in the multinomial probit example of McFadden (1989). Another example is the simulated method of moments (SMM) proposed by Duffie and Singleton (1993) to estimate dynamic models where a long string of simulations from the model, say $\{y_s(\theta) : s = 1, \dots, S\}$, are used to approximate unconditional moments of the model. The resulting estimator is of the minimum-distance type. Creel and Kristensen (2012) generalize the approach of Duffie and Singleton (1993) to the case where a set of conditional moments are used in the estimation; in this

case, similar to SPML, γ_0 is a conditional moment function which is then approximated by combining simulations with kernel regression techniques.

Example 4: Estimation of structural models. Evaluating the value function in dynamic programming models most often requires numerical approximations. Sometimes they involve simulations of moments³, which are stochastic approximators in our nomenclature. They also often resort to interpolation or sieve methods (also referred to as parametric approximations); Rust (1987) or Keane-Wolpin (1994, 1997) are classic examples. Then the approximated value function plays the role of γ_0 .

Similarly, many models used in macroeconomics are so complex that estimation is based on an approximate model, often by linearizing equations close to a steady state. The quality of the model approximation can be improved at a larger computational cost by using a finer grid or by using, for example, more iterations of perturbations or projection methods as advocated by Judd, Kubler and Schmedders (2003). For a first-order theoretical analysis of the impact on the resulting approximate MLE, see Fernández-Villaverde, Rubio-Ramirez and Santos (2006) and Akerberg, Geweke and Hahn (2009).

Another example involves numerical inversion of functions. One example of this arises in the estimation procedure proposed by Berry, Levinsohn and Pakes (1995—hereafter BLP) to estimate discrete choice models in industrial organization. Here, observed market shares (s) are modelled as functions of unobserved (ξ) and observed (z) product characteristics, $s = P(\xi, z; \theta)$ for some choice probability function P . The function P is usually computed by Monte-Carlo integration over unobserved individual preference shocks—a stochastic approximation scheme. The BLP estimation procedure requires the econometrician to compute the unobserved product characteristics given observed market shares by inverting the market share function with respect to its first argument, $\xi(s, z; \theta) = P^{-1}(s, z; \theta)$. Since P^{-1} is not available in closed form, this is done using a numerical fixed-point algorithm. It leads to an approximate solution, $\hat{\xi}_S(s, z; \theta)$, where S captures both the number of draws in the Monte Carlo integral and the number of iterations (and/or the tolerance level used in the algorithm). Our function γ here is simply P^{-1} . Judd and Su (2012) and Dubé, Fox and Su (2012) recently emphasized that the quality of the estimates of θ is very sensitive to errors in the computation of the fixed point. They propose using mathematical programming under equilibrium constraints instead, with an augmented parameter vector (θ, ξ) . This also fits into our framework, and this time it would be natural to use γ_0 to denote the Lagrange multipliers for the market share equations. In Freyberger 2015, the difference between γ and γ_0 comes both from the number of simulation draws and from the finite number of consumers on each market (which we assumed infinite); he uses MPEC and he neglects the approximation in the inversion algorithm itself.

There are many other examples of numerical approximators. For instance, the derivatives

³See e.g. Kristensen and Schjerning (2011) and Norets (2009, 2012).

of the sample objective function are often approximated numerically, either to maximize it or to estimate the asymptotic variance. This often involves computing sample averages of derivatives, e.g. $\sum_{i=1}^n \gamma_k(z_i; \theta) / n$, where $\gamma_k(z; \theta) = \partial q(z; \theta) / (\partial \theta_k)$ for $k = 1, \dots, \dim(\theta)$. We replace $\gamma_k(z; \theta)$ with, for example, $\hat{\gamma}_{S,k}(z; \theta) = [q(z; \theta + \epsilon_S e_k) - q(z; \theta - \epsilon_S e_k)] / (2\epsilon_S)$, where e_k is the k th column of the identity matrix and $\epsilon_S \rightarrow 0$ as $S \rightarrow \infty$. Our theory applies to approximate variance estimators built around numerical derivatives, as well as to estimators built around quasi-Newton iterations that use numerical derivatives:

$$\hat{\theta}_{n,S}^{(k+1)} = \hat{\theta}_{n,S}^{(k)} - H_1(\hat{\theta}_{n,S}^{(k)}) \frac{1}{n} \sum_{i=1}^n \hat{\gamma}_S(z_i; \hat{\theta}_{n,S}^{(k)}), \quad k = 1, \dots, \dim(\theta).$$

* * *

In all of the examples above, approximations reduce the quality of the estimator. Start with our first three examples, where stochastic approximations (i.e. simulations) are used to evaluate a mathematical expectation. The mean of course is an unbiased estimator of the expectation; but in many simulation-based estimation methods the objective function depends nonlinearly on the simulated mean, so that the approximate estimator based on S simulations has an additional bias, along with reduced efficiency. In many cases both are of order $1/S$; this holds for example when the approximator simulates an expectation through a simple average. When using nonparametric techniques such as kernel smoothers or sieve methods in the approximation, the approximator itself is biased, and the objective function will be biased even if the approximator enters linearly.

The simulated method of moments (Example 3) is a special case. This approximate estimator has nicer properties since the objective function is linear in the simulated mean; while the asymptotic efficiency loss still is of order $1/S$, the simulations do not impart bias to the estimator, except when kernel smoothers are employed as in Creel and Kristensen (2012).

Non-stochastic approximations also lead to deteriorations of the properties of the resulting estimators. Take the problem of computing the density $p(y|x; \theta)$ in Example 1 for instance. If the dimension of the integration variable (ε) is small, then instead of simulations the numerical integration may be done by an S point Gaussian quadrature, as in Lee (2001). Because this is a non-stochastic approximation method, the resulting approximate estimator will suffer from additional biases, but its variance will not increase.

As this informal discussion illustrates, the approximate estimator $\hat{\theta}_{n,S}$ often is consistent only if S goes to infinity as n goes to infinity; and \sqrt{n} -consistency requires that S diverges fast enough. In other words⁴, $\|\hat{\theta}_{n,S} - \hat{\theta}_n\| = o_P(1/\sqrt{n})$ as $n \rightarrow \infty$ for some sequence $S = S(n) \rightarrow \infty$, in which case there is no first-order difference between the exact and

⁴Section 3.3 will give more precise statements and regularity conditions.

approximate estimator. For finite S and n , our higher-order expansion allows the researcher to better evaluate the properties of the approximate estimator.

To derive the higher-order expansion of the approximate estimator, we need to impose regularity conditions both on the estimating equation and on the approximators. We present these conditions in the next two subsections.

2.2 Estimating Equation

We restrict our attention to the class of exact estimators $\hat{\theta}_n$ that (asymptotically) satisfy a first order condition of the form

$$G_n(\hat{\theta}_n, \gamma_0) = o_P(1/\sqrt{n}), \text{ where } G_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n g(z_i; \theta, \gamma) \quad (1)$$

for some functional $g(z; \theta, \gamma)$ that depends on data, z_1, \dots, z_n , the parameter of interest, θ , and some unknown “true” function, $\gamma_0(u)$. The function argument, u , is often an observed variable and so part of z , but it could also be an unobserved component that, for example, gets integrated out in the computation of g . The corresponding approximate estimator $\hat{\theta}_{n,S}$ similarly satisfies

$$G_n(\hat{\theta}_{n,S}, \hat{\gamma}_S) = o_P(1/\sqrt{n}). \quad (2)$$

The above framework corresponds to the so-called z -estimators, a large family that includes m -estimators and GMM estimators. All of the examples described in Section 2.1 are z -estimators. For m -estimators, $\hat{\theta}_{n,S} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n q(z_i; \theta, \hat{\gamma}_S) / n$, we can choose $g(z; \theta, \gamma) = \partial q(z; \theta, \gamma) / (\partial \theta)$. For GMM estimators, $\hat{\theta}_{n,S} = \arg \min_{\theta \in \Theta} M_n(\theta, \hat{\gamma}_S) W_n M_n(\theta, \hat{\gamma}_S)$ where $W_n \xrightarrow{P} W > 0$ and $M_n(\theta, \gamma) = \sum_{i=1}^n m(z_i; \theta, \gamma) / n$, we set $g(z_i; \theta, \gamma) = H(\theta, \gamma) W m(z_i; \theta, \gamma)$ where $H(\theta, \gamma) = E[\partial m(z_i; \theta, \gamma) / (\partial \theta)]$.

We assume that the function of interest $\gamma_0 : \mathcal{U} \times \Theta \mapsto \mathbb{R}^p$ belongs to a linear function space Γ equipped with a norm $\|\cdot\|$. In most cases, the norm will be either the sup-norm, $\|\gamma\| = \sup_{u \in \mathcal{U}} \sqrt{(\gamma(u)' \gamma(u))}$, or some L_q -norm induced by the probability measure associated with the data generating process, $\|\gamma\| = E[(\gamma(u)' \gamma(u))^q]^{1/q}$ for some $q \geq 1$. Our analysis will involve the following sample and population averages,

$$H_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n h(z_i; \theta, \gamma), \quad G(\theta, \gamma) = E[g(z_i; \theta, \gamma)], \quad H(\theta, \gamma) = E[h(z_i; \theta, \gamma)], \quad (3)$$

where $h(z_i; \theta, \gamma) = \partial g(z_i; \theta, \gamma) / (\partial \theta)$. We first impose conditions to ensure that the exact, but infeasible estimator is well-behaved:

- A.1** (i) $\hat{\theta}_n \xrightarrow{P} \theta_0$ which lies in the interior of the parameter space Θ ; (ii) $\{z_i\}$ is stationary and geometrically α -mixing; (iii) $E[\|g(z_i; \gamma_0)\|^{2+\delta}] < \infty$ for some $\delta > 0$; (iv) $G(\theta_0, \gamma_0) = 0$.

A.2 For all γ in a neighbourhood \mathcal{N} of γ_0 , $h(z; \theta, \gamma)$ satisfies: (i) $H_0 := H(\theta_0, \gamma_0)$ is positive definite; (ii) for some $\delta > 0$, $E[\sup_{\|\theta - \theta_0\| < \delta} \|h(z_i; \theta, \gamma_0)\|] < \infty$; (iii) for some $\delta, \lambda, \bar{H} > 0$, and for all $\gamma \in \mathcal{N}$,

$$E \left[\sup_{\|\theta - \theta_0\| < \delta} \|h(z_i; \theta, \gamma) - h(z_i; \theta, \gamma_0)\| \right] \leq \bar{H} \|\gamma - \gamma_0\|^\lambda.$$

Assumption A.1(i) requires that the infeasible estimator be consistent; Lemma 1 below provides a set of sufficient conditions. A.1(ii) rules out strongly persistent data, thereby allowing us to obtain standard rates of convergence for the resulting estimators. In particular, A1(ii) and A.1(iii) together imply that a central limit theorem (CLT) applies to $G_n(\theta_0, \gamma_0)$. The geometric mixing condition could be weakened, but this would lead to more complicated results; we refer the reader to Kristensen and Shin (2012) for some results on approximate estimators based on strongly persistent and/or non-stationary data (and thereby estimators with non-standard rates) in the context of SMLE.

Assumption A.2 imposes differentiability of $\theta \mapsto g(z; \theta, \gamma)$. In particular, when γ depends on θ (as is the case for all of our examples), it requires the approximator to be a smooth function of θ . Therefore A.2 rules out discontinuous and non-differentiable approximators, such as the simulated method of moment estimators for discrete choice models proposed in McFadden (1989) and Pakes and Pollard (1989) which involve indicator functions⁵. The Lipschitz condition imposed on $h(z; \theta, \gamma)$ is used to ensure that $H_n(\theta, \hat{\gamma}_S) \xrightarrow{P} H(\theta, \gamma)$ uniformly in θ as $\hat{\gamma}_S \xrightarrow{P} \gamma$.

Since our focus is on higher-order properties of the approximate estimator, we also assume that consistency has already been established, so that we can conduct our analysis locally around θ_0 :

A.3 $\hat{\theta}_{n,S} \xrightarrow{P} \theta_0$ as $n, S \rightarrow \infty$.

A set of sufficient conditions (similar to those in Newey and McFadden, 1994 for consistency of two-step semiparametric estimators) for Assumptions A.1 (i) and A.3 to hold are provided in the following lemma:

Lemma 1 *Suppose that $\hat{\theta}_{n,S} = \arg \max_{\theta \in \Theta} Q_n(\theta, \hat{\gamma}_S)$ where: (i) Θ is compact; (ii) $\hat{\gamma}_S \xrightarrow{P} \gamma_0$; (iii) either $\sup_{\theta \in \Theta, \|\gamma - \gamma_0\| < \delta} |Q_n(\theta, \gamma) - Q(\theta, \gamma)| \xrightarrow{P} 0$ or $|Q_n(\theta, \gamma_1) - Q_n(\theta, \gamma_2)| \leq B_n \|\gamma_1 - \gamma_2\|$ for all γ_1, γ_2 in a neighbourhood of γ_0 where $B_n = O_P(1)$ and $\sup_{\theta \in \Theta} |Q_n(\theta, \gamma_0) - Q(\theta, \gamma_0)| \xrightarrow{P} 0$; (iv) $\theta \mapsto Q(\theta, \gamma_0)$ is continuous and has a unique maximum at θ_0 . Then A.1(i) and A.3 hold.*

⁵These cases could be handled by introducing a smoothed version of the approximators as discussed in McFadden (1989); see also Fermanian and Salanié (2004). Alternatively, one could extend our results by resorting to empirical process theory, as done in the work by Armstrong et al (2012) on simulation-based estimators.

As a first step in our higher-order analysis, we prove in Appendix B (Lemma 7) that under Assumptions A.1 to A.3,

$$\hat{\theta}_{n,S} - \hat{\theta}_n = -H_0^{-1} \{G_n(\theta_0, \hat{\gamma}_S) - G_n(\theta_0, \gamma_0)\} + o_P(1/\sqrt{n}). \quad (4)$$

If γ was a finite-dimensional parameter, we could use a Taylor expansion of $G_n(\theta_0, \hat{\gamma}_S) - G_n(\theta_0, \gamma_0)$ to analyze the additional estimation errors due to the approximator $\hat{\gamma}_S$. However, γ may be a function; for such a functional expansion to be well-defined and for the individual terms in the expansion to be well-behaved, we need to impose some further regularity conditions on $g(z_i, \gamma)$ as a functional of γ . In all of the following, $d\gamma \in \Gamma$ denotes a small change around γ .

A.4(m) There exist functionals $\nabla^k g(z, \theta, \gamma) [d\gamma_1, \dots, d\gamma_k]$ for (θ, γ) in a neighbourhood of (θ_0, γ_0) , and constants $\delta > 0$ and $\bar{G}_k > 0$, $k = 0, \dots, m$, such that: (i) each $\nabla^k g$ is linear in each of its components $d\gamma_i \in \Gamma$, $i = 1, \dots, k$; (ii)

$$E \left[\left\| g(z, \theta, \gamma_0 + d\gamma) - g(z, \theta, \gamma_0) - \sum_{k=1}^m \frac{1}{k!} \nabla^k g(z, \theta, \gamma) [d\gamma, \dots, d\gamma] \right\|^2 \right] \leq \bar{G}_0 \|d\gamma\|^{m+1}, \quad (5)$$

where

$$E \left[\|\nabla g(z, \theta, \gamma) [d\gamma]\|^2 \right] \leq \bar{G}_1 \|d\gamma\|^2,$$

and, for some $\nu > 0$ and for $k = 2, \dots, m$,

$$E \left[\left\| \nabla^k g(z, \theta, \gamma) [d\gamma_1, \dots, d\gamma_k] \right\|^{2+\nu} \right] \leq \bar{G}_k (\|d\gamma_1\| \cdots \|d\gamma_k\|)^{2+\nu}. \quad (6)$$

Assumption A.4(m) restricts $g(z, \theta, \gamma)$ to be m times pathwise differentiable w.r.t. γ with differentials $\nabla^k g(z) [d\gamma_1, \dots, d\gamma_k]$ that are Lipschitz in $d\gamma_1, \dots, d\gamma_k$, $k = 1, \dots, m$. For a given choice of m , this allows us to use an m th order expansion of $G_n(\theta, \gamma)$ w.r.t. γ to evaluate the impact of $\hat{\gamma}_S$. In particular, the difference between the approximate and the exact objective functions can be written as

$$G_n(\theta_0, \hat{\gamma}_S) - G_n(\theta_0, \gamma_0) = \sum_{k=1}^m \frac{1}{k!} \nabla^k G_n(\theta_0, \gamma_0) [\hat{\gamma}_S - \gamma_0, \dots, \hat{\gamma}_S - \gamma_0] + R_{n,S}, \quad (7)$$

where $R_{n,S} = O_P(\|\hat{\gamma}_S - \gamma_0\|^{m+1})$ is the remainder term, and

$$\nabla^k G_n(\theta, \gamma) [d\gamma_1, \dots, d\gamma_k] = \frac{1}{n} \sum_{i=1}^n \nabla^k g(z_i, \theta, \gamma) [d\gamma_1, \dots, d\gamma_k]. \quad (8)$$

To evaluate the higher-order errors due to the approximation, we will study the mean and variance of each of the terms in the sum on the right hand side of (7).

Assumption A.4 holds for a wide range of objective functions, such as those of SPML and SMM. However, in some cases minor modifications of the estimating equations such as trimming will be required.

2.3 Approximators

To analyze the impact of approximations, we need to further specify how the approximator behaves. Let us first introduce two alternative ways of implementing the approximation: Either one common approximator is used across all observations, or a new approximator is used for each observation. To differentiate between the two approximation schemes, we will refer to the approximate estimator based on the first scheme as an *estimator based on common approximators* (ECA) and to the second one as an *estimator based on individual approximators* (EIA). In the first case, the approximate sample moment takes the form

$$\mathbf{ECA} : G_n(\theta, \hat{\gamma}_S) = \frac{1}{n} \sum_{i=1}^n g(z_i; \theta, \hat{\gamma}_S), \quad (9)$$

which uses a single approximator $\hat{\gamma}_S$ in the computation of the moment conditions. In the second case,

$$\mathbf{EIA} : G_n(\theta, \hat{\gamma}_S) = \frac{1}{n} \sum_{i=1}^n g(z_i; \theta, \hat{\gamma}_{i,S}), \quad (10)$$

and n approximators $\hat{\gamma}_{1,S}, \dots, \hat{\gamma}_{n,S}$ are used in the computation. We stress that the ECA and EIA are both targeting the same infeasible estimator; the only difference lies in how the approximators are used in the computation of the objective function.

Take simulation-based estimation for instance. Earlier papers (e.g. Laroque and Salanié 1989, McFadden 1989) used EIAs, and most papers on cross-sectional or panel data still do. ECAs were proposed by Lee (1992) for cross-sectional discrete choice models, but they have been more useful in dynamic models where one long trajectory of the model is simulated and used to compute simulated moments (see Example 3) or densities (see Example 1).

When the number of approximators remains fixed, as in ECAs, the resulting approximate estimator is similar to semiparametric two-step estimators where in the first step a function is nonparametrically estimated, see e.g. Andrews (1994) and Chen et al (2003). In contrast, EIAs employ n approximators—one for each observation. Thus, the number of approximators increases with sample size, and EIAs give rise to an incidental parameters problem. Some of our results for this situation are similar to those found in the literature on higher-order properties and bias-correction of estimators in an incidental parameters setting, see e.g. Hahn and Newey (2004) and Arellano and Hahn (2007).

To provide a streamlined set of regularity conditions that apply to both of these approximation schemes, we let $J \geq 1$ denote the number of approximators used in the computation of $\hat{\theta}_{n,S}$. For ECAs and EIAs, $J = 1$ and $J = n$, respectively. In what follows, it is crucial to separate assumptions on the bias of the approximator (which is common amongst the J approximators),

$$b_S(u; \theta) := E[\hat{\gamma}_{i,S}(u; \theta) | u] - \gamma_0(u; \theta), \quad (11)$$

from assumptions on its stochastic component (which, by definition, is zero for non-stochastic approximators):

$$\psi_{i,S}(u; \theta) := \hat{\gamma}_{i,S}(u; \theta) - E[\hat{\gamma}_{i,S}(u; \theta) | u], \quad i = 1, \dots, J. \quad (12)$$

A.5(p) The approximator(s) lies in Γ and, for any value of u : satisfies⁶:

- (i) The J ($= 1$ or $= n$) random functions $\hat{\gamma}_{1,S}(u; \theta), \dots, \hat{\gamma}_{J,S}(u; \theta)$ are identically distributed, mutually independent and independent of \mathcal{Z}_n .
- (ii) Their common bias b_S is of order $\beta > 0$, $b_S(u; \theta) = S^{-\beta} \bar{b}(u; \theta) + o(S^{-\beta})$.
- (iii) For $2 \leq q \leq p$, the stochastic component satisfies $E[\|\psi_{i,S}(u; \theta)\|^q] = S^{-\alpha_q} v_q(u; \theta) + o(S^{-\alpha_q})$, $i = 1, \dots, J$, for some constant $\alpha_q > 0$.

Assumption A.5 is sufficiently general to cover all of the examples in Section 2 under suitable regularity conditions. First consider A.5(iii). It requires that the approximator have p moments and that each of these vanish at a given rate as $S \rightarrow \infty$. We will choose p in conjunction with the order of the expansion m of Assumption A.4, since we wish to evaluate the mean and variance of each of the higher-order terms. For example, in order to ensure that the variance of $\nabla^k G_n[\hat{\gamma}_S, \dots, \hat{\gamma}_S]$ exists and to evaluate its rate of convergence, we will require A.5(p) to hold with $p = 2k$.

While the need to distinguish between common and independent approximators introduces a subscript i in the $\hat{\gamma}_{i,S}$ and $\psi_{i,S}$ terms, the reader need only remember that for ECA approximators, the i subscript can be disregarded, while for EIA approximators, the i subscript refers to the index i of the observation.

2.3.1 Non-stochastic approximators

First consider an approximation that does not involve any randomness, as with numerical integration, discretization, or numerical inversion of a function. A.5(i) clearly has no bite when non-stochastic approximators are used. Then by construction the conditional variance of the approximator is zero, so that $\alpha_p = +\infty$ for all $p \geq 2$. Non-stochastic approximation imparts a bias, which in leading cases obeys assumption A.5(ii) for some $\beta > 0$. As we will

⁶The $o(\cdot)$ terms in (ii)-(iii) are w.r.t. the function norm of Γ .

see, our general theory also applies in this case. However, with numerical approximators the bias b_S is often hard to correct for using analytical adjustments. This is one reason why we propose an alternative approach in section 4.

2.3.2 Stochastic approximators

Next, let us examine stochastic approximation schemes, which encompass simulation-based inference methods. The typical EIA simulation scheme takes n independent draws from the same distribution; this satisfies A.5(i). Our assumptions do not rule out dependence between the simulated values within each simulated sample, however. Dependence occurs, for example, in SMM and SMLE for time series models where a long trajectory of the model is simulated. In such ECA schemes⁷, only one approximator is used for all observations and so A.5(i) is automatically satisfied. Note that A.5(i) is stated for some *fixed* value of x ; the requirement that the simulations be independent of data is satisfied by most standard simulation schemes.

Monte Carlo schemes are of course the most prominent example of stochastic approximators. As our examples illustrate, many of the most useful approximate estimators rely on Monte Carlo approximators. We will therefore specialize some of our results to the following class of Monte Carlo approximators:

A.6(p) The approximator $\hat{\gamma}_{i,S}(u; \theta)$ takes the form $\hat{\gamma}_{i,S}(u; \theta) = \sum_{s=1}^S w_S(u, \varepsilon_{i,s}; \theta) / S$, $i = 1, \dots, J$, where: (i) $\{\varepsilon_{i,s}\}_{s=1}^S$ is stationary and geometrically β -mixing; (ii) $\{\varepsilon_{i,s}\}_{s=1}^S$ and $\{\varepsilon_{j,s}\}_{s=1}^S$ are independent for $i \neq j$, and they are all independent of the sample; (iii) the function $w_S(u, \varepsilon_{i,s}; \theta)$ satisfies, with expectations being taken w.r.t. $\varepsilon_{i,s}$,

$$\bar{w}_S(u; \theta) := E[w_S(u, \varepsilon_{i,s}; \theta) | u] = \gamma_0(u; \theta) + S^{-\beta} \bar{b}(u; \theta) + o\left(S^{-\beta}\right);$$

and for every $2 \leq q \leq p$

$$E[\|w_S(u, \varepsilon_{i,s}; \theta) - \bar{w}_S(u; \theta)\|^q | u] = O(S^{\mu_q}) \text{ for some } \mu_q < q/2.$$

To our knowledge, the class of approximators that satisfies A.6 includes all simulation-based approximators proposed in the literature, including Markov Chain Monte Carlo methods where $\varepsilon_{i,s}$, $s = 1, 2, \dots$, is a Markov chain that is designed so as to have stationary distribution equal to some unknown target distribution of interest. The assumption of β -mixing is only used in the proof of Theorem 6. It could be weakened to “strongly mixing” elsewhere, but we maintain the assumption of β -mixing throughout to streamline the assumptions. The bias, b_S , and variance, ψ_S , of approximators satisfying A.6 follow directly from those of the simulators w_S .

⁷See the discussion in Example 3.

Using Jensen's inequality, $E[\|X\|^q] \leq E[\|X\|^p]^{q/p}$ for $q \leq p$, we see that $\mu_q \leq q\mu_p/p$ for $2 \leq q \leq p$. Given this inequality, it is easily verified that Assumption A.6 implies A.5 with the same rate β for the bias term and with $\alpha_q = p/2 - \mu_q > 0$ in A.5(iii).

In parametric simulation-based estimation, the simulating function $w_S \equiv w$ in Assumption A.6 is typically independent of the number of simulations, and the approximator has no bias: $b_S \equiv 0$ and so $\beta = \infty$. Moreover, Assumption A.6(iii) typically holds with $\mu_p = 0$, and A.5(iii) with $\alpha_p = p/2$. The class of approximators in Assumption A.6 allows for nonparametric density estimation techniques such as the methods proposed in Fermanian and Salanié (2004) and Kristensen and Shin (2012). These have a bias component, so that $\beta < \infty$, and w_S depends on S through the bandwidth; but A.6 still applies (see section 6). It also allows for nonparametric regression estimators, such as those used in Creel and Kristensen (2012), Kristensen and Scherning (2012) and Norets (2009, 2012). Consider for instance the case where $m(x; \theta) = E_\theta[y_i|x_i = x]$ is used in the estimation but not available on closed form. We approximate m by

$$\hat{m}_S(x; \theta) = \frac{\sum_{s=1}^S y_s(\theta) K_h(x_s(\theta) - x)}{\sum_{s=1}^S K_h(x_s(\theta) - x)},$$

where $(y_s(\theta), x_s(\theta))$ are simulated draws. While \hat{m}_S itself does not satisfy A.6, we can express it as a function of the vector function $\hat{\gamma}_S(x; \theta) = \sum_{s=1}^S \bar{y}_s(\theta) K_h(x_s(\theta) - x) / S$, $\bar{y}_s(\theta) = (y_s(\theta)', 1)'$, where $\hat{\gamma}_S$ does satisfy A.6.

3 Effects of Approximations

We are now ready to derive the leading bias and variance terms of the estimator due to approximation errors. In the following, when we discuss biases and variances and, for example, write $E[\hat{\theta}_{n,S}]$, we refer to the means and variances of the leading terms of a valid stochastic expansion of the estimators. This is a standard approach in the higher-order analysis of estimators; see, for example, Rothenberg (1984) and Newey and Smith (2004, section 3).

3.1 Higher-order Expansion

To state the asymptotic expansion in a compact manner, we introduce some additional notation and moments which will make up the leading bias and variance terms. Let $g_i := g(z_i, \theta_0, \gamma_0)$; $\nabla g_i[d\gamma] := \nabla g(z_i, \theta_0, \gamma_0)[d\gamma]$ and $\nabla^2 g_i[d\gamma, d\gamma] := \nabla^2 g(z_i, \theta_0, \gamma_0)[d\gamma, d\gamma]$ for any function $d\gamma$. As we will see, the leading terms in the bias of the approximate estimator are

$$B_{S,1} = -H_0^{-1} E[\nabla g_i[b_S]] \quad \text{and} \quad B_{S,2} = -\frac{1}{2} H_0^{-1} E[\nabla^2 g_i[\psi_{i,S}, \psi_{i,S}]], \quad (13)$$

where b_S and $\psi_{i,S}$ are defined in eqs. (11)-(12). The first bias term $B_{S,1}$ is zero for unbiased approximators, as in parametric simulation-based inference. The second one, $B_{S,2}$, is zero for non-stochastic approximators of the type found in numerical approximation schemes.

The leading variance term due to the presence of approximations is $\nabla G_n(\theta_0)[\hat{\gamma}_S - \gamma]$. It can be decomposed into two terms. The first one is

$$D_{n,S} = \frac{1}{n} \sum_{i=1}^n d_{i,S}, \quad \text{with } d_{i,S} = \nabla g_i[b_S] - E \nabla g_i[b_S],$$

which is common to the two approximation schemes. The asymptotic properties of the second variance component,

$$E_{n,S} = \frac{1}{n} \sum_{i=1}^n \nabla g_i[\psi_{i,S}]$$

depend on whether we use EIA or ECA, however. The variance components $\psi_{i,S}$ vary across observations for EIAs; as a consequence, one can directly apply a CLT for stationary and mixing sequences to $E_{n,S}$. On the other hand, ECAs only have one ψ_S , which is common across observations; and getting a CLT takes more work and additional assumptions. We start by rewriting $E_{n,S}$ as

$$E_{n,S} = \frac{1}{n} \sum_{i=1}^n \{\nabla g_i[\psi_S] - \nabla G[\psi_S]\} + \nabla G[\psi_S], \quad \text{with } \nabla G[\psi_S] := E[\nabla g_i[\psi_S] | \psi_S].$$

The first term is $O_P(S^{-\alpha_2/2}/\sqrt{n})$, and so is dominated by the second term $\nabla G[\psi_S] = O_P(S^{-\alpha_2/2})$. In general, the large-sample distribution of $\nabla G[\psi_S]$ is not known in closed-form. However, if we strengthen Assumption A.5 to A.6, we can write

$$\nabla G[\psi_S] = \frac{1}{S} \sum_{s=1}^S \nabla G[e_{s,S}], \quad \text{with } e_{s,S}(\varepsilon_s) := w_S(u, \varepsilon_s; \theta_0) - E[w_S(u, \varepsilon_s; \theta_0)], \quad (14)$$

and a CLT can be applied as $S \rightarrow \infty$.

The above terms make up the first-order expansion of the effects of approximations on the estimators:

Theorem 2 *Assume A.1-A.3, A.4(2), and A.5(4). Then:*

$$\begin{aligned} \hat{\theta}_{n,S} - \theta_0 &= B_{S,1} + B_{S,2} + H_0^{-1} \{G_n + D_{n,S} + E_{n,S}\} \\ &\quad + O_P(S^{-3\beta}) + O_P(S^{-\alpha_3}) + o_P(1/\sqrt{n}), \end{aligned} \quad (15)$$

where $G_n := G_n(\theta_0, \gamma_0)$ and the two sequences $(G_n, D_{n,S})$ and $E_{n,S}$ are asymptotically mutually independent. Moreover, the following limit results hold as $n, S \rightarrow \infty$:

- For both EIA and ECA approximators,

$$\sqrt{n}(\Omega_S^{G+D})^{-1/2}\{G_n + D_{n,S}\} \xrightarrow{d} N(0, I_k), \text{ with } \Omega_S^{G+D} = \lim_{n \rightarrow \infty} \frac{1}{n} \text{Var} \left(\sum_{i=1}^n g_i + d_{i,S} \right)$$

$$\text{and } \Omega_S^{G+D} = \Omega^G + O(S^{-2\beta}) \text{ with } \Omega^G = \frac{1}{n} \text{Var}(\sum_{i=1}^n g_i).$$

- The bias terms have orders $B_{S,1} = O(S^{-\beta})$ and $B_{S,2} = O(S^{-\alpha_2})$.
- For EIA approximators, $\text{Var}(E_{n,S}) = O_P(S^{-\alpha_2}n^{-1})$; for ECA approximators, $\text{Var}(E_{n,S}) = O_P(S^{-\alpha_2})$.

For the family of approximators covered by Assumption 6, which covers most parametric simulation-based estimators, we can obtain a more precise characterization of the variance term $E_{n,S}$:

Corollary 3 *Assume that A.1-A.3, A.4(2), A.5(4), and A.6(4) hold, and $w = w_S$ does not depend on S . Then $\alpha_2 = 1$ and*

$$\mathbf{EIA} : \sqrt{nS}E_{n,S} \xrightarrow{d} N(0, \Omega_{EIA}^E), \text{ with } \Omega_{EIA}^E = \lim_{S \rightarrow \infty} \frac{1}{S} \text{Var} \left(\sum_{s=1}^S \nabla g_0[e_s] \right),$$

$$\mathbf{ECA} : \sqrt{S}E_{n,S} \xrightarrow{d} N(0, \Omega_{ECA}^E), \text{ with } \Omega_{ECA}^E = \lim_{S \rightarrow \infty} \frac{1}{S} \text{Var} \left(\sum_{s=1}^S \nabla G[\tilde{e}_s] \right),$$

where $e_s(u) = e_{s,S}(u) = w(u, \varepsilon_s; \theta_0) - E[w(u, \varepsilon_s; \theta_0)]$ is defined in eq. (14).

These expansions allow us to analyze the effects due to approximation errors. In particular, both EIA's and ECA's are normally distributed as $n, S \rightarrow \infty$ with leading bias and variance terms due to approximations given by:

$$E[\hat{\theta}_{n,S} - \theta_0] \simeq B_{S,1} + B_{S,2}, \quad \text{Var}(\hat{\theta}_{n,S} - \theta_0) \simeq H_0^{-1} \left\{ \Omega_S^{G+D}/n + \text{Var}(E_{n,S}) \right\} H_0^{-1}. \quad (16)$$

The bias and the variance of the approximator enter the two leading bias terms of the approximate estimator separately: the bias b_S drives $B_{S,1}$, and the stochastic components $\psi_{j,S}$ drive $B_{S,2}$. When the approximator is a simple unbiased simulated average, $B_{S,1} = 0$ and the leading bias term $B_{S,2} = O(1/S)$; this is a well-known result for specific simulation-based estimators in cross-sectional settings—see e.g. Gouriéroux and Monfort (1996), Lee (1995), or Corradi and Swanson (2007). Our theorem shows that this result holds more generally under weak regularity conditions.

EIA's and ECA's differ regarding the second variance term $E_{n,S}$. In the computation of the ECA, one common approximator is used across all observations; this introduces additional

correlations across observations. In contrast, for EIA, $\psi_{i,S}$ and $\psi_{j,S}$ are independent for $i \neq j$. As a consequence, the variance due to a given number S of simulations is larger for ECA's; and in leading simulation-based inference cases with $\beta = \infty$ and $\alpha_2 = 1$, we need S to go to infinity faster than n to keep the variance from exploding. This seems to suggest that one should prefer EIA to ECA; but statistical efficiency must be traded off with computational efficiency. If for instance $\hat{\gamma}_S$ is costly to implement, it may be convenient to use the same approximator across all observations. As with the bias components, this result generalizes the existing ones for SMLE and similar estimation techniques (see, e.g., Gouriéroux and Monfort (1996), Lee (1995), or Corradi and Swanson (2007)) to allow for more complex simulation-based techniques and for the approximator to enter the estimating equation in a more involved manner. Specifically, their results are as special cases of the above corollary. In Section 6 we show how the general theory can be used to derive a higher-order analysis of NPSMLE (Example 1) and SNLS (Example 2).

3.2 Sharpness

The sharpness of the rates in Theorem 2 depends on the type of approximator being used and how it enters into the objective function; that is, the precise nature of the mapping $\gamma \mapsto g(z, \theta, \gamma)$.

Theorem 4 *Under the assumptions of Theorem 2, if the rates in Assumption A.5 are sharp then*

- *For non-stochastic approximators, all rates listed in the Theorem are sharp.*
- *For EIA's with $\nabla^2 g_i[d\gamma, d\gamma] \neq 0$, the rates of $B_{S,1}$ and $B_{S,2}$ and $D_{n,S}$ and $E_{n,S}$ are sharp. If additionally Assumption A.6(4) holds with $w_S \equiv w$, the same is true for ECA's.*

The proof of Theorem 4 follows from the arguments in the proof of Theorem 2 together with rate results for sample averages. Note that it does not cover nonparametric simulators, for which w_S depends on S through the bandwidth. If for instance $\hat{\gamma}_S$ is a kernel estimator and ECA is used, one can show that $\text{Var}(E_{n,S}) = O(S^{-1})$. Since $\alpha_2 < 1$ in this case, this bound is sharper than the rate stated in the theorem; see Creel and Kristensen (2012) and Kristensen and Shin (2012).

In some special cases, a term in the expansion is zero. In SMM for instance, the function g is linear in the approximator γ . Then $\nabla^2 g_i[d\gamma, d\gamma] = 0$, so that $B_{S,2} = 0$; and our rates are obviously not sharp.

As our discussion of parametric simulation-based estimators showed, often some of the terms in the expansion will be zero. For methods with a nonparametric component, such

as NPSML, all of the terms may be simultaneously nonzero. This follows directly from the coexistence of bias and variance in nonparametric smoothers; see Section 6.1.

3.3 First-order efficiency

Our results allow us to provide rates on the degree of approximation under which the approximate estimator is asymptotically first-order equivalent to the exact estimator; that is, which choices of the sequence $S = S_n$ guarantee $\|\hat{\theta}_{n,S_n} - \hat{\theta}_n\| = o_P(n^{-1/2})$. In general, asymptotic equivalence for ECAs obtain if $n/S^{\min(\alpha_2, 2\beta)} \rightarrow 0$; for EIA's we have a weaker condition, replacing α_2 with $2\alpha_2$.

For parametric simulation-based estimators ($\beta = \infty$ and $\alpha_2 = 1$), this gives the standard result that n/S_n should go to zero for ECA's (Duffie and Singleton, 1993; Lee, 1995, Theorem 1), while \sqrt{n}/S_n should go to zero for EIA's (Laroque and Salanié, 1989; Lee, 1995, Theorem 4). Section 6.1 takes up the more complicated case of nonparametric kernel methods, as used in NPSML.

4 Newton-Raphson Adjustment

We first propose a very general method that can reduce both bias and variance of the approximate estimator in a simple manner. It works with non-stochastic approximations as well as with stochastic approximations. Our proposal builds on the well-known result that a consistent estimator can be made asymptotically efficient by applying one Newton-Raphson (NR) step of the log-likelihood function to it. E.g. if $\hat{\theta}_n$ is a \sqrt{n} -consistent estimator of θ_0 in a model with log-likelihood $L_n(\theta)$, then a single NR-step yields a consistent and asymptotically efficient estimator. We extend this idea to our setting by starting from some initial approximate estimator based on a degree of approximation S , say $\bar{\theta}_{n,S}$. We then define the corrected estimator through one or possibly several Newton-Raphson iterations of an approximate objective function that uses a much finer approximation, $S^* \gg S$. With $H_n(\theta, \gamma) = \partial G_n(\theta, \gamma) / (\partial \theta)$, we define iteratively

$$\hat{\theta}_{n,S}^{(k+1)} = \hat{\theta}_{n,S}^{(k)} - H_n^{-1}(\hat{\theta}_{n,S}^{(k)}, \hat{\gamma}_{S^*}) G_n(\hat{\theta}_{n,S}^{(k)}, \hat{\gamma}_{S^*}), \quad k = 1, 2, 3, \dots \quad (17)$$

where $\hat{\theta}_{n,S}^{(1)} = \bar{\theta}_{n,S}$ is some initial estimator and we use the S^* th order approximator, $\hat{\gamma}_{S^*}$, in the iterations.

Note that the cost of computing each new iterate from the previous one is (very) roughly S^*/S times the cost of one iteration in the minimization of $Q_n(\theta, \hat{\gamma}_{S^*})$. Since the minimization itself can easily require a hundred iterations or so, we can therefore take S^* ten or twenty times larger than S without adding much to the cost of the estimation procedure⁸.

⁸In many cases, many of the dimensions of θ only come into play within some linear indexes $\theta'x$; then the

To evaluate the performance of $\hat{\theta}_{n,S}^{(k+1)}$ relative to $\bar{\theta}_{n,S^*}$, we first note that

$$\|\hat{\theta}_{n,S}^{(k+1)} - \hat{\theta}_n\| \leq \|\hat{\theta}_{n,S}^{(k+1)} - \bar{\theta}_{n,S^*}\| + \|\bar{\theta}_{n,S^*} - \hat{\theta}_n\|.$$

Combining this with Robinson (1988, Theorem 2), we obtain the following theorem:

Theorem 5 *Assume that A.1-A.3, A.4(3) and A.5(6) hold. Let the initial estimate $\bar{\theta}_{n,S}$ be chosen as either $\hat{\theta}_{n,S}$, $\hat{\theta}_{n,S}^{\text{AB}}$, or $\hat{\theta}_{n,S}^{\text{JK}}$. Then the NR-estimator $\hat{\theta}_{n,S}^{(k+1)}$ defined in (17) satisfies:*

$$\|\hat{\theta}_{n,S}^{(k+1)} - \hat{\theta}_n\| = O_P(\|\bar{\theta}_{n,S} - \hat{\theta}_n\|^{2^k}) + O_P(\|\bar{\theta}_{n,S^*} - \hat{\theta}_n\|) \quad (18)$$

as n, S and S^* go to infinity with $S^* > S$.

The above result formalizes the intuition that a large enough number of NR-steps with the score and Hessian evaluated at γ_{S^*} yields an estimator that is equivalent to the extremum estimator obtained from full optimization of the objective function based on γ_{S^*} . This holds irrespective of the convergence rate of the initial estimator. However, the number of NR iterations, k , needed to obtain this result does depend on the precision of the initial estimator. For unadjusted parametric simulation-based estimators in the EIA scheme for instance, we know from Theorem 2 that $\|\bar{\theta}_{n,S} - \hat{\theta}_n\| = O_P(1/S)$. Then the first term on the right-hand side of the inequality in Theorem 5 is asymptotically dominated by the second term if $S^* = o(S^{2^k})$. Taking $k = 1$ and having S^*/S converge to some positive number would be enough in this case.

The above iterative estimator requires computation of the Hessian, $H_n(\theta, \hat{\gamma}_S)$. If this is not feasible or computationally burdensome, an approximation can be employed, e.g. numerical derivatives. This however will slow down the convergence rate and the result of Theorem 5 has to be adjusted, cf. Robinson (1988, Theorem 5). In particular, more iterations are required to obtain a given level of precision.

As a partial alternative to Newton-Raphson iterations, resampling methods could be used⁹. They will in general handle the biases due to both the stochastic and the non-stochastic component of the approximator; and the researcher is not required to derive an expression of the bias. On the other hand, they are computationally more demanding than the analytical bias correction proposed in the next section, and they may lead to an increase in finite-sample variance.

To motivate bias adjustment via resampling, recall from Theorem 2 that $E[\hat{\theta}_{n,S} - \hat{\theta}_n] \simeq b_1 S^{-\beta} + b_2 S^{-\alpha_2}$. First compute two approximators of order S^* which we denote $\hat{\gamma}_{S^*}^{[1]}$ and $\hat{\gamma}_{S^*}^{[2]}$.

trade off is even more favourable since the computation of the second derivative H_n is much simplified.

⁹See Hahn and Newey (2004) and Dhaene and Jochmans (2015) for bias correction using Jackknife in the context of panel models, and Phillips and Yu (2005) for a time series application.

Let $\hat{\theta}_{n,S^*}^{[m]}$ be the estimator based on the same data sample \mathcal{Z}_n but using the m th approximator $\hat{\gamma}_{S^*}^{[m]}$, $m = 1, 2$. Then consider the following jackknife (JK) type estimator:

$$\hat{\theta}_{n,S}^{\text{JK}} := 2\hat{\theta}_{n,S} - \frac{1}{2}\{\hat{\theta}_{n,S^*}^{[1]} + \hat{\theta}_{n,S^*}^{[2]}\}. \quad (19)$$

It is easy to see that

$$\begin{aligned} E[\hat{\theta}_{n,S}^{\text{JK}} - \hat{\theta}_n] &= 2E[\hat{\theta}_{n,S} - \hat{\theta}_n] - \frac{1}{2}\{E[\hat{\theta}_{n,S^*}^{[1]} - \hat{\theta}_n] + E[\hat{\theta}_{n,S^*}^{[2]} - \hat{\theta}_n]\} \\ &\simeq b_1\{2S^{-\beta} - (S^*)^{-\beta}\} + b_2\{2S^{-\alpha} - (S^*)^{-\alpha}\}, \end{aligned}$$

where we ignored higher-order terms. We would now ideally choose S^* such that both of the above bias terms cancel out. However, we can only remove either of the two: By choosing either

$$S^* = \frac{S}{2^{1/\beta}} \text{ or } S^* = \frac{S}{2^{1/\alpha_2}}, \quad (20)$$

we will remove the first or the second term respectively. Obviously, S^* should be chosen so as to remove the bias component that dominates in the expansion. In a previous version we also reported results for this resampling method; and we tested it on the mixed logit model that we explore in section 7. We found that the improvements from resampling were dominated by those obtained with the other methods.

5 Analytical Adjustments

The expansions derived in section 3 naturally suggest correcting the approximate estimators and standard errors to take into account the biases and variances due to approximations. The corrections are obtained by constructing consistent estimators of the leading terms in the formulæ of Theorem 2, and Corollary 3 when applicable.

5.1 Bias Adjustment

The leading bias terms are $B_{S,1}$ and $B_{S,2}$. We mainly focus on the case where $\beta > \alpha_2$. Recall that this includes parametric simulation-based estimation methods, but it excludes most purely non-stochastic approximators. Then $B_{S,1}$ is of lower order and the leading bias component is $B_{S,2} = -\frac{1}{2}H_0^{-1}E[\nabla^2 g(z_i, \theta_0, \gamma_0)[\psi_S, \psi_S]]$.

We wish to adjust the approximate estimator to remove this bias component. The two main approaches to bias adjustment in the econometric literature are “corrective” and “preventive”¹⁰. The *corrective method* first computes the unadjusted estimator, $\hat{\theta}_{n,S}$, obtains

¹⁰See Sections 3 and 4, respectively, in Arellano and Hahn (2007) for a discussion of corrective and preventive adjustments in a panel data setting.

a consistent estimator of the bias, here $B_{S,2}$, and then combines the two to obtain a new, bias-adjusted (BA) estimator

$$\hat{\theta}_{n,S}^{\text{BA}} = \hat{\theta}_{n,S} - \hat{B}_{S,2}.$$

One example of this approach for can be found in Lee (1995) for the special case of SMLE and SNLS in limited dependent variable models. A natural estimator of $\hat{B}_{S,2}$ would be $\hat{B}_{S,2} = -\frac{1}{2}\hat{H}_n^{-1}\nabla^2\hat{G}_n(\hat{\theta}_{n,S})$ for some consistent estimator $\nabla^2\hat{G}_n(\theta)$ of $\nabla^2G(\theta) := E[\nabla^2g(z_i; \theta, \gamma_0)[\psi_{i,S}, \psi_{i,S}]]$. We propose two different estimators depending on whether A.6 holds or not. If A.6 does not apply, the following estimator is available for EIA:

$$\mathbf{EIA} : \nabla^2\hat{G}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla^2g(z_i; \theta, \hat{\gamma}_S)[\hat{\psi}_{i,S}, \hat{\psi}_{i,S}], \quad \hat{\psi}_{i,S} := \hat{\gamma}_{i,S} - \frac{1}{n} \sum_{i=1}^n \hat{\gamma}_{i,S}. \quad (21)$$

For the ECA version, one cannot estimate the variance component of $\hat{\gamma}_S$ without further simulations. One possibility would be to simulate m extra, mutually independent versions, $\hat{\gamma}_{k,S}$, $k = 1, \dots, m$, of $\hat{\gamma}_S$, and then compute $\nabla^2\hat{G}_n(\theta) = \frac{1}{nm} \sum_{i=1}^n \sum_{k=1}^m \nabla^2g(z_i; \theta, \hat{\gamma}_{k,S})[\hat{\psi}_{k,S}, \hat{\psi}_{k,S}]$, where $\hat{\psi}_{k,S} = \hat{\gamma}_{k,S} - \frac{1}{m} \sum_{k=1}^m \hat{\gamma}_{k,S}$. Here, m has to be chosen large enough so that the variance component of $\nabla^2\hat{G}_n(\theta)$ does not dominate the bias that we are trying to remove. This means that the computational cost of this first ECA bias estimator can be large, especially if $\hat{\gamma}_S$ is not easy to compute.

When A.6 also holds, the following alternative estimator is available; and it can be used for both ECA's and EIA's:

$$\nabla^2\hat{G}_n(\theta) = \frac{1}{nS(S-1)} \sum_{i=1}^n \sum_{s=1}^S \nabla^2g(z_i; \theta, \hat{\gamma}_{i,S})[\hat{e}_{i,s,S}, \hat{e}_{i,s,S}]; \quad (22)$$

here, in the case of EIA's, $\hat{e}_{i,s,S}(x; \theta) = w_S(x, \varepsilon_{i,s}; \theta) - \hat{\gamma}_{i,S}(x; \theta)$ while in the case of ECA, $\hat{e}_{i,s,S}(x; \theta) = w_S(x, \varepsilon_s; \theta) - \hat{\gamma}_S(x; \theta)$ and so does not change across observations $i = 1, \dots, n$.

Instead of adjusting the estimator, we can do *preventive correction* where we adjust the estimating equation $G_n(\theta, \hat{\gamma}_S)$ to remove the component leading to the bias $B_{S,2}$. By inspection of the proof of Theorem 2, it is easily seen that the relevant adjustment of $G_n(\theta, \hat{\gamma}_S)$ is $E[\nabla^2g(z_i; \theta)[\psi_{i,S}, \psi_{i,S}]]/2$. This suggests a bias-adjusted estimator $\hat{\theta}_{n,S}^{\text{BA}}$ that solves

$$G_n(\hat{\theta}_{n,S}^{\text{BA}}, \hat{\gamma}_S) - \frac{1}{2}\nabla^2\hat{G}_n(\hat{\theta}_{n,S}^{\text{BA}}) = o_P(1/\sqrt{n}), \quad (23)$$

where $\nabla^2\hat{G}_n(\theta)$ is taken either from eq. (21) or (under A6) from eq. (22). This approach was pursued in the context of SNLS (see Example 2) by Laffont et al (1995).

After either preventive or corrective adjustment, the bias component $B_{S,2}$ is replaced by

$$\tilde{B}_{S,2} := -\frac{1}{2}H_0^{-1}E[\nabla^2G_n(\theta_0)[\psi_S, \psi_S] - \nabla^2\hat{G}(\theta_0)]. \quad (24)$$

The following theorem analyzes the properties of the bias adjusted estimator based on $\nabla^2 \hat{G}_n(\theta)$ given in eq. (22). We expect similar results to hold for any bias adjusted EIA estimator that uses eq. (21).

Theorem 6 *Assume that A.1-A.3, A.4(3), and A.6(8) hold together with*

$$\|\nabla^2 g(z; \theta_0)[e_{is}, e_{it}]\| \leq b(z) \|e_{is}(z)\| \|e_{it}(z)\|,$$

where $E[b^8(z)] < \infty$. Then any $\hat{\theta}_{n,S}^{\text{BA}}$ solving eq. (23) with $\nabla^2 \hat{G}_n(\theta)$ defined in (22) satisfies as $n, S \rightarrow \infty$:

$$\begin{aligned} \hat{\theta}_{n,S} - \theta_0 &= B_{S,1} + \tilde{B}_{S,2} + H_0^{-1} \{G_n + D_{n,S} + E_{n,S}\} \\ &\quad + O_P(S^{-3\beta}) + O_P(S^{-2+\mu_4}) + O(S^{-2+\mu_3}) + o_P(1/\sqrt{n}), \end{aligned}$$

where the new bias term given in eq. (24) satisfies $\tilde{B}_{S,2} = O(S^{-2+\mu_2})$ and $\mu_p, p \geq 2$, is defined in A.6. All other terms in the expansion are exactly as in Theorem 2.

Note that under the assumptions of Theorem 6, $-2 + \mu_4 < 0$, $-2 + \mu_3 < -1/2$ and $-2 + \mu_2 < -1$. The theorem therefore shows that under slightly stronger conditions¹¹ than in Theorem 2, $\tilde{B}_{S,2}$ has a faster rate of convergence than $B_{S,2}$, while the rate of the other leading terms is unchanged. More precisely, compared to Theorem 2, the bias term $B_{S,2} = O(S^{-\alpha_2}) = O(S^{-1+\mu_2})$ has been replaced by $\tilde{B}_{S,2} = O(S^{-2+\mu_2})$. Also note that the higher-order variance component of order $O_P(S^{-\alpha_3})$ that appeared in Theorem 2 has been replaced by $O_P(S^{-2+\mu_4}) + O(S^{-2+\mu_3})$. In the proof, we show that the variance of $\nabla^2 \hat{G}_n(\theta)$, that we use to estimate $B_{S,2}$, is of order $O_P(n^{-1/2}S^{-1+\mu_8/4}) + O_P(n^{-1/2}S^{-1+\alpha_4/2}) = o_P(1/\sqrt{n})$. In particular, the additional variances that we introduce when estimating the bias are of smaller order than the bias being adjusted for and so the bias adjusted estimator dominates the unadjusted one.

With unbiased simulators, we have $\mu_2 = 0$ and $\beta = \infty$, and by Theorem 2 the leading bias term of the unadjusted estimator is of order $O(S^{-1})$. Theorem 6 shows that for the adjusted estimator the leading term of the bias is of order $O(S^{-2})$. The improvement is by a factor S and may be quite large. More generally, the proposed adjustment will remove the largest bias component as long as $\alpha_2 < \beta$. Otherwise the bias term $O_P(S^{-\beta})$ is of a larger order than $O_P(S^{-\alpha_2})$ and the proposed bias adjustment does not remove the leading term anymore. In particular, when non-stochastic approximations are employed the above adjustment does not help. If we could estimate b_S , then $B_{S,1}$ could be taken care of easily by adjusting either estimator or estimating equation using $\nabla \hat{G}_n(\theta) := \sum_{i=1}^n \nabla g_i(\theta, \hat{\gamma}_S)[\hat{b}_S]/n$. However, estimating b_S can be a difficult task.

¹¹The higher order on A.6 is required to ensure that in the asymptotic expansion, the remainder term $R_{n,S}$ is still dominated.

5.2 Adjusting Standard Errors

If the approximator is stochastic, the approximate estimator will not only be biased; it will also contain additional variance terms, c.f. eq. (16). For a given sample size n and number of simulations S , we should adjust inferential tools (such as standard errors and t -statistics) to account for these additional variances. This turns out to be quite straightforward in many cases. To keep the notation simple, we assume in the following that data is i.i.d.¹².

The different terms appearing in the variance expansion in eq. (16) implicitly depend on θ_0 and γ_0 . In standard estimation procedures, one would usually estimate the above variance components by simply replacing θ_0 and γ_0 by $\hat{\theta}_{n,S}$ and $\hat{\gamma}_S$, respectively, in the expressions of Ω^{G+D} , $\text{Var}(E_{n,S})$ and H , and by replacing any population means by their sample counterparts. The variance term Ω_S^{G+D} involves the bias component of the approximator, b_S . This is unknown in most cases, but we know from Theorem 2 that $\Omega_S^{G+D} = \Omega^G + O(S^{-2\beta})$ where $\Omega^G = E[g(z, \gamma_0)g(z, \gamma_0)']$. For large S , a simple estimator would therefore be

$$\hat{\Omega}_S^{G+D} = \hat{\Omega}^G = \frac{1}{n} \sum_{i=1}^n \hat{g}_i \hat{g}_i', \text{ where } \hat{g}_i = g(z_i; \hat{\theta}_{n,S}, \hat{\gamma}_S).$$

However, replacing γ_0 by $\hat{\gamma}_S$ will generate biases. Similarly, if $\hat{\theta}_{n,S}$ has not been bias adjusted, replacing θ_0 by $\hat{\theta}_{n,S}$ will add biases to the variance estimator. Specifically, under suitable regularity conditions and by the same arguments as employed in the proof of Theorem 2,

$$E[\hat{\Omega}^G] = \Omega^G + O(S^{-\beta}) + O(S^{-\alpha_2}). \quad (25)$$

Recall that either $\text{Var}(E_{n,S}) = O_P(S^{-\alpha_2}n^{-1})$ (EIA) or $\text{Var}(E_{n,S}) = O_P(S^{-\alpha_2})$ (ECA), and so the biases in eq. (25) will often be of the same order as the variance components that we are trying to adjust for.

We therefore propose a bias-adjusted estimator of $\hat{\Omega}^G$ to improve on the basic variance estimators in the same way that we bias-adjusted $G_n(\theta, \hat{\gamma}_S)$. We assume in the following that $\hat{\theta}_{n,S}$ has already been bias adjusted so that we only need to adjust any biases due to $\hat{\gamma}_S$. This adjustment takes the form $\tilde{\Omega}_{\text{BA}}^G = \hat{\Omega}^G - \hat{\Delta}_{n,S}^\Omega$ where either, in the case of EIA's with $\hat{\psi}_{i,S} := \hat{\gamma}_{i,S} - \bar{\gamma}_S$,

$$\mathbf{EIA} : \hat{\Delta}_{n,S}^\Omega = \frac{1}{n} \sum_{i=1}^n \left\{ \nabla^2 \hat{g}_i[\hat{\psi}_{i,S}, \hat{\psi}_{i,S}] \hat{g}_i' + 2 \nabla \hat{g}_i[\hat{\psi}_{i,S}] \nabla \hat{g}_i[\hat{\psi}_{i,S}]' + \hat{g}_i \nabla^2 \hat{g}_i[\hat{\psi}_{i,S}, \hat{\psi}_{i,S}]' \right\},$$

¹²Otherwise long-run variance estimators have to be used.

or, under Assumption A.6 for both EIA and ECA,

$$\hat{\Delta}_{n,S}^{\Omega} = \frac{1}{nS(S-1)} \sum_{i=1}^n \sum_{s=1}^S \{ \nabla^2 \hat{g}_i[\hat{e}_{i,s,S}, \hat{e}_{i,s,S}] \hat{g}'_i + 2 \nabla \hat{g}_i[\hat{e}_{i,s,S}] \nabla \hat{g}_i[\hat{e}_{i,s,S}]' + \hat{g}_i \nabla^2 \hat{g}_i[\hat{e}_{i,s,S}, \hat{e}_{i,s,S}]' \};$$

here $\hat{e}_{i,s,S}$ is defined as right after eq. (22). The analysis of this estimator proceeds as in the proof of Theorem 6.

Next consider $\text{Var}(E_{n,S})$. As we know from Theorem 2, the behaviour of this term depends on whether EIA or ECA are used. In the case of EIA, $\text{Var}(E_{n,S}) \simeq \text{Var}(\nabla g_i[\psi_{i,S}])/n$ which can be estimated by $\widehat{\text{Var}}(E_{n,S}) = \sum_{i=1}^n \nabla g_i[\psi_{i,S}] \nabla g_i[\psi_{i,S}]' / n^2$.

When ECA is employed, $\text{Var}(E_{n,S}) \simeq \text{Var}(\nabla G[\psi_S])$ which can be estimated by $\widehat{\text{Var}}(E_{n,S}) = \sum_{k=1}^m \nabla \hat{G}[\psi_{S,k}] \nabla \hat{G}[\psi_{S,k}]' / m$, where $\psi_{S,k} = \hat{\gamma}_{S,k} - \sum_{k=1}^m \hat{\gamma}_{S,k} / m$, $\hat{\gamma}_{S,k}$, $k = 1, \dots, m$, are $m \geq 1$ independent versions of $\hat{\gamma}_S$ distribution of ψ_S , and $\nabla \hat{G}[d\gamma] = \sum_{i=1}^n \nabla \hat{g}_i[d\gamma] / n$. This can be computationally expensive if $\hat{\gamma}_S$ is a costly function.

The proposed estimators will suffer from biases similar to the ones in $\hat{\Omega}^G$, but these biases are of smaller order compared to the variance adjustment that we are making.

If Assumption A.6 holds, better estimates can be obtained since in this case Theorem 2 yields either $\text{Var}(E_{n,S}) \simeq \Omega_{\text{EIA}}^E / (nS)$ (EIA) or $\text{Var}(E_{n,S}) \simeq \Omega_{\text{ECA}}^E / S$ (ECA) where $\Omega_{\text{EIA}}^E = \text{Var}(\nabla g_i[\bar{w}_{i,S}])$ and $\Omega_{\text{ECA}}^E = \text{Var}(\nabla G[\bar{w}_S])$, and we have assumed for simplicity that the simulations are independent across $s = 1, \dots, S$.¹³ This suggests the following simple estimators

$$\hat{\Omega}_{\text{EIA}}^E = \frac{1}{nS} \sum_{i=1}^n \sum_{s=1}^S \nabla \hat{g}_i[\hat{e}_{i,s,S}] \nabla \hat{g}_i[\hat{e}_{i,s,S}]' \text{ and } \hat{\Omega}_{\text{ECA}}^E = \frac{1}{S} \sum_{s=1}^S \nabla \hat{G}[e_{s,S}] \nabla \hat{G}[e_{s,S}]',$$

where $\nabla \hat{G}[\gamma] = \sum_{i=1}^n \nabla g(z_i, \hat{\theta}_{n,S}, \hat{\gamma}_S)[\gamma] / n$ is a consistent estimator of $\nabla G(\theta_0, \gamma_0)[\gamma]$. The estimator $\hat{\Omega}_{\text{ECA}}^E$ is similar to the one proposed in Newey (1994) for semiparametric two-step estimators.

For EIA, two terms cancel out when we combine the bias adjustment $\hat{\Delta}_{n,S}^{\Omega}$ with $\hat{\Omega}_{\text{EIA}}^E$, giving

$$\hat{\Delta}_{n,S}^{\Omega} - \hat{\Omega}_{\text{EIA}}^E = \frac{1}{nS^2} \sum_{i=1}^n \sum_{s=1}^S \{ \nabla^2 \hat{g}_i[\hat{e}_{i,s,S}, \hat{e}_{i,s,S}] \hat{g}'_i + \nabla \hat{g}_i[\hat{e}_{i,s,S}] \nabla \hat{g}_i[\hat{e}_{i,s,S}]' + \hat{g}_i \nabla^2 \hat{g}_i[\hat{e}_{i,s,S}, \hat{e}_{i,s,S}]' \}.$$

Finally, the naive estimator of H_0 takes the form $\hat{H} = \sum_{i=1}^n h(z_i; \hat{\theta}_{n,S}; \hat{\gamma}_{i,S}) / n$. One could bias-adjust this estimator as we did for $\hat{\Omega}^G$. However, note that the approximate estimator

¹³Otherwise a HAC estimator has to be employed.

satisfies:

$$0 = \frac{1}{n} \sum_{i=1}^n g(z_i, \hat{\theta}_{n,S}, \hat{\gamma}_{i,S}) + \left\{ \frac{1}{n} \sum_{i=1}^n h(z_i; \bar{\theta}_{n,S}; \hat{\gamma}_{i,S}) \right\} (\hat{\theta}_{n,S} - \theta_0).$$

So in order to get a precise approximation of the distribution of $\hat{\theta}_{n,S} - \theta_0$, we want to use an estimator that mimics the behaviour of $n^{-1} \sum_{i=1}^n h(z_i; \bar{\theta}_{n,S}; \hat{\gamma}_{i,S})$. This is exactly what \hat{H} does; and we can still use it as an estimator of H_0 .

To sum up, for EIA, we propose the following bias-adjusted variance estimator for $\hat{\theta}_{n,S}$,

$$\hat{H}^{-1} \left(\frac{1}{n} \sum_{i=1}^n \left(\hat{g}_i \hat{g}_i' - \frac{1}{S^2} \sum_{s=1}^S (\nabla^2 \hat{g}_i[\hat{e}_{i,S}, \hat{e}_{i,S}] \hat{g}_i' + \nabla \hat{g}_i[\hat{e}_{i,S}] \nabla \hat{g}_i[\hat{e}_{i,S}]' + \hat{g}_i \nabla^2 \hat{g}_i[\hat{e}_{i,S}, \hat{e}_{i,S}]') \right) \right) \hat{H}^{-1},$$

while for ECA it takes the form $\hat{H}^{-1}(\hat{\Omega}^G - \hat{\Delta}_{n,S}^\Omega + \hat{\Omega}_{\text{ECA}}^E)H^{-1}$.

6 Applications

We now return to the first two examples of section 2.1 to illustrate the application of our expansion and proposed analytical adjustments.

6.1 Example 1: simulated maximum likelihood

For SML we approximate the density $p(z; \theta)$ ($= \gamma_0$) so that $g(z; \theta, p) = \dot{p}(z; \theta) / p(z; \theta)$, where $\dot{p}(z; \theta) = \partial p(z; \theta) / (\partial \theta)$, is the score of the log-likelihood. Then, suppressing dependence on $(z; \theta)$,

$$\nabla g[dp] = \frac{d\dot{p}}{p} - \frac{\dot{p}}{p^2} dp \quad \text{and} \quad \nabla^2 g[dp, dp] = -\frac{2}{p^2} dp d\dot{p} + \frac{2\dot{p}}{p^3} (dp)^2, \quad (26)$$

so that \bar{G}_0 in A.4 involves higher-order moments of $1/p$. If the density $p(z; \theta_0) \rightarrow 0$ as $\|z\| \rightarrow \infty$, these moments may not be finite. One can introduce trimming, replacing the simple simulator $\hat{p}_S(z; \theta)$ described above with $\hat{p}_{a,S}(z; \theta) = \hat{p}_S(z; \theta) \tau_a(\hat{p}_S(z; \theta))$ where $\tau_a(w)$ is a smooth trimming function that satisfies $\tau_a(w) = 1$ for $w \geq 2a$ and $\tau_a(w) = 0$ for $w \leq a$. Then $\bar{G}_{a,0} = O(a^{-(m+1)})$ is finite for any $a > 0$, and the remainder term satisfies $R_{n,S} = O_P(a^{-(m+1)} \|\hat{p}_{a,S} - p\|^{m+1})$. By letting $a = a_S \rightarrow 0$ at a suitable rate as $S \rightarrow \infty$, it is now possible to control the remainder term while the expansion remains valid; see Creel and Kristensen (2012) and Kristensen and Shin (2012) for more details in the context of SMM and SMLE, respectively.

The analytical adjustments are easy to compute when the approximator \hat{p}_S satisfies A.6 with $\beta = \infty$. Assume for instance that it uses independent simulations (the EIA case.)

Denoting $r_{i,s}(\theta) = w_S(x_i, \varepsilon_{i,s}; \theta) - \hat{p}_S(x_i; \theta)$, we obtain

$$\nabla^2 \hat{G}_n(\theta) = \frac{2}{nS} \sum_{i=1}^n \left(\frac{\dot{p}_{i,S}}{\hat{p}_{i,S}^3} \frac{1}{S} \sum_{s=1}^S r_{i,s}^2 - \frac{1}{\hat{p}_{i,S}^2} \frac{1}{S} \sum_{s=1}^S r_{i,s} \dot{r}_{i,s} \right).$$

Our proposed analytical adjustment to the variance of the estimators replaces the standard variance estimator, $n^{-1} \sum_{i=1}^n \dot{p}_{i,S} \dot{p}'_{i,S} / \hat{p}_{i,S}^2$, by

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{\dot{p}_{i,S} \dot{p}'_{i,S}}{\hat{p}_{i,S}^2} - \frac{1}{S^2} \sum_{s=1}^S \left(\frac{\dot{r}_{i,s}^2}{\hat{p}_{i,S}^2} + 9 \frac{r_{i,s}^2 \dot{p}_{i,S} \dot{p}'_{i,S}}{\hat{p}_{i,S}^4} - 4 \frac{r_{i,s} (\dot{r}_{i,s} \dot{p}'_{i,S} + \dot{p}'_{i,S} \dot{r}_{i,s})}{\hat{p}_{i,S}^3} \right) \right).$$

It is sometimes not possible to obtain an unbiased simulator of a density; then the NPSML estimator offers an attractive alternative. Suppose that the model takes the form $y = m(x, \varepsilon, \theta_0)$ and we compute $y_s(x, \theta) = m(x, \varepsilon_s, \theta_0)$, $s = 1, \dots, S$. The nonparametric simulated density then satisfies A.6 with $w_S(y, x, \varepsilon_s; \theta) = K_h(y_s(x, \theta) - y)$ where the bandwidth $h = h(S) \rightarrow 0$ as $S \rightarrow \infty$. Let $d = \dim(y)$ and suppose that we use a kernel of order r . The bias component satisfies $\bar{w}_S(y, x; \theta) - p(y|x; \theta) = h^r \frac{\partial^r p(y|x; \theta)}{\partial y^r} + o(h^r)$. Furthermore, it is easily checked that $E[|K_h(y_s(x, \theta) - x)|^p | x] = O(-h^{d(p-1)})$ for all $p \geq 2$ under suitable regularity conditions. Thus, with a bandwidth of order $h \propto S^{-\delta}$ for some $\delta > 0$, A.6(p) holds with $\beta = r\delta$ and $\mu_p = \delta d(p-1)$ for $p \geq 2$. We only need to choose $\delta < p/(2d(p-1))$ so that $\mu_p < p/2$.

As is well-known, the asymptotic mean integrated squared error is smallest when the bias and variance component are balanced. This occurs when $\delta^* = 1/(2r+d)$, leading to $\beta = r/(2r+d)$. It is easy to check that these values satisfy A6(p) if $r > d(p-2)/(2p)$, which allows for the standard choice of $r = 2$ except in implausibly high-dimensional cases. We recover of course the standard nonparametric rate.¹⁴

Let us now return to first-order efficiency. Using standard arguments from the literature on semiparametric estimation, one can show in great generality that $\Omega_{ECA}^E = O(S^{-1})$ (see Kristensen and Shin, 2012 for further details). Given this result, it easily follows from Theorem 2 that for the NPSMLE based on ECA's to be equivalent to the MLE, we need $\sqrt{nh}^r \rightarrow 0$, $n/S \rightarrow 0$ and $\sqrt{n}/(Sh^d)^2 \rightarrow 0$. For EIA's, $\Omega_{EIA}^E = O(1/(nSh^{d+2}))$ and so $n/S \rightarrow 0$ has to be replaced by $\sqrt{n}Sh^{d+2} \rightarrow \infty$.

We derive in Appendix C the analytical adjustments for such an NPSML estimator when y is scalar ($d = 1$) and the data is i.i.d. Given some additional regularity conditions, we

¹⁴While the standard nonparametric rate is optimal for the approximation of the individual densities that make up the likelihood, this rate does not yield the best NPSML estimators. This is akin to results for semi-parametric two-step estimators where undersmoothing of the first-step nonparametric estimator is normally required for the parametric estimator to be \sqrt{n} -consistent; see Kristensen-Salanié (2010) for details. For example, the optimal rate for NPSML estimation with a kernel of order r turns out to be $\delta^{**} = 1/(r+d+2)$. This involves undersmoothing, except when standard second-order kernels are employed ($r = 2$). Then the rate that minimizes the AMISE of the kernel estimator is also optimal for the MSE of $\hat{\theta}_{n,S}$: $\delta^* = \delta^{**} = 1/(4+d)$.

obtain

$$\begin{aligned} B_{S,1} &\simeq -h^r \frac{\kappa_r}{r!} H_0^{-1} E [b_{1,i}(\theta_0)], \\ B_{S,2} &\simeq \frac{1}{Sh^d} \int K^2(z) dz \times H_0^{-1} E \left[\frac{\dot{p}_i(\theta_0)}{p_i^2(\theta_0)} \right] - \frac{1}{Sh^{d+1}} \int K(z) K'(z) dz \times H_0^{-1} E \left[\frac{\dot{m}_i(\theta_0)}{p_i(\theta_0)} \right], \end{aligned} \quad (27)$$

with $\kappa_r = \int K(z) z^r dz$, $H_0 = E [\dot{p}_i(\theta_0) \dot{p}_i(\theta_0)' / p_i^2(\theta_0)]$,

$$b_{1,i}(\theta) := \frac{\dot{p}_i(\theta)}{p_i^2(\theta)} \frac{\partial^r p_i(\theta)}{\partial y_i^r} - \frac{1}{p_i(\theta)} \frac{\partial^r \dot{p}_i(\theta)}{\partial y_i^r}, \quad (28)$$

and $\dot{m}_i(\theta) = \dot{m}(x_i, r(x_i, y_i); \theta)$. Here, we use “ \simeq ” to indicate that only leading terms are included. From these expressions, we see that the kernel smoother distorts the NPSMLE by an order of magnitude $O(h^r)$ while the simulations, in conjunction with the smoothing, generate additional biases of order $O(1/(Sh^{d+1}))$ and $O(1/(Sh^d))$. If a symmetric kernel is employed, $\int K(z) K'(z) dz = 0$ and the second term in the expression of $B_{S,2}$ drops out. Standard bandwidth selection rules in general imply $Sh^d \rightarrow \infty$, but this is not enough for the bias to vanish with rate \sqrt{n} ; we need to undersmooth so that $\sqrt{n}/(Sh^{d+1}) \rightarrow 0$.

The variance components satisfy $D_{n,S} \simeq -\frac{\kappa_r}{r!} (h^r/n) \sum_{i=1}^n \{b_{1,i}(\theta_0) - E[b_{1,i}(\theta_0)]\}$ and $\text{Var}(E_{n,S}) \simeq (nSh^{d+2})^{-1} E[\dot{\sigma}_i^2(\theta_0)/p_i(\theta_0)] \int K'(z)^2 dz$, where $\dot{\sigma}_i^2(\theta) = \text{Var}(\dot{y}_{i,s}(x_i, \theta)|x_i)$. Note that when EIA is employed, the rate of the correction to the variance is non-standard compared to standard SML, which has an efficiency loss of order $1/S$.

6.2 Example 2: SNLS

We derive the adjustment terms for the SNLS example introduced in Section 2. Recall that for nonlinear least squares,

$$G_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n g_i(\theta, \gamma) \text{ with } g_i(\theta, \gamma) = -(y_i - \gamma_i(\theta)) \dot{\gamma}_i(\theta),$$

where $\dot{\gamma}_i(\theta) = \partial \gamma_i(\theta) / (\partial \theta)$. Its first and second-order pathwise differentials are easily found to be $\nabla g_i[d\gamma] = \dot{\gamma}_i d\gamma - (y_i - \gamma_i) d\dot{\gamma}$ and $\nabla^2 g_i[d\gamma, d\gamma] = -d\dot{\gamma} d\gamma$. Therefore, denoting $r_{i,s} = w_S(x_i, \varepsilon_{i,s}; \theta) - \hat{\gamma}_S(x_i; \theta)$, $\nabla^2 \hat{G}_n(\theta) = -\frac{1}{nS(S-1)} \sum_{i=1}^n \sum_{s=1}^S r_{i,s} \dot{r}_{i,s}$. Note that instead of adjusting $G_n(\theta, \hat{\gamma}_S)$, we could have corrected the nonlinear sum of squares instead and minimized

$$\frac{1}{2n} \sum_{i=1}^n (y_i - \hat{\gamma}_{i,S}(\theta))^2 - \nabla^2 \hat{Q}_n(\theta), \quad \nabla^2 \hat{Q}_n(\theta) := \frac{1}{2nS(S-1)} \sum_{i=1}^n \sum_{s=1}^S r^2(x_i, \varepsilon_{i,s}; \theta).$$

This is exactly the correction proposed in Laffont et al. (1995); and as $\nabla^3 g_i \equiv 0$ in SNLS, all approximation biases are removed.¹⁵

To adjust the variance of the estimators, we start from $\sum_{i=1}^n \hat{g}_i \hat{g}'_i / n = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\gamma}_{i,S})^2 \dot{\gamma}_{i,S} \dot{\gamma}'_{i,S} / n$ and we correct it to

$$\frac{1}{n} \sum_{i=1}^n \left(\hat{g}_i \hat{g}'_i - \frac{1}{S^2} \sum_{s=1}^S ((y_i - \hat{\gamma}_{i,S})^2 \dot{r}_{i,s} \dot{r}'_{i,s} + \dot{\gamma}_{i,S} \dot{\gamma}'_{i,S} r_{i,s}^2) \right)$$

before pre- and post-multiplying by the inverse of $\hat{H} = \sum_{i=1}^n (\dot{\gamma}_{i,S} \dot{\gamma}'_{i,S} - (y_i - \hat{\gamma}_{i,S}) \ddot{\gamma}_{i,S}) / n$.

7 Simulation Study

To explore the performance of our proposed approaches, we set up a small Monte Carlo study of a mixed logit model: the econometrician observes i.i.d. draws of $z_i = (x_i, y_i)$ for $i = 1, \dots, n$, with x_i a centered normal of variance τ^2 and

$$y_i = \mathbf{1}(b + (a + su_i)x_i + e_i > 0)$$

where e_i is standardized type I extreme value and u_i is a centered normal with unit variance, independent of e_i .

We take the true model to have parameters $a = 1, s = 1, b = 0$. In this specification, the mean probability of $y = 1$ is one-half. For $\tau = 1$ (resp. $\tau = 2$) the generalized R^2 is 0.11 (resp. 0.21); in the corresponding simple logit model, which has $s = 0$, the R^2 would be 0.17 (resp. 0.39.)

The mixed logit, in its multinomial form, has become a workhorse in studies of consumer demand (see e.g. the book by Train (2009)); it also figures prominently on the demand side of models of empirical industrial organization. It is usually estimated by simulation-based methods. In empirical IO, the simulated method of moments is commonly used because of endogeneity concerns; since they are absent here, we focus on SML instead.

This is still a very simple model; thus we can use Gaussian quadrature to compute the integral

$$\Pr(y = 1|x) = \int \frac{\phi(u)}{1 + \exp(-(b + (a + su)x))} du. \quad (29)$$

Since Gaussian quadrature achieves almost correct numerical integration in such a regular, one-dimensional case, we can rely on it to do (almost) exact maximum likelihood estimation. By the same token, it is easy to compute the asymptotic variance of the exact ML estimator $\hat{\theta}_n$, and the leading term $B_{S,2}$ of the bias of the SML estimator. Simple calculations¹⁶ give

¹⁵Laroque and Salanié (1993, p. S131) give a formula for preventive bias-adjustment in second-order PML.

¹⁶We used adaptive Gaussian quadrature.

the numbers in Table 1.

| τ | $\sqrt{n\hat{\sigma}}$ | | | S times bias | | |
|--------|------------------------|------|-----|----------------|-------|------|
| | a | s | b | a | s | b |
| 1 | 7.2 | 17.1 | 2.4 | -9.0 | -23.2 | -0.0 |
| 2 | 6.7 | 10.8 | 2.8 | -8.2 | -13.3 | -0.0 |

Table 1: Rescaled asymptotic standard errors and simulation biases

The columns labeled $\sqrt{n\hat{\sigma}}$ give the square roots of the diagonal terms of the inverse of the Fisher information matrix. As can be seen from the values of $\sqrt{n\hat{\sigma}}$, it takes a large number of observations to estimate this model reliably. To take an example, assume that the econometrician would be happy with a modestly precise 95% confidence interval of half-diameter 0.2 for the mean slope a . With $\tau = 1$ it would take about $(7.2 * 1.96/0.2)^2 \simeq 5,200$ observations; and still about 4,500 for $\tau = 2$, even though the generalized R^2 almost doubles. With such sample sizes, the estimate of the size of the heterogeneity s would still be very noisy: its 95% confidence intervals would have half-diameters 0.48 and 0.32, respectively for $\tau = 1$ and $\tau = 2$. We also found that the correlation between the estimators of a and of s is always large and positive—of the order of 0.8. Thus the confidence region for the pair (a, s) is in fact a rather elongated ellipsoid. On the other hand, the estimates of b are reasonably precise, which is not very surprising as b shifts the mean probability of $y = 1$ strongly.

The figures in the columns labeled “ S times bias” refer to the expansions of $\hat{\theta}_{nS} - \hat{\theta}_n$ in our theorems. We will be using SML under the EIA scheme (independent draws across observations). Then we know that the leading term of the bias due to the simulations is $B_{S,2}$ and is of order $1/S$. The figures give our numerical evaluation of $SB_{S,2}$, using our formulæ and Gaussian quadrature again. As appears clearly from Table 1, once again the heterogeneity coefficient s is the harder to estimate, followed by a , while there is hardly any bias on b . With $S = 100$ simulations and $\tau = 1$ for instance, the bias on a is -0.09 , and the bias on s is -0.23 .

We ran experiments for several sets of parameter values, sample sizes n , explanatory power (through τ), and numbers of draws S . Since the results are similar, we only present here those we obtained for a sample of 10,000 observations when the true model has $a = 1, s = 1, b = 0$, and the covariate has standard error $\tau = 1$ or $\tau = 2$.

We present below the results for $S = 50, 100, 200$, and 500 simulations. We ran 5,000 simulations in each case, starting from initial values of the parameters drawn randomly from uniform distributions: $a \sim U[0.5, 1.5]$, $b \sim U[-0.5, 0.5]$, and $s \sim U[0.5, 1.5]$. For each simulated sample with $S \leq 200$, we estimated the model using (i) uncorrected SML, (ii) SML with Newton-Raphson (NR), and (iii) SML with analytic adjustment (AA) for both bias and variance. The AA was done on the objective function. For the NR correction, we use only

$k = 1$ step, with $S^* = 10 \times S$ draws¹⁷.

For each method, we also used several ways of computing the standard errors of the estimates: from the most commonly used, which consists of inverting the outer product of the scores without correcting for the simulations, to the better-grounded sandwich formula which we introduced in Section 6.

We faced very few numerical difficulties. The optimization algorithm sometimes stopped very close to the bounds we had imposed for the heterogeneity parameter, $0.1 \leq s \leq 5$. In some cases it failed to find an optimum, especially for uncorrected SML with 50 draws. Finally, the second derivative of the simulated log-likelihood was sometimes not invertible in one of our sandwich formulæ. Altogether, we had to discard 0.2% to 3% of the 5,000 samples, depending on the run. When a sample fails, it is most often because the uncorrected SML does not converge, or it is hard to evaluate the corresponding standard errors. The corrected SML method appears to be much more robust. The tables and graphs below only refer to the remaining samples.

We focus on a and s since there is little to correct for in the SML estimates of b . We report (Huber) robust means, standard errors and RMSEs. “AA” refers to our analytical bias adjustment.

Tables 2 and 3 report our results for the mean error of our various SML methods. Each row corresponds to a value of the number of simulations S . All numbers in the last three columns of these tables were computed by averaging the “error terms” ($\hat{\theta}_{n,S} - \theta_0$) over the 5,000 samples (minus the small number that were eliminated due to numerical issues.) The standard error of these averages is about 0.001, so that several of the biases from the corrected estimates are statistically insignificant.

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|--------|------------|--------|
| 1 | 50 | -0.133 | -0.089 | 0.004 |
| | 100 | -0.078 | -0.039 | 0.000 |
| | 200 | -0.041 | -0.014 | 0.000 |
| | 500 | -0.017 | - | 0.000 |
| 2 | 50 | -0.133 | -0.051 | 0.010 |
| | 100 | -0.069 | -0.016 | -0.016 |
| | 200 | -0.033 | 0.003 | 0.006 |
| | 500 | -0.010 | - | 0.006 |

Table 2: Mean error on a

The “SML” columns in the tables report the biases of the uncorrected SML estimator. The leading term appears to be a good approximation to the actual size of the bias in these

¹⁷We did not run the NR correction for $S = 500$ as it would have been quite time-consuming, with little benefit.

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|--------|------------|--------|
| 1 | 50 | -0.364 | -0.217 | 0.011 |
| | 100 | -0.206 | -0.093 | 0.000 |
| | 200 | -0.109 | -0.033 | 0.001 |
| | 500 | -0.045 | - | 0.001 |
| 2 | 50 | -0.214 | -0.064 | 0.021 |
| | 100 | -0.110 | -0.018 | -0.018 |
| | 200 | -0.051 | 0.010 | 0.013 |
| | 500 | -0.013 | - | 0.013 |

Table 3: Mean error on s

simulations, and the measured bias is close to proportional to $1/S$. This suggests that our analytical bias adjustment, which focuses on correcting for the leading term of the bias, should work very well. As the last columns show, AA in fact does eliminate most of the bias. The Newton step with ten times more simulations reduces the bias, as expected; but it does not do it as effectively as our analytical bias adjustment. In fact, comparing the SML estimator with $S = 500$ to the Newtonized estimator ($S = 50, S^* = 500$) shows that the Newton method only delivers part of the benefits suggested by the theory. Note that with $S = 500$, there is not that much bias to correct, but what there is AA corrects quite well again.

The discussion above only bears on bias, but one may legitimately be concerned about the possibility that our adjustment procedures introduce more noise into the estimates and perhaps even increase their mean square errors. Tables 4 and 5 show that this concern is unfounded. Correcting the estimates using analytical adjustment or a Newton step reduces the RMSE in all cases. Most often, the reduction in bias dominates and AA works better than Newton. However, for larger number of simulations when $\tau = 2$, bias reduction matters less; and since the Newton method is more effective at reducing dispersion, its RMSE becomes smaller than that of the AA method. This suggests that combining AA and a Newton step could yield an even larger reduction in the RMSE.

Both the bias and the increased variance imparted by the simulations affect the properties of standard tests. Figures 1 and 2 document this for t -tests that a and s , respectively, equal their true values. For such a large sample, we would expect the distributions of the t -statistics to be very close to a standard centered normal; and 95% of the mass should lie between -1.96 and 1.96 . What we observe for the uncorrected SML estimator (“SML”) is quite different: the bias in the estimate skews the distribution to the left, spectacularly so for small number of simulations; and the increased variance flattens the distribution.

Resorting to one Newton-Raphson step (the “SML+Newton” curves) corrects part of the bias and reduces the variance; but except for large number of simulations, the distribution of

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|-------|------------|--------|
| 1 | 50 | 0.139 | 0.095 | 0.041 |
| | 100 | 0.083 | 0.043 | 0.028 |
| | 200 | 0.046 | 0.019 | 0.020 |
| | 500 | 0.021 | – | 0.013 |
| 2 | 50 | 0.136 | 0.053 | 0.032 |
| | 100 | 0.072 | 0.020 | 0.020 |
| | 200 | 0.036 | 0.010 | 0.016 |
| | 500 | 0.014 | – | 0.012 |

Table 4: RMSE on a

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|-------|------------|--------|
| 1 | 50 | 0.379 | 0.234 | 0.107 |
| | 100 | 0.219 | 0.103 | 0.074 |
| | 200 | 0.121 | 0.044 | 0.053 |
| | 500 | 0.056 | – | 0.033 |
| 2 | 50 | 0.219 | 0.068 | 0.056 |
| | 100 | 0.115 | 0.025 | 0.025 |
| | 200 | 0.056 | 0.019 | 0.029 |
| | 500 | 0.021 | – | 0.022 |

Table 5: RMSE on s

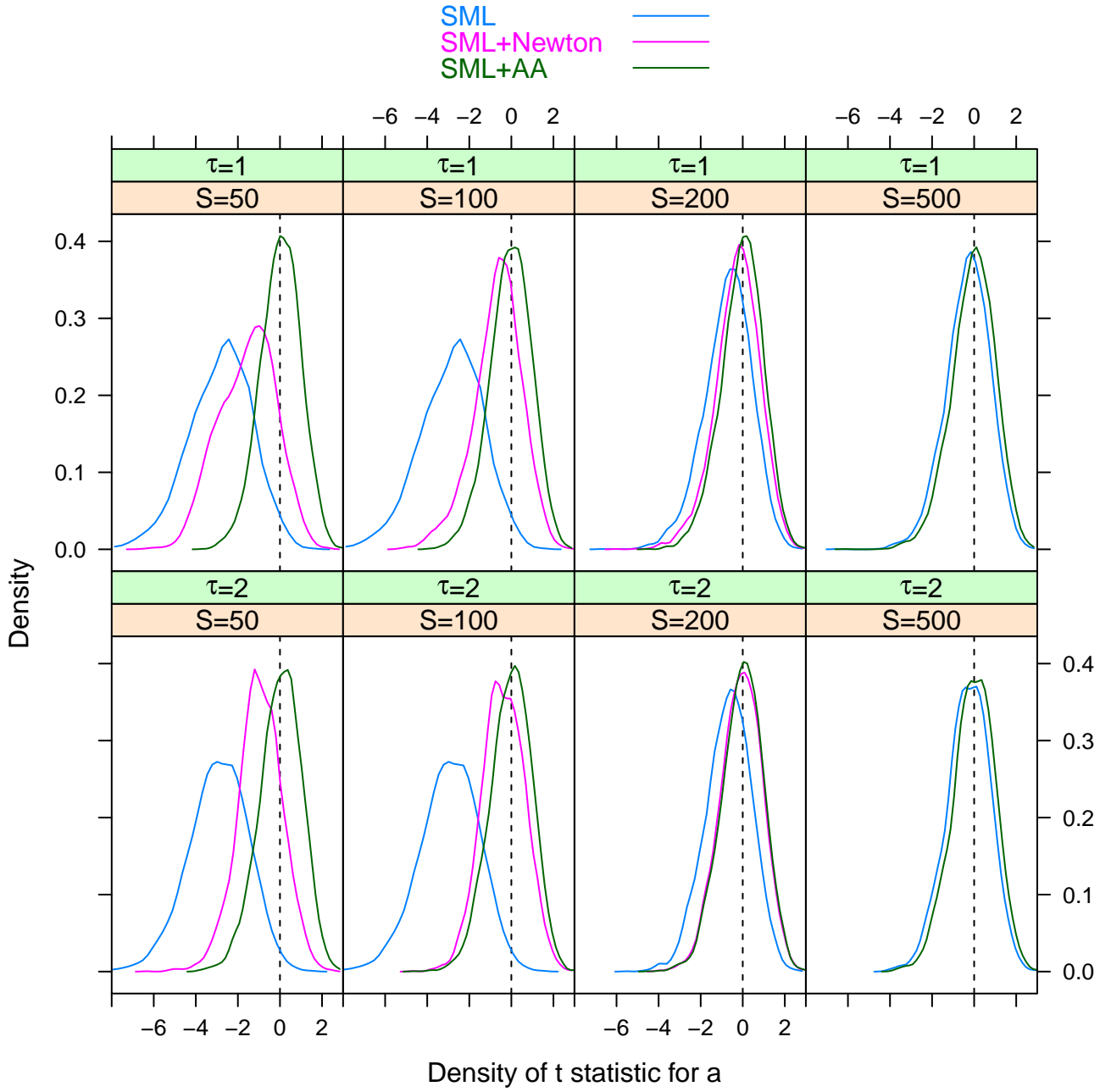


Figure 1: Distributions of the t statistics for $(H_0): a = 1$

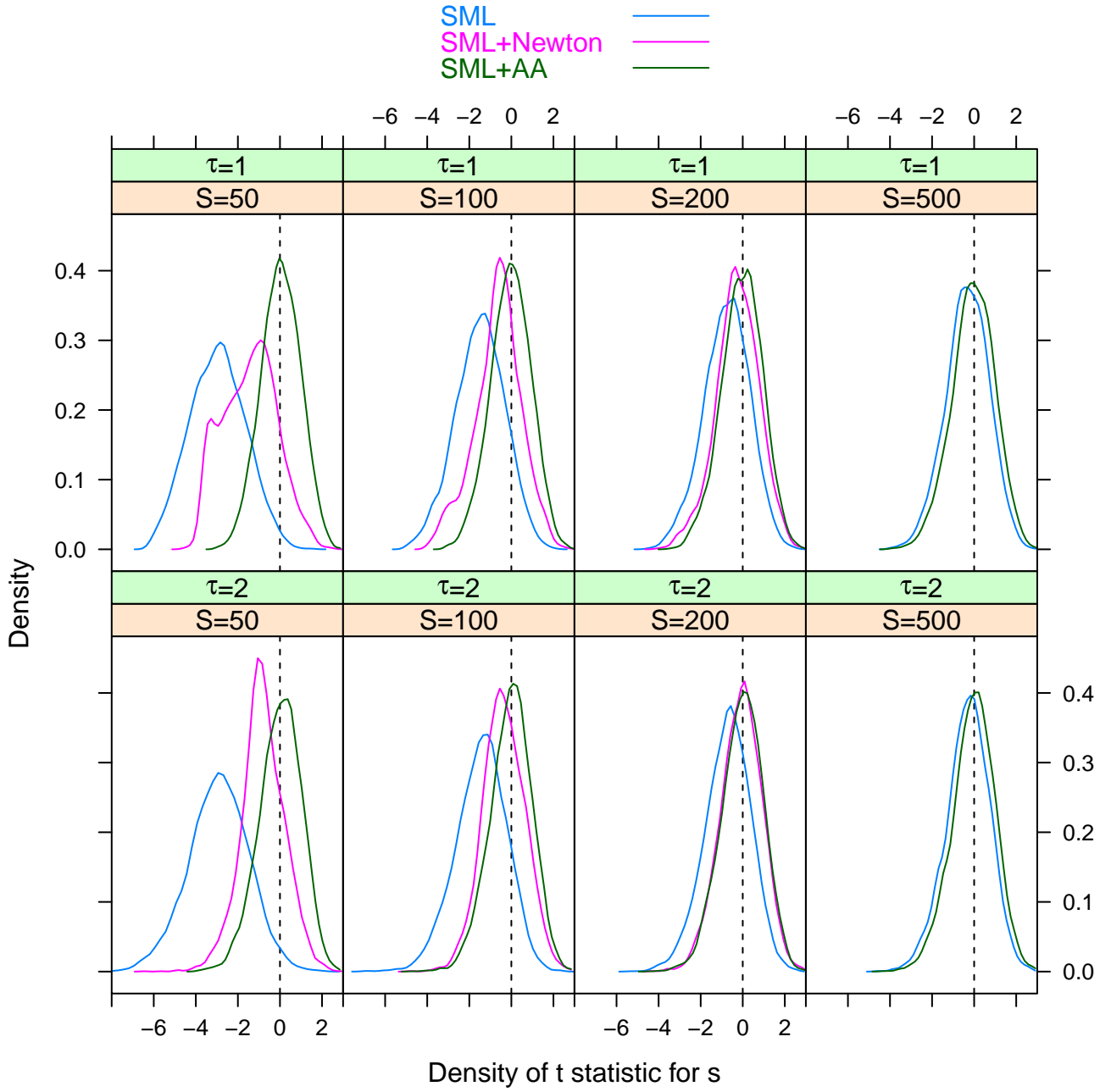


Figure 2: Distributions of the t statistics for $(H_0): s = 1$

the resulting t -statistics is still markedly different from $N(0, 1)$. Using the AA bias-correction and using the proper formula for the variance-covariance matrix (the “SML+AA” curves), on the other hand, produces distributions that are essentially undistinguishable from $N(0, 1)$.

Tables 6 and 7 give the actual coverage probabilities implied by figures 1 and 2. When using uncorrected SML, the nominally 95% confidence intervals undercover very badly, so that the null hypothesis is rejected up to three-quarters of the time when it is in fact true. Our corrections, on the other hand, yield tests that have close to exact coverage.

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|------|------------|--------|
| 1 | 50 | 29.5 | 63.1 | 96.1 |
| | 100 | 68.7 | 86.7 | 95.2 |
| | 200 | 85.9 | 92.6 | 94.9 |
| | 500 | 92.7 | – | 94.5 |
| 2 | 50 | 25.3 | 84.2 | 94.9 |
| | 100 | 69.7 | 93.5 | 95.4 |
| | 200 | 87.2 | 95.4 | 95.6 |
| | 500 | 92.9 | – | 94.6 |

Table 6: Actual coverage probabilities for a

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|------|------------|--------|
| 1 | 50 | 23.2 | 63.7 | 96.7 |
| | 100 | 66.0 | 87.9 | 96.1 |
| | 200 | 86.6 | 93.3 | 95.6 |
| | 500 | 92.9 | – | 94.3 |
| 2 | 50 | 23.6 | 87.9 | 95.1 |
| | 100 | 69.0 | 94.0 | 95.3 |
| | 200 | 87.3 | 95.1 | 95.2 |
| | 500 | 93.0 | – | 94.6 |

Table 7: Actual coverage probabilities for s

Two other considerations are worth mentioning:

- *Ease of implementation:* the analytical bias adjustment wins on that count, since it is usually easy to get a formula for the Δ term and to program it. The Newton method may be more troublesome in models with more than a few parameters, as it requires a reasonably accurate evaluation of the matrix of second derivatives. In our experiment, we relied on the fact that the minimization algorithm itself proceeds by Newton-Raphson steps; after multiplying by ten the number of simulations, we let the algorithm do exactly one iteration of its line search. This appears to work very well, and is very easy to implement.

- *Computer time:* Table 8 reports the mean times per sample. As always, the absolute values are only indicative; we used reasonably optimized C code on a desktop¹⁸. Still, the numbers in the table show that the analytical bias adjustment wins this comparison hands down. For SML for instance, the evaluation of the corrected objective function requires computing the variance of the simulated choice probabilities in addition to their mean, as well as their derivatives—a very small computational cost. Newton adjustment was about five times more costly in our example; it may be more or less time-consuming in other applications, depending on the structure of the model and the care needed to approximate the Hessian. Note however that using 50 simulations to get an estimator and then using a Newton step with 500 simulations is more than twice cheaper than simply working with 500 simulations from the start.

| τ | S | SML | SML+Newton | SML+AA |
|--------|-----|------|------------|--------|
| 1 | 50 | 2.3 | 6.7 | 2.3 |
| | 100 | 3.8 | 11.9 | 3.9 |
| | 200 | 6.6 | 21.6 | 6.8 |
| | 500 | 17.3 | – | 17.7 |
| 2 | 50 | 2.1 | 6.5 | 2.2 |
| | 100 | 3.6 | 11.7 | 3.8 |
| | 200 | 6.3 | 21.3 | 6.6 |
| | 500 | 16.7 | – | 17.2 |

Table 8: Mean CPU time (seconds per sample)

Like any Monte Carlo study, ours can only be illustrative; yet our results are very encouraging. Our analytical corrections for both bias and variance spectacularly improve inference. Using one Newton step, while less effective, can also be a good way to reduce errors.

8 Conclusion

We developed in this paper a unifying framework for the analysis of approximate estimators. We derived a higher-order expansion of the estimators that takes into account additional biases and variances due to approximations. This expansion was then in turn used to develop methods for reducing the bias and the efficiency loss that result from the approximation. Simulations on the mixed logit model confirm that the proposed methods work well in finite samples.

We restricted ourselves to estimators where objective function and approximator (as functions of θ) were both smooth. In principle, one could import the arguments of Chen et al

¹⁸To evaluate derivatives numerically, we used one step of the Ridders-Richardson extrapolation method. We experimented with up to four steps, but the gains in precision were negligible and the results were unchanged.

(2003) to handle non-smooth cases as is done in Armstrong et al (2013). Another approach would be to employ a slight generalization of Robinson (1988, Theorem 1) which in our setting would yield

$$\|\hat{\theta}_{n,S} - \tilde{\theta}_n\| = O_P \left(\sup_{\|\theta - \theta_0\| \leq \delta} \|G_n(\theta, \hat{\gamma}_S) - G_n(\theta, \gamma)\| \right) + o_P(1/\sqrt{n}),$$

for some $\delta > 0$. By strengthening the pointwise bias and variance assumptions to hold uniformly over $\|\theta - \theta_0\| \leq \delta$, we expect our results to remain valid in the non-smooth case.

Also, we require the approximators to be mutually independent, which rules out certain recursive approximation schemes such as particle filtering. Establishing results for this more complicated case would be highly useful. One could here try to use the results of Chen and White (2002) who analyze random dynamic function systems.

We only allowed for one source of approximation in γ . More general situations could have several such terms, possibly with quite different properties. This is for example the case in Kristensen and Scherning (2011), which considers the estimation of dynamic discrete choice models: There, one set of simulations are combined with series regression techniques to approximate the value function (γ_1), and then another set of simulations are used to compute the conditional choice probabilities (γ_2). To cover such situations, Appendix D contains a generalization of Theorem 2 to the case where multiple approximators are employed in the estimation. This is straightforward but tedious, as long as the number of such approximators stays finite; it only requires fairly obvious changes in the assumptions. The expansion can be employed to adjust biases and variances as in the single-approximator case: The analytical bias adjustment will still work when multiple approximators are present, except that we now have to estimate the bias component for each individual approximator. Similarly, the adjustment of standard errors when multiple approximation methods are employed is also relatively straightforward. The Jackknife bias adjustment would on the other hand not be easy to extend to the case where biases vanish at different rates. Finally, the Newton method would remain valid by increasing the degree of approximation for each of the approximators entering the estimation in computing the Newton correction.

Finally, one could interpret an approximate estimator as the exact estimator of a misspecified model. Suppose for instance that we use maximum likelihood to estimate a model with pdf $f(z, \theta)$; and that we suspect that the data may have been generated by a model whose pdf $f^*(z, \theta_0)$ is close to the set of pdfs ($f(\cdot, \theta)$). We can transport all of our results to this problem, with f as γ_S and f^* as γ . In practice we do not know f^* of course; but our methods can be used to explore the likely consequences of any type of (local) misspecification of concern.

9 Acknowledgements

This paper was previously circulated under the title “Higher-order Improvements of Approximate Estimators”. We would like to thank participants at the “Conference on Dynamic Aspects in Economic Decision Making” (University of Copenhagen) and seminar audiences at Northwestern, Ohio State, Toulouse, UCL and Yale for helpful comments and suggestions. Kristensen acknowledges financial support from the Economic and Social Research Council through the ESRC Centre for Microdata Methods and Practice grant RES-589-28-0001, the Danish Social Sciences Research Council (through CREATES) and the European Research Council (Starting grant No. 312474). Parts of this research was conducted while Kristensen visited Princeton University and the University of Copenhagen, whose hospitality is gratefully acknowledged. Salanié gratefully acknowledges the Georges Meyer endowment for its support during a visit at the Toulouse School of Economics.

References

- Ackerberg, D., J. Geweke and J. Hahn (2009) Comments on “Convergence Properties of the Likelihood of Computed Dynamic Models”. *Econometrica* 77, 2009–2017.
- Altissimo, F. and A. Mele (2009) Simulated Nonparametric Estimation of Dynamic Models. *Review of Economic Studies* 76, 413–450.
- Arellano, M. and J. Hahn (2007) Understanding Bias in Nonlinear Panel Models: Some Recent Developments. In *Advances in Economics and Econometrics*, Volume III (eds. R. Blundell, W.K. Newey and T. Persson). Cambridge: Cambridge University Press.
- Andrews, D.W.K. (1994) Asymptotics for Semiparametric Econometric Models Via Stochastic Equicontinuity. *Econometrica* 62, 43–72.
- Andrews, D.W.K. (2002) Higher-order Improvements of a Computationally Attractive k -step Bootstrap for Extremum Estimators. *Econometrica* 70, 119–162.
- Armstrong, T, Gallant, A. R., Hong, H. and L. Huiyu (2012) The Asymptotic Distribution of Estimators with Overlapping Simulation Draws, mimeo Stanford.
- Berry, S., Levinsohn, J., and Pakes, A. (1995) Automobile Prices in Market Equilibrium. *Econometrica* 63, 841–890.
- 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 H. White (2002) Asymptotic Properties of Some Projection-Based Robbins-Monro Procedures in a Hilbert Space. *Studies in Nonlinear Dynamics & Econometrics* 6(1), Article 1.
- Creel, M. and D. Kristensen (2012) Estimation of Dynamic Latent Variable Models Using Simulated Nonparametric Moments. *Econometrics Journal* 15, 490-515.
- Corradi, V. and N.R. Swanson (2007) Evaluation of Dynamic Stochastic General Equilibrium Models Based on Distributional Comparison of Simulated and Historical Data. *Journal of Econometrics* 136, 699-723.
- Dhaene, G. and K. Jochmans (2015) Split-panel Jackknife Estimation of Fixed-effect Models. *Review of Economic Studies*, forthcoming.
- van Dijk, H., A. Monfort and B. Brown (1995) *Econometric Inference Using Simulation Techniques*. John Wiley.
- Dubé, J.P., J. Fox, and C-L. Su (2012) Improving the Numerical Performance of Static and Dynamic Aggregate Discrete Choice Random Coefficients Demand Estimation. *Econometrica* 80, 2231-2267.
- Duffie, D. and K. J. Singleton (1993) Simulated Moments Estimation of Markov Models of Asset Prices. *Econometrica* 61, 929-952.
- Fermanian, J.-D. and B. Salanié (2004) A Nonparametric Simulated Maximum Likelihood Estimation Method. *Econometric Theory* 20, 701-734.
- Fernández-Villaverde, J. and J.F. Rubio-Ramirez (2005) Estimating Dynamic Equilibrium Economies: Linear versus Nonlinear Likelihood. *Journal of Applied Econometrics* 20, 891-910.
- Fernández-Villaverde, J., J.F. Rubio-Ramirez and M. Santos (2006) Convergence Properties of the Likelihood of Computed Dynamic Models. *Econometrica* 74, 93-119.
- Francq, C. and J.-M. Zakoïan (2005) A Central Limit Theorem for Mixing Triangular Arrays of Variables Whose Dependence Is Allowed to Grow with the Sample Size. *Econometric Theory* 21, 1165-1171.
- Freyberger, J. (2015) Asymptotic Theory for Differentiated Products Demand Models with Many Markets. *Journal of Econometrics*, 185, 162-181.
- Gouriéroux, C. and A. Monfort (1996) *Simulation-Based Econometric Methods*. Oxford: Oxford University Press.

- Hahn, J. and W.K. Newey (2004) Jackknife and Analytical Bias Reduction for Nonlinear Panel Models. *Econometrica* 72, 1295–1319.
- Hajivassiliou, V.A. (2000) Some Practical Issues in Maximum Simulated Likelihood. In *Simulation-based Inference in Econometrics* (eds. R. Mariano, T. Schuermann and M.J. Weeks), 71–99. Cambridge: Cambridge University Press.
- Hong, H., A. Mahajan and D. Nekipelov (2015) Extremum Estimation and Numerical Derivatives. *Journal of Econometrics*, forthcoming.
- Judd, K., F. F. Kubler and K. Schmedder (2003) Computational Methods for Dynamic Equilibria with Heterogeneous Agents. In *Advances in Economics and Econometrics* (eds. M. Dewatripont, L.P. Hansen, and S. Turnovsky). Cambridge University Press.
- Judd, K. and C. Su (2012) Constrained Optimization Approaches to Estimation of Structural models. *Econometrica* 80, 2213–2230.
- Keane, M. and K. Wolpin (1994) The solution and estimation of discrete choice dynamic programming models by simulation and interpolation: Monte Carlo evidence. *The Review of Economics and Statistics* 76, 648-672.
- Keane, M. and K. Wolpin (1997) The career decisions of young men. *Journal of Political Economy*, 105, 473-522.
- Kristensen, D. and B. Salanié (2010) Higher Order Improvements for Approximate Estimators. CAM Working Paper 2010-04, University of Copenhagen.
- Kristensen, D. and B. Schjerning (2011) Implementation and Estimation of Discrete Markov Decision Models by Sieve Approximations. Manuscript, University of Copenhagen.
- Kristensen, D. and Y. Shin (2012) Estimation of Dynamic Models with Nonparametric Simulated Maximum Likelihood. *Journal of Econometrics* 167, 76–94.
- Laffont, J.-J., H. Ossard and Q. Vuong (1995) Econometrics of First-Price Auctions. *Econometrica* 63, 953–980.
- Laroque, G. and B. Salanié (1989) Estimation of Multimarket Fix-Price Models: An Application of Pseudo-maximum Likelihood Methods. *Econometrica* 57, 831–860.
- Laroque, G. and B. Salanié (1993) Simulation-based Estimation of Models with Lagged Latent Variables. *Journal of Applied Econometrics* 8, 119–133.
- Laroque, G. and B. Salanié (1994) Estimating the Canonical Disequilibrium Model: Asymptotic Theory and Finite Sample Properties. *Journal of Econometrics* 62, 165–210.

- Lee, L.-F. (1992) On Efficiency of Methods of Simulated Moments and Maximum Simulated Likelihood Estimation of Discrete Response Models. *Econometric Theory* 8, 518–552.
- Lee, L.-F. (1995) Asymptotic Bias in Simulated Maximum Likelihood Estimation of Discrete Choice Models. *Econometric Theory* 11, 437–483.
- Lee, L.-F. (1999) Statistical Inference with Simulated Likelihood Functions. *Econometric Theory* 15, 337–360.
- Lee, L.-F. (2001) Interpolation, Quadrature, and Stochastic Integration. *Econometric Theory* 17, 933–961.
- Mariano, R., T. Schuerman and M. Weeks (2000) *Simulation-based Inference in Econometrics*. Cambridge University Press.
- McFadden, D.F. (1989) A Method of Simulated Moments for Estimation of Discrete Response Models without Numerical Integration. *Econometrica* 57, 995–1026.
- Newey, W.K. (1991) Uniform Convergence in Probability and Stochastic Equicontinuity, *Econometrica* 59, 1161–1167.
- Newey, W.K. (1991) Kernel Estimation of Partial Means and a General Variance Estimator. *Econometric Theory* 10, 233–253.
- Newey, W.K. and D. McFadden (1994) Large Sample Estimation and Hypothesis Testing. In *Handbook of Econometrics*, Vol. 4 (eds. R.F. Engle and D.L. McFadden), Chapter 36. Elsevier Science B.V.
- Newey, W.K. and R. Smith (2004) Higher-order Properties of GMM and Generalized Empirical Likelihood Estimators. *Econometrica* 72, 219–255.
- Norets, A. (2009) Inference in Dynamic Discrete Choice Models with Serially Correlated Unobserved State Variables. *Econometrica* 77, 1665–1682.
- Norets, A. (2012) Estimation of Dynamic Discrete Choice Models Using Artificial Neural Network Approximations. *Econometric Reviews* 31, 84–106.
- Nze, P.A. and P. Doukhan (2004) Weak Dependence: Models and Applications to Econometrics. *Econometric Theory* 20, 995–1045.
- Pakes, A. and D. Pollard (1989) Simulation and the Asymptotics of Optimization Estimators. *Econometrica* 57, 1027–57.
- Rio, E. (1994) Inégalités de moments pour les suites stationnaires et fortement mélangeantes. *Comptes rendus de l'Académie des Sciences* 318, 355–360.

- Robinson, P.M. (1988) The Stochastic Difference Between Econometric Statistics. *Econometrica* 56, 531–548.
- Rothenberg, T.J. (1984) Approximating the Distributions of Econometric Estimators and Test Statistics. In *Handbook of Econometrics*, vol. 2, eds. K. Arrow and M. Intriligator. North Holland.
- Rust, J. (1987) Optimal Replacement of GMC Bus Engines: An Empirical Model of Harold Zurcher *Econometrica* 55, 999–1033.
- Tauchen, G. and R. Hussey (1991) Quadrature-Based Methods for Obtaining Approximate Solutions to Nonlinear Asset Pricing Models. *Econometrica* 59, 371–396.
- Train, K. (2009) *Discrete Choice Methods with Simulation*, Cambridge University Press.
- Yoshihara, K. (1976) Limiting Behaviour of U-Statistics for Stationary, Absolutely Regular Processes. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete* 35, 237–252.

A Proofs of Main Results

Proof of Theorem 2. By Lemma 7,

$$\hat{\theta}_{n,S} - \hat{\theta}_n = -H_0^{-1} \{G_n(\theta_0, \hat{\gamma}_S) - G_n(\theta_0, \gamma_0)\} + o_P(1/\sqrt{n}).$$

Substituting in the expansion given in (7) with $m = 2$ then yields

$$\left\| \hat{\theta}_{n,S} - \hat{\theta}_n \right\| = O_P \left(\left\| \nabla G_n(\theta_0) [\Delta \hat{\gamma}_S] + \frac{1}{2} \nabla^2 G_n(\theta_0) [\Delta \hat{\gamma}_S, \Delta \hat{\gamma}_S] + R_{n,S} \right\| \right) + o_P(1/\sqrt{n}), \quad (30)$$

where $\Delta \hat{\gamma}_{i,S} = \hat{\gamma}_{i,S} - \gamma_0$. We first derive the rate of the remainder term $R_{n,S}$:

$$\begin{aligned} E[\|R_{n,S}\|] &= E \left\| G_n(\theta_0, \hat{\gamma}_S) - G_n(\theta_0, \gamma_0) - \nabla G_n(\theta_0) [\Delta \hat{\gamma}_S] - \frac{1}{2} \nabla^2 G_n(\theta_0) [\Delta \hat{\gamma}_S, \Delta \hat{\gamma}_S] \right\| \\ &\leq \frac{1}{n} \sum_{i=1}^n E \left\| g_i(\theta_0, \hat{\gamma}_{i,S}) - g_i(\theta_0, \gamma_0) - \nabla g_i(\theta_0) [\Delta \hat{\gamma}_{i,S}] - \frac{1}{2} \nabla^2 g_i(\theta_0) [\Delta \hat{\gamma}_{i,S}, \Delta \hat{\gamma}_{i,S}] \right\| \\ &\leq \frac{\bar{G}_0}{n} \sum_{i=1}^n E \left[\|\Delta \hat{\gamma}_{i,S}\|^3 \right], \end{aligned}$$

where we have used A.4(2). Applying first Minkowski's inequality and then the inequality $(a + b)^p \leq 2^{p-1}a^p + 2^{p-1}b^p$ (which holds for all $a, b > 0$ and $p \geq 1$), we obtain—dropping the i index:

$$\begin{aligned} E \left[\|\Delta \hat{\gamma}_S\|^3 \right] &= E \left[\|\psi_S + b_S\|^3 \right] \leq (E [\|\psi_S\|] + E [\|b_S\|])^3 \leq 4E \left[\|\psi_S\|^3 \right] + 4E \left[\|b_S\|^3 \right] \\ &= O(S^{-\alpha_3}) + O(S^{-3\beta}). \end{aligned}$$

The rates of the first and second order functional differentials of $G_n(\theta_0, \gamma)$ are given in Lemmas 10 and 11 depending on whether the ECA approximator of (9) or the EIA approximator of (10) is used. These rates together with the rate of $R_{n,S}$ and (30) yield the higher-order stochastic expansion of the EIA and ECA in equation (15). The rates of the leading bias and variance terms as $S \rightarrow \infty$ also follow from Lemmas 10 and 11.

Finally, the weak convergence of $D_{n,S}$ follows by standard CLT for stationary and mixing triangular arrays. We can, for example, employ Francq and Zakoïan (2005) whose conditions are easily verified given the mixing conditions imposed on data and simulations, which are mutually independent, and the fact that $\text{Var}(\sqrt{n} \{G_n + D_{n,S}\}) = \Omega_S^{G+D} + o(1) = \Omega^G + O(S^{-\beta})$ where $\Omega^G > 0$. ■

Proof of Corollary 3. In the EIA case, $E_{n,S_n} = \sum_{i=1}^n \nabla g_i[\psi_{i,S_n}]$ where (z_i, ψ_{i,S_n}) , for $i = 1, \dots, n$ and $n \geq 1$, is a stationary and mixing triangular array. Under A.6,

$$\begin{aligned} \text{Var}(\sqrt{nS_n} E_{n,S_n}) &= \frac{1}{nS_n} \sum_{i,j=1}^n \sum_{s,t=1}^{S_n} E [\nabla g_i[e_{i,s}] \nabla g_j[e_{j,t}]] = \frac{1}{nS_n} \sum_{i=1}^n \sum_{s,t=1}^{S_n} E [\nabla g_i[e_{i,s}] \nabla g_i[e_{i,t}]] \\ &= \frac{1}{S_n} \sum_{s,t=1}^{S_n} E [\nabla g_0[e_{0,s}] \nabla g_0[e_{0,t}]] = \Omega_{\text{EIA}}^E + o(1), \end{aligned}$$

where we have used that $E [\nabla g_i[e_{i,s}] \nabla g_j[e_{j,t}]] = 0$ for $i \neq j$, and so the claimed result follows from Francq and Zakoïan (2005). For ECA's satisfying A.6, $\sqrt{S} E_{n,S} = \sum_{s=1}^S \nabla G\{\tilde{e}_s\} / \sqrt{S} + o_P(1) \rightarrow^d N(0, \Omega_{\text{ECA}}^E)$, where a CLT for stationary and mixing sequences has been employed. ■

Proof of Theorem 6. We only give a proof for the case of EIA's; the proof for ECA's follows along the same lines. One can easily show that $\sup_{\theta \in \Theta} \|\nabla^2 \hat{G}_n(\theta)\| = o_P(1)$ as $n, S \rightarrow \infty$, and it now follows by the same arguments as in the proof of Theorem 2 that $\hat{\theta}_{n,S}^{\text{BA}}$ is consistent.

Next, we take a Taylor expansion:

$$o_P(n^{-1/2}) = \left\{ G_n(\theta_0, \hat{\gamma}_S) - \frac{1}{2} \nabla^2 \hat{G}_n(\theta_0) \right\} + \left\{ H_n(\bar{\theta}_{n,S}, \hat{\gamma}_S) - \frac{1}{2} \nabla^2 \hat{H}_n(\bar{\theta}_{n,S}) \right\} (\hat{\theta}_{n,S}^{\text{AB}} - \theta_0),$$

where $\nabla^2 \hat{H}_n(\theta) = \partial \nabla^2 \hat{G}_n(\theta) / (\partial \theta)$. From the proof of Theorem 2, $H_n(\bar{\theta}_{n,S}, \hat{\gamma}_S) = H_0 + o_P(1)$, while it is easily shown that $\nabla^2 \hat{H}_n(\bar{\theta}_{n,S}) = o_P(1)$ as $n, S \rightarrow 0$, so that, by the same arguments as in the proof of Theorem 2,

$$\hat{\theta}_{n,S}^{\text{AB}} - \hat{\theta}_n = H_0^{-1} \left\{ G_n(\theta_0, \hat{\gamma}_S) - \frac{1}{2} \nabla^2 \hat{G}_n(\theta_0) - G_n(\theta_0, \gamma) \right\} + o_P(1/\sqrt{n}).$$

Suppressing any dependence on θ_0 , use (7) to write

$$\begin{aligned} G_n(\hat{\gamma}_S) - \frac{1}{2} \nabla^2 \hat{G}_n - G_n(\gamma) &= \frac{1}{2} \left\{ \nabla^2 G_n[\psi_S, \psi_S] - \nabla^2 \hat{G}_n \right\} + \nabla G_n[\hat{\gamma}_S - \gamma] \\ &\quad + \frac{1}{2} \left\{ \nabla^2 G_n[\hat{\gamma}_S - \gamma, \hat{\gamma}_S - \gamma] - \nabla^2 G_n[\psi_S, \psi_S] \right\} + R_{n,S}. \end{aligned} \quad (31)$$

The rates of the second and third terms of (31) are derived in Lemma 11. To ensure that $R_{n,S}$ is negligible, we build on Lemma 12, which uses A.6 to deliver a better rate than that obtained in the proof of Theorem 2.

The crucial term is the first term of (31). Recall $\hat{\gamma}_i(x) = S^{-1} \sum_{s=1}^S w_{is}(x)$, and the definition of $\nabla^2 \hat{G}_n$ in eq. (22). Using the bilinearity of $(d\gamma, d\gamma') \mapsto \nabla^2 g_i[d\gamma, d\gamma']$, and denoting $\bar{w}_i(x) = E[w_{i,s}(x)]$ and $e_{is}(x) = w_{is}(x) - \bar{w}_i(x)$,

$$\begin{aligned} &\nabla^2 G_n[\psi_{n,S}, \psi_{n,S}] - \nabla^2 \hat{G}_n \\ &= \frac{1}{nS^2} \sum_{i=1}^n \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}] + \frac{1}{nS^2} \sum_{i=1}^n \sum_{s=1}^S \nabla^2 g_i[e_{is}, e_{is}] - \frac{1}{nS^2} \sum_{i=1}^n \sum_{s=1}^S \nabla g_i[w_{is} - \hat{\gamma}_i, w_{is} - \hat{\gamma}_i] \\ &= \frac{1}{nS^2} \sum_{i=1}^n \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}] + \frac{1}{nS^2} \sum_{i=1}^n \sum_{s=1}^S \left\{ \nabla^2 g_i[e_{is}, e_{is}] - \nabla g_i[w_{is} - \hat{\gamma}_i, w_{is} - \hat{\gamma}_i] \right\} \\ &= \frac{1}{nS^2} \sum_{i=1}^n \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}] + \frac{1}{nS^2} \sum_{i=1}^n \sum_{s=1}^S \left\{ \nabla^2 g_i[\hat{\gamma}_i - \bar{w}_i, e_{is}] + \nabla^2 g_i[e_{is}, \hat{\gamma}_i - \bar{w}_i] \right\} \\ &= \frac{1}{nS^2} \sum_{i=1}^n \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}] + \frac{2}{nS} \sum_{i=1}^n \nabla^2 g_i[\hat{\gamma}_i - \bar{w}_i, \hat{\gamma}_i - \bar{w}_i], \end{aligned}$$

where the last equality uses the fact that $S^{-1} \sum_{s=1}^S e_{is} = \hat{\gamma}_i - \bar{w}_i$.

Start with the first term, and note that $E[\nabla^2 g_i[e_{is}, e_{it}]] = 0$ when $s \neq t$. Then apply Lemma 8 with $r = 1$ to $W_{i,S} := S^{-2} \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}]$, getting

$$\text{Var} \left(\frac{1}{2nS^2} \sum_{i=1}^n \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}] \right) \leq \frac{C}{n} E \left[\|W_{i,S}\|^{2+\delta} \right]^{2/(2+\delta)}.$$

Now $W_{i,S}$ is a degenerate U -statistic since $E[\nabla^2 g(z_i)[e_{is}, e_{it}] | z_i, e_{it}] = E[\nabla^2 g(z_i)[e_{is}, e_{it}] | z_i, e_{is}] =$

0. Given the conditions imposed on $\{e_{i,s} : 1 \leq s \leq S\}$ in (A.6), we can employ U -statistic results for absolutely regular sequences: Yoshihara (1976, Lemma 3) states that $E \left[\|W_{i,S}\|^4 | z_i \right] = O(S^{-4})$. By inspection of the proof of Yoshihara (1976, Lemma 3), it is easily checked that in fact, for some constant $C > 0$ we have $E \left[\|W_{i,S}\|^4 | z_i \right] \leq CS^{-4} M_S(z_i)$, where $M_S(z_i) := \sup_{s < t} E \left[\|\nabla^2 g(z_i)[e_{is}, e_{it}]\|^{4+\epsilon} | z_i \right]^{4/(4+\epsilon)}$, for some $\epsilon > 0$. Thus, with $\delta = 2$ and using the Lipschitz condition on $\nabla^2 g$, we obtain

$$\begin{aligned}
E \left[\|W_{i,S}\|^4 \right] &\leq CS^{-4} E [M_S(z_i)] \\
&\leq CS^{-4} E \left[\sup_{s < t} E \left[\|\nabla^2 g(z_i)[e_{is}, e_{it}]\|^{4+\epsilon} | z_i \right]^{4/(4+\epsilon)} \right] \\
&\leq CS^{-4} E \left[b^4(z_i) \sup_{s < t} E \left[\|e_{is}(z)\|^{4+\epsilon} \|e_{it}(z)\|^{4+\epsilon} | z_i \right]^{4/(4+\epsilon)} \right] \\
&\leq CS^{-4} E \left[b^4(z_i) E \left[\|e_{is}(z)\|^{8+\epsilon} | z_i \right]^{4/(8+\epsilon)} \right] \\
&\leq CS^{-4} \sqrt{E[b^8(z_i)]} E \left[\|e_{is}\|^{8+2\epsilon} \right]^{4/(8+2\epsilon)} \\
&= O \left(S^{-4+\mu_8/2} \right).
\end{aligned}$$

It follows that $\sum_{i=1}^n \sum_{s \neq t} \nabla^2 g_i[e_{is}, e_{it}] / (nS^2) = O_P(n^{-1/2} S^{-1+\mu_8/4})$. As for the second term, by definition $\hat{\gamma}_i - \bar{w}_i = \psi_{i,S}$; and it follows from Lemma 9 that $E \left[\nabla^2 g_i[\psi_{i,S}, \psi_{i,S}] \right] = O(S^{-\alpha_2})$ and $\frac{1}{n} \sum_{i=1}^n \left(\nabla^2 g_i[\psi_{i,S}, \psi_{i,S}] - E \left[\nabla^2 g_i[\psi_{i,S}, \psi_{i,S}] \right] \right) = O_P(n^{-1/2} S^{-\alpha_4/2})$. Summing up, $\tilde{B}_2 = H_0^{-1} E \left[\nabla^2 G_n[\psi_{n,S}, \psi_{n,S}] - \nabla^2 \hat{G}_n \right] / 2 = O(S^{-2+\mu_2})$ while

$$\text{Var} \left(\nabla^2 G_n[\psi_{n,S}, \psi_{n,S}] - \nabla^2 \hat{G}_n \right) = O(n^{-1} S^{-2+\mu_8/2}) + O(n^{-1} S^{-2+\alpha_4}).$$

This completes the proof. ■

Proof of Theorem 5. To apply the general result in Robinson (1988, Theorem 2), we need to check that his conditions A.1 and A.3 are satisfied in our application. His condition A.1 requires consistency of the approximate estimator for a suitable choice of S , which our assumptions imply. Robinson's condition A.3 also holds, given the smoothness conditions we imposed on $G_n(\theta, \hat{\gamma}_S)$ in our Assumption A.2. ■

B Lemmas

We first derive the expansion in (4):

Lemma 7 *Under Assumptions A.1–A.3, eq. (4) holds.*

Proof. We first take a Taylor expansion of $G_n(\theta, \gamma_0)$ and $G_n(\theta, \hat{\gamma}_S)$ w.r.t. θ :

$$o_P\left(n^{-1/2}\right) = G_n(\hat{\theta}_n, \gamma_0) = G_n(\theta_0, \gamma_0) + H_n(\bar{\theta}_n, \gamma_0)(\hat{\theta}_n - \theta_0), \quad (32)$$

$$o_P\left(n^{-1/2}\right) = G_n(\hat{\theta}_{n,S}, \hat{\gamma}_S) = G_n(\theta_0, \hat{\gamma}_S) + H_n(\tilde{\theta}_{n,S}, \hat{\gamma}_S)(\hat{\theta}_{n,S} - \theta_0), \quad (33)$$

for some $\bar{\theta}_n$ ($\tilde{\theta}_{n,S}$) between $\hat{\theta}_n$ ($\hat{\theta}_{n,S}$) and θ_0 . Since $\hat{\theta}_n$ ($\hat{\theta}_{n,S}$) is consistent, $\bar{\theta}_n$ ($\tilde{\theta}_{n,S}$) $\xrightarrow{P} \theta_0$. By standard arguments together with Assumption A.2,

$$\begin{aligned} \|H_n(\tilde{\theta}_{n,S}, \hat{\gamma}_S) - H_0\| &\leq \|H_n(\tilde{\theta}_{n,S}, \hat{\gamma}_S) - H_n(\tilde{\theta}_{n,S}, \gamma_0)\| + \|H_n(\tilde{\theta}_{n,S}, \gamma_0) - H(\tilde{\theta}_{n,S}, \gamma_0)\| \\ &\quad + \|H(\tilde{\theta}_{n,S}, \gamma_0) - H(\theta_0, \gamma_0)\| \\ &\leq \sup_{\|\theta - \theta_0\| \leq \delta} \|H_n(\theta, \hat{\gamma}_S) - H_n(\theta, \gamma_0)\| + \sup_{\|\theta - \theta_0\| \leq \delta} \|H_n(\theta, \gamma_0) - H(\theta, \gamma_0)\| \\ &\quad + \|H(\tilde{\theta}_{n,S}, \gamma_0) - H(\theta_0, \gamma_0)\| \\ &= o_P(1), \end{aligned}$$

and similar for $H_n(\bar{\theta}_n, \gamma_0)$. Going back to eqs. (32)-(33), we have now shown that

$$\hat{\theta}_{n,S} - \theta_0 = -H_0^{-1}G_n(\theta_0, \hat{\gamma}_S) + o_P(1/\sqrt{n}), \quad \hat{\theta}_n - \theta_0 = -H_0^{-1}G_n(\theta_0, \gamma_0) + o_P(1/\sqrt{n}).$$

Subtracting the second expansion from the first gives the result. ■

To establish the rates for the first and second order differentials, we first establish some useful auxiliary results:

Lemma 8 *Let $\{W_i\}$ be a sequence of α -mixing random variables with $E[W_i] = 0$, $E[\|W_i\|^{2r+\delta}] < \infty$ for some $r \geq 1$ and $\delta > 0$ and its mixing coefficients α_i , $i = 1, 2, \dots$, satisfying $\alpha_i \leq Ai^{-a}$ for some $A > 0$, and $a > 2r + 4r(r-1)/\delta - 2$. Then there exists a constant $C = C(r, a, A) < \infty$ such that:*

$$E\left[\left\|\frac{1}{n}\sum_{i=1}^n W_i\right\|^{2r}\right] \leq n^{-r} \times CE\left[\|W_i\|^{2+\delta}\right]^{2r/(2+\delta)} + o(n^{-r}).$$

Proof. From Rio (1994), we have for $r \geq 1$,

$$E\left[\left\|\frac{1}{n}\sum_{i=1}^n W_i\right\|^{2r}\right] \leq C_r [n^{-r}M_{2,\alpha,n}^r + n^{1-2r}M_{2r,\alpha,n}], \quad (34)$$

where the numbers $M_{p,\alpha,n}$ are defined in Rio (1994). By Nze and Doukhan (2004, p. 1040),

$$M_{p,\alpha,n} \leq \left[E\|W_i\|^{p+\delta}\right]^{p/(p+\delta)} \times \frac{(p+\delta)(p-1)}{\delta} \sum_{n=0}^{\infty} (n+1)^{p+p(p-1)/\delta-2} \alpha_n.$$

Given the bound we imposed on the mixing coefficients, there exists a constant $C(A, a)$ such that

$$\sum_{n=0}^{\infty} (n+1)^{p+p(p-1)/\delta-2} \alpha_n \leq C(A, a) \sum_{n=0}^{\infty} (n+1)^{p+p(p-1)/\delta-2-a} < \infty.$$

In particular, there exist constants $C(r, A, a)$ such that

$$M_{2,\alpha,n}^r \leq C(r, A, a) \left[E \|W_i\|^{2+\delta} \right]^{2r/(2+\delta)}, \quad \text{and} \quad M_{2r,\alpha,n} \leq C(r, A, a) \left[E \|W_i\|^{2r+\delta} \right]^{2r/(2r+\delta)}. \quad (35)$$

The result follows by noting that $n^{1-2r} = o(n^{-r})$ for $r > 1$, and that for $r = 1$ both terms in equation (34) are of order $n^{-1} = n^{-r}$. ■

Lemma 9 *Assume that $\{z_i\}$ satisfies Assumption A.1, and that $\hat{\gamma}_{i,S}$ satisfy Assumption A.5(4) for $i = 1, \dots, J$. Let $m(z; d\gamma)$ be a functional satisfying:*

$$E \left[\|m(z; d\gamma)\|^{2r+\delta} \right] < \infty, \quad E \left[\|m(z; d\gamma)\|^{2+\delta} \right] \leq \bar{M} \|d\gamma\|^{k(2+\delta)}, \quad (36)$$

for some $r, k \geq 1$ and $\delta > 0$. Then, with b_S and ψ_S given in A.5, the following hold:

(i) For EIA's, with $M_S^V := E[m(z_i; \psi_{i,S})]$ and $M_S^B := E[m(z_i; b_{i,S})]$,

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n \{m(z_i; b_{i,S}) - M_S^B\} \right\|^{2r} \right] = O(n^{-r}) \times \left[E \|b_S\|^{k(2+\delta)} \right]^{2r/(2+\delta)},$$

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n \{m(z_i; \psi_{i,S}) - M_S^V\} \right\|^{2r} \right] = O(n^{-r}) \times \left[E \|\psi_S\|^{k(2+\delta)} \right]^{2r/(2+\delta)}.$$

(ii) For ECA's, with $\bar{m}(\gamma) = E[m(z; \gamma)]$ for any fixed γ ,

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n \{m(z_i; b_S) - \bar{m}(b_S)\} \right\|^{2r} \right] = O(n^{-r}) \times \left[E \|\psi_S\|^{k(2+\delta)} \right]^{2r/(2+\delta)},$$

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n \{m(z_i; \psi_S) - \bar{m}(\psi_S)\} \right\|^{2r} \right] = O(n^{-r}) \times \left[E \|\psi_S\|^{k(2+\delta)} \right]^{2r/(2+\delta)}.$$

(iii) The means satisfy $\|M_S^B\| \leq \bar{M} E[\|b_{i,S}\|^k]$, $\|M_S^V\| \leq \bar{M} E[\|\psi_{i,S}\|^k]$, and $E[\|\bar{m}(\psi_S)\|^{2r}] \leq \bar{M} E[\|\psi_S\|^{2kr}]$.

Proof. Define $W_{i,S} := m(z_i; \psi_{i,S}) - M_S(\psi_{i,S})$. By assumptions (A.1) and (A.5), $\{W_{i,S}\}$ is a geometrically mixing process for any given value of S and so its mixing coefficients satisfy the mixing conditions imposed in Lemma 8. Furthermore, (36) implies that $E \left[\|W_{i,S}\|^{2r+\delta} \right] < \infty$.

We can therefore apply Lemma 8

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n \{m(z_i; \psi_{i,S}) - M_S(\psi_{i,S})\} \right\|^{2r} \right] \leq C n^{-r} \left[E \|m(z_i; \psi_{i,S}) - M_S(\psi_{i,S})\|^{2+\delta} \right]^{2r/(2+\delta)} + o(n^{-r})$$

where $C = C(r, a, A)$ only depends on r and the mixing coefficients of $\{z_i\}$ and $\{\psi_{i,S}\}$. By (36), $E \left[\|m(z; \psi_{i,S})\|^{2+\delta} \right] \leq \bar{M} E \left[\|\psi_{i,S}\|^{k(2+\delta)} \right] n^{-r}$ and $\|M_S(\psi_{i,S})\| \leq E \left[\|m(z_i; \psi_{i,S})\| \right] \leq \bar{M} E \left[\|\psi_{i,S}\|^k \right]$. It is easily seen that the above inequalities still go through when replacing $\psi_{i,S}$ with $b_{i,S}$. This prove (i) and (iii).

To derive the second inequality of (ii), now redefine $W_{i,S}$ as $W_{i,S} := m(z_i; \psi_S) - \bar{m}(\psi_S)$. It is easily seen that conditionally on ψ_S , $(W_{i,S})$ satisfies the conditions of Lemma 8, so that

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n W_{i,S} \right\|^{2r} \mid \psi_S \right] \leq C E \left[\|W_{i,S}\|^{2+\delta} \mid \psi_S \right] n^{-r} + o(n^{-r}),$$

where $C = C(r, a, A)$ does not depend on ψ_S . Next, observe that

$$E \left[\|W_{i,S}\|^{2+\delta} \right] \leq C E \left[\|m(z; \psi_S)\|^{2+\delta} \right] \leq C \bar{M} E \left[\|\psi_S\|^{k(2+\delta)} \right];$$

we conclude that

$$E \left[\left\| \frac{1}{n} \sum_{i=1}^n W_{i,S} \right\|^{2r} \right] = E \left[E \left[\left\| \frac{1}{n} \sum_{i=1}^n W_{i,S} \right\|^{2r} \mid \psi_S \right] \right] \leq C E \left[\|\psi_S\|^{k(2+\delta)} \right] n^{-r} + o(n^{-r}).$$

Finally, $E \left[\|\bar{m}(\psi_S)\|^{2r} \right] \leq E \left[\|m(z; \psi_S)\|^{2r} \right] \leq \bar{M} E \left[\|\psi_S\|^{2rk} \right]$. The proof of the first inequality of (ii) follows along the same lines. ■

In the next three lemmas, we suppress the dependence on θ since it is kept fixed at the true value θ_0 .

Lemma 10 *Under A.1-A.3, A.4(2), and A.6(4), the first and second order differentials of $G_n(\theta_0, \hat{\gamma}_S)$ for the ECA yield the rates given in Theorem 2.*

Proof. First consider the EIA case, in which the approximation of $G_n(\gamma)$ is on the form of eq. (10). The functional differentials of G_n are given by

$$\nabla G_n[d\gamma] = \frac{1}{n} \sum_{i=1}^n \nabla g_i[d\gamma], \quad \nabla^2 G_n[d\gamma, d\gamma'] = \frac{1}{n} \sum_{i=1}^n \nabla^2 g_i[d\gamma, d\gamma'],$$

and $d\gamma$ and $d\gamma'$ are the same for all observations $i = 1, \dots, n$. Given A.6(4), the application of the first-order differential to the bias component can be rewritten as

$$\nabla G_n[b_S] = S^{-\beta} \frac{1}{n} \sum_{i=1}^n \nabla g_i [\bar{b}] + \frac{1}{n} \sum_{i=1}^n \nabla g_i [b_S - S^{-\beta} \bar{b}].$$

Now, $E [\sum_{i=1}^n \nabla g_i [\bar{b}] / n] = E [\nabla g_i [\bar{b}]]$, and

$$E \left[\frac{1}{n} \sum_{i=1}^n \left\| \nabla g_i [b_S - S^{-\beta} \bar{b}] \right\| \right] \leq G_1 \|b_S - S^{-\beta} \bar{b}\| = o(S^{-\beta}).$$

By Lemma 9(i) with $m(z; d\gamma) = \nabla g(z) [d\gamma]$, $k = 1$ and $r = 1$, $\text{Var}(\nabla G_n[b_S]) \leq \frac{1}{n} C \|b_S\|^2 = O(S^{-2\beta}/n)$. Since $d\gamma \mapsto \nabla g_i [d\gamma]$ is linear, the conditional mean of the stochastic component of the first-order term is $E[\nabla G_n[\psi_S] | \mathcal{Z}_n] = \frac{1}{n} \sum_{i=1}^n \nabla g_i [E[\psi_S | z_i]] = 0$. Moreover, define $\nabla G[\gamma] = E[\nabla g_i[\gamma]]$ (where expectations are taken w.r.t. the observation z_i); then $\nabla G_n[\psi_S] = \nabla G[\psi_S] + \frac{1}{n} \sum_{i=1}^n \{\nabla g_i[\psi_S] - \nabla G[\psi_S]\}$. Recalling the definition of $\nabla G[\psi_S]$, it follows from Lemma 9(ii) with $m(z; d\gamma) = \nabla g(z) [d\gamma]$ and $k = 2$ that the first term satisfies $\text{Var}(\nabla G[\psi_S]) \leq ME[\|\psi_S\|^2] = O(S^{-\alpha_2})$ while the second term is $O_P(n^{-1/2} S^{-\alpha_2})$.

Regarding the second order differential, its application to the bias component satisfies

$$\nabla^2 G_n[b_S, b_S] = S^{-2\beta} \frac{1}{n} \sum_{i=1}^n \nabla^2 g_i [\bar{b}, \bar{b}] + o_P(S^{-2\beta});$$

moreover, $E[\sum_{i=1}^n \nabla^2 g_i [\bar{b}, \bar{b}] / n] = E[\nabla^2 g_i [\bar{b}, \bar{b}]]$ and, applying Lemma 9(ii) with $m(z; d\gamma) = \nabla^2 g(z) [d\gamma, d\gamma]$, $k = 2$ and $r = 1$, $\text{Var}(\nabla^2 G_n[b_S, b_S]) \leq \frac{1}{n} C \|b_S\|^4 = O(n^{-1} S^{-4\beta})$. To bound the variance component, define $\nabla^2 G[\gamma, \gamma] = E[\nabla^2 g_i[\gamma, \gamma]]$, and write

$$\nabla^2 G_n[\psi_S, \psi_S] = \nabla^2 G[\psi_S, \psi_S] + \frac{1}{n} \sum_{i=1}^n (\nabla^2 g_i[\psi_S, \psi_S] - \nabla^2 G[\psi_S, \psi_S]).$$

Applying Lemma 9(ii) with $m(z; d\gamma) = \nabla^2 g(z) [d\gamma, d\gamma]$ and $r = 1, k = 2$, we obtain that $E\|\nabla^2 G_n[\psi_S, \psi_S]\| = O_P(S^{-2\alpha_2})$.

Finally, by the same arguments as before, $E[\nabla^2 G_n[\psi_S, b_S]] = 0$ while $\text{Var}(\nabla^2 G_n[\psi_S, b_S]) = O(n^{-1} S^{-\alpha_4})$ and $\text{Var}(\nabla^2 G_n[b_S, b_S]) = O(n^{-1} S^{-\alpha_2 - 2\beta})$. ■

Lemma 11 *Under A.1-A.3, A.4(2) and A.5(4), the first and second order differentials of $G_n(\theta_0, \gamma_S)$ for the EIA in (9) yield the rates given in Theorem 2.*

Proof. For the EIA, the first and second order differentials are $\nabla G_n[d\gamma] = \sum_{i=1}^n \nabla g_i [d\gamma_i] / n$ and $\nabla^2 G_n[d\gamma, d\gamma'] = \sum_{i=1}^n \nabla^2 g_i [d\gamma_i, d\gamma'_i] / n$, for any $d\gamma = (d\gamma_1, \dots, d\gamma_n)$ and $d\gamma' = (d\gamma'_1, \dots, d\gamma'_n)$.

It is easily seen that the bias components are the same as those we derived for the ECA in Lemma 10, and so we only consider the variance components. With $\mathcal{Z}_n = (z_1, \dots, z_n)$, the mean of the first-order variance component is zero, $E[\nabla G_n[\psi_S]|\mathcal{Z}_n] = \sum_{i=1}^n \nabla g_i [E[\psi_{i,S}|z_i]]/n = 0$, while its variance satisfies, using Lemma 9.(i) with $m(z, \gamma) = \nabla g(z)[\gamma]$ (in particular, $M_S^V = 0$), $\text{Var}(\nabla G_n[\psi_S]) \leq \frac{1}{n} CE[\|\psi_S\|^2] = O(n^{-1}S^{-\alpha_2})$. Applying Lemma 9(i) and (iii) with $m(z; d\gamma) = \nabla^2 g(z)[d\gamma, d\gamma]$ and $k = 2$, the mean and the variance of the second order differential satisfy

$$E[\nabla^2 G_n[\psi_S, \psi_S]] = E[\nabla^2 g_i[\psi_{i,S}, \psi_{i,S}]] \leq CE[\|\psi_{i,S}\|^2] = O(S^{-\alpha_2}),$$

and $\text{Var}[\nabla^2 G_n[\psi_S, \psi_S]] = O(n^{-1}S^{-\alpha_4})$. The cross term satisfies $E[\nabla^2 G_n[\psi_S, b_S]] = 0$ while $\text{Var}(\nabla^2 G_n[\psi_S, b_S]) = O(n^{-1}S^{-\alpha_2}S^{-2\beta})$, and so we can ignore this term since it is of lower order. ■

Lemma 12 *Assume that A.1-A.3, A.4(3) and A.6(6) hold. Then the rate of the remainder term $R_{n,S}$ can be sharpened to:*

$$R_{n,S} = O_P(S^{-3\beta}) + O_P(S^{-(2-\mu_4)}) + O(S^{-(2-\mu_3)}) + O(n^{-1/2}S^{-(3-\mu_6)/2}).$$

Proof. Since the third-order differential exists, the remainder term in (7) can be further expanded to obtain $R_{n,S} = \nabla^3 G_n[\Delta\hat{\gamma}_S, \Delta\hat{\gamma}_S, \Delta\hat{\gamma}_S]/6 + \bar{R}_{n,S}$ where, by A.4(3) and the same arguments used in the proof of Theorem 2, $E[\|\bar{R}_{n,S}\|] \leq \bar{G}_0 E[\|\Delta\hat{\gamma}_{i,S}\|^4] = O(S^{-4\beta}) + O(S^{-(2-\mu_4)})$. Regarding the third order term, it is easy to check that the bias component is of order $O_P(S^{-3\beta}) + O_P(n^{-1/2}S^{-3\beta})$, by arguments similar to those used in Lemma 10.

This leaves the variance component. In the case of EIA, the variance component can be written as $\nabla^3 G_n[\psi_S, \psi_S, \psi_S] = \sum_{i=1}^n \nabla^3 g_i[\psi_S, \psi_S, \psi_S]/n$. By Lemma 9, we obtain:

$$\nabla^3 G_n[\psi_S, \psi_S, \psi_S] - E[\nabla^3 G_n[\psi_S, \psi_S, \psi_S]] = O(n^{-1/2}S^{-(3-\mu_6)/2});$$

given the independence between simulations,

$$\begin{aligned} |E[\nabla^3 G_n[\psi_S, \psi_S, \psi_S]]| &\leq \frac{1}{S^3} \sum_{s,t,u=1}^S |E[\nabla^3 g_i[e_{i,s}, e_{i,t}, e_{i,u}]]| = \frac{|E[\nabla^3 g_i[e_{i,s}, e_{i,s}, e_{i,s}]]|}{S^2} \\ &\leq \frac{C}{S^2} E[e_{i,s}^3] = O(S^{-(2-\mu_3)}). \end{aligned}$$

In the case of ECA, define $\nabla^3 \bar{g}[\gamma, \gamma, \gamma] = E[\nabla^3 g_i[\gamma, \gamma, \gamma]]$ and write

$$\nabla^3 G_n[\psi_S, \psi_S, \psi_S] = \nabla^3 \bar{g}[\psi_S, \psi_S, \psi_S] + \frac{1}{n} \sum_{i=1}^n \{\nabla^3 g_i[\psi_S, \psi_S, \psi_S] - \nabla^3 \bar{g}[\psi_S, \psi_S, \psi_S]\}.$$

Applying Lemma 9.(ii) with $m(z; d\gamma) = \nabla^3 g(z) [d\gamma, d\gamma, d\gamma]$, the two terms are $O_P(S^{-(3/2-\mu_3)})$ and $O_P(n^{-1/2}S^{-(3-\mu_6)/2})$ respectively. ■

C Expansion of NPSMLE

Let $Q_n(\theta, \gamma_0) = \sum_{i=1}^n \log(p(y_i, x_i; \theta)) / n$ denote the log-likelihood. NPSMLE then maximizes the approximate version where $p (= \gamma_0)$ is replaced by a kernel density estimator based on simulated y_s 's, $\hat{p}_S(y, x; \theta) = \sum_{s=1}^S K_h(y - y_s(x, \theta)) / S$ where $y_s(x, \theta)$, $s = 1, \dots, S$ are i.i.d. draws from $p(\cdot | x; \theta)$. With $\nabla g [dp]$ and $\nabla^2 g [dp, dp]$ given in eq. (26), we here derive explicit expressions for the terms entering the expansion in eq. (15). First note that this expansion is only valid if Eq. (5) holds. It is easily checked that this is the case with $m = 2$ and $\bar{G}_0 := E[\sup_{\theta \in \Theta} \{6 \|\dot{p}_i(\theta)\| / p_i^3(\theta) + 2/p_i^2(\theta)\}]$. To ensure $\bar{G}_0 < \infty$, we either have to assume that the density of covariates is bounded away from zero, or to resort to trimming. Assume for simplicity in the following that the density is bounded away from zero and, moreover, that it is r times differentiable w.r.t. y with its derivatives being integrable, and a r th order kernel is being employed so that $\int z^i K(z) dz = 0$, $i = 1, \dots, r-1$ and $\int z^r K(z) dz < \infty$ for some $r \geq 2$; finally, $\int K'(z)^2 dz < \infty$ and $\int K^2(z) dz < \infty$.

For the analysis of $B_{S,1}$, note that by standard arguments for kernel estimators,

$$\nabla g_i(\theta) [b_S] = \frac{h^r}{r!} \kappa_r \left\{ \frac{\dot{p}_i(\theta)}{p_i^2(\theta)} \frac{\partial^r p_i(\theta)}{\partial y_i^r} - \frac{1}{p_i(\theta)} \frac{\partial^r \dot{p}_i(\theta)}{\partial y_i^r} \right\} + o(h^r),$$

and so, with " \simeq " indicating that only leading terms are included, we obtain from eq. (13) that

$$B_{S,1} = -\frac{\kappa_r}{r!} H_0^{-1} \frac{h^r}{n} \sum_{i=1}^n b_1(y_i, x_i) + o(h^r) = -\frac{\kappa_r}{r!} H_0^{-1} h^r E[b_1(y_i, x_i)] + o(h^r),$$

with $b_1(y_i, x_i)$ defined in eq. (28). This also implies that

$$D_{n,S} = \frac{1}{n} \sum_{i=1}^n d_{i,S} = -\frac{\kappa_r}{r!} \frac{h^r}{n} \sum_{i=1}^n \{b_1(y_i, x_i) - E[b_1(y_i, x_i)]\} + o(h^r).$$

The above analysis is valid irrespectively of whether a single simulation batch (ECA) or n (EIA) simulation batches are used.

Next, we analyze the variance component $E_{n,S}$. First, consider the EIA: By Lemma 9, we obtain that $\text{Var}(\nabla G_n[\psi_S]) = O(1/(nSh^{d+2}))$. More precisely, $E_{n,S} = \nabla G_n[\psi_S] =$

$\frac{1}{n} \sum_{i=1}^n \{a_{S,1,i} + a_{S,2,i}\}$, where $a_{S,1,i}$ and $a_{S,2,i}$, $i = 1, \dots, n$, are i.i.d. sequences given by

$$\begin{aligned} a_{S,1,i} &= \frac{1}{S} \sum_{s=1}^S \frac{\dot{p}_i(\theta)}{p_i^2(\theta)} \{K_h(y_{i,s}(x_i, \theta) - y_i) - E_S [K_h(y_{i,s}(x_i, \theta) - y_i)]\}, \\ a_{S,2,i} &= \frac{1}{S} \sum_{s=1}^S \frac{1}{p_i(\theta)} \{K'_h(y_{i,s}(x_i, \theta) - y_i) \dot{y}_s(x_i, \theta) - E_S [K'_h(y_{i,s}(x_i, \theta) - y_i) \dot{y}_{i,s}(x_i, \theta)]\}. \end{aligned}$$

In particular, with $y_{i,s} = y_s(x_i, \theta_0)$, $p_i = p_i(\theta_0)$ and so forth, we can apply standard results for kernel regression estimators,

$$\begin{aligned} \text{Var}(a_{S,1,i}) &= \frac{1}{S} E \left[\frac{\dot{p}_i \dot{p}'_i}{p_i^2} \{K_h(y_{i,s} - y_i) - E_S [K_h(y_{i,s} - y_i)]\}^2 \right] \\ &= \frac{1}{S} E \left[\frac{\dot{p}_i \dot{p}'_i}{p_i^2} E \left[\{K_h(y_{i,s} - y_i) - E_S [K_h(y_{i,s} - y_i)]\}^2 | y_i \right] \right] \\ &= \frac{1}{Sh^d} E \left[\frac{\dot{p}_i \dot{p}'_i}{p_i} \right] \int K(z)^2 dz, \end{aligned}$$

and, similarly,

$$\begin{aligned} \text{Var}(a_{S,2,i}) &= \frac{1}{S} E \left[\frac{1}{p_i^2} \{K'_h(y_{i,s} - y_i) \dot{y}_{i,s} - E_S [K'_h(y_{i,s} - y_i) \dot{y}_{i,s}]\}^2 \right] \\ &= \frac{1}{S} E \left[\frac{1}{p_i^2} E \left[\{K'_h(y_{i,s} - y_i) \dot{y}_{i,s} - E_S [K'_h(y_{i,s} - y_i) \dot{y}_{i,s}]\}^2 | y_i \right] \right] \\ &= \frac{1}{Sh^{d+2}} E \left[\frac{\dot{\sigma}_i^2}{p_i} \right] \int K'(z)^2 dz. \end{aligned}$$

Thus, $E_{n,S}$ has mean zero and variance $\text{Var}(E_{n,S}) \simeq E[\dot{\sigma}_i^2/p_i] \int K'(z)^2 dz / (nSh^{d+2})$. In the case of ECA, Kristensen and Shin (2012) showed that $\text{Var}(E_{n,S}) = \text{Var}(\nabla G[e_s]) / S + O(1/(nSh^{d+1}))$.

Finally, consider $B_{S,2}$: First note that

$$\nabla^2 g_i[\psi_S, \psi_S] = \frac{2}{p_i^2} \left\{ \frac{\partial \hat{p}_{i,S}}{\partial \theta} - E \left[\frac{\partial \hat{p}_{i,S}}{\partial \theta} \right] \right\} \{\hat{p}_{i,S} - E_S[\hat{p}_{i,S}]\} - \frac{2\dot{p}_i}{p_i^3} \{\hat{p}_{i,S} - E_S[\hat{p}_{i,S}]\}^2.$$

With $m(x, \varepsilon_s) = y_s(x, \theta_0)$ and ε_s being i.i.d. draws from some density $f_\varepsilon(\varepsilon)$ we obtain $p(y|x) = p(y|x; \theta_0) = f_\varepsilon(r(x, y)) |r_y(x, y)|$, where $r(x, y)$ denotes the inverse of $m(x, \varepsilon)$ so

that $\varepsilon = r(x, y)$. Then, for both ECA and EIA,

$$\begin{aligned}
& E_S \left[\left\{ \frac{\partial \hat{p}_{i,S}}{\partial \theta} - E \left[\frac{\partial \hat{p}_{i,S}}{\partial \theta} \right] \right\} \{ \hat{p}_{i,S} - E_S [\hat{p}_{i,S}] \} \right] \\
& \simeq \frac{1}{S} E_S [K'_h(m(x_i, \varepsilon_s) - y_i) K_h(m(x_i, \varepsilon_s) - y_i) \dot{m}(x_i, \varepsilon_s)] \\
& = \frac{1}{S} \int K'_h(m(x_i, \varepsilon) - y_i) K_h(m(x_i, \varepsilon) - y_i) \dot{m}(x_i, \varepsilon) f_\varepsilon(\varepsilon) d\varepsilon \\
& = \frac{1}{S} \int K'_h(y - y_i) K_h(y - y_i) \dot{m}(x_i, r(x_i, y)) p(y|x_i) dy \\
& = \frac{1}{S h^{d+1}} \int K(z) K'(z) dz \times \dot{m}(x_i, r(x_i, y_i)) p(y_i|x_i),
\end{aligned}$$

while, by similar arguments,

$$E_S [\{\hat{p}_{i,S} - E_S [\hat{p}_{i,S}]\}^2] \simeq \frac{1}{S} E [K_h^2(m(x_i, \varepsilon_s) - y_i)] = \frac{1}{S h^d} \int K^2(z) dz \times p_i.$$

Substituting the resulting expression of $E[\nabla^2 g_i[\psi_S, \psi_S]]$ into eq. (13), we obtain the claimed expression in eq. (27).

D Expansion with multiple approximators

We here generalize the theory to handle the case where multiple approximation methods are employed. Let $\hat{\theta}_n$ satisfy a first order condition of the form

$$G_n(\hat{\theta}_n, \gamma_{0,1}, \dots, \gamma_{0,M}) = o_P(1/\sqrt{n}), \quad (37)$$

for some random functional $G_n(\theta, \gamma_1, \dots, \gamma_M)$. The corresponding approximate estimator $\hat{\theta}_{n,S}$ satisfies

$$G_n(\hat{\theta}_{n,S}, \hat{\gamma}_{S_1,1}, \dots, \hat{\gamma}_{S_M,M}) = o_P(1/\sqrt{n}).$$

Here, we allow for γ_m , $m = 1, \dots, M$, being approximated using different methods and with different degrees of approximations, S_m , $m = 1, \dots, M$, which we collect in $S = (S_1, \dots, S_M)$. Collect the approximated functions in $\gamma = (\gamma_1, \dots, \gamma_M)$ and assume that $G_n(\theta, \gamma)$ takes the form of a sample average, $G_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n g(z_i; \theta, \gamma)$. We assume that $\gamma_{0,m}$ belongs to a linear function space Γ_m equipped with a norm $\|\cdot\|_m$, $m = 1, \dots, M$, so that $\gamma \in \Gamma = \Gamma_1 \times \dots \times \Gamma_M$ with norm $\|\gamma\| = \sum_{m=1}^M \|\gamma_m\|$. We maintain the same notation as in the case of one function being approximated and, for example, let $H_n(\theta, \gamma) = \frac{1}{n} \sum_{i=1}^n h(z_i; \theta, \gamma)$, with $h(z_i; \theta, \gamma) = \partial g(z_i; \theta, \gamma) / (\partial \theta)$, denote the first-order derivative of the sample moments. With the same notation, Assumptions A.1-A.3 provided in the main text remain unchanged. Next, we generalize Assumptions A.4-A.5 to:

A.4*(m) Assumption A.4 holds with $\nabla g(z; \theta) [d\gamma] = \sum_{k=1}^M \nabla_m g(z; \theta) [d\gamma_k]$ and $\nabla^2 g(z; \theta) [d\gamma, d\gamma] = \sum_{j,k=1}^M \nabla_{j,k}^2 g(z; \theta) [d\gamma_j, d\gamma_k]$.

A.5*(p) For $m = 1, \dots, M$: The approximator $\hat{\gamma}_{S,m}$ lies in Γ_m and satisfies:

(ii) Its bias $b_{S,m}(z; \theta) = E[\hat{\gamma}_{S,m}(x; \theta)] - \gamma_{0,m}(x; \theta)$ is of order $\beta_m > 0$:

$$b_{S,m}(x; \theta) = S^{-\beta_m} \bar{b}_m(x; \theta) + o(S^{-\beta_m}).$$

(iii) For some $p \geq 2$ and all $2 \leq q \leq p$, there exists $\alpha_{m,q} > 0$ so that $\psi_{S,m}(x; \theta) = \hat{\gamma}_{S,m}(x; \theta) - E[\hat{\gamma}_{S,m}(x; \theta)]$ satisfies:

$$E[\|\psi_{S,m}(x; \theta)\|^q] = S^{-\alpha_{m,q}} v_{m,q}(x; \theta) + o(S^{-\alpha_{m,q}}).$$

Similar to the single approximation scheme case, the leading bias and variance terms take the form

$$B_{S,1} = -H_0^{-1} \sum_{m=1}^M E[\nabla_m g_i[b_{S_m,m}]], \quad B_{S,2} = -\frac{1}{2} H_0^{-1} \sum_{k,m=1}^M E[\nabla_{k,m}^2 g_i[\psi_{S_k,k}, \psi_{S_m,m}]],$$

and

$$\begin{aligned} D_{n,S} &= \frac{1}{n} \sum_{i=1}^n d_{i,S}, \quad d_{i,S} = \sum_{m=1}^M \{\nabla_m g_i[b_{S_m,m}] - E[\nabla_m g_i[b_{S_m,m}]]\}, \\ E_{n,S} &= \frac{1}{n} \sum_{i=1}^n \sum_{m=1}^M \nabla_m g_i[\psi_{S_m,m}]. \end{aligned}$$

We now obtain the following expansion of the approximate estimator, which generalizes the “univariate” version.

Theorem 13 *Assume A.1-A.3, A.4*(2), and A.5*(4). Then,*

$$\hat{\theta}_{n,S} - \theta_0 = B_{S,1} + B_{S,2} + H_0^{-1} \{G_n + D_{n,S} + E_{n,S}\} + O_P \left(\sum_{m=1}^M \{S_m^{-3\beta_m} + S_m^{-\alpha_{m,3}}\} \right) + o_P(1/\sqrt{n}), \quad (38)$$

where $G_n = G_n(\theta_0, \gamma_0)$ and the two sequences $(G_n, D_{n,S})$ and $E_{n,S}$ are asymptotically mutually independent. Moreover, the following limit results hold as $n, S \rightarrow \infty$:

- $\sqrt{n}(\Omega_S^{G+D})^{1/2} \{G_n + D_{n,S}\} \xrightarrow{d} N(0, I_k)$ with $\Omega^G = \sum_{i=-\infty}^{\infty} \text{Cov}(g_0, g_i)$ and

$$\Omega_S^{G+D} = \sum_{i=-\infty}^{\infty} \text{Cov}(g_0 + d_{0,S}, g_i + d_{i,S}) = \Omega^G + O\left(\sum_{m=1}^M S_m^{-2\beta_m}\right).$$

- The bias terms have orders $B_{S,1} = O(\sum_{m=1}^M S_m^{-\beta_m})$ and $B_{S,2} = O(\sum_{k,m=1}^M \sqrt{S_k^{-\alpha_{k,2}} S_m^{-\alpha_{m,2}}})$.
- $\text{Var}(E_{n,S}) = O_P\left(\sum_{m=1}^M S_m^{-\alpha_{m,2}}/n\right)$ (EIA) or $\text{Var}(E_{n,S}) = O_P\left(\sum_{m=1}^M S_m^{-\alpha_{m,2}}\right)$ (ECA).
- If in addition Assumption A.6*(4) holds with $w_{S,m} \equiv w_m$ not depending on S , then $\alpha_{m,2} = 1$, $m = 1, \dots, M$, and

$$\begin{aligned} \mathbf{EIA} & : \left\{ \frac{1}{\sqrt{nS_m}} \sum_{i=1}^n \nabla_m g_i[\psi_{S_m,m,i}] \right\}_{m=1}^M \rightarrow^d N(0, \Omega_{\text{EIA}}^E), \\ \mathbf{EIA} & : \left\{ \frac{1}{\sqrt{S_m}} \sum_{i=1}^n \nabla_m g_i[\psi_{S_m,m}] \right\}_{m=1}^M \rightarrow^d N(0, \Omega_{\text{ECA}}^E) \end{aligned}$$

where $\Omega_{\text{EIA}}^E = \left[\Omega_{\text{EIA},km}^E \right]_{k,m=1}^M$ and $\Omega_{\text{ECA}}^E = \left[\Omega_{\text{ECA},km}^E \right]_{k,m=1}^M$ with

$$\begin{aligned} \Omega_{\text{EIA},km}^E & = \lim_{S \rightarrow \infty} \frac{1}{\sqrt{S_k S_m}} \text{Cov} \left(\sum_{s=1}^{S_k} \nabla_k g_0[e_{k,s}], \sum_{s=1}^{S_m} \nabla_m g_0[e_{m,s}] \right), \\ \Omega_{\text{ECA},km}^E & = \lim_{S \rightarrow \infty} \frac{1}{\sqrt{S_k S_m}} \text{Cov} \left(\sum_{s=1}^{S_k} \nabla_k G[\tilde{e}_{k,s}], \sum_{s=1}^{S_m} \nabla_m G[\tilde{e}_{m,s}] \right). \end{aligned}$$

If the M approximators are mutually independent, the second bias component simplifies to

$$B_{S,2} = -\frac{1}{2} H_0^{-1} \sum_{m=1}^M E \left[\nabla_{m,m}^2 g_i[\psi_{S_m,m}, \psi_{S_m,m}] \right] = O_P \left(\sum_{m=1}^M S_m^{-\alpha_{m,2}} \right),$$

and the off-diagonal elements of the covariance matrices Ω_{EIA}^E and Ω_{ECA}^E become zero.

Proof. By Lemma 7,

$$\hat{\theta}_{n,S} - \hat{\theta}_n = -H_0^{-1} \{G_n(\theta_0, \hat{\gamma}_S) - G_n(\theta_0, \gamma_0)\} + o_P(1/\sqrt{n}).$$

The expansion ($m = 2$) given in eq. (7) then yields

$$\left\| \hat{\theta}_{n,S} - \hat{\theta}_n \right\| = O_P \left(\left\| \nabla G_n(\theta_0) [\Delta \hat{\gamma}_S] + \frac{1}{2} \nabla^2 G_n(\theta_0) [\Delta \hat{\gamma}_S, \Delta \hat{\gamma}_S] + R_{n,S} \right\| \right) + o_P(1/\sqrt{n}), \quad (39)$$

where $\nabla G_n(\theta_0) [d\gamma] = \sum_{m=1}^M \nabla_m G_n(\theta_0) [d\gamma_m]$, $\nabla^2 G_n(\theta_0) [d\gamma, \Delta \gamma] = \sum_{k,m=1}^M \nabla_{k,m}^2 G_n(\theta_0) [d\gamma_k, d\gamma_m]$, and $\Delta \hat{\gamma}_{m,S_m} = \hat{\gamma}_{m,S_m} - \gamma_{m,0}$. The rate of the remainder term $R_{n,S}$ follows by the same arguments as before, $E[\|R_{n,S}\|] \leq \frac{\tilde{G}_0}{n} \sum_{i=1}^n E[\|\Delta \hat{\gamma}_S\|^3]$, where

$$E[\|\Delta \hat{\gamma}_S\|^3] \leq 4 \sum_{m=1}^M \left\{ E[\|\psi_{S_m,m}\|^3] + 4E[\|b_{S_m,m}\|^3] \right\} = O \left(\sum_{m=1}^M \left\{ S_m^{-\alpha_{m,3}} + S_m^{-3\beta_m} \right\} \right).$$

The rate of $\nabla_m G_n(\theta_0) [\Delta \hat{\gamma}_{S_m, m}]$ follow directly from Lemma 10 while $\nabla_{k, m}^2 G_n(\theta_0) [\Delta \hat{\gamma}_{S_k, k}, \Delta \hat{\gamma}_{S_m, m}]$ is analyzed by a simple extension of the arguments employed in the single-approximator case. More specifically, Lemma 9 still applies and yields $\text{Var}(\nabla_{k, m}^2 G_n(\theta_0)[b_{S_k, k}, b_{S_m, m}]) = O(n^{-1} S_k^{-\beta_k} S_m^{-\beta_m})$, for $k, m = 1, \dots, M$. To bound the variance component, apply Lemma 9 to obtain that $E[|\nabla_{k, m}^2 G(\theta_0)[\psi_{S_k, k}, \psi_{S_m, m}]|] = O_P(S_k^{-\alpha_{k, 2}} S_m^{-\alpha_{m, 2}})$, $E[\nabla_{k, m}^2 G_n[\psi_{S_m, m}, b_{S_k, k}]] = 0$ while $\text{Var}(\nabla^2 G_n[\psi_{S_k, k}, b_{S_m, m}]) = O(n^{-1} S_k^{-\alpha_{k, 2}} S_m^{-2\beta_m})$, for $k, m = 1, \dots, M$. The weak convergence results follow by the same arguments as in the proof of Theorem 2. ■

The analytical bias adjustments proposed in Section 5 straightforwardly generalize to the above set-up by simply setting

$$\nabla^2 g(z_i; \theta, \hat{\gamma}_S)[\hat{\psi}_{i, S}, \hat{\psi}_{i, S}] = \sum_{k, m=1}^M \nabla_{k, m}^2 g(z_i; \theta, \hat{\gamma}_S)[\hat{\psi}_{i, S_k, k}, \hat{\psi}_{i, S_m, m}],$$

$$\nabla^2 g(z_i; \theta, \hat{\gamma}_S)[\hat{e}_{i, s, S}, \hat{e}_{i, s, S}] = \sum_{k, m=1}^M \nabla_{k, m}^2 g(z_i; \theta, \hat{\gamma}_S)[\hat{e}_{i, s, S_k, k}, \hat{e}_{i, s, S_m, m}]$$

in eqs. (21) and (22), respectively. If the M approximators are mutually independent, the cross terms in the above double sums can be left out. The adjustments of the standard errors also remain valid when using the definitions of $\nabla g(z; \theta_0)[d\gamma]$ and $\nabla^2 g(z; \theta_0)[d\gamma, d\gamma]$ given in A.4*. Finally, the Newton-Raphson procedure will still work with $S^* = (S_1^*, \dots, S_M^*)$ where $S_m^* > S_m$, $m = 1, \dots, M$.