

An Alternative Way of Computing Efficient Instrumental Variable Estimators*

Xiaohong Chen[†]
(Yale University)

David T. Jacho-Chávez[‡]
(Indiana University)

Oliver Linton[§]
(London School of Economics)

DP No: EM 2009 536

2009

The Suntory Centre
Suntory and Toyota International Centres for
Economics and Related Disciplines
London School of Economics and Political Science
Houghton Street
London WC2A 2AE
Tel: 020 7955 6674

* We would like to thank Javier Hidalgo, Roger Koenker, Guido Kuersteiner, Benno Pötscher, Tom Rothenberg and Pravin Trivedi for helpful comments. We thank STICERD, the NSF and the ESRC for financial support.

[†] Department of Economics, Yale University, PO Box 208281, New Haven CT 06520-8281, USA. E-mail: xiao-hong.chen@yale.edu. Web Page: <http://cowles.econ.yale.edu/faculty/chen.htm>

[‡] Department of Economics, Indiana University, 251 Wylie Hall, 100 South Woodlawn Avenue, Bloomington IN 47403, USA. E-mail: djachoch@indiana.edu. Web Page: <http://mypage.iu.edu/~djachoch/>

[§] Department of Economics, London School of Economics, Houghton Street, London WC2A 2AE, United Kingdom, e-mail: o.linton@lse.ac.uk. This paper was partly written while I was a Universidad Carlos III de Madrid-Banco Santander Chair of Excellence, and I thank them for financial support.

Abstract

A new way of constructing efficient semiparametric instrumental variable estimators is proposed. The method involves the combination of a large number of possibly inefficient estimators rather than combining the instruments into an optimal instrument function. The consistency and asymptotic normality is established for a class of estimators that are linear combinations of a set of \sqrt{n} -consistent estimators whose cardinality increases with sample size. It is shown that the semiparametrically efficient estimator lies in this class. The proofs do not rely on smoothness of underlying criterion functions. Potential use of the estimator can overcome the undersized sample problem. in simultaneous equation system estimation.

JEL Nos.: C12, C13, C14

Keywords: Instrumental Variables; Minimum Distance; Semiparametric Efficiency; Two-Stage Least Squares

1 Introduction

In this paper we derive the properties of an estimator formed by taking linear combinations of an increasing number of inefficient but \sqrt{n} -consistent estimators obtained from conditional moment restrictions. The proposed methodology has the advantage that one can see how much variation there is in the parameter estimates, and how much weight an optimal combination would place on them. In cases where there is truly little variation, the practitioner can presumably do with very simple inference rules. The new estimator is also liable to be useful in situation where asymptotically equivalent alternatives are either impractical or computationally infeasible, i.e. Two Stage Least Squares (2SLS) with an undersized sample problem.

The idea of combining estimates is not new, and has been used to improve finite sample properties of estimators and forecasts. Granger (2000) provides an useful discussion. For example, Sawa (1973) considered combining k-class estimators in simultaneous equations systems, for the reason of improving bias. Breiman (1996,1999) introduced the idea of bagging, which is based on using bootstrap resamples to compute a largeish sample of subsample estimators and then combining them. Watson (2000) and Stock and Watson (1999) propose various methods for combining large numbers of predictors to improve forecasting performance. In the nonparametric literature, Gray and Schucany (1972) and Bierens (1987) have proposed jackknife estimators that combine different kernel smoothers in order to reduce bias. Similarly, Kotlyarova and Zinde-Walsh (2006, 2007) and Schafgans and Zinde-Walsh (2007) have proposed combining kernel smoothers calculated with different bandwidths and kernel functions to construct robust estimator of densities and average derivatives respectively.

Our method is in effect a generalization of the classical method of minimum chi-squared or minimum distance discussed in Rothenberg (1973), which was conceived as a way of imposing equality restrictions in estimation via first estimating an unrestricted model and then finding the best combination of the unrestricted estimators that imposes the restrictions. In a number of cases this strategy is preferable to solving the constrained estimation problem directly. In our case, the best combination is linear with weights that add up to one.

There is a vast literature on estimating models defined through conditional moment restrictions. We just mention one recent paper that is particularly relevant to our study, Koenker and Machado (1999). They considered a similar problem albeit restricted to certain linear models and to a rather specific estimator. They proved that a sufficient condition for the usual asymptotics for generalized method of moments estimation GMM to be valid when the number of moment conditions τ increases with n is that $\tau^3/n \rightarrow 0$. Their results can be interpreted as a warning not to include too many

moment conditions in GMM: that the consequences of so doing are not just that no improvement is made, but that the distributional approximation can potentially break down. Our objective is quite different and we deal with nonlinear models.¹

We first establish consistency and \sqrt{n} -asymptotic normality of a class of estimators that involve finite linear combinations of an infinite dimensional set of estimators, where the cardinality of the linear combinations increases with sample size. The class of estimators considered is allowed to include those computed from discontinuous criterion functions that are nonlinear in the parameters and data. We also establish that a member of our class of estimators achieves the semiparametric efficiency bound for the conditional moment model. We discuss how to estimate the optimal weights and number of estimators to be included. We conclude by presenting results of a Monte Carlo experiment showing how our procedure works in practice.

We use $\|A\| = (\text{tr}(A^\top A))^{1/2}$ for any matrix A . Let $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ denote the smallest and largest eigenvalues of a real symmetric matrix A .

2 The Standard Approach

We observe an independent and identically distributed sample $\{Z_i\}_{i=1}^n$, where $Z_i^\top = (Y_i^\top, X_i^\top)$. We suppose that there is a unique $\theta_0 \in \Theta \subseteq \mathbb{R}^p$ satisfying the conditional moment conditions

$$E[\rho(Z_i, \theta_0) | X_i] = 0$$

with probability one, where $\rho(z, \theta)$ is a scalar residual function.² This implies the unconditional moment conditions

$$E[A(X_i)\rho(Z_i, \theta_0)] = 0, \tag{2.1}$$

for any $p \times 1$ vector $A(X_i)$ [for which the expectation exists]. The sample version of (2.1) is the basis of estimation as described in many previous papers, including Amemiya (1974) and Hansen (1982).

Suppose that $E[\rho(Z_i, \theta_0)^2 | X_i] = \sigma_0^2(X_i)$ is positive with probability one, and that

$$D_0(X_i) = \left(E \left[\frac{\partial \rho}{\partial \theta}(Z_i, \theta) | X_i \right] \right)_{\theta=\theta_0}$$

¹We do not search for the largest value of τ consistent with our asymptotics, although of course the Koenker and Machado *op. cit.* results provide an upper bound.

²The generalization to a system setting is conceptually straightforward; in order to keep the notation simple we shall concentrate on the scalar single equation case.

exists with probability one. In this case, the optimal (instrumental variables) matrix is proportional to $A_{\text{oiv}}(X_i) = D_0(X_i)\sigma_0^{-2}(X_i)$, and the resulting optimal instrumental variables (oiv) –or optimal GMM– estimator $\tilde{\theta}_{\text{oiv}}$ has asymptotic variance $\Sigma_{\text{oiv}} = \{E[\sigma_0^{-2}(X_i)D_0(X_i)D_0(X_i)^\top]\}^{-1}$ - see for example Hansen (1985), Chamberlain (1987) and Newey (1990, 1993).³

Suppose that the optimal matrix $A_{\text{oiv}}(\cdot)$ is of unknown form, but can be represented, in an L_2 sense, by the following series expansion

$$A_{\text{oiv}}(x) = D_0(x)\sigma_0^{-2}(x) = \sum_{j=1}^{\infty} \beta_{j0}\phi_j(x),$$

where $\phi_j(\cdot)$ are known basis functions chosen by the practitioner, while β_{j0} are unknown coefficients determined uniquely by the basis.⁴ For notational convenience we shall allow ϕ_j to be $p \times 1$ vectors; in general, β_{j0} depends on θ_0 and is a $p \times p$ matrix. A common approach here is to estimate the coefficients β_{j0} and then to let

$$\widehat{A}_\theta(x) = \sum_{j=1}^{\tau(n)} \widehat{\beta}_j(\theta)\phi_j(x),$$

where $\tau(n)$ is some truncation sequence that goes to infinity with sample size but at a slow rate.⁵ Then let $\tilde{\theta}_{\text{oiv}}$ be any sequence that satisfies

$$\frac{1}{n} \sum_{i=1}^n \widehat{A}_{\tilde{\theta}_{\text{oiv}}}(X_i)\rho(Z_i, \tilde{\theta}_{\text{oiv}}) = o_p(n^{-1/2}).$$

In current parlance this would be called a continuously updated oiv estimator. An alternative method is to use some preliminary estimator of θ to first construct an estimator of A , and then to solve a similar first order condition with the estimated instrument. Newey (1990, 1993) showed that such an estimator is asymptotically equivalent to the instrumental variable procedure based on knowing the optimal instrument function A_{oiv} and computing solutions $\tilde{\theta}_{\text{oiv}}$ to

$$\frac{1}{n} \sum_{i=1}^n A_{\text{oiv}}(X_i)\rho(Z_i, \tilde{\theta}_{\text{oiv}}) = o_p(n^{-1/2}).$$

³Note that if ρ is not differentiable but the matrix

$$D_0(X_i) = \left(\frac{\partial}{\partial \theta} E[\rho(Z_i, \theta)|X_i] \right)_{\theta=\theta_0}$$

exists, then one might still obtain efficiency by extending the proof in Newey and Powell (1990). We thank Whitney Newey for suggesting this.

⁴In order for the sum to converge, the coefficients β_j must decline as $j \rightarrow \infty$, at least when the basis functions are of fixed magnitude in j .

⁵One can also directly estimate the conditional expectations inside A by nearest neighbor, kernels, or series methods.

See Newey and McFadden (1994) for discussion. There have been a number of alternative suggestions made more recently with a view to improving small sample performance, Newey and Smith (2004) contains an excellent review of this literature.

3 Our Estimation Idea

We take a different approach. Instead of estimating the optimal instrument function we will estimate the optimal way to combine all the available estimators. We consider a sequence of pre-specified basis ($p \times 1$ vector-valued) functions $\{A_j(\cdot)\}$ such that $E[\|A_j(X_i)\|^2] < \infty$; for instance, we may take a uniformly bounded basis such as the B-spline basis. We define the estimators $\hat{\theta}_j$, $j = 1, 2, \dots$, as any sequence that satisfies

$$G_{nj}(\hat{\theta}_j) = \frac{1}{n} \sum_{i=1}^n A_j(X_i) \rho(Z_i, \hat{\theta}_j) = o_p(n^{-1/2}). \quad (3.1)$$

For each j , this problem is completely parametric and will result in a \sqrt{n} -consistent and asymptotically normal estimator $\hat{\theta}_j$ (under standard conditions).⁶ We combine these estimators in a linear fashion to produce a new estimator

$$\hat{\theta} = \sum_{j=1}^{\tau(n)} W_{nj} \hat{\theta}_j, \quad (3.2)$$

where $\tau = \tau(n)$ is a truncation parameter and W_{nj} are some given matrix weights that sum to the identity. This defines a class of estimators \mathcal{E} indexed by the weighting matrices $\{W_{nj}, j = 1, \dots, \tau(n)\}$; as we show below, by an appropriate choice of weights one can achieve the semiparametric efficiency bound for this problem, i.e., the semiparametrically efficient estimator is a member of \mathcal{E} .

Estimator (3.2) is a form of minimum distance where the number of restrictions could increase with sample size.⁷ Even though each criterion function G_{nj} is a nonlinear function of θ , the computational costs of this procedure may not be so great, since one can use the estimates in one step as starting values in the computation of the next step. Additional computational issues arise in connection with the weights W_{nj} but these are discussed below.

Example 1

⁶It is easy to allow the data in (3.1) to depend on j , but we have suppressed this notationally. For example, we could have $Z_i \in \mathbb{R}^\infty$ but with only a finite number of variables in each estimating equation.

⁷See Rothenberg (1973) and Newey and McFadden (1994) for finite fixed τ .

Classical two stage least squares in simultaneous equations. Suppose that

$$y_{1i} = \theta y_{2i} + \varepsilon_i; \quad y_{2i} = \pi_2^\top X_i + u_i,$$

where $(\varepsilon_i, u_i)^\top$ are i.i.d. error terms, $E[\varepsilon_i|X_i] = 0$, $E[u_i|X_i] = 0$ and $X_i \in \mathbb{R}^k$. The two stage least squares estimator is

$$\tilde{\theta} = \frac{\sum_{i=1}^n \hat{y}_{2i} y_{1i}}{\sum_{i=1}^n [\hat{y}_{2i}]^2} = \frac{\sum_{i=1}^n \hat{y}_{2i} y_{1i}}{\sum_{i=1}^n \hat{y}_{2i} y_{2i}}, \quad (3.3)$$

where $\hat{y}_{2i} = \hat{\pi}_2^\top X_i$ and $\hat{\pi}_2$ is the vector of least squares estimates obtained from the reduced form regression of y_{2i} on all the instruments $X_i = (X_{1i}, \dots, X_{ki})^\top$. Our estimator is

$$\hat{\theta} = \sum_{j=1}^k W_{nj} \hat{\theta}_j, \quad (3.4)$$

where

$$\hat{\theta}_j = \frac{\sum_{i=1}^n \hat{y}_{2i}^j y_{1i}}{\sum_{i=1}^n [\hat{y}_{2i}^j]^2} = \frac{\sum_{i=1}^n \hat{y}_{2i}^j y_{1i}}{\sum_{i=1}^n \hat{y}_{2i}^j y_{2i}^j}, \quad (3.5)$$

where $\hat{y}_{2i}^j = \hat{\pi}_{2j}^\top X_{ji}$, and $\hat{\pi}_{2j}$ is the least squares estimates obtained from the reduced form regression of y_{2i} on the single instrument X_{ji} for $j = 1, \dots, k$. Here, W_{nj} are scalar weights that satisfy $\sum_{j=1}^k W_{nj} = 1$. There is a choice of W_{nj} that makes $\hat{\theta}$ asymptotically equivalent to the 2SLS estimator $\tilde{\theta}$. The classical minimum distance estimator (generalized indirect least squares) exploits the relationship between the reduced form coefficients and the structural parameter, i.e., $\pi_{1j}/\pi_{2j} = \theta$, where $\pi_{\ell j} = E[y_{\ell i} X_{ji}] / E[X_{ji}^2]$ are the parameters of the reduced form of $y_{\ell i}$ on X_{ji} for $\ell = 1, 2$ and $j = 1, \dots, k$ [the estimator is a linear combination of $\hat{\pi}_{1j}/\hat{\pi}_{2j}$, where $\hat{\pi}_{\ell j}$ are the corresponding reduced form estimators], see Rothenberg (1973).

Example 2

Now consider the infinite order regression model

$$Y_i = \sum_{k=1}^{\infty} X_{ki} \beta_k(\theta) + \varepsilon_i, \quad (3.6)$$

where θ is some finite dimensional parameter and ε_i is an error term satisfying $E(\varepsilon_i X_{ji}) = 0$, $j = 1, 2, \dots$. Consider the special case that $\beta_k(\theta) = \theta$ for all k . Then, we need at least that $E[(\sum_{k=1}^{\infty} X_{ki})^2] < \infty$ in order for the summation in (3.6) to be well defined; this would be satisfied if $\sigma_k^2 = E(X_{ki})^2$ goes to zero at a rate faster than k^{-1} as $k \rightarrow \infty$. The optimal estimator under homoskedasticity is the OLS estimator of Y_i on $\sum_{k=1}^{\infty} X_{ki}$. If also the regressors are mutually orthogonal, i.e., $E(X_{ji} X_{ki}) = 0$ for all $j \neq k$, the OLS estimators of Y_i on X_{ki} are consistent, and so will any

linear combination thereof, and so we can construct estimators of θ by taking linear combinations of these marginal OLS regressions.⁸

There are two tasks we now pursue. The first is to prove that such an estimator (3.2) is consistent and root-n asymptotically normal under general conditions on the truncation parameter and weighting sequence. The second task is to determine the optimal choice of weights.

4 Large Sample Properties

We begin by defining the sample and population first order conditions. For $j = 1, 2, \dots$, let

$$G_{nj}(\theta) = \frac{1}{n} \sum_{i=1}^n A_j(X_i) \rho(Z_i, \theta) \text{ and } G_j(\theta) = EG_{nj}(\theta)$$

We do not assume that the function $G_{nj}(\theta)$ is differentiable or even continuous, although smoothness conditions are imposed on the expectation $G_j(\theta)$. In this way, we allow also quantile regression estimators (e.g., Koenker and Bassett, 1978), Huber's (1967) M-estimators, and simulation-based estimators (e.g., McFadden (1989) and Pakes and Pollard (1989)). For some of the arguments we only require high level conditions on the sample and population first order conditions, and so our results can apply more generally to any linear combination of estimators that have appropriate expansions.

4.1 Consistency

In this subsection we give our consistency result for the estimator (3.2). We make the following assumptions.

ASSUMPTION A:

(A1) The triangular array $\{W_{nj}\}_{j=1}^{\tau(n)}$, $n = 1, \dots$, satisfies

$$\sum_{j=1}^{\tau(n)} W_{nj} = I_p \text{ and } \sup_n \sum_{j=1}^{\tau(n)} \|W_{nj}\| < \infty, \quad (4.1)$$

with probability tending to one. Here, $\tau(n)$ satisfies $\tau(n) \rightarrow \infty$ as $n \rightarrow \infty$.

⁸By changing variables to X_{ki}/σ_k the parameters become $\theta \cdot \sigma_k$ in which case the problem is more like the instrumental variables regression because the regressors have the same variance but the parameters decline in importance.

(A2) For each j , $\|G_j(\theta_0)\| = 0$.

(A3) For all $\delta > 0$ and $n \geq 1$, there is an $\epsilon_n(\delta) > 0$ (with $\epsilon_n(\delta) \rightarrow 0$) such that

$$\min_{1 \leq j \leq \tau(n)} \inf_{\|\theta - \theta_0\| > \delta} \|G_j(\theta)\| \geq \epsilon_n(\delta) > 0.$$

(A4) For the sequences $\epsilon_n(\delta), \tau(n)$ defined above, there exists a positive sequence α_{1n} with $\sup_n(\alpha_{1n}/\epsilon_n(\delta)) < \infty$ such that

$$\max_{1 \leq j \leq \tau(n)} \left(\|G_{nj}(\hat{\theta}_j)\| - \inf_{\theta \in \Theta} \|G_{nj}(\theta)\| \right) = o_p(\alpha_{1n}),$$

(A5) For the sequences $\epsilon_n(\delta), \tau(n)$ defined above, there exists a positive sequence α_{2n} with $\sup_n(\alpha_{2n}/\epsilon_n(\delta)) < \infty$ such that

$$\max_{1 \leq j \leq \tau(n)} \sup_{\theta \in \Theta} \|G_{nj}(\theta) - G_j(\theta)\| = o_p(\alpha_{2n}).$$

The assumptions on the weights are quite weak and are satisfied by many suitable weighting sequences both random and non-random. For example, equal weighting $W_{nj} = 1/\tau(n)$ satisfies the assumption [A1](#). There are no explicit conditions on the truncation sequence $\tau(n)$ here, but the assumptions [A3–A5](#) may require some restrictions on the rate at which $\tau(n)$ increases with n . Assumption [A4](#) is just a definition of the estimator and is a bit stronger than usual due to the uniformity over j requirement.

The identification Assumption [A3](#) takes account of the fact that each additional moment condition is adding less and less information. The rate at which $\epsilon_n(\delta)$ declines is determined by the sequence $\tau(n)$ and by the sequence A_j , in particular the rate at which $\|E[A_j(X)]\|$ decreases. By choosing $\tau(n)$ to grow very slowly we can compensate for a rapid decline in the moments of the instruments.

The uniform convergence Assumption [A5](#) is easy to verify, although it is slightly stronger than usual due to the $\max_{1 \leq j \leq \tau(n)}$ factor. This factor costs little extra, as can be verified from the Bonferroni and exponential inequalities (see below). Since we must have $\epsilon_n(\delta)$ of larger order than $n^{-1/2}$ in the case of i.i.d. data this puts an upper limit on the rate at which $\tau(n)$ can grow, but no lower limit. If $\tau(n)$ only increases very slowly, say like $\log n$, the stated rate is easy to achieve.

Theorem 1 (i) *Suppose that Assumptions [A1–A5](#) hold. Then $\hat{\theta} - \theta_0 = o_p(1)$.*

For the purpose of obtaining \sqrt{n} -asymptotic normality of $\widehat{\theta}$ in the next subsection, we need to first establish that $\widehat{\theta} - \theta_0 = o_p(n^{-1/4})$ under the following stronger version of Assumption A:

ASSUMPTION A*:

(A*1) **A1** holds.

(A*2) **A2** holds.

(A*3) For all $\delta_n = o(1)$ and $n \geq 1$, there is a positive c_n which could slowly increase to $+\infty$ such that

$$\min_{1 \leq j \leq \tau(n)} \inf_{\|\theta - \theta_0\| > \delta_n} \|G_j(\theta)\| \geq \delta_n c_n > 0.$$

(A*4) For all $\delta_n = o(1)$ and $n \geq 1$,

$$\max_{1 \leq j \leq \tau(n)} \left(\|G_{nj}(\widehat{\theta}_j)\| - \inf_{\|\theta - \theta_0\| \leq \delta_n} \|G_{nj}(\theta)\| \right) = o_p(n^{-1/4}).$$

(A*5) For all $\delta_n = o(1)$ and $n \geq 1$,

$$\max_{1 \leq j \leq \tau(n)} \sup_{\|\theta - \theta_0\| \leq \delta_n} \|G_{nj}(\theta) - G_j(\theta)\| = o_p(n^{-1/4}).$$

Assumption **A*4** is just a definition of the estimator and is a little bit stronger than usual due to the uniformity over j requirement. Assumption **A*5** is stronger than usual, in that we are taking a maximum over an increasing number of first order conditions and requiring a rate at which the resulting random variable goes to zero. However, it is likely to be satisfied in most problems. The uniformity across θ is usually satisfied, indeed we can expect in many cases that $\sup_{\theta \in \Theta} \|G_{nj}(\theta) - G_j(\theta)\| = O_p(1/\sqrt{n})$ for any compact parameter set Θ . Below we provide a Lemma that can be used to verify the uniformity across j condition and may be useful elsewhere.

Theorem 1 (ii) *Suppose that Assumptions **A*1**–**A*5** hold. Then $\widehat{\theta} - \theta_0 = o_p(n^{-1/4})$.*

Of course there are many alternative ways to impose sufficient conditions which lead to convergence rate. We conclude this subsection with a result that is needed in verifying Assumption **A*5** above.

Lemma 1 Let U_{ji} be a triangular array of random variables, $i = 1, \dots, n$, $j = 1, \dots, \tau(n)$, i.i.d. across i for each j with $E(U_{ji}) = 0$ and $E[|U_{ji}|^\kappa] = c_j < \infty$ for some $\kappa \geq 2$. Let $s_{nj}^2 = \sum_{i=1}^n \text{var}(U_{ji}) = n\sigma_j^2$, where $\sigma_j^2 \rightarrow \infty$ as $j \rightarrow \infty$, and let

$$a_n = \left(\max_{1 \leq j \leq \tau(n)} \sigma_j^2 \right) \log \tau(n) + \left(\sum_{j=1}^{\tau(n)} \frac{c_j^2}{\sigma_j^{2\kappa}} \right)^{1/\kappa}. \quad (4.2)$$

Then we have for $\delta_n = a_n \varrho_n$ for any increasing sequence ϱ_n that

$$\max_{1 \leq j \leq \tau(n)} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n U_{ji} \right| = o_p(\delta_n).$$

For example if we take $\kappa = 2$, then $a_n = (\max_{1 \leq j \leq \tau(n)} \sigma_j^2) \log \tau(n) + \sqrt{\tau(n)}$. One application of this Lemma is when $n^{-1/2} \sum_{i=1}^n U_{ji}$ is the leading term of the estimator $\widehat{\theta}_j$, in which case, σ_j^2 would be Γ_j^{-1} (under homoskedasticity). Therefore, the corresponding a_n is of order $\Gamma_{\tau(n)}^{-1} \log \tau(n) + \sqrt{\tau(n)}$. Provided $\tau(n)$ does not increase too rapidly, this is less than $n^{1/4}$ as would be required by assumption **A*5**. Furthermore, it implies that $\max_{1 \leq j \leq \tau(n)} \|\widehat{\theta}_j - \theta_0\|$ goes to zero no slower in probability than $(\Gamma_{\tau(n)}^{-1} \log \tau(n) + \sqrt{\tau(n)})/\sqrt{n}$.

4.2 Asymptotic Normality

In this subsection we derive the asymptotic distribution of our estimator $\widehat{\theta}$, under additional conditions. We strengthen the conditions of Pakes and Pollard (1989) and Newey and McFadden (1994) to accommodate our more general set-up, but again we do not require smoothness conditions on the residual function $\rho(Z_i, \theta)$. Let $g_j(Z_i, \theta) = A_j(X_i)\rho(Z_i, \theta)$ for each j . Then $G_{nj}(\theta) = n^{-1} \sum_{i=1}^n g_j(Z_i, \theta)$ and $\Gamma_j(\theta) = E[g_j(Z_i, \theta)]$. We denote

$$\Gamma_j = \frac{\partial}{\partial \theta^\top} G_j(\theta_0) = \frac{\partial}{\partial \theta^\top} E[A_j(X_i)\rho(Z_i, \theta)] \Big|_{\theta=\theta_0}.$$

If $D_0(X_i) = \{\partial E[\rho(Z_i, \theta)|X_i]/\partial \theta\}|_{\theta=\theta_0}$ exists with probability one, then we have $\Gamma_j = E[A_j(X_i)D_0(X_i)^\top]$.

ASSUMPTION B:

(B1) $\max_{1 \leq j \leq \tau(n)} (\|G_{nj}(\widehat{\theta}_j)\| - \inf_{\|\theta - \theta_0\| \leq \delta_n} \|G_{nj}(\theta)\|) = o_p(1/\sqrt{n})$ for any $\delta_n = o(n^{-1/4})$.

(B2) There exists a finite constant C such that for any θ within a shrinking $(n^{-1/4})$ neighborhood of θ_0

$$\max_{1 \leq j \leq \tau(n)} \|G_j(\theta) - \Gamma_j(\theta - \theta_0)\| \leq C\|\theta - \theta_0\|^2,$$

where Γ_j is of full (column) rank for each j .

(B3) (a) $\max_{1 \leq j \leq \tau(n)} \|\sqrt{n}[G_{nj}(\theta_0) - G_j(\theta_0)]\| = O_p(1)$.

(b) For any $\delta_n = o(n^{-1/4})$,

$$\max_{1 \leq j \leq \tau(n)} \sup_{\|\theta - \theta_0\| \leq \delta_n} \|[G_{nj}(\theta) - G_j(\theta)] - [G_{nj}(\theta_0) - G_j(\theta_0)]\| = o_p(1/\sqrt{n}).$$

(B4) There exists a deterministic sequence of matrices W_{nj}^0 satisfying: (a) $\sum_{j=1}^{\tau(n)} \|(W_{nj} - W_{nj}^0)\Gamma_j^{-1}\| = o_p(1)$; (b) $\limsup_n \sum_{j=1}^{\tau(n)} \|W_{nj}^0 \Gamma_j^{-1}\| < \infty$.

(B5) (a) The matrix $\Sigma_n = \sum_{j=1}^{\tau(n)} \sum_{l=1}^{\tau(n)} W_{nj}^0 V_{jl} W_{nl}^{0\top}$ has a finite positive definite limit Σ , where for all $j, l = 1, \dots, \tau(n)$,

$$V_{jl} = \Gamma_j^{-1} E[g_j(Z_i, \theta_0) g_l(Z_i, \theta_0)^\top] \Gamma_l^{-1\top} = \Gamma_j^{-1} E[A_j(X_i) \sigma_0^2(X_i) A_l(X_i)^\top] \Gamma_l^{-1\top};$$

(b) The triangular array of random variables $f_n(Z_i) = n^{-1/2} \sum_{j=1}^{\tau(n)} c^\top W_{nj}^0 \Gamma_j^{-1} g_j(Z_i, \theta_0)$ satisfies $nE|f_n(Z_i)|^{2+\kappa} \rightarrow 0$ for $\forall c \in \mathbb{R}^p$ and some $\kappa > 0$.

(B6) θ_0 is in the interior of Θ .

(B7) $\max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| = o_p(n^{-1/4})$.

Assumption **B1** is just the definition of the estimator and is a little bit stronger than usual. Assumption **B2** requires essentially two uniformly continuous derivatives for the population moment function at $\theta = \theta_0$ and that the first derivative matrix be of full rank.

For Assumption **B3(b)**, the empirical distribution function satisfies

$$\sup_{|x-x_0| \leq a/n^\alpha} \left| \sqrt{n}[F_n(x) - F(x)] - \sqrt{n}[F_n(x_0) - F(x_0)] \right| = O_p(n^{-\alpha/2})$$

for any $\alpha < 1$ and constant a .⁹ The cost of the additional max is typically no more than an additional factor of order $\sqrt{\tau(n)}$ as is evidenced in the Lemma 1.

In **B4**, we require that if the weights are random that they can be well approximated by some nonrandom sequence with certain summability properties. This condition entails some restrictions on the rate of growth of τ , and these restrictions can be as much as requiring that $\tau^3/n \rightarrow 0$, see Koenker

⁹We are grateful to Benedikt Pötscher for pointing this out to us. This is due to the Hölder continuity of the limiting Brownian bridge process $B(\cdot)$ of $\sqrt{n}[F_n(\cdot) - F(\cdot)]$, i.e., $|B(x) - B(x_0)| \leq c \cdot |x - x_0|^{1/2}$ for some random variable c with bounded moment. The local uniformity [across i] comes at very little extra cost.

and Machado (1999). The restrictions are not so stringent in special cases and really arise out of the nonlinearity of the estimating equation rather combined with the large number of parameters.

Assumption **B5** allows us to apply the Liapounov's central limit theorem for triangular arrays to the leading term. This condition is satisfied for a variety of problems, and it implicitly imposes restrictions on how fast $\tau(n)$ could grow with sample size n . Notice that Assumption **B5**(b) is simply: for some $\kappa > 0$ and for all c ,

$$E \left(\left| \sum_{j=1}^{\tau(n)} c^\top W_{nj}^0 \Gamma_j^{-1} g_j(Z_i, \theta_0) \right|^{2+\kappa} \right) = o(n^{\kappa/2}).$$

For example, suppose we only require that $g_j(Z_i, \theta_0)$ have uniformly bounded fourth moments. Then, by the Cauchy-Schwarz inequality

$$nE[f_n(Z_i)^4] = \frac{1}{n} \sum_{j,k,l,m=1}^{\tau(n)} E[\varphi_{ji}\varphi_{ki}\varphi_{li}\varphi_{mi}] \leq \frac{1}{n\epsilon_n^4} \left(\sup_n \sum_{j=1}^{\tau(n)} \|W_{nj}^0\| \right)^4,$$

where $\varphi_{ji} = c^\top W_{nj}^0 \Gamma_j^{-1} g_j(Z_i, \theta_0)$. It suffices in this case that $n\epsilon_n^4 \rightarrow \infty$. Now suppose that in fact, the scalar $g_j(Z_i, \theta_0)$ are normally distributed with mean zero and variance Γ_j and mutually independent, and that the weights are equal, i.e., $W_{nj}^0 = 1/\tau(n)$ for each j . Then

$$nE[f_n(Z_i)^4] = \frac{1}{n\tau^4} \left(\sum_{j=1}^{\tau(n)} 3\Gamma_j^{-2} + 3 \sum_{j \neq k}^{\tau(n)} \Gamma_j^{-1} \Gamma_k^{-1} \right) \leq \frac{3}{n\tau^2 \epsilon_n^2},$$

which goes to zero provided $n\tau^2 \epsilon_n^2 \rightarrow \infty$. These conditions can be weakened considerably in special cases.

Notice that we can replace Assumptions **B3**(a) and **B5** by the condition that $\{G_{nj}(\theta_0) - G_j(\theta_0) : 1 \leq j \leq \tau(n)\}$ is a Donsker class, i.e., it satisfies the uniform central limit theorem. This kind of assumption has been used in Portnoy (1985) for example.

The condition **B7** that $\max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| = o_p(n^{-1/4})$ follows from our Theorem **1**(ii). It may be possible to prove our result below without a sup-norm convergence result like this, although we have not been able to find a proof based on other convergences like L_p . The usual proofs in other semiparametric estimation problems typically make use of similar results about the convergence of nuisance parameters.

Theorem 2 *Suppose that Assumptions **B1**–**B7** hold. Then $\sqrt{n}(\hat{\theta} - \theta_0) \implies N(0, \Sigma)$.*

The asymptotic variance matrix Σ depends on the weighting scheme and on the class of estimators considered and of course on the underlying distribution of the data. We discuss the nature of the asymptotic variance more in the next section.

To construct consistent estimates of Σ , we would compute

$$\widehat{\Sigma} = \sum_{j=1}^{\tau(n)} \sum_{l=1}^{\tau(n)} \kappa \left(\frac{|j-l|}{\tau} \right) W_{nj} \widehat{V}_{jl} W_{nl}^\top$$

for some weighting function κ , and

$$\widehat{V}_{jl} = \widehat{\Gamma}_j^{-1} \frac{1}{n} \sum_{i=1}^n g_j(Z_i, \widehat{\theta}) g_l(Z_i, \widehat{\theta})^\top \widehat{\Gamma}_l^{-1\top}. \quad (4.3)$$

The estimation of Γ_j is easy when G_{nj} are differentiable. In this case,

$$\widehat{\Gamma}_j = \frac{1}{n} \sum_{i=1}^n \frac{\partial g_j(Z_i, \widehat{\theta})}{\partial \theta} \rightarrow^p \Gamma_j \quad (4.4)$$

under some regularity conditions. The weighting function κ must satisfy some regularity conditions as in Andrews (1991). When G_{nj} are not differentiable, as for example in the LAD case, this method is not feasible. In some cases, one might be able to estimate directly the quantity Γ_j . For example, in the LAD case [with errors independent of covariates], Γ_j is proportional to the density of the errors evaluated at their median. This quantity can be estimated by a variety of nonparametric methods. A general strategy for estimating Γ_j is to use ‘numerical derivatives’, that is, let

$$\widehat{\Gamma}_{j;lk} = \frac{1}{n} \sum_{i=1}^n \frac{g_{jl}(Z_i, \widehat{\theta} + \delta e_k) - g_{jl}(Z_i, \widehat{\theta})}{\delta}, \quad (4.5)$$

where e_k is a vector of zeros with one in the k^{th} position, while δ is a small constant. If we let $\delta(n)$ go to zero at a certain rate as sample size increases, we can show that $\widehat{\Gamma}_{j;lk} \rightarrow^p \Gamma_{j;lk}$, and under stronger conditions $\max_{1 \leq j \leq \tau} \|\widehat{\Gamma}_j - \Gamma_j\| \rightarrow^p 0$. The actual derivative (4.4) makes δ go to zero before n , but our modified estimator (4.5) allows δ to go to zero with n and indeed slower than n .

Example 2 (*cont.*)

Suppose that the errors are homoskedastic and the regressors are mutually orthogonal with $E(X_{ji}^2) = \sigma_j^2$. A necessary and sufficient condition for the \sqrt{n} -rate of convergence is that

$$\limsup_{n \rightarrow \infty} \sum_{j=1}^{\tau} W_{nj}^2 \sigma_j^{-2} < \infty$$

with probability one. Since we also require $\sum_{j=1}^{\infty} \sigma_j^2 < \infty$, this rules out the equal weighting case. Nevertheless, a variety of weighting conditions satisfy the requirement. Furthermore, there is no explicit restriction on τ itself in this case.

5 Choice of τ

In this section, we describe how one would select the truncation parameter τ in a given application with finite sample n . Since each $\hat{\theta}_j$ is a M-estimator obtained by solving (3.1), then under suitable conditions, i.e. G_{nj} is continuously differentiable with respect to θ , Rilstone, Srivastava and Ullah (1996) showed that:

$$\mathfrak{B}_j \equiv \text{Bias}(\hat{\theta}_j) = \frac{1}{n} Q_j(\theta_0) \left[E(v_{j;i}(\theta_0) d_{j;i}(\theta_0)) + \frac{1}{2} H_{j;2}(\theta_0) E(d_{j;i}(\theta_0) \otimes d_{j;i}(\theta_0)) \right]$$

where $Q_j(\theta) = (-E[\partial^2 g_j(Z_i, \theta)/(\partial\theta^\top \otimes \partial\theta^\top)])$, $H_{j;2}(\theta) = -[Q_j(\theta)]^{-1}$, $d_{j;i}(\theta) = Q_j(\theta) g_j(Z_i, \theta)$, and $v_{j;i}(\theta) = \partial g_j(Z_i, \theta)/\partial\theta^\top - E[\partial g_j(Z_i, \theta)/\partial\theta^\top]$. Therefore, we could choose τ by minimizing the trace of a consistent estimator of the approximated mean squared error¹⁰, i.e.

$$\hat{\tau} = \arg \min_{\tau \in \mathbb{Z}_{++}} \text{trace} \left(\sum_{j=1}^{\tau} \sum_{l=1}^{\tau} W_{nj} \left[\hat{\mathfrak{B}}_j \hat{\mathfrak{B}}_l^\top + \hat{V}_{jl} \right] W_{nl}^\top \right), \quad (5.1)$$

where \hat{V}_{jl} was defined in 4.3, and $\hat{\mathfrak{B}}_j$ is constructed as

$$\hat{\mathfrak{B}}_j = \frac{1}{n} \hat{Q}_j(\hat{\theta}_j) \left[\frac{1}{n} \sum_{i=1}^n \hat{v}_{j;i}(\hat{\theta}_j) \hat{d}_{j;i}(\hat{\theta}_j) + \frac{1}{2} \hat{H}_{j;2}(\hat{\theta}_j) \frac{1}{n} \sum_{i=1}^n (\hat{d}_{j;i}(\hat{\theta}_j) \otimes \hat{d}_{j;i}(\hat{\theta}_j)) \right],$$

with $\hat{Q}_j(\theta) = (-n^{-1} \sum_{i=1}^n [\partial^2 g_j(Z_i, \theta)/(\partial\theta^\top \otimes \partial\theta^\top)])$, $\hat{H}_{j;2}(\theta) = -[\hat{Q}_j(\theta)]^{-1}$, $\hat{d}_{j;i}(\theta) = \hat{Q}_j(\theta) g_j(Z_i, \theta)$, and $\hat{v}_{j;i}(\theta) = \partial g_j(Z_i, \theta)/\partial\theta^\top - n^{-1} \sum_{i=1}^n [\partial g_j(Z_i, \theta)/\partial\theta^\top]$.

The minimization problem (5.1) is computationally feasible once $\hat{\mathfrak{B}}_j$ and \hat{V}_{jl} are calculated for $j, l = 1, \dots, \tau$, because it only involves a numerical search over strictly positive integers.

6 Optimal Weights

We now discuss how to choose optimal weights.

¹⁰Further terms, characterizing the higher order efficiency, can be included (see Lemma 3.3 in Rilstone, Srivastava and Ullah (1996), p 377)

6.1 Case 1: Fixed τ

Suppose that we know only that

$$E[A_j(X_i)\rho(Z_i, \theta_0)] = 0, \quad j = 1, \dots, \tau, \quad (6.1)$$

where τ is fixed, and $A_j \in \mathbb{R}^p$. This is a standard unconditional moments estimation problem, and the optimal estimator can be arrived at by several routes:

GMM with optimal combination of the moment conditions:

That is, we minimize the quadratic form

$$G_n^\tau(\theta)^\top W_n G_n^\tau(\theta) \quad (6.2)$$

with respect to θ , where $G_n^\tau(\theta) = n^{-1} \sum_{i=1}^n A^\tau(X_i)\rho(Z_i, \theta)$ with $A^\tau = (A_1^\top, \dots, A_\tau^\top)^\top \in \mathbb{R}^{\tau p}$ (i.e., $G_n^\tau(\theta)$ is the $\tau p \times 1$ vector containing all the sample moments). The asymptotically optimal (opt) weighting matrix is $W_{\text{opt}} = \Psi_\tau^{-1}$, where $\Psi_\tau = E[G_n^\tau(\theta_0)G_n^\tau(\theta_0)^\top] = E[A^\tau(X)\sigma_0^2(X)A^\tau(X)^\top] \in \mathbb{R}^{\tau p \times \tau p}$.

Optimal instrumental variables:

The optimal instrument in this case is simply a linear combination of the $A_j(X_i)$, $j = 1, \dots, \tau$. That is, we solve the equations

$$\Gamma^\tau \Psi_\tau^{-1} G_n^\tau(\hat{\theta}) = 0, \quad (6.3)$$

where $\Gamma^\tau = \partial E[G_n^\tau(\theta_0)]/\partial \theta = E[A^\tau(X)D_0(X)^\top] \in \mathbb{R}^{\tau p \times p}$. These two approaches provide the oiv (optimal GMM) estimator $\tilde{\theta}_{\text{oiv}}^\tau$ of θ_0 for the model (6.1). Specifically, we have $\sqrt{n}(\tilde{\theta}_{\text{oiv}}^\tau - \theta_0) \implies N(0, \Sigma_{\text{oiv}}^\tau)$ as $n \rightarrow \infty$, where the asymptotic variance is given by (see e.g., Hansen (1982) for differentiable ρ , Newey and McFadden (1994) for non-differentiable ρ):

$$\begin{aligned} \Sigma_{\text{oiv}}^\tau &= \left(E[A^\tau(X)D_0(X)^\top]^\top [E(A^\tau(X)\sigma_0^2(X)A^\tau(X)^\top)]^{-1} E[A^\tau(X)D_0(X)^\top] \right)^{-1} \\ &= (\Gamma^\tau \Psi_\tau^{-1} \Gamma^\tau)^{-1} \end{aligned} \quad (6.4)$$

and the optimal instrument for the model (6.1) is:

$$\begin{aligned} A_{\text{oiv}}^\tau(x) &= \Gamma^\tau \Psi_\tau^{-1} A^\tau(x) \\ &= E[A^\tau(X)D_0(X)^\top]^\top [E(A^\tau(X)\sigma_0^2(X)A^\tau(X)^\top)]^{-1} A^\tau(x). \end{aligned}$$

Minimum distance

This approach refers to the minimum distance method described in Rothenberg (1973). In particular, let $\widehat{\theta}_{\text{omd}}^\tau$ minimize the criterion function

$$Q_n(\theta) = \left[\begin{pmatrix} \widehat{\theta}_1 \\ \vdots \\ \widehat{\theta}_\tau \end{pmatrix} - \theta \otimes i_\tau \right]^\top V^{-1} \left[\begin{pmatrix} \widehat{\theta}_1 \\ \vdots \\ \widehat{\theta}_\tau \end{pmatrix} - \theta \otimes i_\tau \right], \quad (6.5)$$

where i_τ is a $\tau \times 1$ vector of ones, and V is the $\tau p \times \tau p$ asymptotic (as $n \rightarrow \infty$ holding τ constant) variance matrix of the vector $(\sqrt{n}(\widehat{\theta}_1 - \theta_0)^\top, \dots, \sqrt{n}(\widehat{\theta}_\tau - \theta_0)^\top)^\top$, i.e., $V = (V_{jl})$, where $V_{jl} = \Gamma_j^{-1} E[A_j(X_i) \sigma_0^2(X_i) A_l(X_i)^\top] \Gamma_l^{-1\top}$ for all $j, l = 1, \dots, \tau$. The first order condition

$$(I_p \otimes i_\tau)^\top V^{-1} \begin{pmatrix} \widehat{\theta}_1 \\ \vdots \\ \widehat{\theta}_\tau \end{pmatrix} = (I_p \otimes i_\tau)^\top V^{-1} (\widehat{\theta} \otimes i_\tau)$$

implies that the optimal estimator $\widehat{\theta}_{\text{omd}}^\tau$ is a linear combination of the $\widehat{\theta}_j$ with

$$\widehat{\theta}_{\text{omd}}^\tau = \sum_{j=1}^{\tau} W_{0j}^{\text{opt}} \widehat{\theta}_j, \quad (6.6)$$

where

$$W_{0j}^{\text{opt}} = \left(\sum_{l=1}^{\tau} B_l \right)^{-1} B_j,$$

$$\text{and } (B_1, \dots, B_\tau) = (I_p \otimes i_\tau)^\top V^{-1}.$$

Furthermore, $\sqrt{n}(\widehat{\theta}_{\text{omd}}^\tau - \theta_0) \Rightarrow N(0, \Sigma_{\text{omd}}^\tau)$, where the asymptotic [as $n \rightarrow \infty$ and τ fixed] variance is

$$\Sigma_{\text{omd}}^\tau = W_{\text{opt}} V W_{\text{opt}}^\top = ((I_p \otimes i_\tau)^\top V^{-1} (I_p \otimes i_\tau))^{-1}.$$

Proposition 1 For each fixed τ , $\widehat{\theta}_{\text{omd}}^\tau$ is asymptotically efficient for (6.1) with $\Sigma_{\text{omd}}^\tau = \Sigma_{\text{oiv}}^\tau$. Moreover the optimal weighting is simply

$$W_{0j}^{\text{oiv}} = - \left(\sum_{j=1}^{\tau} \alpha_j \Gamma_j^\top \right)^{-1} \alpha_j \Gamma_j^\top \text{ for } j = 1, \dots, \tau, \text{ with}$$

$$(\alpha_1, \dots, \alpha_\tau) = \Gamma^\tau{}^\top \Psi_\tau^{-1}.$$

Example 1 (*cont.*)

Recall the optimal GMM estimator in this model [i.e., under homoskedasticity, etc.] is simply the two stage least squares estimator

$$\tilde{\theta} = (Y_2^\top P_X Y_2)^{-1} Y_2^\top P_X Y_1,$$

where $P_X = X(X^\top X)^{-1}X^\top$, $Y_1 = (y_{11}, \dots, y_{1n})^\top$, $Y_2 = (y_{21}, \dots, y_{2n})^\top$, $X = (X_1^\top, \dots, X_n^\top)$, $X_i = (X_{1i}, \dots, X_{ki})^\top$. Within our class of estimators \mathcal{E} , the optimal estimator is

$$\hat{\theta} = \sum_{j=1}^k W_{nj}^{\text{opt}} \hat{\theta}_j = (i_k^\top V^{-1} i_k)^{-1} i_k^\top V^{-1} \begin{bmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_k \end{bmatrix},$$

where $\hat{\theta}_j = (Y_2^\top P_j Y_2)^{-1} Y_2^\top P_j Y_1$ for $j = 1, \dots, k$, where $P_j = X_j(X_j^\top X_j)^{-1}X_j^\top$ and V is the $k \times k$ covariance matrix with $V_{jl} = \text{asy. cov}(\hat{\theta}_j, \hat{\theta}_l)$.¹¹ Suppose that the instruments are mutually orthogonal, then it is easy to see that $\hat{\theta}$ is identically equal to $\tilde{\theta}$.¹² This gives yet another interpretation to 2SLS as being the optimal combination of exactly identified instrumental variables estimators.¹³ Furthermore, $\hat{\theta}$ is computationally feasible even when $k > n$. This problem known as undersized sample problem arises in almost all large macroeconomic models, because there are usually more exogenous variables than time periods of observations (see Klein (1973) for an early account of alternative methods to solve this problem).

6.2 Case 2: Increasing τ

Here we consider the more general case where τ increases with sample size. Let Σ_{oiv} be the asymptotic variance of the optimal instrumental variable (oiv) estimator, and let Σ_{omd} be the asymptotic variance as $n \rightarrow \infty$ and $\tau(n) \rightarrow \infty$ of the optimal minimum distance (omd) estimator. The next theorem establishes that one can basically interchange the operations under additional assumptions.

ASSUMPTION C:

¹¹There is a connection with portfolio theory. Think of the estimators $\hat{\theta}_j$ as being returns on asset j , where each asset has the same expected return but different variances. The optimal weights are the same as the weights for the global minimum variance portfolio, see Campbell, Lo, and MacKinlay (1997, pp 184-185). This is also related to the idea of combining many forecasts, see Stock and Watson (1999) and Granger (2000) for example.

¹²We are grateful to Tom Rothenberg for pointing this out to us.

¹³Interpreting 2SLS in various ways has a long history in econometrics; see Rothenberg (1974) for an early example.

(C1) The matrix $D_0(X_i) = \left(\frac{\partial}{\partial \theta'} E[\rho(Z_i, \theta) | X_i]\right) |_{\theta=\theta_0}$ exists with probability one.

(C2) $E[\sigma_0^{-2}(X_i) D_0(X_i) D_0(X_i)^\top]$ is finite and positive definite.

(C3) $D_0(X_i) = \sum_{j=1}^{\infty} \beta_{j0} \phi_j(X_i) \sigma_0^2(X_i)$, where the sequence $\{\phi_j\}$ is a complete orthonormal basis satisfying:

$$E[\sigma_0^2(X_i) \phi_j(X_i) \phi_l(X_i)^\top] = \begin{cases} 0_p & \text{for } j \neq l, \\ I_p & \text{for } j = l. \end{cases}$$

Theorem 3 *Suppose that $E[\rho(Z_i, \theta_0) | X_i] = 0$ and that Assumptions C1–C3 hold. Then,*

$$\Sigma_{\text{omd}} = \Sigma_{\text{oiv}} = \left(E[\sigma_0^{-2}(X_i) D_0(X_i) D_0(X_i)^\top]\right)^{-1}.$$

The optimal weights in this case are any sequence like

$$W_{nj}^0 = \left(\sum_{l=1}^{\tau(n)} V_{ll}^{-1}\right)^{-1} V_{jj}^{-1},$$

where V_{jj} is the asymptotic variance matrix of $\sqrt{n}(\hat{\theta}_j - \theta_0)$. With such a sequence of weights, $\hat{\theta}$ has the same asymptotic variance as a comparable implementation of $\tilde{\theta}$. Note that in the scalar homoskedastic case, the optimal weights W_{nj}^0 decrease at the same rate as β_{j0}^2 as $j \rightarrow \infty$, while the weights on the basis terms in the estimation of D_0 would decrease like β_{j0} as $j \rightarrow \infty$. This suggests that one needs to combine fewer estimators than instruments to achieve a specified variance.

6.3 Weights Estimation

In practice, one must use estimated weights. In case 1, consistent estimators can be constructed as

$$\widehat{W}_{0j}^{\text{opt}} = \left(\sum_{l=1}^{\tau} \widehat{B}_l\right)^{-1} \widehat{B}_j, \quad (6.7)$$

where $(\widehat{B}_1, \dots, \widehat{B}_\tau) = (I_p \otimes i_\tau)^\top \widehat{V}^{-1}$, and \widehat{V} having (j, l) sub-matrix calculated using formulae (4.3)–(4.4) for $j, l = 2, \dots, \tau$. In the second case, weights W_{nj}^0 are proportional to the inverse of the asymptotic variance; therefore, given consistent estimators \widehat{V}_{jj} of the asymptotic variances $V_{jj} = \text{var}(\sqrt{n}\hat{\theta}_j)$, we can let

$$\widehat{W}_{nj}^0 = \left(\sum_{l=1}^{\tau(n)} \widehat{V}_{ll}^{-1}\right)^{-1} \widehat{V}_{jj}^{-1}. \quad (6.8)$$

The issues surrounding estimating the optimal weights for similar problems have been treated in Newey (1990) and Koenker and Machado (1999). We do not pursue this further here, but refer the reader to these other papers.

7 Monte Carlo

We consider two data generating processes (DGPs). The first one is adapted from Newey (1990) who consider an endogenous dummy variable model with the following specification:

$$\begin{aligned} Y_i &= \beta_{10} + \beta_{20}s_i + \varepsilon_i; \\ \text{DGP1: } s_i &= 1(\alpha_{10} + \alpha_{20}X_i + \eta_i > 0), \\ X_i &\sim N(0, 1); \quad \alpha_{10} = \alpha_{20} = \beta_{10} = \beta_{20} = 1, \end{aligned}$$

where the errors ε_i and η_i are generated as

$$\begin{bmatrix} \varepsilon_i \\ \eta_i \end{bmatrix} \sim N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \varphi \\ \varphi & 1 \end{bmatrix}\right), \quad (7.1)$$

in which $\varphi \in \{0.2, 0.5, 0.8\}$ indicate weak, medium and strong endogeneity respectively. The optimal instrument for s is $\pi(x) = \Pr[s = 1|X = x]$, which makes $D(x) = (1, \pi(x))$.

Tables 1 and 2 report results for two estimators of β_{20} . The first estimator corresponds to Newey's (1990) and the second is ours. To obtain both estimators, we follow Newey (1990) and use the polynomials $A_j(x) = x^{j-1}$ and $A_j(x) = [x/(1 + |x|)]^{j-1}$ as basis. Newey's (1990) estimator becomes

$$\begin{aligned} \tilde{\beta}_{\text{oiv}} &= \begin{pmatrix} \tilde{\beta}_{10; \text{oiv}} \\ \tilde{\beta}_{20; \text{oiv}} \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n s_i \\ \sum_{i=1}^n \hat{\pi}(X_i) & \sum_{i=1}^n \hat{\pi}(X_i)s_i \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n \hat{\pi}(X_i)Y_i \end{pmatrix}, \\ \hat{\pi}(x) &= \sum_{j=1}^{\tau} \hat{\gamma}_j A_j(x). \end{aligned}$$

for series-based estimated weights $\hat{\gamma}_j$. Using the same basis, we calculate our estimator as

$$\begin{aligned} \hat{\beta}_{\text{omd}} &= \sum_{j=2}^{\tau} \widehat{W}_{0j}^{\text{opt}} \hat{\beta}_j, \text{ where} \\ \hat{\beta}_j &= \begin{pmatrix} \hat{\beta}_{10; j} \\ \hat{\beta}_{20; j} \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n s_i \\ \sum_{i=1}^n A_j(X_i) & \sum_{i=1}^n A_j(X_i)s_i \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n A_j(X_i)Y_i \end{pmatrix}, \end{aligned}$$

and $\widehat{W}_{0j}^{\text{opt}}$ calculated as in (6.7). We consider two sample sizes: $n = 100, 200$ and 1000 replications. Simulated bias, standard deviation (Std. Dev.) and root mean squared errors (RMSE) are reported for each estimator.

Tables 1 and 2 show that for each sample size and endogeneity parameter value (φ) under consideration, the RMSE associated with the proposed estimator of β_{20} is roughly comparable to that of Newey's (1990) for low values of τ . While biases are small for both sets of estimates, their variance behave quite differently with relation to τ . In particular, the precision of $\widetilde{\beta}_{\text{oiv}}$ increases with τ , while that of $\widehat{\beta}_{20; \text{omd}}$ actually decreases. This might be caused by the estimation error in $\widehat{W}_{0j}^{\text{opt}}$. Table 2 shows that the proposed estimator outperforms Newey's (1990) for $\tau = 2$ and 3 when $n = 100$ and 200 respectively.

In DGP2 we consider a two-equation system with the following specification:

$$\begin{aligned}
 Y_i &= \beta_{10} + \beta_{20}s_i + \varepsilon_i; \\
 \text{DGP2: } s_i &= \alpha_{0;0} + \sum_{l=1}^k \alpha_{l;0}X_{li} + \eta_i, \\
 X_i &\sim N(0, I_k); \quad \alpha_{0;0} = \alpha_{1;0} = \dots = \alpha_{k;0} = \beta_{10} = \beta_{20} = 1,
 \end{aligned}$$

where $X_i = (X_{1i}, \dots, X_{ki})^\top$ and $(\varepsilon_i, \eta_i)^\top$ are generated as in (7.1). We set $k = 30$ and assess the performance of the 'optimal' estimator here, i.e. weights determined by (6.8) and $\tau = k$, with undersized samples of $n = 15$ and 25. Table 3 shows the results. The proposed estimator shows small biases and decreasing (with respect to sample size) variances for each endogeneity parameter value. Notice that for these sample sizes generic 2SLS cannot be performed. However, for a sample size of 50 observations, the variance and RMSE of the proposed estimator are comparable to that of 2SLS.

8 Conclusions and Extensions

Our approach has an advantage over the traditional approach to semiparametric instrumental variables in that one has a 'distribution' of estimators of the same quantity and one can view the range of values that these estimators take. If that range is not great, then it would appear that achieving efficiency is not going to be worth very much. If the range is considerable, then the efficient estimator may be very much better than any given estimator but at the same time performance might be very sensitive to how it is constructed. This information contained in the spread of the different estimators

is similar to but not necessarily the same as the information contained in the standard error of an efficient estimator.¹⁴ Also, the optimal weighting just requires the estimation of HAC matrices, at least in the orthonormal basis case, about which much has been written in econometrics.

It is quite straightforward to extend our work to produce results for the range¹⁵

$$\mathfrak{R}_n = \max_{1 \leq j \leq \tau(n)} \hat{\theta}_j - \min_{1 \leq j \leq \tau(n)} \hat{\theta}_j$$

using the theory of extreme values for Gaussian processes [as in Bickel and Rosenblatt (1973) for example]. This statistic can be used as another way of measuring whether the observed range is consistent with the underlying model assumptions, i.e., as a model specification test.

Appendix A: Mathematical Proofs

Proof of Lemma 1. We show that

$$\Pr \left[\max_{1 \leq j \leq \tau(n)} \left| \sum_{i=1}^n U_{ji} \right| \geq \lambda_n \right] \rightarrow 0$$

for any $\lambda_n = \delta_n \sqrt{n}$. For an array $\chi_{nj} \rightarrow \infty$ as $n \rightarrow \infty$ for each j , write

$$U_{ji} = U_{ji} \mathbf{1}(|U_{ji}| \leq \chi_{nj}) + U_{ji} \mathbf{1}(|U_{ji}| > \chi_{nj}) = \tilde{U}_{ji} + \tilde{\tilde{U}}_{ji}.$$

We shall assume for simplicity that U_{ji} is symmetric about zero so that $E(\tilde{U}_{ji}) = 0$. Therefore, \tilde{U}_{ji} are i.i.d. for each j with mean zero and are bounded from above by χ_{nj} . By the Bonferroni and Bernstein inequalities

$$\begin{aligned} \Pr \left[\max_{1 \leq j \leq \tau(n)} \left| \sum_{i=1}^n \tilde{U}_{ji} \right| \geq \lambda_n \right] &\leq \sum_{j=1}^{\tau(n)} \Pr \left[\left| \sum_{i=1}^n \tilde{U}_{ji} \right| \geq \lambda_n \right] \\ &\leq \sum_{j=1}^{\tau(n)} \exp \left(\frac{-\lambda_n^2}{s_{nj}^2 + 2\lambda_n \chi_{nj}} \right). \end{aligned} \tag{A-1}$$

¹⁴Actually, if the estimators themselves are mutually independent with the same limiting distribution, then the 95% confidence interval of a single estimator is approximately the same as the inter hemi-decile range, that is the interval $[\hat{\theta}_{0.025 \cdot \tau}, \hat{\theta}_{0.975 \cdot \tau}]$ of the ordered estimators. In fact, it is not possible that the estimators come from the same asymptotic distribution [since the variances must diverge along some trajectory], and so the two intervals do not coincide. Nevertheless, the connection exists.

¹⁵In the multiparameter case, we take the coordinate-wise ranges.

We shall choose λ_n and χ_{nj} below to make this term vanish.

By the Bonferroni and Markov inequalities

$$\begin{aligned}
\Pr \left[\max_{1 \leq j \leq \tau(n)} \left| \sum_{i=1}^n \tilde{U}_{ji} \right| \geq \lambda_n \right] &\leq \sum_{j=1}^{\tau(n)} \Pr \left[\left| \sum_{i=1}^n \tilde{U}_{ji} \right| \geq \lambda_n \right] \\
&\leq \sum_{j=1}^{\tau(n)} \frac{E \left(\left| \sum_{i=1}^n \tilde{U}_{ji} \right|^\kappa \right)}{\lambda_n^\kappa} \\
&\leq \sum_{j=1}^{\tau(n)} \frac{n^\kappa E(|U_{ji}|^\kappa) \Pr[|U_{ji}| > \chi_{nj}]}{\lambda_n^\kappa} \\
&\leq \sum_{j=1}^{\tau(n)} \frac{n^\kappa [E(|U_{ji}|^\kappa)]^2}{\lambda_n^\kappa \chi_{nj}^\kappa} = o(1)
\end{aligned}$$

provided $\sum_{j=1}^{\tau(n)} n^\kappa \chi_{nj}^{-\kappa} \lambda_n^{-\kappa} c_j^2 \rightarrow 0$.

Letting $\lambda_n = \delta_n \sqrt{n}$ and $\chi_{nj} = \sigma_j^2 \sqrt{n}$ we need to show that:

$$\sum_{j=1}^{\tau(n)} \exp \left(\frac{-\delta_n}{\sigma_j^2} \right) \rightarrow 0 \text{ and } \frac{1}{\delta_n^\kappa} \sum_{j=1}^{\tau(n)} \frac{c_j^2}{\sigma_j^{2\kappa}} \rightarrow 0.$$

For the first condition it suffices that

$$\frac{\delta_n}{\max_{1 \leq j \leq \tau(n)} \sigma_j^2 \log \tau(n)} \rightarrow \infty.$$

For the second condition it certainly suffices if

$$\frac{\delta_n}{\left(\sum_{j=1}^{\tau(n)} c_j^2 \sigma_j^{-2\kappa} \right)^{1/\kappa}} \rightarrow \infty.$$

Proof of Theorem 1 (i). From [A3](#), if $\max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| > \delta$, then $\|G_j(\hat{\theta}_j)\| \geq \epsilon_n(\delta)$ for some j . Consequently

$$\Pr \left(\max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| > \delta \right) \leq \Pr \left(\max_{1 \leq j \leq \tau(n)} \|G_j(\hat{\theta}_j)\| \geq \epsilon_n(\delta) \right), \tag{A-2}$$

and it is sufficient to prove that for the given $\epsilon_n(\delta) > 0$, the latter probability goes to zero. But

$$\begin{aligned}
\max_{1 \leq j \leq \tau(n)} \|G_j(\widehat{\theta}_j)\| &\leq \max_{1 \leq j \leq \tau(n)} \|G_j(\widehat{\theta}_j) - G_{nj}(\widehat{\theta}_j)\| + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\widehat{\theta}_j)\| \text{ by the Triangle Inequality,} \\
&\leq \max_{1 \leq j \leq \tau(n)} \sup_{\theta \in \Theta} \|G_j(\theta) - G_{nj}(\theta)\| + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\widehat{\theta}_j)\| \text{ by set inclusion,} \\
&= o_p(\alpha_{2n}) + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\widehat{\theta}_j)\| \text{ by A5,} \\
&= o_p(\alpha_{2n}) + \max_{1 \leq j \leq \tau(n)} \left(\|G_{nj}(\widehat{\theta}_j)\| - \inf_{\theta \in \Theta} \|G_{nj}(\theta)\| \right) + \max_{1 \leq j \leq \tau(n)} \inf_{\theta \in \Theta} \|G_{nj}(\theta)\|, \\
&\leq o_p(\alpha_{2n}) + \max_{1 \leq j \leq \tau(n)} \left(\|G_{nj}(\widehat{\theta}_j)\| - \inf_{\theta \in \Theta} \|G_{nj}(\theta)\| \right) + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\theta_0)\|, \\
&= o_p(\alpha_{2n}) + o_p(\alpha_{1n}) = o_p(\epsilon_n(\delta)) \text{ by A4, A5 and A2.}
\end{aligned}$$

We conclude that $\max_{1 \leq j \leq \tau(n)} \|\widehat{\theta}_j - \theta_0\| = o_p(1)$. Finally,

$$\|\widehat{\theta} - \theta_0\| \leq \sum_{j=1}^{\tau(n)} \|W_{nj}\| \times \max_{1 \leq j \leq \tau(n)} \|\widehat{\theta}_j - \theta_0\| = o_p(1)$$

by A1.

Proof of Theorem 1 (ii). Consistency implies that for every $\epsilon > 0$ there exists a sequence $\{\delta_n\}$, with $\delta_n \rightarrow 0$, and an N such that for all $n \geq N$,

$$\Pr\{\|\widehat{\theta} - \theta_0\| > \delta_n\} \leq \epsilon.$$

The discussion of subsequent properties can confine itself to conditions that need only hold in “shrinking neighbourhoods” of θ_0 ; i.e., neighbourhoods of θ_0 that can get arbitrarily small as n grows large, and still we know that our estimator will have that property with probability tending to one. Using the same proof as that of Theorem 1 (i), we have under our stronger assumption A*5 that with probability tending to one

$$\begin{aligned}
\max_{1 \leq j \leq \tau(n)} \|G_j(\widehat{\theta}_j)\| &\leq \max_{1 \leq j \leq \tau(n)} \|G_j(\widehat{\theta}_j) - G_{nj}(\widehat{\theta}_j)\| + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\widehat{\theta}_j)\| \\
&\leq \max_{1 \leq j \leq \tau(n)} \sup_{\|\theta - \theta_0\| \leq \delta_n} \|G_j(\theta) - G_{nj}(\theta)\| + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\widehat{\theta}_j)\| \\
&= o_p(n^{-1/4}) + \max_{1 \leq j \leq \tau(n)} \|G_{nj}(\widehat{\theta}_j)\| \text{ by A*5,} \\
&= o_p(n^{-1/4}) \text{ by A*4, A*5 and A*2.}
\end{aligned}$$

Therefore, by A*3

$$\begin{aligned} \Pr \left(\max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| > \delta_n \right) &\leq \Pr \left(\max_{1 \leq j \leq \tau(n)} \|G_j(\hat{\theta}_j)\| \geq \delta_n c_n \right) \\ &\rightarrow 0 \quad \text{if} \quad \delta_n c_n = O(n^{-1/4}). \end{aligned}$$

Hence

$$\max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| = o_p(n^{-1/4}),$$

which implies that

$$\|\hat{\theta} - \theta_0\| \leq \sum_{j=1}^{\tau(n)} \|W_{nj}\| \times \max_{1 \leq j \leq \tau(n)} \|\hat{\theta}_j - \theta_0\| = o_p(n^{-1/4})$$

as required, where $\sum_{j=1}^{\tau(n)} \|W_{nj}\|$ is uniformly bounded by [A1](#).

Proof of Theorem 2. Let

$$L_{nj}(\theta) = G_{nj}(\theta_0) + \Gamma_j(\theta - \theta_0)$$

for each $j = 1, 2, \dots$. Then define θ_j^* as the minimizer of $\|L_{nj}(\theta)\|$ over $\theta \in \mathbb{R}^p$ [Note that θ_j^* minimizes over \mathbb{R}^p , and not over Θ . We ignore this difference below because θ_j^* will eventually be in Θ with probability going to one]. The solution satisfies

$$\sqrt{n}(\theta_j^* - \theta_0) = -\Gamma_j^{-1} \sqrt{n} G_{nj}(\theta_0) \tag{A-3}$$

for each j . Therefore,

$$\begin{aligned} \sqrt{n} \sum_{j=1}^{\tau(n)} W_{nj}(\theta_j^* - \theta_0) &= \sqrt{n} \sum_{j=1}^{\tau(n)} W_{nj}^0(\theta_j^* - \theta_0) + \sqrt{n} \sum_{j=1}^{\tau(n)} (W_{nj} - W_{nj}^0)(\theta_j^* - \theta_0) \\ &= \sum_{i=1}^n T_{in} + R_n, \end{aligned}$$

where $R_n = \sqrt{n} \sum_{j=1}^{\tau(n)} (W_{nj} - W_{nj}^0)(\theta_j^* - \theta_0)$ and $T_{in} = \frac{-1}{\sqrt{n}} \sum_{j=1}^{\tau(n)} W_{nj}^0 \Gamma_j^{-1} g_j(Z_i, \theta_0)$.

The result follows after we establish:

- (i) $\sum_{i=1}^n c^\top T_{in} \implies \mathcal{N}(0, c^\top \Sigma c)$ for any $c \in \mathbb{R}^p$ with $\|c\| = 1$;
- (ii) The remainder term $R_n = o_p(1)$;
- (iii) $\sqrt{n} \sum_{j=1}^{\tau(n)} W_{nj}(\theta_j^* - \hat{\theta}_j) = o_p(1)$.

For (i), the triangular array of random variables $c^\top T_{in}$ is mean zero and independent across i for each n . By **B5(a)** we have:

$$\begin{aligned} \sum_{i=1}^n E[c^\top T_{in}]^2 &= E \left[\left(\sum_{j=1}^{\tau(n)} c^\top W_{nj}^0 \Gamma_j^{-1} g_j(Z_i, \theta_0) \right)^2 \right] \\ &= \sum_{j=1}^{\tau(n)} \sum_{l=1}^{\tau(n)} c^\top W_{nj}^0 \Gamma_j^{-1} E [g_j(Z_i, \theta_0) g_l(Z_i, \theta_0)^\top] \Gamma_l^{-1 \top} W_{nl}^{0 \top} c \\ &\rightarrow c^\top \Sigma c . \end{aligned}$$

And by **B5(b)** we have for some $\kappa > 0$,

$$\sum_{i=1}^n E|c^\top T_{in}|^{2+\kappa} \rightarrow 0.$$

Hence we obtain (i) by applying the Liapounov's triangular array central limit theorem.

For (ii), notice that Assumption **B3(a)** and **(A-3)** imply that $\max_{1 \leq j \leq \tau(n)} \|\Gamma_j \sqrt{n}(\theta_j^* - \theta_0)\| = O_p(1)$.

This together with Assumption **B4(a)** imply (ii) because

$$\begin{aligned} \left\| \sqrt{n} \sum_{j=1}^{\tau(n)} (W_{nj} - W_{nj}^0)(\theta_j^* - \theta_0) \right\| &\leq \sqrt{n} \max_{1 \leq j \leq \tau(n)} \|\Gamma_j(\theta_j^* - \theta_0)\| \sum_{j=1}^{\tau(n)} \|(W_{nj} - W_{nj}^0)\Gamma_j^{-1}\| \\ &= O_p(1) \times o_p(1). \end{aligned}$$

For (iii), by the $n^{1/4}$ -consistency result, there exists a positive sequence $\eta_n \rightarrow 0$ such that $\Pr[n^{1/4} \|\hat{\theta} - \theta_0\| > \eta_n] \rightarrow 0$. For each j we have

$$\begin{aligned} G_{nj}(\theta) &= G_{nj}(\theta_0) + G_j(\theta) + G_{nj}(\theta) - G_j(\theta) - G_{nj}(\theta_0) \\ &= L_{nj}(\theta) + O(\|\theta - \theta_0\|^2) + [G_{nj}(\theta) - G_j(\theta)] - G_{nj}(\theta_0) \text{ by } \mathbf{B2}. \end{aligned}$$

Therefore, for the above η_n and constants a and C we have

$$\begin{aligned} &\max_{1 \leq j \leq \tau(n)} \sup_{\|\theta - \theta_0\| \leq a\eta_n/n^{1/4}} \sqrt{n} \|G_{nj}(\theta) - L_{nj}(\theta)\| \\ &\leq C \times \eta_n^2 a^2 + \max_{1 \leq j \leq \tau(n)} \sup_{\|\theta - \theta_0\| \leq a\eta_n/n^{1/4}} \sqrt{n} \|[G_{nj}(\theta) - G_j(\theta)] - G_{nj}(\theta_0)\| \\ &= O_p(\eta_n^2) + o_p(1) = o_p(1) \text{ by } \mathbf{B3(b)}. \end{aligned}$$

Therefore,

$$\max_{1 \leq j \leq \tau(n)} \|\sqrt{n}[L_{nj}(\theta_j^*) - G_{nj}(\theta_j^*)]\| = o_p(1), \text{ and } \max_{1 \leq j \leq \tau(n)} \|\sqrt{n}[L_{nj}(\hat{\theta}_j) - G_{nj}(\hat{\theta}_j)]\| = o_p(1)$$

because θ_j^* is \sqrt{n} -consistent and $\hat{\theta}_j$ is $o(n^{-1/4})$ -consistent. It now follows from the definition of θ_j^* and Assumption **B1** and the triangular inequality that

$$\max_{1 \leq j \leq \tau(n)} \left| \sqrt{n}\|L_{nj}(\theta_j^*)\| - \sqrt{n}\|L_{nj}(\hat{\theta}_j)\| \right| = o_p(1). \quad (\text{A-4})$$

This implies that $\max_{1 \leq j \leq \tau(n)} \|\Gamma_j \sqrt{n}(\theta_j^* - \hat{\theta}_j)\| = o_p(1)$, because of the properties of least squares residuals. Then we have

$$\begin{aligned} \sqrt{n} \sum_{j=1}^{\tau(n)} W_{nj}(\theta_j^* - \hat{\theta}_j) &\leq \sum_{j=1}^{\tau(n)} \|W_{nj} \Gamma_j^{-1}\| \times \max_{1 \leq j \leq \tau(n)} \|\Gamma_j \sqrt{n}(\theta_j^* - \hat{\theta}_j)\| \\ &\leq O_p(1) \times o_p(1) = o_p(1), \end{aligned}$$

where the last inequality is due to Assumption **B4**(a) and (b) since

$$\begin{aligned} \sum_{j=1}^{\tau(n)} \|W_{nj} \Gamma_j^{-1}\| &\leq \sum_{j=1}^{\tau(n)} \|W_{nj}^0 \Gamma_j^{-1}\| + \sum_{j=1}^{\tau(n)} \|(W_{nj} - W_{nj}^0) \Gamma_j^{-1}\| \\ &= O(1) + o_p(1) = O_p(1), \end{aligned}$$

the result (iii) follows.

Proof of Proposition 1. On the one-hand, by the results in Hansen (1982), the optimal GMM (oiv) estimator is asymptotically efficient among all regular \sqrt{n} -asymptotic normal estimators for the moment restrictions (6.1), hence $\Sigma_{\text{oiv}}^\tau \leq \Sigma_{\text{omd}}^\tau$ in the positive semi-definite matrix sense. On the other hand, we notice that the oiv (optimal GMM) estimator has the expansion

$$\sqrt{n}(\tilde{\theta}_{\text{oiv}}^\tau - \theta_0) = -(\Gamma^{\tau\top} \Psi_\tau^{-1} \Gamma^\tau)^{-1} \Gamma^{\tau\top} \Psi_\tau^{-1} \sqrt{n} G_n^\tau(\theta_0) + o_p(1),$$

which can be rewritten as

$$\sqrt{n}(\tilde{\theta}_{\text{oiv}}^\tau - \theta_0) = - \left(\sum_{j=1}^{\tau} \alpha_j \Gamma_j^\top \right)^{-1} \sum_{j=1}^{\tau} \alpha_j \sqrt{n} G_{nj}(\theta_0) + o_p(1), \quad (\text{A-5})$$

where $\Gamma^{\tau\top}\Psi_{\tau}^{-1} = (\alpha_1, \dots, \alpha_{\tau})$ with $\alpha_j \in \mathbb{R}^{p \times p}$, and $\Gamma^{\tau} = (\Gamma_1^{\top}, \dots, \Gamma_{\tau}^{\top})^{\top}$ with $\Gamma_j = E[A_j(X)D_0(X)^{\top}]$, and $G_{nj}(\theta) = \frac{1}{n} \sum_{i=1}^n A_j(X_i)\rho(Z_i, \theta)$ for $j = 1, \dots, \tau$. That is, the optimal GMM (oiv) estimator $\tilde{\theta}_{\text{oiv}}^{\tau}$ belongs to the class of linear combinations of the $\hat{\theta}_j$, $j = 1, \dots, \tau$ with

$$\tilde{\theta}_{\text{oiv}}^{\tau} = \sum_{j=1}^{\tau} W_{0j}^{\text{oiv}} \hat{\theta}_j + o_p(n^{-1/2}),$$

and

$$W_{0j}^{\text{oiv}} = - \left(\sum_{j=1}^{\tau} \alpha_j \Gamma_j^{\top} \right)^{-1} \alpha_j \Gamma_j^{\top} \text{ for } j = 1, \dots, \tau.$$

However, by the results in Rothenberg (1973), $\hat{\theta}_{\text{omd}}^{\tau} = \sum_{j=1}^{\tau} W_{0j}^{\text{opt}} \hat{\theta}_j$ is asymptotically efficient among the regular class of estimators of the form $\sum_{j=1}^{\tau} W_{0j} \hat{\theta}_j$ with $\sum_{j=1}^{\tau} W_{0j} = I_p$, hence $\Sigma_{\text{omd}}^{\tau} \leq \Sigma_{\text{oiv}}^{\tau}$ in the positive semi-definite matrix sense. Therefore $\Sigma_{\text{omd}}^{\tau} = \Sigma_{\text{oiv}}^{\tau}$ in (6.4).

Proof of Theorem 3. Assumption C3 implies that $\beta_{j0} = E[D_0(X_i)\phi_j(X_i)^{\top}]$. We have:

$$\begin{aligned} \Sigma_{\text{oiv}} &= (E[\sigma_0^{-2}(X_i)D_0(X_i)D_0(X_i)^{\top}])^{-1} = \left(E \left[\sum_{j=1}^{\infty} \beta_{j0} \phi_j(X_i) \sigma_0^2(X_i) \sigma_0^{-2}(X_i) D_0(X_i)^{\top} \right] \right)^{-1} \\ &= \left(\sum_{j=1}^{\infty} \beta_{j0} E[\phi_j(X_i)D_0(X_i)^{\top}] \right)^{-1} = \left(\sum_{j=1}^{\infty} \beta_{j0} \beta_{j0}^{\top} \right)^{-1} \\ &= \left(\sum_{j=1}^{\infty} E[D_0(X_i)^{\top} \phi_j(X_i)^{\top}] E[\phi_j(X_i)D_0(X_i)] \right)^{-1}. \end{aligned}$$

Assumptions C2 and C3 imply that $0 < \sum_{j=1}^{\infty} \beta_{j0} \beta_{j0}^{\top} < \infty$.

By Assumptions C1–C3, we have $V_{jj} = \{\Gamma_j^{\top} \Gamma_j\}^{-1} = \{E[\phi_j(X_i)D_0(X_i)]^{\top} E[\phi_j(X_i)D_0(X_i)]\}^{-1}$ and $V_{jl} = 0$ for all $j \neq l$. Therefore

$$\begin{aligned} \Sigma_{\text{omd}} &= \lim_{\tau \rightarrow \infty} ((I_p \otimes i_{\tau})^{\top} V^{-1} (I_p \otimes i_{\tau}))^{-1} \\ &= \lim_{\tau \rightarrow \infty} \left(\sum_{j=1}^{\tau} V_{jj}^{-1} \right)^{-1} \\ &= \lim_{\tau \rightarrow \infty} \left(\sum_{j=1}^{\tau} \{E[\phi_j(X_i)D_0(X_i)]^{\top} E[\phi_j(X_i)D_0(X_i)]\} \right)^{-1} \\ &= \left(\sum_{j=1}^{\infty} \beta_{j0} \beta_{j0}^{\top} \right)^{-1}. \end{aligned}$$

References

- Ai, C., Chen, X., 2003. Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica* 71, 1795–1843.
- Amemiya, T., 1974. The non-linear two-stage least squares estimator. *Journal of Econometrics* 2, 105–110.
- Andrews, D.W.K., 1991. Asymptotic normality of series estimators for nonparametric and semi-parametric regression models. *Econometrica* 59, 307–346.
- Bickel, P.J., Rosenblatt, M., 1973. On some global measures of the deviations of density function estimates. *Annals of Statistics* 1, 1071–1095.
- Bierens, H.J., 1987. Kernel estimators of regression functions. Cambridge University Press: *Advances in Econometrics*.
- Breiman, L., 1996. Bagging Predictors. *Machine Learning* 26, 123–140.
- Breiman, L., 1996. Using Adaptive Bagging to Debias Regressions. University of California Berkeley, Department of Statistics, Technical Report 547.
- Campbell, J., Lo, A., Mackinlay, A.C., 1997. *The Econometrics of Financial Markets*. Princeton University Press: Princeton.
- Chamberlain, G., 1987. Asymptotic efficiency in estimation with conditional moment restrictions. *Journal of Econometrics* 34, 305–334.
- Granger, C.W.J., 2000. Thick modeling. UCSD Department of Economics Working Paper.
- Gray, H.L., Schucany, W.R., 1972. *The generalized jackknife statistic*. New York: Dekker.
- Hansen, L.P., 1982. Large sample properties of generalized method of moments estimators. *Econometrica* 50, 1029–1054.
- Hansen, L.P., 1985. A method for calculating bounds on the asymptotic covariance matrices of generalized method of moments estimators. *Journal of Econometrics* 30, 203–238.

- Huber, P., 1967. The behavior of maximum likelihood estimates under nonstandard conditions. In Fifth Berkeley Symposium on Mathematical Statistics and Probability, 221–233. Berkeley, CA: University of California.
- Klein, L., 1973. The treatment of undersized samples in econometrics. In *Econometric Studies of Macro and Monetary Relations* edited by Alan A. Powell and Ross A. Williams. American Elsevier Publishing Company, Inc.
- Koenker, R., Bassett, G., 1978. Regression quantiles. *Econometrica* 46, 33–50.
- Koenker, R., Machado, J., 1999. GMM inference when the number of moment conditions is large. *Journal of Econometrics* 93, 327–344.
- Kotlyarova, Y., Zinde-Walsh, V., 2007. Robust kernel estimator for densities of unknown smoothness. *Journal of Nonparametric Statistics*, forthcoming.
- Kotlyarova, Y., Zinde-Walsh, V., 2006. Non- and semi-parametric estimation in models with unknown smoothness. *Economics Letters*, 93, 379–386.
- McFadden, D., 1989. A method of simulated moments for estimation of discrete response models without numerical integration. *Econometrica* 57, 995–1026.
- Newey, W.K., 1990. Efficient instrumental variables estimation of nonlinear models. *Econometrica* 58, 809–837.
- Newey, W.K., 1993. Efficient estimation of models with conditional moment restrictions. In *Handbook of Statistics*, vol. 11. G.S. Maddala, C.R. Rao, and H.D. Vinod eds., Amsterdam: North-Holland.
- Newey, W.K., 1994. The asymptotic variance of semiparametric estimators. *Econometrica* 62, 1349–1382.
- Newey, W.K., 2004. Efficient semiparametric estimation via moment restrictions. *Econometrica* 72, 1877–1897.
- Newey, W.K., McFadden, D.F., 1994. Large sample estimation and hypothesis testing. In *Handbook of Econometrics*, vol. IV, eds. D.F. McFadden and R.F. Engle III. North Holland.

- Newey, W.K., Powell, J., 1990. Efficient estimation of linear and type I censored regression models under conditional quantile restrictions. *Econometric Theory* 6, 295–317.
- Newey, W.K., Smith, R.J. 2004. Higher order properties of GMM and Generalized Empirical Likelihood Estimators. *Econometrica* 72, 219–255.
- Olley S., Pakes, A., 1996. The dynamics of productivity in the telecommunications equipment industry. *Econometrica* 64, 1263–1297.
- Pakes, A., Pollard D., 1989. Simulation and the asymptotics of optimization estimators. *Econometrica* 57, 1027–1057.
- Portnoy, S., 1985. Asymptotic behavior of M estimators of p regression parameters when p^2/n is large: II normal approximation. *Annals of Statistics* 17, 382–400.
- Rilstone, P., Srivastava, V. K., Ullah, A., 1996. The second-order bias and mean squared error of nonlinear estimators. *Journal of Econometrics* 75, 369–395.
- Robinson, P. M., 1991. Best nonlinear three-stage least squares estimation of certain econometric models. *Econometrica* 59, 755–786.
- Rothenberg, T.J., 1973. Efficient estimation with a priori information. Cowles monograph.
- Rothenberg, T.J., 1974. A note on two-stage least squares. Unpublished manuscript, UC Berkeley.
- Sawa, T., 1973. Almost unbiased estimator in simultaneous equations systems. *International Economic Review* 14, 97–106.
- Schafgans, M., Zinde-Walsh, V., 2007. Robust Average Derivative Estimation. Unpublished Manuscript.
- Stock, J.W., Watson, M.W., 1999. Forecasting inflation. *Journal of Monetary Economics* 44, 293–335.
- Watson, M.W., 2000. Macroeconomic forecasting using many predictors. Princeton University.

Table 1

DGP1: $A_j(x) = x^{j-1}$, $j = 2, 3, \dots, 7$

τ	(A)			(B)			(A)			(B)		
	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE
	$\varphi = 0.2, n = 100$						$\varphi = 0.2, n = 200$					
2	0.004	0.489	0.488	0.004	0.489	0.488	-0.014	0.319	0.319	-0.014	0.319	0.319
3	0.008	0.451	0.451	0.002	0.487	0.487	-0.009	0.300	0.300	-0.022	0.314	0.314
4	0.007	0.442	0.442	0.019	0.523	0.524	-0.009	0.299	0.299	0.012	0.361	0.361
5	0.008	0.438	0.437	0.012	0.669	0.669	-0.003	0.297	0.297	-0.020	0.449	0.449
6	0.020	0.434	0.434	0.043	0.884	0.885	0.001	0.296	0.296	0.021	0.623	0.623
7	0.023	0.431	0.431	0.042	1.073	1.072	0.004	0.295	0.295	0.015	0.734	0.734
	$\varphi = 0.5, n = 100$						$\varphi = 0.5, n = 200$					
2	0.014	0.470	0.470	0.014	0.470	0.470	0.000	0.327	0.327	0.000	0.327	0.327
3	0.027	0.442	0.441	0.011	0.465	0.465	-0.001	0.309	0.309	0.001	0.324	0.324
4	0.042	0.431	0.431	0.033	0.515	0.514	0.003	0.302	0.302	0.000	0.371	0.371
5	0.057	0.422	0.423	0.010	0.643	0.644	0.006	0.296	0.296	0.009	0.471	0.471
6	0.074	0.418	0.420	-0.040	0.915	0.916	0.014	0.294	0.294	-0.051	0.642	0.643
7	0.092	0.414	0.419	-0.051	1.102	1.103	0.024	0.293	0.294	-0.009	0.749	0.750
	$\varphi = 0.8, n = 100$						$\varphi = 0.8, n = 200$					
2	-0.014	0.503	0.504	-0.014	0.503	0.504	0.001	0.354	0.355	0.001	0.354	0.355
3	0.006	0.473	0.473	0.028	0.493	0.493	0.007	0.337	0.337	0.019	0.317	0.317
4	0.025	0.460	0.460	0.026	0.532	0.533	0.013	0.331	0.331	0.026	0.374	0.374
5	0.047	0.443	0.444	0.027	0.669	0.670	0.028	0.326	0.326	0.040	0.459	0.459
6	0.061	0.436	0.438	-0.022	0.995	0.994	0.034	0.323	0.324	0.058	0.630	0.631
7	0.092	0.427	0.433	0.042	1.203	1.205	0.040	0.321	0.322	0.019	0.727	0.727

Notes: Number of replications = 1000. (A) = $\hat{\beta}_{20; \text{oiv}}$, (B) = $\hat{\beta}_{20; \text{omd}}$

Table 2

DGP1: $A_j(x) = [x/(1 + |x|)]^{j-1}$, $j = 2, 3, \dots, 7$

τ	(A)			(B)			(A)			(B)		
	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE
	$\varphi = 0.2, n = 100$						$\varphi = 0.2, n = 200$					
2	-0.014	0.493	0.494	-0.014	0.493	0.494	-0.017	0.328	0.328	-0.017	0.328	0.328
3	0.005	0.458	0.458	-0.005	0.383	0.383	-0.015	0.304	0.304	-0.012	0.252	0.252
4	0.002	0.448	0.448	0.028	0.384	0.384	-0.002	0.297	0.297	-0.001	0.244	0.244
5	0.022	0.440	0.440	0.023	0.405	0.405	-0.001	0.298	0.298	0.001	0.246	0.246
6	0.027	0.435	0.435	0.037	0.478	0.478	0.000	0.294	0.294	0.007	0.278	0.278
7	0.027	0.430	0.430	0.017	0.489	0.489	0.004	0.293	0.293	0.012	0.282	0.282
	$\varphi = 0.5, n = 100$						$\varphi = 0.5, n = 200$					
2	0.029	0.485	0.485	0.029	0.485	0.485	-0.001	0.333	0.333	-0.001	0.333	0.333
3	0.034	0.451	0.451	0.052	0.375	0.376	-0.003	0.309	0.309	-0.001	0.258	0.258
4	0.047	0.427	0.427	0.040	0.365	0.365	0.012	0.299	0.299	-0.002	0.252	0.252
5	0.074	0.417	0.419	0.026	0.394	0.394	0.013	0.296	0.296	-0.009	0.259	0.259
6	0.085	0.415	0.418	0.041	0.459	0.459	0.024	0.294	0.294	-0.001	0.298	0.298
7	0.102	0.409	0.415	0.031	0.468	0.468	0.035	0.291	0.292	-0.003	0.300	0.300
	$\varphi = 0.8, n = 100$						$\varphi = 0.8, n = 200$					
2	-0.009	0.511	0.513	-0.009	0.511	0.513	-0.001	0.358	0.358	-0.001	0.358	0.359
3	0.014	0.474	0.474	0.023	0.384	0.384	0.004	0.337	0.338	0.012	0.263	0.263
4	0.048	0.454	0.454	0.023	0.386	0.387	0.026	0.330	0.329	0.014	0.258	0.258
5	0.058	0.442	0.443	0.017	0.410	0.410	0.037	0.325	0.325	0.008	0.271	0.271
6	0.079	0.436	0.441	0.030	0.492	0.492	0.044	0.321	0.322	0.022	0.293	0.294
7	0.108	0.426	0.436	0.032	0.502	0.502	0.060	0.315	0.318	0.022	0.296	0.296

Notes: Number of replications = 1000. (A) = $\hat{\beta}_{20; \text{oiv}}$, (B) = $\hat{\beta}_{20; \text{omd}}$

Table 3
DGP2: Undersized Sample Problem

n	β_{10}			β_{20}		
	Bias	Std. Dev.	RMSE	Bias	Std. Dev.	RMSE
$\varphi = 0.2$						
15	0.009	0.272	0.272	0.008	0.056	0.057
25	0.002	0.211	0.211	0.004	0.042	0.042
50	0.000	0.146	0.146	0.004	0.029	0.029
	-0.006*	0.145*	0.145*	0.005*	0.026*	0.026*
$\varphi = 0.5$						
15	-0.003	0.264	0.264	0.013	0.055	0.057
25	-0.002	0.204	0.204	0.011	0.043	0.044
50	-0.006	0.144	0.144	0.007	0.030	0.031
	-0.016*	0.144*	0.145*	0.009*	0.027*	0.028*
$\varphi = 0.8$						
15	-0.016	0.276	0.276	0.020	0.054	0.058
25	-0.016	0.205	0.205	0.018	0.042	0.046
50	-0.012	0.144	0.145	0.015	0.029	0.032
	-0.013*	0.144*	0.145*	0.016*	0.027*	0.031*

Notes: Number of replications = 1000. (*) = 2SLS estimator.