

Optimal Inference for Instrumental Variables Regression with non-Gaussian Errors*

MATIAS D. CATTANEO

DEPARTMENT OF ECONOMICS, UNIVERSITY OF MICHIGAN

RICHARD K. CRUMP

DEPARTMENT OF ECONOMICS, UC BERKELEY

MICHAEL JANSSON

DEPARTMENT OF ECONOMICS, UC BERKELEY AND *CREATES*

February 2, 2009

ABSTRACT. This paper is concerned with inference on the coefficient on the endogenous regressor in a linear instrumental variables model with a single endogenous regressor, nonrandom exogenous regressors and instruments, and *i.i.d.* errors whose distribution is unknown. It is shown that under mild smoothness conditions on the error distribution it is possible to develop tests which are “nearly” efficient (in the sense of Andrews, Moreira, and Stock (2006)) when identification is weak and consistent and asymptotically optimal when identification is strong. In addition, an estimator is presented which can be used in the usual way to construct valid (indeed, optimal) confidence intervals when identification is strong. The estimator is of the two stage least squares variety and is asymptotically efficient under strong identification whether or not the errors are normal.

1. INTRODUCTION

This paper is concerned with inference on the coefficient on the endogenous regressor in a linear instrumental variables (IVs) model with a single endogenous regressor, nonrandom exogenous regressors and IVs, and *i.i.d.* errors. Models of this type have been studied intensively in recent years, with particular attention being devoted to the case where the IVs are weak (in the terminology of Staiger and Stock (1997)).¹ Analyzing such a model in which the *i.i.d.* errors are furthermore assumed to be

*The authors thank Don Andrews, Bryan Graham, Jim Powell, Tom Rothenberg, Paul Ruud, Hal White, two anonymous referees, and seminar participants at Aarhus, Berkeley, Brown, Columbia, Michigan, NYU, Princeton, Rochester, Stanford, Yale, UCSD, UTDT, and UdeSA for comments. The third author gratefully acknowledges the research support of *CREATES* (funded by the Danish National Research Foundation).

¹For reviews of the weak IV literature, see e.g. Dufour (2003) and Andrews and Stock (2007).

Gaussian, Andrews, Moreira, and Stock (2006, henceforth AMS) find that the conditional likelihood ratio test proposed by Moreira (2003) is “nearly” efficient when identification is weak and asymptotically efficient when identification is strong.

The purpose of the present paper is to explore the consequences of relaxing the assumption of normality on the part of the *i.i.d.* errors in a model which is otherwise identical to the model studied by AMS (and others). Recent work by Andrews and Marmar (2007) and Andrews and Soares (2007) shows that departures from normality can be exploited for power purposes when the errors satisfy a certain symmetry condition. Although these papers do not establish optimality results on the part of the rank-based testing procedures proposed therein, the findings of the papers imply in particular that for certain classes of error distributions the conditional likelihood test ceases to be (“nearly”) optimal once the assumption of normality is relaxed. This paper addresses the issue of optimality and shows that under mild smoothness conditions on the (otherwise unknown) error distribution it is possible to develop tests which are (“nearly”) optimal whether or not the errors are Gaussian.

The asymptotic optimality theory developed herein treats the distribution of the *i.i.d.* errors as an unknown nuisance parameter and is therefore of the semiparametric variety. In fact, under the assumption that the model contains an intercept (an assumption which we maintain throughout), we establish adaptation results, namely that one can construct procedures which perform asymptotically as well as procedures which (optimally) utilize knowledge of the error distribution. This adaptation result bears more than a superficial resemblance to Bickel’s (1982) celebrated result on adaptive estimation of the slope coefficients in a regression model. Specifically, it turns out that the problem of conducting inference in an IV model with an unknown error distribution can be decomposed into two separate problems, each of which is well understood (in isolation) from the works of Bickel (1982) and AMS, respectively. The first of these problems concerns efficient estimation of the slope coefficients in the reduced form of the IV model. That problem is a bivariate version of the problem addressed by Bickel (1982) and can be solved in essentially the same way. Because efficient estimators of the slope coefficients turn out to be asymptotically sufficient statistics for the relevant parameters of the IV model, the problem of conducting optimal inference can be reduced to the problem of optimally extracting information from the efficient estimators of the reduced form regression coefficients. The mathematical structure of that problem turns out to be the same whether or not the errors are Gaussian, implying that we can utilize the results of AMS to construct test statistics which combine the efficient estimators of the reduced form regression coefficients in a “nearly” optimal way.

Our construction of feasible inference procedures proceeds in several steps, culminating with a procedure which is “nearly” efficient when identification is weak and consistent and asymptotically optimal when identification is strong. The result-

ing procedure is of the conditional likelihood ratio variety, but being optimal (or “nearly” so, depending on the strength of identification) it is of necessity different from Moreira’s (2003) procedure. Analogously to Moreira’s (2003) procedure, a potential drawback of our procedure is that although it enjoys optimality properties when identification is strong, it is somewhat tedious to invert it in order to obtain confidence intervals in strongly identified models. To address this issue, we present an estimator (and an accompanying standard error formula) which can be used in the usual way to construct valid (indeed, optimal) confidence intervals when identification is strong. The estimator, which would appear to be new, is of the two stage least squares (2SLS) variety and is asymptotically efficient (under strong identification) whether or not the errors are normal.

The paper proceeds as follows. Section 2 presents the model and the assumptions under which the asymptotic analysis will proceed. Section 3 is concerned with asymptotic inference under the assumptions that the error distribution is known and identification is weak. The counterfactual assumption that the error distribution is known is dispensed with in Section 4, where it is also shown how strong identification can be accommodated. Section 5 presents some simulation results, while mathematical derivations have been relegated to an Appendix.

2. THE MODEL

We consider a model given by

$$\begin{aligned} y_{1i} &= \Gamma_1' x_i + \beta y_{2i} + u_i, \\ y_{2i} &= \gamma_2' x_i + \pi' z_i + v_{2i} \quad (i = 1, \dots, n), \end{aligned} \quad (1)$$

where $y_{1i}, y_{2i} \in \mathbb{R}$, $x_i \in \mathbb{R}^p$, and $z_i \in \mathbb{R}^q$ are observed variables; $u_i, v_{2i} \in \mathbb{R}$ are unobserved errors; and $\beta \in \mathbb{R}$, $\pi \in \mathbb{R}^q$, and $\Gamma_1, \gamma_2 \in \mathbb{R}^p$ are parameters. The exogenous variables x_i and z_i are fixed (i.e., nonrandom) and the first element of x_i is assumed to equal unity. The errors (u_i, v_{2i}) are *i.i.d.* from a continuous distribution with zero mean and finite variance.

It turns out to be convenient to work with the reduced form of the model. The reduced form is given by the pair of equations

$$\begin{aligned} y_{1i} &= \gamma_1' x_i + \beta \pi' z_i + v_{1i}, \\ y_{2i} &= \gamma_2' x_i + \pi' z_i + v_{2i} \quad (i = 1, \dots, n), \end{aligned} \quad (2)$$

where $\gamma_1 = \Gamma_1 + \gamma_2 \beta$ and $v_{1i} = v_{2i} \beta + u_i$. The parameters of the reduced form are β ,

π , $\gamma = (\gamma'_1, \gamma'_2)'$, and f , the Lebesgue density of $v_i = (v_{1i}, v_{2i})'$. The analysis of the reduced form is facilitated by the fact that it can be embedded in the model

$$\begin{aligned} y_{1i} &= \gamma'_1 x_i + \delta'_1 z_i + v_{1i}, \\ y_{2i} &= \gamma'_2 x_i + \delta'_2 z_i + v_{2i} \quad (i = 1, \dots, n), \end{aligned} \quad (3)$$

where $\delta_1, \delta_2 \in \mathbb{R}^q$ and the other parameters are as in (2). (The model (3) reduces to (2) when $\delta = (\delta_1, \delta_2)' = (\beta\pi', \pi)'$.) Indeed, the main results of this paper can and will be derived as relatively simple consequences of results concerning the bivariate regression model (3), which itself can be analyzed by means of fairly standard tools.

Our goal is to develop powerful tests of

$$H_0 : \beta = \beta_0 \quad \text{vs.} \quad H_1 : \beta \neq \beta_0,$$

treating π , γ , and f as unknown nuisance parameters. (Testing problems of this type are of interest partly because the duality between hypothesis testing and interval estimation implies that confidence intervals for β can be obtained by test inversion.) Replacing y_{1i} by $y_{1i} - \beta_0 y_{2i}$ if necessary, we assume without loss of generality that $\beta_0 = 0$.

The analysis proceeds under the following assumptions.

Assumption 1. (a) $Q_{zz,n} = n^{-1} \sum_{i=1}^n z_i z_i' \rightarrow Q_{zz} > 0$ and $\max_{1 \leq i \leq n} \|z_i\| / \sqrt{n} \rightarrow 0$.
 (b) $Q_{xx,n} = n^{-1} \sum_{i=1}^n x_i x_i' \rightarrow Q_{xx} > 0$ and $\max_{1 \leq i \leq n} \|x_i\| / \sqrt{n} \rightarrow 0$.

Assumption 2. The density f admits a function \dot{f} such that

- (a) for almost every $v \in \mathbb{R}^2$, f is differentiable at v , with (total) derivative \dot{f} .
- (b) for every $v \in \mathbb{R}^2$,

$$f(v + \theta) - f(v) = \theta' \int_0^1 \dot{f}(v + \theta t) dt, \quad \forall \theta \in \mathbb{R}^2.$$

- (c) $\int_{\mathbb{R}^2} \|\ell(v)\|^2 f(v) dv < \infty$, where

$$\ell(v) = -\frac{\dot{f}(v)}{f(v)} \mathbf{1}[f(v) > 0].$$

Assumption 3. $Q_{zx,n} = n^{-1} \sum_{i=1}^n z_i x_i' \rightarrow 0$.

Remarks. (i) In Assumption 1 and elsewhere in the paper, $\|\cdot\|$ is the Euclidean norm and limits are taken as $n \rightarrow \infty$, except where otherwise noted.

(ii) Assumption 1 is a fairly standard assumption concerning the exogenous variables. As in Bickel (1982), the assumption that the exogenous variables $(x'_i, z'_i)'$ are nonrandom can be relaxed (and the main results of this paper will remain valid) provided the errors $\{v_i\}$ are assumed to be independent of $\{(x'_i, z'_i)'\}$. Moreover, it seems plausible that certain forms of heteroskedasticity can be accommodated by adapting the methods of Schick (1997), but no attempts to do so will be made in this paper.

(iii) The assumption that (first and) second moments of the errors exist serves three purposes. First, it implies that the Fisher information matrix \mathcal{I} defined in (4) is nonsingular. Second, it implies that the \sqrt{n} -consistency requirements of Assumptions 6-8 are met by OLS estimators. Finally, it is required for the validity of the statements concerning procedures based on the Gaussian (quasi-)likelihood that are made throughout the paper. As in Bickel (1982), the main results of this paper are valid (and Assumptions 6-8 are met by suitable estimators) even without moment assumptions provided it is assumed that $\mathcal{I} > 0$.

(iv) Assumption 2 is a relatively mild smoothness condition on the error density. Parts (a) and (b) of Assumption 2 hold if, but do not require that, f is continuously differentiable. In particular, Assumption 2 accommodates mild departures from continuous differentiability, such as that which occurs when the elements of v_i (or some rotation thereof) are independent and double exponentially distributed.

(v) If Assumption 1 holds and $Q_{zx} = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n z_i x'_i$ exists, Assumption 3 is a normalization (i.e., it entails no loss of generality). Specifically, replacing z_i by $z_i^* = z_i - Q_{zx} Q_{xx}^{-1} x_i$ has no effect on the value of (β, π) and guarantees validity of Assumption 3. Our main results depend on $\{(x'_i, z'_i)'\}$ only through $\{z_i^*\}$, so Assumption 3 is convenient insofar as it enables us to simplify the notation by eliminating the distinction between $\{z_i\}$ and $\{z_i^*\}$.

(vi) Throughout this paper the endogenous regressor y_{2i} is assumed to be scalar. Most of our distributional results should generalize straightforwardly to models with multiple endogenous regressors, as should the optimality results reported in Section 4.4. On the other hand, analogues of the “near” optimality results (established by AMS for Moreira (2003)-type inference procedures in models with weak instruments and a scalar endogenous regressor) that underlie some of the efficiency claims made in other sections of the paper do not seem to be available for models with multiple endogenous regressors.

An immediate implication of Assumptions 1(a) and 2 is that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \ell(v_i) \otimes z_i \rightarrow_d \mathcal{N}(0, \mathcal{I} \otimes Q_{zz}),$$

where

$$\mathcal{I} = \int_{\mathbb{R}^2} \ell(v) \ell(v)' f(v) dv \quad (4)$$

is the Fisher information for the location family generated by f . Assumption 2 furthermore implies that the model (3) is differentiable in quadratic mean at any (γ, δ) (see (5) in the proof of Theorem A.1 in the Appendix) and enables nonparametric estimation of ℓ (as demonstrated by Theorem A.2 in the Appendix). In other words, the roles played by parts (a) and (b) of Assumption 2 are analogous to those played by the assumption of absolute continuity routinely invoked in regression models with scalar errors. In fact, the natural scalar counterpart of Assumption 2(b) is the assumption of absolute continuity.

As mentioned in remark (v), Assumption 3 is a normalization which greatly simplifies the derivation and statements of asymptotic results. Specifically, because the limit of $Q_{zx,n}$ is a zero matrix under Assumption 3, the parameters (β, π) and γ are orthogonal (in the sense of Cox and Reid (1987)). This fact, which is an immediate consequence of the fact that $\delta = (\delta_1, \delta_2)'$ and γ are orthogonal in (3), implies that the analysis can proceed under the “as if” assumption that γ is known. Similarly, the fact that $n^{-1} \sum_{i=1}^n z_i \rightarrow 0$ under Assumption 3 (because the first element of x_i equals unity) implies that the analysis can proceed under the “as if” assumption that f is known. This is so because δ in (3) can be estimated adaptively, the latter fact essentially following from Bickel’s (1982) result on adaptive estimation of slope coefficients in a regression model.

In other words, Assumption 3 implies that π is the only nuisance parameter which matters asymptotically. Concerning π , particular attention will be devoted to the weakly identified case where π is “close” to zero in the sense of the following assumption.

Assumption 4W. $\pi = c/\sqrt{n}$ for some constant $c \in \mathbb{R}^q$ and β is a constant.

Under the local-to-zero parameterization of π specified by Assumption 4W, contiguous alternatives to H_0 are of the form $\beta = \beta_0 + O(1)$. Accordingly, β is modeled as a constant in the weakly identified case. Although our main emphasis is on the weakly identified case, we shall on occasion employ one of the following (strong identification) assumptions.

Assumption 4SC. π is a nonzero constant and $\beta = b/\sqrt{n}$ for some constant $b \in \mathbb{R}$.

Assumption 4SF. π is a nonzero constant and β is a constant.

When π is a nonzero constant, identification is strong and contiguous alternatives to H_0 are of the form $\beta = \beta_0 + O(1/\sqrt{n})$. Assumption 4SC covers that case and is appropriate when studying local asymptotic power properties under strong identification. In contrast, Assumption 4SF assumes strong identification and furthermore holds β fixed. This combination of strong identification and fixed alternatives is appropriate when studying the consistency properties of various tests. Moreover, Assumption 4SF is useful when studying the properties of point estimators of β under strong identification.

Assumptions 4W, 4SC, and 4SF are nonnested, but it seems natural to study them in the order indicated above. This is so because the assumptions impose decreasingly strong upper bounds on the magnitude of the parameters δ_1 and δ_2 of (3). Specifically, Assumption 4W implies that $\delta_1 = O(1/\sqrt{n})$ and $\delta_2 = O(1/\sqrt{n})$. Relative to Assumption 4W, Assumption 4SC removes the requirement $\delta_2 = O(1/\sqrt{n})$ and Assumption 4SF furthermore relaxes the requirement $\delta_1 = O(1/\sqrt{n})$. In this paper, these differences are important because the feasible inference procedures constructed in Section 4 employ one-step estimators of δ . As usual, one-step estimators utilize initial estimators that are required to be \sqrt{n} -consistent. Under Assumption 4W, this requirement is met by the zero vector, while Assumption 4SC and 4SF imply that nondegenerate initial estimators of δ_2 and (δ_1, δ_2) , respectively, are required in order to guarantee that one-step estimators of δ are well behaved. Accordingly, the three constructions presented in Section 4 differ in terms of (and only in terms of) the nature of the initial estimators of δ being employed.

3. THE LIMITING EXPERIMENT WHEN IDENTIFICATION IS WEAK

This section is concerned with asymptotic inference under the assumptions that (i) the nuisance parameters γ and f are known and (ii) identification is weak. As mentioned in the previous section, Assumption 3 ensures that (i) can be dispensed with. Precise statements to that effect will be provided in the next section, where it is also shown how departures from (ii) can be accommodated.

When f is Gaussian and the reduced form variance

$$\Omega = \int_{\mathbb{R}^2} vv' f(v) dv$$

is known, the problem of testing $\beta = \beta_0$ vs. $\beta \neq \beta_0$ is nonstandard, but amenable to finite sample analysis using the theory of curved exponential families (e.g., Moreira (2003) and AMS). This feature is lost, in general, when f is not Gaussian. On the other hand, the testing problem remains amenable to asymptotic analysis using the limits of experiments approach even when f is non-Gaussian.² In fact, it turns out

²For an exposition of the elements of the theory of limits of experiments employed in this paper,

that the family of limiting experiments associated with non-Gaussian error distributions coincides with the family of limiting experiments for the Gaussian case.

In the Gaussian case, the limiting experiment is that of a single observation from the $\mathcal{N}[\mu(\beta, c), \Omega \otimes Q_{zz}^{-1}]$ distribution, where

$$\mu(\beta, c) = \begin{pmatrix} \beta \\ 1 \end{pmatrix} \otimes c.$$

Equivalently, because $\Omega = \mathcal{I}^{-1}$ when f is Gaussian, the limiting experiment in the Gaussian case is that of a single observation from the $\mathcal{N}[\mu(\beta, c), \mathcal{I}^{-1} \otimes Q_{zz}^{-1}]$ distribution. As it turns out, the latter characterization generalizes readily to non-Gaussian error distributions.

To give a precise statement, we proceed in the spirit of van der Vaart (1998, Section 7.6). Define the log likelihood ratio function

$$\begin{aligned} L_n(\beta, c) &= \sum_{i=1}^n \log f(y_{1i} - \gamma'_1 x_i - \beta c' z_i / \sqrt{n}, y_{2i} - \gamma'_2 x_i - c' z_i / \sqrt{n}) \\ &\quad - \sum_{i=1}^n \log f(y_{1i} - \gamma'_1 x_i, y_{2i} - \gamma'_2 x_i) \end{aligned}$$

and let “ $o_{p_0}(1)$ ” and “ \rightarrow_{d_0} ” be shorthand for “ $o_p(1)$ under the distributions associated with $(\beta, \pi) = (0, 0)$ ” and “ \rightarrow_d under the distributions associated with $(\beta, \pi) = (0, 0)$ ”, respectively.

Theorem 1. *If Assumptions 1(a) and 2 hold, then*

$$L_n(\beta, c) = \mu(\beta, c)' (\mathcal{I} \otimes Q_{zz}) \Delta_n - \frac{1}{2} \mu(\beta, c)' (\mathcal{I} \otimes Q_{zz}) \mu(\beta, c) + o_{p_0}(1)$$

for every (β, c) , where

$$\Delta_n = (\mathcal{I}^{-1} \otimes Q_{zz}^{-1}) \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell(y_{1i} - \gamma'_1 x_i, y_{2i} - \gamma'_2 x_i) \otimes z_i \rightarrow_{d_0} \mathcal{N}(0, \mathcal{I}^{-1} \otimes Q_{zz}^{-1}).$$

Theorem 1 is a special case of a local asymptotic normality (LAN) result for the model (3). The general LAN result is given in Theorem A.1 in the Appendix.

see e.g. van der Vaart (1998).

As in van der Vaart (1998, Section 9.3), Theorem 1 and Le Cam’s third lemma can be used to show that if Assumptions 1(a), 2, and 4W hold, then the asymptotically sufficient statistic Δ_n satisfies

$$\Delta_n \rightarrow_d \mathcal{N} [\mu(\beta, c), \mathcal{I}^{-1} \otimes Q_{zz}^{-1}],$$

implying in particular that the limiting experiment is that of a single observation from the $\mathcal{N}[\mu(\beta, c), \mathcal{I}^{-1} \otimes Q_{zz}^{-1}]$ distribution whether or not the errors are Gaussian.

Under the same assumptions, the quasi-sufficient (i.e., sufficient when the errors are Gaussian) statistic

$$\bar{\Delta}_n = (\Omega \otimes Q_{zz,n}^{-1}) \frac{1}{\sqrt{n}} \sum_{i=1}^n \bar{\ell}(y_{1i} - \gamma'_1 x_i, y_{2i} - \gamma'_2 x_i) \otimes z_i, \quad \bar{\ell}(v) = \Omega^{-1} v,$$

obtained from the Gaussian quasi-likelihood satisfies

$$\bar{\Delta}_n \rightarrow_d \mathcal{N} [\mu(\beta, c), \Omega \otimes Q_{zz}^{-1}].$$

The Cauchy-Schwarz inequality can be used to show that $\mathcal{I}^{-1} \leq \Omega$, with equality if and only if $\ell(v)$ is linear in v (on the support of f). By implication, procedures based on the Gaussian quasi-likelihood are asymptotically inefficient in general. More specifically, any test based on a “smooth” (e.g., almost everywhere continuous) function of $\bar{\Delta}_n$, such as those proposed by Anderson and Rubin (1949), Kleibergen (2002), and Moreira (2003), will be dominated (under weak identification and whenever the inequality $\mathcal{I}^{-1} \leq \Omega$ is strict) by a test which is efficient (or “nearly” so) under the assumptions of Theorem 1. (Section 4 will exhibit tests which are “nearly” efficient under the assumptions of Theorem 1.)

Nevertheless, the results obtained under the assumption of Gaussian errors are of considerable relevance also in models with non-Gaussian errors. This is so because the limiting experiments (indexed by $\mathcal{I}^{-1} \otimes Q_{zz}^{-1}$) in the general case are isomorphic to the limiting experiments (indexed by $\Omega \otimes Q_{zz}^{-1}$) associated with Gaussian errors, a very convenient result because it implies that the insights concerning the relative merits of various testing procedures obtained under the assumption of normality are directly applicable in the general case.

To be specific, let $S_n, T_n \in \mathbb{R}^q$ be given by

$$\begin{pmatrix} S_n \\ T_n \end{pmatrix} = [\mathcal{I}^{1/2'} \otimes Q_{zz}^{1/2'}] \Delta_n,$$

where $M^{1/2}$ denotes the upper triangular Cholesky factor of a (symmetric, positive semi-definite) matrix M ; that is, $M = M^{1/2}M^{1/2'}$, where $M^{1/2}$ is upper triangular.³ The pair (S_n, T_n) is a non-Gaussian counterpart of

$$\begin{pmatrix} \bar{S}_n \\ \bar{T}_n \end{pmatrix} = \left[(\Omega^{-1})^{1/2'} \otimes Q_{zz,n}^{1/2'} \right] \bar{\Delta}_n,$$

which features prominently in the work by Moreira (2003), AMS, and others.

In terms of (\bar{S}_n, \bar{T}_n) , the (known Ω) Anderson-Rubin, Lagrange multiplier, and likelihood ratio test statistics popularized by Anderson and Rubin (1949), Kleibergen (2002), and Moreira (2003), respectively, can be expressed as

$$\begin{aligned} \overline{AR}_n &= \bar{S}'_n \bar{S}_n, & \overline{LM}_n &= \frac{(\bar{S}'_n \bar{T}_n)^2}{\bar{T}'_n \bar{T}_n}, \\ \overline{LR}_n &= \frac{1}{2} \left(\bar{S}'_n \bar{S}_n - \bar{T}'_n \bar{T}_n + \sqrt{(\bar{S}'_n \bar{S}_n - \bar{T}'_n \bar{T}_n)^2 + 4(\bar{S}'_n \bar{T}_n)^2} \right). \end{aligned}$$

In perfect analogy with the Gaussian case, let

$$\begin{aligned} AR_n &= S'_n S_n, & LM_n &= \frac{(S'_n T_n)^2}{T'_n T_n}, \\ LR_n &= \frac{1}{2} \left(S'_n S_n - T'_n T_n + \sqrt{(S'_n S_n - T'_n T_n)^2 + 4(S'_n T_n)^2} \right). \end{aligned}$$

The tests which reject H_0 when $AR_n > \chi_\alpha^2(q)$, $LM_n > \chi_\alpha^2(1)$, and $LR_n > \kappa_\alpha(T_n)$ have asymptotic size α , where $\chi_\alpha^2(d)$ is the $1 - \alpha$ quantile of the χ^2 distribution with d degrees of freedom and $\kappa_\alpha(t)$ is the $1 - \alpha$ quantile of the distribution of $\frac{1}{2} \left(\mathcal{Z}' \mathcal{Z} - t't + \sqrt{(\mathcal{Z}' \mathcal{Z} - t't)^2 + 4(\mathcal{Z}' t)^2} \right)$, where $\mathcal{Z} \sim \mathcal{N}(0, I_q)$. Because of the isomorphism between the Gaussian case and the general case, the relative merits of these testing procedures are well understood from the numerical work of AMS. In particular, it follows from AMS that the test which rejects when $LR_n > \kappa_\alpha(T_n)$ is “nearly efficient” in the sense that its power function is “close” to the two-sided power envelope for invariant similar tests.

³In particular, letting \mathcal{I}_{ij} denote element (i, j) of \mathcal{I} , we have:

$$\mathcal{I}^{1/2} = \begin{pmatrix} \sqrt{\mathcal{I}_{11.2}} & \mathcal{I}_{12}/\sqrt{\mathcal{I}_{22}} \\ 0 & \sqrt{\mathcal{I}_{22}} \end{pmatrix}, \quad \mathcal{I}_{11.2} = \mathcal{I}_{11} - \mathcal{I}_{12}^2/\mathcal{I}_{22}.$$

Remarks. (i) As shown by Moreira (2003), $\kappa_\alpha(t)$ depends on t only through $\|t\|$, is monotonically decreasing in $\|t\|$, and satisfies $\lim_{\|t\| \rightarrow \infty} \kappa_\alpha(t) = \chi_\alpha^2(1)$. The latter result will be utilized when studying the behavior of the test based on LR_n under strong identification.

(ii) The existence of tests which are equivalent to procedures based on the Gaussian quasi-likelihood under the assumption of normality and enjoy improved power properties for certain non-Gaussian error distributions has been pointed out by Andrews and Marmor (2007) and Andrews and Soares (2007). The rank-based tests proposed in those papers, while superior to tests based on the Gaussian quasi-likelihood under some conditions, are also inefficient in general (even for those error distributions for which they dominate tests based on the Gaussian quasi-likelihood).

4. FEASIBLE INFERENCE PROCEDURES

The results of the previous section were obtained under the (tacit) assumption that γ and f are known. In addition, it was assumed to be known that identification is weak (i.e., that π is “close” to zero). This section relaxes these assumptions.

4.1. Inference without knowledge of γ and f . First, consider the problem of conducting inference under weak identification without knowledge of the nuisance parameters γ and f . Doing so is easy provided we can find a pair $(\hat{\Delta}_n, \hat{\mathcal{I}}_n)$ which is asymptotically equivalent to (Δ_n, \mathcal{I}) under weak identification and can be computed without knowledge of (γ, f) . To that end, let

$$\hat{\Delta}_n = \left(\hat{\mathcal{I}}_n^{-1} \otimes Q_{zz,n}^{-1} \right) \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\ell}_{i,n} \otimes z_i, \quad \hat{\mathcal{I}}_n = n^{-1} \sum_{i=1}^n \hat{\ell}_{i,n} \hat{\ell}'_{i,n},$$

where $\hat{\ell}_{i,n}$ is an estimator of $\ell(y_{1i} - \gamma'_1 x_i, y_{2i} - \gamma'_2 x_i)$. In the spirit of Schick (1987), we assume that $\hat{\ell}_{i,n} = \hat{\ell}_n(\hat{v}_i)$, where $\hat{v}_i = (y_{1i} - \hat{\gamma}'_{1n} x_i, y_{2i} - \hat{\gamma}'_{2n} x_i)'$ for some estimator $\hat{\gamma}_n = (\hat{\gamma}'_{1n}, \hat{\gamma}'_{2n})'$ of γ and

$$\hat{\ell}_n(v) = -\frac{\partial \hat{f}_n(v) / \partial v}{\hat{f}_n(v) + a_n}, \quad \hat{f}_n(v) = \frac{1}{nh_n^2} \sum_{i=1}^n K\left(\frac{v - \hat{v}_i}{h_n}\right),$$

where K is a kernel and a_n and h_n are positive sequences. Theorem 2 shows that this construction, which does not involve sample splitting, works when the following assumptions hold.

Assumption 5. (a) $K(s_1, s_2) = k(s_1)k(s_2)$, where k is a bounded, symmetric, continuously differentiable density satisfying

$$\int_{\mathbb{R}} r^2 k(r) dr < \infty \quad \text{and} \quad \sup_{r \in \mathbb{R}} |k'(r)|/k(r) < \infty.$$

(b) $a_n \rightarrow 0$, $h_n \rightarrow 0$, and $na_n^2 h_n^4 \rightarrow \infty$.

Assumption 6. $\hat{\gamma}_n$ is discrete and $\sqrt{n}(\hat{\gamma}_n - \gamma) = O_p(1)$.

Remarks. (i) The nonparametric estimation method used here involves two smoothing parameters, h_n and a_n , of which the former is a bandwidth sequence whereas the latter enables us to avoid trimming when handling the density estimator \hat{f}_n appearing in the denominator of $\hat{\ell}_n$. Choosing smoothing parameters is well known to be a difficult problem in practice, but it is beyond the scope of this paper to provide guidance about how to choose h_n and a_n .

(ii) If the variances of v_1 and v_2 are suspected to be of different magnitude it may be desirable to let K be a product kernel of the form $K(s_1, s_2) = \sigma_1^{-1} \sigma_2^{-1} k(s_1/\sigma_1) k(s_2/\sigma_2)$, where σ_1 and σ_2 are positive constants and k is as in Assumption 5(a). All results (and their proofs) remain valid if Assumption 5(a) is modified in this way.

(iii) Assumption 5(a) holds if k is the logistic density, but not if k is the standard normal density, the reason being that the normal density violates the condition $\sup_{r \in \mathbb{R}} |k'(r)|/k(r) < \infty$. As explained in remark (ii) following the proof of Theorem A.2 in the Appendix, it is possible to accommodate the normal kernel provided the error density f is such that $\sup_{v \in \mathbb{R}^2} \|\dot{f}(v)\| < \infty$ (and provided the requirement $\overline{\lim}_{n \rightarrow \infty} h_n/a_n < \infty$ is added to Assumption 5(b)).

(iv) In Assumption 6, the statement “ $\hat{\gamma}_n$ is discrete” is shorthand for the assumption that $\hat{\gamma}_n$ takes only values in the grid $\{\varkappa Z/\sqrt{n} : Z \in \mathbb{Z}^{2p}\}$, where \varkappa is some (constant) $2p \times 2p$ matrix. (Assuming discreteness on the part of an initial estimator is technically convenient and it seems plausible that this assumption can be dropped if additional smoothness is assumed on the part of f .) If $\tilde{\gamma}_n$ is a \sqrt{n} -consistent estimator of γ , then Assumption 6 will be satisfied by $\hat{\gamma}_n = \lfloor \sqrt{n} \tilde{\gamma}_n \rfloor / \sqrt{n}$, where $\lfloor \cdot \rfloor$ denotes the integer part of the argument (defined element-by-element). A similar remark applies to Assumptions 7 and 8.

(v) Assumption 6 is satisfied (under both weak and strong identification) by a discretized version of

$$\hat{\gamma}_n^{OLS} = \begin{bmatrix} (\sum_{i=1}^n x_i x_i')^{-1} (\sum_{i=1}^n x_i y_{1i}) \\ (\sum_{i=1}^n x_i x_i')^{-1} (\sum_{i=1}^n x_i y_{2i}) \end{bmatrix},$$

the OLS estimator of γ .

Theorem 2. *If Assumptions 1-3, 4W, and 5-6 hold, then*

$$\left(\hat{\Delta}_n, \hat{\mathcal{I}}_n \right) = (\Delta_n, \mathcal{I}) + o_p(1).$$

In the model (3), the statistic $\hat{\Delta}_n/\sqrt{n}$ can be interpreted as a one-step estimator of δ which uses the zero vector as an initial estimator. As a consequence, Theorem 2 can (and will) be derived as a special case of a general adaptation result, Theorem A.2 in the Appendix, for one-step estimators of δ in the model (3). Theorem A.2 assumes existence of a (discrete) \sqrt{n} -consistent initial estimator of δ . This requirement is easily met, especially so under weak identification because the zero vector can serve as a (discrete) \sqrt{n} -consistent estimator of δ in that case. (The full force of Theorem A.2 will be needed when Assumption 4W is replaced by Assumption 4SC or 4SF.) Somewhat surprisingly, perhaps, some aspects of conducting inference are therefore simplified by the assumption of weak identification.

Theorem 2 (and the continuous mapping theorem) can be used to show that if identification is weak, then the local asymptotic power properties of the tests based on AR_n , LM_n , and LR_n are matched by those of the tests based on

$$\widehat{AR}_n = \hat{S}'_n \hat{S}_n, \quad \widehat{LM}_n = \frac{\left(\hat{S}'_n \hat{T}_n \right)^2}{\hat{T}'_n \hat{T}_n},$$

and

$$\widehat{LR}_n = \frac{1}{2} \left(\hat{S}'_n \hat{S}_n - \hat{T}'_n \hat{T}_n + \sqrt{\left(\hat{S}'_n \hat{S}_n - \hat{T}'_n \hat{T}_n \right)^2 + 4 \left(\hat{S}'_n \hat{T}_n \right)^2} \right),$$

respectively, where $\left(\hat{S}'_n, \hat{T}'_n \right)' = \left[\hat{\mathcal{I}}_n^{1/2'} \otimes Q_{zz,n}^{1/2'} \right] \hat{\Delta}_n$. More specifically, we have the following corollary, which implies in particular that the (feasible) test which rejects when $\widehat{LR}_n > \kappa_\alpha \left(\hat{T}_n \right)$ is “nearly efficient” when identification is weak.

Corollary 3. *If Assumptions 1-3, 4W, and 5-6 hold, then*

$$\left[\widehat{AR}_n, \widehat{LM}_n, \widehat{LR}_n, \kappa_\alpha \left(\hat{T}_n \right) \right] = \left[AR_n, LM_n, LR_n, \kappa_\alpha (T_n) \right] + o_p(1).$$

4.2. Inference when identification may be strong. Next, consider the consequences of relaxing the assumption that identification is known to be weak. We are interested in finding a pair of statistics, computable without knowledge of (γ, f) , which is asymptotically equivalent to (Δ_n, \mathcal{I}) under weak identification and is “well behaved” also when identification is strong.

When Assumptions 1-3 and 4SC hold, the quasi-sufficient statistic $\bar{\Delta}_n$ obtained from the Gaussian quasi-likelihood satisfies

$$\bar{\Delta}_n - \sqrt{n} \begin{pmatrix} 0 \\ \pi \end{pmatrix} \rightarrow_d \mathcal{N} \left[\begin{pmatrix} b\pi \\ 0 \end{pmatrix}, \Omega \otimes Q_{zz}^{-1} \right].$$

It follows immediately from this result that if Assumptions 1-3 and 4SC holds, then

$$\overline{AR}_n \rightarrow_d \chi^2 (q; b^2 \omega_{11}^{-1} \pi' Q_{zz} \pi)$$

and

$$\overline{LM}_n = \overline{LR}_n + o_p(1) = \frac{\left(\bar{S}'_n Q_{zz}^{1/2} \pi \right)^2}{\pi' Q_{zz} \pi} + o_p(1) \rightarrow_d \chi^2 (1; b^2 \omega_{11}^{-1} \pi' Q_{zz} \pi),$$

where ω_{11} is element (1, 1) of Ω and $\chi^2(d; \lambda)$ denotes the noncentral χ^2 distribution with d degrees of freedom and noncentrality parameter λ . (Moreover, the properties of κ_α mentioned in remark (i) of Section 3 can be used to show that $\kappa_\alpha(\bar{T}_n) = \chi_\alpha^2(1) + o_p(1)$.) The convergence result for $\bar{\Delta}_n$ derives in part from the linearity of $\bar{\ell}$ and an analogous result will typically fail to hold for Δ_n and/or $\hat{\Delta}_n$. Indeed, at the present level of generality very little can be said about the asymptotic null properties of statistics such as \widehat{LR}_n under strong identification. This observation motivates the search for a statistic which is asymptotically equivalent to Δ_n under weak identification and exhibits behavior qualitatively similar to that of $\bar{\Delta}_n$ under Assumption 4SC.

Theorem 4 gives conditions under which this property is enjoyed by

$$\hat{\Delta}_n^* = \begin{pmatrix} 0 \\ \sqrt{n} \hat{\pi}_n \end{pmatrix} + \left(\hat{\mathcal{I}}_n^{*-1} \otimes Q_{zz,n}^{-1} \right) \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\ell}_{i,n}^* \otimes z_i, \quad \hat{\mathcal{I}}_n^* = n^{-1} \sum_{i=1}^n \hat{\ell}_{i,n}^* \hat{\ell}_{i,n}^{*'},$$

with $\hat{\ell}_{i,n}^* = \hat{\ell}_n^*(\hat{v}_i^*)$, where $\hat{v}_i^* = (y_{1i} - \hat{\gamma}'_{1n} x_i, y_{2i} - \hat{\gamma}'_{2n} x_i - \hat{\pi}'_n z_i)'$ for some estimators $(\hat{\gamma}_n, \hat{\pi}_n)$ of (γ, π) , and

$$\hat{\ell}_n^*(v) = -\frac{\partial \hat{f}_n^*(v) / \partial v}{\hat{f}_n^*(v) + a_n}, \quad \hat{f}_n^*(v) = \frac{1}{nh_n^2} \sum_{i=1}^n K \left(\frac{v - \hat{v}_i^*}{h_n} \right).$$

As defined, $\hat{\Delta}_n^*/\sqrt{n}$ is a one-step estimator of δ (in the model (3)) which uses $(0', \hat{\pi}'_n)'$ as an initial estimator of δ . This initial estimator is \sqrt{n} -consistent under Assumption 4SC provided $\hat{\pi}_n$ satisfies the following condition.

Assumption 7. $\hat{\pi}_n$ is discrete and $\sqrt{n}(\hat{\pi}_n - \pi) = O_p(1)$.

Assumption 7 holds (under both weak and strong identification) if $\hat{\pi}_n$ is a discretized version of $\hat{\pi}_n^{OLS} = (\sum_{i=1}^n z_i z_i')^{-1} (\sum_{i=1}^n z_i y_{2i})$.

Theorem 4. (a) If Assumptions 1-3, 4W, and 5-7 hold, then

$$\left(\hat{\Delta}_n^*, \hat{\mathcal{I}}_n^* \right) = (\Delta_n, \mathcal{I}) + o_p(1).$$

(b) If Assumptions 1-3, 4SC, and 5-7 hold, then $\hat{\mathcal{I}}_n^* = \mathcal{I} + o_p(1)$ and

$$\hat{\Delta}_n^* - \sqrt{n} \begin{pmatrix} 0 \\ \pi \end{pmatrix} \rightarrow_d \mathcal{N} \left[\begin{pmatrix} b\pi \\ 0 \end{pmatrix}, \mathcal{I}^{-1} \otimes Q_{zz}^{-1} \right].$$

Remark. The pair $(\hat{\Delta}_n^*, \hat{\mathcal{I}}_n^*)$ reduces to $(\hat{\Delta}_n, \hat{\mathcal{I}}_n)$ when $\hat{\pi}_n = 0$. Moreover, $\sqrt{n}\pi = O(1)$ under weak identification, so Theorem 2 is a special case of Theorem 4 (a).

As a consequence of Theorem 4, we have the following result concerning the statistics

$$\begin{aligned} \widehat{AR}_n^* &= \hat{S}_n^{*'} \hat{S}_n^*, & \widehat{LM}_n^* &= \frac{(\hat{S}_n^{*'} \hat{T}_n^*)^2}{\hat{T}_n^{*'} \hat{T}_n^*}, \\ \widehat{LR}_n^* &= \frac{1}{2} \left(\hat{S}_n^{*'} \hat{S}_n^* - \hat{T}_n^{*'} \hat{T}_n^* + \sqrt{\left(\hat{S}_n^{*'} \hat{S}_n^* - \hat{T}_n^{*'} \hat{T}_n^* \right)^2 + 4 \left(\hat{S}_n^{*'} \hat{T}_n^* \right)^2} \right), \end{aligned}$$

where $(\hat{S}_n^{*'}, \hat{T}_n^{*'})' = [\hat{\mathcal{I}}_n^{*1/2'} \otimes Q_{zz,n}^{1/2'}] \hat{\Delta}_n^*$.

Corollary 5. (a) If Assumptions 1-3, 4W, and 5-7 hold, then

$$\left[\widehat{AR}_n^*, \widehat{LM}_n^*, \widehat{LR}_n^*, \kappa_\alpha(\hat{T}_n^*) \right] = \left[AR_n, LM_n, LR_n, \kappa_\alpha(T_n) \right] + o_p(1).$$

(b) If Assumptions 1-3, 4SC, and 5-7 hold, then

$$\begin{aligned} \widehat{AR}_n^* &= AR_n + o_p(1) \rightarrow_d \chi^2(q; b^2 \mathcal{I}_{11.2} \pi' Q_{zz} \pi), \\ \widehat{LM}_n^* &= \widehat{LR}_n^* + o_p(1) = \frac{(S_n' Q_{zz}^{1/2'} \pi)^2}{\pi' Q_{zz} \pi} + o_p(1) \rightarrow_d \chi^2(1; b^2 \mathcal{I}_{11.2} \pi' Q_{zz} \pi), \end{aligned}$$

and $\kappa_\alpha(\hat{T}_n^*) = \chi_\alpha^2(1) + o_p(1)$.

It follows from Corollary 5(a) that the test which rejects when $\widehat{LR}_n^* > \kappa_\alpha(\widehat{T}_n^*)$ is “nearly efficient” when identification is weak. Moreover, Theorem A.1 in the Appendix and Choi, Hall, and Schick (1996, Theorem 2) can be used to show that the test which rejects for large values of $\left(S_n' Q_{zz}^{1/2} \pi\right)^2 / (\pi' Q_{zz} \pi)$ is asymptotically uniformly most powerful unbiased (in the terminology of Choi, Hall, and Schick (1996, Section 4)) under the assumptions of Corollary 5(b). As a consequence, Corollary 5(b) implies that the test which rejects when $\widehat{LR}_n^* > \kappa_\alpha(\widehat{T}_n^*)$ enjoys demonstrable optimality properties under strong identification, as does the test which rejects when $\widehat{LM}_n^* > \chi_\alpha^2(1)$. In particular, under strong identification these (asymptotically equivalent) tests are superior to the tests based on the statistics \overline{AR}_n , \overline{LM}_n , \overline{LR}_n (and Andrews and Soares’s (2007) rank-based analogues thereof).

4.3. Consistency. Finally, we address the issue of test consistency under strong identification. The tests based on \overline{AR}_n , \overline{LM}_n , and \overline{LR}_n are all consistent because $(\kappa_\alpha(\cdot))$ is bounded and

$$n^{-1} \overline{AR}_n = n^{-1} \overline{LM}_n + o_p(1) = n^{-1} \overline{LR}_n + o_p(1) = \beta^2 \omega_{11}^{-1} \pi' Q_{zz} \pi + o_p(1)$$

under Assumptions 1-3 and 4SF, the displayed results following almost immediately from the fact that if Assumptions 1-3 and 4SF hold, then

$$\bar{\Delta}_n - \sqrt{n} \begin{pmatrix} \beta\pi \\ \pi \end{pmatrix} \rightarrow_d \mathcal{N} \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Omega \otimes Q_{zz}^{-1} \right].$$

Once again, this convergence result for $\bar{\Delta}_n$ derives in part from the linearity of $\bar{\ell}$ and an analogous result will typically fail to hold for Δ_n , $\hat{\Delta}_n$ and/or $\hat{\Delta}_n^*$. In fact, at the present level of generality there is no guarantee that the tests based on \widehat{AR}_n^* , \widehat{LM}_n^* , and \widehat{LR}_n^* are consistent under strong identification.

Fortunately this potential problem is easily avoided. Indeed, let

$$\hat{\Delta}_n^{**} = \begin{pmatrix} \sqrt{n} \hat{\Pi}_n \\ \sqrt{n} \hat{\pi}_n \end{pmatrix} + \left(\hat{\mathcal{I}}_n^{**} \otimes Q_{zz}^{-1} \right) \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\ell}_{i,n}^{**} \otimes z_i, \quad \hat{\mathcal{I}}_n^{**} = n^{-1} \sum_{i=1}^n \hat{\ell}_{i,n}^{**} \hat{\ell}_{i,n}^{**'}$$

with $\hat{\ell}_{i,n}^{**} = \hat{\ell}_n^{**}(\hat{v}_i^{**})$, where $\hat{v}_i^{**} = \left(y_{1i} - \hat{\gamma}'_{1n} x_i - \hat{\Pi}'_n z_i, y_{2i} - \hat{\gamma}'_{2n} x_i - \hat{\pi}'_n z_i \right)'$ for some estimators $\left(\hat{\gamma}_n, \hat{\pi}_n, \hat{\Pi}_n \right)$ of $(\gamma, \pi, \beta\pi)$,

$$\hat{\ell}_n^{**}(v) = -\frac{\partial \hat{f}_n^{**}(v)/\partial v}{\hat{f}_n^{**}(v) + a_n}, \quad \hat{f}_n^{**}(v) = \frac{1}{nh_n^2} \sum_{i=1}^n K\left(\frac{v - \hat{v}_i^{**}}{h_n}\right),$$

and $\hat{\Pi}_n$ is assumed to satisfy the following condition, which holds (under weak and strong identification) if $\hat{\Pi}_n$ is a discretized version of $\hat{\Pi}_n^{OLS} = (\sum_{i=1}^n z_i z_i')^{-1} (\sum_{i=1}^n z_i y_{1i})$.

Assumption 8. $\hat{\Pi}_n$ is discrete and $\sqrt{n}(\hat{\Pi}_n - \beta\pi) = O_p(1)$.

Once again, $\hat{\Delta}_n^{**}/\sqrt{n}$ can be interpreted as a one-step estimator of δ in (3). Unlike $\hat{\Delta}_n/\sqrt{n}$ and $\hat{\Delta}_n^*/\sqrt{n}$, $\hat{\Delta}_n^{**}/\sqrt{n}$ employs an initial estimator of δ with global \sqrt{n} -consistency properties. This feature is utilized in the proof of part (c) of the following result, which in turn can be used to establish consistency of tests based on $\hat{\Delta}_n^{**}$.

Theorem 6. (a) If Assumptions 1-3, 4W, and 5-8 hold, then

$$\left(\hat{\Delta}_n^{**}, \hat{\mathcal{I}}_n^{**}\right) = (\Delta_n, \mathcal{I}) + o_p(1).$$

(b) If Assumptions 1-3, 4SC, and 5-8 hold, then $\hat{\mathcal{I}}_n^{**} = \mathcal{I} + o_p(1)$ and

$$\hat{\Delta}_n^{**} - \sqrt{n} \begin{pmatrix} 0 \\ \pi \end{pmatrix} \rightarrow_d \mathcal{N} \left[\begin{pmatrix} b\pi \\ 0 \end{pmatrix}, \mathcal{I}^{-1} \otimes Q_{zz}^{-1} \right].$$

(c) If Assumptions 1-3, 4SF, and 5-8 hold, then $\hat{\mathcal{I}}_n^{**} = \mathcal{I} + o_p(1)$ and

$$\hat{\Delta}_n^{**} - \sqrt{n} \begin{pmatrix} \beta\pi \\ \pi \end{pmatrix} \rightarrow_d \mathcal{N} \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \mathcal{I}^{-1} \otimes Q_{zz}^{-1} \right].$$

Remark. The pair $(\hat{\Delta}_n^{**}, \hat{\mathcal{I}}_n^{**})$ reduces to $(\hat{\Delta}_n^*, \hat{\mathcal{I}}_n^*)$ when $\hat{\Pi}_n = 0$. Moreover, $\sqrt{n}\beta\pi = O(1)$ under the assumptions of Theorem 4, so the latter theorem is a special case of Theorem 6(a)-(b).

Let $(\hat{S}_n^{**'}, \hat{T}_n^{**'})' = [\hat{\mathcal{I}}_n^{**'}]^{1/2'} \otimes Q_{zz,n}^{1/2'} \hat{\Delta}_n^{**}$ and define

$$\widehat{AR}_n^{**} = \hat{S}_n^{**'} \hat{S}_n^{**}, \quad \widehat{LM}_n^{**} = \frac{(\hat{S}_n^{**'} \hat{T}_n^{**})^2}{\hat{T}_n^{**'} \hat{T}_n^{**}},$$

$$\widehat{LR}_n^{**} = \frac{1}{2} \left(\hat{S}_n^{**'} \hat{S}_n^{**} - \hat{T}_n^{**'} \hat{T}_n^{**} + \sqrt{(\hat{S}_n^{**'} \hat{S}_n^{**} - \hat{T}_n^{**'} \hat{T}_n^{**})^2 + 4(\hat{S}_n^{**'} \hat{T}_n^{**})^2} \right).$$

The salient properties of these statistics are characterized in the following corollary to Theorem 6.

Corollary 7. (a) *If Assumptions 1-3, 4W, and 5-8 hold, then*

$$\left[\widehat{AR}_n^{**}, \widehat{LM}_n^{**}, \widehat{LR}_n^{**}, \kappa_\alpha \left(\widehat{T}_n^{**} \right) \right] = \left[AR_n, LM_n, LR_n, \kappa_\alpha (T_n) \right] + o_p(1).$$

(b) *If Assumptions 1-3, 4SC, and 5-8 hold, then*

$$\begin{aligned} \widehat{AR}_n^{**} &= AR_n + o_p(1) \rightarrow_d \chi^2(q; b^2 \mathcal{I}_{11.2} \pi' Q_{zz} \pi), \\ \widehat{LM}_n^{**} &= \widehat{LR}_n^{**} + o_p(1) = \frac{\left(S_n' Q_{zz}^{1/2} \pi \right)^2}{\pi' Q_{zz} \pi} + o_p(1) \rightarrow_d \chi^2(1; b^2 \mathcal{I}_{11.2} \pi' Q_{zz} \pi), \end{aligned}$$

and $\kappa_\alpha \left(\widehat{T}_n^{**} \right) = \chi_\alpha^2(1) + o_p(1)$.

(c) *If Assumptions 1-3, 4SF, and 5-8 hold, then*

$$n^{-1} \widehat{AR}_n^{**} = n^{-1} \widehat{LM}_n^{**} + o_p(1) = n^{-1} \widehat{LR}_n^{**} + o_p(1) = \beta^2 \mathcal{I}_{11.2} \pi' Q_{zz} \pi + o_p(1).$$

In perfect analogy with Corollary 5, parts (a) and (b) of Corollary 7 imply that the test which rejects when $\widehat{LR}_n^{**} > \kappa_\alpha \left(\widehat{T}_n^{**} \right)$ is “nearly” optimal when identification is weak and demonstrably optimal when identification is strong. Relative to Corollary 5, which establishes analogous results for the test which rejects when $\widehat{LR}_n^* > \kappa_\alpha \left(\widehat{T}_n^* \right)$, the additional property that can be claimed on the part of the test based on \widehat{LR}_n^{**} is that of consistency under strong identification. This, and the analogous consistency results about the tests based on \widehat{AR}_n^{**} and \widehat{LM}_n^{**} , is the content of Corollary 7(c).

4.4. Inference when identification is strong. If identification is strong, then the usual duality between estimation and testing holds, implying in particular that the asymptotic optimality properties of the tests based on \widehat{LR}_n^{**} and \widehat{LM}_n^{**} are shared by a Wald test based on an asymptotically efficient estimator of β .

Let

$$\hat{\beta}_n^{**} = \frac{\hat{\Delta}_{1,n}^{**'} Q_{zz,n} \hat{\Delta}_{2,n}^{**}}{\hat{\Delta}_{2,n}^{**'} Q_{zz,n} \hat{\Delta}_{2,n}^{**}},$$

where $\hat{\Delta}_n^{**} = \left(\hat{\Delta}_{1,n}^{**'}, \hat{\Delta}_{2,n}^{**'} \right)'$ and partitioning is after the q th row. The estimator $\hat{\beta}_n^{**}$

can be interpreted as a non-Gaussian counterpart of the 2SLS estimator of β , the latter being given by

$$\bar{\beta}_n = \frac{\bar{\Delta}'_{1,n} Q_{zz,n} \bar{\Delta}_{2,n}}{\bar{\Delta}'_{2,n} Q_{zz,n} \bar{\Delta}_{2,n}},$$

where $\bar{\Delta}_n = (\bar{\Delta}'_{1,n}, \bar{\Delta}'_{2,n})'$ and partitioning is after the q th row. The estimators $\hat{\beta}_n^{**}$ and $\bar{\beta}_n$ are both obtained by means of a generalized least squares (GLS) regression of an estimator of δ_1 onto an estimator of δ_2 (in (3)). The GLS regressions utilize identical weighting matrices, but differ in terms of the estimators of δ being employed, with $\hat{\beta}_n^{**}$ being based on an asymptotically efficient estimator (namely $\hat{\Delta}_n^{**}/\sqrt{n}$) and $\bar{\beta}_n$ being based on the OLS estimator $\bar{\Delta}_n/\sqrt{n}$.

If Assumptions 1-3 and 4SF hold, then

$$\sqrt{n}(\bar{\beta}_n - \beta) \rightarrow_d \mathcal{N}(0, \bar{\Sigma}_\beta), \quad \bar{\Sigma}_\beta = \left[\begin{pmatrix} 1 \\ -\beta \end{pmatrix}' \Omega \begin{pmatrix} 1 \\ -\beta \end{pmatrix} \right] (\pi' Q_{zz} \pi)^{-1}.$$

The next result, which follows from Theorem 6(c) and the delta method, gives the corresponding result for $\hat{\beta}_n^{**}$.

Corollary 8. *If Assumptions 1-3, 4SF, and 5-8 hold, then*

$$\sqrt{n}(\hat{\beta}_n^{**} - \beta) \rightarrow_d \mathcal{N}(0, \Sigma_\beta), \quad \Sigma_\beta = \left[\begin{pmatrix} 1 \\ -\beta \end{pmatrix}' \mathcal{I}^{-1} \begin{pmatrix} 1 \\ -\beta \end{pmatrix} \right] (\pi' Q_{zz} \pi)^{-1}.$$

Under normality the convergence result in Corollary 8 agrees with that for the 2SLS estimator of β (and its asymptotic equivalents, such as the limited information maximum likelihood (LIML) estimator and Fuller's (1977) modification thereof). With non-Gaussian errors, on the other hand, the estimator $\hat{\beta}_n^{**}$ compares favorably with $\bar{\beta}_n$ whenever the inequality $\mathcal{I}^{-1} \leq \Omega$ is strict.

The existence of estimators which outperform 2SLS for certain non-Gaussian error distributions has been known at least since Amemiya (1982) and Powell (1983). For the purposes of relating $\hat{\beta}_n^{**}$ to the two-stage least absolute deviations (2SLAD) and double 2SLAD (D2SLAD) estimators studied in those papers, define

$$\tilde{\beta}_n(\lambda_1, \lambda_2) = \frac{\hat{\Pi}_n(\lambda_1)' Q_{zz,n} \hat{\pi}_n(\lambda_2)}{\hat{\pi}_n(\lambda_2)' Q_{zz,n} \hat{\pi}_n(\lambda_2)}, \quad (\lambda_1, \lambda_2)' \in \mathbb{R}^2,$$

where

$$\hat{\Pi}_n(\lambda_1) = \lambda_1 \hat{\Pi}_n^{LAD} + (1 - \lambda_1) \hat{\Pi}_n^{OLS}, \quad \hat{\pi}_n(\lambda_2) = \lambda_2 \hat{\pi}_n^{LAD} + (1 - \lambda_2) \hat{\pi}_n^{OLS},$$

$$\left(\hat{\gamma}_1^{LAD}, \hat{\Pi}_n^{LAD} \right) = \arg \min_{(\gamma_1, \Pi)} \sum_{i=1}^n |y_{1i} - \gamma_1' x_i - \Pi' z_i|,$$

$$\left(\hat{\gamma}_2^{LAD}, \hat{\pi}_n^{LAD} \right) = \arg \min_{(\gamma_2, \pi)} \sum_{i=1}^n |y_{2i} - \gamma_2' x_i - \pi' z_i|.$$

In this notation $\tilde{\beta}_n(0, 0)$ is the 2SLS estimator, while nonzero pairs (λ_1, λ_2) give rise to estimators that are asymptotically distinct from the 2SLS estimator. The Bahadur representation of any $\tilde{\beta}_n(\lambda_1, \lambda_2)$ is readily obtained (by means of the delta method) from the Bahadur representations of $\hat{\Pi}_n^{LAD}$, $\hat{\Pi}_n^{OLS}$, $\hat{\pi}_n^{LAD}$, and $\hat{\pi}_n^{OLS}$. Utilizing these Bahadur representations it can be shown that $\tilde{\beta}_n(\lambda_1, 0)$ is asymptotically equivalent to the 2SLAD(λ_1) estimator and that $\tilde{\beta}_n(1, 1)$ is asymptotically equivalent to the D2SLAD estimator(s).

Because $(\hat{\Delta}_{1,n}^{**}/\sqrt{n}, \hat{\Delta}_{2,n}^{**}/\sqrt{n})$ is an asymptotically efficient estimator of (δ_1, δ_2) in (3), it compares favorably with $(\hat{\Pi}_n(\lambda_1), \hat{\pi}_n(\lambda_2))$ for any value of (λ_1, λ_2) . This superiority is inherited by $\hat{\beta}_n^{**}$, which compares favorably with all estimators of the form $\tilde{\beta}_n(\lambda_1, \lambda_2)$ (and their asymptotic equivalents, such as the 2SLAD and D2SLAD estimators). In fact, Theorems A.1 and A.2 can be used to show that $\hat{\beta}_n^{**}$ is an asymptotically efficient (i.e., best regular) estimator of β under strong identification.

As a consequence, one would expect the strong identification local asymptotic power properties of the tests based on \widehat{LR}_n^{**} and \widehat{LM}_n^{**} to be matched by those of the test which rejects when $\widehat{W}_n^{**} > \chi_\alpha^2(1)$, where

$$\widehat{W}_n^{**} = \frac{(\hat{\beta}_n^{**})^2}{\hat{\Sigma}_\beta^{**}/n}, \quad \hat{\Sigma}_\beta^{**} = \left[\begin{pmatrix} 1 \\ -\hat{\beta}_n^{**} \end{pmatrix}' \hat{\mathcal{I}}_n^{** - 1} \begin{pmatrix} 1 \\ -\hat{\beta}_n^{**} \end{pmatrix} \right] (\hat{\pi}_n' Q_{zz,n} \hat{\pi}_n)^{-1}.$$

The next result, which follows from Theorem 6(b) and the delta method, verifies that conjecture.

Corollary 9. *If Assumptions 1-3, 4SC, and 5-8 hold, then*

$$\widehat{W}_n^{**} = \frac{(S_n' Q_{zz}^{1/2} \pi)^2}{\pi' Q_{zz} \pi} + o_p(1) \rightarrow_d \chi^2(1; b^2 \mathcal{I}_{11.2} \pi' Q_{zz} \pi).$$

An attractive feature of \widehat{W}_n^{**} is that its ingredients, $\hat{\beta}_n^{**}$ and $\hat{\Sigma}_\beta^{**}$, can be combined in the usual way to form a Wald test of any null hypothesis regarding β , not just the null hypothesis that $\beta = 0$. This feature is particularly convenient when hypothesis tests are used to construct confidence intervals by inversion, as it implies that valid (indeed, optimal) confidence intervals are trivial to construct. Indeed, a confidence interval with asymptotic coverage probability $1 - \alpha$ is given by

$$\left(\hat{\beta}_n^{**} - \sqrt{\chi_\alpha^2(1) \frac{\hat{\Sigma}_\beta^{**}}{n}}, \hat{\beta}_n^{**} + \sqrt{\chi_\alpha^2(1) \frac{\hat{\Sigma}_\beta^{**}}{n}} \right).$$

It should be emphasized, however, that the displayed confidence interval is invalid (i.e., does not have asymptotic coverage probability $1 - \alpha$) under weak identification. As a consequence, while the computational simplicity of \widehat{W}_n^{**} makes it an attractive competitor to \widehat{LM}_n^{**} and \widehat{LR}_n^{**} under strong identification, the Wald statistic does not enjoy the robustness (and, in the case of \widehat{LR}_n^{**} , “near” optimality) properties under weak identification that Corollary 7(a) establishes on the part of \widehat{LM}_n^{**} and \widehat{LR}_n^{**} .

Remark. The LIMLK (i.e., LIML with known Ω) estimator of β is given by

$$\arg \min_\beta \frac{(1, -\beta) (\bar{\Delta}_{1,n}, \bar{\Delta}_{2,n}) Q_{zz,n} (\bar{\Delta}_{1,n}, \bar{\Delta}_{2,n})' (1, -\beta)'}{(1, -\beta) \Omega (1, -\beta)'}$$

This estimator is asymptotically equivalent to the 2SLS estimator $\bar{\beta}_n$ when identification is strong, but enjoys certain advantages over $\bar{\beta}_n$ when identification is weak (e.g., Staiger and Stock (1997)). Analogously, the following non-Gaussian counterpart of the LIMLK estimator of β is asymptotically equivalent (superior) to $\hat{\beta}_n^{**}$ under strong (weak) identification:

$$\arg \min_\beta \frac{(1, -\beta) (\hat{\Delta}_{1,n}^{**}, \hat{\Delta}_{2,n}^{**}) Q_{zz,n} (\hat{\Delta}_{1,n}^{**}, \hat{\Delta}_{2,n}^{**})' (1, -\beta)'}{(1, -\beta) \hat{\mathcal{I}}_n^{**^{-1}} (1, -\beta)'}$$

5. SIMULATIONS

This section presents the results of a simulation study investigating the finite-sample performance of the procedure considered in this paper. Although we primarily focus on power properties of the tests based on \widehat{AR}_n^{**} , \widehat{LM}_n^{**} , and \widehat{LR}_n^{**} , we also discuss the properties of the point estimator $\hat{\beta}_n^{**}$ under strong identification.

5.1. Model Setup. The data are generated by the model (2). Specifically, we set $x_i = 1$ and set q , the dimension of the instrumental variable, equal to 4. The instruments are randomly generated from a standard Gaussian distribution, demeaned, and then kept fixed throughout the experiment. For the errors we consider two different specifications, based on (i) the standard normal distribution and (ii) the t distribution with 3 degrees of freedom ($t(3)$), respectively. (The Fisher information for the location model generated by the $t(3)$ distribution is $2/3$, twice the inverse of the variance of the $t(3)$ distribution.). The probability densities associated with the distributions are depicted in Figure 1.

FIGURE 1 ABOUT HERE

We generate $2n$ independent (studentized) errors $\tilde{v}_i = (\tilde{v}_{1i}, \tilde{v}_{2i})'$ from each distribution and define

$$v_{1i} = \tilde{v}_{1i} \quad \text{and} \quad v_{2i} = \sqrt{1 - \rho^2} \tilde{v}_{2i} + \rho \tilde{v}_{1i},$$

hereby inducing a correlation of ρ between the errors v_{1i} and v_{2i} . Consistent with the previous discussion, we take $\beta_0 = 0$. The 4×1 vector π is given by

$$\pi = \iota \cdot \sqrt{\frac{\zeta}{\iota' Z' Z \iota}},$$

where ι is a 4×1 vector of ones, Z is the $n \times 4$ matrix of instruments, and ζ is the concentration parameter $\pi' Z' Z \pi / q$, which determines the “strength” of the instruments. For the simulations, we chose $n = 1,000$ as the sample size, $S = 5,000$ as the number of simulations, $\rho = 0.5$, and ζ taking on the values 1 and 10. In addition we chose $\alpha = 0.05$ for the size of our tests. (We obtained qualitatively similar results for other choices of n , S , ρ , and ζ , but omit these to conserve space.)

5.2. Implementation. The new procedures are compared to three benchmark procedures. The first of these is the Gaussian procedure constructed using a feasible version of the quasi-sufficient statistics (\bar{S}_n, \bar{T}_n) employing the OLS estimator

$$\hat{\Omega}_n^{OLS} = \frac{1}{n - p - q} \sum_{i=1}^n \hat{v}_i^{OLS} \hat{v}_i^{OLS'} = \frac{1}{n - 5} \sum_{i=1}^n \hat{v}_i^{OLS} \hat{v}_i^{OLS'}$$

of Ω , where \hat{v}_i^{OLS} are the OLS residuals. We will refer to this technique as “OLS” for simplicity.

As a second benchmark procedure we compute the Normal Scores Rank Tests introduced by Andrews and Soares (2007). We refer to this procedure as “RNK”

for brevity. These tests are seen to have superior power properties to those denoted OLS herein and are recommended by the authors based on both asymptotic and finite-sample results. However, based on our asymptotic results, our procedures are expected to have superior power properties over the corresponding RNK tests.

Finally, the third benchmark procedure utilizes an “oracle” version of $\hat{\Delta}_n^{**}$. Specifically, using the true ℓ instead of its estimate, we obtain

$$\hat{\Delta}_n^{MLE} = \begin{pmatrix} \sqrt{n}\hat{\Pi}_n^{OLS} \\ \sqrt{n}\hat{\pi}_n^{OLS} \end{pmatrix} + \left(\hat{\mathcal{I}}_\ell^{-1} \otimes Q_{zz,n}^{-1} \right) \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell(\hat{v}_i^{OLS}) \otimes z_i,$$

where

$$\hat{\mathcal{I}}_\ell = \frac{1}{n} \sum_{i=1}^n \ell(\hat{v}_i^{OLS}) \ell(\hat{v}_i^{OLS})'.$$

It should be noted that this is not a true “oracle” procedure in the sense that it uses the estimated error terms rather than their true values and also relies on an estimate of the information matrix. We include this additional benchmark in an effort to identify the effects on performance of using nonparametric estimates of the score function. Although a slight abuse of notation, we will refer to this technique as “MLE” for simplicity.

The (feasible) adaptive procedure based on $\hat{\Delta}_n^{**}$ is referred to as “ADP” for notational simplicity. This procedure is fully data-driven, but requires the additional choice of three parameters: the kernel k , the trimming parameter a , and the smoothing parameter h . For specificity we set k equal to a standard Gaussian kernel. Since our explicit goal is to explore the extent to which the asymptotic optimality properties of adaptive procedures are inherited at least partially in finite samples, we chose values of a and h that deliver tests with actual size close to nominal size in our simulations. (Power curves are easier to interpret and compare when competing tests have common size.) In our simulations, this was achieved by setting $a = 0$ and employing bandwidth estimators of the form $h_1 = c\sqrt{\hat{\omega}_{11}^{OLS}}/\sqrt[8]{n}$ and $h_2 = c\sqrt{\hat{\omega}_{22}^{OLS}}/\sqrt[8]{n}$, where h_1 and h_2 are the bandwidth choice for the first and second dimension of the nonparametric score estimator, respectively, and the constant c is chosen from a grid of possible values to obtain approximately correct empirical size. As expected, these grid values were sensitive to the true distribution of the error terms. In our simulations, we used $c = 0.65$ and $c = 0.55$ for the Gaussian and $t(3)$ model, respectively, as these values produced testing procedures with good size properties.

Remarks. (i) The size of the tests was found to be quite sensitive to the choice

of a and h , being a decreasing function (ranging from .20 to 0 over the range of values of a and h considered in our simulations) of each of these tuning parameters. The choice $a = 0$ violates Assumption 5(b), but was made for simplicity and concreteness because the qualitative results seemed to be more sensitive to the choice of h than to the choice of a . Regarding the choice of h , we experimented with a variety of procedures and specifications. In terms of procedures we considered both first-generation and second-generation bandwidth selection procedures for both univariate density and derivative estimation and bivariate density and derivative estimation (e.g., Ichimura and Todd (2007)). In terms of specifications, we considered a common bandwidth as well as different combinations of alternative bandwidths for densities and partial derivatives. Unfortunately, but unsurprisingly in light of previous Monte Carlo results on adaptive estimation in the univariate case (e.g., Steigerwald (1992)), our preliminary findings showed that these procedures have disappointing size properties for modest sample sizes. In the end we therefore opted for the above-mentioned procedure, which is a simple re-scaling of a rule of thumb choice for bivariate density estimation.

(ii) A possible explanation for the poor size performance of tests using existing bandwidth selection techniques is that these techniques are engineered to minimize a criterion function related to the (pointwise) distance between the estimated and true function and therefore do not necessarily lead to the optimal bandwidth choice for the adaptive procedure considered in this paper. Moreover to our knowledge the only results available concerning the choice of smoothing parameters for adaptive estimators are those of Linton and Xiao (2001), who derived a second-order approximation to the distribution of an adaptive regression estimator. Their results are particular to the case of univariate densities and are derived under distributional assumptions that are violated in the designs considered here.

5.3. Results. Figure 2 presents the power graphs for the AR, LM, and CLR tests for the case where the reduced form errors are generated from a Gaussian distribution.

FIGURE 2 ABOUT HERE

The strength of the instruments is equal to 1 and 10 in the first and second rows of graphs, respectively. In this particular case, the OLS and MLE estimators of the linear coefficients coincide, while the second-moment matrices are equal up to a constant multiple which converges to 1 with the sample size. As a consequence, the power curves of the tests based on these two procedures are virtually equivalent. The RNK tests also appear to reach the power curve generated by OLS and MLE. Because the adaptive procedures employ a nonparametric estimator of ℓ , we would expect them to have reduced finite sample power relative to the “oracle” procedures

and this does indeed seem to be the case. Nevertheless, the power loss is encouragingly small and the findings suggest that the ADP procedures can dominate the OLS and RNK procedures when the errors are non-Gaussian.

FIGURE 3 ABOUT HERE

Figure 3 presents the results for the case when the errors are generated from a non-Gaussian distribution, a $t(3)$ distribution in this case. Again, the first and second rows differ by the choice of the strength of the instruments. The results in this case are consistent with the theoretical predictions. The tests based on the MLE estimator have superior power relative to the test statistics based on the OLS estimator, while the test statistics based on the RNK procedure and ADP estimator have power curves which reside in between the other two. Moreover, tests based on ADP appear to (non-strictly) dominate the corresponding RNK tests. As expected, the MLE estimator delivers important power gains when compared to the RNK procedure. Presumably the difference between the MLE and ADP power curves can be attributed to the fact that ADP employs a nonparametric estimator of the nuisance parameter ℓ . In other words, the asymptotic theory probably overstates the extent to which departures from Gaussianity can be exploited in finite samples. On the other hand, the qualitative predictions of the asymptotic theory are borne out in the simulations insofar as Figure 3 clearly suggests that even in finite samples the ADP procedures can enjoy power advantages over the OLS procedures when the errors are non-Gaussian.

FIGURE 4 ABOUT HERE

Finally, in Figure 4 we present (kernel density estimators of) the sampling distributions of the estimators of β using each procedure when instruments are “strong.” The sampling distribution of the ADP estimator $\hat{\beta}_n^{**}$ is more concentrated than that of the “OLS” estimator $\bar{\beta}_n$ and less concentrated than that of the “oracle” estimator. This is also consistent with the theoretical predictions. (Similar results were obtained for the ADP estimator of the reduced-form coefficients. We omit the results to conserve space.)

In our view, the Monte Carlo results provide evidence in favor of the procedure(s) developed in this paper. The key potential drawback of the new procedure(s), which is common to all nonparametric procedures, is the fact that no firm guidance on the choice of the smoothing parameter is available. As discussed above, the lack of formal theory in this area led us to consider a simple rule-of-thumb procedure which includes a constant parameter that needs to be chosen and depends on the underlying design. Although this constant is arbitrary and set in advance for each design, we showed that if this constant is chosen so that the empirical size is approximately correct, then important power gains are realized by constructing test statistics based on the ADP

estimates of the sufficient statistics. This would appear to suggest that in practice a bootstrap procedure (in the spirit of Hsieh and Manski (1987), but targeted at the size properties of the test rather than the mean squared error of the estimator) is likely to produce tests with good size and power.

6. APPENDIX: PROOFS

The main results of the paper will follow from two facts, Theorems A.1 and A.2, about the model (3). Neither result is particularly surprising, but we have been unable to find statements of these results in the literature.

Theorem A.1 is an LAN result. To state it, let

$$\begin{aligned} \mathcal{L}_n(d, g) &= \sum_{i=1}^n \log f [y_{1i} - \gamma_{1n}(g_1)' x_i - \delta_{1n}(d_1)' z_i, y_{2i} - \gamma_{2n}(g_2)' x_i - \delta_{2n}(d_2)' z_i] \\ &\quad - \sum_{i=1}^n \log f [y_{1i} - \gamma_1' x_i - \delta_1' z_i, y_{2i} - \gamma_2' x_i - \delta_2' z_i] \end{aligned}$$

denote the log likelihood ratio function associated with the local reparameterization

$$\gamma = \begin{bmatrix} \gamma_{1n}(g_1) \\ \gamma_{2n}(g_2) \end{bmatrix} = \begin{bmatrix} \gamma_1 + g_1/\sqrt{n} \\ \gamma_2 + g_2/\sqrt{n} \end{bmatrix}, \quad \delta = \begin{bmatrix} \delta_{1n}(d_1) \\ \delta_{2n}(d_2) \end{bmatrix} = \begin{bmatrix} \delta_1 + d_1/\sqrt{n} \\ \delta_2 + d_2/\sqrt{n} \end{bmatrix},$$

let “ $o_{p_{\delta, \gamma}}(1)$ ” and “ $\rightarrow_{d_{\delta, \gamma}}$ ” be shorthand for “ $o_p(1)$ under the distributions associated with $(d, g) = (0, 0)$ ” and “ \rightarrow_d under the distributions associated with $(d, g) = (0, 0)$ ”, respectively, and let

$$\ell_i = \ell(y_{1i} - \gamma_1' x_i - \delta_1' z_i, y_{2i} - \gamma_2' x_i - \delta_2' z_i).$$

Theorem A.1. *Suppose (y_{1i}, y_{2i}) is generated by (3).*

(a) *If Assumptions 1(a) and 2 hold and d_n is a bounded sequence, then*

$$\mathcal{L}_n(d_n, 0) = \mathcal{L}_n^\delta(d_n) + o_{p_{\delta, \gamma}}(1),$$

where

$$\mathcal{L}_n^\delta(d_n) = d_n' (\mathcal{I} \otimes Q_{zz}) \Delta_n^\delta - \frac{1}{2} d_n' (\mathcal{I} \otimes Q_{zz}) d_n, \quad \Delta_n^\delta = (\mathcal{I}^{-1} \otimes Q_{zz}^{-1}) \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell_i \otimes z_i,$$

$$\Delta_n^\delta \rightarrow_{d_{\delta,\gamma}} \mathcal{N}(0, \mathcal{I}^{-1} \otimes Q_{zz}^{-1}).$$

(b) If, moreover, Assumptions 1(b) and 3 hold and g_n is a bounded sequence, then

$$\mathcal{L}_n(d_n, g_n) = \mathcal{L}_n^\delta(d_n) + \mathcal{L}_n^\gamma(g_n) + o_{p_{\delta,\gamma}}(1),$$

where

$$\begin{aligned} \mathcal{L}_n^\gamma(g_n) &= g_n' (\mathcal{I} \otimes Q_{xx}) \Delta_n^\gamma - \frac{1}{2} g_n' (\mathcal{I} \otimes Q_{xx}) g_n, & \Delta_n^\gamma &= (\mathcal{I}^{-1} \otimes Q_{xx}^{-1}) \frac{1}{\sqrt{n}} \sum_{i=1}^n \ell_i \otimes x_i, \\ \left(\begin{array}{c} \Delta_n^\delta \\ \Delta_n^\gamma \end{array} \right) &\rightarrow_{d_{\delta,\gamma}} \mathcal{N} \left[\left(\begin{array}{c} 0 \\ 0 \end{array} \right), \left(\begin{array}{cc} \mathcal{I}^{-1} \otimes Q_{zz}^{-1} & 0 \\ 0 & \mathcal{I}^{-1} \otimes Q_{xx}^{-1} \end{array} \right) \right]. \end{aligned}$$

Theorem A.2 is an adaptation result for one-step estimators of δ . Given initial estimators $\hat{\delta}_n = (\hat{\delta}'_{1n}, \hat{\delta}'_{2n})'$ and $\hat{\gamma}_n = (\hat{\gamma}'_{1n}, \hat{\gamma}'_{2n})'$ of δ and γ , let

$$\tilde{\delta}_n(\hat{\delta}_n, \hat{\gamma}_n) = \hat{\delta}_n + \frac{1}{\sqrt{n}} \hat{\Delta}_n^\delta(\hat{\delta}_n, \hat{\gamma}_n),$$

where

$$\begin{aligned} \hat{\Delta}_n^\delta(\hat{\delta}_n, \hat{\gamma}_n) &= \left[\tilde{\mathcal{I}}_n(\hat{\delta}_n, \hat{\gamma}_n)^{-1} \otimes Q_{zz,n}^{-1} \right] \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\ell}_n(\hat{v}_i) \otimes z_i, \\ \tilde{\mathcal{I}}_n(\hat{\delta}_n, \hat{\gamma}_n) &= n^{-1} \sum_{i=1}^n \hat{\ell}_n(\hat{v}_i) \hat{\ell}_n(\hat{v}_i)', & \hat{v}_i &= \begin{pmatrix} y_{1i} - \hat{\gamma}'_{1n} x_i - \hat{\delta}'_{1n} z_i \\ y_{2i} - \hat{\gamma}'_{2n} x_i - \hat{\delta}'_{2n} z_i \end{pmatrix}, \end{aligned}$$

and

$$\hat{\ell}_n(v) = -\frac{\partial \hat{f}_n(v) / \partial v}{\hat{f}_n(v) + a_n}, \quad \hat{f}_n(v) = \frac{1}{nh_n^2} \sum_{i=1}^n K\left(\frac{v - \hat{v}_i}{h_n}\right).$$

Theorem A.2. Suppose (y_{1i}, y_{2i}) is generated by (3). If Assumptions 1-3 and 5 hold, $(\hat{\delta}_n, \hat{\gamma}_n)$ is discrete, and $\sqrt{n}(\hat{\delta}_n - \delta, \hat{\gamma}_n - \gamma) = O_p(1)$, then

$$\tilde{\mathcal{I}}_n(\hat{\delta}_n, \hat{\gamma}_n) = \mathcal{I} + o_{p_{\delta, \gamma}}(1)$$

and

$$\sqrt{n} \left[\tilde{\delta}_n(\hat{\delta}_n, \hat{\gamma}_n) - \delta \right] = \Delta_n^\delta + o_{p_{\delta, \gamma}}(1).$$

Proof of Theorem 1. Apply Theorem A.1(a) with $\delta = 0$ and $d_n = \mu(\beta, c)$. ■

Proof of Theorems 2, 4, and 6. Theorems 2 and 4 (a) are special cases of Theorem 6 (a) and Theorem 4 (b) is a special case of Theorem 6 (b), so it suffices to prove Theorem 6.

Theorem 6 can be derived with the help of Theorem A.2 because

$$\hat{\Delta}_n^{**} = \sqrt{n} \tilde{\delta}_n(\hat{\delta}_n, \hat{\gamma}_n), \quad \hat{\mathcal{I}}_n^{**} = \tilde{\mathcal{I}}_n(\hat{\delta}_n, \hat{\gamma}_n),$$

where $\hat{\delta}_n = (\hat{\Pi}'_n, \hat{\pi}'_n)'$ and $\hat{\gamma}_n$ is as in the main text.

Proof of Theorem 6(a). If $c = 0$ in Assumption 4W, then the result can be obtained by applying Theorem A.2 with $\delta = (0', 0)'$. The result for $c \neq 0$ follows by the contiguity property implied by Theorem A.1(a).

Proof of Theorem 6(b). If $b = 0$ in Assumption 4SC, then the result can be obtained by applying Theorem A.2 with $\delta = (0', \pi)'$. The result for $b \neq 0$ follows by applying Theorem A.1(a) with $d_n = (b\pi', 0)'$ and using Le Cam's third lemma.

Proof of Theorem 6(c). Apply Theorem A.2 with $\delta = (\beta\pi', \pi)'$. ■

Proof of Theorem A.1. Define

$$R(v, \theta) = 2 \left[\sqrt{\frac{f(v - \theta)}{f(v)}} - 1 - \frac{1}{2} \theta' \ell(v) \right] 1[f(v) > 0], \quad v, \theta \in \mathbb{R}^2,$$

and

$$\bar{R}(\theta) = \frac{1}{4} \theta' \mathcal{I} \theta + \int_{\mathbb{R}^2} R(v, \theta) f(v) dv, \quad \theta \in \mathbb{R}^2.$$

If Assumption 2 holds, then

$$\sqrt{f(v-\theta)} - \sqrt{f(v)} = \frac{1}{2}\theta' \int_0^1 \ell(v-\theta t) \sqrt{f(v-\theta t)} dt, \quad \forall v, \theta \in \mathbb{R}^2$$

and for almost every $v \in \mathbb{R}^2$, \sqrt{f} is differentiable at v , with (total) derivative $-\frac{1}{2}\ell\sqrt{f}$. Using these facts and proceeding as in the proof of van der Vaart (1998, Lemma 7.6), it can be shown that if Assumption 2 holds, then

$$\lim_{\eta \downarrow 0} V(\eta) = 0, \quad V(\eta) = \sup_{\|\theta\| \leq \eta, \theta \neq 0} \|\theta\|^{-2} \int_{\mathbb{R}^2} R(v, \theta)^2 f(v) dv. \quad (5)$$

It follows from this result and Lemma 1 of Pollard (1997) that

$$\lim_{\eta \downarrow 0} \bar{V}(\eta) = 0, \quad \bar{V}(\eta) = \sup_{\|\theta\| \leq \eta, \theta \neq 0} \|\theta\|^{-2} \bar{R}(\theta). \quad (6)$$

The proofs of parts (a) and (b) are completely analogous, so to conserve space we only establish part (a). The log likelihood ratio $\mathcal{L}_n(d_n, 0)$ admits the expansion

$$\begin{aligned} \mathcal{L}_n(d_n, 0) &= d_n' (\mathcal{I} \otimes Q_{zz}) \Delta_n^\delta + \sum_{i=1}^n R_{i,n} \\ &\quad - \frac{1}{4} \sum_{i=1}^n \left[d_n' \frac{\ell_i \otimes z_i}{\sqrt{n}} + R_{i,n} \right]^2 (1 + \xi_{i,n}), \end{aligned}$$

where

$$R_{i,n} = R \left[\begin{pmatrix} y_{1i} - \gamma_1' x_i - \delta_1' z_i \\ y_{2i} - \gamma_2' x_i - \delta_2' z_i \end{pmatrix}, \begin{pmatrix} d_{1n}' z_i / \sqrt{n} \\ d_{2n}' z_i / \sqrt{n} \end{pmatrix} \right], \quad \xi_{i,n} = \xi \left[d_n' \frac{\ell_i \otimes z_i}{\sqrt{n}} + R_{i,n} \right],$$

and the defining property of $\xi(\cdot)$ is $\log(1+t) = t - \frac{1}{2}t^2 [1 + \xi(2t)]$.

It suffices to show that the following conditions hold:

$$\sum_{i=1}^n R_{i,n} = -\frac{1}{4} d_n' (\mathcal{I} \otimes Q_{zz}) d_n + o_{p_{\delta, \gamma}}(1), \quad (7)$$

$$\max_{1 \leq i \leq n} |\xi_{i,n}| = o_{p_{\delta, \gamma}}(1), \quad (8)$$

$$\sum_{i=1}^n \left[d_n' \frac{\ell_i \otimes z_i}{\sqrt{n}} + R_{i,n} \right]^2 = d_n' (\mathcal{I} \otimes Q_{zz}) d_n + o_{p_{\delta, \gamma}}(1). \quad (9)$$

To do so, suppose $(d, g) = (0, 0)$.

Proof of (7). The random variables $R_{1,n}, \dots, R_{n,n}$ are independent and satisfy

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}(R_{i,n}^2) &\leq \sum_{i=1}^n V \left(\left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \right) \left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \\ &\leq \max_{1 \leq i \leq n} V \left(\left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \right) \sum_{i=1}^n \left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \\ &= o(1) O(1) = o(1), \end{aligned}$$

where the penultimate equality uses (5) and Assumption 1(a). As a consequence,

$$\sum_{i=1}^n R_{i,n} = \sum_{i=1}^n \mathbb{E}(R_{i,n}) + o_p(1),$$

where

$$\sum_{i=1}^n \mathbb{E}(R_{i,n}) = -\frac{1}{4} d'_n (\mathcal{I} \otimes Q_{zz,n}) d_n + \sum_{i=1}^n \bar{R} \left[\begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right].$$

By Assumption 1(a),

$$d'_n (\mathcal{I} \otimes Q_{zz,n}) d_n = d'_n (\mathcal{I} \otimes Q_{zz}) d_n + o(1).$$

Moreover,

$$\begin{aligned} \left| \sum_{i=1}^n \bar{R} \left[\begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right] \right| &\leq \sum_{i=1}^n \left| \bar{R} \left[\begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right] \right| \\ &\leq \sum_{i=1}^n \bar{V} \left(\left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \right) \left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \\ &\leq \max_{1 \leq i \leq n} \bar{V} \left(\left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \right) \sum_{i=1}^n \left\| \begin{pmatrix} d'_{1n} z_i / \sqrt{n} \\ d'_{2n} z_i / \sqrt{n} \end{pmatrix} \right\|^2 \\ &= o(1) O(1) = o(1), \end{aligned}$$

where the penultimate equality uses (6) and Assumption 1(a).

Proof of (8). Because $\lim_{t \rightarrow 0} \xi(t) = 0$ (by Taylor's Theorem), the result follows from the fact that

$$\max_{1 \leq i \leq n} \left\| \frac{\ell_i \otimes z_i}{\sqrt{n}} \right\| = o_p(1)$$

and

$$\max_{1 \leq i \leq n} |R_{i,n}| \leq \sqrt{\sum_{i=1}^n R_{i,n}^2} = o_p(1),$$

where the first convergence result uses $\ell_i \sim i.i.d.$ $(0, \mathcal{I})$ and Assumption 1(a), while the second convergence result uses the relation $\mathbb{E}(\sum_{i=1}^n R_{i,n}^2) = o(1)$ established in the proof of (7).

Proof of (9). Because $\sum_{i=1}^n R_{i,n}^2 = o_p(1)$ and

$$\sum_{i=1}^n \left[d'_n \frac{\ell_i \otimes z_i}{\sqrt{n}} \right]^2 = d'_n \left(\frac{1}{n} \sum_{i=1}^n \ell_i \ell'_i \otimes z_i z'_i \right) d_n,$$

it suffices to show that

$$\frac{1}{n} \sum_{i=1}^n \ell_i \ell'_i \otimes z_i z'_i = \mathcal{I} \otimes Q_{zz} + o_p(1).$$

The latter result can be established using $\ell_i \sim i.i.d.$ $(0, \mathcal{I})$ and Assumption 1(a). ■

Proof of Theorem A.2. The proof uses Schick's (1987) approach.

First, it follows from Theorem A.1(b) and the properties of $(\hat{\delta}_n, \hat{\gamma}_n)$ that we may assume $(\hat{\delta}_n, \hat{\gamma}_n) = (\delta, \gamma)$. (This is so because Theorem 6.2 of Bickel (1982) can be used to verify that Condition A of Schick's (1987) Method 3 holds.) In other words, it suffices to show that

$$\check{\Delta}_n^\delta = [\check{\mathcal{I}}_n^{-1} \otimes Q_{zz,n}^{-1}] \frac{1}{\sqrt{n}} \sum_{i=1}^n \check{\ell}_n(v_i) \otimes z_i = \Delta_n^\delta + o_p(1) \quad (10)$$

and

$$\check{\mathcal{I}}_n = n^{-1} \sum_{i=1}^n \check{\ell}_n(v_i) \check{\ell}_n(v_i)' = \mathcal{I} + o_p(1), \quad (11)$$

where

$$\check{\ell}_n(v) = -\frac{\partial \check{f}_n(v) / \partial v}{\check{f}_n(v) + a_n}, \quad \check{f}_n(v) = \frac{1}{nh_n^2} \sum_{i=1}^n K\left(\frac{v - v_i}{h_n}\right).$$

To do so, let $\check{\ell}_{n,i}(\cdot)$ denote the leave-one-out version of $\check{\ell}_n(\cdot)$ given by

$$\check{\ell}_{n,i}(v) = -\frac{\partial \check{f}_{n,i}(v) / \partial v}{\check{f}_{n,i}(v) + a_n}, \quad \check{f}_{n,i}(v) = \check{f}_n(v) - \frac{1}{nh_n^2} \left[K\left(\frac{v - v_i}{h_n}\right) - K(0) \right].$$

It follows from (the proof of) Lemma 3.1 and Remark 3.2 of Schick (1987) that condition (10) is implied by condition (11), Assumptions 1(a) and 2, and the following conditions:

$$\mathbb{E} \left[\int_{\mathbb{R}^2} \|\check{\ell}_n(v) - \ell(v)\|^2 f(v) dv \right] = o(1), \quad (12)$$

$$\max_{1 \leq i \leq n} \mathbb{E} \left[\int_{\mathbb{R}^2} \|\check{\ell}_n(v) - \check{\ell}_{n,i}(v)\|^2 f(v) dv \right] = o\left(\frac{1}{n}\right). \quad (13)$$

Utilizing Assumptions 2 and 5 and proceeding as in Schick (1987, p. 100), it can be shown that

$$\int_{\mathbb{R}^2} \left\| -\frac{\partial f_n(v) / \partial v}{f_n(v) + a_n} - \ell(v) \right\|^2 f(v) dv = o(1), \quad (14)$$

where $f_n(v) = \int_{\mathbb{R}^2} f(v - h_n r) K(r) dr = \mathbb{E}[\check{f}_n(v)]$. It follows from this result that

$$\int_{\mathbb{R}^2} \left\| \frac{\partial f_n(v) / \partial v}{f_n(v) + a_n} \right\|^2 f(v) dv = O(1). \quad (15)$$

Now, using Assumptions 2 and 5, we have

$$\sup_{v \in \mathbb{R}^2} \mathbb{E} \left[\|\check{f}_n(v) - f_n(v)\|^2 \right] = O\left(\frac{1}{nh_n^2}\right)$$

and

$$\sup_{v \in \mathbb{R}^2} \mathbb{E} \left[\left\| \frac{\partial \check{f}_n(v)}{\partial v} - \frac{\partial f_n(v)}{\partial v} \right\|^2 \right] = O \left(\frac{1}{nh_n^4} \right).$$

Utilizing these facts, (15), and the decomposition

$$\frac{\frac{\partial \check{f}_n(v)}{\partial v}}{\check{f}_n(v) + a_n} - \frac{\frac{\partial f_n(v)}{\partial v}}{f_n(v) + a_n} = \frac{\frac{\partial f_n(v)}{\partial v} \check{f}_n(v) - f_n(v)}{f_n(v) + a_n} + \frac{\frac{\partial \check{f}_n(v)}{\partial v} - \frac{\partial f_n(v)}{\partial v}}{\check{f}_n(v) + a_n},$$

it is easily shown that

$$\int_{\mathbb{R}^2} \mathbb{E} \left[\left\| \frac{\frac{\partial \check{f}_n(v)}{\partial v}}{\check{f}_n(v) + a_n} - \frac{\frac{\partial f_n(v)}{\partial v}}{f_n(v) + a_n} \right\|^2 \right] f(v) dv = O \left(\frac{1}{na_n^2 h_n^4} \right) = o(1), \quad (16)$$

a result which can be combined with (14) to yield (12).

It follows from (15) – (16) that

$$\int_{\mathbb{R}^2} \mathbb{E} \left[\left\| \frac{\frac{\partial \check{f}_n(v)}{\partial v}}{\check{f}_n(v) + a_n} \right\|^2 \right] f(v) dv = O(1).$$

Utilizing this fact, Assumption 5, and the decomposition

$$\check{\ell}_n(v) - \check{\ell}_{n,i}(v) = -\frac{\frac{\partial \check{f}_n(v)}{\partial v} \check{f}_n(v) - \check{f}_{n,i}(v)}{\check{f}_n(v) + a_n} + \frac{\frac{\partial \check{f}_n(v)}{\partial v} - \frac{\partial \check{f}_{n,i}(v)}{\partial v}}{\check{f}_{n,i}(v) + a_n},$$

it is easily shown that (13) holds.

Finally, condition (11) holds because

$$\check{\mathcal{I}}_n = n^{-1} \sum_{i=1}^n \check{\ell}_{n,i}(v_i) \check{\ell}_{n,i}(v_i)' = n^{-1} \sum_{i=1}^n \ell_i \ell_i' + o_p(1) = \mathcal{I} + o_p(1),$$

where the first equality uses the fact that $\check{\ell}_{n,i}(v_i) = \check{\ell}_n(v_i)$ for each i and the second equality uses (14) and (16). ■

Remarks. (i) Conditions (12) and (13) are counterparts of Schick's (1987) conditions (3.2) and (3.6). No counterpart of Schick's (1987) condition (3.1) is needed because $n^{-1} \sum_{i=1}^n z_i \rightarrow 0$. Also, the present definition of $\check{\ell}_{n,i}$ ensures that $\check{\ell}_{n,i}(v_i) = \check{\ell}_n(v_i)$

for every i , implying in particular that the natural counterpart of Schick's (1987) condition (3.5) is satisfied.

(ii) With the possible exception of (14), all steps in the proof of Theorem A.2 remain valid if the condition $\sup_{r \in \mathbb{R}} |k'(r)|/k(r) < \infty$ of Assumption 5(a) is replaced by the condition $\int_{\mathbb{R}} k'(r)^2 dr < \infty$. The latter condition, which is implied by Assumption 5(a), is satisfied by the normal kernel. Furthermore, if the error density f is such that $\sup_{v \in \mathbb{R}^2} \|\dot{f}(v)\| < \infty$, then (14) is satisfied (for any kernel) provided $\overline{\lim}_{n \rightarrow \infty} h_n/a_n < \infty$. This is so because

$$\begin{aligned}
 & \int_{\mathbb{R}^2} \left\| -\frac{\partial f_n(v)/\partial v|_{r=v}}{f_n(v) + a_n} - \ell(v) \right\|^2 f(v) dv \\
 & \leq 2 \int_{\mathbb{S}_f} \left\| \frac{\partial f_n(v)/\partial v}{f_n(v) + a_n} \right\|^2 [\sqrt{f(v)} - \sqrt{f_n(v)}]^2 dv \\
 & \quad + 2 \int_{\mathbb{S}_f} \left\| \frac{\partial f_n(v)/\partial v}{f_n(v) + a_n} \sqrt{f_n(v)} - \frac{\dot{f}(v)}{f(v)} \sqrt{f(v)} \right\|^2 dv \\
 & = \left(\sup_{v \in \mathbb{R}^2} \|\dot{f}(v)\| \right)^2 o(h_n^2/a_n^2) + o(1),
 \end{aligned}$$

where $\mathbb{S}_f = \{v \in \mathbb{R}^2 : f(v) > 0\}$ and the last equality uses

$$\int_{\mathbb{S}_f} [\sqrt{f(v)} - \sqrt{f_n(v)}]^2 dv = o(h_n^2), \tag{17}$$

$$\int_{\mathbb{S}_f} \left\| \frac{\partial f_n(v)/\partial v}{f_n(v) + a_n} \sqrt{f_n(v)} - \frac{\dot{f}(v)}{f(v)} \sqrt{f(v)} \right\|^2 dv = o(1), \tag{18}$$

and the bound

$$\sup_{v \in \mathbb{R}^2} \left\| \frac{\partial f_n(v)/\partial v}{f_n(v) + a_n} \right\|^2 \leq \left(\sup_{v \in \mathbb{R}^2} \|\dot{f}(v)\| \right)^2 / a_n^2.$$

The result (17) can be shown by means of Proposition A.7 of Koul and Schick (1996), while (18) can be established using Vitali's theorem, the L^1 -continuity theorem, and arguments analogous to those used in the proof of Lemma 6.2 of Bickel (1982).

REFERENCES

- AMEMIYA, T. (1982): "Two Stage Least Absolute Deviations Estimators," *Econometrica*, 50, 689–711.
- ANDERSON, T. W., AND H. RUBIN (1949): "Estimation of the Parameters of a Single Equation in a Complete Set of Stochastic Equations," *Annals of Mathematical Statistics*, 20, 46–63.
- ANDREWS, D. W. K., AND V. MARMER (2007): "Exactly Distribution-Free Inference in Instrumental Variables Regression with Possibly Weak Instruments," *Journal of Econometrics*, 142, 183–200.
- ANDREWS, D. W. K., M. J. MOREIRA, AND J. H. STOCK (2006): "Optimal Two-sided Invariant Similar Tests for Instrumental Variables Regression," *Econometrica*, 74, 715–752.
- ANDREWS, D. W. K., AND G. SOARES (2007): "Rank Tests for Instrumental Variables Regression with Weak Instruments," *Econometric Theory*, 23, 1033–1082.
- ANDREWS, D. W. K., AND J. H. STOCK (2007): "Inference with Weak Instruments," in *Advances in Economics and Econometrics: Theory and Applications, Ninth World Congress, Volume III*, ed. by R. Blundell, W. K. Newey, and T. Persson. New York: Cambridge University Press, 122–173.
- BICKEL, P. J. (1982): "On Adaptive Estimation," *Annals of Statistics*, 10, 647–671.
- CHOI, S., W. J. HALL, AND A. SCHICK (1996): "Asymptotically Uniformly Most Powerful Tests in Parametric and Semiparametric Models," *Annals of Statistics*, 24, 841–861.
- COX, D. R., AND N. R. REID (1987): "Parameter Orthogonality and Approximate Conditional Inference (with Discussion)," *Journal of the Royal Statistical Society, Series B*, 49, 1–39.
- DUFOUR, J.-M. (2003): "Identification, Weak Instruments, and Statistical Inference in Econometrics," *Canadian Journal of Economics*, 36, 767–808.
- FULLER, W. A. (1977): "Some Properties of a Modification of the Limited Information Estimator," *Econometrica*, 45, 939–953.
- HSIEH, D. A., AND C. F. MANSKI (1987): "Monte Carlo Evidence on Adaptive Maximum Likelihood Estimation of a Regression," *Annals of Statistics*, 15, 541–551.

- ICHIMURA, H., AND P. E. TODD (2007): “Implementing Nonparametric and Semiparametric Estimators,” in *Handbook of Econometrics, Volume 6*, ed. by J. J. Heckman, and E. E. Leamer. New York: North Holland, 5369-5468.
- KLEIBERGEN, F. (2002): “Pivotal Statistics for Testing Structural Parameters in Instrumental Variables Regression,” *Econometrica*, 70, 1781–1803.
- KOUL, H. L., AND A. SCHICK (1996): “Adaptive Estimation in a Random Coefficient Autoregressive Model,” *Annals of Statistics*, 24, 1025–1052.
- LINTON, O., AND Z. XIAO (2001): “Second-Order Approximation for Adaptive Regression Estimators,” *Econometric Theory*, 17, 984–1024.
- MOREIRA, M. J. (2003): “A Conditional Likelihood Ratio Test for Structural Models,” *Econometrica*, 71, 1027–1048.
- POLLARD, D. (1997): “Another Look at Differentiability in Quadratic Mean,” in *Festschrift for Lucien Le Cam: Research Papers in Probability and Statistics*, ed. by D. Pollard, E. Torgersen, and G. L. Yang. New York: Springer-Verlag, 305-314.
- POWELL, J. L. (1983): “The Asymptotic Normality of Two-Stage Least Absolute Deviations Estimators,” *Econometrica*, 51, 1569–1575.
- SCHICK, A. (1987): “A Note on the Construction of Asymptotically Linear Estimators,” *Journal of Statistical Planning and Inference*, 16, 89–105.
- (1997): “Efficient Estimates in Linear and Nonlinear Regression with Heteroscedastic Errors,” *Journal of Statistical Planning and Inference*, 58, 371–387.
- STAIGER, D., AND J. H. STOCK (1997): “Instrumental Variables Estimation with Weak Instruments,” *Econometrica*, 65, 557–586.
- STEIGERWALD, D. G. (1992): “On the Finite Sample Behavior of Adaptive Estimators,” *Journal of Econometrics*, 54, 371–400.
- VAN DER VAART, A. W. (1998): *Asymptotic Statistics*. New York: Cambridge University Press.

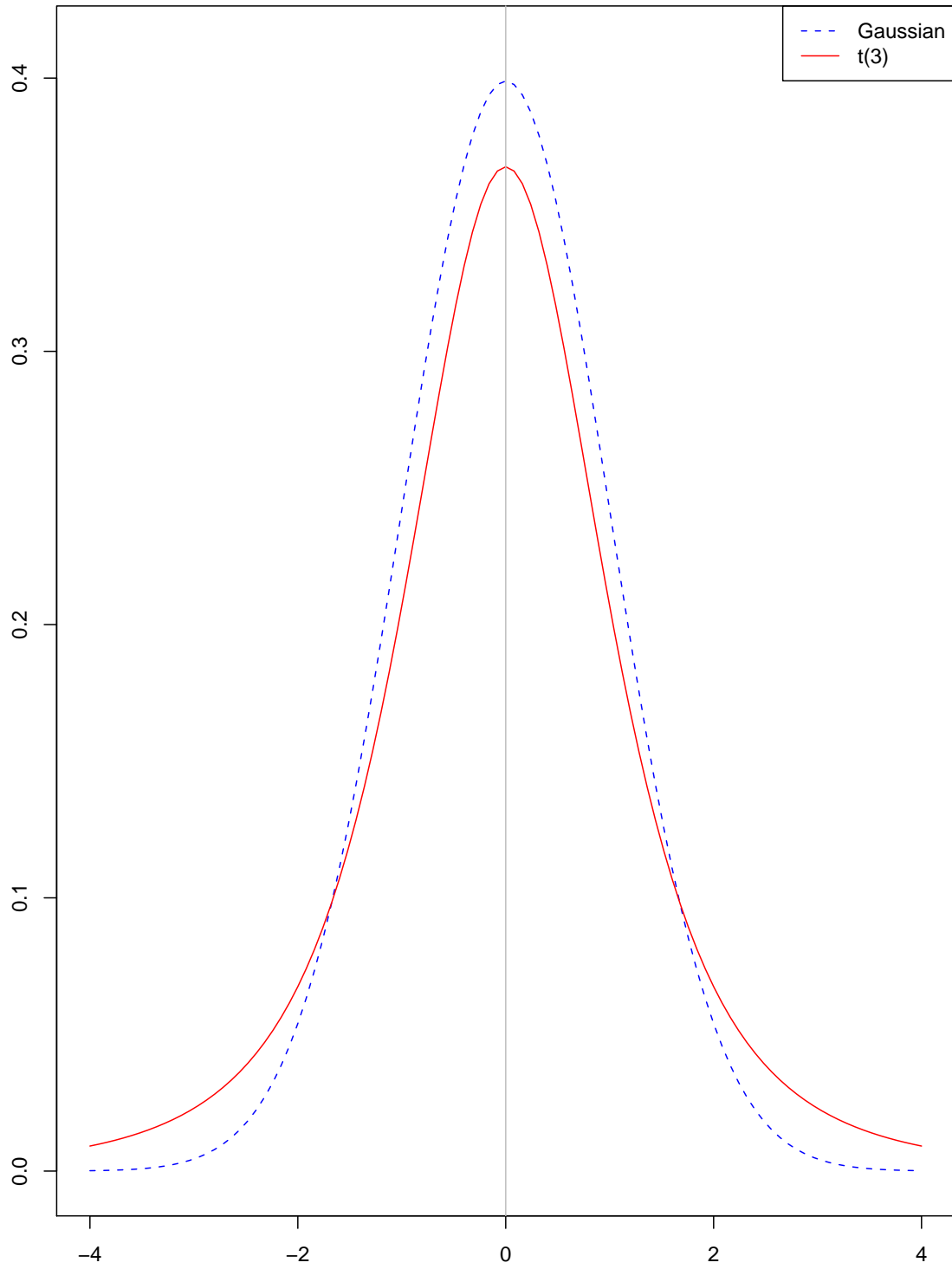


Figure 1: Probability Densities

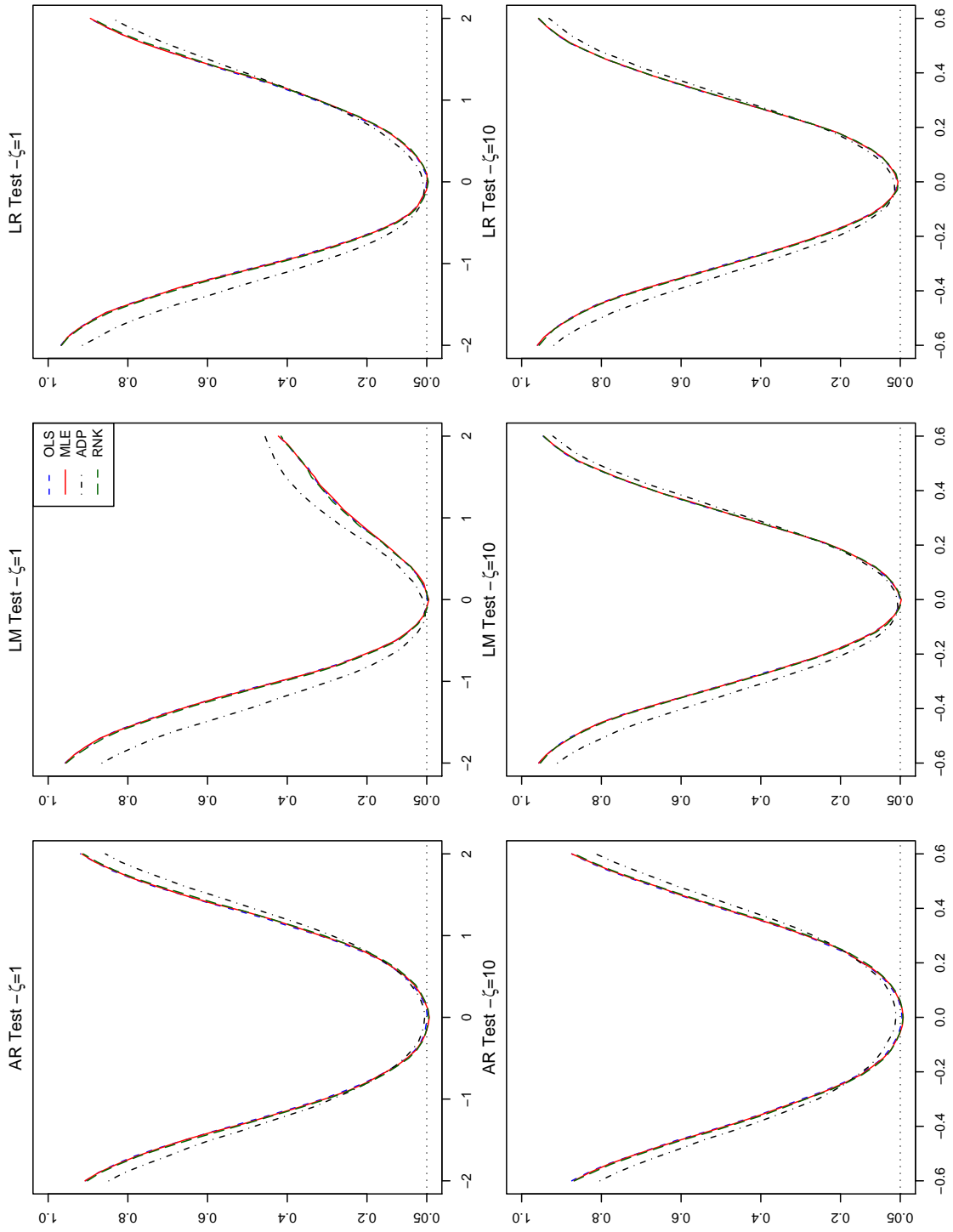


Figure 2: Power Curves, Gaussian Errors

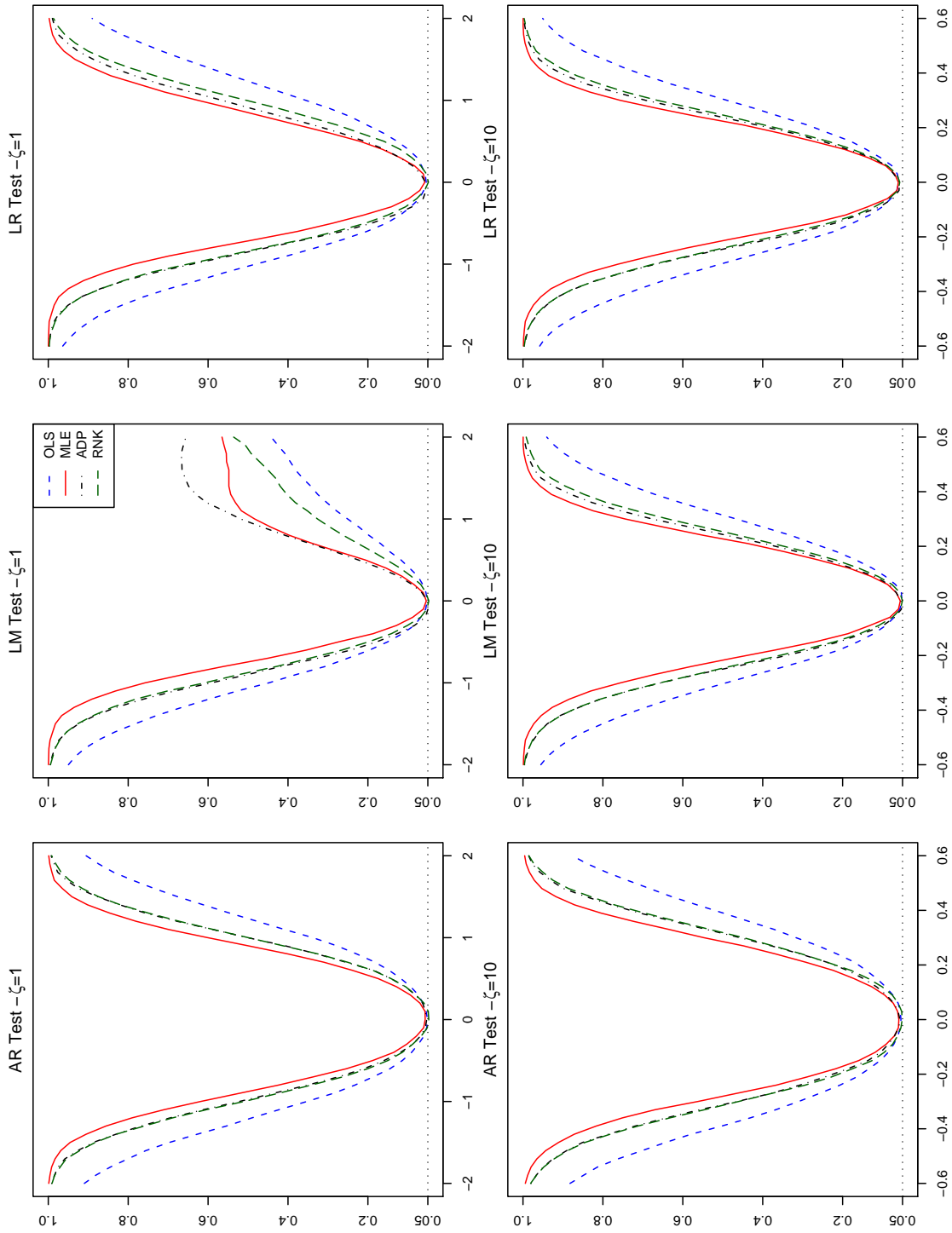


Figure 3: Power Curves, Non-Gaussian Errors

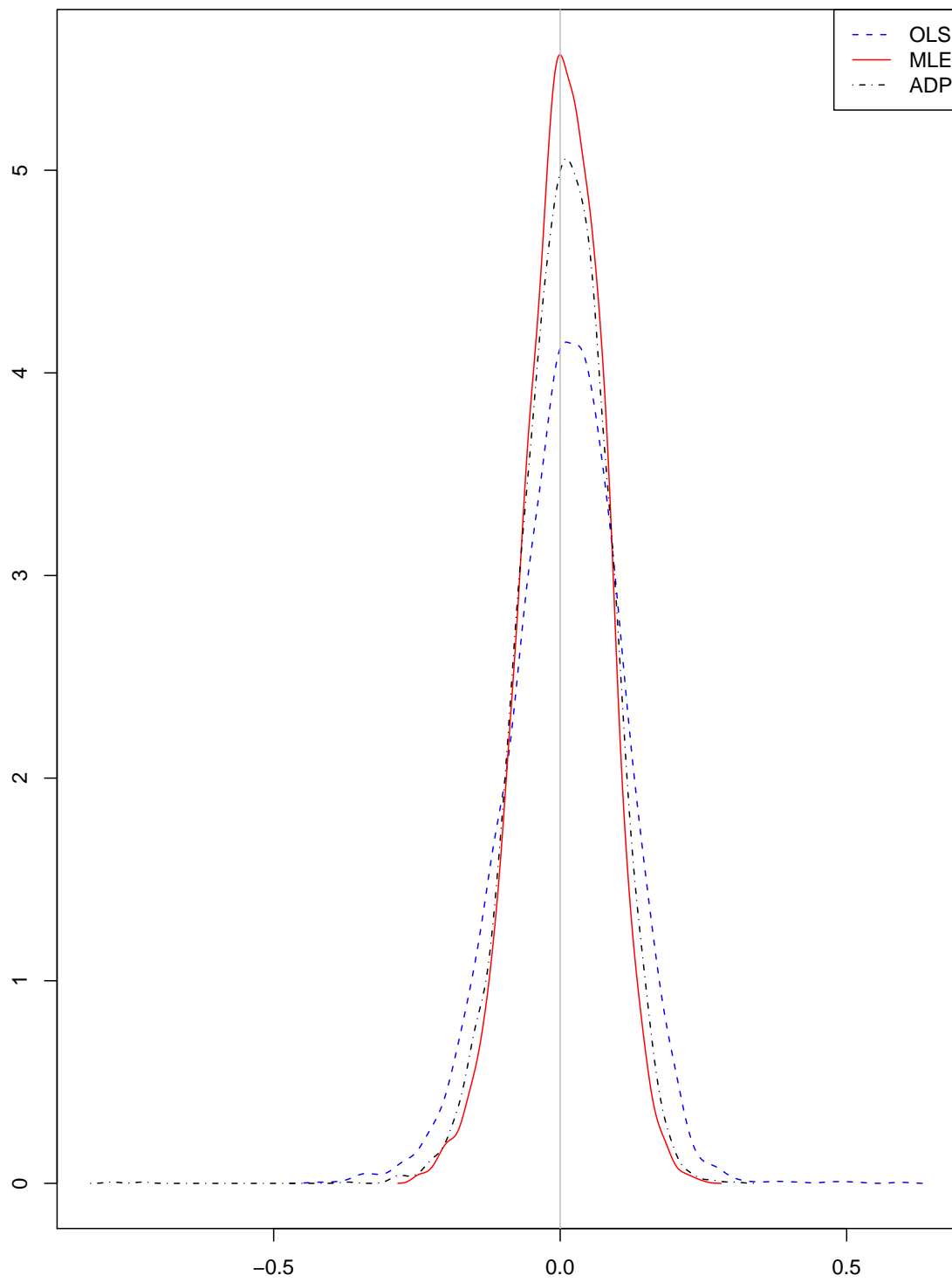


Figure 4: Estimators of β