

EFFICIENT MULTIVARIATE QUANTILE REGRESSION ESTIMATION

Sung Jae Jun* and Joris Pinkse†
The Pennsylvania State University

November 2006

Abstract

We propose an efficient semiparametric estimator for the multivariate linear quantile regression model in which the conditional joint distribution of errors given regressors is unknown. The procedure can be used to estimate multiple conditional quantiles of the same regression relationship. The proposed estimator is asymptotically as efficient as if the conditional distribution were known. Simulation results suggest that the estimation procedure works well in practice and dominates an equation-by-equation efficiency correction if the errors are dependent conditional on the regressors.

*(corresponding author) Department of Economics, The Pennsylvania State University, 608 Kern Graduate Building, University Park PA 16802, sjun@psu.edu

†joris@psu.edu We thank participants at the Carnegie Mellon departmental seminar and Ari Kang for their useful suggestions.

1 Introduction

We propose an efficient semiparametric estimator for the multivariate linear quantile regression model in which the conditional joint distribution of errors given regressors is unknown. The procedure can be used to estimate multiple conditional quantiles of the same regression relationship. The proposed estimator is asymptotically as efficient as if the conditional distribution were known. Simulation results suggest that the estimation procedure works well in practice and dominates an equation-by-equation efficiency correction if the errors are dependent conditional on the regressors.

The proposed method entails the nonparametric estimation of optimal instruments for a set of moment conditions corresponding to the conditional quantiles of interest and subsequently using these estimated optimal instruments to obtain the efficient quantile estimates. Two-step efficiency corrections like this go back to at least Aitken (1935), but semiparametric corrections like ours have been around for a while, also. Carroll (1982), Delgado (1992) and Robinson (1987) achieve full GLS¹ asymptotic efficiency by estimating the conditional error variance function nonparametrically. Newey (1990, 1993) proposes methods for estimating optimal instruments nonparametrically, thereby allowing for multivariate regressions and ones with endogenous regressors. Pinkse (2006) introduces a method which addresses the *curse of dimensionality* associated with the nonparametric estimation of functions with many arguments. Finally, Zhao (2001), Whang (2006) and Komunjer and Vuong (2006) propose efficiency corrections for the univariate quantile regression model.

The multivariate quantile regression case is of interest for applied work for several reasons. First, even absent dependence between errors and regressors quantile regression estimators tend to have greater asymptotic variances than mean regression ones² and efficiency improvements are hence more valuable. Further, an optimal parametric correction in the mean regression model requires one to guess the correct parametric form of the conditional variance function (matrix-valued in the multivariate case), which is

¹ Generalized Least Squares

² The asymptotic relative efficiency for a median regression estimator versus a mean regression estimator for a model with normally distributed errors is $2/\pi$. There are cases in which the median regression estimator has a smaller asymptotic variance, e.g. for the Laplace distribution. Please note that median and mean regression estimators typically estimate different coefficients.

difficult since little reliable information may be available as to its shape. In the quantile case, one would need to know the marginal conditional error densities at zero plus, in the multivariate case, the probability for each pair of errors that both are negative, conditional on the regressors. It is even more unrealistic for an empirical researcher to possess that much information; incorrect guesses will lead to inefficient estimators, quite possibly to ones that have lesser asymptotic efficiency than uncorrected ones.³ Unless the errors are independent conditional on the regressors *and* there are no cross-equation restrictions on the regression coefficients, multivariate efficiency corrections are moreover generally more efficient than univariate ones. Finally, with quantile estimation it is possible to estimate multiple quantiles of the same regression relationship (i.e. the same dependent variable and the same regressors) simultaneously, which would imply strong dependence between the corresponding errors and hence more scope for efficiency improvements.

Like all of the above semiparametric estimators ours relies on the availability of a \sqrt{n} -consistent first round estimator; a natural choice is the standard quantile regression estimator. A problem with such a two-step procedure is that the first round estimation error, while asymptotically absent, can be such that correction is not worthwhile in small samples. This is especially true when the number of regressors is large due to the fact that nonparametric estimators of high-dimensional functions are notoriously inaccurate. Please note however, that our correction does not require (nor do we establish) pointwise consistent estimation of the optimal instruments. Further, since the uncorrected estimates are special cases of the correction procedure for particular values of the input parameters of the semiparametric procedure, the semiparametric procedure is in principle never worse irrespective of the sample size. Please note, however, that we offer no procedure for the optimal selection of the input parameters. Earlier work (Pinkse, 2006) and some experimentation (not reported) suggest that our procedure is comparatively insensitive to the choice of input parameter.

This paper contains several theoretical innovations. While Newey (1990, 1993) allows for multiple equations to be estimated jointly, his results do not cover the current case because of nondifferentiability issues.

³ In the univariate case it can be reasonable to assume that errors factor as the product of a function of regressors and some error independent of the regressors, see e.g. Koenker (2005), section 5.3.2, and Koenker and Zhao (1994).

Zhao (2001), Whang (2006) and Komunjer and Vuong (2006) propose estimators for the single equation case. In the single equation case the nuisance function is just conditional error density at zero instead of the product of a matrix and the inverse of another matrix, as is the case here. Whang (2006) and Komunjer and Vuong (2006) achieve the semiparametric efficiency bound (the latter for time series) by optimizing an objective function involving a series expansion of the nuisance function; the nondifferentiability problems we solve do not arise then.

Our paper is closer to Zhao (2001) in that we use a nonparametric plugin estimator, but they differ in several dimensions. Zhao's results only cover the single-equation case and are not readily generalized to the multivariate case. Further, Zhao uses a less primitive technical condition on the construction of the weights, while we have specifically opted for nearest neighbor estimation; we believe that nearest neighbor weights satisfy Zhao's condition. A final difference concerns the way in which the first step estimation error is addressed. For both methods (Zhao's and ours) first step estimates enter the second step via a nondifferentiable function. Zhao proposes two distinct procedures to address this problem. The first procedure entails *sample splitting*, i.e. using half the data to get the first step estimator to be used as a plug-in in the second step for the second half of the data and vice versa. This procedure does not make a difference asymptotically, but is less attractive due to its inherent (finite sample) inefficiency. His second procedure assumes that the first step estimator has a certain *Bahadur representation*, which we believe is likely to hold in practice.

We do not know whether Zhao's methods can be extended to the multivariate case. Instead, we follow a new line of proof which neither entails sample-splitting nor does it require any assumptions on the first step estimator beyond a convergence rate. The new proof (contained in the last two lemmas of Appendix C and using L1 of Appendix A) entails ratcheting up of the established uniform convergence rate of the feasible estimator of the moment condition and the feasible estimator of the parameter vector of interest alternately. This method of proof has uses that go well beyond the particular problem at hand or indeed differentiability problems or ones involving nonparametric estimation.

To compute our estimates we propose procedures involving solving standard linear programming problems possibly combined with taking a Newton step . The procedure is guaranteed to yield estimates

satisfying our constraints — we prove this — and does so fast; computing the nonparametric weights takes the most time. The reason that computation here is simple, in contrast to e.g. Chernozhukov and Hansen’s (2006) estimator, is that we have an initial easily computable \sqrt{n} -consistent but inefficient estimator at our disposal, namely the standard least absolute deviations estimator. The Matlab code is available from the authors on request.

The outline of the paper is as follows. In section 2 we introduce the setup and define our estimator. Section 3 contains the theoretical results for our estimator, whose computation and performance are studied in section 4. Section 5 concludes.

2 Model and Estimator

Let $\{y_i, X_i\}$ be an i.i.d. sequence for which

$$Q(y_i|X_i) = X_i'\theta_0 \text{ a.s.}, \quad i = 1, \dots, n, \quad (1)$$

or equivalently,

$$y_i = X_i'\theta_0 + u_i, \quad Q(u_i|X_i) = 0 \text{ a.s.}, \quad i = 1, \dots, n, \quad (2)$$

where $y_i \in \mathbb{R}^d$, $X_i \in \mathbb{R}^{K \times d}$ and Q denotes the vector of quantiles of interest.

The formulations in (1) and (2) allow for several possibilities. The restriction that the regression coefficients are the same in all regression equations is not restrictive because we can make the choices

$$X_i = \begin{bmatrix} x_{i1} & & & \\ & \ddots & & \\ & & & x_{id} \end{bmatrix}, \quad \theta_0 = \begin{bmatrix} \theta_{01} \\ \vdots \\ \theta_{0d} \end{bmatrix},$$

resulting in

$$y_{ij} = x'_{ij}\theta_{0j} + u_{ij}, \quad i = 1, \dots, n; \quad j = 1, \dots, d. \quad (3)$$

So (1) allows for arbitrary amounts of overlap between the vectors of regression coefficients across equations. An assumption implicit in (1) is that all regressors in equation ℓ which enter the conditional quantile

function in equation j are also regressors in equation j . This is where part of the efficiency gain originates; it is akin to an orthogonality condition between regressors and errors across equations in the mean regression case.⁴ It is possible to choose $y_{ij} = y_{i\ell}$, $x_{ij} = x_{i\ell}$, $j \neq \ell$, for all i in (3) if different regression quantiles of the same regression relationship are desired. Assuming multiple quantiles of the same relationship to all be linear, however, imposes strong restrictions on the types of dependence between errors and regressors that can be accommodated and a procedure that exploits such restrictions will likely work better in practice than the more general procedure proposed here; a more fruitful avenue would be to estimate the median and mean jointly, a possibility not covered by our results.

We now formulate an infeasible efficient estimation procedure for θ_0 . Let $s_i(\theta) = I(y_i \leq X_i'\theta) - \tau$, where τ is the vector indicating which quantiles are desired (a vector with values 0.5 in case of the median) and I is the *indicator function*, where for any $v \in \mathbb{R}^{d_v}$, $I(v) = [I(v_1), \dots, I(v_{d_v})]'$. Then the conditional moment condition is ($s_i = s_i(\theta_0)$)

$$E(s_i|X_i) = 0 \text{ a.s..}$$

The corresponding optimal unconditional moment conditions are

$$E(A_i s_i) = 0, \tag{4}$$

where $A_i = S_i' T_i^{-1}$ with

$$S_i = F_i X_i', \quad F_i = \begin{bmatrix} f_{u_{i1}|X_i}(0) & & \\ & \ddots & \\ & & f_{u_{id}|X_i}(0) \end{bmatrix}, \quad T_i = E(s_i s_i' | X_i). \tag{5}$$

The asymptotic variance of an infeasible estimator $\hat{\theta}_I$ based on (4) will later be shown to be V^{-1} with

$$V = E(A_1 s_1 s_1' A_1') = E(S_1' T_1^{-1} S_1). \tag{6}$$

The proposed procedure yields a natural efficiency improvement over equation-by-equation estimation

⁴ It is possible to obtain efficiency improvements when the conditional quantiles do not depend on some but not all of the regressors in another equation; this possibility can be accommodated in our setup by a judicious choice of y and X .

when there is overlap between the regression coefficients across equations. Absent such overlap, the asymptotic variance of $\hat{\theta}_{I1}$, the infeasible estimator of the first subvector θ_{01} , is for $d = 2$ equal to

$$V_{I1} = \left(E[t_i^{11} f_{i1}^2 x_{i1} x_{i1}'] - E[t_i^{12} f_{i1} f_{i2} x_{i1} x_{i2}'] \{ E[t_i^{22} f_{i2}^2 x_{i2} x_{i2}'] \}^{-1} E[t_i^{12} f_{i1} f_{i2} x_{i2} x_{i1}'] \right)^{-1},$$

where $f_{ij} = f_{u_{ij}|X_i}(0)$ and $t_i^{j\ell}$ is the (j, ℓ) -element of

$$T_i^{-1} = \frac{1}{t_{i11}t_{i22} - t_{i12}^2} \begin{bmatrix} t_{i22} & -t_{i12} \\ -t_{i12} & t_{i11} \end{bmatrix}; \quad T_i = \begin{bmatrix} t_{i11} & t_{i12} \\ t_{i12} & t_{i22} \end{bmatrix}.$$

The corresponding asymptotic variances for the inefficient and efficient single equation estimators are

$$V_{SI1} = \tau_1(1 - \tau_1)(E[\tilde{f}_{i1} x_{i1} x_{i1}'])^{-1} E(x_{i1} x_{i1}') (E(\tilde{f}_{i1} x_{i1} x_{i1}'))^{-1}, \quad V_{SE1} = \tau_1(1 - \tau_1)(E[\tilde{f}_{i1}^2 x_{i1} x_{i1}'])^{-1},$$

where $\tilde{f}_{i1} = f_{u_{i1}|x_{i1}}(0)$. It is necessarily true that $V_{I1} \leq V_{SE1} \leq V_{SI1}$; we now discuss when they are equal. When \tilde{f}_{i1} does not depend on x_{i1} , $V_{SI1} = V_{SE1}$; otherwise equality only occurs in exceptional cases.

Using information from different equations is useful because one can exploit (i) the information that regressors in equation 2 do not impact the conditional quantile of equation 1 and (ii) the fact that u_{i1} and u_{i2} are not necessarily independent conditional on X_i . Consideration (i) can be accommodated in the single equation case (as in Zhao (2001)) by extending the conditioning set to regressors outside of the equation being estimated; in the multivariate case the conditioning vector can likewise be extended (and the efficiency thereby improved) by including variables from outside the system. But (ii) cannot be used in the single equation setup.

So even if the regressors in both equations are the same and $\tilde{f}_{i1} = f_{i1}$, there is still an efficiency gain from our method unless u_{i1}, u_{i2} are independent conditional on X_i ,⁵ in which case $t_i^{12} = t_{i12} = 0$, or if u_{i1}, u_{i2} do not depend on X_i . Conversely, even if u_{i1}, u_{i2} are independent of X_i there is still an efficiency gain unless $x_{i1} = x_{i2}$. All of this is similar to a SUR model with random regressors where no efficiency gain obtains from joint estimation if the errors are uncorrelated conditional on the regressors or if the regressors

⁵ Or more precisely: if $I(u_{i1} \leq 0)$ and $I(u_{i2} \leq 0)$ are independent conditional on X_i .

are identical and independent of the errors.⁶ Table 1 in the appendix contains the full details of when efficiency improvements obtain for the various estimators.

If the errors are known to be independent of the regressors, then no nonparametric correction is needed since only the joint distribution of $I(u_{ij} \leq 0)$ with $I(u_{i\ell} \leq 0)$ for all j, ℓ is needed, and this distribution entails only $d(d-1)/2$ unknowns. The types of dependence between errors and regressors that lead to efficiency improvements is different from the mean regression case. In the mean regression case efficiency improvements obtain only if $\Sigma(X_i) = V(u_i|X_i)$ varies with X_i whereas in the quantile regression case improvements obtain if the conditional error densities at zero vary with X_i or if $P(u_{ij} \leq 0, u_{i\ell} \leq 0|X_i)$ varies with X_i for some j, ℓ . Neither situation implies the other, except in special models like

$$u_i = (\Sigma(X_i))^{1/2} e_i, \quad (7)$$

where the elements of e_i are independent with unit variances and $\Sigma_i = \Sigma(X_i)$ is some positive definite matrix. The problem with (7) is that quantiles are generally not invariant to linear transformations, e.g. $\text{Med}(a+b) \neq \text{Med}(a) + \text{Med}(b)$. If the e_i 's are mean zero normal, however, then so are the u_i 's and their conditional median is zero.⁷ With (7), $f_{u_i|X_i}(0) = f_{e_i}(0)/\sqrt{|\Sigma_i|}$ and hence varies with X_i unless Σ_i is constant.

We now proceed with the formulation of our estimators. We begin with the infeasible estimator $\hat{\theta}_I$ which is defined as any estimator satisfying

$$m_n(\hat{\theta}_I) = o_p(n^{-1/2}), \quad \text{where } m_n(\theta) = n^{-1} \sum_{i=1}^n A_i s_i(\theta). \quad (8)$$

We do not set m_n equal to zero in (8) because no value of θ may exist that satisfies $m_n(\theta) = 0$ since s_i involves an indicator function. m_n converges to m with

$$m(\theta) = E[A_1 s_1(\theta)].$$

⁶ In the classical SUR model errors are assumed independent of the regressors, in which case no efficiency gain arises when the regressors are identical or the errors are uncorrelated.

⁷ This holds for any class of multivariate distributions that is closed to linear transformations and which are element-wise even.

$\hat{\theta}_I$ is infeasible since the A_i 's in (8) are unknown. We will estimate them and using their estimates \hat{A}_i we can define $\hat{\hat{\theta}}$ as any value satisfying

$$\hat{m}_n(\hat{\hat{\theta}}) = o_p(n^{-1/2}), \quad \text{where } \hat{m}_n(\theta) = n^{-1} \sum_{i=1}^n \hat{A}_i s_i(\theta). \quad (9)$$

The only remaining question is how to estimate A_i . Let $\hat{\theta}$ be any \sqrt{n} -consistent first stage estimator of θ_0 , e.g. based on single equation quantile estimation. We estimate T_i, S_i separately using KNN estimators

$$\hat{T}_i = n^{-1} \sum_{j=1}^n w_{ij} \hat{s}_j \hat{s}_j', \quad \hat{S}_i = n^{-1} \sum_{j=1}^n w_{ij} \hat{F}_j X_j' \quad (10)$$

where $\hat{s}_i = I(\hat{u}_i \leq 0) - \tau$, $\hat{F}_i = \text{diag}(I(|\hat{u}_i| \leq \beta_n \iota) / (2\beta_n))$ with ι a vector of ones, β_n a *bandwidth* parameter, $\hat{u}_i = y_i - X_i' \hat{\theta}$ and w_{ij} a KNN weight,⁸ setting $\hat{A}_i = \hat{S}_i' \hat{T}_i^{-1}$.

The KNN weights are all nonnegative and w_{ij} is positive only if observation j is among observation i 's k_n closest neighbors in terms of the distance between X_i and X_j ; ties only occur when all regressors are discrete and can be resolved by randomizing among the tying observations. The only other constraints we impose are upper and lower bounds to their values and conditions on the rate at which the number of neighbors should increase.

3 Results

We now discuss our main result, formulated in [T3](#), which shows that the feasible estimator $\hat{\hat{\theta}}$ has a limiting normal distribution with variance V^{-1} . For our main result, we need the following assumptions.⁹

A1 θ_0 is an interior point of the compact parameter space Θ .

A2 For some $C_T > 0$, $P(\lambda_{\min}(T_1) \geq C_T) = 1$.

A3 $E(X_1 X_1') > 0$.

⁸ See Newey and Powell (1990) for a similar use of \hat{F}_i .

⁹ We have not separated the assumptions by theorem since we are mostly concerned with [T3](#).

A4 For some $0 < C_f < \infty$, and all $j = 1, \dots, d$, $P(f_{u_{1j}|X_1}(0) \geq 1/C_f) > 0$, $P(f_{u_{1j}|X_1}(0) \leq C_f) = 1$, $P(\sup_t |f'_{u_{1j}|X_1}(t)| \leq C_f) = 1$ and $P(\sup_t |f''_{u_{1j}|X_1}(t)| \leq C_f) = 1$.

A5 $\forall \theta \in \Theta : m(\theta) = 0 \Leftrightarrow \theta = \theta_0$.

A6 The weights w_{ij} are nonnegative and all k_n nonzero weights take values in the range $[1/(C_w k_n); C_w/k_n]$.

A7 Let for any $p > 0$, $\zeta_{npT} = n^{1/p_x - 1/2} + n^{1/p} k_n^{-1/2}$ and $\zeta_{npS} = n^{1/p_x} k_n^{-1/2} \beta_n^{1/p_x - 1} + n^{1/p_x} \beta_n^2 + n^{1/2} k_n^{-1} \beta_n$. Then for some $p < \infty$, $\sqrt{n} \zeta_{npT}^2 \rightarrow 0$, $\sqrt{n} \zeta_{npT} \zeta_{npS} \rightarrow 0$ and $k_n/n \rightarrow 0$, as $n \rightarrow \infty$.

A1 and **A3** are standard. **A2** essentially says that $\text{Corr}[I(u_{i1} \leq 0), I(u_{i2} \leq 0)|X_i]$ should be a.s. bounded away from ± 1 ; this is reasonable and similar to a condition used in Pinkse (2006). The assumption (**A4**) that the conditional error densities have two uniformly bounded derivatives excludes distributions like the Laplace distribution, but is otherwise reasonable within the context of nonparametric estimation.¹⁰ The assumption that the conditional densities at zero are bounded away from zero with positive probability is needed for the invertibility of V . Further, **A6** is not a restriction on the model, but rather on how to choose the nearest neighbor weights and is hence innocuous.

That leaves **A5** and **A7**. **A5** is not primitive. It is a necessary and sufficient condition to ensure identification. In the univariate case **A5** is implied by **A2**, **A3** and **A4**, but we have failed to find a natural and primitive sufficient condition in the multivariate case. Finally, **A7** deals with the rate at which k_n increases. As long as a sequence exists that satisfies the restrictions, **A7** is merely a prescription on how to choose k_n . **A7** is for instance satisfied when $p_x = 6$, $\beta_n \sim k_n^{-3/17}$ and $k_n \sim n^{35/36}$. It can be shown that **A7** can only be satisfied for values of p_x greater than $3 + \sqrt{8}$. However, if an expansion taken in **L21** and **L22** in the appendix is taken beyond the second order the requirements would improve but would never be better than $\sqrt{n} \zeta_{npT}^o \rightarrow 0$, $\sqrt{n} \zeta_{npT}^{o-1} \zeta_{npS} \rightarrow 0$ where o denotes the order of the expansion. Since with cross-sectional data fat regressor tails are rarely an issue and the extension would merely involve a repetition of the same arguments, we have omitted it in the interest of brevity.

We now state our theorems.

¹⁰ The Laplace distribution could be accommodated since its density has bounded first left and right derivatives at zero, but this would come at the expense of longer proofs, stronger conditions on the value of p_x and more restrictive choices of $\{k_n\}$.

T1 For any estimator $\hat{\theta}_I$ satisfying (8), $\hat{\theta}_I \xrightarrow{p} \theta_0$.

T2 For any estimator $\hat{\theta}_I$ satisfying (8), $\sqrt{n}(\hat{\theta}_I - \theta_0) \xrightarrow{d} N(0, V^{-1})$.

T3 For any estimator $\hat{\theta}$ satisfying (9), $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, V^{-1})$.

For the purpose of hypothesis testing the matrix V needs to be estimated. The assumptions made are amply sufficient to guarantee convergence of our estimator \hat{V} of V .

T4 $\hat{V} = n^{-1} \sum_{i=1}^n \hat{A}_i \hat{S}_i \xrightarrow{p} V$.

4 Computation and Simulations

In this section, we report the results of a small simulation study and we discuss issues of computation.

We begin by outlining a simple method for the computation of estimates $\hat{\theta}$ that satisfy (9). This procedure entails taking one Newton step from any \sqrt{n} -consistent starting value, e.g. $\hat{\theta}_{(0)} = \hat{\theta}$, i.e. computing

$$\hat{\theta}_{(1)} = \hat{\theta}_{(0)} - \hat{V}^{-1}(\hat{\theta}_{(0)}) \hat{m}_n(\hat{\theta}_{(0)}),$$

This is a familiar procedure, where only the nondifferentiability issues provide minor complications; complications which were largely addressed in the earlier theorems.

T5 $\hat{\theta}_{(1)}$ solves (9).

Experience based on our simulations suggests that the above-described procedure often leads to an (undesirable) increase in the value of $\|\hat{m}_n\|$. We therefore propose an alternative procedure which can be based on $\hat{\theta}_{(1)}$ (or whichever of $\hat{\theta}_{(1)}$ and $\hat{\theta}_{(0)}$ which yields the smallest $\|\hat{m}_n\|$ value) to ensure that (9) remains satisfied, but it only works in case there is no overlap in the estimating equations.

In this case computing our estimator for a bivariate model entails solving the estimating equations ($\hat{\theta} = (\hat{\beta}', \hat{\gamma})'$ and ditto for similar symbols)

$$\begin{bmatrix} \hat{m}_{1n}(\hat{\beta}, \hat{\gamma}) \\ \hat{m}_{2n}(\hat{\beta}, \hat{\gamma}) \end{bmatrix} = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n \hat{a}_i x_{i1} (I\{y_{i1} \leq x'_{i1} \hat{\beta}\} - \tau_1) + \hat{b}_i x_{i1} (I\{y_{i2} \leq x'_{i2} \hat{\gamma}\} - \tau_2) \\ \frac{1}{n} \sum_{i=1}^n \hat{c}_i x_{i2} (I\{y_{i2} \leq x'_{i2} \hat{\gamma}\} - \tau_2) + \hat{d}_i x_{i2} (I\{y_{i1} \leq x'_{i1} \hat{\beta}\} - \tau_1) \end{bmatrix} = o_p(n^{-\frac{1}{2}}), \quad (11)$$

where $\hat{a}_i, \hat{b}_i, \hat{c}_i$, and \hat{d}_i are taken from the estimated optimal instruments as

$$\begin{bmatrix} \hat{a}_i & \hat{b}_i \\ \hat{d}_i & \hat{c}_i \end{bmatrix} = \begin{bmatrix} \hat{f}_{1i}(0|x_{i1}, x_{i2}) & 0 \\ 0 & \hat{f}_{2i}(0|x_{i1}, x_{i2}) \end{bmatrix} \hat{T}_i^{-1}. \quad (12)$$

When we have two median restrictions, we set $\tau_1 = \tau_2 = 0.5$. Note that the solutions are not unique, and that the estimating equation does not have to be exactly satisfied.

In order to find a solution that satisfies this equation, we consider the following strategy. Suppose that we start from $(\hat{\beta}_{(t)}, \hat{\gamma}_{(t)})$. Then, define $\hat{\beta}_{(t+1)}$ and $\hat{\gamma}_{(t+1)}$ to be solutions to the following linear programming (LP) problems:

$$\begin{cases} \hat{\beta}_{(t+1)} \equiv \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n \hat{a}_i \rho_{\tau_1}(y_{i1} - x'_{i1} \beta) + 2\hat{b}_i (I\{y_{i2} \leq x'_{i2} \hat{\gamma}_{(t)}\} - \tau_2) x'_{i1} \beta \\ \hat{\gamma}_{(t+1)} \equiv \arg \min_{\gamma} \frac{1}{n} \sum_{i=1}^n \hat{c}_i \rho_{\tau_2}(y_{i2} - x'_{i2} \gamma) + 2\hat{d}_i (I\{y_{i1} \leq x'_{i1} \hat{\beta}_{(t)}\} - \tau_1) x'_{i2} \gamma, \end{cases} \quad (13)$$

where $\rho_{\tau}(s) = |s| + (2\tau - 1)s$, which is reduced to $\rho_{\tau}(s) = |s|$ when $\tau = 0.5$. Note that the LP problems in (13) have the asymptotic first order conditions

$$\begin{cases} \|\hat{m}_{1n}(\hat{\beta}_{(t+1)}, \hat{\gamma}_{(t)})\| = \delta_{n,t} = o_p(n^{-1/2}) \\ \|\hat{m}_{2n}(\hat{\beta}_{(t)}, \hat{\gamma}_{(t+1)})\| = \eta_{n,t} = o_p(n^{-1/2}) \end{cases} \quad (14)$$

for $t = 1, 2, 3, 4, \dots$. We may therefore choose $(\hat{\beta}_{(t+1)}, \hat{\gamma}_{(t+1)})$ as a solution to equation 11, if

$$\begin{cases} \|\hat{m}_{1n}(\hat{\beta}_{(t+1)}, \hat{\gamma}_{(t+1)})\| \leq M\delta_{n,s} \\ \|\hat{m}_{2n}(\hat{\beta}_{(t+1)}, \hat{\gamma}_{(t+1)})\| \leq M\eta_{n,s'} \end{cases} \quad (15)$$

for some $s, s' = 1, 2, 3, \dots$, where M is a prespecified constant. Although the starting point $(\hat{\beta}_{(1)}, \hat{\gamma}_{(1)})$ does not matter in this algorithm, the first step estimators $\hat{\beta}, \hat{\gamma}$ are the most natural choice. The results reported in our experiments, used the second step estimators $(\hat{\hat{\beta}}, \hat{\hat{\gamma}}) = (\hat{\beta}_{(2)}, \hat{\gamma}_{(2)})$. This computational strategy naturally generalizes to the case with more than two equations.

The design of our experiment follows Zhao (2001), i.e.

$$\begin{cases} y_{i1} = \beta_{00} + x_{i11}\beta_{01} + x_{i21}\beta_{02} + u_{i1} \\ y_{i2} = \gamma_{00} + x_{i12}\gamma_{01} + x_{i22}\gamma_{02} + u_{i2}, \end{cases} \quad (16)$$

where $\beta_0 = [\beta_{00} : \beta_{01} : \beta_{02}]' = \gamma_0 = [\gamma_{00} : \gamma_{01} : \gamma_{02}]' = [10 : -4 : 2]'$. Like Zhao (2001) we generated the regressors by

$$\begin{aligned} x_{i11} &= N_{1i} + 0.2U_{i1}, \text{ and } x_{i21} = 0.2N_{1i} + U_{i1} \\ x_{i12} &= N_{2i} + 0.2U_{i2}, \text{ and } x_{i22} = 0.2N_{2i} + U_{i2}, \end{aligned}$$

where $N_{1i}, N_{2i} \sim N(5, 9)$, and $U_{i1}, U_{i2} \sim U(0, 4)$, and the errors follow

$$\begin{cases} u_{i1} = h(x_{i1}, x_{i2})\epsilon_{i1} \\ u_{i2} = g(x_{i1}, x_{i2})\epsilon_{i2}, \end{cases} \quad (17)$$

where ϵ_{i1} and ϵ_{i2} are jointly normally distributed with mean zero, variance one, and correlation ρ ; h and g are specifying different types of heteroskedasticity. Note that the regressors in one equation are allowed to enter the conditional error variance function of the other equation, but that in our experiments the correlation is fixed. The functional forms we used for h and g in the experiments are provided in Table 2.

For each scenario we computed regular quantile regression estimates (Unweighted LAD), within-equation efficient quantile estimation (à la Zhao (2001)) and our own methodology. The within-equation efficient method ignores the possibility that the conditional density function of one equation can depend on the regressors of the other equation. In particular, the weights for equation $t \in \{1, 2\}$ were computed by

$$\hat{f}_{it}^s = \sum_{j=1}^n w_{tij} I\{|\hat{\epsilon}_{jt}| \leq h_n\} \frac{1}{2h_n},$$

where h_n is a bandwidth, $\hat{\epsilon}_{tj}$ is a residual from the first step, and w_{tij} is a k -nearest neighbor weight based on the regressors of equation t . These weights are what we would use, if we tried efficient estimation within each equation separately. For our own methodology, weights were based on the stacked regressors from both equations. Specifically, we used

$$\begin{cases} \hat{f}_{it}^J = \sum_{j=1}^n w_{ij} I\{|\hat{\epsilon}_{jt}| \leq h_n\} \frac{1}{2h_n} \\ \hat{T}_i = \sum_{j=1}^n w_{ij} \hat{s}_j \hat{s}_j' \end{cases}$$

for $t \in \{1, 2\}$, where w_{ij} is the k -nearest neighbor weight based on the regressors of both equations, and $\hat{s}_j = [I\{y_{i1} \leq x'_{i1}\hat{\beta}\} - \tau_1 \mid I\{y_{i2} \leq x'_{i2}\hat{\beta}\} - \tau_2]'$. For the number of nearest neighbors and the bandwidth, we chose $k_n \propto n^{\frac{4}{5}}$, and $h_n \propto k_n^{-\frac{1}{6}}$.

Tables 3–8 contain the results. We are especially interested in the accuracy with which the slope coefficients are estimated. In table 3 all estimators are equally efficient, but the regular quantile regression estimator does somewhat better owing to the fact that it does not require the overhead of first round nonparametric estimation. When the errors are correlated, the proposed method dominates the others even in small samples, as illustrated in 4, with the uncorrected LAD marginally beating out the single equation efficient procedure.

The results in Table 5 imply that single equation–efficient estimation is preferable in moderate size samples when both nonparametric methods are fully efficient due to the additional overhead involved with joint nonparametric estimation. When the heteroskedastic form depends on regressors from both equations joint estimation appears to do somewhat better; see table 7.

The real benefit from estimating equations jointly, however, appears to materialize when there is error correlation. This is true absent heteroskedasticity, as discussed above, but continues to be the case in the presence of heteroskedasticity of unknown form, see tables 6 and 8. This conclusion, however, is likely due to the amount and nature of the correlation and heteroskedasticity specified.

5 Conclusions

We have proposed a new estimator for the efficient semiparametric estimation of systems of quantile regression estimation. Our simulation study shows that the new estimation method works well in practice.

Appendices

A Infeasible Estimator

Proof of T1: Consider the following class of functions:

$$\mathcal{F} \equiv \left\{ c' A_1 s_1(\theta) = \sum_{j=1}^d c' A_{1j} s_{1j}(\theta) : \theta \in \Theta \subset \mathbb{R}^D \right\},$$

where $c = [c_1, c_2, \dots, c_d]'$ is an arbitrary vector and A_{1j} is the j^{th} column vector of A_1 . Since $\mathcal{G}_j \equiv \{1(y_{1j} \leq X'_{1j}\theta) : \theta \in \Theta \subset \mathbb{R}^p\}$ is a Vapnik Červonenkis subgraph class (or simply VČ class),¹¹ it follows that $\mathcal{F}_j \equiv \{c' A_{1j} s_{1j}(\theta) : \theta \in \Theta \subset \mathbb{R}^p\}$ is also a VČ class by lemma 2.6.18 of van der Vaart and Wellner (1996). Since a VČ class is Euclidean for every envelope function (Pakes and Pollard (PP), 1989, lemma 2.12), we know that \mathcal{F}_j is Euclidean with envelope function $\mathcal{E}_j = c' A_{1j}$. Therefore, by lemma 2.14 of PP, \mathcal{F} is Euclidean with envelope function $\mathcal{E} = \sum_{j=1}^n \mathcal{E}_j$. Since $E(\mathcal{E}) < \infty$ by A3 and A4, it follows from lemma 2.8 of PP that

$$\sup_{\theta \in \Theta} |c' m_n(\theta) - c' m(\theta)| = o_p(1).$$

Since c is arbitrary, we have $\sup_{\theta \in \Theta} \|m_n(\theta) - m(\theta)\| = o_p(1)$. Now, by the triangle inequality

$$\|m(\hat{\theta}_I)\| \leq \|m_n(\hat{\theta}_I)\| + \|m(\theta) - m_n(\theta)\| = o_p(n^{-1/2}) + o_p(1) = o_p(1).$$

Hence, by assumptions A1, A4 and A5, $\hat{\theta}_I - \theta_0 = o_p(1)$. ■

L1 For any positive sequence $\{r_n\}$ and a consistent estimator θ_n , $m_n(\theta_n) = o_p(r_n)$ implies $\|\theta_n - \theta_0\| = O_p(n^{-1/2}) + o_p(r_n)$.

¹¹ See problem 14 on page 152 of van der Vaart and Wellner (1996).

Proof: Let $\{\delta_n\}$ be a sequence such that $P(\|\theta_n - \theta_0\| > \delta_n) = o(1)$. Then, recalling that $A_i s_i(\theta)$ is VČ,

$$\begin{aligned} \|m(\theta_n)\| &\stackrel{\text{triangle}}{\leq} \|m_n(\theta_n) - m(\theta_n)\| + \|m_n(\theta_n)\| \lesssim \sup_{\|\theta - \theta_0\| < \delta_n} \|m_n(\theta) - m(\theta)\| + o_p(r_n) \\ &\leq \sup_{\|\theta - \theta_0\| < \delta_n} \|m_n(\theta) - m(\theta) - m_n(\theta_0) + m(\theta_0)\| + \|m_n(\theta_0)\| + o_p(r_n) \\ &= o_p(n^{-1/2}) + O_p(n^{-1/2}) + o_p(r_n). \end{aligned} \quad (18)$$

A2, A3 and A4 imply that

$$m(\theta) = V(\theta - \theta_0) + o(\|\theta - \theta_0\|). \quad (19)$$

Hence

$$\lambda_{\min}(V)\|\theta_n - \theta_0\| \leq \|V(\theta_n - \theta_0)\| \leq \|m(\theta_n)\| + o_p(\|\theta_n - \theta_0\|),$$

which, together with the consistency of θ_n , implies that

$$(\lambda_{\min}(V) - o_p(1))\|\theta_n - \theta_0\| \leq \|m(\theta_n)\| = O_p(n^{-1/2}) + o_p(r_n).$$

Since V is positive definite, $\|\theta_n - \theta_0\| = O_p(n^{-1/2}) + o_p(r_n)$. ■

Proof of T2: First, recall that \mathcal{F} is a Euclidean class with envelope function $\mathcal{E} = \sum_{j=1}^d \mathcal{E}_j = \sum_{j=1}^d c' A_{1j}$. Note also that $E(\mathcal{E}^2) = c' \{\sum_{j=1}^d \sum_{t=1}^d E(A_{1j} A'_{1t})\} c < \infty$. Therefore, it follows from lemma 2.17 of PP that

$$\sup_{\|\theta - \theta_0\| < \delta_n} |\sqrt{n}(c' m_n(\theta) - c' m(\theta)) - \sqrt{n}(c' m_n(\theta_0) - c' m(\theta_0))| = o_p(1)$$

for any sequence $\{\delta_n\}$ with $\delta_n = o(1)$. Since c is arbitrary, it implies that

$$\sup_{\|\theta - \theta_0\| < \delta_n} \|\sqrt{n}(m_n(\theta) - m(\theta)) - \sqrt{n}(m_n(\theta_0) - m(\theta_0))\| = o_p(1)$$

The asserted result now follows from theorem 3.3 of PP. Specifically, note that by lemma L1, $\hat{\theta}_I - \theta_0 = O_p(n^{-1/2})$. Using derivations similar to those in (18) and (19) we have

$$\begin{aligned} o_p(n^{-1/2}) &= m_n(\hat{\theta}_I) = (m_n(\hat{\theta}_I) - m(\hat{\theta}_I) - m_n(\theta_0) + m(\theta_0)) + m(\hat{\theta}_I) + m_n(\theta_0) \\ &= o_p(n^{-1/2}) + V(\hat{\theta}_I - \theta_0) + o_p(n^{-1/2}) + m_n(\theta_0) = m_n(\theta_0) + V(\hat{\theta}_I - \theta_0) + o_p(n^{-1/2}). \end{aligned}$$

Hence since $E(A_1 s_1 s_1' A_1') = E(A_1 T_1 A_1') = E(A_1 F_1 T_1^{-1} F_1 A_1') = V > 0$,

$$\sqrt{n}(\hat{\theta}_I - \theta_0) = -V^{-1} \sqrt{n} m_n(\theta_0) + o_p(1) \xrightarrow{d} N(0, V^{-1}). \quad \blacksquare$$

B Nonparametric Approximation

In addition to $\hat{T}_i, T_i, \hat{S}_i, S_i$ we define

$$\tilde{T}_i = \sum_{j=1}^n w_{ij} s_j s_j', \quad \bar{T}_i = \sum_{j=1}^n w_{ij} T_j, \quad \tilde{S}_i = \sum_{j=1}^n w_{ij} \tilde{F}_j X_j', \quad \bar{S}_i = \sum_{j=1}^n w_{ij} S_j,$$

where $\tilde{F}_j = \text{diag}(I(|u_{jt}| \leq \beta_n)) / (2\beta_n)$.

B.1 Lemmas showing that $\max_i \|\hat{A}_i - \bar{A}_i\| = o_p(1)$.

Note that

$$\begin{aligned} \hat{A}_i - \bar{A}_i &= (\hat{S}_i' - \bar{A}_i \hat{T}_i) \hat{T}_i^{-1} = ((\hat{S}_i - \bar{S}_i)' - \bar{A}_i (\hat{T}_i - \bar{T}_i)) (\bar{T}_i^{-1} + (\hat{T}_i^{-1} - \bar{T}_i^{-1})) \\ &= \left((\hat{S}_i - \tilde{S}_i)' + (\tilde{S}_i - \bar{S}_i)' - \bar{A}_i ((\hat{T}_i - \tilde{T}_i) + (\tilde{T}_i - \bar{T}_i)) \right) \left(\bar{T}_i^{-1} + (\hat{T}_i^{-1} - \bar{T}_i^{-1}) \right). \end{aligned} \quad (20)$$

We deal with the uniform convergence of the differences in turn and then find a bound on the growth of \bar{A}_i .

B.1.1 $\tilde{T}_i - \bar{T}_i$

L2 $\exists \epsilon > 0 : \forall n : P(\min_i \lambda_{\min}(\bar{T}_i) < \epsilon) = 0$.

Proof:

$$P(\min_i \lambda_{\min}(\bar{T}_i) < \epsilon) \leq P(\min_i \lambda_{\min}(T_i) < \epsilon) = 0,$$

by **A2**. \blacksquare

L3 For any $p > 2$ for which $E(R_{ni}|X_i) = 0$ a.s. and $\limsup E\|R_{ni}\|^p < \infty$, $E\|\sum_{j=1}^n w_{ij}R_{nj}\|^p = O(k_n^{-p/2})$.

Proof: This is a special case of Pinkse (2006), L3, which was inspired by Robinson (1987), lemma 7. ■

L4 For any $\{\xi_{ni}\}$ for which $E\|\xi_{ni}\|^p < \infty$ for all i, n and any $\epsilon > 0$,

$$P(\max_i \|\xi_{ni}\| \geq \epsilon) \leq \epsilon^{-p} \sum_{i=1}^n E\|\xi_{ni}\|^p.$$

Proof: The LHS is bounded by $\sum_i P(\|\xi_{ni}\| \geq \epsilon)$ which is bounded by the RHS by the Markov inequality. ■

L5 For any $p > 2$ for which $E(R_{ni}|X_i) = 0$ a.s. and $\limsup E\|R_{ni}\|^p < \infty$, $\max_i \|\sum_{j=1}^n w_{ij}R_{nj}\| = O_p(n^{1/p}k_n^{-1/2})$.

Proof: Take $\xi_{ni} = n^{-1/p}k_n^{1/2} \sum_j w_{ij}R_j$ in L4 to obtain

$$P\left(\max_i \left\| n^{-1/p}k_n^{1/2} \sum_{j=1}^n w_{ij}R_j \right\| \geq \epsilon\right) \leq n^{-1}k_n^{p/2} \epsilon^{-p} \sum_{i=1}^n E\left\| \sum_{j=1}^n w_{ij}R_j \right\|^p \stackrel{\text{L3}}{=} O(1)\epsilon^{-p} \rightarrow 0,$$

as $\epsilon \rightarrow \infty$. ■

L6 For all values of $p > 2$, $\max_i \|\tilde{T}_i - \bar{T}_i\| = O_p(k_n^{-1/2}n^{1/p})$.

Proof: Use L5 with $R_i = s_i s'_i - T_i$. ■

B.1.2 $\hat{T}_i - \tilde{T}_i$

We will make frequent use of the inequality

$$\|\hat{s}_j \hat{s}'_j - s_j s'_j\| \leq \|\hat{s}_j - s_j\|^2 + \|s_j\| \cdot \|\hat{s}_j - s_j\| \leq C_s \|\hat{s}_j - s_j\|, \quad (21)$$

which holds for some $0 < C_s < \infty$ since both s_j and \hat{s}_j are vectors of zeroes and ones. We will also make multiple use of the inequality

$$\begin{aligned} \|\hat{s}_j - s_j\| &= \left| I(u_j \leq X'_j(\hat{\theta} - \theta_0)) - I(u_j \leq 0) \right| \leq \left| I(|u_j| \leq \|X_j\| \cdot \|\hat{\theta} - \theta_0\|) \right| \\ &\leq \left| I(|u_j| \leq \|X_j\| r_n) \right| + I(\|\hat{\theta} - \theta_0\| > r_n) = \|\alpha_{jr_n}\| + I(\|\hat{\theta} - \theta_0\| > r_n), \end{aligned} \quad (22)$$

which holds for any sequence $\{r_n\}$.

L7 For some $C > 0$ and any $r \geq 0$, $E(\|\alpha_{ir}\| | X_i) \leq C\|X_i\|r$ a.s.

Proof: Note that

$$0 \leq E(\alpha_{irj} | X_i) = P(|u_{ij}| \leq r \|X_i\| | X_i) = F_{u_{ij}|X_i}(r \|X_i\|) - F_{u_{ij}|X_i}(-r \|X_i\|) \stackrel{\text{A4}}{\leq} 2C_f \|X_i\| r. \quad \blacksquare$$

L8 For any $p > 0$, $\max_i \|\hat{T}_i - \tilde{T}_i\| = O_p(\zeta_{npT})$.

Proof: First,

$$\begin{aligned} C_s^{-1} \|\hat{T}_i - \tilde{T}_i\| &= C_s^{-1} \left\| \sum_{j=1}^n w_{ij} (\hat{s}_j \hat{s}'_j - s_j s'_j) \right\| \stackrel{(21)}{\leq} \sum_{j=1}^n w_{ij} \|\hat{s}_j - s_j\| \\ &\stackrel{(22)}{\leq} \sum_{j=1}^n w_{ij} (\|\alpha_{jr_n}\| - E(\|\alpha_{jr_n}\| | X_j)) + \sum_{j=1}^n w_{ij} E(\|\alpha_{jr_n}\| | X_j) + I(\|\hat{\theta} - \theta_0\| > r_n). \end{aligned} \quad (23)$$

Take $r_n = 1/(\sqrt{n} - \log n)$. Since $e^{-1/t}$ is an increasing function of t and for arbitrary positive a, b $I(a > b) \leq g(a)/g(b)$ for any increasing function g ,

$$I(\|\hat{\theta} - \theta_0\| > r_n) \leq e^{1/r_n} e^{-1/\|\hat{\theta} - \theta_0\|} = O_p(e^{1/r_n - \sqrt{n}}) = O_p(e^{-\log n}) = O_p(n^{-1}). \quad (24)$$

For the second RHS term in (23), note that

$$\max_i \sum_{j=1}^n w_{ij} E(\|\alpha_{jr_n}\| | X_j) \stackrel{\text{L7}}{\leq} C_\alpha r_n \max_i \sum_{j=1}^n w_{ij} \|X_j\| \leq C_\alpha r_n \max_i \|X_i\| \stackrel{\text{L4}}{=} O_p(r_n n^{1/p_x}) = O_p(n^{1/p_x - 1/2}).$$

Finally, noting that the $\|\alpha_{jr_n}\|$'s are uniformly bounded, **L5** implies that for any $p > 0$,

$$\max_i \left\| \sum_{j=1}^n w_{ij} (\|\alpha_{jr_n}\| - E(\|\alpha_{jr_n}\| | X_j)) \right\| = O_p(n^{1/p} k_n^{-1/2}),$$

which takes care of the first RHS term in (23). \blacksquare

L9 For any $p > 0$, $\max_i \|\hat{T}_i^{-1} - \bar{T}_i^{-1}\| = O_p(\zeta_{npT})$.

Proof: Since $\hat{T}_i^{-1} = \bar{T}_i^{-1} (I + (\hat{T}_i - \bar{T}_i) \bar{T}_i^{-1})^{-1}$, the result follows from lemmas **L2**, **L6** and **L8**. \blacksquare

B.1.3 $\tilde{S}_i - \bar{S}_i$

L10 $\max_i \|\bar{S}_i\| = O_p(n^{1/p_x})$ and $\max_i \|\bar{A}_i\| = O_p(n^{1/p_x})$.

Proof: Note that for some $0 < C < \infty$,

$$\max_i \|\bar{A}_i\| \leq \max_i \|\bar{S}_i\| \max_i \|\bar{T}_i^{-1}\| \stackrel{\text{L2}}{\leq} C \max_i \|\bar{S}_i\| \leq C \max_i \|S_i\| \stackrel{\text{A4}}{\leq} CC_f \max_i \|X_i\| \stackrel{\text{L4}}{=} O_p(n^{1/p_x}). \quad \blacksquare$$

L11 $\max_i \|\tilde{S}_i - \bar{S}_i\| = O_p(n^{1/p_x}(k_n^{-1/2}\beta_n^{1/p_x-1} + \beta_n^2))$.

Proof: Note that

$$\tilde{S}_i - \bar{S}_i = \sum_{j=1}^n w_{ij}(\tilde{F}_j - E(\tilde{F}_j|X_j))X_j' + \sum_{j=1}^n w_{ij}(E(\tilde{F}_j|X_j) - F_j)X_j'. \quad (25)$$

Take $R_{nj} = \beta_n^{1-1/p_x}(\tilde{F}_j - E(\tilde{F}_j|X_j))X_j'$ in **L5** to obtain the rate $O_p(n^{1/p_x}k_n^{-1/2}\beta_n^{1/p_x-1})$ for the first RHS term in (25). For the second RHS term note that by the mean value theorem for all $t = 1, \dots, d$,

$$\|E(\tilde{F}_{jt}|X_j) - F_{jt}\| = \|6^{-1}\beta_n^2 f''_{u_{jt}|X_j}(\cdot)\| \stackrel{\text{A4}}{\leq} 6^{-1}C_f\beta_n^2. \quad (26)$$

Hence the second RHS term in (25) is bounded by

$$6^{-1}C_f\beta_n^2 \max_i \sum_{j=1}^n w_{ij}\|X_j\| \leq 6^{-1}C_f\beta_n^2 \max_i \|X_i\| = O_p(n^{1/p_x}\beta_n^2). \quad \blacksquare$$

B.1.4 $\hat{S}_i - \tilde{S}_i$

L12 $\max_i \|\hat{S}_i - \tilde{S}_i\| = O_p(n^{1/2}k_n^{-1}\beta_n^{-1})$.

Proof: Let $r_n = 1/(\sqrt{n} - \log n)$. Now,

$$\max_i \|\hat{S}_i - \tilde{S}_i\| = \max_i \|\hat{S}_i - \tilde{S}_i\|I(\|\hat{\theta} - \theta_0\| \leq r_n) + \max_i \|\hat{S}_i - \tilde{S}_i\|I(\|\hat{\theta} - \theta_0\| > r_n). \quad (27)$$

By (24) $I(\|\hat{\theta} - \theta_0\| > r_n) = O_p(n^{-1})$, such that the second RHS term in (27) converges faster than the first. Now the first RHS term in (27). Using the inequality (for generic a, b, t)

$$|I(|a| \leq t) - I(|b| \leq t)| \leq I(|b| \leq t + |a - b|) - I(|b| \leq t - |a - b|),$$

it follows that

$$\|\hat{F}_j - \tilde{F}_j\| I(\|\hat{\theta} - \theta_0\| \leq r_n) \leq \left\| I(|u_j| \leq (\beta_n + \|X_j\|r_n)\iota) - I(|u_j| \leq (\beta_n - r_n\|X_j\|)\iota) \right\|, \quad (28)$$

and hence

$$\begin{aligned} \max_i \|\hat{S}_i - \tilde{S}_i\| I(\|\hat{\theta} - \theta_0\| \leq r_n) &\leq \max_i \sum_{j=1}^n w_{ij} \|X_j\| \cdot \|\hat{F}_j - \tilde{F}_j\| \\ &\leq \beta_n^{-1} \max_i \|X_i\| \max_i \sum_{j=1}^n w_{ij} \left\| I(|u_j| \leq (\beta_n + \|X_j\|r_n)\iota) - I(|u_j| \leq (\beta_n - r_n\|X_j\|)\iota) \right\| \\ &\stackrel{\text{A6}}{\leq} C_w (k_n \beta_n)^{-1} \sum_{j=1}^n \left\| I(|u_j| \leq (\beta_n + r_n\|X_j\|)\iota) - I(|u_j| \leq (\beta_n - r_n\|X_j\|)\iota) \right\|. \end{aligned} \quad (29)$$

Since for all $t = 1, \dots, d$,

$$\begin{aligned} E\left(I(|u_{jt}| \leq (\beta_n + r_n\|X_j\|)) - I(|u_{jt}| \leq (\beta_n - r_n\|X_j\|)) \mid X_j\right) &= \mathcal{F}_{u_{jt}|X_j}(\beta_n + r_n\|X_j\|) - \mathcal{F}_{u_{jt}|X_j}(\beta_n - r_n\|X_j\|) \\ &= f_{u_{jt}|X_j}(\cdot) \|X_j\| r_n \leq C_f r_n \|X_j\|, \end{aligned} \quad (30)$$

the unconditional expectation of (29) is bounded by

$$d C_w C_f r_n (k_n \beta_n)^{-1} \sum_{j=1}^n E \|X_j\|^2 = O(nr_n (k_n \beta_n)^{-1}) = O(n^{1/2} k_n^{-1} \beta_n^{-1}). \quad \blacksquare$$

L13 $\max_i \|\hat{A}_i - \bar{A}_i\| = o_p(1)$.

Proof: Using **L2**, **L6**, **L8**, **L9**, **L10**, **L11** and **L12** in (20) yields

$$\hat{A}_i - \bar{A}_i = O_p((n^{1/p_x} \zeta_{npT} + \zeta_{npS})(1 + \zeta_{npT}) = o_p(1),$$

by **A7**. \blacksquare

B.2 $\sqrt{n}(\hat{m}_n(\theta_0) - m_n(\theta_0))$

Observe that

$$\sqrt{n}(\hat{m}_n(\theta_0) - m_n(\theta_0)) = n^{-1/2} \sum_{i=1}^n (\hat{A}_i - A_i) s_i = n^{-1/2} \sum_{i=1}^n (\hat{A}_i - \bar{A}_i) s_i + n^{-1/2} \sum_{i=1}^n (\bar{A}_i - A_i) s_i. \quad (31)$$

We use the expansion in (20) to deal with the first RHS term and show the following results.

$$n^{-1/2} \sum_{i=1}^n \bar{A}_i(\hat{T}_i - \tilde{T}_i) \bar{T}_i^{-1} s_i = o_p(1), \quad (32)$$

$$n^{-1/2} \sum_{i=1}^n \bar{A}_i(\tilde{T}_i - \bar{T}_i) \bar{T}_i^{-1} s_i = o_p(1), \quad (33)$$

$$n^{-1/2} \sum_{i=1}^n (\hat{S}_i - \tilde{S}_i)' \bar{T}_i^{-1} s_i = o_p(1), \quad (34)$$

$$n^{-1/2} \sum_{i=1}^n (\tilde{S}_i - \bar{S}_i)' \bar{T}_i^{-1} s_i = o_p(1), \quad (35)$$

$$n^{-1/2} \sum_{i=1}^n \bar{A}_i(\hat{T}_i - \bar{T}_i)(\hat{T}_i^{-1} - \bar{T}_i^{-1}) = o_p(1), \quad (36)$$

$$n^{-1/2} \sum_{i=1}^n (\hat{S}_i - \bar{S}_i)'(\hat{T}_i^{-1} - \bar{T}_i^{-1}) = o_p(1), \quad (37)$$

$$n^{-1/2} \sum_{i=1}^n (\bar{A}_i - A_i) s_i = o_p(1). \quad (38)$$

B.2.1 (32)

L14 Let $\{\xi_i\}$ be an i.i.d. random sequence for which $E(\xi_i|X) = 0$ a.s. and $\text{ess sup}(|\xi_i|) \leq 1$. Then

$$\max_j \left\| \sum_{i=1}^n w_{ij} \xi_i \right\| = o_p(\sqrt{n \log n / k_n}).$$

Proof: Let $\epsilon_n = C_w \sqrt{3n \log n / k_n}$. Then

$$P\left(\max_j \left\| \sum_i w_{ij} \xi_i \right\| \geq 2\epsilon_n\right) \leq P\left(\max_j \left\| \sum_{i \neq j} w_{ij} \xi_i \right\| \geq \epsilon_n\right) + P\left(\max_j \|w_{jj} \xi_j\| \geq \epsilon_n\right). \quad (39)$$

The second RHS term in (39) is bounded by $I(C_w/k_n \geq \epsilon_n)$, which equals zero for sufficiently large n . We now deal with the first RHS term in (39). By the *Hoeffding inequality*,¹² noting that $\|w_{ij} \xi_i\| \leq C_w/k_n$ for

¹² The Hoeffding inequality says that if $\{\mu_i\}$ is an independent sequence of mean zero random variables taking values on $[a_i, b_i]$, then $P(\|\sum_i \mu_i\| > \epsilon_n) \leq \exp[-2\epsilon_n^2 / \sum_{i=1}^n (b_i - a_i)^2]$.

all i, j ,

$$\begin{aligned} P\left(\max_j \left\| \sum_{i \neq j} w_{ij} \xi_i \right\| \geq \epsilon_n | X_j, \xi_j\right) &\leq \sum_{j=1}^n P\left(\left\| \sum_{i \neq j} w_{ij} \xi_i \right\| \geq \epsilon_n | X_j, \xi_j\right) \\ &\leq \sum_{j=1}^n \exp\left(-\frac{\epsilon_n^2 k_n^2}{2n C_w^2}\right) = n \exp(- (3/2) \log n) = n^{-1/2} = o(1). \quad \blacksquare \end{aligned}$$

L15 Let $\{\xi_i\}$ be as in [L14](#) and let $\xi_{ni} = \Xi_{ni}(X)\xi_i$, where for some $p_\Xi > 0$, $\limsup E\|\Xi_{ni}(X)\|^{p_\Xi} < \infty$.

Then

$$\max_j \left\| \sum_{i=1}^n w_{ij} \xi_{ni} \right\| = o_p(n^{1/p_\Xi+1/2} k_n^{-1} \log n).$$

Proof: Let $\epsilon_n^* = n^{1/p_\Xi} \sqrt{\log n}$, $\epsilon_n = \sqrt{3} C_w n^{1/p_\Xi+1/2} \log n / k_n$ and $\xi_{ni}^* = \xi_{ni} I(\|\Xi_{ni}(X)\| \leq \epsilon_n^*) / \epsilon_n^*$. Then

$$\begin{aligned} P\left(\max_j \left\| \sum_{i=1}^n w_{ij} \xi_{ni} \right\| \geq 2\epsilon_n\right) &= P\left(\max_j \left\| \sum_{i=1}^n w_{ij} (\epsilon_n^* \xi_{ni}^* + \xi_{ni} I(\|\Xi_{ni}(X)\| > \epsilon_n^*)) \right\| \geq 2\epsilon_n\right) \\ &\leq P\left(\max_j \left\| \sum_{i=1}^n w_{ij} \xi_{ni}^* \right\| \geq 2\frac{\epsilon_n}{\epsilon_n^*}\right) + P\left(\max_i \|\Xi_{ni}(X)\| \geq \epsilon_n^*\right). \quad (40) \end{aligned}$$

The second RHS term in (40) is by [L4](#) bounded by

$$(\epsilon_n^*)^{-p_\Xi} \sum_{i=1}^n E\|\Xi_{ni}\|^{p_\Xi} = O((\log n)^{-p_\Xi/2}) = o(1).$$

The first RHS term in (40) is also $o(1)$ because $\text{ess sup} \|\xi_{ni}^*\| \leq 1$ by construction and since

$$\frac{\epsilon_n}{\epsilon_n^*} = \frac{\sqrt{3} C_w n^{\frac{p_\Xi+2}{2p_\Xi}} \log n / k_n}{n^{1/p_\Xi} \sqrt{\log n}} = \frac{C_w \sqrt{3n \log n}}{k_n},$$

[L14](#) can be applied. \blacksquare

L16 $n^{-1/2} \sum_i \bar{A}_i (\hat{T}_i - \tilde{T}_i) \bar{T}_i^{-1} s_i = o_p(1)$.

Proof: The LHS is

$$\begin{aligned} \left\| n^{-1/2} \sum_{j=1}^n \sum_{i=1}^n w_{ij} \bar{A}_i (\hat{s}_j \hat{s}'_j - s_j s'_j) \bar{T}_i^{-1} s_i \right\| &\stackrel{\text{L15}}{\leq} \sum_{j=1}^n \|\hat{s}_j \hat{s}'_j - s_j s'_j\| \times o_p(n^{1/p_x} k_n^{-1} \log n) \\ &\stackrel{(21)}{\leq} C_s \sum_{j=1}^n \|\hat{s}_j - s_j\| \times o_p(n^{1/p_x} k_n^{-1} \log n). \quad (41) \end{aligned}$$

Set $r_n = 1/(\sqrt{n} - \log n)$. Now,

$$\sum_{j=1}^n \|\hat{s}_j - s_j\| \stackrel{(22)}{\leq} \sum_{j=1}^n (\|\alpha_{jr_n}\| - E(\|\alpha_{jr_n}\| | X)) + \sum_{j=1}^n E(\|\alpha_{jr_n}\| | X) + nI(\|\hat{\theta} - \theta_0\| > r_n). \quad (42)$$

The third RHS term is $O_p(1)$ by (24) and the second RHS term is by L7 bounded by $C_\alpha r_n \sum_{j=1}^n \|X_j\| = O_p(nr_n) = O_p(n^{1/2})$. Squaring the first RHS term and taking its expectation yields

$$\sum_{j=1}^n E(\|\alpha_{jr_n}\| - E(\|\alpha_{jr_n}\| | X))^2 \stackrel{L7}{\leq} Cnr_n = O(nr_n).$$

Hence the RHS in (42) is $O_p(\sqrt{nr_n}) + O_p(\sqrt{n}) + O_p(1) = O_p(\sqrt{n})$, which implies that the RHS in (41) is $o_p(n^{1/p_x+1/2}k_n^{-1} \log n) = o_p(1)$ by A7. ■

B.2.2 (33)

L17 Let $\xi_{nij} = \xi_n(u_i, u_j; X)$ be such that $E(\xi_{nij} | u_i, X) = E(\xi_{nij} | u_j, X) = 0$ a.s. for all i, j and $\max_{i,j} E\|\xi_{nij}\|^2 = O(1)$. Then $n^{-1} \sum_{i,j=1}^n w_{ij} \xi_{nij} = O_p(k_n^{-1})$.

Proof: Square the LHS and take the expectation to obtain

$$n^{-2} \sum_{i,j=1}^n \left(E(w_{ij}^2 \|\xi_{nij}\|^2) + E(w_{ij} w_{ji} \xi_{nij}' \xi_{nji}) \right) \stackrel{A6}{\leq} 2C_w^2 k_n^{-2} \max_{i,j} E\|\xi_{nij}\|^2 = O(k_n^{-2}). \quad \blacksquare$$

L18 $n^{-1/2} \sum_{i=1}^n \bar{A}_i (\tilde{T}_i - \bar{T}_i) \bar{T}_i^{-1} s_i = o_p(1)$.

Proof: In L17, take $\xi_{nij} = \bar{A}_i (s_j s_j' - T_j) \bar{T}_i^{-1} s_i$ to obtain a convergence rate of $O_p(n^{1/2} k_n^{-1}) = o_p(1)$. ■

B.2.3 (34) and (35)

L19 $n^{-1/2} \sum_i (\hat{S}_i - \tilde{S}_i)' \bar{T}_i^{-1} s_i = o_p(1)$.

Proof: The norm of the LHS is

$$\begin{aligned} \left\| n^{-1/2} \sum_{j=1}^n (\hat{F}_j - \tilde{F}_j) X_j' \sum_{i=1}^n w_{ij} \bar{T}_i^{-1} s_i \right\| &\leq \max_j \left\| n^{-1/2} \sum_{i=1}^n w_{ij} \bar{T}_i^{-1} s_i \right\| \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| \\ &\stackrel{L14}{=} O_p(k_n^{-1} \sqrt{\log n}) \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\|. \end{aligned}$$

Let (as in L12) $r_n = 1/(\sqrt{n} - \log n)$. Then

$$\begin{aligned} \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| &= \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| I(\|\hat{\theta} - \theta_0\| \leq r_n) + \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| I(\|\hat{\theta} - \theta_0\| > r_n) \\ &\stackrel{(24)}{=} \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| I(\|\hat{\theta} - \theta_0\| \leq r_n) + \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| \times O_p(n^{-1}). \end{aligned}$$

Finally,

$$\frac{\sqrt{\log n}}{k_n} \sum_{j=1}^n \|\hat{F}_j - \tilde{F}_j\| \times \|X_j\| I(\|\hat{\theta} - \theta_0\| \leq r_n) \stackrel{(28),(30)}{\leq} \frac{C_f d r_n \sqrt{\log n}}{k_n \beta_n} \sum_{j=1}^n \|X_j\|^2 = O_p\left(\frac{\sqrt{n \log n}}{k_n \beta_n}\right) = o_p(1),$$

by A7. ■

L20 $n^{-1/2} \sum_{i=1}^n (\tilde{S}_i - \bar{S}_i)' \bar{T}_i^{-1} s_i = o_p(1)$.

Proof: The LHS is

$$n^{-1/2} \sum_{i,j=1}^n w_{ij} (\tilde{F}_j - E(\tilde{F}_j | X_j)) X_j' \bar{T}_i^{-1} s_i + n^{-1/2} \sum_{i,j=1}^n w_{ij} (E(\tilde{F}_j | X_j) - F_j) X_j' \bar{T}_i^{-1} s_i. \quad (43)$$

The first RHS term is $O_p(n^{1/2} \beta_n^{-1/2} k_n^{-1}) = o_p(1)$ by L17. The norm of the second RHS term is bounded by

$$\begin{aligned} n^{-1/2} \max_j \left\| \sum_{i=1}^n w_{ij} \bar{T}_i^{-1} s_i \right\| \left\| \sum_{j=1}^n w_{ij} \|E(\tilde{F}_j | X_j) - F_j\| \right\| \times \|X_j\| \\ \leq \stackrel{\text{L14,(26)}}{=} O_p(k_n^{-1} \sqrt{\log n}) 6^{-1} C_f \beta_n^2 \sum_j \|X_j\| = O_p(n k_n^{-1} \beta_n^2 \sqrt{\log n}) = o_p(1), \end{aligned}$$

by A7.

B.2.4 (36) and (37)

L21 $n^{-1/2} \sum_{i=1}^n \bar{A}_i(\hat{T}_i - \bar{T}_i)(\hat{T}_i^{-1} - \bar{T}_i^{-1})s_i = o_p(1)$.

Proof: Note that

$$\begin{aligned} & \left\| n^{-1/2} \sum_{i=1}^n \bar{A}_i(\hat{T}_i - \bar{T}_i)(\hat{T}_i^{-1} - \bar{T}_i^{-1})s_i \right\| \\ & \leq \max_i \|\hat{T}_i - \bar{T}_i\| \times \|\hat{T}_i^{-1} - \bar{T}_i^{-1}\| \times n^{-1/2} \sum_{i=1}^n \|\bar{A}_i\| \times \|s_i\| = O_p(\sqrt{n}\zeta_{npT}^2) = o_p(1), \end{aligned}$$

by L8, L9 and A7. ■

L22 $n^{-1/2} \sum_{i=1}^n (\hat{S}_i - \bar{S}_i)'(\hat{T}_i^{-1} - \bar{T}_i^{-1})s_i = o_p(1)$.

Proof: Use a similar inequality to the one used in L21 to obtain a rate of $n^{1/2}\zeta_{npS}\zeta_{npT} = o(1)$ by A7. ■

B.2.5 (38)

L23 $E\|\bar{A}_i - A_i\|^2 = o(1)$.

Proof: The square of the LHS is bounded by

$$C\left(E\|A_i\|^4 E\|\bar{T}_i - T_i\|^4 + (E\|\bar{S}_i - S_i\|^2)^2\right) = o(1),$$

by theorem 1 of Stone (1977). ■

L24 $n^{-1/2} \sum_{i=1}^n (\bar{A}_i - A_i)s_i = o_p(1)$.

Proof:

$$E\left\| n^{-1/2} \sum_{i=1}^n (\bar{A}_i - A_i)s_i \right\|^2 \leq E\|\bar{A}_i - A_i\|^2 = o(1),$$

by L23. ■

L25 $\hat{m}_n(\theta_0) - m_n(\theta_0) = o_p(n^{-1/2})$.

Proof: Using the expansion in (31) and (32)–(38), the stated result follows from lemmas L16, L18, L19, L20, L21, L22, and L24. ■

C Feasible Estimator

L26 *There exists a positive sequence $\{\mu_{1n}\}$ with $\mu_{1n} = o(1)$ such that for any positive sequence $\{r_n\}$, $n^{-1} \sum_{i=1}^n \|\bar{A}_i - A_i\| \|\alpha_{ir_n}\| = o_p(r_n \mu_{1n})$.*

Proof: Let μ_{1n} be such that $\mu_{1n} = o(1)$ and $E\|\bar{A}_i - A_i\|^2 = o(\mu_{1n}^2)$; such μ_{1n} exist by lemma L23. Now,

$$E(\|\bar{A}_i - A_i\| \|\alpha_{ir_n}\|) \stackrel{\text{L7}}{\leq} Cr_n E(\|\bar{A}_i - A_i\| \|X_i\|) \stackrel{\text{Schwarz}}{\leq} Cr_n \sqrt{E(\|\bar{A}_i - A_i\|^2)} \sqrt{E\|X_i\|^2} = o(r_n \mu_{1n}). \quad \blacksquare$$

Let $\Theta_r = \{\theta \in \Theta : \|\theta - \theta_0\| < r\}$.

L27 *There exists a positive sequence $\{\mu_n\}$ with $\mu_n = o(1)$ such that for any positive sequence $\{r_n\}$,*

$$\sup_{\theta \in \Theta_{r_n}} \|\hat{m}_n(\theta) - m_n(\theta)\| = o_p(r_n \mu_n + n^{-1/2}).$$

Proof: First note that

$$\begin{aligned} \sup_{\theta \in \Theta_{r_n}} \|\hat{m}_n(\theta) - m_n(\theta)\| &\stackrel{\text{triangle}}{\leq} \sup_{\theta \in \Theta_{r_n}} \|\hat{m}_n(\theta) - m_n(\theta) - \hat{m}_n(\theta_0) + m_n(\theta_0)\| + \|\hat{m}_n(\theta_0) - m_n(\theta_0)\| \\ &\stackrel{\text{L25}}{\leq} \sup_{\theta \in \Theta_{r_n}} n^{-1} \sum_{i=1}^n \|\hat{A}_i - A_i\| \|s_i(\theta) - s_i(\theta_0)\| + o_p(n^{-1/2}) \\ &\leq n^{-1} \sum_{i=1}^n \|\hat{A}_i - A_i\| \|\alpha_{ir_n}\| + o_p(n^{-1/2}). \end{aligned}$$

Now, let μ_{2n} be such that $\max_i \|\hat{A}_i - A_i\| = o_p(\mu_{2n})$ and $\mu_{2n} = o(1)$; such μ_{2n} exist by L13. Then by the triangle inequality,

$$\begin{aligned} n^{-1} \sum_{i=1}^n \|\hat{A}_i - A_i\| \|\alpha_{ir_n}\| &\leq n^{-1} \sum_{i=1}^n \|\hat{A}_i - \bar{A}_i\| \|\alpha_{ir_n}\| + n^{-1} \sum_{i=1}^n \|\bar{A}_i - A_i\| \|\alpha_{ir_n}\| \\ &\leq \max_i \|\hat{A}_i - \bar{A}_i\| n^{-1} \sum_{i=1}^n \|\alpha_{ir_n}\| + n^{-1} \sum_{i=1}^n \|\bar{A}_i - A_i\| \|\alpha_{ir_n}\| \stackrel{\text{L7, L13, L26}}{=} o_p(\mu_{2n}) O_p(r_n) + o_p(\mu_{1n} r_n) = o_p((\mu_{1n} + \mu_{2n}) r_n), \end{aligned}$$

Take $\mu_n = \mu_{1n} + \mu_{2n}$. \blacksquare

L28 $m_n(\hat{\theta}) = o_p(n^{-1/2})$.

Proof: Let $\{\psi_n\}$ be such that $\|\hat{\theta} - \theta_0\| = O_p(\psi_n)$ but $\|\hat{\theta} - \theta_0\| \neq o_p(\psi_n)$. Let μ_n be as in L27. Then for $r_n = \psi_n/\sqrt{\mu_n}$ we have

$$\|m_n(\hat{\theta})\| \stackrel{\text{triangle}}{\leq} \|m_n(\hat{\theta}) - \hat{m}_n(\hat{\theta})\| + \|\hat{m}_n(\hat{\theta})\| \lesssim \sup_{\theta \in \Theta_{r_n}} \|m_n(\theta) - \hat{m}_n(\theta)\| + o_p(n^{-1/2}) \stackrel{\text{L27}}{=} o_p(\psi_n \sqrt{\mu_n}) + o_p(n^{-1/2}).$$

So by L1, $\|\hat{\theta} - \theta_0\| = o_p(\psi_n) + O_p(n^{-1/2})$. Hence $\psi_n \sim n^{-1/2}$. Apply L27 with $r_n = n^{-1/2}$. ■

Proof of T3: By L28, $\hat{\theta}$ satisfies (8). ■

D Covariance Matrix Estimation

Let $\bar{V} = n^{-1} \sum_{i=1}^n \bar{A}_i \bar{S}_i$.

L29 $\hat{V} - \bar{V} = o_p(1)$.

Proof: Using the expansion

$$\hat{V} - \bar{V} = n^{-1} \sum_{i=1}^n (\hat{A}_i - \bar{A}_i)(\hat{S}_i - \bar{S}_i) + n^{-1} \sum_{i=1}^n (\hat{A}_i - \bar{A}_i)\bar{S}_i + n^{-1} \sum_{i=1}^n \bar{A}_i(\hat{S}_i - \bar{S}_i),$$

the stated result follows from L11, L12 and L13. ■

L30 $\bar{V} - V = o_p(1)$.

Proof: Using a similar expansion to the one in L29, we have

$$\begin{aligned} E\|\bar{V} - V\| &= E\left\|n^{-1} \sum_{i=1}^n (\bar{A}_i \bar{S}_i - A_i S_i)\right\| \\ &\leq E\left(\|\bar{A}_i - A_i\| \times \|\bar{S}_i - S_i\|\right) + E\left(\|A_i\| \times \|\bar{S}_i - S_i\|\right) + E\left(\|\bar{A}_i - A_i\| \times \|S_i\|\right) \\ &\stackrel{\text{Schwarz}}{\leq} \sqrt{E\|\bar{A}_i - A_i\|^2} \sqrt{E\|\bar{S}_i - S_i\|^2} + \sqrt{E\|A_i\|^2} \sqrt{E\|\bar{S}_i - S_i\|^2} + \sqrt{E\|\bar{A}_i - A_i\|^2} \sqrt{E\|S_i\|^2}. \end{aligned}$$

Apply L23, theorem 1 of Stone (1977) and the fact that $E\|A_i\|^2, E\|S_i\|^2 < \infty$ by assumption. ■

Proof of T4: Combine the previous two lemmas. ■

E Computation

Proof of T5: By L27 and T4 it follows that $\hat{\theta}_{(1)} = O_p(n^{-1/2})$. Hence by L1, $\hat{m}_n(\hat{\theta}_{(j)}) - m_n(\hat{\theta}_{(j)}) = o_p(n^{-1/2})$ for $j = 0, 1$. Because $\{A_i s_i\}$ is a VC class (see (18)), it follows that

$$\left| m_n(\hat{\theta}_{(1)}) - m_n(\hat{\theta}_{(0)}) - m(\hat{\theta}_{(1)}) + m(\hat{\theta}_{(0)}) \right| = o_p(n^{-1/2}).$$

Since $m(\hat{\theta}_{(1)}) - m(\hat{\theta}_{(0)}) = V(\hat{\theta}_{(1)} - \hat{\theta}_{(0)}) + o_p(n^{-1/2})$ (see (19)), it follows that

$$\begin{aligned} \hat{m}_n(\hat{\theta}_{(1)}) - \hat{m}_n(\hat{\theta}_{(0)}) &= m_n(\hat{\theta}_{(1)}) - m_n(\hat{\theta}_{(0)}) + o_p(n^{-1/2}) = m(\hat{\theta}_{(1)}) - m(\hat{\theta}_{(0)}) + o_p(n^{-1/2}) \\ &= V(\hat{\theta}_{(1)} - \hat{\theta}_{(0)}) + o_p(n^{-1/2}) = -V\hat{V}^{-1}(\hat{\theta}_{(0)})\hat{m}_n(\hat{\theta}_{(0)}) + o_p(n^{-1/2}) \stackrel{\text{T4}}{=} -\hat{m}_n(\hat{\theta}_{(0)}) + o_p(n^{-1/2}). \end{aligned}$$

So $\hat{m}_n(\hat{\theta}_{(1)}) = o_p(n^{-1/2})$ and (9) is satisfied. ■

References Cited

- Aitken, Alexander C. (1935) "On least squares and linear combination of observations," *Proceedings of the Royal Society of Edinburgh* 55, 42–48.
- Carroll, Raymond J. (1982) "Adapting for heteroscedasticity in linear models," *Annals of Statistics* 10, 1224–1233.
- Chernozhukov, Victor and Christian Hansen (2006) "Instrumental quantile regression inference for structural and treatment effect models," *Journal of Econometrics* 132, 491–525.
- Delgado, Miguel (1992) "Semiparametric generalized least squares estimation in the multivariate nonlinear regression model," *Econometric Theory* 8, 203–222.
- Koenker, Roger (2005) "Quantile regression," Cambridge University Press (New York).
- Koenker, Roger and Quanshui Zhao (1994) "L-estimation for linear heteroscedastic models," *Journal of Nonparametric Statistics* 3, 223–235.
- Komunjer, Ivana and Quang Vuong (2006) "Efficient Conditional Quantile Estimation: The Time Series Case," UCSD working paper.
- Newey, Whitney K. (1990) "Efficient instrumental variables estimation of nonlinear models," *Econometrica* 58, 809–837.
- Newey, Whitney K. (1993) "Efficient estimation of models with conditional moment restrictions," in *Handbook of Statistics* 11, G.S. Maddala, C.R. Rao and H.D. Vinod, eds., North Holland, Amsterdam.
- Newey, Whitney K. and James L. Powell (1990) "Efficient estimation of linear and type I censored regression models under conditional quantile restrictions," *Econometric Theory* 6, 295–317.
- Pakes, Ariel and David Pollard (1989) "Simulation and the asymptotics of optimization estimators," *Econometrica* 57, 1027–1057.
- Pinkse, Joris (2006) "Heteroskedasticity correction and dimension reduction," Pennsylvania State University working paper.
- Robinson, Peter M. (1987) "Asymptotically efficient estimation in the presence of heteroskedasticity of unknown form," *Econometrica* 55, 875–891.

- Stone, Charles J. (1977) “Consistent nonparametric regression,” *Annals of Statistics* 5, 595–645.
- van der Vaart, Aad, and Jon A. Wellner (1996) “Weak convergence and empirical processes: with applications to statistics,” Springer (New York).
- Wang, Yoonjae (2006) “Smoothed Empirical Likelihood Methods for Quantile Regression Models,” *Econometric Theory* 22, 173–205.
- Zhao, Quanshui (2001) “Asymptotically efficient median regression in the presence of heteroskedasticity of unknown form,” *Econometric Theory* 17, 765–784.

		No Overlap		Overlap in θ_{01}, θ_{02}	
		$x_{i1} = x_{i2}$	$x_{i1} \neq x_{i2}$	$x_{i1} = x_{i2}$	$x_{i1} \neq x_{i2}$
$x_{i1}, x_{i2} \perp\!\!\!\perp u_{i1}, u_{i2}$	$u_{i1} \perp\!\!\!\perp u_{i2}$	all same	all same	$J \succcurlyeq S^*SO$	$J \succcurlyeq S^*SO$
	$u_{i1} \not\perp\!\!\!\perp u_{i2}$	all same	$J \succcurlyeq S^*SO$	$J \succcurlyeq S^*SO$	$J \succcurlyeq S^*SO$
$x_{ij} \perp\!\!\!\perp u_{ij^*}; j \neq j^*$	$u_{i1} \perp\!\!\!\perp u_{i2} x_{i1}, x_{i2}$	all same	$JS^*S \succcurlyeq O$	$J \succcurlyeq S^*SO$	$J \succcurlyeq S^*S \succcurlyeq O$
	$u_{i1} \not\perp\!\!\!\perp u_{i2} x_{i1}, x_{i2}$	all same	$J \succcurlyeq S^*S \succcurlyeq O$	$J \succcurlyeq S^*SO$	$J \succcurlyeq S^*S \succcurlyeq O$
$x_{ij} \not\perp\!\!\!\perp u_{ij^*}$	$u_{i1} \perp\!\!\!\perp u_{i2} x_{i1}, x_{i2}$	$JS^*S \succcurlyeq O$	$JS^* \succcurlyeq S \succcurlyeq O$	$J \succcurlyeq S^*S \succcurlyeq O$	$J \succcurlyeq S^* \succcurlyeq S \succcurlyeq O$
	$u_{i1} \not\perp\!\!\!\perp u_{i2} x_{i1}, x_{i2}$	$J \succcurlyeq S^*S \succcurlyeq O$	$J \succcurlyeq S^* \succcurlyeq S \succcurlyeq O$	$J \succcurlyeq S^*S \succcurlyeq O$	$J \succcurlyeq S^* \succcurlyeq S \succcurlyeq O$

The entries indicate which methods are preferable to others in terms of asymptotic efficiency in various situations. ‘ $\perp\!\!\!\perp$ ’ denotes independence and ‘ \succcurlyeq ’ means “is typically more efficient but never less efficient than.” **J**=joint estimation (new methodology), **S**=separate estimation (Zhao’s method), **S***=separate estimation using the regressors from both equations (Zhao’s results can be used for this) and **O**=no efficiency correction.

Please note: when errors are independent of each other and of the regressors *and* the coefficient vectors do not overlap, then equation by equation adaptive (to error distribution) estimation dominates all of the other estimation methods mentioned here.

This comparison applies equally to mean and quantile regressions.

Table 1: Asymptotic Efficiency Comparison of Semiparametric Methods

	$h(x_{i1}, x_{i2})$	$g(x_{i1}, x_{i2})$
Homo	1	1
Hetero I	$\exp(0.1 x'_{i1}\beta_0)$	$1 + 3 \exp(-(x'_{i1}\beta_0 + 5)^2/100)$
Hetero II	$\exp(0.1 x'_{i1}\beta_0 + x'_{i2}\gamma_0)$	$1 + 3 \exp(-(x'_{i1}\beta_0 + x'_{i2}\gamma_0 + 10)^2/100)$

Table 2: Heteroskedasticity Designs

Methods		β_{00}	β_{01}	β_{02}	γ_{00}	γ_{01}	γ_{02}
Joint Estimation	$n = 100$	0.1189	0.0027	0.0131	0.1093	0.0026	0.0144
	$n = 300$	0.0380	0.0008	0.0047	0.0404	0.0009	0.0046
	$n = 500$	0.0225	0.0004	0.0026	0.0207	0.0005	0.0027
Single Equation	$n = 100$	0.1137	0.0027	0.0130	0.1074	0.0026	0.0139
	$n = 300$	0.0373	0.0008	0.0046	0.0394	0.0009	0.0046
	$n = 500$	0.0225	0.0004	0.0026	0.0206	0.0005	0.0027
Unweighted LAD	$n = 100$	0.1123	0.0026	0.0125	0.1056	0.0025	0.0137
	$n = 300$	0.0373	0.0008	0.0046	0.0391	0.0009	0.0045
	$n = 500$	0.0222	0.0004	0.0025	0.0203	0.0005	0.0026

Table 3: MSE by Monte Carlo: Homo with $\rho = 0$

Methods		β_{00}	β_{01}	β_{02}	γ_{00}	γ_{01}	γ_{02}
Joint Estimation	$n = 100$	0.0929	0.0020	0.0100	0.0880	0.0021	0.0111
	$n = 300$	0.0285	0.0007	0.0033	0.0327	0.0007	0.0035
	$n = 500$	0.0162	0.0004	0.0020	0.0169	0.0004	0.0021
Single Equation	$n = 100$	0.1109	0.0025	0.0128	0.1125	0.0027	0.0149
	$n = 300$	0.0360	0.0009	0.0044	0.0404	0.0009	0.0045
	$n = 500$	0.0204	0.0005	0.0027	0.0209	0.0005	0.0027
Unweighted LAD	$n = 100$	0.1113	0.0023	0.0126	0.1115	0.0026	0.0145
	$n = 300$	0.0355	0.0009	0.0044	0.0400	0.0009	0.0045
	$n = 500$	0.0200	0.0005	0.0027	0.0205	0.0005	0.0027

Table 4: MSE by Monte Carlo: Homo with $\rho = 0.7$

Methods		β_{00}	β_{01}	β_{02}	γ_{00}	γ_{01}	γ_{02}
Joint Estimation	$n = 100$	0.7641	0.0475	0.0721	0.4289	0.0060	0.0699
	$n = 300$	0.2188	0.0140	0.0213	0.1369	0.0018	0.0179
	$n = 500$	0.1368	0.0086	0.0130	0.0822	0.0010	0.0115
Single Equation	$n = 100$	0.7155	0.0443	0.0669	0.3987	0.0057	0.0655
	$n = 300$	0.2087	0.0131	0.0201	0.1312	0.0017	0.0172
	$n = 500$	0.1295	0.0080	0.0122	0.0780	0.0010	0.0109
Unweighted LAD	$n = 100$	0.7959	0.0521	0.0750	0.4121	0.0059	0.0682
	$n = 300$	0.2323	0.0156	0.0223	0.1382	0.0017	0.0183
	$n = 500$	0.1474	0.0099	0.0142	0.0866	0.0010	0.0120

Table 5: MSE by Monte Carlo: No Cross-Equation Hetero (Hetero I) with $\rho = 0$

Methods		β_{00}	β_{01}	β_{02}	γ_{00}	γ_{01}	γ_{02}
Joint Estimation	$n = 100$	0.5316	0.0366	0.0582	0.3675	0.0044	0.0500
	$n = 300$	0.1749	0.0117	0.0180	0.1075	0.0012	0.0148
	$n = 500$	0.0971	0.0069	0.0095	0.0653	0.0008	0.0088
Single Equation	$n = 100$	0.6787	0.0418	0.0652	0.4298	0.0051	0.0616
	$n = 300$	0.2136	0.0133	0.0210	0.1312	0.0015	0.0183
	$n = 500$	0.1245	0.0079	0.0115	0.0800	0.0010	0.0111
Unweighted LAD	$n = 100$	0.7453	0.0479	0.0722	0.4450	0.0053	0.0634
	$n = 300$	0.2325	0.0160	0.0242	0.1408	0.0016	0.0199
	$n = 500$	0.1408	0.0097	0.0132	0.0847	0.0011	0.0120

Table 6: MSE by Monte Carlo: No Cross-Equation Hetero (Hetero I) with $\rho = 0.7$

Methods		β_{00}	β_{01}	β_{02}	γ_{00}	γ_{01}	γ_{02}
Joint Estimation	$n = 100$	1.0014	0.0456	0.1249	0.3149	0.0055	0.0428
	$n = 300$	0.2792	0.0142	0.0360	0.0955	0.0016	0.0125
	$n = 500$	0.1626	0.0084	0.0233	0.0600	0.0009	0.0076
Single Equation	$n = 100$	1.0344	0.0456	0.1266	0.3106	0.0054	0.0413
	$n = 300$	0.2784	0.0147	0.0371	0.0964	0.0016	0.0125
	$n = 500$	0.1693	0.0087	0.0239	0.0610	0.0009	0.0077
Unweighted LAD	$n = 100$	1.1002	0.0496	0.1326	0.3159	0.0054	0.0418
	$n = 300$	0.2899	0.0163	0.0390	0.0977	0.0016	0.0126
	$n = 500$	0.1805	0.0097	0.0255	0.0607	0.0009	0.0078

Table 7: MSE by Monte Carlo: Cross-Equation Hetero (Hetero II) with $\rho = 0$

Methods		β_{00}	β_{01}	β_{02}	γ_{00}	γ_{01}	γ_{02}
Joint Estimation	$n = 100$	0.8263	0.0397	0.0974	0.2943	0.0048	0.0340
	$n = 300$	0.2260	0.0126	0.0288	0.0845	0.0014	0.0104
	$n = 500$	0.1284	0.0062	0.0166	0.0522	0.0009	0.0061
Single Equation	$n = 100$	0.9934	0.0475	0.1176	0.3046	0.0054	0.0415
	$n = 300$	0.3060	0.0157	0.0376	0.0989	0.0017	0.0141
	$n = 500$	0.1645	0.0081	0.0222	0.0620	0.0011	0.0085
Unweighted LAD	$n = 100$	1.0346	0.0516	0.1220	0.3021	0.0055	0.0415
	$n = 300$	0.3210	0.0171	0.0400	0.0988	0.0017	0.0142
	$n = 500$	0.1758	0.0091	0.0237	0.0635	0.0011	0.0086

Table 8: MSE by Monte Carlo: Cross-Equation Hetero (Hetero II) with $\rho = 0.7$