

# TESTING UNDER WEAK IDENTIFICATION WITH CONDITIONAL MOMENT RESTRICTIONS\*

Sung Jae Jun and Joris Pinkse

Center for Auctions, Procurements, and Competition Policy  
Department of Economics  
The Pennsylvania State University

September 2011

We propose a semiparametric test for the value of coefficients in models with conditional moment restrictions which has correct size regardless of identification strength. The test is in essence an Anderson–Rubin (AR) test using nonparametrically estimated instruments to which we apply a standard error correction. We show that the test is (1) always size-correct, (2) consistent when identification is not too weak, and (3) asymptotically equivalent to an infeasible AR test when identification is sufficiently strong. We moreover prove that under homoskedasticity and strong identification our test has a limiting noncentral chi-square distribution under a sequence of local alternatives, where the noncentrality parameter is given by a quadratic form of the inverse of the semiparametric efficiency bound.

**Key Words:** Conditional Moment Restrictions, Weak Identification, Optimal Instruments, k-Nearest Neighbors

**JEL Classification Codes:** C12, C14, C21, C31

---

\*This paper is in part based on research supported by NSF grant SES-0922127. We thank seminar participants at Syracuse University, Columbia University, Arizona State University, the University of Arizona, The Pennsylvania State University, McGill University, and The Festschrift Conference for Tony Lancaster at Brown University, Yuichi Kiramura, Whitney Newey and four anonymous referees for their comments. We thank the Human Capital Foundation ([www.hcfoundation.ru](http://www.hcfoundation.ru)), and especially Andrey P. Vavilov, for their support of the Department of Economics, the Center for the Study of Auctions, Procurements, and Competition Policy (CAPCP, <http://capcp.psu.edu/>), and the Center for Research in International Financial and Energy Security (CRIFES, <http://crifes.psu.edu/>) at Penn State University. Sung Jae Jun is the Strumpf early career professor of economics. Joris Pinkse is an extramural fellow of the Center for Economic Research at Tilburg University.

## 1. INTRODUCTION

It is now well understood that standard asymptotic approximations can be misleading when instruments are *weak* (e.g. Phillips (1989); Staiger and Stock (1997); Stock, Wright, and Yogo (2002)). There is a substantial amount of work establishing limit distributions with alternative asymptotic variances of estimators in which the number of weak instruments is allowed to increase with the sample size, e.g. Newey (2004); Chao and Swanson (2005); Stock and Yogo (2005); Han and Phillips (2006); Newey and Windmeijer (2009). However, even  $t$ -tests based on such alternative asymptotic variances do not allow for correct inferences in the sense of Dufour (1997); see also Newey and Windmeijer (2009). There are testing-based inference procedures that are robust to weak identification but all of them are designed in a parametric setup (e.g. Anderson and Rubin, 1949; Kleibergen, 2002; Moreira, 2003; Kleibergen, 2005). In this paper we propose a semiparametric test for the value of coefficients in models with conditional moment restrictions that always has correct size. We show that the proposed test is consistent and has optimal local power under homoskedasticity when identification is sufficiently strong; i.e. under a sequence of local alternatives the test statistic then has a noncentral chi-square distribution whose noncentrality parameter is equal to a quadratic form of the inverse of the semiparametric efficiency bound.<sup>1</sup>

We begin our discussion with a simple example in which there is a single linear structural equation, a scalar parameter of interest  $\theta_0$  and no nuisance parameters, i.e.  $y_i = Y_i\theta_0 + u_i$ . In this example we have a vector of instruments  $z_i$  for which  $\mathbb{E}[u_i|z_i] = 0$  and  $\text{Var}[u_i|z_i] = \sigma^2$ . We use this simple example because it highlights the issues and differences across testing procedures with the least number of caveats and exceptions. The optimal instruments in this model are given by  $g_i = g(z_i) = \mathbb{E}[Y_i|z_i]$ , which is the derivative of  $M_i(\theta) = \mathbb{E}[Y_i\theta - y_i|z_i] = \mathbb{E}[Y_i|z_i](\theta - \theta_0)$ . Note that  $g$  is an unknown, but nonparametrically estimable, function of  $z_i$ .

We let  $g_i = \lambda\tilde{g}_i$  with  $0 < \mathbb{E}\tilde{g}_i^2 < \infty$ , where  $\lambda$  may depend on the sample size  $n$ . This construction is a technical but common tool to study and improve the behavior of statistics in finite samples, which was first introduced by Staiger and Stock (1997) who assumed that  $g(z) = z^\top\pi_0 = z^\top\tilde{\pi}_0\lambda$ . Staiger and Stock showed that when  $\lambda \preceq 1/\sqrt{n}$ , the two stage least squares (2SLS) estimator is inconsistent for  $\theta_0$  and its limiting distribution is nonnormal, where  $a \preceq b$  means that  $a/b \rightarrow C$  for some constant  $C$ ; when  $C = 0$ , we write  $a \prec b$ .<sup>2</sup> Given this fact, it is not surprising that standard semiparametric estimation theory also breaks down when  $\lambda$  decreases too fast.

Indeed, consider the semiparametric estimator  $\hat{\theta}$  based on  $k$  nearest neighbor (knn) estimates of  $g_i$ , which was also considered by Newey (1990);  $k$  is a smoothing parameter chosen by the researcher (see also Stone, 1977; Robinson, 1987; Newey, 1993). Although standard theory says that  $\hat{\theta}$  is asymptotically normal with asymptotic variance equal to the semiparametric efficiency bound, this occurs only when identification is sufficiently strong. For instance, if  $\lambda \preceq 1/\sqrt[4]{nk}$ , the

<sup>1</sup>We also propose a general statistic in the presence of heteroskedasticity of unknown form, and we show that the dominant term of the test statistic under the alternative has an expansion around the truth that is equal to a quadratic form of the inverse of the semiparametric efficiency bound.

<sup>2</sup>We will similarly use  $\succ, \succeq$  in this paper.

estimator  $\hat{\theta}$  ceases to be consistent for  $\theta_0$ . In fact, the limiting distribution implied by standard theory requires that  $\lambda \succ 1/\sqrt{k}$ . When  $1/\sqrt[4]{nk} \prec \lambda \preceq 1/\sqrt{k}$ ,  $\hat{\theta}$  is consistent and asymptotically normal but its asymptotic variance exceeds the one implied by standard theory; see example I in section 2.<sup>3</sup> This phenomenon is not specific to the knn method and similar conclusions obtain if other nonparametric estimators of  $g_i$  are used.<sup>4</sup>

The focus of this paper is on testing rather than estimation. As noted before, our objective is to obtain correct inferences in the sense of [Dufour \(1997\)](#) while maintaining performance in more favorable circumstances. More specifically, we propose a semiparametric test of the hypothesis  $H_0 : \theta_0 = \theta_H$ , which is always size-correct, is consistent when  $\lambda \succ 1/\sqrt[4]{nk}$ , and is equivalent in terms of local power to a  $t$ -test based on the efficient estimator  $\hat{\theta}$  when  $\lambda \succ 1/\sqrt{k}$ .

One way of dealing with a conditional moment restrictions model is to make a choice of instruments to obtain a set of unconditional moment conditions. When the conditional moment condition  $\mathbb{E}[Y_i\theta_0 - y_i|z_i] = 0$  a.s. is replaced with an unconditional one with a naïve choice of instruments, e.g.  $\mathbb{E}[z_i(Y_i\theta_0 - y_i)] = \mathbb{E}[z_i g_i](\theta_0 - \theta_0) = 0$ , two things can go wrong: (i)  $g_i$  is so small that the model is intrinsically difficult to estimate and/or (ii) the instruments  $z_i$  are chosen poorly in the sense that  $\mathbb{E}[z_i g_i]$  is close to zero. In both cases generalized methods of moments (GMM) estimators are known to behave poorly and  $t$ -tests based on them suffer from serious size distortions (e.g. [Dufour, 1997](#); [Stock and Wright, 2000](#)). Problem (ii) can be addressed by using the optimal instrument  $g_i$  instead of  $z_i$ , albeit that even unconditional moment conditions using optimal instruments are not generally equivalent to the original conditional moment restrictions ([Domínguez and Lobato, 2004](#), example 2).

There are several methods that address weak identification (caused by either (i) or (ii)) in the context of unconditional moment conditions models such as the Anderson–Rubin (1949, AR) procedure, the Kleibergen (2002, 2005; K) test, and the Moreira (2003, M) test.<sup>5</sup> [Andrews, Moreira, and Stock \(2004\)](#) found in a simulation study that the M-test approximates optimal average power in the above model with linear  $g$  and normally distributed errors which are independent of the instruments.

Although these testing procedures are robust to potential identification failure, all of them are fully parametric and designed for the case of a given choice of instruments. Therefore, when poor identification is due to a poor choice of instruments, inference based on these procedures can be improved. In particular, we will show that a robust test providing improved inference can be constructed by using nonparametrically estimated  $g_i$ .

<sup>3</sup>Note that as  $k$  has a rate that is arbitrarily close to  $n$ , both  $1/\sqrt[4]{nk}$  and  $1/\sqrt{k}$  are arbitrarily close to the corresponding parametric cut-off rate  $1/\sqrt{n}$ .

<sup>4</sup>In the case of series estimation similar conclusions can be deduced from the many-instrument literature; e.g. LIML estimators are known to be consistent and asymptotically normal with a nonstandard asymptotic variance when the number of instruments is allowed to increase slowly so that the concentration parameter increases sufficiently fast. See [Bekker \(1994\)](#); [Chao and Swanson \(2005\)](#); [Hausman, Newey, Woutersen, and Chao \(2007\)](#); [Newey and Windmeijer \(2009\)](#) for related work.

<sup>5</sup>See also [Staiger and Stock \(1997\)](#) for the AR-test and [Andrews, Moreira, and Stock \(2004, 2006\)](#) for the M-test.

The meaning of improved inference is twofold. Since the optimal instrument  $g_i$  is used, the problem due to a poor choice of instruments does not arise. Therefore, the proposed test may have nontrivial power even when other tests, not using optimal instruments, have no power. When identification is sufficiently strong, a further benefit of using optimal instruments is local power optimality. Indeed, the local power of the proposed test is equal to that of a  $t$ -test based on the efficient estimator  $\hat{\theta}$ , where the square of the mean of the local power distribution equals the inverse of the semiparametric efficiency bound (e.g. [Newey \(1993\)](#)) times the square of the local drift parameter.

The above discussion can be extended in several directions. When the parameter of interest  $\theta_0$  is vector-valued, a semiparametric chi-square statistic can be constructed without much difficulty. In this case, our test becomes comparable with a Wald test based on a semiparametric efficient estimator.<sup>6</sup> Further, nuisance parameters arising from the presence of additional regressors (endogenous or exogenous) whose coefficients are not of interest, can be partialled out without changing our results, provided that such nuisance parameters are identified if  $\theta_0$  is known (see also [Kleibergen, 2002, 2005](#)).

When the parameter of interest is scalar-valued (with or without nuisance parameters), the asymptotic null distribution of our test statistic is standard normal. Consequently, testing one-sided alternatives is straightforward. [Andrews, Moreira, and Stock \(2004\)](#) pointed out that no one-sided version of the  $K$ -test is consistent and they proposed a one-sided modification of the  $M$ -test. They also showed that the power function of the one-sided version of the  $M$ -test is close to the Gaussian power envelope when errors are normal and  $g_i$  is indeed linear in  $z_i$ . However, the  $M$ -test has no such property if errors are nonnormal or  $g_i$  is nonlinear. Indeed, our simulation experiments show that the one-sided version of our test can have better power than the corresponding  $M$ -test when  $g_i$  is nonlinear. A further advantage of an asymptotic normal or chi-square distribution is that the usual critical values can be used, whereas the critical values for the  $M$ -test must be simulated.

The presence of heteroskedasticity of unknown form does not affect the asymptotic validity (i.e. correct size) or consistency of (the simple form of) our test but it causes it to lose local power optimality because  $g_i$  is then no longer the optimal instrument. This feature is comparable to the 2SLS estimator in a parametric setup, which is efficient only under homoskedasticity but (provided that identification is sufficiently strong) is consistent under heteroskedasticity, also; the same applies to the  $M$  test.<sup>7</sup> We hence generalize our results to the case of multiple equations with heteroskedasticity of unknown form in section 3.

As mentioned before, the number  $k$  is an input parameter. Having a sample-size-dependent input parameter is common for nonparametric procedures. Although we state rate conditions on  $k$ , these are of limited practical use in a single sample of finite size. In estimation problems, an

<sup>6</sup>In this paper we will refer to both the  $t$ -test and the Wald test based on an efficient estimator as the  $N$ -test.

<sup>7</sup>[Andrews, Moreira, and Stock \(2004\)](#) proposed a modified  $M$ -test that is designed to be robust to heteroskedasticity. However, our simulation experiments show that the heteroskedasticity-robust version of the  $M$ -test tends to have size distortions in small samples when identification is weak and many instruments are used.

‘optimal’ choice of  $k$  can be motivated by minimizing a criterion like the (integrated) mean square error. With hypothesis testing, it is unclear what optimality criterion one would use. Indeed, the choice of  $k$  that minimizes estimation error of the nonparametric nuisance function is unlikely to equal the value optimizing test performance, however defined. Fortunately, we have found the performance of the test to be relatively insensitive to the choice of  $k$  in extensive Monte Carlo experiments.<sup>8</sup>

Finally, knn is not the only nonparametric method available, but it is attractive in this context because knn estimators have smaller variances in the sparse density area of  $z_i$  than e.g. kernel estimators, thereby obviating the need for an unnatural trimming procedure and (in settings like ours) knn estimators require minimal conditions.

The paper is organized as follows. In section 2 we propose a simple statistic for the single equation case that does not require estimation of conditional variances. The size, power, and local power properties of the test are carefully discussed. Section 3 extends the discussion to a more general case with multiple equations, where we use optimal instruments under heteroskedasticity. In section 4 we report the results of our Monte Carlo experiments. Section 5 concludes.

## 2. INFERENCE WITH A SINGLE EQUATION AND WEAK IDENTIFICATION

**2.1. Setup and Motivation.** We consider the following model:

$$\mathbb{E}[m(\omega_i, \theta_0)|z_i] = \mathbb{E}[m_i(\theta_0)|z_i] = 0 \quad a.s., \quad (1)$$

where  $m : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$  is twice continuously differentiable in its second argument with  $\mathcal{X}$  the support of  $\omega_i$  and  $\Theta \subset \mathbb{R}^{d_\theta}$  the parameter space; so  $\theta_0$  is a  $d_\theta$ -dimensional parameter vector of interest.

Identification strength will be measured by  $\lambda$ , which is a nonincreasing (with the sample size  $n$ ) sequence of numbers such that for any  $\theta \in \Theta \in \mathbb{R}^{d_\theta}$ ,

$$M(z_i, \theta) = M_i(\theta) = \mathbb{E}[m_i(\theta)|z_i] = \lambda \tilde{M}(z_i, \theta) = \lambda \tilde{M}_i(\theta) \quad a.s., \quad (2)$$

where  $\tilde{M} : \mathcal{Z} \times \Theta \rightarrow \mathbb{R}$  is a function that does not depend on  $n$ , where  $\mathcal{Z}$  is the support of  $z_i$ . By letting  $\mathbb{E}[m_i(\theta)|z_i]$  vary with  $\lambda$  (and hence with  $n$ ), we follow the setup of e.g. [Staiger and Stock \(1997\)](#), albeit that Staiger and Stock’s focus was on the case  $\lambda \sim n^{-1/2}$ . Note that the above construction is artificial in that it is not intended to describe the true data generating process; moment conditions do not generally vary with the sample size in real applications. The sequence  $\lambda$  is introduced as a tool to study and improve on the properties of statistical procedures in finite samples by analyzing the asymptotic effect of allowing identification strength to deteriorate as the sample size increases.

The following example shows that semiparametric estimators are not consistent when  $\lambda$  decreases too fast, i.e. identification is too weak. An implication is that standard hypothesis tests

---

<sup>8</sup>Not all experiments are reported here.

based on such estimators are invalid. Although this is a simple example, we note that the problem applies generally.

**Example I.** Recall the simple model from the introduction, where  $m_i(\theta) = y_i - Y_i\theta$ ,  $M_i(\theta) = \mathbb{E}(Y_i|z_i)(\theta - \theta_0)$ , and  $g_i = \mathbb{E}(Y_i|z_i) = \lambda\tilde{g}_i$ .<sup>9</sup> Let  $\hat{\theta}_I = \sum_i g_i y_i / \sum_i g_i Y_i$  and  $\hat{\theta} = \sum_i \hat{g}_i y_i / \sum_i \hat{g}_i Y_i$ , where  $\hat{g}_i$  is a knn estimator with  $k \prec n$ .<sup>10</sup> So,  $\hat{\theta}_I$  and  $\hat{\theta}$  are the infeasible and feasible semiparametric estimators, respectively, which are known to be efficient under homoskedasticity and strong identification (see [Newey \(1990, 1993\)](#)). However, as  $\lambda$  is allowed to vary with  $n$ , we have

$$\begin{cases} \hat{\theta}_I - \theta_0 & \xrightarrow{d} \Psi_{\tilde{g}u} / (\tilde{C}\mathbb{E}\tilde{g}_i^2 + \Psi_{\tilde{g}v}), & \text{if } 0 < \sqrt{n}\lambda \rightarrow \tilde{C} < \infty, \\ \sqrt{n}\lambda(\hat{\theta}_I - \theta_0) & \xrightarrow{d} \Psi_{\tilde{g}u} / \mathbb{E}\tilde{g}_i^2, & \text{if } \sqrt{n}\lambda \rightarrow \infty, \end{cases} \quad (3)$$

$$\begin{cases} \hat{\theta} - \theta_0 & \xrightarrow{d} \Psi_{vu} / (\Psi_{vv} + \tilde{C}^2\mathbb{E}\tilde{g}_1^2), & \text{if } \sqrt[4]{nk}\lambda \rightarrow \tilde{C} < \infty, \\ \sqrt{nk}\lambda^2(\hat{\theta} - \theta_0) & \xrightarrow{d} \Psi_{vu} / \mathbb{E}\tilde{g}_1^2, & \text{if } 1/\sqrt[4]{nk} \prec \lambda \prec 1/\sqrt{k}, \\ \sqrt{n}\lambda(\hat{\theta} - \theta_0) & \xrightarrow{d} (\tilde{C}\Psi_{vu} + \Psi_{\tilde{g}u}) / \mathbb{E}\tilde{g}_1^2, & \text{if } 1/(\sqrt{k}\lambda) \rightarrow \tilde{C} < \infty, \end{cases} \quad (4)$$

where all  $\Psi$ 's together have a joint normal distribution with zero mean and finite variance.

An informal justification for the above results is in appendix G.<sup>11</sup> □

The behavior of the infeasible estimator as a function of  $\lambda$  is effectively the same as that of a parametric estimator such as 2SLS. The behavior of the feasible estimator as a function of  $\lambda$ , however, is highly unusual and the asymptotic distribution of the two estimators coincides only when  $\lambda \succ 1/\sqrt{k}$ . When  $1/\sqrt[4]{nk} \prec \lambda \leq 1/\sqrt{k}$ , the feasible estimator is consistent and asymptotically normal, but its variance is nonstandard. The feasible estimator loses consistency when  $\lambda \leq 1/\sqrt[4]{nk}$ , compared with  $\lambda \leq 1/\sqrt{n}$  for the infeasible version. The results for the feasible estimator above are similar to those obtained in [Chao and Swanson \(2005\)](#) for the many weak instrument case.

A GMM estimator based on ‘linear’ instruments  $z_i$  will be inconsistent both when  $\lambda \leq 1/\sqrt{n}$  and at all levels of  $\lambda$  when  $g_i$  is orthogonal to  $z_i$ ; to see this point, note that  $\mathbb{E}[z_i m_i(\theta)] = \mathbb{E}[z_i g_i](\theta - \theta_0)$ . Since no prior information is available for  $g_i$  there is no way of knowing ex-ante if one’s instruments are sufficiently strong for consistency to obtain. Note also that the asymptotic distribution of  $\hat{\theta}$  in example I becomes nonstandard for larger  $\lambda$  than with an estimator based on a naïve choice of instruments such as 2SLS; i.e. semiparametrics requires stronger identification. Therefore, finding a testing method which is robust to weak identification is more valuable when the first stage equation is left unspecified.

We now describe our approach here. Note first that the optimal instruments are given by  $p^*(z_i) = g_i(\theta_0)/\sigma_i^2(\theta_0)$ , where  $g_i(\theta_0) = \mathbb{E}[\partial m_i(\theta_0)/\partial \theta|z_i]$  and  $\sigma_i^2(\theta_0) = \text{Var}[m_i(\theta_0)|z_i]$ . Since the

<sup>9</sup>In the parametric analysis of [Staiger and Stock \(1997\)](#),  $g_i = z_i^\top \pi_0 = z_i^\top \tilde{\pi}_0/\sqrt{n}$  so that  $\lambda = 1/\sqrt{n}$  and  $\tilde{g}_i = z_i^\top \tilde{\pi}_0$ .

<sup>10</sup>For a formal definition of a knn estimator, please see e.g. (6)

<sup>11</sup>The justification uses results established and notation and choices used in (the proof of) theorem 1, so it is better to finish reading this section prior to exploring the justification.

conditional variance can be ignored under homoskedasticity,<sup>12</sup> we start with the moment condition

$$\mathbb{E}[g_i(\theta_0)m_i(\theta_0)] = 0. \quad (5)$$

Our semiparametric procedure tests the null hypothesis  $H_0 : \theta_0 = \theta_H$  based on a nonparametric estimator of the expectation in (5) evaluated at  $\theta_H$ . We establish that it has correct size irrespective of  $\lambda$  and discuss its power properties as a function of  $\lambda$ .

In the present section we ignore the conditional variance and focus on (5) because it leads to a simpler (yet optimal under homoskedasticity) procedure, postponing the vector-valued case with estimated variance matrix to section 3.

**2.2. The Test Statistic.** We assume i.i.d. data throughout. We begin by describing our knn method. We define the *nearest neighbor* weights  $w_{ij}$ 's as follows. Let  $k$  be such that  $n^{3/5+c_k} \prec k \prec n$  for some  $c_k > 0$ . For each observation  $i$ , observation  $j$  is assigned a nearest neighbor weight  $w_{ij}$ , which is positive only if  $z_j$  belongs to the  $k$  nearest neighbors of  $z_i$ , where  $j \neq i$  and the distance between  $z_i$  and  $z_j$  is measured by the Euclidean distance.<sup>13</sup> In the case of ties positive weights are assigned randomly among the ties. Further, the weights  $\{w_{ij}\}$  are chosen such that  $w_{ii} = 0$ ,  $\sum_{j=1}^n w_{ij} = 1$ , and for some fixed  $C_w^-, C_w$ , either  $w_{ij} = 0$  or  $0 < C_w^- \leq kw_{ij} \leq C_w$ .

One can for instance use uniform weights  $\{w_{ij}\}$ , where  $w_{ij} = 1/k$  for any observation  $j(\neq i)$  whose  $z_j$  belongs to the  $k$  closest observations to  $z_i$  in terms of the Euclidean norm. The condition  $w_{ii} = 0$  removes the bias arising from the correlation between the moment function and its derivative. Since the weights are chosen by the researcher and there are always weights that satisfy the above requirements, the conditions imposed on the weights are innocuous. For further discussion on nearest neighbor weights, see e.g. [Stone \(1977\)](#); [Robinson \(1987\)](#); [Jun and Pinkse \(2009\)](#).

To define our test statistic based on (5), we first estimate  $g_i(\theta)$  by the knn method; i.e.

$$\hat{g}_i(\theta) = \sum_{j=1}^n w_{ij} m_{\theta_i}^T(\theta), \quad (6)$$

where  $m_{\theta_i}(\theta) = \partial m_i(\theta) / \partial \theta^T$ . Our statistic for testing  $H_0 : \theta_0 = \theta_H$  is then given by  $S(\theta_H) = \|\mathbf{t}(\theta_H)\|^2$ ,<sup>14</sup> where

$$\mathbf{t}(\theta) = \hat{D}_s^{-1}(\theta) \hat{N}_s(\theta), \quad (7)$$

with (omitting the  $\theta$ -argument)

$$\begin{cases} \hat{D}_s^2 &= \sum_i \hat{g}_i \hat{g}_i^T m_i^2 - n^{-1} (\sum_i \hat{g}_i m_i) (\sum_i \hat{g}_i^T m_i) + \sum_{ij} w_{ij} w_{ji} m_{\theta_i}^T m_{\theta_j} m_i m_j, \\ \hat{N}_s &= \sum_{i=1}^n \hat{g}_i m_i. \end{cases}$$

<sup>12</sup>I.e.,  $\text{Var}[m_i(\theta_0)|z_i]$  does not depend on  $z_i$ .

<sup>13</sup>Alternatively, the Mahalanobis distance  $(z_i - z_j)^T (\sum_{s=1}^n z_s z_s^T)^{-1} (z_i - z_j)$  can be used to address scaling issues. Note however that taking the Mahalanobis transformation before applying the Euclidean distance is asymptotically irrelevant, albeit that we do not specifically address the issue of 'estimated' transformations.

<sup>14</sup>When  $d_\theta = 1$ ,  $\mathbf{t}(\theta_H)$  can be used as a statistic instead of  $S(\theta_H)$ , which can be useful for one-sided testing.

Note that  $n^{-1} \sum_i \hat{g}_i(\theta_H) m_i(\theta_H)$  is an estimator of  $\mathbb{E}[g_1(\theta_H) m_1(\theta_H)]$ , which equals zero under the null hypothesis. So,  $S(\theta_H)$  is a semiparametric version of the Anderson–Rubin statistic (e.g. [Anderson and Rubin, 1949](#); [Stock and Wright, 2000](#)). However, the statistic has an unusual feature;  $\hat{D}_S^2(\theta)$  has an additional correction term, which is asymptotically irrelevant if identification is strong.

**2.3. Size, Power, and Local Power Optimality.** We now discuss the asymptotic properties of the statistic  $S(\theta_H)$ . Note that

$$\sum_{i=1}^n \hat{g}_i(\theta_H) m_i(\theta_H) = \sum_{i=1}^n (\hat{g}_i(\theta_H) - g_i(\theta_H)) m_i(\theta_H) + \lambda \sum_{i=1}^n \tilde{g}_i(\theta_H) m_i(\theta_H), \quad (8)$$

where  $\tilde{g}_i(\theta) = \partial \tilde{M}_i(\theta) / \partial \theta$ . If  $\lambda \neq 0$  is fixed, the first term on the right hand side of (8) can be shown to be asymptotically negligible relative to the second term. However, if  $\lambda$  is very small, the estimation error in  $\hat{g}_i$  can dominate the second term, in which case the nonparametric estimation error must be taken into account. In what follows we study the behavior of  $t(\theta_H)$  as a function of  $\lambda$ , allowing for such estimation error.

Let  $u_i(\theta) = m_i(\theta) - M_i(\theta)$  and  $v_i(\theta) = m_{\theta_i}^\top(\theta) - g_i(\theta)$ . We make the following assumptions.

**Assumption A.**  $\mathbb{E}[\mathbb{E}(m_1^2(\theta_H)|z_1)\mathbb{E}(v_1^2(\theta_H)|z_1) - \{\mathbb{E}(m_1(\theta_H)v_1(\theta_H)|z_1)\}^2] > 0$ .<sup>15</sup>

If  $\text{Var}[m_1(\theta_H)|z_1 = z], \text{Var}[v_1(\theta_H)|z_1 = z]$  are bounded away from zero on the support of  $z_1$ , as is often assumed (see e.g. [Robinson, 1987](#)), then assumption A is implied by

$$\mathbb{P}[\text{Corr}(m_1(\theta_H), v_1(\theta_H)|z_1) = 1] < 1,$$

which is unlikely ever to be violated in empirical work.

**Assumption B.**  $\mathbb{E}[|\tilde{M}_1(\theta_H)|^4] < \infty, \mathbb{E}[\|\tilde{g}_1(\theta_H)\|^4] < \infty$ , and  $\mathbb{E}[\{\mathbb{E}(\|u_1(\theta_H)\|^4|z_1)\}^4] < \infty$ . Further,  $\mathbb{E}[\{\mathbb{E}(|v_1(\theta_H)|^4|z_1)\}^4] < \infty$  and  $\text{Var}[\tilde{g}_1(\theta_H)m_1(\theta_H)] > 0$ .

Assumption B excludes fat-tailed error distributions, which is a common regularity assumption in cross section analysis.

**Theorem 1.** *Suppose that assumptions A and B hold. Under the null of  $H_0 : \theta_0 = \theta_H$ ,  $t(\theta_H) \xrightarrow{d} N(0, I_{d_\theta})$  and hence  $S(\theta_H) \xrightarrow{d} \chi^2(d_\theta)$  regardless of the (nonincreasing)  $\lambda$ -sequence.*

Please note that we have not assumed homoskedasticity for theorem 1; it holds even under heteroskedasticity of unknown form.

Theorem 1 establishes our first and most important desideratum, i.e. asymptotically correct size regardless of the quality of identification. Given that our test has correct size, we now proceed with an analysis of its power.

<sup>15</sup>We implicitly assume here that only the conditional means vary with  $\lambda$  (and hence  $n$ ). One can allow the conditional (co)variances depend on  $\lambda$  by taking infimums over  $n$ .



We begin with a discussion of consistency against a fixed alternative  $H_1 : \theta_0 \neq \theta_H$ . For this we need to discuss the issue of identification since consistency cannot be obtained without (sufficiently strong) identification.

**Assumption C.**  $\mathbb{E}[\tilde{g}_1(\theta_H)\tilde{M}_1(\theta_H)] \neq 0 \Leftrightarrow \theta_H \neq \theta_0$ .

Assumption C means that  $\lambda$  provides a complete description of identification quality; without assumption C, consistency against  $H_1 : \theta_0 \neq \theta_H$  does not obtain even for fixed  $\lambda \neq 0$ . But consistency moreover requires sufficient identification quality, i.e.  $\lambda$  not going to zero too fast. For parametric weak identification–robust tests, consistency requires that  $\lambda \succ 1/\sqrt{n}$ . In the semiparametric case,  $\lambda \succ 1/\sqrt[4]{nk}$  is the best achievable rate. To see why this is so, consider

$$\sum_{i=1}^n \hat{g}_i(\theta_H)m_i(\theta_H) = \sum_{i=1}^n \hat{g}_i(\theta_H)u_i(\theta_H) + \lambda \sum_{i=1}^n (\hat{g}_i(\theta_H) - g_i(\theta_H))\tilde{M}_i(\theta_H) + \lambda^2 \sum_{i=1}^n \tilde{g}_i(\theta_H)\tilde{M}_i(\theta_H). \quad (9)$$

In the proof of theorem 1 we establish that the first term on the right hand side of (9) is  $O_p(\rho)$  with  $\rho = \sqrt{n}(\lambda + 1/\sqrt{k})$ . The middle term on the right hand side of (9) can be shown to be negligible relative to the other terms and the last term is  $O_p(n\lambda^2)$ . Because the first term on the right hand side of (9) has a limiting normal distribution (after renorming), the last right hand side term must dominate the first for consistency to obtain. Such dominance only occurs when  $n\lambda^2 \succ \sqrt{n/k}$ , i.e. when  $\lambda \succ 1/\sqrt[4]{nk}$ . This rate is exactly the cutoff rate for consistency of the semiparametric estimator discussed in example I, which is no coincidence; if a parameter cannot be estimated consistently it cannot be tested consistently.

In the parametric case with  $g_i = z_i^\top \tilde{\pi}_0 \lambda$  and  $\hat{g}_i = z_i^\top \hat{\pi}$  with  $\hat{\pi} = \lambda \tilde{\pi}_0 + O_p(1/\sqrt{n})$ , the first term on the right hand side of (9) becomes  $\hat{\pi}^\top \sum_i z_i u_i(\theta_H) = O_p(\sqrt{n}\lambda + 1)$  such that consistency requires that  $n\lambda^2 \succ \sqrt{n}\lambda + 1$  or  $\lambda \succ 1/\sqrt{n}$ . So the cause for the discrepancy is the fact that the nonparametric estimator of  $g_i$  converges more slowly than does the parametric one. It should be noted that the rate at which  $k$  increases can be chosen to be close to  $n$ , in which case the difference is minor.<sup>16</sup>

**Theorem 2.** *Let assumptions A to C hold and let  $\lambda \succ 1/\sqrt[4]{nk}$ . Under  $H_1 : \theta_0 \neq \theta_H$ ,*

*(i)  $\forall C < \infty : \lim_{n \rightarrow \infty} \mathbb{P}[\mathbf{S}(\theta) > C] = 1$ , (ii) if  $\lambda \succ 1/\sqrt{k}$ ,  $\mathbf{S}(\theta_H) = \|\mathcal{B}_n(\theta_H)\|^2 + o_p(n\lambda^2)$ , where  $\mathcal{B}_n(\theta_H) = \sqrt{n}\lambda(\text{Var}[\tilde{g}_1(\theta_H)m_1(\theta_H)])^{-1/2}\mathbb{E}[\tilde{g}_1(\theta_H)\tilde{M}_1(\theta_H)]$ .*

Although not stated in theorem 2 it should be clear that  $t(\theta_H)$  can be used for one–sided testing when  $d_\theta = 1$ . For instance, if  $\mathbb{E}[\tilde{g}_1(\theta)\tilde{M}_1(\theta)]$  is continuous in  $\theta$  and zero at  $\theta_0$  then, in view of assumption C, it is positive when  $\theta_H > \theta_0$  and negative when  $\theta_H < \theta_0$ .

Part (ii) of theorem 2 implies that if  $\lambda \succ 1/\sqrt{k}$  then  $\mathbf{S}(\theta_H)$  is asymptotically equivalent to an infeasible AR–test based on (5). Recall that the power of a traditional Wald test based on a semiparametric estimator (i.e. N–test) comes from  $\mathcal{B}_n^o(\theta_H) = \sqrt{n}\lambda(\theta_0 - \theta_H)\mathcal{B}$  where  $\mathcal{B} = \{\text{Var}[\tilde{g}_1(\theta_0)m_1(\theta_0)]\}^{-1/2}\mathbb{E}[\tilde{g}_1(\theta_0)\tilde{g}_1^\top(\theta_0)]$ , which is comparable but not generally equal to  $\mathcal{B}_n(\theta_H)$ .

<sup>16</sup>However, using  $k$  too large causes size distortions in practice although it helps power under strong identification.

Comparing  $\mathcal{B}_n(\theta_H)$  to  $\mathcal{B}_n^o(\theta_H)$  allows us to draw some conclusions about the local power of our test. For the sequence  $\theta_H = \theta_n = \theta_0 + \Delta/(\sqrt{n}\lambda)$ ,  $\mathcal{B}_n(\theta_n) = \mathcal{B}\Delta + o(1)$  while  $\mathcal{B}_n^o(\theta_n) = -\mathcal{B}\Delta$ . Therefore, when identification is sufficiently strong (i.e.  $\lambda \succ 1/\sqrt{k}$ ), our test has the same local power as the usual Wald test based on a semiparametric estimator using  $\hat{g}_i$ .<sup>17</sup> Theorem 3 formalizes this idea.

**Theorem 3.** Consider  $H_{1L} : \theta_H = \theta_n = \theta_0 + \Delta/(\sqrt{n}\lambda)$  and let assumptions A to C hold. Let further  $\lambda \succ 1/\sqrt{k}$  and  $k \prec n/(\log n)^4$ . Further, suppose that  $\mathbb{E}[\sup_{\mathcal{N}^*} |u_i(\theta)|^{p_u}]$ ,  $\mathbb{E}[\sup_{\mathcal{N}^*} |v_i(\theta)|^{p_v}] < \infty$  for some  $1/p_u + 1/p_v < 1/2$  in an open neighborhood  $\mathcal{N}^*$  of  $\theta_0$ . Then under  $H_{1L}$  we have

$$\mathbf{t}(\theta_n) \xrightarrow{d} N(\mathcal{B}\Delta, I_{d_\theta}). \quad (10)$$

Consequently,  $S(\theta_n)$  converges to a noncentral  $\chi^2$ -distribution with  $d_\theta$  degrees of freedom and noncentrality parameter  $\mathcal{B}\Delta$ ; note that  $\lambda\mathcal{B}^\top\mathcal{B}$  is the inverse of the semiparametric efficiency bound under homoskedasticity with fixed  $\lambda \neq 0$ .

The rate at which the local alternative approximates the null is  $1/(\sqrt{n}\lambda)$ , not the usual  $1/\sqrt{n}$ -rate. The  $1/(\sqrt{n}\lambda)$ -rate is not specific to our test but applies to all tests under weak identification; see e.g. Andrews, Moreira, and Stock (2004, 2006).

**2.4. Nuisance Parameters.** Up to now, we have only considered testing the entire parameter vector of the model. It is, however, more realistic to consider a model that contains both the parameter(s) of interest and a vector of nuisance parameters. We will write  $\theta_0$  and  $\beta_0$  for the subvector of interest and the nuisance parameter vector, respectively. We could for instance have  $m_i(\theta_0, \beta_0) = y_i - Y_i\theta_0 - x_i^\top\beta_0$ . Under additional conditions discussed below, all results derived earlier go through in the presence of nuisance parameters provided that the nuisance parameters can be consistently estimated under the null hypothesis.

We define  $\beta(\theta)$  as

$$\beta(\theta) = \underset{\beta}{\operatorname{argmin}} \mathbb{E}[\mathbb{E}\{m_1(\theta, \beta)|z_1\}]^2, \quad (11)$$

which is assumed to be unique. Other definitions of  $\beta(\theta)$  are conceivable and it is possible that  $\beta(\theta)$  is not uniquely determined by (11). In linear models this would occur when the number of good instruments is less than the total number of regressors minus  $d_\theta$ . If there are only  $d_\theta$  endogenous regressors with coefficient vector  $\theta$  then uniqueness of  $\beta(\theta)$  requires the absence of multicollinearity of the exogenous regressors. With more than  $d_\theta$  endogenous regressors the assumption of the availability of good instruments can be unrealistic and goes against the spirit of the weak identification literature, but there are sufficiently many applications in which it is warranted to allow for identification of nuisance parameters once the parameters of interest are known; see e.g. Kleibergen (2002, 2005). For nonlinear models conditions under which  $\beta(\theta)$  is uniquely defined are discussed at length in the GMM literature.

<sup>17</sup>The power of our test against a fixed alternative can be greater or less than that of the Wald test depending on the alternative.

Let  $h_i(\theta, \beta) = E[m_{\beta_i}^\top(\theta, \beta)|z_i]$  and let  $h_i(\theta) = h_i(\theta, \beta(\theta))$ . All other symbols that depend on both  $\theta$  and  $\beta$  are similarly defined. In particular,  $M_i(\theta) = M_i(\theta, \beta(\theta)) = \lambda \tilde{M}_i(\theta) = \lambda \tilde{M}_i(\theta, \beta(\theta))$ . Note that  $\beta(\theta)$  satisfies the first order condition

$$\mathbb{E}[h_1(\theta, \beta(\theta))m_1(\theta, \beta(\theta))] = 0, \quad (12)$$

which is natural since  $h_i(\theta_0)$  is exactly the vector of optimal instruments under homoskedasticity. In fact, using  $h_i$  as instruments is necessary to maintain equality of the local power of our test and the N-test.

We propose to estimate  $\beta(\theta)$  on the basis of the moment condition (12), which requires its identification.

**Assumption D.** For all  $\theta \in \Theta$  there is a unique  $\beta = \beta(\theta) \in \mathcal{B}$ , with  $\mathcal{B}$  compact, such that (12) is satisfied. Further, at this unique  $\beta = \beta(\theta)$  the Jacobian  $Q = Q(\theta) = \mathbb{E}[h_1(\theta)h_1^\top(\theta) + M_1(\theta)m_{\beta\beta 1}(\theta)]$  is invertible.

Note that the invertibility requirement is automatically fulfilled in sufficiently large samples if  $\lambda \prec 1$  provided that  $\mathbb{E}[h_1(\theta)h_1^\top(\theta)] > 0$  for all  $\theta \in \Theta$ .

There are situations in which  $h$  is known, e.g. when the structural equation is linear, in which case the nuisance parameters can be estimated parametrically. Here we focus on the more challenging case in which  $h_i$  itself is estimated. We define  $\hat{\beta}(\theta)$  as a solution to

$$\sum_i \hat{h}_i(\theta, \hat{\beta}(\theta))m_i(\theta, \hat{\beta}(\theta)) = 0, \quad (13)$$

where  $\hat{h}$  is a knn estimator of  $h$ . Under conditions to be outlined below, we will show that  $\hat{\beta}(\theta)$  is a consistent estimator of  $\beta(\theta)$ . In fact, under the null hypothesis  $\hat{\beta}(\theta) = \hat{\beta}(\theta_0)$  is a  $\sqrt{n}$ -consistent estimator achieving the semiparametric efficiency bound for the estimation of  $\beta(\theta_0)$  under homoskedasticity. Under the alternative hypothesis, however,  $\hat{\beta}(\theta)$  generally converges at a rate slower than  $\sqrt{n}$ . This is innocuous since it converges fast enough to ensure that the consistency properties of our test statistic are not affected provided that  $\lambda \succ 1/\sqrt[4]{nk}$ , as was assumed in theorem 2. The discrepancy in convergence rates under the null and alternative hypotheses is due to the fact that  $M_1(\theta_0) = \mathbb{E}[m_1(\theta_0, \beta(\theta_0))|z_1] = 0$  a.s. but  $M_1(\theta)$  is not generally zero for other values of  $\theta$ .<sup>18</sup> Since we maintain the setup of (2) here, the only exception arises when  $\lambda = 0$ .

The moment condition we tested previously, i.e. (5), can continue to be used here, albeit that now  $g_i(\theta) = g_i(\theta, \beta(\theta))$  and  $m_i(\theta) = m_i(\theta, \beta(\theta))$ . However, it is more convenient to replace (5) with an equivalent one from which the effect of the estimation of the nuisance parameters has been isolated, i.e. letting  $q_i(\theta) = q_i(\theta, \beta(\theta)) = g_i(\theta) - \kappa^\top(\theta)h_i(\theta)$ ,

$$\mathbb{E}[q_1(\theta_0)m_1(\theta_0)] = 0, \quad (14)$$

<sup>18</sup> $n^{-1} \sum_i (\hat{h}_i - h_i)m_i = o_p(n^{-1/2})$  if  $\mathbb{E}[m_i|z_i] = 0$  a.s. but this is not true if  $h_i$  is merely orthogonal to  $m_i$ .

where  $\kappa(\theta) = \kappa(\theta, \beta(\theta)) = (\mathbb{E}[h_1(\theta)h_1^\top(\theta)])^{-1}\mathbb{E}[h_1(\theta)g_1^\top(\theta)]$ . Let  $\hat{q}_i(\theta) = \hat{g}_i(\theta) - \kappa^\top(\theta)\hat{h}_i(\theta)$  and further  $\hat{q}_i(\theta) = \hat{g}_i(\theta) - \hat{\kappa}(\theta)^\top\hat{h}_i(\theta)$ , where  $\hat{g}_i(\theta) = \hat{g}_i(\theta, \hat{\beta}(\theta))$ ,  $\hat{h}_i(\theta) = \hat{h}_i(\theta, \hat{\beta}(\theta))$ , and  $\hat{\kappa}(\theta) = (\sum_{i=1}^n \hat{h}_i(\theta)\hat{h}_i^\top(\theta))^{-1}(\sum_{i=1}^n \hat{h}_i(\theta)\hat{g}_i^\top(\theta))$ . Condition (14) has the appealing feature that for  $\hat{m}_i(\theta) = m_i(\theta, \hat{\beta}(\theta))$  the difference between  $\sum_i \hat{q}_i(\theta)\hat{m}_i(\theta)$  and  $\sum_i \hat{q}_i(\theta)m_i(\theta)$  is asymptotically negligible. The transition from  $g$  to  $q$  is similar to ‘partialing out’ other regressors in a linear regression model.

Our test statistic then becomes

$$\hat{\mathbf{t}}(\theta) = \hat{\mathcal{D}}_N^{-1}(\theta) \sum_i \hat{q}_i(\theta)\hat{m}_i(\theta), \quad (15)$$

where  $\hat{m}_{\theta i}(\theta) = m_{\theta i}(\theta, \hat{\beta}(\theta))$  and (omitting the  $\theta$  argument)

$$\hat{\mathcal{D}}_N^2 = \sum_i \hat{q}_i \hat{q}_i^\top \hat{m}_i^2 - n^{-1} \left( \sum_i \hat{q}_i \hat{m}_i \right) \left( \sum_i \hat{q}_i^\top \hat{m}_i \right) + \sum_{ij} w_{ij} w_{ji} \hat{m}_j \hat{m}_{\theta j}^\top \hat{m}_i \hat{m}_{\theta i}.$$

The only difference between  $\mathbf{t}$  and  $\hat{\mathbf{t}}$  is that we now use  $q$  instead of  $g$  and that the right hand side quantities in (15) depend on  $\hat{\beta}$ . All assumptions made earlier will hence now be applied to  $q$  instead of  $g$ .

Previously  $g$ , being  $M$ 's derivative, varied proportionally to  $\lambda$ . Now the entire  $q$ -function varies with  $\lambda$ . To see this, note that  $M_{\theta i}(\theta) = g_i^\top(\theta) + h_i^\top(\theta)\beta_\theta(\theta)$ , and that by the implicit function theorem it follows from (12) that for  $Q(\theta) = E[h_1(\theta)h_1^\top(\theta) + M_1(\theta)m_{\beta\beta 1}(\theta)]$ ,

$$\begin{aligned} \beta_\theta(\theta) &= -(Q(\theta))^{-1} \mathbb{E}[h_1(\theta)g_1^\top(\theta) + M_1(\theta)m_{\beta\theta 1}(\theta)] \\ &= -\kappa(\theta) - (\mathbb{E}[h_1(\theta)h_1^\top(\theta)])^{-1} \mathbb{E}[M_1(\theta)(m_{\beta\theta 1}(\theta) - m_{\beta\beta 1}(\theta)\kappa(\theta))] + O(\lambda^2). \end{aligned}$$

Hence since  $M_{\theta i}(\theta) = q_i^\top(\theta) + h_i^\top(\theta)[\beta_\theta(\theta) + \kappa(\theta)]$ ,  $q_i$  varies proportionally to  $\lambda$  up to terms of order  $\lambda^2$ .

We are now in a position to state our assumptions and formulate our nuisance parameter theorem. Most assumptions made previously carry over albeit that they are now made with respect to  $q$  instead of  $g$ . The main implication of this shift is that assumption C requires (i) that there is a unique solution  $\beta(\theta)$  to (12), which was assumed in assumption D, and (ii) that there is a unique combination  $(\theta_0, \beta_0)$  that zeroes both  $\mathbb{E}[\hat{g}_1(\theta, \beta)m_1(\theta, \beta)]$  and  $\mathbb{E}[h_1(\theta, \beta)m_1(\theta, \beta)]$ . Therefore, (ii) requires that  $(\theta_0, \beta_0)$  are identified for fixed  $\lambda$  if  $g_i, h_i$  are used as instruments. The only additional assumptions we make here relate to the smoothness of  $m$ .

**Assumption E.** For  $f = m, m_\theta, m_\beta, m_{\theta\beta}$  we have  $\mathbb{E}[\mathbb{E}(\sup_{\beta \in \mathcal{B}} \|f_1(\theta, \beta)\|^2 | z_1)]^2 < \infty$ . Further,  $\mathbb{E}[\mathbb{E}(\sup_{\beta \in \mathcal{B}} \|m_{\beta\beta 1}(\theta, \beta)\|^2 | z_1)]^2 < \infty$ .

Assumption E is strong. It assumes the existence of at least two partial derivatives and moreover assumes that a uniform bound of these derivatives is finite in expectation. The assumption of the existence of derivatives excludes interesting applications such as quantile regression models. There are also nonpathological situations in which assumption E would be violated when the derivatives do exist; an example can be found in van der Vaart (1998, pp. 48–49).

**Theorem 4.** *If assumptions A to C hold when  $g$  is replaced with  $q$  and moreover assumptions D and E are satisfied then theorems 1 and 2 hold when  $g$  is replaced with  $q$  and  $t$  with  $\hat{t}$ .*

### 3. THE GENERAL CASE

In this section we provide a generalization of our test statistic in two respects: we allow for a vector-valued moment function  $m_i$ , as one would have in a demand/supply system, and we consider optimal instruments allowing for heteroskedasticity of unknown form using nonparametric estimation of the conditional variance function.

Consider the model in (1), where  $m_i$  now takes values in  $\mathbb{R}^{d_m}$  for  $d_m \geq 1$ . Then, the unconditional moment equations based on optimal instruments are given by

$$\mathbb{E}[G_i^\top(\theta_0)\mathcal{V}_i^{-1}(\theta_0)m_i(\theta_0)] = 0, \quad (16)$$

where  $G_i(\theta) = E[m_{\theta i}(\theta)|z_i]$  and  $\mathcal{V}_i(\theta) = \text{Var}[m_i(\theta)|z_i]$ .<sup>19</sup> A semiparametric version of the AR statistic can be constructed based on (16).

Let  $\hat{G}_i$  and  $\hat{\mathcal{V}}_i$  be the (knn) nonparametric estimators for  $G_i$  and  $\mathcal{V}_i$ , respectively; i.e.

$$\hat{G}_i(\theta) = \sum_{j=1}^n w_{ij}m_{\theta j}(\theta) \quad \text{and} \quad \hat{\mathcal{V}}_i(\theta) = \sum_{j=1}^n w_{ij}(m_j(\theta) - \hat{\mu}_j(\theta))(m_j(\theta) - \hat{\mu}_j(\theta))^\top,$$

where  $\hat{\mu}_i(\theta) = \sum_{j=1}^n w_{ij}m_j(\theta)$ . We then define a statistic  $T(\theta_H)$  for testing  $H_0 : \theta_0 = \theta_H$ , where

$$T(\theta) = \left( \sum_i \hat{G}_i^\top(\theta)\hat{\mathcal{V}}_i^{-1}(\theta)m_i(\theta) \right)^\top \hat{D}^{-2}(\theta) \left( \sum_i \hat{G}_i^\top(\theta)\hat{\mathcal{V}}_i^{-1}(\theta)m_i(\theta) \right), \quad (17)$$

where (omitting the  $\theta$ -argument)

$$\hat{D}^2 = \sum_i \hat{G}_i^\top \hat{\mathcal{V}}_i^{-1} m_i m_i^\top \hat{\mathcal{V}}_i^{-1} \hat{G}_i - n^{-1} \left( \sum_i \hat{G}_i^\top \hat{\mathcal{V}}_i^{-1} m_i \right) \left( \sum_i \hat{G}_i^\top \hat{\mathcal{V}}_i^{-1} m_i \right)^\top + \sum_{ij} w_{ij} w_{ji} m_{\theta j}^\top \hat{\mathcal{V}}_i^{-1} m_i m_j^\top \hat{\mathcal{V}}_j^{-1} m_{\theta i}.$$

$T(\theta_H)$  is similar to, but more complicated than,  $S(\theta_H)$  in that  $T(\theta_H)$  involves nonparametric estimation of conditional variances.

The asymptotic properties of  $T(\theta_H)$  are similar to those of  $S(\theta_H)$ , but the proofs are more complicated. Let  $u_i(\theta) = m_i(\theta) - M_i(\theta)$  as before and let  $V_i(\theta) = m_{\theta i}(\theta) - G_i(\theta)$ . Further, let  $\underline{e\mathbf{v}}, \overline{e\mathbf{v}}$  denote the smallest and largest eigenvalues of their arguments, respectively. We will use  $\text{tr}$  to denote the trace of a matrix. For any matrix  $A$ , let  $\|A\| = \sqrt{\overline{e\mathbf{v}}(A^\top A)}$ . The following assumptions correspond to assumptions A to C. The comments on assumptions A to C also apply to the following assumptions.

**Assumption AG.**  $\inf_z \underline{e\mathbf{v}}(\text{Var}[u_1(\theta_H)|z_1 = z]) > 0$  and for all  $c \in \mathbb{R}^{d_\theta}$  with  $\|c\| = 1$ ,

$$\mathbb{E} \text{tr} \left( \mathbb{E} \left\{ \text{Var}[u_1(\theta_H)|z_1] \text{Var}[V_1(\theta_H)c|z_1] - \text{Cov}[u_1(\theta_H), V_1(\theta_H)c|z_1] \text{Cov}[u_1(\theta_H), V_1(\theta_H)c|z_1]^\top \right\} \right) > 0.$$

<sup>19</sup>Although  $G_i^\top = g_i$  when  $m_i$  is scalar-valued, we use an upper-case letter to emphasize that the Jacobian is potentially matrix-valued in this section.

**Assumption BG.**  $\mathbb{E}[\|\tilde{M}_1(\theta_H)\|^4] < \infty$ ,  $\mathbb{E}[\|\tilde{G}_1(\theta_H)\|^4] < \infty$ , and  $\mathbb{E}[\{\mathbb{E}(\|u_1(\theta_H)\|^4|z_1)\}^4] < \infty$ . Further,  $\mathbb{E}[\{\mathbb{E}(\|V_1(\theta_H)\|^4|z_1)\}^4] < \infty$ , and  $\mathbb{E}[\tilde{G}_1^\top(\theta_H)\mathcal{V}_1^{-1}(\theta_H)\tilde{G}_1(\theta_H)] > 0$ .

**Assumption CG.**  $\mathbb{E}[\tilde{G}_1^\top(\theta_H)\mathcal{V}_1^{-1}(\theta_H)M_1(\theta_H)] \neq 0 \Leftrightarrow \theta_H \neq \theta_0$ .

**Theorem 5.** Suppose that assumptions AG and BG hold. Under  $H_0 : \theta_0 = \theta_H$ ,  $T(\theta_H) \xrightarrow{d} \chi^2(d_\theta)$  regardless of the (nonincreasing)  $\lambda$ -sequence.

**Theorem 6.** Let assumptions AG to CG hold and let  $\lambda \succ 1/\sqrt[4]{nk}$ . Under  $H_1 : \theta_0 \neq \theta_H$ , (i)  $\forall C < \infty : \lim_{n \rightarrow \infty} \mathbb{P}[T(\theta_H) > C] = 1$ , (ii) if  $\lambda \succ 1/\sqrt{k}$ ,  $T(\theta_H) = \|\mathcal{B}_n^*(\theta_H)\|^2 + o_p(n\lambda^2)$ , where  $\mathcal{B}_n^*(\theta) = \sqrt{n}\lambda \{\text{Var}[\tilde{G}_1^\top(\theta)\mathcal{V}_1^{-1}(\theta)m_1(\theta)]\}^{-1/2} \mathbb{E}[\tilde{G}_1^\top(\theta)\mathcal{V}_1^{-1}(\theta)\tilde{M}_1(\theta)]$ .

Theorem 5 shows that  $T(\theta_H)$  is robust to weak identification. Although it does not formally derive the limiting distribution of  $T(\theta_H)$  under local alternatives, part (ii) of theorem 6 shows that the dominant term of  $T(\theta_H)$  under the alternative is closely related with the semiparametric efficiency bound when identification is sufficiently strong. In particular, when  $\theta_H = \theta_n = \theta_0 + \Delta/(\sqrt{n}\lambda)$ , the Taylor expansion of  $\mathcal{B}_n^*(\theta_H)$  shows that  $\mathcal{B}_n^*(\theta_H) = \mathcal{B}^*\Delta + o(1)$ , where

$$\mathcal{B}^* = \{\text{Var}[\tilde{G}_1^\top(\theta_0)\mathcal{V}_1^{-1}(\theta_0)m_1(\theta_0)]\}^{-1/2} \mathbb{E}[\tilde{G}_1^\top(\theta_0)\mathcal{V}_1^{-1}(\theta_0)\tilde{G}_1(\theta_0)].$$

Note that  $\lambda\mathcal{B}^{*\top}\mathcal{B}^*$  is the inverse of the semiparametric efficiency bound of a conditional moment restrictions model with fixed  $\lambda \neq 0$ . We do not provide a full generalization of theorem 3 here.

#### 4. SIMULATIONS

We now compare several tests proposed in the literature with ours using simulation experiments. Because not all tests can be used when the structural equation is nonlinear, we use the simple model from the introduction, i.e.

$$\begin{cases} y_i = Y_i\theta_0 + u_i \\ Y_i = g(z_i) + v_i \end{cases}.$$

All eight instruments ( $z_i \in \mathbb{R}^8$ ) are independent standard normals. We consider the following data generating processes (DGP's). DGP1 was motivated by Heckman (1978).

**DGP1: Nonlinear IV:**  $\begin{cases} y_i = Y_i\theta_0 + u_i \\ Y_i = I\{\epsilon_i \leq \frac{1}{2} + \Phi(z_i^\top \iota)\lambda\} - \frac{1}{2}, \end{cases}$

where  $u_i = 5(\epsilon_i - \frac{1}{2}) + \eta_i$  with  $\epsilon_i$  and  $\eta_i$  drawn from standard uniform and normal distributions, respectively.<sup>20</sup> Note that the correlation between  $u_i$  and  $\epsilon_i$  is about 0.8, and hence  $Y_i$  is highly endogenous.  $Y_i$  is a binary variable taking values  $\pm 1/2$  with conditional mean (given  $z_i$ )

$$g(z) \equiv \mathbb{E}[Y_i|z_i = z] = \Phi(z_i^\top \iota)\lambda I\{\Phi(z_i^\top \iota)\lambda \leq 1/2\} + I\{\Phi(z_i^\top \iota)\lambda > 1/2\}/2,$$

which are the optimal instruments. For sufficiently small  $\lambda$  (i.e. less than 1/2), we note that

$$g(z) = \Phi(z_i^\top \iota)\lambda \rightarrow 0 \text{ as } \lambda \rightarrow 0,$$

<sup>20</sup> $\Phi(\cdot)$  is the distribution function of  $N(0,1)$ , and  $\iota$  is a vector of ones.

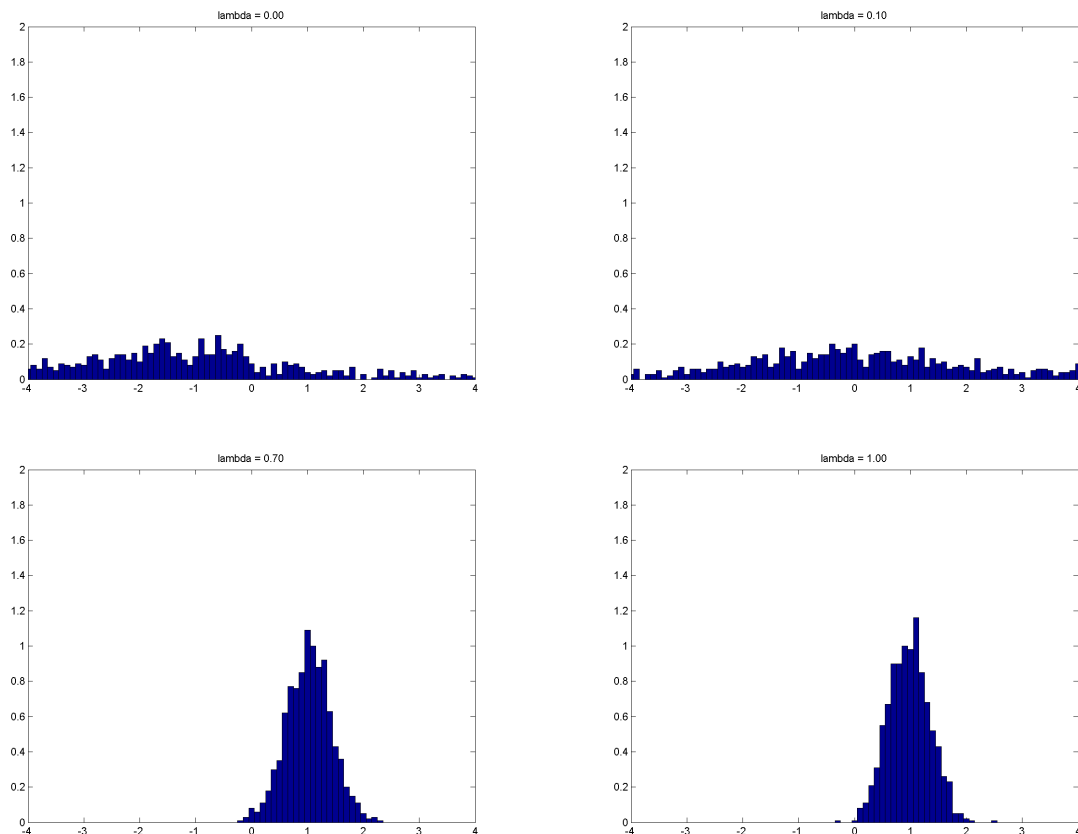


FIGURE 1. Histograms of the semiparametric (knn) estimator  $\hat{\theta}$  under DGP1

which causes identification failure. When  $\lambda$  is big, the nonlinearity of  $g$  makes using a linear probability model for the estimation of  $g$  inefficient.

**DGP2: Linear IV:** 
$$\begin{cases} y_i = Y_i\theta_0 + u_i \\ Y_i = z_i^\top \lambda + v_i, \end{cases}$$

where  $\lambda \geq 0$ , and  $u_i, v_i$  are drawn from the joint normal distribution with mean, variance, covariance equal to 0, 1, and 0.8, respectively. This is the case where the linear specification for the endogenous regressor is correct, which means that the ‘optimality’ results of Andrews, Moreira, and Stock (2004, 2006) apply.

Recall that the parameter  $\lambda$  indicates the strength of identification; we considered the following values for  $\lambda$ : 0, 0.05, 0.1, 0.3, 0.5, 0.7, 1 and 2. In all cases we test  $H_0 : \theta_0 = \theta_H$  versus  $H_1 : \theta_0 \neq \theta_H$  (or  $H_0 : \theta_0 = \theta_H$  versus  $H_1 : \theta_0 < \theta_H$  for one-sided alternatives). The primary goal is for the tests to have correct size properties irrespective of identification strength (as measured by  $\lambda$ ), but we also want them to have good power if identification is sufficiently strong. We used heteroskedasticity-robust versions of all statistics (except for the N-test) in the experiments; our test is still valid

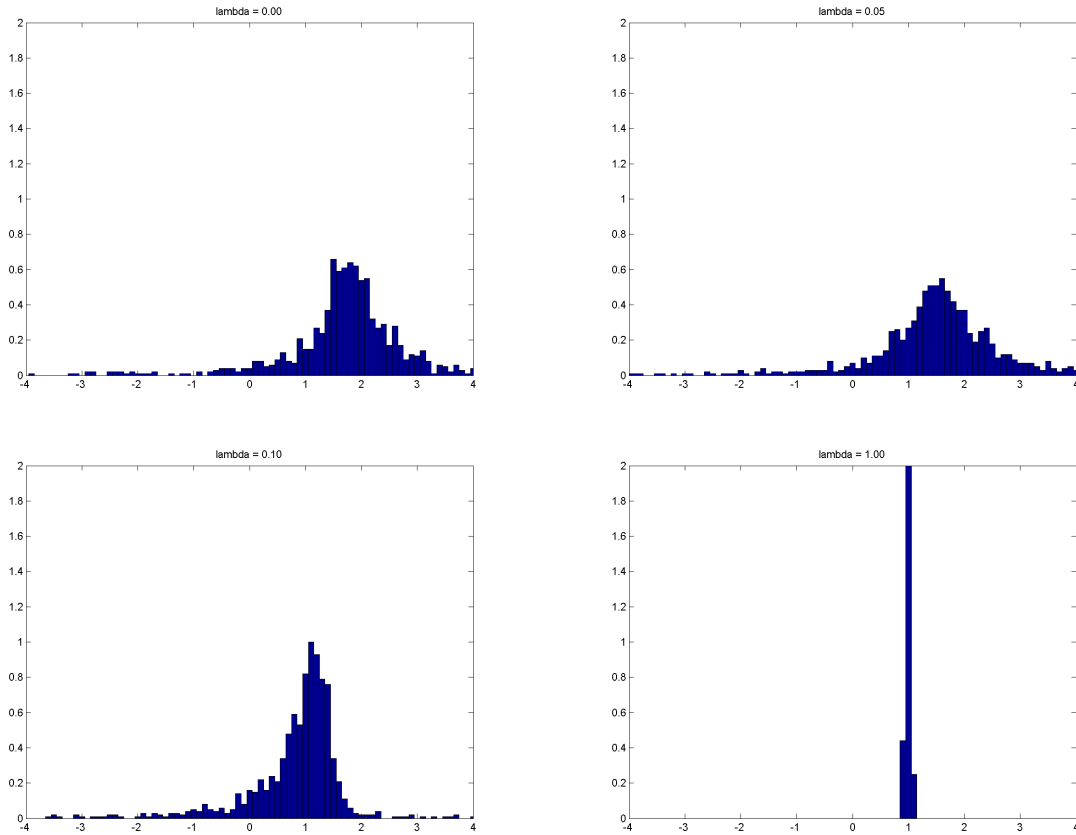


FIGURE 2. Histograms of the semiparametric (knn) estimator  $\hat{\theta}$  under DGP2

in the presence of heteroskedasticity.<sup>21</sup> In all experiments, we used  $n = 200, k = 70$ , and 1,000 replications.

We now illustrate our earlier finding that the semiparametric estimator  $\hat{\theta}$  does not behave as strong identification asymptotics would indicate when  $\lambda$  is small; the results are depicted in figures 1 and 2, which contain histograms of its distribution. For both DGP's, the histograms are not nicely concentrated around the true  $\theta_0 = 1$ , when the value of  $\lambda$  is small. This is not surprising in view of the inconsistency of  $\hat{\theta}$  when  $\lambda \leq 1/\sqrt[4]{nk}$ . We repeated this experiment for the HLIM estimator<sup>22</sup> of Hausman, Newey, Woutersen, and Chao (2007); the graphs (not shown) are similar to those in figures 1 and 2. See also figure 5.

Figures 3 and 4 show the rejection rates of a number of test statistics as a function of  $\theta_H$  for various values of  $\lambda$ ; figure 3 for DGP1 and figure 4 for DGP2. In all cases  $\theta_0 = 1$  and in all cases except for the N-test we used heteroskedasticity-robust versions of the test statistics. The nominal

<sup>21</sup>For the heteroskedasticity-robust version of the M-test, see Andrews, Moreira, and Stock (2004).

<sup>22</sup>Heteroskedasticity robust Limited Information Maximum Likelihood estimators. These are LIML type estimators but designed to be robust to heteroskedasticity by jackknife methods. They are consistent under "many-weak" instruments asymptotics that allows the concentration parameter to increase more slowly than  $n$ .



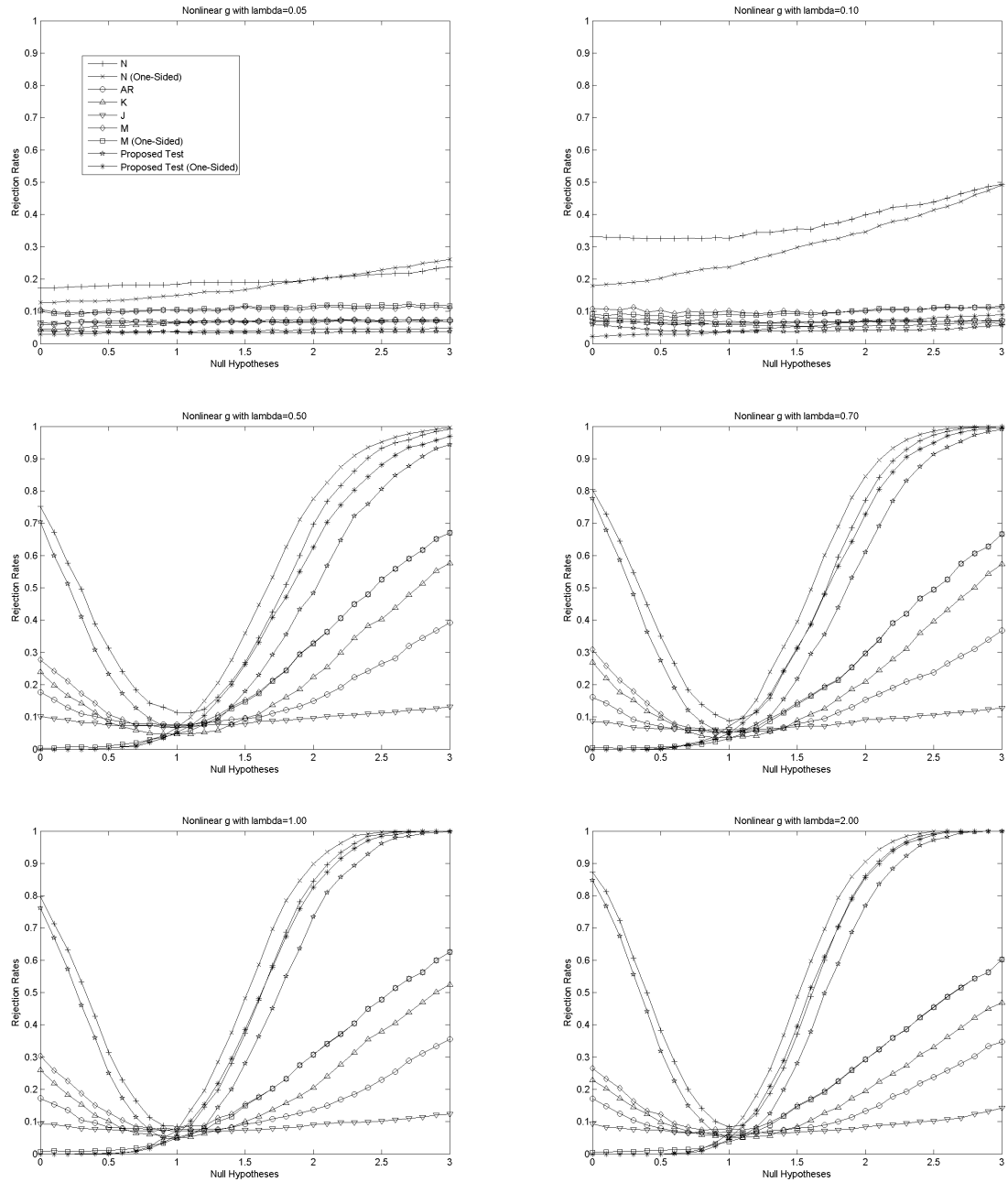


FIGURE 3. Monte Carlo rejection rates for nonlinear  $g$  (DGP1)

size of all tests is 0.05. As expected, the  $t$ -test based on an efficient point estimator suffers from serious size distortions, when the values of  $\lambda$  are small. But note that the power curves of the proposed tests are close to those of the  $t$ -tests based on the efficient estimator when  $\lambda$  is relatively large.

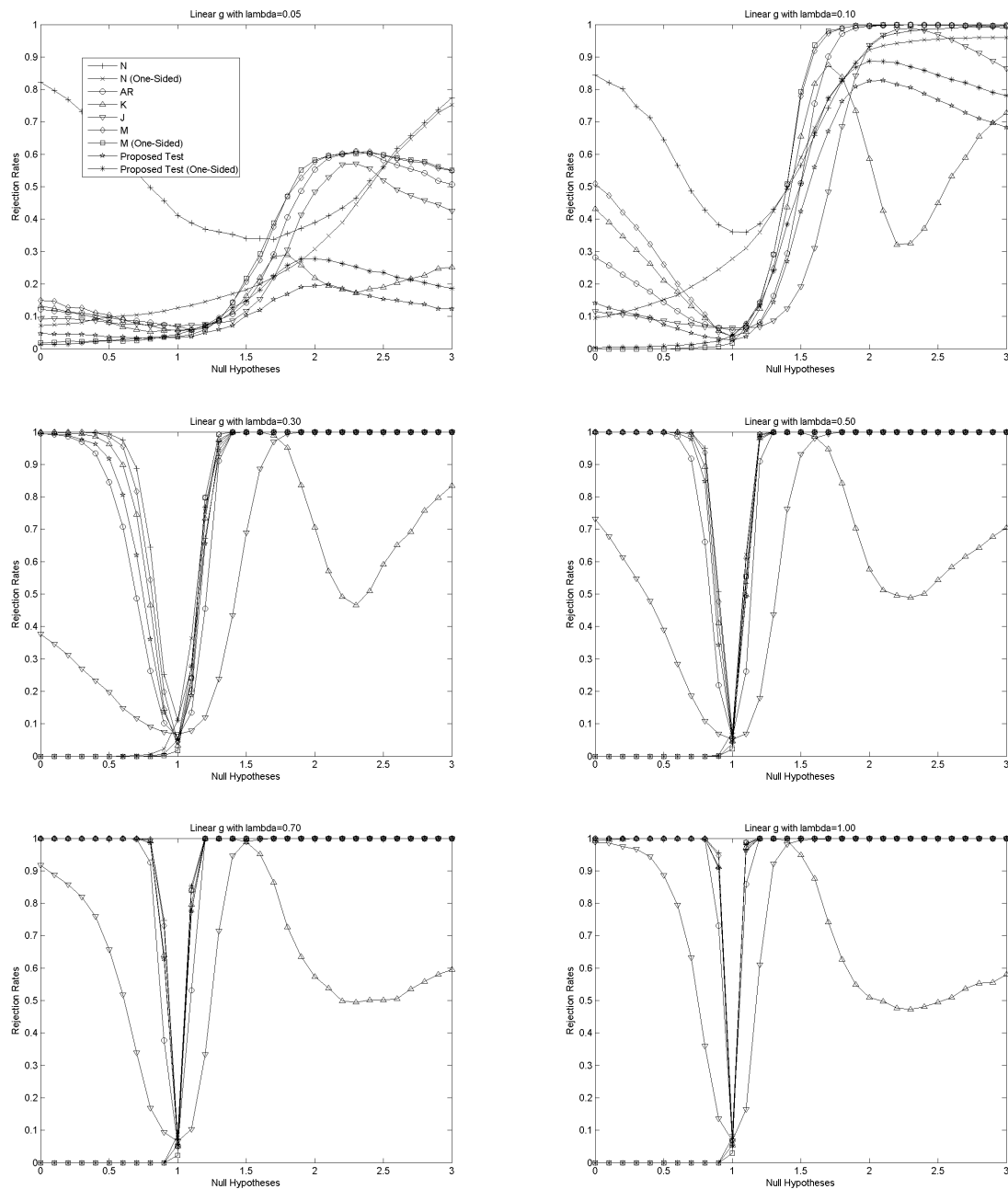


FIGURE 4. Monte Carlo rejection rates for linear  $g$  (DGP2)

The size of the M-test is moderately distorted, which may be attributable to the use of the heteroskedasticity-robust version of the test and would gradually disappear with an increase of the sample size. We feel encouraged that our test (which is also robust to heteroskedasticity) appears to have better size properties in these experiments. As we increase  $\lambda$ , the gain from

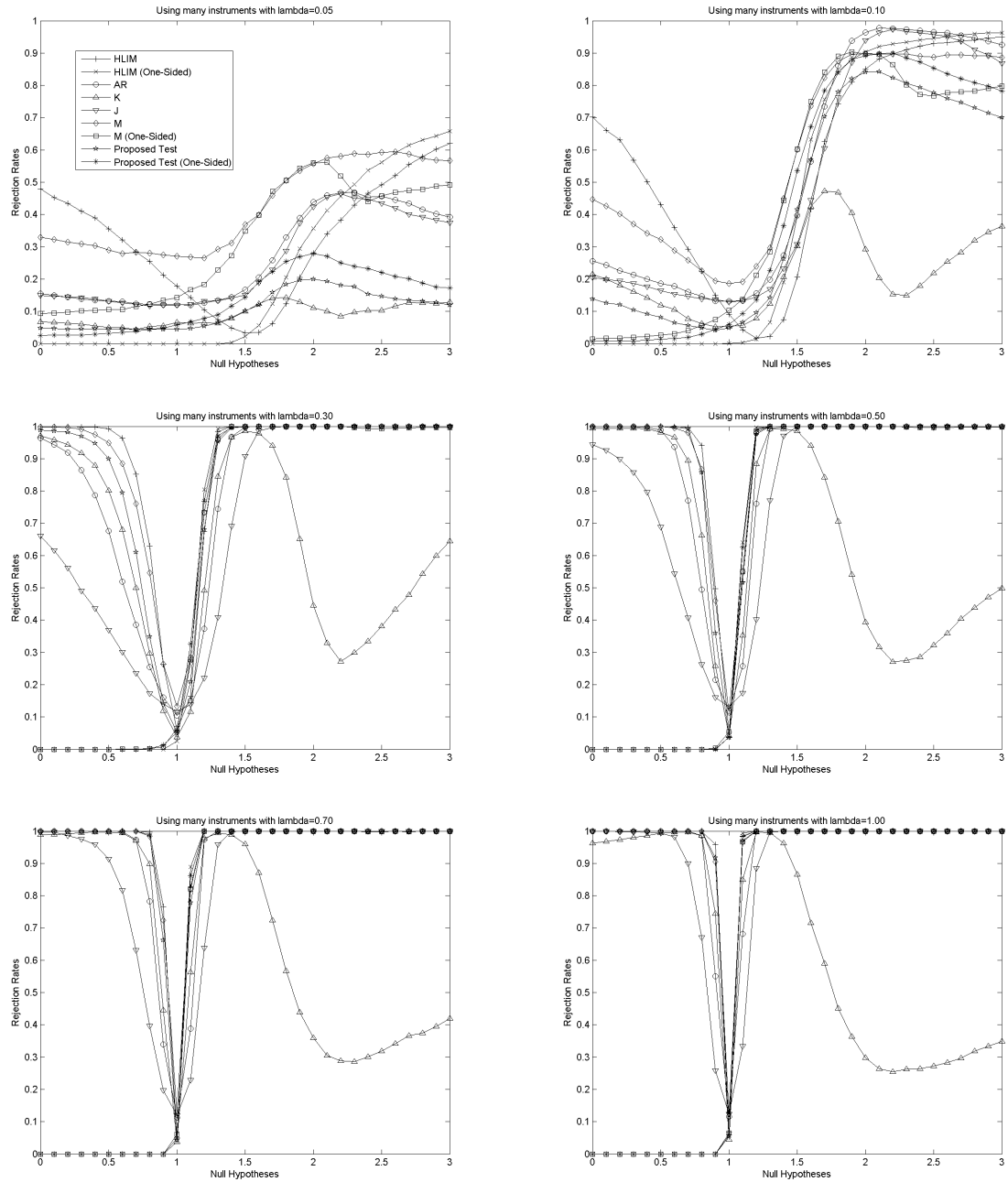


FIGURE 5. Monte Carlo rejection rates of the statistics using many instruments for linear  $g$  (DGP2)

using nonparametrically estimated instruments becomes apparent under DGP1 and our test is comparable to the N-test in this case.

Figure 4 shows that, again as expected, our test has less power than the M-test in a fully linear specification. When  $\lambda$  increases, however, the power curves of both tests are almost the same, as

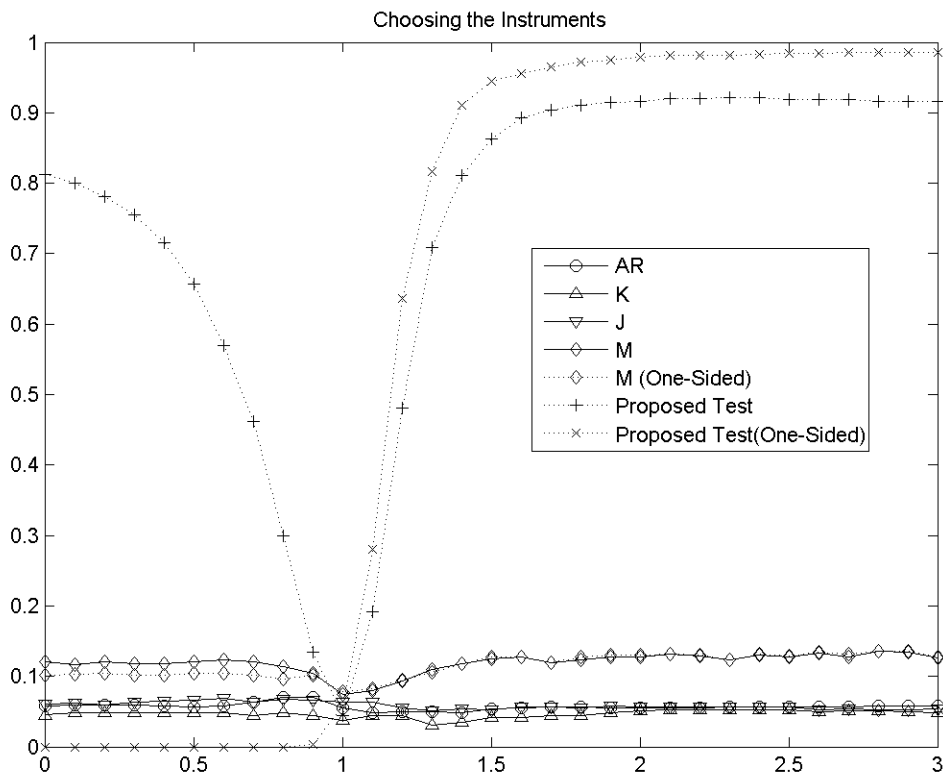


FIGURE 6. Monte Carlo rejection rates for quadratic  $g$  (DGP3)

our theoretical results would indicate. Note that the power curves of several tests are decreasing in the right tail and that the power curve of the  $K$ -test has a peculiar shape. It is straightforward to show that the  $AR$ -test has a (theoretical) nonmonotone power curve and since our test is in effect a semiparametric version of the  $AR$ -test, it inherits this feature. As for the  $K$ -test, it is well-known that when  $\theta_H - \theta_0 = \sigma_{uu}/\sigma_{uv}$  ( $= 1/0.8 \approx 1.25$  in DGP2), the  $K$ -test suffers from spurious power-declines.

In figure 5 we explore the effect of using more instruments, namely  $z_i$ ,  $z_i^2$  and  $z_i^3$  on the rejection rates of the  $AR$ -, the  $K$ -, the  $M$ -tests, and of the  $t$ -test based on HLIM estimators.. The test based on HLIM is incorrectly centered when  $\lambda$  is small, which is due to the fact that the HLIM estimator is then inconsistent. Note that the  $AR$ - and the  $M$ -tests get size-distorted as more instruments are used and identification becomes weaker. Although these distortions disappear asymptotically,<sup>23</sup> they are a concern in moderate-size samples and do not occur with our statistic.

There is any number of nonlinear specifications we could have chosen for DGP1 and depending on our choice the relative performance of the tests considered will vary. We feel that DGP1 is

<sup>23</sup>Simulations (not graphed) suggest that distortions of the  $AR$ - and the  $M$ -tests indeed become substantially smaller as the sample size gets bigger, even when  $\lambda = 0$ .

favoring the parametric tests somewhat since  $g$  is a function of a linear index. To emphasize this point we now consider an extreme example to illustrate the fact that using arbitrary unconditional moment conditions generated from unconditional moment conditions can lead to poor inference.

$$\text{DGP3: Nonlinear IV 2: } \begin{cases} y_i = Y_i\theta_0 + u_i \\ Y_i = \|z_i\|^2 - 8 + v_i, \end{cases}$$

where  $u_i, v_i$  are drawn from the same joint normal distribution as in DGP2. Note that here  $\mathbb{E}[z_i(y_i - Y_i\theta)] = -\mathbb{E}[z_i(\|z_i\|^2 - 8)](\theta - \theta_0) = 0$  for any  $\theta \in \mathbb{R}$  because the distribution of the  $z_i$ 's is even. Therefore, the unconditional moment condition using linear instruments,  $\mathbb{E}[z_i(y_i - Y_i\theta)]$ , does not identify the parameter of interest, while using  $g_i = \|z_i\|^2 - 8$  does identify  $\theta_0$ . Figure 6 shows that the power curves of the parametric tests are indeed almost flat, while the proposed test performs well. It is of course true that a parametric test using a polynomial expansion would have worked well in this particular case. But the whole point of using a nonparametric procedure is that one does not know the correct form of  $g$ , such that problems like the one illustrated here arise for any parametric choice for  $g$ .

In summary, parametric tests appear to perform better when  $g$  is linear and our test is preferable when  $g$  is nonlinear.

## 5. CONCLUSIONS

We have proposed a weak identification–robust test which can be applied in models with conditional moment restrictions. The test has attractive properties and requires weak assumptions, which do however exclude nonsmooth moment conditions and fat–tailed distributions. We have shown it to be asymptotically size–correct regardless of identification strength yet to have the same local power as the corresponding semiparametric t–test or Wald–test (at least under homoskedasticity) if identification is sufficiently strong. Our simulation experiments suggest that the test works well in practice.

## REFERENCES CITED

- ANDERSON, T. W., AND H. RUBIN (1949): “Estimation of the parameters of a single equation in a complete system of stochastic equations,” *Annals of Mathematical Statistics*, 20, 46–63.
- ANDREWS, D. W., M. J. MOREIRA, AND J. H. STOCK (2004): “Optimal invariant similar tests for instrumental variables regression,” Discussion paper, Cowles Foundation.
- (2006): “Optimal two–sided invariant similar tests for instrumental variables regression,” *Econometrica*, 74, 715–752.
- BEKKER, P. A. (1994): “Alternative approximations to the distribution of instrumental variables estimators,” *Econometrica*, 62, 657–681.
- CHAO, J. C., AND N. R. SWANSON (2005): “Consistent estimation with weak instruments,” *Econometrica*, 73, 1673–1692.
- DOMÍNGUEZ, M., AND I. LOBATO (2004): “Consistent estimation of models defined by conditional moment restrictions,” *Econometrica*, 72(5), 1601–1615.

- DUFOUR, J.-M. (1997): "Some impossibility theorems in econometrics with applications to structural and dynamic models," *Econometrica*, 65, 1365–1387.
- HAN, C., AND P. C. PHILLIPS (2006): "GMM with many moment conditions," *Econometrica*, 74, 147–192.
- HAUSMAN, J. A., W. K. NEWEY, T. WOUTERSEN, AND J. C. CHAO (2007): "Instrumental variable estimation with heteroskedasticity and many instruments," Discussion paper, MIT.
- JUN, S. J., AND J. PINKSE (2009): "Semiparametric tests of conditional moment restrictions under weak or partial identification," *Journal of Econometrics*, 152, 3–18.
- KLEIBERGEN, F. (2002): "Pivotal statistics for testing structural parameters in instrumental variables regression," *Econometrica*, 70, 1781–1804.
- (2005): "Testing parameters in GMM without assuming that they are identified," *Econometrica*, 73, 1103–1123.
- MOREIRA, M. J. (2003): "A conditional likelihood ratio test for structural models," *Econometrica*, 71, 1027–1048.
- NEWEY, W. K. (1990): "Efficient instrumental variables estimation of nonlinear models," *Econometrica*, 58, 809–837.
- (1993): *Efficient estimation of models with conditional moment restrictions* pp. 419–454. North Holland.
- (2004): "Many weak moment asymptotics for the continuously updated GMM estimator," Discussion paper, MIT.
- NEWEY, W. K., AND F. WINDMEIJER (2009): "GMM with many weak moment conditions," *Econometrica*, 77, 687–719.
- PHILLIPS, P. (1989): "Partially identified econometric models," *Econometric Theory*, 5(02), 181–240.
- ROBINSON, P. M. (1987): "Asymptotically efficient estimation in the presence of heteroskedasticity of unknown form," *Econometrica*, 55, 875–891.
- STAIGER, D., AND J. H. STOCK (1997): "Instrumental variables regression with weak instruments," *Econometrica*, 65, 557–586.
- STOCK, J. H., AND J. H. WRIGHT (2000): "GMM with weak identification," *Econometrica*, 68, 1055–1096.
- STOCK, J. H., J. H. WRIGHT, AND M. YOGO (2002): "A survey of weak instruments and weak identification in generalized method of moments," *Journal of Business and Economic Statistics*, 20, 518–529.
- STOCK, J. H., AND M. YOGO (2005): "Asymptotic distributions of instrumental variables statistics with many instruments," in *Identification and inference for econometric models: essays in honor of Thomas Rothenberg*, pp. 109–120. Cambridge University Press.
- STONE, C. J. (1977): "Consistent nonparametric regression," *Annals of Statistics*, 5, 595–620.
- VAN DER VAART, A. W. (1998): *Asymptotic statistics*. Cambridge University Press.

## APPENDIX A. PROOFS OF THEOREMS

Throughout the appendix we use the abbreviations RHS and LHS for *right hand side* and *left hand side*, respectively. Further, RHS<sub>1</sub> means the first right hand side term, and similarly for RHS<sub>2</sub>, RHS<sub>3</sub>, etcetera. Further, we let  $\mathcal{Z} = \{z_1, \dots, z_n\}$ .

**A.1. Proof of theorems 1 and 5.** We focus on  $T(\theta_0)$  defined in (17); the derivations for  $t(\theta_0)$  defined in (7) can be similarly shown by assuming that  $\mathcal{V}_i(\theta_0)$  is known. Omitting the  $\theta_0$  argument, write  $T = \hat{N}^\top \hat{D}^{-2} \hat{N}$ , where  $\hat{N} = \sum_i \hat{G}_i^\top \hat{\mathcal{V}}_i^{-1} m_i$ . Let  $\tilde{N}$  be defined as  $\hat{N}$  but with  $\mathcal{V}_i$  in lieu of  $\hat{\mathcal{V}}_i$ . Let  $m_i^* = \mathcal{V}_i^{-1} m_i$  and  $m_{\theta i}^* = \mathcal{V}_i^{-1} m_{\theta i}$ . Let  $\rho = \sqrt{n}(\lambda + 1/\sqrt{k})$ . Define

$$\tilde{D}^2 = n \text{Var}[G_1^\top m_1^*] + \sum_{ij} w_{ij}^2 \mathbb{E}[V_j^\top m_i^* m_i^{*\top} V_j | z_i, z_j] + \sum_{ij} w_{ij} w_{ji} \mathbb{E}[V_j^\top m_i^* m_j^{*\top} V_i | z_i, z_j].$$

Let further  $B_i = B_i(\theta_0) = \zeta_i + \sum_{j < i} \zeta_{ij}$  with  $\zeta_i = G_i^\top m_i^*$  and  $\zeta_{ij} = \ell_{ij} + \ell_{ji}$ , where  $\ell_{ij} = w_{ij} V_j^\top m_i^*$ . Let  $\tilde{N} = \sum_i B_i$ , where  $\{\tilde{D}^{-1} B_i\}$  forms a martingale difference array due to the fact that  $\mathbb{E}[B_i | B_{i-1}, \mathcal{Z}] = 0$  a.s..

- (1) We first show that  $\tilde{D}^{-1} \hat{N} \xrightarrow{d} N(0, I_{d_\theta})$ .
  - (a) We show that  $\|\hat{N} - \tilde{N}\| = o_p(\rho)$ , which is established in lemma C8.
  - (b) We show that  $\|\tilde{N} - \tilde{N}\| = o_p(\rho)$ 
    - (i) Since  $M_i(\theta_0) = 0$  a.s., lemma C15 establishes that  $\hat{N} = \sum_i \zeta_i + \sum_{ij} \ell_{ij} + o_p(\rho)$ .
    - (ii) The result then follows from the fact that  $\sum_i \sum_{j \neq i} \ell_{ij} = \sum_i \sum_{j < i} (\ell_{ij} + \ell_{ji})$ .
  - (c) We show that  $\tilde{D}^{-1} \tilde{N} \xrightarrow{d} N(0, I_{d_\theta})$ .
    - (i) Lemmas D1 to D4 imply that  $\forall \|c\| = 1 : c^\top \tilde{D}^{-1} \tilde{N} \xrightarrow{d} N(0, 1)$ .
    - (ii) The result then follows from the Cramér–Wold device.
  - (d) Finally, note that  $\|\tilde{D}^{-1}(\hat{N} - \tilde{N})\| \leq \|\tilde{D}^{-1}\| \|\hat{N} - \tilde{N}\| = O_p(\rho^{-1}) o_p(\rho) = o_p(1)$ .
- (2) We show that  $\tilde{D}^{-1} \hat{D} \xrightarrow{p} I_{d_\theta}$ .
  - (a)  $\hat{D}^2 - \tilde{D}^2 = o_p(\rho^2)$  is established in lemmas C14 and C22.
  - (b) Since  $\tilde{D}^{-1} = O_p(\rho^{-1})$  by lemma C25, we have  $\|\tilde{D}^{-1} \hat{D} - I_{d_\theta}\| = \|\tilde{D}^{-1}(\hat{D} - \tilde{D})\| \leq \|\tilde{D}^{-1}\| \sqrt{\|\hat{D}^2 - \tilde{D}^2\|} = o_p(1)$ .<sup>24</sup>
- (3) So, the result follows from combining the results established in items (1) and (2).

**A.2. Proof of theorems 2 and 6.** Again, we focus on  $T(\theta)$  because  $t(\theta)$  can be similarly dealt with by assuming  $\mathcal{V}_i(\theta)$  is known. We start with part (i).

- (1) We show that  $\hat{N}$  diverges at a rate of  $n\lambda^2$ .
  - (a) We show that  $\hat{N} = n\lambda^2 \mathbb{E}[\tilde{G}_1^\top \mathcal{V}_1^{-1} \tilde{M}_1] + o_p(\rho) + o_p(n\lambda^2)$ , which follows from lemmas C8 and C17, noting that  $n\lambda^2 \succ \rho$  by assumption.
  - (b)  $\mathbb{E}[\tilde{G}_1^\top \mathcal{V}_1^{-1} \tilde{M}_1] \neq 0$  by assumption CG.
- (2) We show that  $\hat{D}^2 = o_p(n^2\lambda^4)$ , which implies that  $\hat{D} = o_p(n\lambda^2)$ .

<sup>24</sup>  $\|\hat{D} - \tilde{D}\|^2 \leq \|\hat{D}^2 - \tilde{D}^2\|$  because if  $c$  is an eigenvector of  $\hat{D} - \tilde{D}$  with eigenvalue  $\nu$  then if  $\nu > 0$ ,  $c^\top (\hat{D} - \tilde{D})^2 c = c^\top (\hat{D}^2 - \tilde{D}^2) c - 2\nu c^\top \tilde{D} c \leq c^\top (\hat{D}^2 - \tilde{D}^2) c$  and if  $\nu < 0$  then  $c^\top (\hat{D} - \tilde{D})^2 c = c^\top (\tilde{D}^2 - \hat{D}^2) c + 2\nu c^\top \hat{D} c \leq c^\top (\tilde{D}^2 - \hat{D}^2) c$ .

- (a) Recall that  $\|\hat{D}^2 - \tilde{D}^2\| = o_p(\rho^2)$  by lemmas C14 and C22.  
 (b) Since  $\|\tilde{D}^2\| = O_p(\rho^2)$  by lemma C24, we have  $\|\hat{D}^2\|/n^2\lambda^4 \leq o_p(\rho^2/n^2\lambda^4) + O_p(\rho^2/n^2\lambda^4) = o_p(1)$ .  
 (3) Since  $T^{-1} = (\hat{N}^\top \hat{D}^{-2} \hat{N})^{-1} \leq \{\|\hat{N}\|^2 \underline{\text{ev}}(\hat{D}^{-2})\}^{-1} = \|\hat{D}^2\|/\|\hat{N}\|^2 = o_p(1)$ , (i) hence holds.  
 (4) Now part (ii). By lemma C23,  $\tilde{D}^2 = n \text{Var}[\tilde{G}_1^\top \mathcal{V}_1^{-1} m_1] + o_p(n\lambda^2)$ , and hence

$$\begin{aligned} \hat{D}^{-1} \hat{N} &= \left( n\lambda^2 \text{Var}[\tilde{G}_1^\top \mathcal{V}_1^{-1} m_1] + o_p(n\lambda^2) \right)^{-1/2} \left( n\lambda^2 \mathbb{E}[\tilde{G}_1^\top \mathcal{V}_1^{-1} \tilde{M}_1] + o_p(n\lambda^2) \right) \\ &= \sqrt{n\lambda} \left( \text{Var}[\tilde{G}_1^\top \mathcal{V}_1^{-1} m_1] \right)^{-1/2} \mathbb{E}[\tilde{G}_1^\top \mathcal{V}_1^{-1} \tilde{M}_1] + o_p(\sqrt{n\lambda}). \end{aligned}$$

### A.3. Proof of theorem 3.

- (1) We first have by lemmas E5 and E6,

$$\sum_i \hat{g}_i(\theta_n) m_i(\theta_n) = \sum_i \hat{g}_i m_i + \sqrt{n\lambda} \mathbb{E} \tilde{g}_1 \tilde{g}_1^\top \Delta + o_p(\rho). \quad (18)$$

- (2) Also,

$$\hat{D}_s^2(\theta_n) = \left\{ \sum_i \hat{g}_i \hat{g}_i^\top m_i^2 - n^{-1} \left( \sum_i \hat{g}_i m_i \right) \left( \sum_i \hat{g}_i^\top m_i \right) + \sum_{ij} w_{ij} w_{ji} m_j m_j^\top m_i m_{\theta i} \right\} + o_p(\rho^2), \quad (19)$$

by lemmas E5 to E8.

- (3) Combining the results of (18) and (19) with those of theorem 1, we obtain the stated result.

A.4. **Proof of theorem 4.** Express the test statistic  $\hat{t}$  as  $\hat{\mathcal{G}}_N^{-1} \hat{\mathcal{N}}_N$  and let  $\mathbf{t} = \hat{\mathcal{G}}_s^{-1} \hat{\mathcal{N}}_s$ , where  $\hat{\mathcal{G}}_s, \hat{\mathcal{N}}_s$  are defined as  $\hat{N}_s$  and  $\hat{D}_s$  with  $\hat{g}_i$  replaced with  $\hat{q}_i$ ; then, the results for  $\hat{N}_s$  and  $\hat{D}_s$  also hold for  $\hat{\mathcal{N}}_s$  and  $\hat{\mathcal{G}}_s$ . We proceed as follows. The  $\theta_0$ -argument is dropped in the list below.

- (1) First asymptotic validity.

- (a) Lemma F4 shows that  $\hat{\mathcal{N}}_N - \hat{\mathcal{N}}_s = o_p(\rho)$ .  
 (b) Lemma F8 establishes that  $\hat{\mathcal{G}}_N^2 - \hat{\mathcal{G}}_s^2 = o_p(\rho^2)$ .  
 (c) Consequently,

$$\left\{ (\hat{\mathcal{G}}_N - \hat{\mathcal{G}}_s)/\rho + \hat{\mathcal{G}}_s/\rho \right\}^{-1} \left\{ (\hat{\mathcal{N}}_N - \hat{\mathcal{N}}_s)/\rho + \hat{\mathcal{N}}_s/\rho \right\} = \mathbf{t} + o_p(1).$$

- (d) Apply theorem 1.

- (2) For consistency:

- (a) Lemma F4 shows that  $\hat{\mathcal{N}}_N - \hat{\mathcal{N}}_s = O_p(\rho) + o_p(n\lambda^2)$ . Since  $n\lambda^2 \succ \rho$  by assumption,  $\hat{\mathcal{N}}_N = n\lambda^2 \mathbb{E}[\tilde{q}_i \tilde{M}_1] + o_p(n\lambda^2)$ , which follows from applying lemma C17 to  $\hat{\mathcal{N}}_s$ , replacing  $g_i$  with  $q_i$ .  
 (b) Lemma F8 establishes that  $\hat{\mathcal{G}}_N^2 - \hat{\mathcal{G}}_s^2 = o_p(\rho^2)$ . Further, from the proof of theorem 2 with  $\hat{\mathcal{G}}_s$  in lieu of  $\hat{D}_s$ , it follows that  $\hat{\mathcal{G}}_s^2 = o_p(n^2\lambda^4)$ . Hence  $\hat{\mathcal{G}}_s = o_p(n\lambda^2)$ .



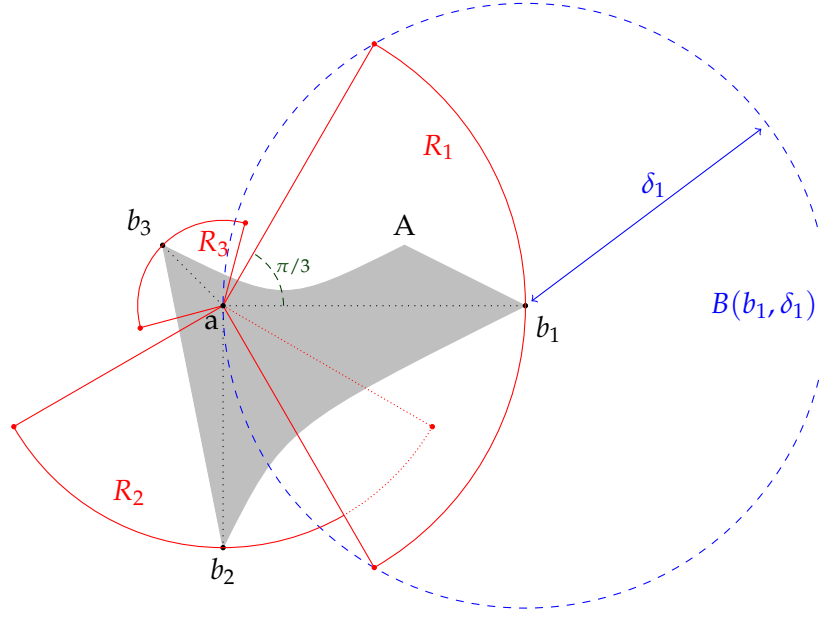


FIGURE 7. Proof of lemma B2

## APPENDIX B. TECHNICAL LEMMAS

**Lemma B1.** For  $p_1 > 0$  let  $f_i = f(z_i)$  be such that  $\mathbb{E}|f_i|^{p_1 p_2} < \infty$  and  $a_i > 0$  be such that  $\mathbb{E}a_i^{p_2/(p_2-1)} < \infty$  for some  $p_2 \geq 1$ .<sup>25</sup> Then

$$\mathbb{E}[w_{12}a_1|f_1 - f_2|^{p_1}] = o(n^{-1}). \quad (20)$$

*Proof.* Note first that by the Jensen inequality,

$$\left(\sum_j w_{ij}|f_i - f_j|^{p_1}\right)^{p_2} \leq \sum_j w_{ij}|f_i - f_j|^{p_1 p_2}. \quad (21)$$

The expectation of the RHS in (21) is  $o(1)$  by lemma 1 of Robinson (1987). Hence the LHS in (20) is by the Hölder inequality bounded by

$$\left(\mathbb{E}|a_1|^{p_2/(p_2-1)}\right)^{(p_2-1)/p_2} \left[\mathbb{E}\left(\sum_j w_{ij}|f_i - f_j|^{p_1}\right)^{p_2}\right]^{1/p_2} = O(1)o(1) = o(1). \quad \square$$

Let  $\tau(a, b) = \mathbb{P}(\|z_1 - b\| \leq \|a - b\|)$  and  $A(a) = \{b : \tau(a, b) \leq 2k/n\}$ .

**Lemma B2.**  $\sup_a \mathbb{P}[z_1 \in A(a)] = O(k/n)$ .

*Proof.* We establish the result for  $z_i \in \mathbb{R}^2$ ; the case in which  $z_i$  is scalar-valued is simple and the case in which  $z_i$  has dimension greater than two is only more complicated in terms of exposition.

Let  $A = A(a)$  and  $b(t, \omega) = (a_1 + t \cos \omega, a_2 + t \sin \omega)$  and further  $\delta_1 = \sup_{\omega, t: b(t, \omega) \in A} t$ ,  $\omega_1$  a corresponding  $\omega$  and  $b_1 = b(\delta_1, \omega_1)$ ; see figure 7. If  $\delta_1 = 0$  there is nothing to prove so suppose that  $\delta_1 > 0$ . Then let  $R_1 = \{b = b(t, \omega) : |\omega - \omega_1| < \pi/3, 0 \leq t \leq \delta_1\}$ . So  $R_1$  is a slice of a pie

<sup>25</sup> $p_2 = 1$  means that  $a_i$  is bounded a.s..

with radius  $\delta_1$  extending  $\pi/3$  in both directions from  $\omega$ ; see again figure 7. By construction for all  $|\omega - \omega_1| < \pi/3$  and  $t \geq 0$ :  $b(t, \omega) \in A \Rightarrow b(t, \omega) \in R_1$ .

Now let  $\delta_2 = \sup_{\omega, t: b(t, \omega) \in A \setminus R_1} t$  and define  $\omega_2, b_2, R_2$  accordingly. Repeat this procedure until  $A$  is covered by at most 10  $R_j$ -sets. By symmetry it suffices to show that  $\mathbb{P}(z_2 \in R_1) \leq 2k/n$ . Now let  $A_1 = \{b(t, \omega_1) : 0 \leq t < \delta_1\}$ . Then if  $\mathcal{B}$  denotes a closed ball and  $B(b_1, \delta_1) = \bigcup_{b \in A_1} \mathcal{B}(b, \|b - a\|)$ , then because  $A_1 \subset A$ ,

$$\mathbb{P}(z_2 \in R_1) \leq \mathbb{P}[z_2 \in B(b_1, \delta_1)] = \lim_{t \uparrow \delta_1} \mathbb{P}[z_2 \in \mathcal{B}(b(t, \omega_1), t)] \leq \lim_{t \uparrow \delta_1} 2k/n = 2k/n. \quad \square$$

**Lemma B3.**  $\sup_a \mathbb{E}[w_{21}w_{31}|z_1 = a] = O(n^{-2})$ .

*Proof.* Let  $\mathcal{N}_i$  denote the set of neighbors of  $i$ . Since for any  $a$ ,

$$\mathbb{E}[w_{21}w_{31}|z_1 = a] \leq \frac{C_w^2}{k^2} \mathbb{E}[I(1 \in \mathcal{N}_2)I(1 \in \mathcal{N}_3)|z_1 = a],$$

it suffices to show that  $\sup_a \mathbb{E}[I(1 \in \mathcal{N}_2)I(1 \in \mathcal{N}_3)|z_1 = a] = O(k^2/n^2)$ . First, suppose that  $\mathbb{P}(z_1 = a) = c > 0$ . Let  $S_j^* = \sum_{i \neq j, 1} I(z_i = a)$ . Then using the randomization scheme for assigning weights in the case of ties we have

$$\begin{aligned} & \mathbb{E}[I(1 \in \mathcal{N}_2)I(1 \in \mathcal{N}_3)|z_1 = a] \\ &= \mathbb{E}[I(1 \in \mathcal{N}_2)I(1 \in \mathcal{N}_3)\{I(S_2^* \leq c(n-2)/2) + I(S_2^* > c(n-2)/2)\}|z_1 = a] \\ &\leq \mathbb{P}[S_2^* \leq c(n-2)/2] + \frac{2k}{c(n-2)} \mathbb{E}[I(1 \in \mathcal{N}_3)|z_1 = a] \\ &\leq \mathbb{P}[S_2^* \leq c(n-2)/2] + \frac{2k}{c(n-2)} \mathbb{P}[S_3^* \leq c(n-2)/2] + \frac{4k^2}{c^2(n-2)^2}. \quad (22) \end{aligned}$$

Now by the Hoeffding inequality,

$$\mathbb{P}(S_2^* \leq c(n-2)/2|z_1 = a) \leq \mathbb{P}(|S_2^* - c(n-2)| \geq c(n-2)/2) \leq \exp(-c^2(n-2)/2),$$

which decreases exponentially fast. Now suppose that  $\mathbb{P}(z_1 = a) = 0$ . Let  $S_j(a, b) = \sum_{i \neq j, 1} I(\|z_i - b\| \leq \|a - b\|)$  and let  $\tau$  be as defined prior to lemma B2. Then because

$$\begin{aligned} I(1 \in \mathcal{N}_2) &\leq I(\tau(z_1, z_2) \leq 2k/n) + I(1 \in \mathcal{N}_2)I(\tau(z_1, z_2) > 2k/n) \\ &\leq I(\tau(z_1, z_2) \leq 2k/n) + I(S_2^*(z_1, z_2) < k)I(\tau(z_1, z_2) > 2k/n) \\ &\leq I(\tau(z_1, z_2) \leq 2k/n) + I(|S_2(z_1, z_2) - (n-2)\tau(z_1, z_2)| > k), \end{aligned}$$

we have

$$\begin{aligned} \mathbb{E}[I(1 \in \mathcal{N}_2)I(1 \in \mathcal{N}_3)|z_1 = a] &\leq \mathbb{E}[I(\tau(a, z_2) \leq 2k/n)I(\tau(a, z_3) \leq 2k/n)|z_1 = a] \\ &\quad + \mathbb{P}[|S_2(a, z_2) - (n-2)\tau(a, z_2)| > k] + \mathbb{P}[|S_3(a, z_3) - (n-2)\tau(a, z_3)| > k]. \quad (23) \end{aligned}$$

RHS1 in (23) is  $\sup_a [\mathbb{P}(z_1 \in A(a))]^2 = O(k^2/n^2)$  by lemma B2. RHS2 and RHS3 in (23) are the same and are by the Hoeffding inequality bounded by  $\exp(-2k^2/(n-2))$ , which goes to zero at an exponential rate.  $\square$

**Lemma B4.**  $\max_j \sum_i w_{ij} = O_p(1)$ .

*Proof.* Since the proof is similar to that of lemma B3, we only establish the result for the case in which the probability of ties is zero. Let  $S_{ij}(a, b) = \sum_{t \neq i, j} I(\|z_t - b\| \leq \|a - b\|)$ . Now, for any  $0 < C^* < \infty$  and  $C = C/3C_w$ ,

$$\begin{aligned} & \mathbb{P} \left[ \max_j \sum_i w_{ij} > C^* \right] \leq \mathbb{P} \left[ \max_j \sum_{i \neq j} I(j \in \mathcal{N}_i) > 3Ck \right] \\ & \stackrel{(22)}{\leq} \mathbb{P} \left[ \max_j \sum_{i \neq j} I(\tau(z_j, z_i) \leq 2k/n) > 2Ck \right] + \mathbb{P} \left[ \max_j \sum_{i \neq j} I(|S_{ij}(z_j, z_i) - (n-2)\tau(z_j, z_i)| > k) > Ck \right] \\ & \leq n \sup_a \mathbb{P} \left[ \left| \sum_i \{I(\tau(a, z_i) \leq 2k/n) - \mathbb{P}[\tau(a, z_i) \leq 2k/n]\} \right| > Ck \right] + \sup_a I(\mathbb{P}[\tau(a, z_1) \leq 2k/n] > Ck/n) \\ & \quad + n^2 \sup_{a, b} \mathbb{P} \left[ \left| \sum_t (I(\|z_t - b\| \leq \|a - b\|) - \tau(a, b)) \right| > k \right]. \quad (24) \end{aligned}$$

RHS1 and RHS3 in (24) are by theorem 1.1 of de la Peña (1999) respectively bounded by  $2n \exp(-C^2k/(4+2C)) = o(1)$  and  $2n^2 \exp(-k^2/4n) = o(1)$ . For RHS2, let  $n \rightarrow \infty$  followed by  $C \rightarrow \infty$ .  $\square$

**Lemma B5.** *Suppose that for some  $p \geq 1$ ,  $\mathbb{E}|\mathbb{E}[b_1|z_1]|^p < \infty$  and  $\mathbb{E}|\mathbb{E}[a_1|z_1]|^{p/(p-1)} < \infty$ . Then*

(i)  $\mathbb{E}[w_{12}a_1b_2] = n^{-1}\mathbb{E}[a_1\mathbb{E}[b_1|z_1]] + o(n^{-1}) = O(n^{-1})$ ,

(ii) for any  $1 < p^* < \infty$ :  $\mathbb{E}[w_{12}^{p^*}a_1b_2] = O(n^{-1}k^{1-p^*})$ .

If  $\mathbb{E}|\mathbb{E}[a_1|z_1]|^{2p} < \infty$ ,  $\mathbb{E}|\mathbb{E}[b_1|z_1]|^{2p} < \infty$  and  $\mathbb{E}|\mathbb{E}[c_1|z_1]|^{p/(p-1)} < \infty$  for some  $p \geq 1$  then

(iii)  $\mathbb{E}[w_{12}w_{13}c_1a_2b_3] = O(n^{-2})$ ,

(iv)  $\sum_t \mathbb{E}[w_{12}w_{1t}c_1a_2b_t] = O(n^{-1})$ .

Further, if for some  $p_1 > 4$  and  $p_2 \geq 1$ ,  $\mathbb{E}|\mathbb{E}[a_1|z_1]|^{p_1 p_2} + \mathbb{E}|\mathbb{E}[b_1|z_1]|^{p_1 p_2 / (p_2 - 1)} + \mathbb{E}|c_1| < \infty$ , then

(v)  $\mathbb{E}[w_{21}w_{31}c_1a_2b_3] = o(n^{-3/2}k^{-1/2})$ , and

(vi)  $\sum_t \mathbb{E}[w_{21}w_{t1}c_1a_2b_t] = o(1/\sqrt{nk})$ .

*Proof.* First (i). Note that

$$\begin{aligned} \mathbb{E}[w_{12}a_1b_2] &= \mathbb{E}[w_{12}\mathbb{E}[a_1|z_1]\mathbb{E}[b_2|z_2]] \\ &= \mathbb{E}[w_{12}\mathbb{E}[a_1|z_1]\mathbb{E}[b_1|z_1]] + \mathbb{E}[w_{12}\mathbb{E}[a_1|z_1]\{\mathbb{E}[b_2|z_2] - \mathbb{E}[b_1|z_1]\}] \\ &= n^{-1}\mathbb{E}[\mathbb{E}[a_1|z_1]\mathbb{E}[b_1|z_1]] + o(n^{-1}), \end{aligned}$$

by lemma B1. (ii) follows from (i) noting that  $w_{12}^{p^*-1} \leq C_w^{p^*-1}k^{1-p^*}$  by construction. For (iii), note that

$$|\mathbb{E}[w_{12}w_{13}c_1a_2b_3]| \leq \mathbb{E}[w_{12}w_{13}|c_1|(a_2^2 + b_3^2)] = n^{-1}\mathbb{E}[w_{12}|c_1|a_2^2] + n^{-1}\mathbb{E}[w_{12}|c_1|b_3^2] = O(n^{-2}),$$

by (i). For (iv) note that

$$\sum_t \mathbb{E}[w_{12}w_{1t}c_1a_2b_t] = \mathbb{E}[w_{12}^2c_1a_2b_2] + (n-1)\mathbb{E}[w_{12}w_{13}c_1a_2b_3] = O(n^{-1}k^{-1}) + O(n^{-1}) = O(n^{-1}),$$

by parts (ii) and (iii). Further, (v) follows since for  $\tilde{a}_i = \mathbb{E}[a_i|z_i]$ ,  $\tilde{b}_i = \mathbb{E}[b_i|z_i]$ ,

$$\begin{aligned} |\mathbb{E}[w_{21}w_{31}c_1a_2b_3]| &= |\mathbb{E}[w_{21}w_{31}c_1\tilde{a}_2\tilde{b}_3]| = |\mathbb{E}\{\mathbb{E}[w_{21}w_{31}\tilde{a}_2\tilde{b}_3|z_1]c_1\}| \\ &\stackrel{\text{H\"older}}{\leq} \mathbb{E}\{(\mathbb{E}[(w_{21}w_{31})^{p_1/(p_1-1)}|z_1])^{(p_1-1)/p_1}c_1\}(\mathbb{E}|\tilde{a}_2\tilde{b}_3|_1^p)^{1/p_1} \\ &\stackrel{\text{B}_3, \text{H\"older}}{\leq} C_w^2 n^{-2} (n/k)^{2/p_1} \mathbb{E}|c_1| (\mathbb{E}|\tilde{a}_2|^{p_1 p_2})^{1/(p_1 p_2)} (\mathbb{E}|\tilde{b}_3|^{p_1 p_2/(p_2-1)})^{(p_2-1)/(p_1 p_2)} \prec n^{-3/2} k^{-1/2}. \end{aligned}$$

where the first equality follows from lemma B3 and the second from part (i). Finally, (vi) follows from (ii) and (v).  $\square$

**Lemma B6.** Let  $a_i, b_i$  be such that for some  $p \geq 1$ ,  $\mathbb{E}|\mathbb{E}[b_1^2|z_1]|^p < \infty$  and  $\mathbb{E}|\mathbb{E}[a_1^2|z_1]|^{p/(p-1)} < \infty$  and let  $p_1, p_2$  be such that  $\min(p_1, p_2) \geq 0$  and  $\max(p_1, p_2) \geq 1$ . Let furthermore  $|c_{ij}| \leq 1$  be constants. Then

- (i) If  $\mathbb{E}[a_i|z_i] = \mathbb{E}[b_i|z_i] = 0$  a.s. then  $\sum_{ij} w_{ij}^{p_1} w_{ji}^{p_2} c_{ij} a_i b_j = O_p(n^{1/2} k^{1/2 - p_1 - p_2})$ ,
  - (ii) If  $\mathbb{E}[a_i|z_i] = 0$  a.s. and  $p_1 \geq 1$  then  $\sum_{ij} w_{ij}^{p_1} w_{ji}^{p_2} c_{ij} a_i b_j = O_p(n^{1/2} k^{1 - p_1 - p_2})$ ,
  - (iii) if  $\max(p_1, p_2) \geq 1$  then  $\sum_{ij} w_{ij}^{p_1} w_{ji}^{p_2} c_{ij} [a_i b_j - \mathbb{E}[a_i|z_i] \mathbb{E}[b_j|z_j]] = O_p(n^{1/2} k^{1 - p_1 - p_2})$ ,
- If furthermore, for some  $p_3 > 8$ ,  $\mathbb{E}|\mathbb{E}[b_1|z_1]|^{p_3} < \infty$  then
- (iv) if  $\mathbb{E}[a_i|z_i] = 0$  a.s. then  $\sum_{ij} w_{ij}^{p_2} c_{ij} a_i b_j = o_p(n^{3/4} k^{3/4 - p_2})$ ,
  - (v)  $\sum_{ij} w_{ij}^{p_2} c_{ij} [a_i b_j - \mathbb{E}[a_i|z_i] \mathbb{E}[b_j|z_j]] = o_p(n^{3/4} k^{3/4 - p_2})$ ,
  - (vi)  $\sum_{ij} w_{ij}^{p_1} w_{ji}^{p_2} a_i b_j = \sum_{ij} w_{ij}^{p_1} w_{ji}^{p_2} \mathbb{E}[a_i|z_i] \mathbb{E}[b_j|z_j] + o_p(nk^{1 - p_1 - p_2})$ .

*Proof.* For (i), suppose that  $p_1 \geq 1$  (the case  $p_2 \geq 1$  is symmetric), square and take expectations to obtain

$$\begin{aligned} \sum_{ij} c_{ij} \mathbb{E}[w_{ij}^{2p_1} w_{ji}^{2p_2} a_i^2 b_j^2 + (w_{ij} w_{ji})^{p_1 + p_2} a_i b_i a_j b_j] \\ \leq (C_w/k)^{2p_1 + 2p_2 - 1} \sum_{ij} w_{ij} (a_i^2 b_j^2 + |a_i b_i a_j b_j|) \stackrel{\text{B}_5}{=} O(nk^{1 - 2p_1 - 2p_2}). \end{aligned}$$

For (ii) likewise square and take expectations which yields

$$\sum_{ijt} c_{ij} c_{it} \mathbb{E}[(w_{ij} w_{it})^{p_1} (w_{ji} w_{ti})^{p_2} a_i^2 b_j b_t] \leq (C_w/k)^{2p_1 + 2p_2 - 2} \sum_{ijt} \mathbb{E}[w_{ij} w_{it} a_i^2 b_j b_t] \stackrel{\text{B}_5}{=} O(nk^{2(1 - p_1 - p_2)}).$$

For (iii) write  $a_i b_j - \mathbb{E}[a_i|z_i] \mathbb{E}[b_j|z_j] = (a_i - \mathbb{E}[a_i|z_i]) (b_j - \mathbb{E}[b_j|z_j]) + \mathbb{E}[a_i|z_i] (b_j - \mathbb{E}[b_j|z_j]) + (a_i - \mathbb{E}[a_i|z_i]) \mathbb{E}[b_j|z_j]$  and apply parts (i) and (ii). For (iv) we get

$$\sum_{ijt} c_{ij} c_{it} \mathbb{E}[(w_{ji} w_{ti})^{p_2} a_i^2 b_j b_t] \leq (C_w/k)^{2p_2 - 2} \sum_{ijt} \mathbb{E}[w_{ji} w_{ti} a_i^2 b_j b_t] \stackrel{\text{B}_5}{=} o(n^{3/2} k^{3/2 - 2p_2}).$$

For (v) use the same expansion as for (iii) but apply parts (i) and (iv). Finally, for (vi) use either (iii) or (v) and (supposing  $p_1 \geq 1$ ;  $p_2 \geq 1$  is similar)

$$\begin{aligned} & \left| \sum_{ij} w_{ij}^{p_1} w_{ji}^{p_2} \mathbb{E}[a_i | b_i] [\mathbb{E}[b_j | z_j] - \mathbb{E}[b_i | z_i]] \right| \\ & \leq n^2 (C_w/k)^{p_1+p_2-1} \mathbb{E}[w_{12} \mathbb{E}[a_1 | z_1] (\mathbb{E}[b_2 | z_2] - \mathbb{E}[b_1 | z_1])] \stackrel{\text{B1}}{=} o(nk^{1-p_1-p_2}). \quad \square \end{aligned}$$

**Lemma B7.** Let  $\{b_j\}$  be i.i.d. conditional on  $\mathcal{Z}$  with  $\mathbb{E}[b_j | \mathcal{Z}] = 0$  a.s., and let  $c_{ij} = c(z_i, z_j, \mathcal{Z})$ . Let further  $\mathcal{S}_{ni} = \sum_j c_{ij} b_j$  and  $\mathcal{Q}_{ni} = \sum_j c_{ij}^2 \mathbb{E}[b_j^2 | \mathcal{Z}]$ . Then if  $\mathbb{E}|b_j|^p < \infty$  for some  $2 \leq p < \infty$  and  $\max_i \mathcal{Q}_{ni} = O_p(C_n^*)$ , then for  $C_n = (\sqrt{C_n^*} + n^{1/p}/k)(\log n)^2$ ,  $\max_i |\mathcal{S}_{ni}| = o_p(C_n)$ .

*Proof.* First, for some  $0 < C, C^* < \infty$ ,

$$\begin{aligned} \mathbb{P}[\max_i |\mathcal{S}_{ni}| > C_n] & \leq \sum_i \mathbb{P}[|\mathcal{S}_{ni}| > C_n, \max_i \mathcal{Q}_{ni} \leq C^* C_n^*, \max_j |b_j| \leq C_n^{1/p}] \\ & \quad + \mathbb{P}[\max_i \mathcal{Q}_{ni} > C^* C_n^*] + \sum_i \mathbb{P}[|b_i| > C_n^{1/p}]. \quad (25) \end{aligned}$$

RHS1 in (25) is by theorem 1.1 of de la Peña (1999) bounded by  $2n \exp(-C_n^2/2(C^* C_n^* + C_n^{1/p} C_n/k)) = o(1)$  by the definition of  $C_n$ . RHS2 is also  $o(1)$ , which follows by taking the limit for  $n \rightarrow \infty$  followed by  $C^* \rightarrow \infty$ . Finally, RHS3 is by the Markov inequality bounded by  $\mathbb{E}|b_i|^p/C^p$ . Take  $C \rightarrow \infty$ .  $\square$

**Lemma B8.** If  $\{b_j\}$  is i.i.d. and  $\mathbb{E}|b_j|^p < \infty$  for some  $1 \leq p < \infty$  then  $\max_i |\sum_j w_{ij} b_j| = O_p((n/k)^{1/p})$ .

*Proof.* By the Hölder inequality for any  $i$  and some  $C < \infty$  independent of  $i$ ,

$$|\sum_j w_{ij} b_j| \leq (\sum_j w_{ij}^{p/(p-1)})^{1-1/p} (\sum_j |b_j|^p)^{1/p} \leq C (\sum_j |b_j|^p)^{1/p} / k^{1/p}. \quad \square$$

**Lemma B9.** Let  $\{b_j\}$  be such that  $\{(b_j, z_j)\}$  is i.i.d. with  $\mathbb{E}[b_j | z_j] = 0$  a.s. and  $\mathbb{E}|b_j|^p < \infty$  for some  $2 \leq p < \infty$  and  $\mathbb{E}[\mathbb{E}[b_j^2 | z_j]]^{p^*} < \infty$  for some  $1 \leq p^* < \infty$ . Then  $\max_i |\sum_j w_{ij} b_j| = o_p((n^{1/p}/k + n^{1/2p^*}/k^{(p^*+1)/2p^*})(\log n)^2)$ .

*Proof.* Take  $c_{ij} = w_{ij}$  in B7 and note that  $\max_i \mathcal{Q}_{ni} = O_p(n^{1/p^*}/k^{1+1/p^*})$  by B8.  $\square$

**Lemma B10.** For an arbitrary random sequence  $\{b_n\}$ , if  $\mathbb{E}[|b_n| | \mathcal{Z}] = O_p(1)$ , then  $b_n = O_p(1)$ .

*Proof.* Choose an arbitrary  $\epsilon > 0$ . We show that for some  $C_1 > 0$ ,  $\sup_n \mathbb{P}[|b_n| > C_1] < \epsilon$ . Since  $a_n = \mathbb{E}[|b_n| | \mathcal{Z}] = O_p(1)$ , there exists a  $C_2 > 0$  such that  $\sup_n \mathbb{P}[a_n > C_2] < \epsilon/2$ . For this value of  $C_2$ , choose  $C_1 > 0$  such that  $C_2/C_1 \leq \epsilon/2$ . We then have

$$\begin{aligned} \sup_n \mathbb{P}[|b_n| > C_1] & = \sup_n \int \mathbb{P}[|b_n| > C_1 | \mathcal{Z}] (I(a_n > C_2) + I(a_n \leq C_2)) dP_{\mathcal{Z}} \\ & \leq \sup_n \mathbb{P}[a_n > C_2] + \frac{1}{C_1} \sup_n \int a_n I(a_n \leq C_2) dP_{\mathcal{Z}} \leq \frac{\epsilon}{2} + \frac{C_2}{C_1} \leq \epsilon. \quad \square \end{aligned}$$

#### APPENDIX C. APPROXIMATIONS

Suppressing  $\theta$ , let  $u_i = m_i - M_i$ ,  $u_i^* = \mathcal{V}_i^{-1} u_i$ ,  $G_i^* = \mathcal{V}_i^{-1} m_{\theta i}$ , and  $V_i^* = m_{\theta i}^* - G_i^* = \mathcal{V}_i^{-1} V_i$ .

**C.1.  $\check{\mathcal{V}}_i$  and  $\mathcal{V}_i$ .** Let  $\hat{u}_i(\theta) = \sum_j w_{ij} u_j(\theta)$ ,  $\hat{M}_i(\theta) = \sum_j w_{ij} M_j(\theta)$ ,  $\check{\mathcal{V}}_i(\theta) = \sum_j w_{ij} (M_j(\theta) - \hat{M}_j(\theta)) (M_j(\theta) - \hat{M}_j(\theta))^\top + \sum_j w_{ij} \mathcal{V}_j(\theta)$ . In the lemmas below we will frequently drop the argument  $\theta$  when this can be done without losing clarity.

**Lemma C1.**  $\exists \epsilon > 0 : \mathbb{P}[\max_i \|\check{\mathcal{V}}_i^{-1}\| > 1/\epsilon] = 0$ .

*Proof.* For sufficiently small  $\epsilon > 0$ ,  $\|\check{\mathcal{V}}_i^{-1}\| \leq \|(\sum_j w_{ij} \mathcal{V}_j)^{-1}\| \leq 1/\underline{\text{ev}}(\mathcal{V}_i) < 1/\epsilon$  a.s. by assumption **A**.  $\square$

**Lemma C2.**  $\max_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\| = o_p((\log n)^2 n^{1/2}/k)$ .

*Proof.* Note that

$$\begin{aligned} \hat{\mathcal{V}}_i - \check{\mathcal{V}}_i &= \sum_j w_{ij} (u_j u_j^\top - \mathcal{V}_j) - \sum_j w_{ij} (u_j \hat{u}_j^\top + \hat{u}_j u_j^\top) + \sum_j w_{ij} \hat{u}_j \hat{u}_j^\top \\ &\quad + \sum_j w_{ij} ((u_j - \hat{u}_j)(M_j - \hat{M}_j)^\top + (M_j - \hat{M}_j)(u_j - \hat{u}_j)^\top). \end{aligned} \quad (26)$$

We show that the desired uniform rate holds for RHS<sub>1</sub> and RHS<sub>2</sub> in (26), where similar arguments can be used for RHS<sub>3</sub>. For RHS<sub>1</sub>, take  $p = 2$ ,  $p^* = 4$  in **B9**. For RHS<sub>2</sub>, note that

$$\begin{aligned} \max_i \left\| \sum_j w_{ij} (u_j \hat{u}_j^\top + \hat{u}_j u_j^\top) \right\| &\leq 2 \max_i |\hat{u}_i| \max_i \sum_j w_{ij} |u_j| = o_p((n^{1/4}/k)(n/k)^{1/4}(\log n)^2) \\ &= o_p(n^{1/2} k^{-5/4} (\log n)^2), \end{aligned}$$

where the penultimate equality follows by taking  $p = 4$ ,  $p^* = 8$  in **B8** and **B9**.  $\square$

**Lemma C3.**  $\max_i \|\hat{\mathcal{V}}_i^{-1}\| = O_p(1)$ .

*Proof.* Follows from lemmas **C1** and **C2** since  $\underline{\text{ev}}(\hat{\mathcal{V}}_i) \geq \underline{\text{ev}}(\check{\mathcal{V}}_i) - \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|$ .  $\square$

**Lemma C4.** (i) For  $1 \leq p \leq 2$ ,  $\sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^p = O_p(n/k^{p/2})$  and (ii) for  $p > 2$ ,  $\sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^p = o_p((\log n)^{2p-4} n^{p/2}/k^{p-1})$ .

*Proof.* For (i) each RHS term in the expansion (26) can be dealt with separately. We only work on the dominant term, i.e. RHS<sub>1</sub>, where the other terms follow similarly. We have for some  $C < \infty$ ,

$$\sum_i \mathbb{E} \left\| \sum_j w_{ij} (u_j u_j^\top - \mathcal{V}_j) \right\|^p \stackrel{\text{Burkholder}^{26}}{\leq} C \sum_i \left( \sum_j \mathbb{E}[w_{ij}^2 \|u_j\|^4] \right)^{p/2} \leq \frac{C^2}{k^{p/2}} \sum_i \left( \sum_j \mathbb{E}[w_{ij} \|u_j\|^4] \right)^{p/2} = O\left(\frac{n}{k^{p/2}}\right).$$

(ii) then follows from (i) and **C3**.  $\square$

**Lemma C5.**  $\sum_i \hat{G}_i^\top (\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}) m_i = \begin{cases} o_p(\rho), & \theta = \theta_0, \\ o_p(\rho + n\lambda^2), & \theta \neq \theta_0. \end{cases}$

<sup>26</sup>See theorem 15.18 of Davidson (1994) and comments following it.

*Proof.* It suffices to show that

$$\sum_{ij} w_{ij} V_j^\top (\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}) u_i = o_p(\rho), \quad (27)$$

$$\sum_{ij} w_{ij} G_j^\top (\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}) u_i = o_p(\rho), \quad (28)$$

$$\sum_{ij} w_{ij} V_j^\top (\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}) M_i = o_p(\rho), \quad (29)$$

$$\sum_{ij} w_{ij} G_j^\top (\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}) M_i = o_p(n\lambda^2). \quad (30)$$

The LHS in (30) is bounded by  $\lambda^2 \max_i \|\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}\| \sum_{ij} w_{ij} \|\tilde{G}_j\| \|\tilde{M}_i\| = o_p(n\lambda^2)$  by lemmas **C1**, **C2**, and **C3**. The LHS of (29) is by the Schwarz inequality plus lemmas **C1**, **C3** bounded by

$$O_p(\lambda) \sqrt{\sum_i \left\| \sum_j w_{ij} V_j \right\|^2 \|\tilde{M}_i\|^2 \sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^2} = O_p(\lambda \sqrt{n/k} \sqrt{n/k}) = O_p(\rho \sqrt{n/k}) = o_p(\rho),$$

where the first equality follows from lemma **C4** and from taking expectations of the first factor under the square root. Letting  $\hat{G}_i = \sum_j w_{ij} \tilde{G}_j$ , the LHS in (28) equals

$$-\lambda \sum_i \hat{G}_i^\top \check{\mathcal{V}}_i^{-1} (\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1} u_i + \lambda \sum_i \hat{G}_i^\top \hat{\mathcal{V}}_i^{-1} (\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1} (\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1} u_i, \quad (31)$$

The second term in (31) is by **C1**, **C3** bounded in norm by a constant times

$$\begin{aligned} \lambda \sum_i \|\hat{G}_i\| \|u_i\| \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^2 &\stackrel{\text{H\"older}}{\leq} \lambda \left( \sum_i (\|\hat{G}_i\| \|u_i\|)^{p/(p-1)} \right)^{1-1/p} \left( \sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^{2p} \right)^{1/p} \\ &\stackrel{\text{C4}}{=} o_p(\lambda n^{1-1/p} (\log n)^{4-4/p} n/k^{2-1/p}) = o_p(\rho n^{3/2-1/p} (\log n)^{4-4/p} / k^{2-1/p}). \end{aligned}$$

Take  $p = 4/3$ . A bound for the norm of the first term in (31) can be found by expanding  $\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i$  using (26). For the first term in the expansion (26), we get

$$\lambda^2 \mathbb{E} \left\| \sum_{ij} w_{ij} \hat{G}_i^\top \check{\mathcal{V}}_i^{-1} (u_j u_j^\top - \mathcal{V}_j) \check{\mathcal{V}}_i^{-1} u_i \right\|^2 \stackrel{\text{C1}}{\leq} C \lambda^2 \sum_{ij} \mathbb{E} \left[ w_{ij}^2 \|\hat{G}_i\|^2 \|u_j\|^4 \|u_i\|^2 \right] = O(n\lambda^2/k) = o(\rho^2).$$

The other expansion terms can be dealt with similarly. Finally, the LHS in (27) equals

$$-\sum_{ij} w_{ij} V_j^\top \check{\mathcal{V}}_i^{-1} (\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1} u_i + \sum_{ij} w_{ij} V_j^\top \hat{\mathcal{V}}_i^{-1} (\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1} (\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1} u_i. \quad (32)$$

The second term in (32) is by the Hölder inequality bounded in norm by a constant times

$$\begin{aligned} &\left\{ \sum_i \left( \left\| \sum_j w_{ij} V_j \right\| \|u_i\| \right)^{p/(p-1)} \right\}^{1-1/p} \left\{ \sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^{2p} \right\}^{1/p} \\ &= o_p((n/k)^{1-1/p} (n/k^{2-1/p}) (\log n)^{4-4/p}) = o_p(n^{2-1/p} (\log n)^{4-4/p} / k^{3-2/p}) = o_p(\rho), \end{aligned}$$

for sufficiently large  $p$ , where the first equality follows with a Burkholder type argument like the one used earlier and lemma **C4**. For the first term in (32) we can once again use the expansion (26)

and deal with each RHS term in the expansion separately. For RHS<sub>1</sub> in (26) we get

$$\begin{aligned} & \mathbb{E} \left\| \sum_{ijt} w_{ij} w_{it} V_j^\top \check{\mathcal{V}}_i^{-1} (u_t u_t^\top - \mathcal{V}_t) \check{\mathcal{V}}_i^{-1} u_i \right\| \\ & \leq \mathbb{E} \left\| \sum_{ij} w_{ij}^2 V_j^\top \check{\mathcal{V}}_i^{-1} (u_j u_j^\top - \mathcal{V}_j) \check{\mathcal{V}}_i^{-1} u_i \right\| + \mathbb{E} \left\| \sum_{ij:t \neq j} w_{ij} w_{it} V_j^\top \check{\mathcal{V}}_i^{-1} (u_t u_t^\top - \mathcal{V}_t) \check{\mathcal{V}}_i^{-1} u_i \right\|. \end{aligned} \quad (33)$$

By the Schwarz inequality, the square of RHS<sub>2</sub> in (33) is bounded by a constant times  $\sum_{ij:t \neq j} \mathbb{E} [w_{ij}^2 w_{it}^2 \|V_j\|^2 \|u_t\|^4 \|u_i\|^2] = O(n/k^2) = o(\rho^2)$ . For RHS<sub>1</sub> in (33) note that for  $p = 4/3$ , some  $C < \infty$ , and  $\mathcal{R}_{ni} = \sum_{j < i} w_{ij}^2 V_j^\top \check{\mathcal{V}}_i^{-1} (u_j u_j^\top - \mathcal{V}_j) \check{\mathcal{V}}_i^{-1}$  (the  $\sum_{j > i}$  case is analogous), by the Burkholder and Loève inequalities,

$$\begin{aligned} & \mathbb{E} \left\| \sum_i \mathcal{R}_{ni} u_i \right\|^p \leq C \mathbb{E} \left[ \sum_i \|\mathcal{R}_{ni} u_i\|^2 \right]^{p/2} \leq C^2 \sum_i \mathbb{E} \|\mathcal{R}_{ni} u_i\|^p \leq \frac{C^3}{k^p} \sum_i \sum_{j < i} \mathbb{E} [w_{ij} \|V_j\|^p \|u_j\|^{2p} \|u_i\|^p] \\ & \stackrel{\text{Hölder}}{\leq} \frac{C^3}{k^p} \sum_i \sum_{j < i} \mathbb{E} \left[ w_{ij} (\mathbb{E} [\|V_j\|^4 | z_j])^{1/3} (\mathbb{E} [\|u_j\|^4 | z_j])^{2/3} \mathbb{E} [\|u_i\|^{4/3} | z_i] \right] = O(n/k^p) = o(\rho^p). \quad \square \end{aligned}$$

**Lemma C6.** For any  $p \geq 1$ ,  $\sum_i \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^p = o_p(n)$ .

*Proof.* Since  $\|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^p = \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^{p-1} \|\check{\mathcal{V}}_i^{-1} (\check{\mathcal{V}}_i - \mathcal{V}_i) \mathcal{V}_i^{-1}\|$  it suffices by lemma C<sub>1</sub> to establish the result for  $p = 1$ . For some  $C < \infty$ ,

$$\begin{aligned} \frac{1}{C} \sum_i \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\| & \leq \sum_i \|\check{\mathcal{V}}_i - \mathcal{V}_i\| \leq \sum_{ij} w_{ij} \|\mathcal{V}_j - \mathcal{V}_i\| + \sum_{ij} w_{ij} \|M_j - \hat{M}_j\|^2 \\ & \leq o_p(n) + \sum_{ijt} w_{ij} w_{jt} \|M_j - M_t\|^2 = o_p(n), \end{aligned}$$

by lemmas B<sub>1</sub> and B<sub>4</sub>. □

**Lemma C7.**  $\sum_i \hat{G}_i^\top (\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) m_i = \begin{cases} o_p(\rho), & \theta = \theta_0, \\ o_p(\rho + n\lambda^2), & \theta \neq \theta_0. \end{cases}$

*Proof.* It suffices to show that

$$\sum_{ij} w_{ij} V_j^\top (\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) u_i = o_p(\rho), \quad (34)$$

$$\sum_{ij} w_{ij} G_j^\top (\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) u_i = o_p(\rho), \quad (35)$$

$$\sum_{ij} w_{ij} V_j^\top (\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) M_i = o_p(\rho), \quad (36)$$

$$\sum_{ij} w_{ij} G_j^\top (\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) M_i = o_p(n\lambda^2). \quad (37)$$

The LHS in (37) is by the Schwarz inequality and lemma C<sub>1</sub> bounded in norm by

$\lambda^2 \sqrt{\sum_i \|\hat{G}_i\|^2 \|\tilde{M}_i\|^2 \sum_i \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|} = o_p(n\lambda^2)$  by lemma C6. The expectation of the norm of the



LHS in (36) is by repeated application of the Schwarz inequality bounded by

$$\lambda \left\{ \sum_i \left( \sum_{jt} w_{ij} w_{tj} \|V_j\|^2 \|\tilde{M}_i\|^2 \right) \sum_i \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^4 \right\}^{1/4} = o_p(\sqrt{n}\lambda) = o_p(\rho),$$

by lemmas B4 and C6. The derivation for (34) and (35) is similar to that for (36).  $\square$

**Lemma C8.**  $\sum_i \hat{G}_i^\top (\hat{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) m_i = \begin{cases} o_p(\rho), & \theta = \theta_0, \\ o_p(\rho + n\lambda^2), & \theta \neq \theta_0. \end{cases}$

*Proof.* Combine lemmas C5 and C7.  $\square$

Let  $\hat{\hat{\alpha}}_i = \hat{G}_i^\top \hat{\mathcal{V}}_i^{-1} m_i$ ,  $\hat{\alpha}_i = \hat{G}_i^\top \mathcal{V}_i^{-1} m_i$ .

**Lemma C9.**  $n^{-1} \|(\sum_i \hat{\hat{\alpha}}_i)(\sum_i \hat{\hat{\alpha}}_i^\top) - (\sum_i \hat{\alpha}_i)(\sum_i \hat{\alpha}_i^\top)\| = o_p(\rho^2)$ .

*Proof.* The LHS is bounded by  $n^{-1} \|\sum_i (\hat{\hat{\alpha}}_i - \hat{\alpha}_i)\|^2 + 2n^{-1} \|\sum_i (\hat{\hat{\alpha}}_i - \hat{\alpha}_i)\| \|\sum_i \hat{\alpha}_i\| = o_p(n^{-1}\rho^2 + n\lambda^4) + o_p(\rho + n\lambda^4) = o_p(\rho^2)$  by lemma C8.  $\square$

**Lemma C10.**  $\sum_i \|\hat{\hat{\alpha}}_i\|^2 = O_p(\rho^2)$ .

*Proof.* The LHS is bounded by a constant times  $\lambda^2 \sum_i \|\hat{G}_i^\top \mathcal{V}_i^{-1} m_i\|^2 + \sum_i \|\sum_j w_{ij} V_j^\top \mathcal{V}_i^{-1} m_i\|^2 = O_p(n\lambda^2 + n/k) = O_p(\rho^2)$ .  $\square$

**Lemma C11.**  $\sum_i \|\hat{\hat{\alpha}}_i - \hat{\alpha}_i\|^2 = o_p(\rho^2)$ .

*Proof.* The LHS is bounded by a constant times

$$\lambda^2 \sum_i \|\hat{G}_i^\top (\hat{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) m_i\|^2 + \sum_i \left\| \sum_j w_{ij} V_j^\top (\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}) m_i \right\|^2 + \sum_i \left\| \sum_j w_{ij} V_j^\top (\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}) m_i \right\|^2. \quad (38)$$

The first term in (38) is bounded by  $\lambda^2 \sqrt{\sum_i \|\hat{G}_i\|^2 \|m_i\|^2 \sum_i \|\hat{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^2} = o_p(n\lambda^2)$  by lemmas C1, C4, C6. The expectation of RHS3 conditional on  $\mathcal{Z}$  is bounded by

$$\begin{aligned} & \sum_{ij} w_{ij}^2 \mathbb{E}[\|V_j\|^2 | z_j] \mathbb{E}[\|m_i\|^2 | z_i] \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^2 \\ & \leq \stackrel{\text{H\"older, Jensen}}{\leq} \frac{C_w}{k} \left( \sum_{ij} w_{ij} (\mathbb{E}[\|V_j\|^2 | z_j] \mathbb{E}[\|m_i\|^2 | z_i])^p \right)^{1/p} \left( \sum_i \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^{2p/(p-1)} \right)^{1-1/p}. \end{aligned}$$

Pick  $p$  no greater than two to obtain a rate of  $o_p(n/k)$  by lemmas B1, C6. RHS2 in (38) is by the Schwarz inequality bounded by  $\sqrt{\sum_i \|\sum_j w_{ij} V_j\|^4 \|m_i\|^4 \sum_i \|\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}\|^4}$ . By the Jensen inequality,  $\sum_i \mathbb{E} \|\sum_j w_{ij} V_j\|^4 \leq (C_w/k)^2 \sum_{ij} \mathbb{E} [w_{ij} \|V_j\|^4] = O(n/k^2)$  by lemma B1. Finally, using  $\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1} = \check{\mathcal{V}}_i^{-1} (\check{\mathcal{V}}_i - \hat{\mathcal{V}}_i) \check{\mathcal{V}}_i^{-1}$ , we have that  $\sum_i \|\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}\|^4 \leq \max_i \|\hat{\mathcal{V}}_i^{-1}\|^4 \max_i \|\check{\mathcal{V}}_i^{-1}\|^4 \sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^4 = o_p(n^2(\log n)^4/k^3)$ , by lemmas C1, C3, and C4.  $\square$

**Lemma C12.**  $\sum_i (\hat{\hat{\alpha}}_i \hat{\hat{\alpha}}_i^\top - \hat{\alpha}_i \hat{\alpha}_i^\top) = o_p(\rho^2)$ .

*Proof.* The LHS is bounded in norm by a constant times  $\sum_i \|\hat{\hat{\alpha}}_i - \hat{\alpha}_i\|^2 + \sqrt{\sum_i \|\hat{\hat{\alpha}}_i\|^2 \sum_i \|\hat{\alpha}_i - \hat{\hat{\alpha}}_i\|^2} = o_p(\rho^2)$  by lemmas C10, C11.  $\square$

Let  $\hat{\alpha}_{ij}^* = w_{ij} m_{\theta_j}^\top \hat{\mathcal{V}}_i^{-1} m_i$ ,  $\alpha_{ij}^* = w_{ij} m_{\theta_j}^\top \mathcal{V}_i^{-1} m_i$ .

**Lemma C13.**  $\sum_{ij} (\hat{\alpha}_{ij}^* \hat{\alpha}_{ji}^{*\top} - \alpha_{ij}^* \alpha_{ji}^{*\top}) = o_p(\rho^2)$ .

*Proof.* It suffices to establish that (i)  $\sum_{ij} \|\alpha_{ij}^*\|^2 = O_p(\rho^2)$  and that (ii)  $\sum_{ij} \|\hat{\alpha}_{ij}^* - \alpha_{ij}^*\|^2 = o_p(\rho^2)$ . The LHS in (i) is bounded by  $(C_w/k) \sum_{ij} w_{ij} \|m_{\theta_j}\|^2 \|\mathcal{V}_i^{-1}\|^2 \|m_i\|^2 = O_p(n/k)$ . The LHS in (ii) is bounded by

$$\begin{aligned} \frac{2C_w}{k} \sqrt{\sum_i \left( \sum_j w_{ij} \|m_{\theta_j}\|^2 \right)^2 \|m_i\|^4 \left( \sum_i \|\hat{\mathcal{V}}_i^{-1} - \check{\mathcal{V}}_i^{-1}\|^2 + \sum_i \|\check{\mathcal{V}}_i^{-1} - \mathcal{V}_i^{-1}\|^2 \right)} &\stackrel{\text{Schwarz, C6}}{\leq} \\ \frac{2C_w}{k} \sqrt{\sum_{ij} w_{ij} \|m_{\theta_j}\|^4 \|m_i\|^4 \left( \max_i \|\hat{\mathcal{V}}_i^{-1}\|^2 \max_i \|\check{\mathcal{V}}_i^{-1}\|^2 \sum_i \|\hat{\mathcal{V}}_i - \check{\mathcal{V}}_i\|^2 + o_p(n) \right)} &= o_p\left(\frac{n}{k}\right), \end{aligned}$$

by lemmas B1, C1, C3, C4. □

**Lemma C14.**  $\hat{D}^2 - \hat{D}^2 = o_p(\rho^2)$ .

*Proof.* Follows from lemmas C9, C12, C13. □

C.2.  $\mathcal{V}_i$ . Let  $M_i^*(\theta) = \mathcal{V}_i^{-1}(\theta) M_i(\theta)$ ,  $\tilde{M}_i^*(\theta) = \check{\mathcal{V}}_i^{-1}(\theta) \tilde{M}_i(\theta)$ , and  $u_i^*(\theta) = m_i^*(\theta) - M_i^*(\theta)$ .

**Lemma C15.**  $\sum_i \hat{G}_i^\top(\theta) u_i^*(\theta) = \sum_i G_i^\top(\theta) u_i^*(\theta) + \sum_{ij} w_{ij} V_j^\top(\theta) u_i^*(\theta) + o_p(\sqrt{n\lambda})$ .

*Proof.* Since  $\hat{G}_i = \sum_j w_{ij} m_{\theta_j} = G_i + \sum_j w_{ij} (G_j - G_i) + \sum_j w_{ij} V_j$ , we consider (omitting  $\theta$ )

$$\begin{aligned} \mathbb{E} \left\| \sum_{ij} w_{ij} (G_j - G_i)^\top u_i^* \right\|^2 &\leq \sum_{ijs} \mathbb{E} [w_{ij} w_{is} \|G_j - G_i\| \|G_s - G_i\| \|u_i^*\|^2] \\ &\leq \sum_{ijs} \mathbb{E} [w_{ij} w_{is} (\|G_j - G_i\|^2 + \|G_s - G_i\|^2) \|u_i^*\|^2] = 2 \sum_{ij} \mathbb{E} [w_{ij} \|G_j - G_i\|^2 \|u_i^*\|^2] \stackrel{\text{B1}}{=} o(n\lambda^2). \end{aligned} \quad \square$$

**Lemma C16.**  $\sum_i \hat{G}_i^\top(\theta) M_i^*(\theta) = n \mathbb{E} [G_1^\top(\theta) M_1^*(\theta)] + o_p(n\lambda^2 + \rho)$ .

*Proof.* The LHS equals (omitting  $\theta$ )

$$\lambda^2 \sum_i (\tilde{G}_i^\top \tilde{M}_i^* - \mathbb{E}[\tilde{G}_1^\top \tilde{M}_1^*]) + \lambda \sum_{ij} w_{ij} V_j^\top \tilde{M}_i^* + \lambda^2 \sum_{ij} w_{ij} (\check{G}_j - \check{G}_i)^\top \tilde{M}_i^*. \quad (39)$$

RHS1 in (39) is  $O_p(\sqrt{n\lambda^2})$  and RHS3 is  $o_p(n\lambda^2)$  by lemma B1. Finally, take the norm of RHS2, square and take its expectation conditional on  $\mathcal{Z}$  (using lemma B10) to obtain a bound of

$$\lambda^2 \sum_{ijt} w_{ij} w_{tj} \mathbb{E} [\|V_j\|^2 |z_j] \|\tilde{M}_i^*\| \|\tilde{M}_t^*\| \leq 2\lambda^2 \max_j \left( \sum_t w_{tj} \right) \sum_{ij} w_{ij} \mathbb{E} [\|V_j\|^2 |z_j] \|\tilde{M}_i^*\|^2 \stackrel{\text{B4}}{=} O_p(n\lambda^2). \quad \square$$

Recall the definition of  $\hat{\alpha}_i, \hat{\alpha}_i$ , defined prior to lemma C9 and let  $\alpha_i(\theta) = G_i^\top(\theta) \mathcal{V}_i^{-1}(\theta) m_i(\theta)$ .

**Lemma C17.**  $\sum_i \hat{\alpha}_i(\theta) - n \mathbb{E} \alpha_1(\theta) = o_p(n\lambda^2) + O_p(\rho)$ .

*Proof.* In view of lemmas C15, C16, it suffices to show that (omitting  $\theta$ ) (i)  $\sum_i G_i^\top u_i^* = O_p(\rho)$  and that (ii)  $\sum_{ij} w_{ij} V_j^\top u_i^* = O_p(\rho)$ . The expectation of the squared norm of the LHS in (i) is  $\lambda^2 \sum_i \mathbb{E} \|\tilde{G}_i^\top u_i^*\|^2 = O(n\lambda^2) = O(\rho^2)$  and the expectation of the squared norm of (ii) is bounded by twice  $\sum_{ij} \mathbb{E} [w_{ij}^2 \|V_j\|^2 \|u_i^*\|^2] = O(n/k) = O(\rho^2)$  by lemma B1.  $\square$

**Lemma C18.**  $\sum_i \hat{\alpha}_i(\theta) \sum_i \hat{\alpha}_i^\top(\theta) - n^2 \mathbb{E} \alpha_1(\theta) \mathbb{E} \alpha_1^\top(\theta) = o_p(\rho^4)$ .

*Proof.* Omitting  $\theta$ , the LHS is bounded in norm by  $\sum_i \|\hat{\alpha}_i - \mathbb{E} \alpha_1\|^2 + 2 \|\mathbb{E} \alpha_1\| \sum_i \|\hat{\alpha}_i - \mathbb{E} \alpha_1\| = o_p(\rho^4)$  by lemma C17 and the fact that  $\|\mathbb{E} \alpha_1\| = O(\lambda^2)$ .  $\square$

Let  $\hat{\alpha}_{V_i}(\theta) = \sum_j w_{ij} V_j^\top(\theta) m_i^*(\theta)$  and  $\hat{\alpha}_i^*(\theta) = (G_i(\theta) + \sum_j w_{ij} V_j(\theta))^\top m_i^*(\theta)$ .

**Lemma C19.**  $\sum_i (\hat{\alpha}_i(\theta) \hat{\alpha}_i^\top(\theta) - \alpha_i(\theta) \alpha_i^\top(\theta) - \hat{\alpha}_{V_i}(\theta) \hat{\alpha}_{V_i}^\top(\theta)) = o_p(\rho^2)$ .

*Proof.* We show separately that (i)  $\sum_i (\hat{\alpha}_i \hat{\alpha}_i^\top - \hat{\alpha}_i^* \hat{\alpha}_i^{*\top}) = o_p(\rho^2)$  and that (ii)  $\sum_i (\hat{\alpha}_i^* \hat{\alpha}_i^{*\top} - \alpha_i \alpha_i^\top - \hat{\alpha}_{V_i} \hat{\alpha}_{V_i}^\top) = o_p(\rho^2)$ . The LHS in (i) is bounded by twice  $\sum_i \|\hat{\alpha}_i - \hat{\alpha}_i^*\|^2 + \sqrt{\sum_i \|\hat{\alpha}_i - \hat{\alpha}_i^*\|^2 \sum_i \|\hat{\alpha}_i^*\|^2}$ . Now,  $\sum_i \|\hat{\alpha}_i^*\|^2 / 2 \leq \lambda^2 \sum_i \|\tilde{G}_i\|^2 \|m_i^*\|^2 + \sum_i \|\sum_j w_{ij} V_j\|^2 \|m_i^*\|^2 = O_p(\rho^2)$ . Further, by the Jensen inequality  $\sum_i \|\hat{\alpha}_i - \hat{\alpha}_i^*\|^2 \leq \lambda^2 \sum_{ij} w_{ij} \|\tilde{G}_j - \tilde{G}_i\|^2 \|m_i^*\|^2 = o_p(\rho^2)$  by lemma B1. The norm of the LHS in (ii) is twice  $\|\sum_i \alpha_i \hat{\alpha}_{V_i}^\top\|$ . Thus, using lemma B10,

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_i \alpha_i \hat{\alpha}_{V_i}^\top \right\|^2 \middle| \mathcal{Z} \right] &= \lambda^2 \mathbb{E} \left[ \left\| \sum_{ij} w_{ij} \tilde{G}_i^\top m_i^* m_i^{*\top} V_j \right\|^2 \middle| \mathcal{Z} \right] \\ &\leq \lambda^2 \sum_{ijt} w_{ij} w_{tj} \|\tilde{G}_i\| \|\tilde{G}_t\| \mathbb{E} [\|\tilde{m}_i^*\|^2 \|\tilde{m}_t^*\|^2 | z_i, z_t] \mathbb{E} [\|V_j\|^2 | z_j] \\ &\leq 2\lambda^2 \sum_{ijt} w_{ij} w_{tj} \|\tilde{G}_i\|^2 \mathbb{E} [\|\tilde{m}_i^*\|^4 | z_i] \mathbb{E} [\|V_j\|^2 | z_j] \stackrel{\text{B4}}{=} O_p(\lambda^2 n) = o_p(\rho^4). \quad \square \end{aligned}$$

**Lemma C20.**  $\sum_i \hat{\alpha}_{V_i}(\theta) \hat{\alpha}_{V_i}^\top(\theta) - \sum_{ij} w_{ij}^2 \mathbb{E} [V_j^\top(\theta) m_i^*(\theta) m_i^{*\top}(\theta) V_j(\theta) | z_i, z_j] = o_p(\rho^2)$ .

*Proof.* Note first that the LHS is

$$\left( \sum_{ij} w_{ij}^2 (V_j^\top m_i^* m_i^{*\top} V_j - \mathbb{E} [V_j^\top m_i^* m_i^{*\top} V_j | z_i, z_j]) \right) + \sum_{ij} \sum_{t \neq j} w_{ij} w_{it} V_j^\top m_i^* m_i^{*\top} V_t. \quad (40)$$

Using lemma B10, we have for the second term in (40) that for some  $C < \infty$ ,

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{ij} \sum_{t \neq j} w_{ij} w_{it} V_j^\top m_i^* m_i^{*\top} V_t \right\|^2 \middle| \mathcal{Z} \right] &\leq C \sum_{ijts} w_{ij} w_{it} w_{sj} w_{st} \mathbb{E} [\|V_j\|^2 | z_j] \mathbb{E} [\|V_t\|^2 | z_t] \mathbb{E} [\|m_i^*\|^2 \|m_s^*\|^2 | z_i, z_s] \\ &\stackrel{\text{B4}}{\leq} O_p(1/k) \sum_{ij} w_{ij} \mathbb{E} [\|V_j\|^2 | z_j]^2 \mathbb{E} [\|m_i^*\|^4 | z_i] = O_p(n/k) = o_p(\rho^4). \end{aligned}$$

Again by lemma B10, we have that the expectation of the squared norm of the first term in (40) given  $\mathcal{Z}$  is for some  $C < \infty$  bounded by

$$C \sum_{ijt} w_{ij}^2 w_{it}^2 \mathbb{E} [\|V_j\|^2 | z_j] \mathbb{E} [\|V_t\|^2 | z_t] \mathbb{E} [\|m_i^*\|^4 | z_i] \leq \frac{C^2}{k^2} \sum_{ij} w_{ij} (\mathbb{E} [\|V_j\|^2 | z_j])^2 \mathbb{E} [\|m_i^*\|^4 | z_i] = o_p(\rho^4). \quad \square$$

**Lemma C21.**  $\sum_{ij} w_{ij} w_{ji} m_{\theta_j}^\top(\theta) m_i^*(\theta) m_j^{*\top}(\theta) m_{\theta_i}(\theta) - \sum_{ij} w_{ij} w_{ji} \mathbb{E}[V_j^\top(\theta) m_i^*(\theta) m_j^{*\top}(\theta) V_i(\theta) | z_i, z_j] = o_p(\rho^2)$ .

*Proof.* Write the LHS as

$$\begin{aligned} & \sum_{ij} w_{ij} w_{ji} (m_{\theta_j}^\top m_i^* m_j^{*\top} m_{\theta_i} - \mathbb{E}[m_{\theta_j}^\top m_i^* m_j^{*\top} m_{\theta_i} | z_i, z_j]) + \lambda \sum_{ij} w_{ij} w_{ji} \mathbb{E}[V_j^\top m_i^* m_j^{*\top} \tilde{G}_i | z_i, z_j] \\ & \quad + \lambda \sum_{ij} w_{ij} w_{ji} \mathbb{E}[\tilde{G}_j^\top m_i^* m_j^{*\top} V_i | z_i, z_j] + \lambda^4 \sum_{ij} w_{ij} w_{ji} \tilde{G}_j^\top \tilde{M}_i^* \tilde{M}_j^{*\top} \tilde{G}_i. \end{aligned} \quad (41)$$

By lemma B6(iii) the first term in (41) is  $O_p(n^{1/2}k^{-1}) = o_p(\rho^2)$ . The expectation of the fourth term in (41) is bounded by  $(n^4 \lambda^4 C_w / k) \mathbb{E}[w_{12} \|\tilde{M}_1^*\| \|\tilde{M}_2^*\| \|\tilde{G}_1\| \|\tilde{G}_2\|] \stackrel{B5}{=} O(n\lambda^4/k) = o(\rho^2)$ . Finally, the second and third terms in (41) are  $O_p(n\lambda^2/k) = o_p(\rho^2)$ , which can be similarly established.  $\square$

**Lemma C22.**  $\hat{D}^2(\theta) - \tilde{D}^2(\theta) = o_p(\rho^2)$ .

*Proof.* Omitting the  $\theta$  argument, the LHS is

$$\begin{aligned} & \underbrace{\sum_i (\hat{\alpha}_i \hat{\alpha}_i^\top - \alpha_i \alpha_i^\top - \hat{\alpha}_{V_i} \hat{\alpha}_{V_i}^\top)}_{\stackrel{C19}{=} o_p(\rho^2)} - n^{-1} \underbrace{\left( \sum_i \hat{\alpha}_i \sum_i \hat{\alpha}_i^\top - n^2 \mathbb{E} \alpha_1 \mathbb{E} \alpha_1^\top \right)}_{\stackrel{C18}{=} o_p(\rho^4/n) = o_p(\rho^2)} + \underbrace{\sum_i (\alpha_i \alpha_i^\top - \mathbb{E}[\alpha_1 \alpha_1^\top])}_{= O_p(\sqrt{n}\lambda^2)} \\ & + \underbrace{\left( \sum_i \hat{\alpha}_{V_i} \hat{\alpha}_{V_i}^\top - \sum_{ij} w_{ij}^2 \mathbb{E}[V_j m_i^* m_i^{*\top} V_j^\top | z_i, z_j] \right)}_{\stackrel{C20}{=} o_p(\rho^2)} + \underbrace{\sum_{ij} w_{ij} w_{ji} (m_{\theta_j}^\top m_i^* m_j^{*\top} m_{\theta_i} - \mathbb{E}[V_j^\top m_i^* m_j^{*\top} V_i | z_i, z_j])}_{\stackrel{C21}{=} o_p(\rho^2)} \\ & = o_p(\rho^2). \quad \square \quad (42) \end{aligned}$$

**Lemma C23.**  $\tilde{D}^2(\theta) - n \text{Var}[\alpha_1(\theta)] = O_p(n/k)$ .

*Proof.* The expectation of the norm of the LHS is bounded by a constant times

$$\sum_{ij} w_{ij}^2 \mathbb{E}[\|V_j\|^2 | z_j] \mathbb{E}[\|m_i^*\|^2 | z_i] = O_p(n/k) \text{ by lemma B5.} \quad \square$$

**Lemma C24.**  $\tilde{D}(\theta) = O_p(\rho)$ .

*Proof.* Apply lemma C23 and note that  $n \text{Var}[\alpha_1] = O(n\lambda^2)$ , such that  $\tilde{D}^2 = O_p(\rho^2)$ .  $\square$

**Lemma C25.**  $\tilde{D}^{-1} = O_p(\rho^{-1})$ .

*Proof.* Recall that  $\ell_{ij} = w_{ij} V_j^\top m_i^*$ , such that  $\tilde{D}^2 = n \text{Var}[\alpha_1] + \sum_{ij} \mathbb{E}[\ell_{ij} \ell_{ij}^\top | \mathcal{Z}] + \sum_{ij} \mathbb{E}[\ell_{ij} \ell_{ji}^\top | \mathcal{Z}]$ . Now,

$$\frac{1}{\|\tilde{D}^{-2}\|} \geq \max \left\{ \underline{\text{ev}}(n \text{Var}[\alpha_1]), \underline{\text{ev}} \left( \sum_{ij} \mathbb{E}[\ell_{ij} \ell_{ij}^\top | \mathcal{Z}] + \sum_{ij} \mathbb{E}[\ell_{ij} \ell_{ji}^\top | \mathcal{Z}] \right) \right\}$$

Now,  $\underline{\text{ev}}(n \text{Var}[\alpha_1]) \geq n\lambda^2 \underline{\text{ev}}(\mathbb{E}[\tilde{G}_1^\top \mathcal{V}_1^{-1} \tilde{G}_1]) \geq Cn\lambda^2$  for some  $C > 0$  by assumption B. We show that for any fixed  $\|c\| = 1$  and some  $C > 0$ ,

$$\sum_{ij} c^\top \mathbb{E}[\ell_{ij} \ell_{ij}^\top | \mathcal{Z}] c + \sum_{ij} c^\top \mathbb{E}[\ell_{ij} \ell_{ji}^\top | \mathcal{Z}] c \geq Cn/k - o_p(n/k). \quad (43)$$

First, letting  $A_{ij} = w_{ij} \mathbb{E}[m_i^* c^\top V_i^\top | z_i]$ , note that the second term in (43) is equal to

$$\begin{aligned} \sum_{ij} w_{ij} w_{ji} \mathbb{E}[c^\top V_j^\top m_i^* m_j^{*\top} V_i c | z_i, z_j] &= \sum_{ij} \text{tr}(A_{ij} A_{ji}) = \frac{1}{2} \sum_{ij} \text{tr}(A_{ij} A_{ji} + A_{ji}^\top A_{ij}^\top) \\ &= \frac{1}{2} \sum_{ij} \text{tr}\{(A_{ij} + A_{ji}^\top)(A_{ij} + A_{ji}^\top)^\top - (A_{ij} A_{ij}^\top + A_{ji} A_{ji}^\top)\} \geq - \sum_{ij} \text{tr}(A_{ij} A_{ij}^\top) \\ &= - \sum_{ij} w_{ij}^2 \text{tr}(\mathbb{E}[m_i^* c^\top V_i^\top | z_i] \mathbb{E}[V_i c m_i^{*\top} | z_i]). \end{aligned} \quad (44)$$

The first term in (43) equals

$$\begin{aligned} \sum_{ij} w_{ij}^2 \mathbb{E}[(c^\top V_j^\top m_i^*)^2 | z_i, z_j] &= \text{tr}\left(\sum_{ij} w_{ij}^2 \mathbb{E}[m_i^* m_i^{*\top} V_j c c^\top V_j^\top | z_i, z_j]\right) \\ &\stackrel{\text{B1, B6(iii)}}{=} \sum_{ij} w_{ij}^2 \text{tr}(\mathbb{E}[m_i^* m_i^{*\top} | z_i] \mathbb{E}[V_j c c^\top V_j^\top | z_i]) + o_p(n/k). \end{aligned} \quad (45)$$

Noting that  $m_i^* = u_i^* + \lambda \tilde{M}_i^*$ , it follows that the sum of the RHS in (44) and that in (45) is bounded below by

$$\begin{aligned} \sum_{ij} w_{ij}^2 \text{tr}(\text{Var}[u_i^* | z_i] \text{Var}[V_i c | z_i] - \text{Cov}[u_i^*, V_i c | z_i] \text{Cov}[V_i c, u_i^* | z_i]) &+ o_p(n/k) \\ &= \sum_{ij} w_{ij}^2 \text{tr}\{\mathcal{V}_i^{-2} (\text{Var}[u_i | z_i] \text{Var}[V_i c | z_i] - \text{Cov}[u_i, V_i c | z_i] \text{Cov}[V_i c, u_i | z_i])\} + o_p(n/k) \\ &\geq \frac{C_w n}{k} \mathbb{E} \text{tr}\{\mathcal{V}_1^{-2} (\text{Var}[u_1 | z_1] \text{Var}[V_1 c | z_1] - \text{Cov}[u_1, V_1 c | z_1] \text{Cov}[V_1 c, u_1 | z_1])\} - o_p(n/k) \\ &\geq Cn/k - o_p(n/k), \end{aligned}$$

for some  $C > 0$  by assumption A.  $\square$

#### APPENDIX D. MARTINGALE DIFFERENCES

**Lemma D1.** Let  $\{X_{ni}, \mathcal{F}_{ni}\}$  be a martingale difference array for which

(i)  $\sup_n \mathbb{E}[\max_{i \leq n} X_{ni}^2] < \infty$ , (ii)  $\sum_{i=1}^n X_{ni}^2 \xrightarrow{p} 1$  and (iii)  $\max_{i \leq n} |X_{ni}| \xrightarrow{p} 0$ . Then  $\sum_{i=1}^n X_{ni} \xrightarrow{d} N(0, 1)$ .

*Proof.* This is Davidson (1994), theorem 24.3, where the only difference is that Davidson's sufficient condition that the unconditional variances sum to one is replaced with condition (i). Since in Davidson's proof of theorem 24.3 the finite unconditional variances requirement is only used to verify (i), the lemma statement holds trivially.  $\square$

**Lemma D2.**  $\max_i \|\tilde{D}^{-1} B_i\| = o_p(1)$ .

*Proof.* Since  $\|\tilde{D}^{-1}\| = O_p(\rho^{-1})$  by lemma C25, it suffices to show that  $\max_i \|B_i\| = o_p(\rho)$ . Pick any  $\epsilon > 0$ . Then for any  $2 < p \leq 4$  (suppressing  $\theta_0$ ),

$$\mathbb{P}(\max_i \|B_i\| > \epsilon \rho) \leq \sum_i \mathbb{P}(\|B_i\| > \epsilon \rho) \stackrel{\text{Markov}}{\leq} \rho^{-p} \epsilon^{-p} \sum_i \mathbb{E} \|B_i\|^p.$$

So by the Minkowski inequality it suffices to show that

$$n\mathbb{E}\|\zeta_1\|^p = o(\rho^p), \quad \sum_i \mathbb{E}\left\|\sum_{j<i} \ell_{ij}\right\|^p = o(\rho^p), \quad \sum_i \mathbb{E}\left\|\sum_{j<i} \ell_{ji}\right\|^p = o(\rho^p). \quad (46)$$

The LHS in the first condition in (46) is  $O(n\lambda^p) = o(\rho^p)$ . The second and third conditions in (46) are similar; we only show the second. We have for some fixed  $C_p < \infty$ ,

$$\begin{aligned} \sum_i \mathbb{E}\left\|\sum_{j<i} \ell_{ij}\right\|^p &= \sum_i \mathbb{E}\left\|\sum_{j<i} w_{ij} V_j^\top m_i^*\right\|^p \stackrel{\text{Burkholder}}{\leq} C_p \sum_i \mathbb{E}\left[\sum_{j<i} w_{ij}^2 \|m_i^*\|^2 \|V_j\|^2\right]^{p/2} \\ &\stackrel{\text{Jensen}}{\leq} C_p \sum_i \mathbb{E}\left[i^{p/2-1} \sum_{j<i} w_{ij}^p \|m_i^*\|^p \|V_j\|^p\right] \leq C_p n^{p/2+1} \mathbb{E}[w_{12}^p \|m_1^*\|^p \|V_2\|^p] \stackrel{\text{B5}}{=} O(n^{p/2} k^{1-p}) = o(\rho^p). \end{aligned}$$

□

**Lemma D3.**  $\sum_i B_i B_i^\top - \tilde{D}^2(\theta_0) = o_p(\rho^2)$ .

*Proof.* Noting that  $\tilde{D}^2(\theta_0) = n\mathbb{E}[\zeta_1 \zeta_1^\top] + \sum_i \sum_{j<i} \mathbb{E}[\tilde{\zeta}_{ij} \tilde{\zeta}_{ij}^\top | \mathcal{Z}]$ , we have (omitting the  $\theta_0$ -argument)

$$\begin{aligned} \sum_i B_i B_i^\top &= \tilde{D}^2 + \sum_i (\zeta_i \zeta_i^\top - \mathbb{E}[\zeta_i \zeta_i^\top]) + \sum_i \sum_{j<i} (\zeta_i \tilde{\zeta}_{ij}^\top + \tilde{\zeta}_{ij} \zeta_i^\top) \\ &\quad + \sum_i \sum_{j<i} \{\tilde{\zeta}_{ij} \tilde{\zeta}_{ij}^\top - \mathbb{E}[\tilde{\zeta}_{ij} \tilde{\zeta}_{ij}^\top | \mathcal{Z}]\} + 2 \sum_i \sum_{j<i} \sum_{t<j} \tilde{\zeta}_{ij} \tilde{\zeta}_{it}^\top. \end{aligned} \quad (47)$$

For RHS2 in (47), note that  $\sum_i \text{Var}[\text{vec}(\zeta_i \zeta_i^\top)] = O(n\lambda^4) = o(\rho^4)$ . For RHS3, note that (the case  $\sum_{j>i}$  is symmetric)

$$\sum_i \sum_{j<i} w_{ij} \tilde{\zeta}_i m_i^{*\top} V_j + \sum_i \sum_{j<i} w_{ji} \tilde{\zeta}_i m_j^{*\top} V_i \stackrel{\text{B6(ii),(iv)}}{=} o_p(n^{3/4} k^{-1/4} \lambda) + O_p(n^{1/2} \lambda) = o_p(\rho^2),$$

where lemma B6 is applied element-wise with the  $c_{ij}$ -weights taken as  $c_{ij} = I(j < i)$ . RHS4 in (47) is  $o_p(\rho^2)$  because using lemma B6(iii),

$$\begin{aligned} \sum_i \sum_{j<i} (\ell_{ij} \ell_{ij}^\top - \mathbb{E}[\ell_{ij} \ell_{ij}^\top | \mathcal{Z}]) &= \sum_i \sum_{j<i} w_{ij}^2 (V_j^\top m_i^* m_i^{*\top} V_j - \mathbb{E}[V_j^\top m_i^* m_i^{*\top} V_j | z_i, z_j]) = O_p(n^{1/2}/k) = o_p(\rho^2), \\ \sum_i \sum_{j<i} (\ell_{ij} \ell_{ji}^\top - \mathbb{E}[\ell_{ij} \ell_{ji}^\top | \mathcal{Z}]) &= \sum_i \sum_{j<i} w_{ij} w_{ji} (V_j^\top m_i^* m_j^{*\top} V_i - \mathbb{E}[V_j^\top m_i^* m_j^{*\top} V_i | z_i, z_j]) = O_p(n^{1/2}/k) = o_p(\rho^2). \end{aligned}$$

Finally, RHS5 in (47) is  $o_p(\rho^2)$  because for arbitrary elements  $\tilde{\zeta}_{ij}^*, \tilde{\zeta}_{it}^{**}$  of  $\tilde{\zeta}_{ij}, \tilde{\zeta}_{it}$ ,  $\mathbb{E}[\sum_i \sum_{j<i} \sum_{t<j} \tilde{\zeta}_{ij}^* \tilde{\zeta}_{it}^{**}]^2 = \sum_i \sum_{j<i} \sum_{t<j} \mathbb{E}[\tilde{\zeta}_{ij}^{*2} \tilde{\zeta}_{it}^{**2}] = o(\rho^4)$ , by lemma B5. □

**Lemma D4.**  $\sup_n \mathbb{E}[\max_{i \leq n} \|B_i^\top(\theta_0) \tilde{D}^{-1}(\theta_0)\|^2] < \infty$ .

*Proof.* Omit  $\theta_0$  throughout. Note that

$$\begin{aligned} \tilde{D}^{-1} \sum_i \mathbb{E}[B_i B_i^\top | \mathcal{Z}] \tilde{D}^{-1} &= \tilde{D}^{-1} \left( \tilde{D}^2 + \sum_i (G_i^\top \mathcal{V}_i^{-1} G_i - \mathbb{E}[G_1^\top \mathcal{V}_1^{-1} G_1]) \right) \tilde{D}^{-1} \\ &= I + \sum_i \tilde{D}^{-1} (G_i^\top \mathcal{V}_i^{-1} G_i - \mathbb{E}[G_1^\top \mathcal{V}_1^{-1} G_1]) \tilde{D}^{-1}. \end{aligned}$$

For  $\lambda = 0$  the result is immediate. For  $\lambda > 0$ , note that for some  $C < \infty$ ,  $\|\tilde{D}^{-2}\| \leq \|(nV[G_1 m_1^*])^{-1}\| \leq C/n\lambda^2$ . Thus,

$$\max_{i \leq n} \|B_i^\top \tilde{D}^{-1}\|^2 \leq \sum_i \|B_i^\top \tilde{D}^{-1}\|^2 \leq 1 + \frac{C}{n} \left\| \sum_i (\tilde{G}_i^\top \mathcal{V}_i^{-1} \tilde{G}_i - \mathbb{E}[\tilde{G}_1^\top \mathcal{V}_1^{-1} \tilde{G}_1]) \right\| \leq 1 + C\mathbb{E}\|\mathcal{V}_1^{-1} \tilde{G}_1\|^4 < \infty. \quad \square$$

## APPENDIX E. LOCAL ALTERNATIVES

**E.1. Uniform Lemmas.** Let  $\mathcal{N}$  be a convex compact set including  $\theta_0$  and  $|\mathcal{N}| = \sup_{\theta \in \mathcal{N}} \|\theta - \theta_0\|$ .

**Lemma E1.** For  $p_1 > 0$  let  $f_i = f_i(\theta) = f(z_i, \theta)$  be such that  $\mathbb{E}|\sup_{\mathcal{N}} f_i|^{p_1 p_2} < \infty$  and  $a_i = a_i(\theta) > 0$  be such that  $\mathbb{E} \sup_{\mathcal{N}} a_i^{p_2/(p_2-1)} < \infty$  for some  $p_2 \geq 1$ . If  $\lim_{|\mathcal{N}| \rightarrow 0} \mathbb{E}[\sup_{\mathcal{N}} f_i - \inf_{\mathcal{N}} f_i]^{p_2} = 0$  then

$$\lim_{|\mathcal{N}| \rightarrow 0} n\mathbb{E} \left[ w_{12} \sup_{\theta \in \mathcal{N}} (a_1 |f_1 - f_2|) \right] = 0.$$

*Proof.* Let sup and inf be taken over  $\theta \in \mathcal{N}$  for the remainder of this lemma. Now,

$$n\mathbb{E} [w_{12} \sup (a_1 |f_1 - f_2|)] \leq n\mathbb{E} [w_{12} \sup a_1 | \sup f_1 - \inf f_2 |] + n\mathbb{E} [w_{12} \sup a_1 | \inf f_1 - \sup f_2 |]. \quad (48)$$

Both RHS terms in (48) can be dealt with using the same argument, so we concentrate on RHS1 which is by the triangle inequality bounded by

$$\begin{aligned} & \mathbb{E} [\sup a_1 | \sup f_1 - \inf f_1 |] + n\mathbb{E} [w_{12} \sup a_1 | \inf f_1 - \inf f_2 |] \stackrel{\text{B1}}{=} \mathbb{E} [\sup a_1 | \sup f_1 - \inf f_1 |] + o(1) \\ & \leq \stackrel{\text{Hölder}}{=} (\mathbb{E} [\sup a_1^{p_2/(p_2-1)}])^{1-1/p_2} (\mathbb{E} |\sup f_1 - \inf f_1|^{p_2})^{1/p_2} + o(1) \rightarrow 0 \text{ as } |\mathcal{N}| \rightarrow 0. \quad \square \end{aligned}$$

**Lemma E2.** Consider a process  $\mathcal{T}_n$  which for any polynomial  $\mathcal{P}_n$  and any  $1 = O(\epsilon_n)$ ,  $\sup_{\theta} \mathbb{P}[|\mathcal{T}_n(\theta)| > \epsilon_n] = o(\mathcal{P}_n^{-1})$  and for some polynomial  $\mathcal{P}_n^*$  satisfies  $\sup_{\theta} \|\mathcal{T}_{n\theta}(\theta)\| = o_p(\mathcal{P}_n^*)$  then  $\sup_{\theta} |\mathcal{T}_n(\theta)| = o_p(\epsilon_n)$ .

*Proof.* Divide  $\Theta$  into convex sets  $\Theta_1, \dots, \Theta_{\mathcal{P}_n}$  such that for all  $t = 1, \dots, \mathcal{P}_n$  and all  $\theta, \tilde{\theta} \in \Theta_t$ :  $\|\theta - \tilde{\theta}\| \leq \epsilon_n / \mathcal{P}_n^*$ . Then if  $\theta_t \in \Theta_t$ ,

$$\sup_{\theta} |\mathcal{T}_n(\theta)| \leq \max_{t=1, \dots, \mathcal{P}_n} \sup_{\theta \in \Theta_t} |\mathcal{T}_n(\theta) - \mathcal{T}_n(\theta_t)| + \max_{t=1, \dots, \mathcal{P}_n} |\mathcal{T}_n(\theta_t)|. \quad (49)$$

RHS1 in (49) is by the mean value theorem bounded by  $\sup_{\theta} \|\mathcal{T}_{n\theta}(\theta)\| \epsilon_n / \mathcal{P}_n^* = o_p(\epsilon_n)$ . For RHS2 note that

$$\mathbb{P} \left[ \max_{t=1, \dots, \mathcal{P}_n} |\mathcal{T}_n(\theta_t)| > \epsilon_n \right] \leq \mathcal{P}_n \sup_{\theta} \mathbb{P} [|\mathcal{T}_n(\theta)| > \epsilon_n] = o(1). \quad \square$$

**Lemma E3.** Let  $\{a_i, b_i\}$  be an i.i.d. sequence such that for some  $p_a, p_b > 2$  with  $1/p_a + 1/p_b < 1/2$ ,  $\mathbb{E} \sup_{\theta} |a_i(\theta)|^{p_a} < \infty$ ,  $\mathbb{E} \sup_{\theta} |b_i(\theta)|^{p_b} < \infty$ ,  $\mathbb{E} \sup_{\theta} \|a_{\theta i}(\theta)\| < \infty$ ,  $\mathbb{E} \sup_{\theta} \|b_{\theta i}(\theta)\| < \infty$ , and further that  $\mathbb{E}[a_i(\theta) | \mathcal{Z}^c] = 0$  a.s.,  $b_i(\theta) = b(z_i, \theta)$ . Then  $\sup_{\theta} |\sum_{ij} w_{ij} b_j(\theta) a_i(\theta)| = o_p(\sqrt{n} \log n)$ .

*Proof.* Let  $\tilde{a}_i = a_i I(|a_i| \leq Cn^{1/p_a}) - \mathbb{E}[a_i I(|a_i| \leq Cn^{1/p_a}) | z_i]$ ,  $\sigma_{ai}^2 = \text{Var}[\tilde{a}_i | z_i]$ , and  $\epsilon_n = \sqrt{n} \log n$ . We first show that  $\sup_{\theta} |\tilde{\mathcal{S}}_n - \mathcal{S}_n| = o_p(\epsilon_n)$  for  $\tilde{\mathcal{S}}_n = \sum_{ij} w_{ij} b_j \tilde{a}_i$  and  $\mathcal{S}_n = \sum_{ij} w_{ij} b_j a_i$ . Now, since

for any two random variables  $\mathbb{E}_1, \mathbb{E}_2$ ,  $\mathbb{P}[|\mathbb{E}_1| + |\mathbb{E}_2| > \epsilon_n] \leq \mathbb{P}[\mathbb{E}_1 \neq 0] + \mathbb{P}[|\mathbb{E}_2| > \epsilon_n]$  and noting that  $\mathbb{E}[a_i I(|a_i| > Cn^{1/p_a}) | z_i] = -\mathbb{E}[a_i I(|a_i| \leq Cn^{1/p_a}) | z_i]$  a.s.,

$$\begin{aligned} & \mathbb{P}\left[\sup_{\theta} |\mathcal{S}_n - \tilde{\mathcal{S}}_n| > \epsilon_n\right] \\ & \leq \mathbb{P}\left[\max_i \sup_{\theta} |a_i| > Cn^{1/p_a}\right] + \mathbb{P}\left[\sum_{ij} w_{ij} \sup_{\theta} |b_j| \sup_{\theta} \mathbb{E}[|a_i| I(|a_i| > Cn^{1/p_a}) | z_i] > \epsilon_n\right]. \end{aligned} \quad (50)$$

RHS1 in (50) is by the Markov inequality bounded by  $\mathbb{E} \sup |a_i|^{p_a} / C^{p_a} \rightarrow 0$  as  $C \rightarrow \infty$ . For RHS2 in (50), note first that  $\max_j \sup_{\theta} |b_j| = O_p(n^{1/p_b})$ , such that

$$\sum_{ij} w_{ij} \sup_{\theta} |b_j| \sup_{\theta} \mathbb{E}[|a_i| I(|a_i| > Cn^{1/p_a}) | z_i] = O_p(n^{1/p_b}) \sum_i \sup_{\theta} \mathbb{E}[|a_i| I(|a_i| > Cn^{1/p_a}) | z_i]. \quad (51)$$

Since for  $a_i^* = \sup_{\theta} |a_i|$  by the Hölder and Markov inequalities,

$$\mathbb{E}[a_i^* I(a_i^* > Cn^{1/p_a})] \leq (\mathbb{E}|a_i^*|^{p_a})^{1/p_a} (\mathbb{P}[a_i^* > Cn^{1/p_a}])^{1-1/p_a} \leq C^{1-p_a} n^{1/p_a-1} \mathbb{E}|a_i^*|^{p_a},$$

it follows that the RHS in (51) is  $O_p(n^{1/p_b+1/p_a}) = o_p(\sqrt{n})$  such that RHS2 in (50) is  $o(1)$ .

It remains to be shown that  $\sup_{\theta} |\tilde{\mathcal{S}}_n| = o_p(\epsilon_n)$ . Let  $\mathcal{Q}_n = \sum_{ij} w_{ij} b_j^2 \sigma_{a_i}^2$ . Then

$$\begin{aligned} \mathbb{P}\left[\sup_{\theta} |\tilde{\mathcal{S}}_n| > \epsilon_n\right] & \leq \mathbb{P}\left[\sup_{\theta} |\tilde{\mathcal{S}}_n| > \epsilon_n, \sup_{\theta} \mathcal{Q}_n \leq Cn, \max_i \sup_{\theta} |b_i| \leq Cn^{1/p_b}\right] \\ & \quad + \mathbb{P}\left[\sup_{\theta} \mathcal{Q}_n > Cn\right] + \mathbb{P}\left[\max_i \sup_{\theta} |b_i| > Cn^{1/p_b}\right]. \end{aligned} \quad (52)$$

Note that by taking  $C$  to  $\infty$  after taking the supremum over  $n$ , RHS3 in (52) tends to zero because  $\max_i \sup_{\theta} |b_i| = O_p(n^{1/p_b})$ . The same is true for RHS2 in (52) because  $\sup_{\theta} \mathcal{Q}_n \leq \sum_{ij} w_{ij} \sup_{\theta} b_j^2 \sup_{\theta} \sigma_{a_i}^2 = O_p(n)$  by lemma B1 and the moment conditions.

Finally RHS1 in (52). Note first that by theorem 1.1 of de la Peña (1999),

$$\begin{aligned} & \sup_{\theta} \mathbb{P}\left[|\tilde{\mathcal{S}}_n| > \epsilon_n \mid \mathcal{Z}, \sup_{\theta} \mathcal{Q}_n \leq Cn, \max_i \sup_{\theta} |b_i| \leq Cn^{1/p_b}\right] \\ & \leq 2 \exp\left(-\frac{\epsilon_n^2/2}{Cn + C^2 n^{1/p_a+1/p_b} \epsilon_n}\right) \leq 2 \exp\left(-\frac{(\log n)^2}{4C}\right), \end{aligned}$$

which converges faster than any polynomial of  $n$ . Finally, since by the conditions of this lemma,  $\sum_{ij} w_{ij} (\sup_{\theta} \|b_{\theta j}\| \sup_{\theta} |a_i| + \sup_{\theta} |b_j| \sup_{\theta} \|a_{\theta i}\|)$  increases at a polynomial rate, both conditions of lemma E2 are satisfied for  $\mathcal{T}_n(\theta) = \tilde{\mathcal{S}}_n(\theta) I(\max_i \sup_{\tilde{\theta}} |a_i(\tilde{\theta})| \leq Cn^{1/p_a})$ .<sup>27</sup> Since by the Markov inequality,

$$\mathbb{P}\left[\sup_{\theta} I(\mathcal{T}_n \neq \tilde{\mathcal{S}}_n) = 1\right] = \mathbb{P}\left[\max_i \sup_{\theta} |a_i| > Cn^{1/p_a}\right] \leq C^{-p_a} \mathbb{E} \sup_{\theta} |a_i|^{p_a} \rightarrow 0,$$

as  $C \rightarrow \infty$ , the lemma statement holds.  $\square$

<sup>27</sup>When  $\max_i \sup_{\theta} |a_i(\theta)| \leq Cn^{1/p_a}$ ,  $\tilde{\mathcal{S}}_n = \mathcal{S}_n$ , and hence  $\mathcal{T}_n$  is differentiable.



**Lemma E4.** Let  $\{a_i, b_i\}$  be an i.i.d. sequence such that for some  $p_a, p_b > 2$  with  $1/p_a + 1/p_b < 1/2$ ,  $n^{2/p_b}/k = o(1)$ ,  $\mathbb{E} \sup_{\theta} |a_i(\theta)|^{p_a} < \infty$ ,  $\mathbb{E} \sup_{\theta} |b_i(\theta)|^{p_b} < \infty$ ,  $\mathbb{E} \sup_{\theta} \|a_{\theta i}(\theta)\| < \infty$ ,  $\mathbb{E} \sup_{\theta} \|b_{\theta i}(\theta)\| < \infty$ , and further that  $\mathbb{E}[a_i(\theta)|\mathcal{Z}] = \mathbb{E}[b_i(\theta)|\mathcal{Z}] = 0$  a.s.. Then  $\sup_{\theta} |\sum_{ij} w_{ij} b_j(\theta) a_i(\theta)| = o_p(\sqrt{n/k}(\log n)^2)$ .

*Proof.* The arguments surrounding truncation of  $b_j, a_i$  are the same as those for  $a_i$  in lemma E3, so they are omitted here. We moreover take  $\sigma_{\tilde{a}_i}^2 = \text{Var}[\tilde{a}_i|z_i]$  and  $\sigma_{\tilde{b}_i}^2 = \text{Var}[\tilde{b}_i|z_i]$  to have an a.s. upper bound of one, uniformly in  $\theta, n$ , because the same truncation argument can be used otherwise. The discussion of the application of lemma E2 is similar to that in lemma E3 and is omitted here, also. We hence limit ourselves to showing that for any polynomial  $\mathcal{P}_n$ ,

$$\mathcal{P}_n \sup_{\theta} \mathbb{P} \left[ \left| \sum_{ij} w_{ij} \tilde{a}_i \tilde{b}_j \right| > 2\sqrt{n/k}(\log n)^2 \right] = o(1), \quad (53)$$

where  $\tilde{a}_i = \bar{a}_i - \mathbb{E}[\bar{a}_i|\mathcal{Z}]$  with  $\bar{a}_i = a_i I(|a_i| \leq Cn^{1/p_a})$  and  $\tilde{b}_i$  is similarly defined. Instead of (53) we in fact show that

$$\mathcal{P}_n \sup_{\theta} \mathbb{P} \left[ \left| \sum_i \sum_{j < i} w_{ij} \tilde{b}_j \tilde{a}_i \right| > \sqrt{n/k}(\log n)^2 \right] = o(1), \quad (54)$$

where the derivation for  $\sum_{j > i}$  obtains by reversing the sum. Let  $\mathcal{A}_i = \sum_{j < i} w_{ij} \tilde{b}_j \tilde{a}_i$ , such that  $\{\mathcal{A}_i\}$  is a martingale sequence. Then for  $\epsilon_{An} = n^{1/p_a+1/p_b}(\log n)/\sqrt{k}$  by theorem 1.1 of de la Peña (1999),

$$\begin{aligned} \sup_{\theta} \mathbb{P} \left[ \max_i |\mathcal{A}_i| > \epsilon_{An} \right] &\leq \sum_{i=1}^n \sup_{\theta} \mathbb{E} \left[ \mathbb{P} \left[ \left| \sum_j w_{ij} \tilde{b}_j \tilde{a}_i \right| > \epsilon_{An} \mid \mathcal{Z}, \tilde{a}_i \right] \right] \\ &\leq 2n \exp \left( - \frac{\epsilon_{An}^2/2}{C_w C^4 n^{2/p_a+2/p_b} + \epsilon_{An} C_w C^2 n^{1/p_a+1/p_b}} \right) = o(\mathcal{P}_n^{-1}). \end{aligned} \quad (55)$$

Further, for  $\mathcal{Q}_n = \sum_i \text{Var}[\mathcal{A}_i | \mathcal{Z}, \tilde{b}_1, \dots, \tilde{b}_{i-1}] = \sum_i (\sum_{j < i} w_{ij} \tilde{b}_j)^2 \sigma_{\tilde{a}_i}^2$  and  $\epsilon_{Qn} = n(\log n)^2/k$ , we have

$$\begin{aligned} \sup_{\theta} \mathbb{P} [\mathcal{Q}_n > \epsilon_{Qn}] &\leq \sum_{i=1}^n \sup_{\theta} \mathbb{P} \left[ \left| \sum_{j < i} w_{ij} \tilde{b}_j \right| > \sqrt{\frac{\epsilon_{Qn}}{n}} \right] \\ &\leq 2n \mathbb{E} \exp \left( - \frac{(\log n)^2/(2k)}{\sum_{j < i} w_{ij}^2 + C_w C n^{1/p_b}(\log n)/k^{3/2}} \right) = o(\mathcal{P}_n^{-1}). \end{aligned} \quad (56)$$

Finally, letting  $\epsilon_n = \sqrt{n/k}(\log n)^2$ ,

$$\begin{aligned} \sup_{\theta} \mathbb{P} \left[ \left| \sum_i \mathcal{A}_i \right| > \epsilon_n \right] &\leq \\ \sup_{\theta} \mathbb{P} \left[ \left| \sum_i \mathcal{A}_i \right| > \epsilon_n, \mathcal{Q}_n \leq \epsilon_{Qn}, \max_i |\mathcal{A}_i| \leq \epsilon_{An} \right] &+ \sup_{\theta} \mathbb{P} [\mathcal{Q}_n > \epsilon_{Qn}] + \sup_{\theta} \mathbb{P} \left[ \max_i |\mathcal{A}_i| > \epsilon_{An} \right] \\ &\leq 2 \exp \left( - \frac{\epsilon_n^2/2}{\epsilon_{Qn} + \epsilon_{An} \epsilon_n} \right) + o(\mathcal{P}_n^{-1}) + o(\mathcal{P}_n^{-1}) = o(\mathcal{P}_n^{-1}), \end{aligned} \quad (57)$$

where the last inequality in (57) follows from theorem 1.2A of de la Peña (1999) plus (55), (56), and the equality by the definition of  $\epsilon_n$ . So (54) holds.  $\square$

**E.2. Specifics.** Let below  $\mathcal{N} = \{\theta : \|\theta - \theta_0\| \leq \|\theta_n - \theta_0\|\}$ .

**Lemma E5.**  $\sum_{ij} w_{ij} g_j(\theta_n) m_i(\theta_n) = \sum_{ij} w_{ij} g_j m_i + \sqrt{n} \lambda \mathbb{E} \tilde{g}_1 \tilde{g}_1^\top \Delta + o_p(\rho)$ .

*Proof.* Using the mean value theorem,

$$\begin{aligned} & \left\| \sum_{ij} w_{ij} (g_j(\theta_n) m_i(\theta_n) - g_j m_i) - \sqrt{n} \lambda \mathbb{E} \tilde{g}_1 \tilde{g}_1^\top \Delta \right\| \\ & \leq \frac{\|\Delta\|}{\sqrt{n}} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} \tilde{M}_{\theta\theta j}(\theta) u_i(\theta) \right\| + \frac{\|\Delta\|}{\sqrt{n}} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} \tilde{g}_j(\theta) v_i^\top(\theta) \right\| \\ & \quad + \frac{\|\Delta\|}{\sqrt{n}} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} \tilde{g}_j(\theta) M_i(\theta) \right\| + \frac{\|\Delta\| \lambda}{\sqrt{n}} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} \tilde{g}_j(\theta) \tilde{g}_i^\top(\theta) - n \mathbb{E} \tilde{g}_1 \tilde{g}_1^\top \right\|. \end{aligned} \quad (58)$$

RHS1 and RHS2 in (58) are  $o_p(\log n) = o_p(\rho)$  by lemma E3. By the mean value theorem RHS3 is bounded by  $\|\Delta\|^2 n^{-1} \sup_{\theta, \theta^* \in \mathcal{N}} \left\| \sum_{ij} w_{ij} \tilde{g}_j(\theta) \tilde{g}_i^\top(\theta^*) \right\| = O_p(1) = o_p(\rho)$ . Finally RHS4 is by lemma E1 equal to  $o_p(\sqrt{n} \lambda) = o_p(\rho)$ .  $\square$

**Lemma E6.**  $\sum_{ij} w_{ij} v_j(\theta_n) m_i(\theta_n) - \sum_{ij} w_{ij} v_j m_i = o_p(\rho)$ .

*Proof.* The norm of the LHS is by the mean value theorem bounded by

$$\begin{aligned} & \frac{\|\Delta\|}{\sqrt{n} \lambda} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} (m_{\theta\theta j}(\theta) - M_{\theta\theta j}(\theta)) u_i(\theta) \right\| + \frac{\|\Delta\|}{\sqrt{n} \lambda} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} v_j(\theta) v_i^\top(\theta) \right\| \\ & \quad + \frac{\|\Delta\|}{\sqrt{n}} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} (m_{\theta\theta j}(\theta) - M_{\theta\theta j}(\theta)) \tilde{M}_i(\theta) \right\| + \frac{\|\Delta\|}{\sqrt{n}} \sup_{\mathcal{N}} \left\| \sum_{ij} w_{ij} v_j(\theta) \tilde{g}_i^\top(\theta) \right\|. \end{aligned} \quad (59)$$

The first two terms in (59) are  $o_p((\log n)^2 / (\sqrt{k} \lambda)) = o_p((\log n)^2) = o_p(\rho)$  by lemma E4. The last two terms in (59) are  $o_p(\log n) = o_p(\rho)$  by lemma E3, making minor adjustments for the fact that  $a, b$  are interchanged.  $\square$

**Lemma E7.**  $\sum_i \hat{g}_i(\theta_n) \hat{g}_i^\top(\theta_n) m_i^2(\theta_n) = \sum_i \hat{g}_i \hat{g}_i^\top m_i^2 + o_p(\rho)$ .

*Proof.* By the mean value theorem,

$$\begin{aligned} & \left\| \sum_i \hat{g}_i(\theta_n) \hat{g}_i^\top(\theta_n) m_i^2(\theta_n) - \sum_i \hat{g}_i \hat{g}_i^\top m_i^2 \right\| \leq \frac{\|\Delta\|}{\sqrt{n} \lambda} \sup_{\mathcal{N}} \sum_i \|\hat{g}_{\theta i}(\theta)\| \|\hat{g}_i(\theta)\| m_i^2(\theta) \\ & \quad + \frac{\|\Delta\|}{\sqrt{n} \lambda} \sup_{\mathcal{N}} \sum_i \|\hat{g}_i(\theta)\| \|\hat{g}_{\theta i}(\theta)\| m_i^2(\theta) + \frac{2\|\Delta\|}{\sqrt{n} \lambda} \sup_{\mathcal{N}} \left\| \sum_i \hat{g}_i(\theta) \hat{g}_i(\theta) m_{\theta i}(\theta) m_i(\theta) \right\|. \end{aligned} \quad (60)$$

RHS1, RHS2, and RHS3 in (60) are similar, and we only obtain a result for RHS3, which is bounded by

$$\frac{2\|\Delta\| \lambda}{\sqrt{n}} \sum_i \left\| \sum_j w_{ij} \tilde{g}_j(\theta) \right\|^2 \|m_{\theta i}(\theta) m_i(\theta)\| + \frac{2\|\Delta\|}{\sqrt{n} \lambda} \sum_i \left\| \sum_j w_{ij} v_j(\theta) \right\|^2 \|m_{\theta i}(\theta) m_i(\theta)\|. \quad (61)$$

The first term in (61) is  $O_p(\sqrt{n} \lambda) = o_p(\rho^2)$ . For the second term, note that by arguments similar to (but simpler than) those made in lemma E3,  $\sup_{\mathcal{N}} \max_i \left\| \sum_j w_{ij} v_j \right\| = o_p((\log n) / \sqrt{k})$ . Hence the second term in (61) is  $o_p((\log n)^2 \sqrt{n} / (k \lambda)) = o_p(\rho^2)$ .  $\square$

**Lemma E8.**  $\sum_{ij} w_{ij} w_{ji} m_j(\theta_n) m_{\theta_j}^\top(\theta_n) m_i(\theta_n) m_{\theta_i}(\theta_n) - \sum_{ij} w_{ij} w_{ji} m_j m_{\theta_j}^\top m_i m_{\theta_i} = o_p(\rho^2)$ .

*Proof.* By the mean value theorem,

$$\begin{aligned} & \left\| \sum_{ij} w_{ij} w_{ji} m_j(\theta_n) m_{\theta_j}^\top(\theta_n) m_i(\theta_n) m_{\theta_i}(\theta_n) - \sum_{ij} w_{ij} w_{ji} m_j m_{\theta_j}^\top m_i m_{\theta_i} \right\| \\ & \leq \frac{\|\Delta\|}{\sqrt{n\lambda}} \sup_{\mathcal{N}} \left\{ \sum_{ij} w_{ij} w_{ji} \|m_{\theta_j}(\theta)\|^2 \|m_i(\theta) m_{\theta_i}(\theta)\| + \sum_{ij} w_{ij} w_{ji} \|m_j(\theta) m_{\theta_j}(\theta)\| \|m_i(\theta) m_{\theta_i}(\theta)\| \right. \\ & \quad \left. + \sum_{ij} w_{ij} w_{ji} \|m_j(\theta) m_{\theta_j}(\theta)\| \|m_{\theta_i}(\theta)\|^2 + \sum_{ij} w_{ij} w_{ji} \|m_j(\theta) m_{\theta_j}(\theta)\| \|m_i(\theta) m_{\theta_i}(\theta)\| \right\} \\ & = O_p\left(\frac{\sqrt{n}}{k\lambda}\right) = o_p(\rho) = o_p(\rho^2). \quad \square \end{aligned}$$

#### APPENDIX F. NUISANCE PARAMETERS

**Lemma F1.** If  $f, f^*$  are such that for a neighborhood  $\mathfrak{N}$  of  $\tilde{\beta}$ ,  $\mathbb{E}(\mathbb{E}[\sup_{\beta \in \mathfrak{N}} \|f_1(\theta, \beta)\| | z_1])^2 < \infty$  and  $\mathbb{E}(\mathbb{E}[\sup_{\beta \in \mathfrak{N}} \|f_1^*(\theta, \beta)\| | z_1])^2 < \infty$ , then for any  $\beta^*$  converging to  $\tilde{\beta}$ ,  $\sum_{ij} w_{ij} f_i^*(\theta, \beta^*) f_j(\theta, \beta^*) = O_p(n)$ .

*Proof.* Let  $n$  be big enough to ensure that  $\beta^* \in \mathfrak{N}$ . Then

$$\begin{aligned} \left\| \sum_{ij} w_{ij} f_j(\theta, \beta^*) f_i^*(\theta, \beta^*) \right\| & \leq \sup_{\beta \in \mathfrak{N}} \left\| \sum_{ij} w_{ij} f_j(\theta, \beta) f_i^*(\theta, \beta) \right\| \\ & \leq \sum_{ij} w_{ij} \sup_{\beta \in \mathfrak{N}} \|f_j(\theta, \beta)\| \sup_{\beta \in \mathfrak{N}} \|f_i^*(\theta, \beta)\| = O_p(n), \end{aligned}$$

by lemma B5. □

**Lemma F2.**  $\hat{\beta}(\theta) \xrightarrow{p} \beta(\theta)$ .

*Proof.* Since  $\mathbb{E}[h_1 m_1]$  is continuous and yields a unique solution  $\beta(\theta)$  by assumption D, we are left to show uniform convergence (in  $\beta$  at  $\theta$ ) of  $n^{-1} \sum_i \hat{h}_i m_i$  to  $\mathbb{E}[h_1 m_1]$ . We first show pointwise convergence and then stochastic equicontinuity. By lemma B6(vi) at fixed  $\theta, \beta$ ,

$$n^{-1} \sum_i \hat{h}_i m_i = n^{-1} \sum_{ij} w_{ij} m_{\beta j} m_i = n^{-1} \sum_i h_i \mathbb{E}[m_i | z_i] + o_p(1) = \mathbb{E}[h_1 m_1] + o_p(1).$$

Stochastic equicontinuity then follows from the fact that

$$\begin{aligned} & \sup_{\tilde{\beta}} \sup_{\beta: \|\beta - \tilde{\beta}\| \leq \delta} \left\| n^{-1} \sum_i [\hat{h}_i(\theta, \tilde{\beta}) m_i(\theta, \tilde{\beta}) - \hat{h}_i(\theta, \beta) m_i(\theta, \beta)] \right\| \\ & \leq \left( n^{-1} \sum_{ij} w_{ij} \sup_{\beta} \|m_{\beta j}(\theta, \beta)\| \sup_{\beta} \|m_i(\theta, \beta)\| \right. \\ & \quad \left. + n^{-1} \sum_{ij} w_{ij} \sup_{\beta} \|m_{\beta j}(\theta, \beta)\| \sup_{\beta} \|m_{\beta i}(\theta, \beta)\| \right) \times \sup_{\tilde{\beta}} \sup_{\beta: \|\beta - \tilde{\beta}\| \leq \delta} \|\tilde{\beta} - \beta\| \stackrel{\text{B5.E}}{=} O_p(\delta). \end{aligned}$$

Choose  $\delta$  proportional to  $\varepsilon$  in (21.42) of Davidson (1994). □

**Lemma F3.**  $\hat{\beta}(\theta) - \beta(\theta) = \begin{cases} O_p(n^{-1/2}), & \theta = \theta_0, \\ O_p(n^{-1/2}) + o_p(\lambda), & \theta \neq \theta_0. \end{cases}$

*Proof.* Let  $S_n(\theta, \beta^*)$  be the vector with  $t$ -th element

$$\begin{aligned} & (\hat{\beta}(\theta) - \beta(\theta))^\top \sum_i \left( \hat{h}_{\beta\beta it}(\theta, \beta^*) m_i(\theta, \beta^*) + 2\hat{h}_{\beta it}(\theta, \beta^*) m_{\beta i}^\top(\theta, \beta^*) \right. \\ & \qquad \qquad \qquad \left. + \hat{h}_{it}(\theta, \beta^*) m_{\beta\beta i}(\theta, \beta^*) \right) (\hat{\beta}(\theta) - \beta(\theta)) / 2, \end{aligned}$$

where  $\hat{h}_{it}$  is the  $t$ -th element of  $\hat{h}_i$  and  $\hat{h}_{\beta it}, \hat{h}_{\beta\beta it}$  are its first and second partial derivatives with respect to  $\beta$ . Then by lemmas F1 and F2,  $S_n(\theta, \beta^*) = O_p(n \|\hat{\beta}(\theta) - \beta(\theta)\|^2)$  for any  $\beta^*$  between  $\hat{\beta}(\theta)$  and  $\beta(\theta)$ . Hence by the mean value theorem (omitting  $\theta$  from hereon),

$$0 = \sum_i \hat{h}_i \hat{m}_i = \sum_i \hat{h}_i m_i + \sum_i m_i \hat{h}_{\beta i}^\top (\hat{\beta} - \beta) + \sum_i \hat{h}_i m_{\beta i}^\top (\hat{\beta} - \beta) + O_p(n \|\hat{\beta} - \beta\|^2). \quad (62)$$

Now, noting that  $M_i(\theta_0) = 0$  a.s. and that  $\mathbb{E}[h_1 M_1] = \mathbb{E}[h_1 m_1] = 0$  by definition,

$$\sum_i \hat{h}_i m_i = \underbrace{\sum_i \hat{h}_i (m_i - M_i)}_{\stackrel{\text{B6(ii)}}{=} O_p(\sqrt{n})} + \underbrace{\lambda \sum_i (\hat{h}_i - h_i) \tilde{M}_i}_{\stackrel{\text{B6(vi)}}{=} o_p(n\lambda)} + \underbrace{\lambda \sum_i h_i \tilde{M}_i}_{=O_p(\sqrt{n}\lambda)} = \begin{cases} O_p(\sqrt{n}), & \theta = \theta_0, \\ O_p(\sqrt{n}) + o_p(n\lambda), & \theta \neq \theta_0. \end{cases} \quad (63)$$

Further,

$$\begin{aligned} \sum_i \hat{h}_{\beta i} m_i &= \underbrace{\sum_i \hat{h}_{\beta i} (m_i - M_i)}_{\stackrel{\text{B6(ii)}}{=} O_p(\sqrt{n})} + \underbrace{\lambda \sum_i (\hat{h}_{\beta i} - h_{\beta i}) \tilde{M}_i}_{\stackrel{\text{B6(vi)}}{=} o_p(n\lambda)} + \underbrace{\lambda \sum_i (h_{\beta i} \tilde{M}_i - \mathbb{E}[h_{\beta 1} \tilde{M}_1])}_{=O_p(\sqrt{n}\lambda)} + n\mathbb{E}[h_{\beta 1} M_1] \\ &= n\mathbb{E}[m_{\beta\beta 1} M_1] + \begin{cases} O_p(\sqrt{n}), & \theta = \theta_0, \\ O_p(\sqrt{n}) + o_p(n\lambda), & \theta \neq \theta_0. \end{cases} \end{aligned} \quad (64)$$

And

$$\sum_i \hat{h}_i m_{\beta i}^\top \stackrel{\text{B6(vi)}}{=} \sum_i h_i h_i^\top + o_p(n) = n\mathbb{E}[h_1 h_1^\top] + o_p(n). \quad (65)$$

Plugging (63)–(65) into (62) we get

$$n(Q + o_p(1))(\hat{\beta} - \beta) = \begin{cases} O_p(\sqrt{n}), & \theta = \theta_0, \\ O_p(\sqrt{n}) + o_p(n\lambda), & \theta \neq \theta_0. \end{cases}$$

The stated result then follows from the fact that  $Q$  is invertible by assumption D.  $\square$

**Lemma F4.**  $\sum_i \hat{q}_i(\theta) \hat{m}_i(\theta) = \sum_i \hat{q}_i(\theta) m_i(\theta) + \begin{cases} o_p(\rho), & \theta = \theta_0, \\ O_p(\rho) + o_p(n\lambda^2), & \theta \neq \theta_0. \end{cases}$

*Proof.* Since we can consider each element of  $\hat{q}_i$  separately, we may assume that  $\hat{q}_i$  is a scalar without loss of generality. Note then that for any  $\beta^*$  between  $\hat{\beta}(\theta)$  and  $\beta(\theta)$ ,

$$\left\| \left( \hat{\beta}(\theta) - \beta(\theta) \right)^\top \sum_i \left[ \hat{g}_{\beta\beta i}(\theta, \beta^*) m_i(\theta, \beta^*) + 2\hat{g}_{\beta i}(\theta, \beta^*) m_{\beta i}^\top(\theta, \beta^*) + \hat{g}_i(\theta, \beta^*) m_{\beta\beta i}(\theta, \beta^*) \right] \left( \hat{\beta}(\theta) - \beta(\theta) \right) \right\| \stackrel{\text{F1}}{=} O_p(n \|\hat{\beta}(\theta) - \beta(\theta)\|^2).$$

Hence by the mean value theorem (omitting arguments from hereon)

$$\begin{aligned} \sum_i \hat{q}_i \hat{m}_i &= \sum_i \hat{g}_i \hat{m}_i = \sum_i \hat{g}_i [m_i + m_{\beta i}^\top(\hat{\beta} - \beta)] + \sum_i m_i \hat{g}_{\beta i}^\top(\hat{\beta} - \beta) + O_p(n \|\hat{\beta} - \beta\|^2) \\ &= \sum_i \hat{q}_i m_i + \sum_i \hat{q}_i m_{\beta i}^\top(\hat{\beta} - \beta) + \kappa^\top \sum_i \hat{h}_i [m_i + m_{\beta i}^\top(\hat{\beta} - \beta)] \\ &\quad + \sum_i m_i \hat{g}_{\beta i}^\top(\hat{\beta} - \beta) + \begin{cases} O_p(1), & \theta = \theta_0, \\ O_p(1) + o_p(n\lambda^2), & \theta \neq \theta_0. \end{cases} \end{aligned} \quad (66)$$

We deal with each of the RHS2–RHS4 terms separately. For RHS3 in (66) note that by (62),

$$\begin{aligned} \left\| \sum_i \hat{h}_i (m_i + m_{\beta i}^\top(\hat{\beta} - \beta)) \right\| &\leq \left\| \sum_i m_i \hat{h}_{\beta i} \right\| \|\hat{\beta} - \beta\| + O_p(n \|\hat{\beta} - \beta\|^2) \\ &\stackrel{\text{F3}, (64)}{=} \begin{cases} O_p(n^{1/2}) O_p(n^{-1/2}) + O_p(1), & \theta = \theta_0, \\ \{O_p(n^{1/2}) + o_p(n\lambda)\} \{O_p(n^{-1/2}) + o_p(\lambda)\} + O_p(1) + o_p(n\lambda^2), & \theta \neq \theta_0. \end{cases} \end{aligned}$$

which is  $\begin{cases} O_p(1), & \theta = \theta_0, \\ o_p(\rho + n\lambda^2), & \theta \neq \theta_0. \end{cases}$ . For RHS4 in (66) we have

$$\sum_i \hat{g}_{\beta i} m_i = \underbrace{\sum_i \hat{g}_{\beta i} (m_i - M_i)}_{\stackrel{\text{B6(ii)}}{=} O_p(\sqrt{n})} + \underbrace{\lambda \sum_i (\hat{g}_{\beta i} - g_{\beta i}) \tilde{M}_i}_{\stackrel{\text{B6(vi)}}{=} o_p(n\lambda)} + \underbrace{\lambda \sum_i g_{\beta i} \tilde{M}_i}_{=O_p(n\lambda)} = \begin{cases} O_p(\sqrt{n}), & \theta = \theta_0, \\ O_p(\sqrt{n} + n\lambda), & \theta \neq \theta_0. \end{cases} \quad (67)$$

Now combine (67) with lemma F3 to obtain rates for RHS4 in (66) of  $O_p(1)$  if  $\theta = \theta_0$  and  $O_p(\rho) + o_p(n\lambda^2)$  if  $\theta \neq \theta_0$ . Finally RHS2 in (66). We have

$$\begin{aligned} \sum_i \hat{q}_i m_{\beta i} &\stackrel{\text{B6(ii)}}{=} \sum_i \hat{q}_i h_i + O_p(\sqrt{n}) \stackrel{\mathbb{E}[q_1 h_1] = 0}{=} \sum_i (\hat{q}_i - q_i) h_i + O_p(\sqrt{n}) \\ &= \underbrace{\lambda \sum_{ij} w_{ij} (\tilde{q}_j - \tilde{q}_i) h_i}_{\stackrel{\text{B6(vi)}}{=} o_p(n\lambda)} + \underbrace{\sum_{ij} w_{ij} (m_{\theta j}^\top - \kappa^\top m_{\beta j} - q_j) h_i}_{\stackrel{\text{B6(iv)}}{=} o_p(n^{3/4} k^{-1/4})} + O_p(\sqrt{n}) \\ &= O_p(\sqrt{n}) + o_p(\sqrt{n\rho}) + o_p(n\lambda), \end{aligned}$$

such that RHS2 in (66) is  $\begin{cases} o_p(\rho), & \theta = \theta_0, \\ o_p(\rho + n\lambda^2), & \theta \neq \theta_0. \end{cases}$   $\square$

**Lemma F5.**  $\hat{\kappa}(\theta) - \kappa(\theta) = o_p(\rho/\sqrt{n})$ .

*Proof.* Note that (omitting  $\theta$ )

$$\hat{\kappa} - \kappa = \left( n^{-1} \sum_i \hat{h}_i \hat{h}_i^\top \right)^{-1} n^{-1} \sum_i \hat{h}_i (\hat{g}_i^\top - \hat{h}_i^\top \kappa).$$

Since we can consider each column of  $\hat{\kappa} - \kappa$  separately, we may assume that  $\hat{g}_i$  (and hence  $q_i$ ) is a scalar without loss of generality. By Slutsky it then suffices to show that (i)  $n^{-1} \sum_i \hat{h}_i \hat{h}_i^\top - \mathbb{E}[h_1 h_1^\top] = o_p(1)$  and (ii)  $n^{-1} \sum_i \hat{h}_i (\hat{g}_i - \hat{h}_i^\top \kappa) = o_p(\rho/\sqrt{n})$ . Establishing (i) is similar to but simpler than showing (ii) so we only show (ii) here. We show that

$$\sum_i [\hat{h}_i (\hat{g}_i - \hat{h}_i^\top \kappa) - \hat{h}_i \hat{q}_i] = o_p(\sqrt{n}\rho), \quad (68)$$

$$\sum_i (\hat{h}_i \hat{q}_i - h_i q_i) = o_p(\sqrt{n}\rho), \quad (69)$$

$$\sum_i h_i q_i = o_p(\sqrt{n}\rho). \quad (70)$$

Since  $\mathbb{E}[h_1 q_1] = 0$  by construction, the LHS in (70) is  $O_p(\sqrt{n}\lambda)$ . For (69) note that

$$\sum_i (\hat{h}_i \hat{q}_i - h_i q_i) = \sum_i (\hat{h}_i - h_i) (\hat{q}_i - q_i) + \sum_i (\hat{h}_i - h_i) q_i + \sum_i h_i (\hat{q}_i - q_i). \quad (71)$$

First RHS<sub>3</sub> in (71). Let  $\epsilon_i = m_{\theta i}^\top - \kappa^\top m_{\beta i} - q_i$ . Then

$$\sum_i h_i (\hat{q}_i - q_i) = \lambda \sum_{ij} w_{ij} h_i (\hat{q}_j - \tilde{q}_j) + \sum_{ij} w_{ij} h_i \epsilon_j \stackrel{\text{B1}}{=} o_p(n\lambda) + \sum_{ij} w_{ij} h_i \epsilon_j.$$

But

$$\mathbb{E} \left\| \sum_{ij} w_{ij} h_i \epsilon_j \right\|^2 \leq \sum_{ij} \mathbb{E}[w_{ij}^2 |h_i|^2 |\epsilon_j|^2] + \sum_{ij} \sum_{t \neq i} \mathbb{E}[w_{ij} w_{it} |h_i| |h_t| |\epsilon_j|^2] \stackrel{\text{B5}}{=} o(n^{3/2} k^{-1/2}) = o(n\rho).$$

The derivation for RHS<sub>2</sub> in (71) is similar to that for RHS<sub>3</sub> and is omitted. By the Schwarz inequality, sufficient conditions for RHS<sub>1</sub> in (71) to be  $o_p(\sqrt{n}\rho)$  are that (a)  $\sum_i \|\hat{h}_i - h_i\|^2 = o_p(n)$  and (b)  $\sum_i \|\hat{q}_i - q_i\|^2 = O_p(\rho^2)$ . First (b). Since  $\hat{q}_i - q_i = \sum_j w_{ij} (q_j - q_i + \epsilon_j)$ , we have

$$\begin{aligned} 2^{-1} \sum_i \mathbb{E} \|\hat{q}_i - q_i\|^2 &\leq \sum_i \mathbb{E} \left\| \sum_j w_{ij} (q_j - q_i) \right\|^2 + \sum_i \mathbb{E} \left\| \sum_j w_{ij} \epsilon_j \right\|^2 \\ &\stackrel{\text{Schwarz}}{\leq} n^2 \lambda^2 \mathbb{E}[w_{12} |\tilde{q}_2 - \tilde{q}_1|^2] + n^2 \mathbb{E}[w_{12}^2 |\epsilon_2|^2] \stackrel{\text{B1, B5}}{=} o(n\lambda^2) + O(n/k) = O(\rho^2). \end{aligned}$$

Requirement (a) follows with the same derivation resulting in a rate of  $o(n) + O(n/k) = o(n)$ . Finally apply the mean value theorem to the LHS in (68) and obtain for some  $\beta^*$  between  $\hat{\beta}$  and  $\beta$

$$\sum_i \left[ \hat{h}_{\beta i} (\hat{g}_i - \hat{h}_i^\top \kappa) + \hat{h}_i (\hat{g}_{\beta i} - \hat{h}_{\beta i}^\top \kappa)^\top \right] (\hat{\beta} - \beta) \quad \text{all at } (\theta, \beta^*). \quad (72)$$

By lemma F3,  $\|\hat{\beta} - \beta\| = o_p(\rho/\sqrt{n})$ . It is hence sufficient to show that the sum in (72) is  $O_p(n)$ , which follows from repeated application of lemma F1.  $\square$

**Lemma F6.**  $\sum_i \hat{q}_i(\theta) \hat{q}_i^\top(\theta) \hat{m}_i^2(\theta) = \sum_i \hat{q}_i(\theta) \hat{q}_i^\top(\theta) m_i^2(\theta) + o_p(\rho^2)$ .

*Proof.* Since for generic  $a, b$ ,  $\|aa^\top - bb^\top\| \leq \|a - b\|^2 + 2\|b\|\|a - b\|$ , and  $\sum_i \|\hat{q}_i m_i\|^2 = O_p(\rho^2)$  by lemmas C22 and C24 (substituting  $\hat{q}$  for  $\hat{g}$ ), we need to show that  $\sum_i \|\hat{q}_i \hat{m}_i - \hat{q}_i m_i\|^2 = o_p(\rho^2)$ . Since we can consider each element of  $\hat{q}_i \hat{m}_i - \hat{q}_i m_i$  separately, we may assume that  $\hat{q}_i$  and  $\hat{q}_i$  are scalars without loss of generality, and we only need to show  $\sum_i (\hat{q}_i \hat{m}_i - \hat{q}_i m_i)^2 = o_p(\rho)$ . Now,

$$\hat{q}_i \hat{m}_i - \hat{q}_i m_i = (\hat{g}_i \hat{m}_i - \hat{g}_i m_i) - (\hat{\kappa} - \kappa)^\top (\hat{h}_i \hat{m}_i - \hat{h}_i m_i) - \kappa^\top (\hat{h}_i \hat{m}_i - \hat{h}_i m_i) - (\hat{\kappa} - \kappa)^\top \hat{h}_i m_i.$$

By lemma F5 it hence suffices to show that

$$\sum_i (\hat{g}_i \hat{m}_i - \hat{g}_i m_i)^2 = o_p(\rho^2), \quad \sum_i \|\hat{h}_i \hat{m}_i - \hat{h}_i m_i\|^2 = o_p(\rho^2), \quad \sum_i \|\hat{h}_i m_i\|^2 = O_p(n). \quad (73)$$

The third condition in (73) follows from lemma B6. Now the first condition. We have for some  $\beta^*$  between  $\hat{\beta}$  and  $\beta$ ,

$$\begin{aligned} \sum_i (\hat{g}_i \hat{m}_i - \hat{g}_i m_i)^2 &= \sum_i \left( \sum_j w_{ij} (\hat{m}_{\theta_j} \hat{m}_i - m_{\theta_j} m_i) \right)^2 \\ &\leq \|\hat{\beta} - \beta\|^2 \sum_{ijt} w_{ij} w_{it} \left| m_t(\theta, \beta^*) m_{\theta\beta_j}(\theta, \beta^*) + m_{\theta_j}(\theta, \beta^*) m_{\beta t}(\theta, \beta^*) \right|. \end{aligned} \quad (74)$$

Now since for generic  $a_i, b_i$

$$\sum_{ijt} w_{ij} w_{it} |a_j b_t| \leq \sum_{ijt} w_{ij} w_{it} (a_j^2 + b_t^2) = \sum_{ij} w_{ij} (a_j^2 + b_j^2),$$

it follows from lemmas F1 and F3 that the RHS in (74) is  $o_p(\rho^2/n) O_p(n) = o_p(\rho^2)$ . So we have established the first condition in (73), where the derivation for the second condition follows similarly.  $\square$

**Lemma F7.**  $\sum_{ij} w_{ij} w_{ji} \hat{m}_j(\theta) \hat{m}_{\theta_j}^\top(\theta) \hat{m}_i(\theta) \hat{m}_{\theta_i}(\theta) = \sum_{ij} w_{ij} w_{ji} m_j(\theta) m_{\theta_j}^\top(\theta) m_i(\theta) m_{\theta_i}(\theta) + o_p(\rho^2)$ .

*Proof.* By the mean value theorem, we have

$$\begin{aligned} &\left| \sum_{ij} w_{ij} w_{ji} (\hat{m}_j \hat{m}_{\theta_j}^\top \hat{m}_i \hat{m}_{\theta_i}^\top - m_j m_{\theta_j}^\top m_i m_{\theta_i}^\top) \right| \\ &\leq (C_w/k) \|\hat{\beta} - \beta\| \left| \sum_{ij} w_{ij} \left| m_{\beta_j}(\theta, \beta^*) m_{\theta_j}^\top(\theta, \beta^*) m_i(\theta, \beta^*) m_{\theta_i}(\theta, \beta^*) + \dots + \right. \right. \\ &\quad \left. \left. m_j(\theta, \beta^*) m_{\theta_j}^\top(\theta, \beta^*) m_i(\theta, \beta^*) m_{\theta_{\beta i}}(\theta, \beta^*) \right| \right|. \end{aligned}$$

Apply the triangle inequality and lemmas F1 and F3 to obtain a rate of  $o_p(\rho/\sqrt{n}) O_p(n/k) = o_p(\rho^2)$ .  $\square$

**Lemma F8.**  $\hat{\mathcal{D}}_N^2(\theta) = \hat{\mathcal{D}}_N^2(\theta) + o_p(\rho^2)$ .

*Proof.* By lemmas F4, F6 and F7 (omitting  $\theta$ ),

$$\hat{\mathcal{D}}_N^2 - \hat{\mathcal{D}}_N^2 = \underbrace{\left[ \sum_i \hat{q}_i \hat{q}_i^\top \hat{m}_i^2 - \sum_i \hat{q}_i \hat{q}_i^\top m_i^2 \right]}_{\stackrel{\text{F6}}{=} o_p(\rho^2)} - n^{-1} \left[ \left( \sum_i \hat{q}_i \hat{m}_i \right) \left( \sum_i \hat{q}_i^\top \hat{m}_i \right) - \left( \sum_i \hat{q}_i m_i \right) \left( \sum_i \hat{q}_i^\top m_i \right) \right]$$

$$+ \underbrace{\left[ \sum_{ij} w_{ij} w_{ji} \hat{m}_j \hat{m}_j^\top \hat{m}_i \hat{m}_i^\top - \sum_{ij} w_{ij} w_{ji} m_j m_j^\top m_i m_i^\top \right]}_{\stackrel{\text{Fz}}{=} o_p(\rho^2)} \quad (75)$$

Finally, the norm of RHS2 in (75) is bounded by

$$\underbrace{n^{-1} \left\| \sum_i (\hat{q}_i \hat{m}_i - \hat{q}_i m_i) \right\|^2}_{\stackrel{\text{F4}}{=} O_p(\rho^2/n) + o_p(n\lambda^4)} + 2n^{-1} \underbrace{\left\| \sum_i \hat{q}_i m_i \right\|}_{=O_p(\rho)} \underbrace{\left\| \sum_i (\hat{q}_i \hat{m}_i - \hat{q}_i m_i) \right\|}_{\stackrel{\text{F4}}{=} O_p(\rho) + o_p(n\lambda^2)} = o_p(\rho^2). \quad \square$$

## APPENDIX G. JUSTIFICATION OF EXAMPLE I

We first justify that

$$\begin{bmatrix} n^{-1/2} \sum_i \tilde{g}_i u_i \\ n^{-1/2} \sum_i \tilde{g}_i v_i \\ (n/k)^{-1/2} \sum_{ij} w_{ij} v_j u_i \\ (n/k)^{-1/2} \sum_{ij} w_{ij} v_j v_i \end{bmatrix} \xrightarrow{d} \begin{bmatrix} \Psi_{\tilde{g}u} \\ \Psi_{\tilde{g}v} \\ \Psi_{vu} \\ \Psi_{vv} \end{bmatrix} \sim N(0, V_\Psi), \quad (76)$$

for some  $0 < V_\Psi < \infty$ . By the Cramér–Wold device, it suffices to show the normality of an arbitrary linear combination, which can be expressed as

$$S_n = \sum_i \left( \frac{\Xi_i}{\sqrt{n}} + \frac{\sqrt{k}}{\sqrt{n}} \sum_{j<i} Y_{ij} \right),$$

where  $\Xi_i = a_1 \tilde{g}_i u_i + a_2 \tilde{g}_i v_i$ , and  $Y_{ij} = w_{ij}(a_3 v_j u_i + a_4 v_j v_i) + w_{ji}(a_3 v_i u_j + a_4 v_i v_j)$  with  $a_1, a_2, a_3, a_4$  being arbitrary constants. Since  $X_i = \frac{\Xi_i}{\sqrt{n}} + \frac{\sqrt{k}}{\sqrt{n}} \sum_{j<i} Y_{ij}$  is a martingale difference array, so is  $X_i/\zeta$  where  $\zeta^2 = \sum_i \mathbb{E} X_i^2$ . The necessary conditions of theorem 24.3 of Davidson (1994) can be verified similar to the way in which the conditions of lemma D1 were verified in the proof of theorem 1.

We now analyze  $\hat{\theta}_I - \theta_0 = \hat{N}_I / \hat{D}_I$  for implicitly defined  $\hat{N}_I, \hat{D}_I$ . We have

$$n^{-1/2} \lambda^{-1} \hat{N}_I = n^{-1/2} \sum_i \tilde{g}_i u_i, \quad (77)$$

$$\begin{aligned} n^{-1/2} \lambda^{-1} \hat{D}_I &= n^{-1/2} \sum_i \tilde{g}_i Y_i = \sqrt{n} \lambda n^{-1} \sum_i \tilde{g}_i^2 + n^{-1/2} \sum_i \tilde{g}_i v_i \\ &= (\sqrt{n} \lambda) \mathbb{E} \tilde{g}_i^2 + n^{-1/2} \sum_i \tilde{g}_i v_i + O_p(\lambda), \end{aligned} \quad (78)$$

where we used the fact that  $n^{-1} \sum_i \tilde{g}_i^2 = \mathbb{E}[\tilde{g}_i^2] + O_p(n^{-1/2})$ . The result stated in (3) then follows directly from combining (76)–(78) and applying the continuous mapping theorem. We now turn to  $\hat{\theta} - \theta_0 = \hat{N} / \hat{D}$ , for which

$$\hat{N} = \sum_i \hat{g}_i u_i = \sum_{ij} w_{ij} Y_j u_i = \lambda \sum_{ij} w_{ij} \tilde{g}_j u_i + \sum_{ij} w_{ij} v_j u_i = \lambda \sum_{ij} w_{ij} (\tilde{g}_j - \tilde{g}_i) u_i$$



$$+ \lambda \sum_i \tilde{g}_i u_i + \sum_{ij} w_{ij} v_j u_i = \lambda \underbrace{\sum_i \tilde{g}_i u_i}_{O_p(\sqrt{n}\lambda)} + \underbrace{\sum_{ij} w_{ij} v_j u_i}_{O_p(\sqrt{n/k})} + o_p(\sqrt{n}\lambda). \quad (79)$$

Similarly,

$$\begin{aligned} \hat{D} = \sum_i \hat{g}_i Y_i &= \lambda \sum_i \hat{g}_i \tilde{g}_i + \sum_i \hat{g}_i v_i = \lambda^2 \underbrace{\sum_i \tilde{g}_i^2}_{O_p(n\lambda^2)} + \lambda^2 \underbrace{\sum_{ij} w_{ij} (\tilde{g}_j - \tilde{g}_i) \tilde{g}_i}_{o_p(n\lambda^2)} + \underbrace{\lambda \sum_{ij} w_{ij} \tilde{g}_i v_j}_{o_p(n\lambda^2 + \sqrt{n}\lambda + \sqrt{n/k})} \\ &\quad + \lambda \underbrace{\sum_i \tilde{g}_i v_i}_{O_p(\sqrt{n}\lambda)} + \lambda \underbrace{\sum_{ij} w_{ij} (\tilde{g}_j - \tilde{g}_i) v_i}_{o_p(\sqrt{n}\lambda)} + \underbrace{\sum_{ij} w_{ij} v_i v_j}_{O_p(\sqrt{n/k})} \end{aligned} \quad (80)$$

For RHS<sub>3</sub> in (80), note that by lemma B5,

$$\mathbb{E} \left( \sum_{ij} w_{ij} \tilde{g}_i v_j \right)^2 = \sum_{ijt} \mathbb{E} [w_{ij} w_{tj} \tilde{g}_i \tilde{g}_t v_j^2] = o(n^{3/2} k^{-1/2}),$$

which implies

$$\lambda \sum_{ij} w_{ij} \tilde{g}_i v_j = o_p(\lambda n^{3/4} k^{-1/4}) = o_p(\sqrt{n}\lambda\sqrt{\rho}) = o_p(n\lambda^2 + \rho)$$

where the last equality used the generic inequality  $2ab \leq a^2 + b^2$ . Therefore, we have

$$\hat{D} = \lambda^2 \underbrace{\sum_i \tilde{g}_i^2}_{O_p(n\lambda^2)} + \lambda \underbrace{\sum_i \tilde{g}_i v_i}_{O_p(\sqrt{n}\lambda)} + \underbrace{\sum_{ij} w_{ij} v_i v_j}_{O_p(\sqrt{n/k})} + o_p(n\lambda^2 + \sqrt{n}\lambda + \sqrt{n/k}).$$

The dominant numerator terms converge at the same rate if  $\lambda \sim 1/\sqrt{k}$  and the denominator terms when  $\lambda \sim 1/\sqrt[4]{nk}$ . Combining (79) and (80) with (76) and applying the continuous mapping theorem for each of the three cases in (4) gives the results stated therein.  $\square$