

QUADERNI DEL DIPARTIMENTO DI ECONOMIA POLITICA E STATISTICA

Federico Crudu Giovanni Mellace Zsolt Sándor

Inference in instrumental variables models with heteroskedasticity and many instruments

n. 821 – Febbraio 2020 versione aggiornata del Quaderno n.761 novembre 2017



Inference in instrumental variables models with heteroskedasticity and many instruments^{*}

Federico Crudu[†]

Giovanni Mellace[‡]

Università di Siena and CRENoS

University of Southern Denmark

Zsolt Sándor[§]

Sapientia Hungarian University of Transylvania

November 2019

Abstract

This paper proposes novel inference procedures for instrumental variable models in the presence of many, potentially weak instruments that are robust to the presence of heteroskedasticity. First, we provide an Anderson-Rubin-type test for the entire parameter vector that is valid under assumptions weaker than previously proposed Anderson-Rubin-type tests. Second, we consider the case of testing a subset of parameters under the assumption that a consistent estimator for the parameters not under test exists. We show that under the null the proposed statistics have Gaussian limiting distributions and derive alternative chi square approximations. An extensive simulation study shows the competitive finite sample properties in terms of size and power of our procedures. Finally, we provide an empirical application using college proximity instruments to estimate the returns to education.

Key words: Instrumental variables, heteroskedasticity, many instruments, jackknife, inference. *JEL classification*: C12, C13, C23.

^{*}We are grateful to Stanislav Anatolyev, Samuele Centorrino, and Neil Davies for valuable help. F. Crudu thanks financial support from the Chilean government through CONICYT's grant FONDECYT Iniciacion n. 11140433. Z. Sándor thanks financial support from grant PN-II-ID-PCE-2012-4-0066 of the Romanian Ministry of National Education, CNCS-UEFISCDI.

 $^{^\}dagger \mathrm{Department}$ of Economics and Statistics, Piazza San Francesco 7/8, 53100 Siena, Italy, federico.crudu@unisi.it

[‡]Department of Business and Economics, Campusvej 55, 5230 Odense M, Denmark, giome@sam.sdu.dk

 $^{^{\$}}$ Department of Business Sciences, Piata Libertătii 1, 530104 Miercurea Ciuc, Romania, sandorzsolt@cs.sapientia.ro

1 Introduction

The performance of test statistics based on instrumental variable (IV) models crucially depends on the quality and quantity of said IVs. In the presence of weak instruments, standard test statistics tend to deliver unreliable results. It is also well known that the number of instruments used in the construction of such tests plays a key role (see, e.g., Kleibergen, 2002, and references therein).

The Anderson-Rubin test (Anderson and Rubin, 1949, henceforth AR) is one of the most widely used statistics in the context of IV. Notoriously, this approach has the advantage of being robust to the presence of weak instruments. However, when the number of instruments grows larger than the number of parameters, the performance of the AR test starts deteriorating (e.g., Anatolyev and Gospodinov, 2011). The presence of hetero-skedasticity may exacerbate the problem.

Over the years a number of improvements on the basic formulation of the AR test have been introduced (see, e.g., Staiger and Stock, 1997; Wang and Zivot, 1998; Zivot *et al.*, 1998; Kleibergen, 2002; Stock *et al.*, 2002; Andrews *et al.*, 2006; Moreira, 2009; Andrews *et al.*, 2019). However, those tests do not consider the framework when the number of instruments grows with the sample size.

Anatolyev and Gospodinov (2011) study the limiting behavior of the AR and Sargan statistic under Bekker's many instruments framework (Bekker, 1994). Under conditional homoskedasticity they find that their test statistics are asymptotically normal and that the resulting limiting distributions depend on $\lambda = \lim_{n\to\infty} k/n$, $0 < \lambda < 1$ where k is the number of instruments and n is the sample size. Since the tests may display some size distortion when λ is close to either zero or one, the authors propose a suitable chi square approximation. Donald *et al.* (2003) and Andrews and Stock (2007) obtain similar results where the instruments are allowed to grow at slower rates.

Probably, the paper closest to ours is that of Chao *et al.* (2014), where the authors propose an overidentification test for many (weak) instruments and heteroskedasticity that

exploits the properties of the jackknife IV estimator (see Hausman *et al.*, 2012). The framework in Chao *et al.* (2014) is sufficiently general to include the Bekker's many instruments case, the many weak instruments case of Chao and Swanson (2005) and instruments that are either weak or strong. Furthermore, the ratio k/n is bounded and the number of instruments cannot grow faster than the square of the concentration parameter.

Newey and Windmeijer (2009) study generalized empirical likelihood and generalized method of moments methods in a model with moment restrictions and show that the tests of Guggenberger and Smith (2005) and Kleibergen (2005) have canonical chi square limits even when the number of instruments goes to infinity. However, the rate of growth of the instruments is slower than that in Chao *et al.* (2014). In a recent paper Bun *et al.* (2018) propose a general version of the AR test based on an Edgeworth expansion argument both for k fixed and, in the homoskedastic linear model case, for $k \to \infty$. We are not aware of any other studies that generalize the AR test to the case of many instrumental variables and heteroskedasticity.

The objective of this paper is to construct test statistics for the parameter vector of a linear IV model in the presence of many, potentially weak instruments and heteroskedasticity. The starting point of our work is the paper by Bekker and Crudu (2015). The analysis is closely related to the papers by Hausman *et al.* (2012) and Chao *et al.* (2014).

First of all, we show that the many-instrument results in Anatolyev and Gospodinov (2011) are no longer valid under heteroskedasticity. Then, we propose a test statistic to test null hypotheses on the full vector of parameters associated to both endogenous and exogenous variables and a test statistic to test null hypotheses on a subset of the parameters of the model (see e.g., Guggenberger *et al.*, 2012). In the latter case we assume the existence of a plug-in estimator that is consistent under the null hypothesis. We also allow for heteroskedasticity of unknown form. In this sense, our test statistics may be seen as generalizations of the AR test. The first statistic we introduce refers to the whole parameter vector and, under the null, does not explicitly depend on the convergence properties of the concentration parameter. On the other hand, the second test statistic is

built to test a subset of the parameter vector and relies on a consistent plug-in estimator. In this case, when the plug-in is an IV estimator, the concentration parameter plays a role in the limiting properties of the test. The assumptions on the concentration parameter match those in Bekker and Crudu (2015) and are rather similar to those in Hausman *et al.* (2012). To the best of our knowledge there is no other test on only a subset of the parameter vector in IV models with many, potentially weak instruments that allow for the presence of heteroskedasticity (for the fixed instruments model see e.g., Guggenberger *et al.*, 2012,0).

The plan of the paper is as follows. Section 2 introduces the model, Section 3 describes the test statistics, the main asymptotic results and the associated assumptions. Section 4 and Section 5 contain the simulation results and an empirical application using the college proximity instruments of Card (1995), respectively. Section 6 concludes the paper. Proofs, auxiliary results and figures are relegated to the Appendix and some additional material is available in an online Supplemental Appendix.

2 The IV model

Let us consider the model

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{1}$$

$$\boldsymbol{X} = \boldsymbol{Z}\boldsymbol{\Pi} + \boldsymbol{U} \tag{2}$$

where \boldsymbol{y} is a vector of dimension n and \boldsymbol{X} is a $n \times g$ matrix. Throughout the paper it is assumed that the $n \times k$ matrix of instruments \boldsymbol{Z} is nonstochastic and $\mathbf{E}[\boldsymbol{X}] = \boldsymbol{Z}\boldsymbol{\Pi}$, where the components of $\boldsymbol{\Pi}$ are allowed to vary with the sample size n. Such assumptions are made for convenience and may be generalized.¹ The rows of the disturbance couple ($\boldsymbol{\varepsilon}, \boldsymbol{U}$),

¹We may, for example, consider Z to be stochastic and in this case E[X] should be interpreted as a conditional expectation with respect to Z. The linearity of E[X] may also be relaxed as suggested in, e.g., Bekker (1994) and Chao *et al.* (2014).

say (ε_i, U'_i) $i = 1, \ldots, n$, are independent with zero mean and covariance matrices

$$\boldsymbol{\Sigma}_{i} = \begin{pmatrix} \sigma_{i}^{2} & \boldsymbol{\sigma}_{i12} \\ \boldsymbol{\sigma}_{i21} & \boldsymbol{\Sigma}_{i22} \end{pmatrix}$$
(3)

while the covariance matrix of the rows (y_i, X'_i) are

$$\boldsymbol{\Omega}_{i} = \begin{pmatrix} 1 & \boldsymbol{\beta}' \\ \mathbf{0} & \boldsymbol{I}_{g} \end{pmatrix} \boldsymbol{\Sigma}_{i} \begin{pmatrix} 1 & \mathbf{0} \\ \boldsymbol{\beta} & \boldsymbol{I}_{g} \end{pmatrix}.$$
(4)

3 Asymptotic results

In this section we introduce a set of assumptions that are used to prove our asymptotic results. Furthermore, we generalize a result due to Anatolyev and Gospodinov (2011) to the heteroskedastic case and we introduce our main results. In addition to that, we compare our assumptions with those introduced in other related papers and we comment on the behavior of the proposed tests when some critical assumptions are violated.

The assumptions we use are similar to those in Bekker and Crudu (2015). Additional assumptions are included to generalize some results due to Anatolyev and Gospodinov (2011). In what follows it is understood that the generic positive constant c_u may be different in different situations.

Assumption 1. The generic diagonal element P_{ii} of the projection matrix $\mathbf{P} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$ satisfies $\max_i P_{ii} \leq 1 - 1/c_u$, with $1 < c_u < \infty$. In addition, $k \to \infty$ as $n \to \infty$.

Assumption 2. The variances satisfy $\sigma_i^2 \geq \underline{\sigma}^2$ with $0 < \underline{\sigma}^2 < \infty$, for any *i*.

Assumption 3. $E[\varepsilon_i^4] \leq c_u$ and $E[||\boldsymbol{U}_i||^4] \leq c_u$ with $0 < c_u < \infty$, for any *i*.

Assumption 1 is a technical condition on the projection matrix P. It requires the main diagonal elements of P to be bounded away from 1. This assumption is rather standard in the literature (e.g., Hausman *et al.*, 2012; Bekker and Crudu, 2015) and is strictly weaker than the so called asymptotic balanced design (see Anatolyev, 2018) imposed, for example, in Anatolyev and Gospodinov (2011) and Bun *et al.* (2018) according to which all the diagonal elements of the projection matrix converge to the same constant. The assumption $k \to \infty$ as $n \to \infty$ formalizes the many instruments idea in a way that is known as Bekker asymptotics. Assumption 2 and Assumption 3 are standard regularity conditions; the former bounds variances of the disturbances away from zero, while the latter bounds the fourth moments of the errors.

3.1 The AR test under heteroskedasticity

In this section we study the limiting distribution of the AR test statistics in the presence of heteroskedasticity. In addition, our derivation implies that the test statistics we propose in Section 3.2 are also valid under homoskedasticity.

The AR statistic is a popular choice to test a null hypothesis defined as $H_0: \beta = \beta_0$. The statistic is defined as

$$AR = (n-k) \frac{\varepsilon_0' P \varepsilon_0}{\varepsilon_0' (I_n - P) \varepsilon_0}$$
(5)

and, under certain assumptions, it is asymptotically chi square distributed with k degrees of freedom. In the many instruments context and in the presence of homoskedasticity, the behavior of the AR test has been studied by Andrews and Stock (2007) and Anatolyev and Gospodinov (2011), among others. The following result generalizes the results in Lemma 1 of Anatolyev and Gospodinov (2011) to the heteroskedastic case. Let us define $\overline{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$ and $W_n = \frac{2}{k} \sum_{i \neq j} P_{ij}^2 \sigma_i^2 \sigma_j^2$.

Proposition 1. Suppose that Assumption 3 is satisfied, $\lambda = \lim_{n \to \infty} \frac{k}{n} < 1$ exists and $\frac{1}{k} \sum_{i=1}^{n} (P_{ii} - \frac{k}{n})^2 \to 0, \frac{1}{\sqrt{k}} \sum_{i=1}^{n} (P_{ii} - \frac{k}{n}) \sigma_i^2 \to 0$ hold.² In addition, assume that $\lim_{n \to \infty} \overline{\sigma}_n^2 = 1$

²The assumption $\frac{1}{\sqrt{k}} \sum_{i=1}^{n} (P_{ii} - \frac{k}{n}) \sigma_i^2 \to 0$ is needed here in order for the expected value of the statistic to converge to 0 because this does not always hold. In Example B.1 in the Supplemental Appendix we provide an instance when this property is violated in the context of indicator instruments (Bekker and Van der Ploeg, 2005).

 σ_0^2 and $\lim_{n\to\infty} W_n = W_0$ exist. Then the statistic $AR_{AG} = \sqrt{k} \left(\frac{AR}{k} - 1\right)$ proposed by Anatolyev and Gospodinov (2011) has the limit ³

$$AR_{AG} \xrightarrow{d} \mathcal{N}\left(0, \frac{W_0}{\sigma_0^4 \left(1-\lambda\right)^2}\right).$$

Remark 1. The asymptotic distribution result in Proposition 1 has two important implications. First, the asymptotic size of this test is

$$\Pr\left(AR_{AG} > \Phi^{-1}\left(1-\alpha\right)\right) = \Pr\left(\frac{\sigma_0^2\left(1-\lambda\right)}{\sqrt{W_0}}AR_{AG} < \frac{\sigma_0^2\left(1-\lambda\right)}{\sqrt{W_0}}\Phi^{-1}\left(\alpha\right)\right)$$
$$\to \Phi\left(\frac{\sigma_0^2\left(1-\lambda\right)}{\sqrt{W_0}}\Phi^{-1}\left(\alpha\right)\right).$$

Second, the test statistic T_1 proposed in Section 3.2 has broader applicability than that proposed by Anatolyev and Gospodinov even under homoskedasticity. This is because its asymptotic distribution requires the assumption that the main diagonal elements P_{ii} , i =1, ..., n, of the projection matrix \mathbf{P} should be bounded away from 1. The test statistic proposed by Anatolyev and Gospodinov (2011) requires the stronger assumption that the main diagonal elements of \mathbf{P} converge to λ . This difference in the assumptions comes from the fact that the former test statistic does not involve the diagonal elements of \mathbf{P} while the latter statistic does. The following example clarifies this concept.

Example 1. Consider indicator instruments with unequal group sizes (Bekker and Van der Ploeg, 2005). Anatolyev and Yaskov (2017, Section 5.1) show that in this case the main diagonal elements of \mathbf{P} do not converge to λ . In the Supplemental Appendix we show that under homoskedasticity the convergence in distribution $\sqrt{k} \left(\frac{AR}{k} - 1\right) \xrightarrow{d} \mathcal{N}\left(0, \frac{2}{1-\lambda}\right)$ is violated.

$$\frac{\overline{\sigma}_n^2}{\sqrt{W_n}} \to \frac{\sigma^2}{\sqrt{2(1-\lambda)}\sigma^2} = \frac{1}{\sqrt{2(1-\lambda)}}, \text{ so } \sqrt{k}\left(\frac{AR}{k} - 1\right) \stackrel{d}{\to} \mathcal{N}\left(0, \frac{2}{1-\lambda}\right),$$

which is exactly as in Lemma 1 of Anatolyev and Gospodinov (2011).

³We note that under homoskedasticity

3.2 Inference with heteroskedasticity and many instruments

In this section we present the main results. First, we present our test on the entire parameter vector. Then, we consider the more challenging case where we test the null on a subset of the coefficients; in this case we assume that a consistent plug-in estimator exists for the parameters not under test. Furthermore, we study our tests when the number of instruments is fixed. Finally, we briefly discuss the behavior of our "subset" test for some commonly encountered specific plug-in estimators and in some pathological situations.

The test statistics proposed in this paper are related to the symmetric jackknife instrumental variable estimator (SJIVE) proposed by Bekker and Crudu (2015). The SJIVE estimates consistently, in the many (weak) instruments sense, the parameter vector $\boldsymbol{\beta}$ and it is defined as

$$\widehat{\boldsymbol{\beta}}_{SJIVE} = \operatorname*{arg\,min}_{\boldsymbol{\beta}} Q_{SJIVE}(\boldsymbol{\beta}) = \operatorname*{arg\,min}_{\boldsymbol{\beta}} \frac{(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})' \boldsymbol{C}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})}{(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})' \boldsymbol{B}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})}$$
(6)

and, given the projection matrix P and the diagonal matrix D containing the diagonal elements of P,

$$oldsymbol{C} = oldsymbol{A} - oldsymbol{B}, \qquad oldsymbol{A} = oldsymbol{P} + oldsymbol{\Delta}, \qquad oldsymbol{B} = (oldsymbol{I}_n - oldsymbol{P}) oldsymbol{D} (oldsymbol{I}_n - oldsymbol{D})^{-1} (oldsymbol{I}_n - oldsymbol{P})^{-1} (oldsymbol{I}_n - oldsymbol{I}_n - oldsymbol{I$$

Consider now testing the null hypothesis $H_0: \beta = \beta_0$, where β is the true parameter vector.

The test statistic we propose is based on the numerator of the objective function in equation (6), namely,

$$Q(\boldsymbol{\beta}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta})'\boldsymbol{C}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}), \tag{7}$$

and it is defined as

$$T_1 = \frac{1}{\sqrt{k}} \frac{(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_0)' \boldsymbol{C}(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}_0)}{\sqrt{\widehat{V}(\boldsymbol{\beta}_0)}}, \quad \widehat{V}(\boldsymbol{\beta}_0) = \frac{2}{k} \boldsymbol{\varepsilon}_0^{(2)'} \boldsymbol{C}^{(2)} \boldsymbol{\varepsilon}_0^{(2)}$$
(8)

where $\varepsilon_0 = y - X\beta_0$ and the superscript "⁽²⁾" indicates the elementwise product of two conformable matrices or vectors. The following theorem provides the asymptotic distribution of the T_1 test statistic.⁴

Theorem 1. If Assumptions 1, 2, 3 are satisfied, then under H_0 : $\boldsymbol{\beta} = \boldsymbol{\beta}_0$ we have $T_1 \xrightarrow{d} \mathcal{N}(0,1).$

Let us now consider a nominal level α and let z_{α} be the α -th quantile of the normal distribution. Then, the null hypothesis is rejected if $T_1 \geq z_{1-\alpha}$.

Sometimes one is interested only in performing inference on a subset of parameters. In particular, we would like to test the coefficients associated to the endogenous variables. Let us now define the parameter vector as $\boldsymbol{\beta} = (\boldsymbol{\beta}_1', \boldsymbol{\beta}_2')'$ and suppose we want to test the following null hypothesis

$$H_0: \boldsymbol{\beta}_1 = \boldsymbol{\beta}_{10} \tag{9}$$

in the model

$$oldsymbol{y} = oldsymbol{X}oldsymbol{eta} + oldsymbol{arepsilon} = oldsymbol{X}_1oldsymbol{eta}_1 + oldsymbol{X}_2oldsymbol{eta}_2 + oldsymbol{arepsilon}$$

where the dimensions of X_1 and X_2 are $n \times g_1$ and $n \times g_2$ respectively with $g = g_1 + g_2$. Let $y_0 = y - X_1 \beta_1$ so that under the null hypothesis the model becomes

 $\boldsymbol{y}_0 = \boldsymbol{X}_2 \boldsymbol{eta}_2 + \boldsymbol{\varepsilon}.$

⁴We could apply the same type of analysis by replacing C with P - D as in Chao *et al.* (2014). We do not pursue that avenue since, as suggested in Bekker and Crudu (2015), C allows us to retain the whole signal matrix.

Accordingly, the reduced form model corresponding to X_2 is

$$\boldsymbol{X}_2 = \boldsymbol{Z}\boldsymbol{\Pi}_2 + \boldsymbol{U}_2,$$

where Π_2 and U_2 both have g_2 columns. Further, let

$$oldsymbol{H} = oldsymbol{\Pi}'oldsymbol{Z}'oldsymbol{Z}oldsymbol{\Pi} = \left(egin{array}{cc}oldsymbol{H}_{11} & oldsymbol{H}_{12}\ oldsymbol{H}_{12} & oldsymbol{H}_{22}\end{array}
ight),$$

denote the signal matrix and let $H_{22} = \Pi'_2 Z' Z \Pi_2$, which has dimension $g_2 \times g_2$.

We assume that a consistent estimator for β_2 , say $\tilde{\beta}_2$, exists. If the variables associated to β_2 are exogenous, the OLS estimator is a valid candidate. However, if this is not the case, we need a suitable IV estimator. Under the null, a consistent estimator is, for example, the SJIVE. For the null hypothesis $H_0: \beta_1 = \beta_{10}$ consider $\tilde{\beta} = (\beta'_{10}, \tilde{\beta}'_2)', \tilde{\epsilon} = y - X\tilde{\beta}$ and let the modified test statistic, denoted as T_2 , be

$$T_{2} = \frac{1}{\sqrt{k}} \frac{\widetilde{\boldsymbol{\varepsilon}}' \boldsymbol{C} \widetilde{\boldsymbol{\varepsilon}}}{\sqrt{\widehat{V}(\widetilde{\boldsymbol{\beta}})}}, \text{ where } \widehat{V}(\widetilde{\boldsymbol{\beta}}) = \frac{2}{k} \widetilde{\boldsymbol{\varepsilon}}^{(2)'} \boldsymbol{C}^{(2)} \widetilde{\boldsymbol{\varepsilon}}^{(2)}.$$
(10)

Let now $r_{\min} = \lambda_{\min}(\mathbf{H}_{22})$ and $r_{\max} = \lambda_{\max}(\mathbf{H}_{22})$ be the smallest eigenvalue and the largest eigenvalue of \mathbf{H}_{22} , respectively. Moreover, let us define a generic constant κ such that $0 \leq \kappa < \infty$.

Assumption 4. $k/r_{\min} \rightarrow \kappa$, $r_{\max}/k \rightarrow \kappa$ when $n \rightarrow \infty$.

Assumption 5. $r_{\text{max}}/k \to \kappa$, $r_{\text{min}}/k \to 0$, $\sqrt{k}/r_{\text{min}} \to 0$ when $n \to \infty$.

We have two remarks on these assumptions. First, Assumptions 4 and 5 are used in conjunction with Assumption 1 (specifically, $k \to \infty$ as $n \to \infty$), and therefore, either of them implies that $r_{\min} \to \infty$ and $r_{\max} \to \infty$ as $n \to \infty$. Second, Assumptions 4 and 5 regulate the convergence of the plug-in IV estimator. When the growth rates of r_{\min} and r_{\max} are the same, we are either in the many instruments framework of Bekker (1994) or in the many weak instruments framework of Chao and Swanson (2005). As in Chao *et al.* (2014), the growth rates of r_{\min} and r_{\max} are allowed to vary.

The following theorem provides the asymptotic distribution of the T_2 test.

Theorem 2. If Assumptions 1, 2, 3 and either 4 (many strong instruments case) or 5 (many weak instruments case) are satisfied, then $T_2 \xrightarrow{d} \mathcal{N}(0,1)$.

Analogously to the T_1 case, the null hypothesis is rejected if $T_2 \ge z_{1-\alpha}$.

It is important to derive the limiting distribution of our tests in case the number of instruments does not grow with the sample size. The following theorem provides the limiting distribution of T_1 and T_2 under the assumption that k is fixed and the error couple (ε, U) is homoskedastic.

Theorem 3. Let the disturbance couple (ε, U) be zero mean and homoskedastic and let Assumption 3 hold. Furthermore, assume (i) k fixed and $n \to \infty$, (ii) as $n \to \infty$, $\frac{Z'Z}{n} \to \Sigma_{ZZ}$ a full rank non stochastic matrix, (iii) as $n \to \infty$, $\frac{X'Z}{n} \to_p \Sigma_{XZ}$ a non stochastic matrix with rank $(\Sigma_{XZ}) = g$, (iv) as $n \to \infty$, $\frac{Z'\varepsilon}{\sqrt{n}} \to_d \mathcal{N}(\mathbf{0}, \sigma^2 \Sigma_{ZZ})$. Then $\sqrt{2kT_1} + k \to_d \chi_{k}^2$. If $\widetilde{\boldsymbol{\beta}}_2$ is the two-stage least squares estimator, then $\sqrt{2kT_2} + k \to_d \chi_{k-g_2}^2$.

While the Gaussian approximation may work well in finite samples, it does not allow us to control for the number of instruments. This, as stressed in Anatolyev and Gospodinov (2011), may be an important issue. The following corollary shows how to obtain a chi square approximation for T_1 and T_2 .

Corollary 1. If the assumptions of Theorem 1 hold true, then (i) $\sqrt{k}T_1 + k \rightarrow_d \chi_k^2$. If the assumptions of Theorem 2 hold true, then (ii) $\sqrt{k}T_2 + k \rightarrow_d \chi_k^2$, (iii) $\sqrt{k}T_2 + k \rightarrow_d \chi_{k-g_2}^2$ or (iv) $\sqrt{k-g_2}T_2 + k - g_2 \rightarrow_d \chi_{k-g_2}^2$.

Corollary 1 shows that there are different possible chi square approximations for T_2 . While approximation (*iii*) seems to be a natural candidate, also because it matches the result in Theorem 3, it may not deliver the best results in finite samples. We expect, for example, that, when k is small, approximations (ii) and (iv) enjoy better finite sample properties.⁵

The convergence properties of T_1 are determined by the behavior of the diagonal elements of P and by the properties of the disturbances. The T_2 test also depends on the properties of the plug-in estimator of the parameters not under test. When the OLS estimator $\widetilde{\boldsymbol{\beta}}_2 = (\boldsymbol{X}_2' \boldsymbol{X}_2)^{-1} \boldsymbol{X}_2' \boldsymbol{y}_0$ is a consistent plug-in estimator, that is, \boldsymbol{X}_2 is exogenous, it is easy to show that T_2 converges to a standard normal basically under the same assumptions as those imposed in Theorem 1. No further assumptions on the strength of the instruments need to be imposed (see Theorem B.1 in the Supplemental Appendix for a formal treatment). In the Supplemental Appendix (see Theorem B.2) we derive the limiting distribution of T_2 for the case when X_2 is endogenous and the JIV1 estimator is used as plug-in. The convenient expression of the JIV1 estimator allows us to explain why underrejection of the null hypothesis occurs in most cases (see Remark B.1 in the Supplemental Appendix) and to better link the weak instrument assumption $\sqrt{k}/r_{\min} \rightarrow 0$ to the asymptotic distribution of $T_{\rm 2}$ (see Remark B.2 in the Supplemental Appendix). Our derivation suggests that the assumption $\sqrt{k}/r_{\rm min} \rightarrow 0$ is likely to be necessary for the asymptotic standard normality of the statistic T_2 (see Remark B.2). In Section 4.2 below we further discuss the behavior of T_2 in relation with the limiting behavior of $\sqrt{k}/r_{\rm min}$. Specifically, we illustrate that for relatively large $\sqrt{k}/r_{\rm min}$ the histogram of T_2 differs substantially from the standard normal density. Finally, in the case when the plug-in estimator converges slowly to the true value we find that the density of T_2 is shifted to the right causing the test to overreject (see Example B.2 in the Supplemental Appendix).

3.3 Comparison with other tests

In this section we compare our set of assumptions with those used in some closely related papers. Some papers provide a broad range of results and a certain degree of variation in the specification of the assumptions. Therefore, for ease of presentation, some assumptions

 $^{^5\}mathrm{See}$ Section C in the Supplemental Appendix for some Monte Carlo evidence.

considered here are stronger than necessary (e.g., Newey and Windmeijer, 2009). In Table 1, we report the different assumptions imposed on the rate of convergence of the number of instruments and the concentration parameter, and whether or not they are robust to heteroskedasticity. Moreover, we distinguish between test statistics that consider null hypotheses on the full set of parameters or on a subset. We also report whether they allow for instruments with unbalanced design. Finally, we only consider the case where the model contains one endogenous variable. Thus, $r_{\min} = r_{\max} = r$ and $r = \pi' Z' Z \pi$ where r is the scalar version of the signal matrix H and is proportional to the concentration parameter.

There is a certain degree of heterogeneity in the type of assumptions that we show in Table 1. For example, Anatolyev and Gospodinov (2011) and Bun et al. (2018) use Bekker's framework. Our assumptions, on the other hand, are more in line with those in Chao et al. (2014), with the difference that in our case r is bounded by k, while in Chao et al. (2014) it is bounded by $n.^6$ The assumptions in Andrews and Stock (2007) and Newey and Windmeijer (2009) are to some extent similar to ours but generally their rates tend to be slower. We also notice that only Newey and Windmeijer (2009) consider AR-type tests that are robust to heteroskedasticity. Moreover, no test other than T_2 seems to explicitly consider the subset null hypothesis presented in Equation (9).

Monte Carlo simulations 4

In this Section we study the finite sample properties of the T_1 and T_2 tests in terms of size and power (see Figures 1 and 2 for the results on size and Figures 3 to 6 for the results on power).⁷ Further Monte Carlo results may be found in the Supplemental Appendix. We make inference on the full parameter vector and on the sole parameter associated to the endogenous variable. The proposed tests are compared to the version of the ARtest proposed by Anatolyev and Gospodinov (2011), denoted as AR_{AG} , and the AR test

⁶They assume either $\frac{n}{r} \to \kappa$ or $\frac{n}{r} \to 0$ and $\frac{\sqrt{k}}{r} \to 0$. ⁷The size properties of T_1 and T_2 are investigated by means of PP-plots as described in Davidson and MacKinnon (1998).

	Subset	k/n	r	Heteroskedasticity	Unbalanced
	Sabbot	10/10		meterosmetabereney	instruments
Anatolyev and Gospodinov (2011)	No	$\tfrac{k}{n} \to \lambda, 0 < \lambda < 1$	$\tfrac{r}{n} \to \kappa, \kappa \in (0,\infty)$	No	No
Andrews and Stock (2007)	No	$\frac{k^3}{n} \to 0$	$\begin{array}{c} \frac{r}{k^{\zeta}} \rightarrow \kappa_{\zeta}, \kappa_{\zeta} \in [0,\infty) \\ \zeta \in (0,\infty) \end{array}$	No	Yes
Bun <i>et al.</i> (2018)	No	$\tfrac{k}{n} \to \lambda, 0 < \lambda < 1$	_	No	No
Newey and Windmeijer (2009)	No	$\frac{k^2}{n} \rightarrow 0 \text{ or } \frac{k^3}{n} \rightarrow 0$	$\frac{\frac{n}{r} \to \kappa \text{ or } \frac{r}{n} \to 0}{\frac{k}{r} \text{ bounded}}$	Yes	Yes
T_1	No	$\frac{k}{n}$ bounded	_	Yes	Yes
T_2	Yes	$\frac{k}{n}$ bounded	$\frac{\frac{k}{r} \to \kappa \text{ or}}{\frac{r}{k} \to 0, \frac{\sqrt{k}}{r} \to 0}$	Yes	Yes

 Table 1: Comparison of assumptions in the many instruments framework.

Notes: For simplicity we refer to the single endogenous variable case where $r = \pi' Z' Z \pi = r_{\min} = r_{\max}$ and restrict ourselves to tests that use $k \to \infty$. Bun *et al.* (2018) also propose tests for the fixed k case that are robust to heteroskedasticity. Andrews and Stock (2007) and Newey and Windmeijer (2009) impose different set of assumptions depending on the problem considered and the ones reported here might be stronger than necessary in some cases.

introduced in Bun *et al.* (2018) and defined as

$$\widetilde{AR}_{df} = n\widehat{g}(\boldsymbol{\beta})'\widetilde{\boldsymbol{\Omega}}_{df}(\boldsymbol{\beta})^{-1}\widehat{g}(\boldsymbol{\beta})$$
(11)

where $\widetilde{\Omega}_{df}(\beta) = \frac{n}{n-k}\widetilde{\Omega}(\beta)$, $\widetilde{\Omega}(\beta) = \widehat{\Omega}(\beta) - \widehat{g}(\beta)\widehat{g}(\beta)'$ and $\widehat{\Omega}(\beta) = \frac{1}{n}\sum_{i=1}^{n} g(\beta)g(\beta)'$. The moment condition model is defined as $g_i(\beta) = Z_i(y_i - X'_i\beta)$ and $\widehat{g}(\beta) = \frac{1}{n}\sum_{i=1}^{n} g_i(\beta)$. In the case of T_1, T_2 and AR_{AG} , we use the corresponding chi square asymptotic distribution.⁸

This comparison is interesting for a number of reasons. First, we get a clearer idea on how much we gain by using our tests in a heteroskedastic context. Second, we get some important insights on how the considered test statistics work in the extreme cases where $\frac{k}{n} \approx 0$ and $\frac{k}{n} \approx 1$. A priori, we may expect the AR_{AG} to work well under homoskedasticity and for moderately large values of $\frac{k}{n}$, while it is probable that \widetilde{AR}_{df} performs well also in the heteroskedastic case.

Furthermore, we introduce a two parameter model with two endogenous regressors; this model is used to study the role played by the boundary condition $\sqrt{k}/r_{\min} \to 0$ and by the

⁸Due to the results in Figure C.2 in the Supplementary Appendix, for T_2 we use approximation (*ii*) in Corollary 1.

plug-in estimator in determining the behavior of T_2 (see Figure 7 and Figure 8 in Appendix B).

4.1 Data generating processes

Let us consider the Monte Carlo set up of Hausman *et al.* (2012). One of the features of this experiment is that the sum of the diagonal elements of \boldsymbol{P} does not converge to $\lambda = \lim \frac{k}{n}$, as shown in Anatolyev and Yaskov (2017). The DGP is given by

$$y = \iota \gamma + x\beta + \varepsilon$$
(12)
$$x = z\pi + v$$

where $\gamma = \beta = 1$, while $\pi = 0.1$ in the analysis of size and $\pi \in \{0.1, 1\}$ in the analysis of power. The sample size is n = 800, $\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_n)$ and independently $\boldsymbol{v} \sim \mathcal{N}(\boldsymbol{0}, 0.1^2 \times \boldsymbol{I}_n)$. The disturbances vector $\boldsymbol{\varepsilon}$ is generated as

$$\boldsymbol{\varepsilon} = \rho \boldsymbol{v} + \sqrt{\frac{1 - \rho^2}{\phi^2 + \psi^4}} (\phi \boldsymbol{w}_1 + \psi \boldsymbol{w}_2), \tag{13}$$

where $\rho = 0.3$, $\psi = 0.86$ and conditional on z, independent of v, $w_1 \sim \mathcal{N}(0, \operatorname{Diag}(z)^2)$ where $\operatorname{Diag}(z)$ is a diagonal matrix where the diagonal elements are the elements of zand $w_2 \sim \mathcal{N}(0, \psi^2 I_n)$. Notice that, $\phi = 1.38072$ implies heteroskedasticity, while $\phi = 0$ corresponds to the homoskedastic case. The instrument matrix Z is given by matrices with rows $(1, z_i, z_i^2, z_i^3, z_i^4)$ and $(1, z_i, z_i^2, z_i^3, z_i^4, z_i b_{1i}, \ldots, z_i b_{\ell i})$, $\ell = 95,695$, where, independent of other random variables, the elements $b_{1i}, \ldots, b_{\ell i}$ are i.i.d. Bernoulli distributed with $p = 1/2.^9$ We consider also two rather extreme situations: k = 2 and k = 700. We replicate our experiments 5000 times. When using the T_1 test and the T_2 test we consider $H_0: (\gamma, \beta)' = (1, 1)'$ and $H_0: \beta = 1$ respectively.¹⁰

⁹The same set of instruments is used throughout the various repetitions.

¹⁰We computed results also for $\ell = 5, 15, 35, 55, 75$ and we noticed that the p-value curves would converge from the p-value curve associated to k = 5 to the p-value curve with k = 100. This result replicates in all cases, including the power curves.

The following DGP is used to explore the properties of the T_2 test when the boundary condition $\sqrt{k}/r_{\min} \rightarrow 0$ is violated and when the plug-in estimator is inconsistent. Let us consider the following model

$$y = x\beta + w\gamma + \varepsilon$$

$$x = Z\pi_x + u_x, \quad w = Z\pi_w + u_w.$$
(14)

Let us now suppose we want to test the null H_0 : $\beta = \beta_0$, define $\eta_i = (\varepsilon_i, u_{xi}, u_{wi})'$ and assume that

$$\boldsymbol{\eta}_{i} \sim \mathcal{N}\left(\begin{pmatrix} 0\\0\\0 \end{pmatrix}, \begin{pmatrix} 1 & \rho & \xi\\\rho & 1 & 0\\\xi & 0 & 1 \end{pmatrix}\right), \quad \boldsymbol{Z}_{i} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_{k}), i = 1, \dots, n.$$
(15)

We consider two cases.

- 1. In the first case we consider n = 400. In addition to that, we set $\rho = 0.2$, $\xi = 0.3$, $\pi_x = \pi_w = \sqrt{\frac{R^2}{k(1-R^2)}} \iota_k$, k = n/2 and R^2 is chosen in such a way that $\frac{\sqrt{k}}{n\pi'_w\pi_w} = 0.1$ and $\frac{\sqrt{k}}{n\pi'_w\pi_w} = 1.^{11}$ Finally, as a plug-in estimator we use the JIV1 estimator. The number of replications is 5000.
- 2. The sample size is set to n = 400. Moreover, $\rho = 0.2, \xi \in \{0, 0.1, 0.2, 0.3\}, k \in \{2, 20, 200\}, \pi_x = \pi_w = \sqrt{\frac{R^2}{k(1-R^2)}} \iota_k$ with $R^2 = 0.2$. Finally, as a plug-in estimator we use the OLS estimator. The number of replications is 5000.

4.2 Simulation results

We first provide some interpretation of the simulations by separately analyzing the results on size and power. Then we discuss the behavior of T_2 when an inconsistent plug-in is

¹¹The condition $\frac{\sqrt{k}}{n\pi'_w\pi_w} = 1$ replicates the idea that the boundary condition $\sqrt{k}/r_{\min} \to 0$ is violated. We did run simulations also for $\frac{\sqrt{k}}{n\pi'_w\pi_w} = 10$ and n = 100,200 finding similar results.

used.

Size. Analyzing Figure 1 and Figure 2 we observe that, in general, T_1 and T_2 work well in all the considered cases.¹² The \widetilde{AR}_{df} test, on the other hand, works well for most of the cases but tends overreject when k is large. Finally, as expected, the AR_{AG} test overrejects for any value of k and its performance deteriorates as k increases.

Power. The power properties of the various test statistics display some interesting patterns. When k = 2, 5, the T_1 and the T_2 tests along with the AR_{df} test are able to discriminate among alternatives (Figure 3 to Figure 6 panels (a) and (b)). To some extent the same could be said about the case where k = 100 (Figure 3 to Figure 6 panel (c)). Finally, when $k = 700, \pi = 0.1$, the T_1 and the T_2 tests are unable to discriminate among alternatives. More precisely, no test statistic among those considered seems to work well in this case. However, when $\pi = 1$, the T_1 and T_2 tests tend to outperform their competitors (Figure 3 to Figure 6 panel (d)). It is interesting to notice that the properties of T_1 and T_2 are affected by a trade off between size and power with respect to k: as k grows the empirical size approaches the nominal size, but the power curves tend to get wider. This may be a problem when the instruments are weak as the tests may eventually have no power for k large. When the instruments are stronger the effect of such a trade-off is less severe and our tests work well even in the extreme case with k = 700.

Over/underrejection. The comparison of the histograms and QQ-plots in Figure 7 displays how the T_2 test behaves when the boundary condition $\sqrt{k}/r_{\min} \rightarrow 0$ is violated. In particular, we notice that the (empirical) density tends to be more leptokurtic with respect to its asymptotic counterpart. This feature induces the test to underreject. On the other hand, the plots in Figure 8 show the behavior of the T_2 statistic when a slow plug-in estimator is used. We notice that the use of OLS instead of a more appropriate IV estimator causes T_2 to overreject. In particular, T_2 overrejects more as ξ increases. Furthermore, we notice that the overrejection tendency is mitigated by the increased number of instruments.

¹²It is worth noticing that, in general, for the hypothesis $H_0: \beta_1 = \beta_{10}$ all the tests tend to underreject for small values of k.

Violation of the boundary condition. Figure 7 illustrates the behavior of the T_2 test with a small \sqrt{k}/r_{\min} (Figure 7(a)) and a large \sqrt{k}/r_{\min} (Figure 7(b)). In the latter case, the histogram of T_2 differs substantially from the one of a standard normal density, suggesting that $\sqrt{k}/r_{\min} \rightarrow 0$ is important for the asymptotic normality of our test.

5 Empirical application

In this section we apply our methods to the data from the National Longitudinal Survey of Young Men (NLSYM) used by Card (1995) to estimate the returns to education. The data set includes 3010 observations and 35 variables.¹³

We consider two different models to estimate the returns of education. Both models assume that the log of wages (*wage*) is a linear function of education measured in years of schooling (*school*) and a set of exogenous variables \boldsymbol{x} , namely

 $\log(wage_i) = \beta school_i + \boldsymbol{x}'_i \boldsymbol{\gamma} + \varepsilon_i.$

Similar to Kleibergen (2004), \boldsymbol{x} includes a constant and binary variables for race, residence in a metropolitan area, and residence in the south of the United States as well as IQ test score. As experience is measured simply as age - school - 6 in this data, we do not use it as a control variable in our models. ¹⁴ For the instruments, following once again Kleibergen (2004), in our first specifications we use age and age square and two variables that indicate college proximity. The exogeneity of the college proximity instruments is somewhat questionable for several reasons. For example, areas with a high prevalence of people with high unobserved ability may be more likely to have a college nearby. Card (1995) argues that including other observable characteristics, as we do, should mitigate this issue. However, we cannot completely exclude the potential endogeneity of our instruments.

¹³The data are from the R package ivmodel of Jiang *et al.* (2016).

¹⁴Another reason not to control for experience when estimating returns of education, at least in this data, is that experience is mechanically an outcome of education and it is therefore a bad control as discussed for example in Angrist and Pischke (2008).

In our second specification, we generate additional excluded instruments by interacting age, age squared, and the two college proximity variables with the geographical indicators and race. In the first specification, the instrument set includes four variables, while in the second it includes fourteen variables.

It is very likely that the variance of the error depends on the exogenous variables which motivates the use of our T_2 test for inference. For example, it appears very plausible that the conditional variance of the unobservables driving wages differ by college proximity, location as well as race. We run our T_2 statistic, using both the chi square (T_2) and Gaussian (T_2^{gauss}) approximations, the AR_{AG} statistic of Anatolyev and Gospodinov (2011), the \widetilde{AR}_{df} statistics of Bun *et al.* (2018), and the standard AR statistic to test 301 equidistant values in the interval [0, 3] for the coefficient of education β . The results for the model with four instruments are reported in Figure 9. With only 4 excluded instruments all the tests give very similar results, in particular they are not able to reject values in the (approximate) interval [.5, 1.5], at the 5% significance level, and in the (approximate) interval [.5, .1], at the 10% significance level. The implied effects of education on wages are much higher than the one found in Card's study who, however, includes experience (which is arguably a "bad control") in his model. However, the large effects implied by our models are in line with the one found in Imbens and Rubin (1997).

The results with fourteen instruments are reported in Figure 10. Probably due to the presence of heteroskedasticity, adding instruments deteriorates the performances of the AR, AR_{AG} , and \widehat{AR}_{df} tests, which reject every single value of β at the 10% significance level. On the other hand, increasing the number of instruments does not seem to have a big impact on our T_2 and T_2^{gauss} tests. These results are in line with what we find in our simulation study where, with strong heteroskedasticity, the performances of the AR_{AG} deteriorates when we increase the number of instruments and the \widehat{AR}_{df} tends to overreject with many instruments.

6 Conclusion

This paper introduces two test statistics for the parameters of a linear model in the presence of endogeneity, heteroskedasticity and many, potentially weak, instruments. The tests are easy to build as they are based on the numerator of the SJIVE estimator proposed by Bekker and Crudu (2015). We prove that, after appropriate rescaling, the limiting distribution of the test statistics are standard normal. Moreover, simulation evidence shows that, in finite samples, the proposed tests generally outperform their competitors, such as the AR tests proposed in Anatolyev and Gospodinov (2011) and in Bun *et al.* (2018).

In our empirical application, the standard AR test and its modification by Anatolyev and Gospodinov (2011), probably due to the presence of heteroskedasticity, reject every single value chosen for the null when we increase the number of instruments from four to fourteen. On the other hand, our proposed statistic provides similar results independently of the number of instruments used.

The tests we propose can be applied broadly to any linear overidentified IV model and they are particularly appealing for the growing literature using genetic markers as instruments, see for example Von Hinke *et al.* (2016). In this literature, the number of instruments is potentially very large and the instruments are typically weak, a framework where our tests potentially outperform existing methods. Another potential field of application for our tests is the framework of Kang *et al.* (2016) and Windmeijer *et al.* (2017) where inference is carried out after a potentially large set of valid instruments is selected via LASSO.

Appendix

A Proofs of main results

This section contains the proofs of the main theorems and some auxiliary results. In what follows it is understood that O is a conformable matrix of zeros and that the abbreviations LLN, CLT and IID stand for law of large numbers, central limit theorem and independently and identically distributed respectively. In addition to that, $\sum_{i \neq j}$ is a double sum for $i, j = 1, \ldots, n$ that excludes the same index elements and $\sum_{i,j,k,\ell}$ replaces the quadruple sum $\sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{\ell=1}^{n}$. Triple sums are defined similarly.

Proof of Proposition 1. Under $H_0: \boldsymbol{\beta} = \boldsymbol{\beta}_0$ we have

$$\sqrt{k}\left(\frac{AR}{k}-1\right) = \frac{\frac{1}{\sqrt{k}}\left(\frac{n-k}{k}\boldsymbol{\varepsilon}'\boldsymbol{P}\boldsymbol{\varepsilon}-\boldsymbol{\varepsilon}'\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\varepsilon}\right)}{\frac{1}{k}\boldsymbol{\varepsilon}'\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\varepsilon}} = \frac{n}{k}\frac{\frac{1}{\sqrt{k}}\left(\boldsymbol{\varepsilon}'\boldsymbol{P}\boldsymbol{\varepsilon}-\frac{k}{n}\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}\right)}{\frac{1}{k}\boldsymbol{\varepsilon}'\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\varepsilon}}.$$
(16)

Note that

$$\frac{1}{\sqrt{k}}\left(\boldsymbol{\varepsilon}'\boldsymbol{P}\boldsymbol{\varepsilon} - \frac{k}{n}\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}\right) = \frac{1}{\sqrt{k}}\sum_{i\neq j}P_{ij}\varepsilon_i\varepsilon_j + \frac{1}{\sqrt{k}}\sum_{i=1}^n\left(P_{ii} - \frac{k}{n}\right)\varepsilon_i^2 \equiv E_1 + E_2.$$
 (17)

We can apply the CLT from (Chao et al., 2012, Lemma A2) to the quadratic form

$$R = \sum_{i \neq j} P_{ij} \varepsilon_i \varepsilon_j$$

involved in E_1 . We obtain that

$$\frac{R}{\sqrt{kW_n}} \stackrel{d}{\to} \mathcal{N}(0,1) \,,$$

where

$$W_n = \frac{\operatorname{Var}[R]}{k} = \frac{2}{k} \sum_{i \neq j} P_{ij}^2 \sigma_i^2 \sigma_j^2$$

with the property that

$$\frac{1}{k}\operatorname{Var}[R] = \frac{2}{k}\sum_{i\neq j}P_{ij}^2\sigma_i^2\sigma_j^2 \ge \frac{2\underline{\sigma}^4}{k}\sum_{i\neq j}P_{ij}^2 \ge \frac{2\underline{\sigma}^4}{c_u},$$

(the latter inequality comes from (22)), which is bounded away from 0. Consequently, W_n is bounded between two positive numbers. We obtain that $E_1/\sqrt{W_n} \xrightarrow{d} \mathcal{N}(0,1)$. Regarding E_2 , by the assumption $\frac{1}{\sqrt{k}} \sum_{i=1}^n (P_{ii} - \frac{k}{n}) \sigma_i^2 \to 0$ we have

$$\mathbf{E}[E_2] = \frac{1}{\sqrt{k}} \sum_{i=1}^n \left(P_{ii} - \frac{k}{n} \right) \sigma_i^2 \to 0.$$

Further, by Assumption 3

$$\operatorname{Var}[E_2] = \frac{1}{k} \sum_{i} \left(P_{ii} - \frac{k}{n} \right)^2 \operatorname{Var}\left[\varepsilon_i^2 \right] \le \frac{c_u}{k} \sum_{i=1}^n \left(P_{ii} - \frac{k}{n} \right)^2.$$

Using the assumption $\frac{1}{k} \sum_{i=1}^{n} (P_{ii} - \frac{k}{n})^2 \to 0$, we obtain that $\operatorname{Var}[E_2] = o(1)$. Then by Chebyshev's inequality $E_2 = o_p(1)$. Therefore,

$$\frac{E_1 + E_2}{\sqrt{W_n}} \xrightarrow{d} \mathcal{N}(0, 1) \,. \tag{18}$$

Regarding the denominator involved in (16) we observe that

$$\frac{1}{k}\boldsymbol{\varepsilon}'\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\varepsilon} = \frac{1}{k}\left(1-\frac{k}{n}\right)\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} - \frac{1}{k}\boldsymbol{\varepsilon}'\left(\boldsymbol{P}-\frac{k}{n}\boldsymbol{I}\right)\boldsymbol{\varepsilon}.$$

The second term is just the expression from (17) divided by \sqrt{k} , that is,

$$\frac{1}{k}\boldsymbol{\varepsilon}'\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\varepsilon} = \frac{1}{k}\left(1-\frac{k}{n}\right)\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} - \frac{1}{\sqrt{k}}\left(E_1+E_2\right) = \frac{1}{k}\left(1-\frac{k}{n}\right)\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} + O_p\left(\frac{1}{\sqrt{k}}\right).$$

Using Assumption 3 and the LLN, using the notation

$$\overline{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n \sigma_i^2$$

we have that

$$\frac{1}{n}\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} - \overline{\sigma}_n^2 = O_p\left(\frac{1}{\sqrt{k}}\right). \tag{19}$$

Consequently,

$$\frac{1}{k}\boldsymbol{\varepsilon}'\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\varepsilon}=\frac{n}{k}\left(1-\frac{k}{n}\right)\overline{\sigma}_{n}^{2}+O_{p}\left(\frac{1}{\sqrt{k}}\right).$$

Now, from equation (16) and the fact that $\frac{n}{k}\left(1-\frac{k}{n}\right)\overline{\sigma}_n^2$ is bounded between two positive numbers, we have

$$\begin{split} \sqrt{k} \left(\frac{AR}{k} - 1\right) &= \frac{n}{k} \frac{\frac{1}{\sqrt{k}} \left(\boldsymbol{\varepsilon}' \boldsymbol{P} \boldsymbol{\varepsilon} - \frac{k}{n} \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}\right)}{\frac{n}{k} \left(1 - \frac{k}{n}\right) \overline{\sigma}_n^2} + \frac{n}{k} \frac{\frac{1}{\sqrt{k}} \left(\boldsymbol{\varepsilon}' \boldsymbol{P} \boldsymbol{\varepsilon} - \frac{k}{n} \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon}\right)}{\frac{k}{k} \left(1 - \frac{k}{n}\right) \overline{\sigma}_n^2} \left(\frac{\frac{n}{k} \left(1 - \frac{k}{n}\right) \overline{\sigma}_n^2}{\frac{1}{k} \boldsymbol{\varepsilon}' \left(\boldsymbol{I} - \boldsymbol{P}\right) \boldsymbol{\varepsilon}} - 1\right) \\ &= \frac{E_1 + E_2}{\left(1 - \frac{k}{n}\right) \overline{\sigma}_n^2} + o_p\left(1\right). \end{split}$$

Therefore, collecting the above results we obtain that

$$\left(1-\frac{k}{n}\right)\frac{\overline{\sigma}_{n}^{2}}{\sqrt{W_{n}}}\sqrt{k}\left(\frac{AR}{k}-1\right) = \frac{E_{1}+E_{2}}{\sqrt{W_{n}}} + o_{p}\left(1\right),$$

which by (18) implies that

$$\left(1 - \frac{k}{n}\right) \frac{\overline{\sigma}_n^2}{\sqrt{W_n}} \sqrt{k} \left(\frac{AR}{k} - 1\right) \stackrel{d}{\to} \mathcal{N}(0, 1) \,.$$

$$(20)$$

Since we assume that $\lim_{n\to\infty} \overline{\sigma}_n^2 = \sigma_0^2$ and $\lim_{n\to\infty} W_n = W_0$ exist, we obtain the result. \Box

In the proof of Theorem 1 we use the following CLT, which, as argued by Bekker and Crudu (2015, Appendix A.4) can be proved in a way similar to Lemma A2 from Chao *et al.*

(2012).

Lemma A.1. Consider the quadratic form $Q = \sum_{i \neq j} C_{ij} \varepsilon_i \varepsilon_j$, where C_{ij} is the (i, j) element of matrix C that is symmetric and has zero main diagonal elements. Suppose that there is a matrix P that is symmetric, idempotent, $P_{ii} \leq c_u < 1$, $|C_{ij}| \leq c_u |P_{ij}|$ for any $i \neq j$, and rank(P) = k, where $k \to \infty$ as $n \to \infty$, and the following properties hold: (a) $E[\varepsilon_i] = 0$ and $\varepsilon_1, ..., \varepsilon_n$ are independent; (b) $E[\varepsilon_i^4] < \infty$; (c) $\frac{1}{k} Var[Q] \geq c_u > 0$. Then,

$$\frac{Q}{\sqrt{\operatorname{Var}[Q]}} \xrightarrow{d} \mathcal{N}(0,1) \,.$$

Lemma A.2. Let $\widehat{V}(\boldsymbol{\beta}_0) = \frac{2}{k} \boldsymbol{\varepsilon}_0^{(2)'} \boldsymbol{C}^{(2)} \boldsymbol{\varepsilon}_0^{(2)}$. If Assumptions 1, 3 hold, $\widehat{V}(\boldsymbol{\beta}_0) - V_n = O_p\left(\frac{1}{\sqrt{k}}\right)$; consequently $\widehat{V}(\boldsymbol{\beta}_0) - V_n \xrightarrow{p} 0$, where

$$V_{n} = \frac{2}{k} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} \sigma_{i}^{2} \sigma_{j}^{2}.$$

(For a proof see Supplemental Appendix A.)

Proof of Theorem 1. Under the null hypothesis we have

$$\mathbf{E} \left[\boldsymbol{\varepsilon}_{0}^{\prime} \boldsymbol{C} \boldsymbol{\varepsilon}_{0} \right] = 0,$$

$$\operatorname{Var} \left[\boldsymbol{\varepsilon}_{0}^{\prime} \boldsymbol{C} \boldsymbol{\varepsilon}_{0} \right] = \mathbf{E} \left[\left(\boldsymbol{\varepsilon}_{0}^{\prime} \boldsymbol{C} \boldsymbol{\varepsilon}_{0} \right)^{2} \right] = 2 \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} \sigma_{i}^{2} \sigma_{j}^{2} \equiv k V_{n}.$$

We verify the conditions of the CLT stated in Lemma A.1 for C and P defined in Section 2. The properties of C and P hold by definition, Assumption 1 and the fact that $|C_{ij}| = \frac{|P_{ij}|}{2} \left(\frac{1}{1-P_{ii}} + \frac{1}{1-P_{jj}}\right) \leq c_u |P_{ij}| \text{ for any } i, j, \text{ (see the proof of Lemma A.2 in the Supplemental Appendix).}$

Further, (a) is clearly satisfied; (b) is satisfied due to Assumption 3. Regarding (c) note that by Assumption 2

$$\frac{1}{k} \operatorname{Var} [Q] \equiv V_n = \frac{2}{k} \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \sigma_i^2 \sigma_j^2 \ge \frac{2\underline{\sigma}^4}{k} \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2,$$

where

$$\sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} = \sum_{i \neq j} \frac{P_{ij}^{2}}{4} \left(\frac{1}{1 - P_{ii}} + \frac{1}{1 - P_{jj}} \right)^{2} \ge \sum_{i \neq j} \frac{P_{ij}^{2}}{4} (1 + 1)^{2} = \sum_{i \neq j}^{n} P_{ij}^{2}$$
$$= \sum_{i=1}^{n} \sum_{j=1}^{n} P_{ij}^{2} - \sum_{i=1}^{n} P_{ii}^{2} = tr\left(\mathbf{P}\right) - \sum_{i=1}^{n} P_{ii}^{2} = k - \sum_{i=1}^{n} P_{ii}^{2}.$$
(21)

By Assumption 1

$$\sum_{i=1}^{n} P_{ii}^{2} \le \max P_{ii} \sum_{i=1}^{n} P_{ii} \le (1 - 1/c_{u}) tr \left(\boldsymbol{P}\right) = (1 - 1/c_{u}) k.$$
(22)

 So

$$\sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \ge k/c_u$$

therefore,

$$\frac{1}{k} \operatorname{Var}\left[Q\right] \ge \frac{2\underline{\sigma}^4}{c_u},$$

which is bounded away from 0. In this case we can apply the CLT in Lemma A.1 and complete the proof by using Lemma A.2. $\hfill \Box$

For the proof of Theorem 2 we need the following results (for proofs see Supplemental Appendix A).

Lemma A.3. Let $\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) = \frac{2}{k} \widetilde{\boldsymbol{\epsilon}}^{(2)\prime} \boldsymbol{C}^{(2)} \widetilde{\boldsymbol{\epsilon}}^{(2)}$. If $\widetilde{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$ and Assumptions 1, 3 hold, then $\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) - V_n \xrightarrow{p} 0$.

Lemma A.4. Under Assumptions 1, 3

1. $\operatorname{E} [\mathbf{X}_{2}'\mathbf{C}\mathbf{X}_{2}] = \mathbf{H}_{22}, \operatorname{Var} [\mathbf{X}_{2}'\mathbf{C}\mathbf{X}_{2}] \leq c_{u}\mathbf{H}_{22} + c_{u}k\mathbf{I}_{g_{2}} + c_{u}r_{\max}\mathbf{I}_{g_{2}},$ 2. $\operatorname{E} [\mathbf{X}_{2}'\mathbf{C}\boldsymbol{\varepsilon}] = \mathbf{0}, \operatorname{Var} [\mathbf{X}_{2}'\mathbf{C}\boldsymbol{\varepsilon}] \leq c_{u}\mathbf{H}_{22} + c_{u}k\mathbf{I}_{g_{2}}.$

Before proceeding to the proof of Theorem 2 we present some general facts that are used in the proofs of several results below. Consider a plug-in estimator $\tilde{\beta}_2$ of β_2 and, as above, let $\widetilde{\boldsymbol{\beta}} = \left(\boldsymbol{\beta}_1', \widetilde{\boldsymbol{\beta}}_2'\right)'$. Notice that under $H_0: \boldsymbol{\beta}_1 = \boldsymbol{\beta}_{10}$ it holds that

$$\boldsymbol{X}\left(\widetilde{\boldsymbol{eta}}-\boldsymbol{eta}
ight)=\boldsymbol{X}_{1}\left(\boldsymbol{eta}_{10}-\boldsymbol{eta}_{1}
ight)+\boldsymbol{X}_{2}\left(\widetilde{\boldsymbol{eta}}_{2}-\boldsymbol{eta}_{2}
ight)=\boldsymbol{X}_{2}\left(\widetilde{\boldsymbol{eta}}_{2}-\boldsymbol{eta}_{2}
ight),$$

so the residual vector can be written as

$$\widetilde{oldsymbol{arepsilon}} = oldsymbol{y} - oldsymbol{X} \widetilde{oldsymbol{eta}} = oldsymbol{arepsilon} - oldsymbol{X} \left(\widetilde{oldsymbol{eta}} - oldsymbol{eta}
ight) = oldsymbol{arepsilon} - oldsymbol{X}_2 \left(\widetilde{oldsymbol{eta}}_2 - oldsymbol{eta}_2
ight).$$

The statistic T_2 can be rewritten as

$$T_{2} = \frac{1}{\sqrt{k}} \frac{\varepsilon' C\varepsilon}{\sqrt{\widehat{V}(\beta)}} \left(\frac{\sqrt{\widehat{V}(\beta)}}{\sqrt{\widehat{V}(\widehat{\beta})}} - 1 \right) + \frac{1}{\sqrt{k}} \frac{\Delta}{\sqrt{\widehat{V}(\widehat{\beta})}} + \frac{1}{\sqrt{k}} \frac{\varepsilon' C\varepsilon}{\sqrt{\widehat{V}(\beta)}}$$

$$\equiv B_{1} + B_{2} + B_{3},$$
(23)

where

$$\Delta = \left(\widetilde{\boldsymbol{\beta}}_{2} - \boldsymbol{\beta}_{2}\right)' \boldsymbol{X}_{2}' \boldsymbol{C} \boldsymbol{X}_{2} \left(\widetilde{\boldsymbol{\beta}}_{2} - \boldsymbol{\beta}_{2}\right) - 2 \left(\widetilde{\boldsymbol{\beta}}_{2} - \boldsymbol{\beta}_{2}\right)' \boldsymbol{X}_{2}' \boldsymbol{C} \boldsymbol{\varepsilon}.$$
(24)

The first term is equal to

$$B_{1} = \frac{1}{\sqrt{k}} \frac{\boldsymbol{\varepsilon}' \boldsymbol{C} \boldsymbol{\varepsilon}}{\sqrt{\widehat{V}(\boldsymbol{\beta})}} \frac{\sqrt{\widehat{V}(\boldsymbol{\beta})} - \sqrt{\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right)}}{\sqrt{\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right)}},$$

where from Lemma A.3 and the consistency of the plug-in it follows that $\sqrt{\hat{V}(\beta)} - \sqrt{\hat{V}(\tilde{\beta})} = o_p(1)$, while since V_n is bounded away from 0 by Assumption 3, it follows that $1/\sqrt{\hat{V}(\tilde{\beta})} = O_p(1)$. Theorem 1 implies that $\frac{1}{\sqrt{k}} \frac{\epsilon' C \epsilon}{\sqrt{\hat{V}(\beta)}} = O_p(1)$, so $B_1 = o_p(1)$. Regarding B_3 , from Theorem 1 we have that $B_3 \to_d \mathcal{N}(0, 1)$.

Consequently, if the plug-in estimator $\widetilde{\beta}$ is consistent then under Assumptions 1, 3 we

have $B_1 = o_p(1)$ and $B_3 \to_d \mathcal{N}(0, 1)$. In order to derive the asymptotic distribution of T_2 we need to study the term B_2 .

Proof of Theorem 2. Note that the first term from Δ in (24) is

$$\left(\widetilde{\boldsymbol{\beta}}_{2}-\boldsymbol{\beta}_{2}\right)'\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{X}_{2}\left(\widetilde{\boldsymbol{\beta}}_{2}-\boldsymbol{\beta}_{2}\right)=\left(\widetilde{\boldsymbol{\beta}}_{2}-\boldsymbol{\beta}_{2}\right)'\boldsymbol{H}_{22}^{1/2}\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{X}_{2}\boldsymbol{H}_{22}^{-1/2}\boldsymbol{H}_{22}^{1/2}\left(\widetilde{\boldsymbol{\beta}}_{2}-\boldsymbol{\beta}_{2}\right).$$
(25)

First we show that $H_{22}^{-1/2}X'_2CX_2H_{22}^{-1/2} \xrightarrow{p} I_{g_2}$. Lemma A.4 implies that

$$\mathbb{E}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\boldsymbol{H}_{22}^{-1/2}\right] = \boldsymbol{I}_{g_{2}}$$

$$\tag{26}$$

and

$$\operatorname{Var}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\boldsymbol{H}_{22}^{-1/2}\right] \leq \frac{1}{r_{\min}}\boldsymbol{H}_{22}^{-1/2}\left(c_{u}\boldsymbol{H}_{22} + c_{u}k\boldsymbol{I}_{g_{2}} + c_{u}r_{\max}\boldsymbol{I}_{g_{2}}\right)\boldsymbol{H}_{22}^{-1/2}$$
$$= \frac{1}{r_{\min}}\left(c_{u} + c_{u}\frac{k}{r_{\min}} + c_{u}\frac{r_{\max}}{r_{\min}}\right)\boldsymbol{I}_{g_{2}},$$

where the inequality holds due to $\boldsymbol{H}_{22}^{-1} \leq \frac{1}{r_{\min}} \boldsymbol{I}_{g_2}$. Therefore, Assumption 4 (many strong instruments case) implies that $\operatorname{Var} \left[\boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \boldsymbol{H}_{22}^{-1/2} \right] = O\left(\frac{1}{k}\right)$ while Assumption 5 (many weak instruments case) implies that $\operatorname{Var} \left[\boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \boldsymbol{H}_{22}^{-1/2} \right] = o(1)$. In either case we obtain that $\operatorname{Var} \left[\boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \boldsymbol{H}_{22}^{-1/2} \right] \to 0$, which together with (26) implies that $\boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \boldsymbol{H}_{22}^{-1/2} \right] \to 0$, which together with (26) implies that $\boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \boldsymbol{H}_{22}^{-1/2} \to \boldsymbol{I}_{g_2}$.

Note that under Assumption 4 $\boldsymbol{H}_{22}^{1/2}\left(\boldsymbol{\tilde{\beta}}_2-\boldsymbol{\beta}_2\right) = O_p(1)$ while under Assumption 5 $\frac{1}{\sqrt{k}}\boldsymbol{H}_{22}\left(\boldsymbol{\tilde{\beta}}_2-\boldsymbol{\beta}_2\right) = O_p(1)$ (see Section 4 in Bekker and Crudu, 2015). Therefore, under either Assumption 4 or Assumption 5, from (25) we conclude that

$$\frac{1}{\sqrt{k}} \left(\widetilde{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \right)' \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \left(\widetilde{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \right) = o_p \left(1 \right).$$
(27)

The second term from Δ in (24) involves

$$\left(\widetilde{oldsymbol{eta}}_2-oldsymbol{eta}_2
ight)'oldsymbol{X}_2'oldsymbol{C}oldsymbol{arepsilon}=\left(\widetilde{oldsymbol{eta}}_2-oldsymbol{eta}_2
ight)'oldsymbol{H}_{22}oldsymbol{H}_{22}^{-1}oldsymbol{X}_2'oldsymbol{C}oldsymbol{arepsilon}.$$

Lemma A.4 implies that

$$\mathbf{E}\left[\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] = \boldsymbol{0} \tag{28}$$

and

$$\operatorname{Var}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] \leq c_{u}\left(\frac{1}{r_{\min}} + \frac{k}{r_{\min}^{2}}\right)\boldsymbol{I}_{g_{2}},\tag{29}$$

where the latter inequality is due to $H_{22}^{-1} \leq \frac{1}{r_{\min}} I_{g_2}$. We also obtain that

$$\operatorname{Var}\left[\frac{1}{\sqrt{k}}\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] \leq c_{u}\left(\frac{1}{k} + \frac{1}{r_{\min}}\right)\boldsymbol{I}_{g_{2}}.$$
(30)

Under Assumption 4 (many strong instruments case) we get $\operatorname{Var}\left[\frac{1}{\sqrt{k}}\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{\varepsilon}\right] = O\left(\frac{1}{k}\right)$, which together with (28) implies that $\frac{1}{\sqrt{k}}\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{\varepsilon} = o_{p}\left(1\right)$. Since $\boldsymbol{H}_{22}^{1/2}\left(\widetilde{\boldsymbol{\beta}}_{2} - \boldsymbol{\beta}_{2}\right) = O_{p}\left(1\right)$ holds, we obtain

$$rac{1}{\sqrt{k}}\left(\widetilde{oldsymbol{eta}}_2-oldsymbol{eta}_2
ight)'oldsymbol{X}_2'oldsymbol{C}oldsymbol{arepsilon}=o_p\left(1
ight).$$

Under Assumption 5 (many weak instruments case) (29) implies $\operatorname{Var}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{\varepsilon}\right] = o(1)$, which together with (28) implies that $\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{\varepsilon} = o_{p}(1)$. Since $\frac{1}{\sqrt{k}}\boldsymbol{H}_{22}\left(\boldsymbol{\widetilde{\beta}}_{2}-\boldsymbol{\beta}_{2}\right) = O_{p}(1)$ holds, we obtain

$$\frac{1}{\sqrt{k}}\left(\widetilde{\boldsymbol{\beta}}_{2}-\boldsymbol{\beta}_{2}\right)'\boldsymbol{X}_{2}'\boldsymbol{C}\boldsymbol{\varepsilon}=o_{p}\left(1\right).$$

Regarding B_3 , from Theorem 1 we have that $B_3 \xrightarrow{d} \mathcal{N}(0,1)$.

Proof of Theorem 3. Note that

$$\sum_{i \neq j} C_{ij} \varepsilon_i \varepsilon_j = \sum_{i \neq j} \frac{P_{ij}}{2} \left(\frac{1}{1 - P_{ii}} + \frac{1}{1 - P_{jj}} \right) \varepsilon_i \varepsilon_j = (1 + o(1)) \sum_{i \neq j} P_{ij} \varepsilon_i \varepsilon_j$$

as $\max_i P_{ii} \to 0$. Further,

$$\sum_{i \neq j} C_{ij} \varepsilon_i \varepsilon_j = (1 + o(1)) \boldsymbol{\varepsilon}' \boldsymbol{P} \boldsymbol{\varepsilon} - (1 + o(1)) \sum_{i=1}^n P_{ii} \varepsilon_i^2.$$

By assumptions (*ii*) and (*iv*), $\varepsilon' P \varepsilon \rightarrow_d \sigma^2 \chi_k^2$. Moreover, by independence of ε_i , Assumption 3 and the properties of P_{ii}

$$\mathbf{E}\left[\left(\sum_{i=1}^{n} P_{ii}\varepsilon_{i}^{2} - k\sigma^{2}\right)^{2}\right] \to 0,$$

which implies $\sum_{i=1}^{n} P_{ii} \varepsilon_i^2 \rightarrow_p k\sigma^2$ (see Chao *et al.*, 2014). Consider now

$$V_n = \frac{2}{k} \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \sigma_i^2 \sigma_j^2 = \frac{2\sigma^4}{k} \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 = \frac{2\sigma^4}{k} \sum_{i=1}^n \sum_{j=1}^n \frac{P_{ij}^2}{4} \left(\frac{1}{1 - P_{ii}} + \frac{1}{1 - P_{jj}}\right)^2.$$

Since $\sum_{i=1}^{n} \sum_{j=1}^{n} P_{ij}^2 = k$ is fixed and $\max_i P_{ii} \to 0, V_n \to 2\sigma^4$. Hence, by Lemma A.2, $\widehat{V}(\boldsymbol{\beta}_0) \to_p 2\sigma^4$. Finally,

$$T_1 = \frac{1}{\sqrt{k}} \frac{\boldsymbol{\varepsilon}_0' \boldsymbol{C} \boldsymbol{\varepsilon}_0}{\sqrt{\frac{2}{k} \boldsymbol{\varepsilon}_0^{(2)'} \boldsymbol{C}^{(2)} \boldsymbol{\varepsilon}_0^{(2)}}} \to_d \frac{\chi_k^2 - k}{\sqrt{2k}}.$$

Thus, $\sqrt{2k}T_1 + k \rightarrow_d \chi_k^2$. Let us consider now the T_2 statistic. Notice that $\tilde{\boldsymbol{\varepsilon}} = \boldsymbol{y} - \boldsymbol{X}_1 \boldsymbol{\beta}_{10} - \boldsymbol{X}_2 \boldsymbol{\beta}_2$ where $\boldsymbol{\beta}_2 = (\boldsymbol{X}_2' \boldsymbol{P} \boldsymbol{X}_2)^{-1} \boldsymbol{X}_2' \boldsymbol{P} (\boldsymbol{y} - \boldsymbol{X}_1 \boldsymbol{\beta}_{10})$, the two-stage least squares under the null. By standard manipulations, CLT and Slutsky's theorem we get

$$\widetilde{\boldsymbol{\varepsilon}}'\boldsymbol{P}\widetilde{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon}'\left(\boldsymbol{P} - \boldsymbol{P}\boldsymbol{X}_{2}(\boldsymbol{X}_{2}'\boldsymbol{P}\boldsymbol{X}_{2})^{-1}\boldsymbol{X}_{2}'\boldsymbol{P}\right)\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}'\boldsymbol{Z}^{*}\left(\boldsymbol{I}_{k} - \boldsymbol{P}_{\boldsymbol{X}_{2}'\boldsymbol{Z}^{*}}\right)\boldsymbol{Z}^{*'}\boldsymbol{\varepsilon} \rightarrow_{d} \sigma^{2}\chi_{k-g_{2}}^{2}$$

where $\mathbf{Z}^* = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1/2}$. Using the results in Lemma A.3 we get $\sum_{i=1}^n P_{ii} \tilde{\varepsilon}_i^2 \rightarrow_p k\sigma^2$ and

$$\widehat{V}(\widetilde{\boldsymbol{\beta}}) \to_p 2\sigma^4$$
. So, by the usual standard arguments $\sqrt{2kT_2} + k \to_d \chi^2_{k-g_2}$.

Proof of Corollary 1. The proof mimics that of Theorem 1 in Chao *et al.* (2014). Let $q_{\alpha}^{\chi_k^2}$ be the generic α -th quantile of the chi square distribution with k degrees of freedom. As $k \to \infty$, $\frac{q_{\alpha}^{\chi_k^2} - k}{\sqrt{2k}} \to z_{\alpha}$, where z_{α} is the generic α -th quantile of the standard normal distribution. This proves parts (*i*) and (*ii*). With respect to part (*iii*) and part (*iv*), notice that $\sqrt{\frac{k-g_2}{k}} \frac{q_{\alpha}^{\chi_k^2} - g_2}{\sqrt{2(k-g_2)}} - \frac{g_2}{\sqrt{2k}} \to z_{\alpha}$ as $k \to \infty$.

B Figures



Figure 1: PP-plots with heteroskedasticity, $H_0: \boldsymbol{\beta} = \boldsymbol{\beta}_0$.



Figure 2: PP-plots with heteroskedasticity, $H_0: \beta_1 = \beta_{10}$.



Figure 3: Power curves with heteroskedasticity and $\pi = 0.1$, $H_0: \beta = \beta_0$.



Figure 4: Power curves with heteroskedasticity and $\pi = 0.1$, $H_0: \beta_1 = \beta_{10}$.



Figure 5: Power curves with heteroskedasticity and $\pi = 1$, $H_0: \beta = \beta_0$.


Figure 6: Power curves with heteroskedasticity and $\pi = 1$, $H_0: \beta_1 = \beta_{10}$.



Figure 7: Histograms and QQ-plots for T_2 and JIV1 plug-in. The blue curve superimposed on the histograms is the standard normal distribution.



Figure 8: PP-plots for T_2 with a slow (inconsistent) plug-in.



Figure 9: T_2 , AR_{AG} and AR P-values for different values of β for the model with four instruments. 38



Figure 10: T_2 , AR_{AG} and AR P-values for different values of β for the model with fourteen instruments. 39

References

- Anatolyev, S. (2018) Many Instruments and/or Regressors: A Friendly Guide. Journal of Economic Surveys 33, 689–726.
- Anatolyev, S. and Gospodinov, N. (2011) Specification Testing in Models with Many Instruments. *Econometric Theory* 27, 427–441.
- Anatolyev, S. and Yaskov, P. (2017) Asymptotics of diagonal elements of projection matrices under many instruments/regressors. *Econometric Theory* 33, 717–738.
- Anderson, T.W. and Rubin, H. (1949) Estimators of the Parameters of a Single Equation in a Complete Set of Stochastic Equations. *The Annals of Mathematical Statistics* 21, 570–582.
- Andrews, D.W.K., Marmer, V. and Yu, Z. (2019) On optimal inference in the linear IV model. *Quantitative Economics* 10, 457–485.
- Andrews, D.W.K., Moreira, M.J. and Stock, J.H. (2006) Optimal Two-Sided Invariant Similar Tests for Instrumental Variable Regression. *Econometrica* 73, 715–752.
- Andrews, D.W.K. and Stock, J. (2005) Inference with Weak Instruments. In R. Blundell, W.K. Newey and T. Persson (eds.), Advances in Economics and Econometrics, Theory and Applications: Ninth World Congress of the Econometric Society, volume III, Cambridge University Press, Cambridge.
- Andrews, D.W.K. and Stock, J.H. (2007) Testing with Many Weak Instruments. The Journal of Econometrics 138, 24–46.
- Angrist, J.D. and Krueger, A. (1991) Does compulsory school attendance affect schooling and earnings? The Quarterly Journal of Economics 106, 979–1014.
- Angrist, J.D. and Pischke, J.S. (2008) Mostly Harmless Econometrics: An Empiricist's Companion. Princeton University Press.

- Bekker, P.A. (1994) Alternative approximations to the distributions of instrumental variable estimators. *Econometrica* 54, 657–682.
- Bekker, P.A. and Crudu, F. (2015) Jackknife Instrumental Variable Estimation with Heteroskedasticity. *The Journal of Econometrics* 185, 332–342.
- Bekker, P.A. and Van der Ploeg, J. (2005) Instrumental variable estimation based on grouped data. *Statistica Neerlandica* 59, 239–267.
- Bound, J., Jaeger, D.A. and Baker, R.M. (1995) Problems with Instrumental Variables Estimation When the Correlation between the Instruments and the Endogenous Explanatory Variable Is Weak. *Journal of the American Statistical Association* 90, 443–450.
- Bun, M., Farbmacher, H. and Poldermans, R. (2018) Finite sample properties of the Anderson and Rubin (1949) test. working paper .
- Card, D. (1995) Using Geographic Variation in College Proximity to Estimate the Return to Schooling. In L. Christofides, E. Grant and R. Swidinsky (eds.), Aspects of Labor Market Behaviour: Essays in Honour of John Vanderkamp, University of Toronto Press, Toronto, 201–222.
- Chao, J.C., Hausman, J.A., Newey, W.K., Swanson, N.R. and Woutersen, T. (2014) Testing Overidentifying Restrictions with Many Instruments and Heteroskedasticity. *The Journal* of Econometrics 178, 15–21.
- Chao, J.C. and Swanson, N.R. (2005) Consistent estimation with a large number of weak instruments. *Econometrica* 73, 1673–1692.
- Chao, J.C., Swanson, N.R., Hausman, J.A., Newey, W.K. and Woutersen, T. (2012) Asymptotic distribution of JIVE in a heteroskedastic IV regression with many instruments. *Econometric Theory* 28, 42–86.
- Davidson, R. and MacKinnon, J.G. (1998) Graphical Methods for Investigating the Size and Power of Hypothesis Tests. *The Manchester School* 66, 1–26.

- Donald, S.G., Imbens, G.W. and Newey, W.K. (2003) Empirical likelihood estimation and consistent tests with conditional moment restrictions. *Journal of Econometrics* 117, 55–93.
- Guggenberger, P., Kleibergen, F. and Mavroeidis, S. (2019) A more powerful subvector Anderson Rubin test in linear instrumental variables regression. *Quantitative Economics* 10, 487–526.
- Guggenberger, P., Kleibergen, F., Mavroeidis, S. and Chen, L. (2012) On the asymptotic sizes of subset Anderson–Rubin and Lagrange multiplier tests in linear instrumental variables regression. *Econometrica* 80, 2649–2666.
- Guggenberger, P. and Smith, R.J. (2005) Generalized Empirical Likelihood Estimators and Tests Under Partial, Weak, and Strong Identification. *Econometric Theory* 21, 667–709.
- Hausman, J.A., Newey, W.K., Woutersen, T., Chao, J.C. and Swanson, N.R. (2012) Instrumental variable estimation with heteroskedasticity and many instruments. *Quantitative Economics* 3, 211–255.
- Imbens, G.W. (2014) Instrumental Variables: An Econometrician's Perspective. Statistical Science 29, 323–358.
- Imbens, G.W. and Rubin, D. (1997) Estimating outcome distributions for compliers in instrumental variables models. *Review of Economic Studies* 64, 555–574.
- Jiang, Y., Kang, H. and Small, D. (2016) ivmodel: Statistical Inference and Sensitivity Analysis for Instrumental Variables Model. URL https://CRAN.R-project.org/ package=ivmodel, r package version 1.2.
- Kang, H., Zhang, A., Cai, T.T. and Small, D.S. (2016) Instrumental Variables Estimation With Some Invalid Instruments and its Application to Mendelian Randomization. *Journal of the American Statistical Association* 111, 132–144.

- Kleibergen, F. (2002) Pivotal Statistics for Testing Structural Parameters in Instrumental Variables Regression. *Econometrica* 70, 1781–1803.
- Kleibergen, F. (2004) Testing Subsets of Structural Parameters in the Instrumental Variables. The Review of Economics and Statistics 86, 418–423.
- Kleibergen, F. (2005) Testing Parameters in GMM Without Assuming They Are Identified. Econometrica 73, 1103–1123.
- Lee, Y. and Okui, R. (2012) Hahn–Hausman test as a specification test. *Journal of Econometrics* 167, 133–139.
- Moreira, M.J. (2003) A Conditional Likelihood Ratio Test for Structural Models. *Econo*metrica 71, 1027–1048.
- Moreira, M.J. (2009) Tests with Correct Size When Instruments Can Be Arbitrarily Weak. The Journal of Econometrics 152, 131–140.
- Newey, W.K. and Windmeijer, F. (2009) Generalized method of moments with many weak moment conditions. *Econometrica* 77, 687–719.
- Staiger, D. and Stock, J.H. (1997) Instrumental Variables Regression with Weak Instruments. *Econometrica* 65, 557–586.
- Stock, J.H., Wright, J.H. and Yogo, M. (2002) A Survey of Weak Instruments and Weak Identification in Generalized Method of Moments. *Journal of Business and Economic Statistics* 20, 518–529.
- Van Hasselt, M. (2010) Many instruments asymptotic approximations under nonnormal error distributions. *Econometric Theory* 26, 633–645.
- Von Hinke, S., Davey Smith, G., Lawlor, D.A., Propper, C. and Windmeijer, F. (2016) Genetic markers as instrumental variables. *Journal of Health Economics* 45, 131–148.

- Wang, J. and Zivot, E. (1998) Inference on a Structural Parameter in Instrumental Variables Regression with Weak Instruments. *Econometrica* 66, 1389–1404.
- Windmeijer, F., Farbmacher, H., Davies, N. and Davey Smith, G. (2017) On the Use of the Lasso for Instrumental Variables Estimation with Some Invalid Instruments. Bristol economics discussion papers, Department of Economics, University of Bristol, UK.
- Zivot, E., Startz, R. and Nelson, C.R. (1998) Valid Confidence Intervals and Inference in the Presence of Weak Instruments. *International Economic Review* 39, 1119–1144.

Supplement to "Inference in instrumental variables models with heteroskedasticity and many instruments"

Federico Crudu*

Giovanni Mellace[†]

Università di Siena and CRENoS

University of Southern Denmark

Zsolt Sándor[‡]

Sapientia Hungarian University of Transylvania

November 2019

Abstract

This supplement contains the proofs of the auxiliary lemmas, some additional theoretical results and further Monte Carlo experiments that complement the results in the main text.

^{*}Department of Economics and Statistics, Piazza San Francesco 7/8, 53100 Siena, Italy, federico.crudu@unisi.it

[†]Department of Business and Economics, Campusvej 55, 5230 Odense M, Denmark, giome@sam.sdu.dk

 $^{^{\}ddagger} \mathrm{Department}$ of Business Sciences, Piata Libertătii 1, 530104 Miercurea Ciuc, Romania, sandorz-solt@cs.sapientia.ro

A Proofs of Lemmas

This Section contains some auxiliary lemmas that are useful to prove the main results of the paper.

Lemma A.2. Let $\widehat{V}(\boldsymbol{\beta}_0) = \frac{2}{k} \boldsymbol{\varepsilon}_0^{(2)} \boldsymbol{C}^{(2)} \boldsymbol{\varepsilon}_0^{(2)}$. If Assumptions 1, 3 hold, $\widehat{V}(\boldsymbol{\beta}_0) - V_n = O_p\left(\frac{1}{\sqrt{k}}\right)$; consequently $\widehat{V}(\boldsymbol{\beta}_0) - V_n \xrightarrow{p} 0$, where

$$V_n = \frac{2}{k} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^2 \sigma_i^2 \sigma_j^2.$$

Proof. Let $\eta_i = \varepsilon_i^2 - \sigma_i^2$; then

$$\widehat{V}(\boldsymbol{\beta}_{0}) - V_{n} = \frac{2}{k} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} \left(\varepsilon_{i}^{2} \varepsilon_{j}^{2} - \sigma_{i}^{2} \sigma_{j}^{2} \right) = \frac{2}{k} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} \left(\eta_{i} \eta_{j} + \sigma_{i}^{2} \eta_{j} + \sigma_{j}^{2} \eta_{i} \right).$$

 So

$$\left| V_n - \widehat{V}(\boldsymbol{\beta}_0) \right| \le \frac{2}{k} \left| \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \eta_i \eta_j \right| + \frac{2}{k} \left| \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \sigma_i^2 \eta_j \right| + \frac{2}{k} \left| \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \sigma_j^2 \eta_i \right|$$
$$\equiv A_1 + A_2 + A_3.$$

Since

$$\mathbf{E}\left[\eta_{i}^{2}\right] = \mathbf{E}\left[\varepsilon_{i}^{4}\right] - \sigma_{i}^{4},$$

from Assumption 3 we have $\mathbf{E}\left[\eta_i^2\right] \leq c_u$. So

$$\mathbf{E}\left[A_{1}^{2}\right] = \frac{8}{k^{2}} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{4} \mathbf{E}\left[\eta_{i}^{2}\right] \mathbf{E}\left[\eta_{j}^{2}\right] \le \frac{c_{u}}{k^{2}} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{4}.$$

Note that for $i \neq j$ we have

$$C_{ij} = \frac{P_{ij}}{2} \left(\frac{1}{1 - P_{ii}} + \frac{1}{1 - P_{jj}} \right),$$

which from Assumption 1 implies

$$|C_{ij}| = \frac{|P_{ij}|}{2} \left(\frac{1}{1 - P_{ii}} + \frac{1}{1 - P_{jj}} \right) \le c_u |P_{ij}| \quad \text{for any } i, j,$$
(A.1)

 \mathbf{SO}

$$\mathbb{E}\left[A_{1}^{2}\right] \leq \frac{c_{u}}{k^{2}} \sum_{i=1}^{n} \sum_{j=1}^{n} P_{ij}^{4}.$$

From Assumption 1 and the fact that $\boldsymbol{P} = \boldsymbol{P}^2$, we have

$$P_{hh} \ge P_{hh}^2 = \left(\sum_{i=1}^n P_{hi}^2\right)^2 = \sum_{i=1}^n \sum_{j=1}^n P_{hi}^2 P_{hj}^2 \ge \sum_{i=1}^n P_{hi}^4,$$

 \mathbf{SO}

$$\sum_{i=1}^{n} \sum_{j=1}^{n} P_{ij}^{4} \le tr\left(\mathbf{P}\right) = k \text{ and } \sum_{i,j,h} P_{hi}^{2} P_{hj}^{2} \le k.$$
(A.2)

Therefore,

$$\operatorname{E}\left[A_1^2\right] \le \frac{c_u}{k}.$$

Now, by Cauchy-Schwarz $(E[\varepsilon_i^2])^2 \leq E[\varepsilon_i^4]$, thus $\sigma_i^2 \leq c_u$, so from Assumption 3, (A.1) and (A.2)

$$\mathbb{E}\left[A_{2}^{2}\right] = \frac{4}{k^{2}} \sum_{i,j,h} C_{hi}^{2} C_{ij}^{2} \sigma_{h}^{2} \sigma_{j}^{2} \mathbb{E}\left[\eta_{i}^{2}\right] \le \frac{4c_{u}^{2}}{k^{2}} \sum_{i,j,h} C_{hi}^{2} C_{ij}^{2} \le \frac{c_{u}}{k^{2}} \sum_{i,j,h} P_{hi}^{2} P_{ij}^{2} \le \frac{c_{u}}{k}.$$

We can obtain a similar inequality for A_3 , so by the Markov and triangle inequalities we obtain that $\hat{V}(\boldsymbol{\beta}_0) - V_n = O_p\left(\frac{1}{\sqrt{k}}\right)$, therefore, $\hat{V}(\boldsymbol{\beta}_0) - V_n \stackrel{p}{\to} 0$.

Lemma A.3. Let $\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) = \frac{2}{k}\widetilde{\boldsymbol{\epsilon}}^{(2)\prime}\boldsymbol{C}^{(2)}\widetilde{\boldsymbol{\epsilon}}^{(2)}$. If $\widetilde{\boldsymbol{\beta}} \xrightarrow{p} \boldsymbol{\beta}$ and Assumptions 1, 3 hold, then $\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) - V_n \xrightarrow{p} 0$.

Proof. Let

$$\widetilde{V}_n = \frac{2}{k} \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \varepsilon_i^2 \varepsilon_j^2.$$

Then

$$\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) - \widetilde{V}_n = \frac{2}{k} \sum_{i=1}^n \sum_{j=1}^n C_{ij}^2 \left(\widetilde{\varepsilon}_i^2 \widetilde{\varepsilon}_j^2 - \varepsilon_i^2 \varepsilon_j^2\right).$$

Note that

$$\begin{split} \left| \widetilde{\varepsilon}_{i}^{2} \widetilde{\varepsilon}_{j}^{2} - \varepsilon_{i}^{2} \varepsilon_{j}^{2} \right| &\leq \widetilde{\varepsilon}_{j}^{2} \left| \widetilde{\varepsilon}_{i}^{2} - \varepsilon_{i}^{2} \right| + \varepsilon_{i}^{2} \left| \widetilde{\varepsilon}_{j}^{2} - \varepsilon_{j}^{2} \right|, \\ \left| \widetilde{\varepsilon}_{i} + \varepsilon_{i} \right| &\leq \left| \widetilde{\varepsilon}_{i} - \varepsilon_{i} \right| + 2 \left| \varepsilon_{i} \right|, \\ \widetilde{\varepsilon}_{j}^{2} &= \left| \widetilde{\varepsilon}_{j}^{2} - \varepsilon_{j}^{2} + \varepsilon_{j}^{2} \right| \leq \left| \widetilde{\varepsilon}_{j}^{2} - \varepsilon_{j}^{2} \right| + \varepsilon_{j}^{2} \end{split}$$

and

$$\begin{aligned} \left| \widetilde{\varepsilon}_{i}^{2} - \varepsilon_{i}^{2} \right| &= \left| \widetilde{\varepsilon}_{i} - \varepsilon_{i} \right| \cdot \left| \widetilde{\varepsilon}_{i} + \varepsilon_{i} \right| = \left| \mathbf{X}_{i}^{\prime} \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right| \cdot \left| \widetilde{\varepsilon}_{i} + \varepsilon_{i} \right| \leq \left| \mathbf{X}_{i}^{\prime} \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right| \cdot \left(\left| \mathbf{X}_{i}^{\prime} \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right| + 2 \left| \varepsilon_{i} \right| \right) \\ &\leq \left\| \mathbf{X}_{i} \right\| \left\| \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right\| \cdot \left(\left\| \mathbf{X}_{i} \right\| \left\| \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right\| + 2 \left| \varepsilon_{i} \right| \right) \equiv d_{i} \left\| \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right\|, \end{aligned}$$

where

$$d_{i} = \left\| \boldsymbol{X}_{i} \right\| \left(\left\| \boldsymbol{X}_{i} \right\| + 2 \left| \varepsilon_{i} \right| \right)$$

because $\left\| \left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}} \right) \right\| \leq 1$ with probability approaching 1 as $n \to \infty$. So

$$\left|\widetilde{\varepsilon}_{i}^{2}\widetilde{\varepsilon}_{j}^{2}-\varepsilon_{i}^{2}\varepsilon_{j}^{2}\right|\leq d_{i}d_{j}\left\|\left(\boldsymbol{\beta}-\widetilde{\boldsymbol{\beta}}\right)\right\|^{2}+\left(d_{i}\varepsilon_{j}^{2}+d_{j}\varepsilon_{i}^{2}\right)\left\|\left(\boldsymbol{\beta}-\widetilde{\boldsymbol{\beta}}\right)\right\|,$$

therefore,

$$\left|\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) - \widetilde{V}_{n}\right| \leq \frac{2\left\|\left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}}\right)\right\|^{2}}{k} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} d_{i} d_{j} + \frac{2\left\|\left(\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}}\right)\right\|}{k} \sum_{i=1}^{n} \sum_{j=1}^{n} C_{ij}^{2} d_{i} \varepsilon_{j}^{2}.$$
(A.3)

Note that by the Cauchy-Schwarz inequality

$$\operatorname{E}\left[d_{i}d_{j}\right] \leq \sqrt{\operatorname{E}\left[d_{i}^{2}\right]\operatorname{E}\left[d_{j}^{2}\right]},$$

where

Assumption 3 and Minkowski's inequality imply $\mathbf{E}\left[\|\mathbf{X}_i\|^4\right] \leq c_u$. Hence

$$\operatorname{E}\left[d_{i}d_{j}\right] \leq c_{u}, \quad \operatorname{E}\left[d_{i}^{2}\right] \leq c_{u},$$

so by Assumption 1 and (A.1)

$$\mathbb{E}\left[\frac{1}{k}\sum_{i=1}^{n}\sum_{j=1}^{n}C_{ij}^{2}d_{i}d_{j}\right] \leq \frac{1}{k}\sum_{i=1}^{n}\sum_{j=1}^{n}C_{ij}^{2}\mathbb{E}\left[d_{i}d_{j}\right] \leq c_{u}\left(\frac{1}{k}\sum_{i=1}^{n}\sum_{j=1}^{n}C_{ij}^{2}\right) \leq c_{u}$$

and

$$\operatorname{E}\left[\frac{1}{k}\sum_{i=1}^{n}\sum_{j=1}^{n}C_{ij}^{2}d_{i}^{2}\right] \leq c_{u}.$$

Then by Markov's and the triangle inequalities $\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right) - V_n \xrightarrow{p} 0.$

Lemma A.4. Under Assumptions 1, 3

1. $\operatorname{E} [\mathbf{X}_{2}^{\prime} \mathbf{C} \mathbf{X}_{2}] = \mathbf{H}_{22}, \operatorname{Var} [\mathbf{X}_{2}^{\prime} \mathbf{C} \mathbf{X}_{2}] \leq c_{u} \mathbf{H}_{22} + c_{u} k \mathbf{I}_{g_{2}} + c_{u} r_{\max} \mathbf{I}_{g_{2}},$ 2. $\operatorname{E} [\mathbf{X}_{2}^{\prime} \mathbf{C} \boldsymbol{\varepsilon}] = \mathbf{0}, \operatorname{Var} [\mathbf{X}_{2}^{\prime} \mathbf{C} \boldsymbol{\varepsilon}] \leq c_{u} \mathbf{H}_{22} + c_{u} k \mathbf{I}_{g_{2}}.$ *Proof.* The model $X_2 = Z\Pi_2 + U_2$ implies that

$$X_{2}'CX_{2} = (Z\Pi_{2})'CZ\Pi_{2} + (Z\Pi_{2})'CU_{2} + U_{2}'CZ\Pi_{2} + U_{2}'CU_{2}.$$
 (A.4)

Therefore,

$$\mathbb{E}\left[oldsymbol{X}_2'oldsymbol{C}oldsymbol{X}_2
ight] = oldsymbol{\Pi}_2'oldsymbol{Z}'oldsymbol{C}oldsymbol{Z}oldsymbol{\Pi}_2 + oldsymbol{E}\left[oldsymbol{U}_2'oldsymbol{C}oldsymbol{U}_2
ight].$$

Since Z'CZ = Z'Z we have that

$$\boldsymbol{\Pi}_{2}^{\prime}\boldsymbol{Z}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}=\boldsymbol{H}_{22}.\tag{A.5}$$

Also,

$$\mathbb{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] = \sum_{i=1}^{n}\sum_{j=1}^{n}E\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{e}_{i}\boldsymbol{e}_{i}^{\prime}\boldsymbol{C}\boldsymbol{e}_{j}\boldsymbol{e}_{j}^{\prime}\boldsymbol{U}_{2}\right] = \sum_{i=1}^{n}\sum_{j=1}^{n}C_{ij}E\left[\boldsymbol{U}_{2i}\boldsymbol{U}_{2j}^{\prime}\right] = \boldsymbol{O}$$

because the main diagonal elements of C are 0, so $E[X'_2CX_2] = H_{22}$. Further,

$$\begin{aligned} \operatorname{Var}\left[\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\right] &= \operatorname{E}\left[\left\{\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2} + \boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2} + \boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right\}\left\{\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2} + \boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2} + \boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right\}^{\prime}\right] \\ &\leq 3\operatorname{E}\left[\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right] + 3\operatorname{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] + 3\operatorname{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] \end{aligned}$$

by the Cauchy-Schwarz inequality. By Assumption 3 $\operatorname{E}[U_2U'_2] \leq c_u I_n$ and from the definition of C it holds that (see Bekker and Crudu (2015), p.337)

$$Z'C^2Z = Z'\left\{I_n + \frac{1}{4}\left(I_n - D\right)^{-1}\left(I_n - P\right)\left(I_n - D\right)^{-1}\right\}Z.$$

Further, by Assumption 1 $(\boldsymbol{I}_n - \boldsymbol{D})^{-1} \leq c_u \boldsymbol{I}_n$, and therefore,

$$\mathbf{Z}'\mathbf{C}^2\mathbf{Z} \le c_u \mathbf{Z}'\mathbf{Z},\tag{A.6}$$

so the first expectation is

$$\mathbb{E}\left[\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right] \leq c_{u}\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}^{2}\boldsymbol{Z}\boldsymbol{\Pi}_{2} \leq c_{u}\boldsymbol{H}_{22}.$$
(A.7)

The second expectation is

$$\mathbb{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right]=\sum_{i=1}^{n}a_{ii}\mathbb{E}\left[\boldsymbol{U}_{2i}\boldsymbol{U}_{2i}^{\prime}\right],$$

where a_{ii} denotes the *i*-th main diagonal component of $CZ\Pi_2(Z\Pi_2)'C$ and U'_{2i} is the *i*-th row of U_2 . By Assumption 3 and (A.6) we obtain that

$$\mathbb{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] \leq c_{u}\operatorname{tr}\left(\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}\right)\boldsymbol{I}_{g_{2}} = c_{u}\operatorname{tr}\left(\left(\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)^{\prime}\boldsymbol{C}^{2}\boldsymbol{Z}\boldsymbol{\Pi}_{2}\right)\boldsymbol{I}_{g_{2}}$$

$$\leq c_{u}\operatorname{tr}\left(\boldsymbol{H}_{22}\right)\boldsymbol{I}_{g_{2}} \leq c_{u}r_{\max}\boldsymbol{I}_{g_{2}}.$$

$$(A.8)$$

$$\begin{split} \operatorname{E} \left[\boldsymbol{U}_{2}^{\prime} \boldsymbol{C} \boldsymbol{U}_{2} \boldsymbol{U}_{2}^{\prime} \boldsymbol{C} \boldsymbol{U}_{2} \right] &= \sum_{i,j,k,\ell} \operatorname{E} \left[\boldsymbol{U}_{2}^{\prime} \boldsymbol{e}_{i} \boldsymbol{e}_{i}^{\prime} \boldsymbol{C} \boldsymbol{e}_{j} \boldsymbol{e}_{j}^{\prime} \boldsymbol{U}_{2} \boldsymbol{U}_{2}^{\prime} \boldsymbol{e}_{k} \boldsymbol{e}_{k}^{\prime} \boldsymbol{C} \boldsymbol{e}_{\ell} \boldsymbol{e}_{\ell}^{\prime} \boldsymbol{U}_{2} \right] \\ &= \sum_{i,j,k,\ell} C_{ij} C_{k\ell} \operatorname{E} \left[\boldsymbol{U}_{2i} \boldsymbol{U}_{2j}^{\prime} \boldsymbol{U}_{2k} \boldsymbol{U}_{2\ell}^{\prime} \right] \\ &= \sum_{i \neq j} C_{ij}^{2} \operatorname{E} \left[\boldsymbol{U}_{2i} \boldsymbol{U}_{2j}^{\prime} \boldsymbol{U}_{2i} \boldsymbol{U}_{2j}^{\prime} \right] + \sum_{i \neq j} C_{ij}^{2} \operatorname{E} \left[\boldsymbol{U}_{2i} \boldsymbol{U}_{2j}^{\prime} \boldsymbol{U}_{2j} \boldsymbol{U}_{2i} \right]. \end{split}$$

By Assumption 3, the Cauchy-Schwarz inequality and (21) we obtain that

$$\mathbb{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] \leq c_{u}\left(\sum_{i\neq j}C_{ij}^{2}\right)\boldsymbol{I}_{g_{2}} \leq c_{u}k\boldsymbol{I}_{g_{2}}.$$
(A.9)

By collecting the results from (A.7), (A.8), (A.9) we obtain the first result.

2. The model $X_2 = Z\Pi_2 + U_2$ implies $X'_2 C \varepsilon = \Pi'_2 Z' C \varepsilon + U'_2 C \varepsilon$. Similar to part 1., since the main diagonal elements of C are 0, we have $E[X'_2 C \varepsilon] = 0$. Regarding the

variance we have

$$\operatorname{Var}\left[\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] = (\boldsymbol{Z}\boldsymbol{\Pi}_{2})^{\prime}\boldsymbol{C}\operatorname{E}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\right]\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2} + \operatorname{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\right]\boldsymbol{C}\boldsymbol{Z}\boldsymbol{\Pi}_{2} + (\boldsymbol{Z}\boldsymbol{\Pi}_{2})^{\prime}\boldsymbol{C}\operatorname{E}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] \\ + \operatorname{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right].$$
(A.10)

By Assumption 3 and (A.6) the first term is

$$\boldsymbol{H}_{22}^{-1/2} \left(\boldsymbol{Z} \boldsymbol{\Pi}_{2} \right)' \boldsymbol{C} \operatorname{E} \left[\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}' \right] \boldsymbol{C} \boldsymbol{Z} \boldsymbol{\Pi}_{2} \boldsymbol{H}_{22}^{-1/2} \leq c_{u} \boldsymbol{H}_{22}^{-1/2} \left(\boldsymbol{Z} \boldsymbol{\Pi}_{2} \right)' \boldsymbol{Z} \boldsymbol{\Pi}_{2} \boldsymbol{H}_{22}^{-1/2} = c_{u} \boldsymbol{I}_{g_{2}}.$$
(A.11)

The second and third terms from (A.10) are 0. This is because

$$E\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\right]\boldsymbol{C}=\sum_{i,j,k}\mathrm{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{e}_{i}\boldsymbol{e}_{i}^{\prime}\boldsymbol{C}\boldsymbol{e}_{j}\boldsymbol{\varepsilon}_{j}^{\prime}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\boldsymbol{e}_{k}\right]\boldsymbol{e}_{k}^{\prime}\boldsymbol{C}=\sum_{i,j,k}\mathrm{E}\left[\boldsymbol{U}_{2i}C_{ij}\varepsilon_{j}\varepsilon_{k}\right]\boldsymbol{e}_{k}^{\prime}\boldsymbol{C}.$$

Since the main diagonal elements of C are 0, the expectations from the above sum are 0. Consequently, $C \to [\varepsilon \varepsilon' C U_2] = O$ as well. The fourth term from the expression in (A.10) is

$$E\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] = \sum_{i,j,k,\ell} E\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{e}_{i}\boldsymbol{e}_{i}^{\prime}\boldsymbol{C}\boldsymbol{e}_{j}\boldsymbol{e}_{j}^{\prime}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\boldsymbol{e}_{k}\boldsymbol{e}_{k}^{\prime}\boldsymbol{C}\boldsymbol{e}_{\ell}\boldsymbol{e}_{\ell}^{\prime}\boldsymbol{U}_{2}\right] = \sum_{i,j,k,\ell} E\left[\boldsymbol{U}_{2i}C_{ij}\varepsilon_{j}\varepsilon_{k}C_{k\ell}\boldsymbol{U}_{2\ell}^{\prime}\right]$$
$$= \sum_{i\neq j} C_{ij}^{2}\left(E\left[\varepsilon_{j}^{2}\boldsymbol{U}_{2i}\boldsymbol{U}_{2i}^{\prime}\right] + E\left[\varepsilon_{i}\boldsymbol{U}_{2i}\varepsilon_{j}\boldsymbol{U}_{2j}^{\prime}\right]\right) = \sum_{i\neq j} C_{ij}^{2}\left(\sigma_{j}^{2}\boldsymbol{\Sigma}_{22i} + \boldsymbol{\sigma}_{12i}\boldsymbol{\sigma}_{12j}^{\prime}\right)$$

By the Cauchy-Schwarz inequality, Assumption 3, and Equation (21) in the main text we obtain that

$$\mathbb{E}\left[\boldsymbol{U}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^{\prime}\boldsymbol{C}\boldsymbol{U}_{2}\right] \leq c_{u}\sum_{i\neq j}C_{ij}^{2}\boldsymbol{I}_{g_{2}} \leq c_{u}k\boldsymbol{I}_{g_{2}}.$$
(A.12)

Collecting the results from (A.11) and (A.12), we obtain the result. \Box

B Auxiliary Results

This Section includes a set of theorems, examples, remarks and propositions associated to the main results of the paper. First we discuss the case when the OLS estimator $\tilde{\beta}_2 = (X'_2 X_2)^{-1} X'_2 y_0$ can be used as a plug-in estimator. This occurs in the practically relevant situation where the null hypothesis contains all parameters corresponding to endogenous variables. The case when there is a single endogenous regressor in the model and the null hypothesis contains exactly its coefficient is a common example.

Theorem B.1. Suppose that \mathbf{X}_2 is exogenous, $\mathbf{H}_{22}/n = O(1)$ and $(\mathbf{X}'_2\mathbf{X}_2/n)^{-1} = O_p(1)$. Then under Assumptions 1, 2, 3 we have that $T_2 \rightarrow_d \mathcal{N}(0, 1)$.

Proof. We need to show that $\frac{\Delta}{\sqrt{k}} = o_p(1)$ where Δ is defined in Equation (24). First note that $\operatorname{E}\left[\frac{\mathbf{X}_2'\varepsilon}{\sqrt{n}}\right] = \mathbf{0}$ because \mathbf{X}_2 is exogenous, and

$$\operatorname{Var}\left[\frac{\boldsymbol{X}_{2}^{\prime}\boldsymbol{\varepsilon}}{\sqrt{n}}\right] \leq c_{u}\frac{\boldsymbol{H}_{22}}{n} + \frac{1}{n}\sum_{i=1}^{n}\sigma_{i}^{2}\boldsymbol{\Sigma}_{i22} = O\left(1\right)$$

due to Assumption 3 and $H_{22}/n = O(1)$. Consequently,

$$\sqrt{n}\left(\widetilde{\boldsymbol{\beta}}_{2}-\boldsymbol{\beta}_{2}\right)=\left(\boldsymbol{X}_{2}^{\prime}\boldsymbol{X}_{2}/n\right)^{-1}\left(\boldsymbol{X}_{2}^{\prime}\boldsymbol{\varepsilon}/\sqrt{n}\right)^{-1}=O_{p}\left(1\right).$$
(B.1)

Using $H_{22}/n = O(1)$ and by Lemma A.4 we have

$$\mathbf{E}\left[\frac{\mathbf{X}_{2}^{\prime}\mathbf{C}\mathbf{X}_{2}}{n\sqrt{k}}\right] = o\left(1\right), \quad \operatorname{Var}\left[\frac{\mathbf{X}_{2}^{\prime}\mathbf{C}\mathbf{X}_{2}}{n\sqrt{k}}\right] = o\left(1\right)$$

and

$$\mathbb{E}\left[\frac{\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}}{\sqrt{kn}}\right] = \boldsymbol{0}, \quad \operatorname{Var}\left[\frac{\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}}{\sqrt{kn}}\right] = o\left(1\right).$$

Therefore,

$$\frac{1}{\sqrt{k}}\Delta = \frac{1}{\sqrt{k}} \left(\widetilde{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \right)' \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \left(\widetilde{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \right) - \frac{2}{\sqrt{k}} \left(\widetilde{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_2 \right)' \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{\varepsilon} = o_p \left(1 \right),$$

so B_2 in (23) is $o_p(1)$. B_1 from (23) is $o_p(1)$ due to the consistency of $\tilde{\beta}_2$ from (B.1) and the result follows.

The following theorem provides sufficient conditions for the asymptotic distribution of the T_2 statistic when the JIV1 estimator is used as plug-in. Recall that $r_{\min} = \lambda_{\min}(\mathbf{H}_{22})$ and $r_{\max} = \lambda_{\max}(\mathbf{H}_{22})$.

Theorem B.2. If Assumptions 1, 2, 3 and $\sqrt{k}/r_{\min} \to 0$, $r_{\max}/k = O(1)$ are satisfied, then the JIV1 estimator $\widetilde{\boldsymbol{\beta}}_2 = (\boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2)^{-1} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{y}_0$ is consistent and $T_2 \to_d \mathcal{N}(0, 1)$.

Proof. First we show consistency, that is, $\tilde{\beta}_2 - \beta_2 = (X'_2 C X_2)^{-1} X'_2 C \varepsilon = o_p(1)$. From Lemma A.4 it follows that

$$\mathbb{E}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\right] = \boldsymbol{I}_{g_{2}}, \text{ Var}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\right] \leq c_{u}\boldsymbol{H}_{22}^{-1} + c_{u}k\boldsymbol{H}_{22}^{-2} + c_{u}r_{\max}\boldsymbol{H}_{22}^{-2}$$

Since $\boldsymbol{H}_{22}^{-1} \leq \frac{1}{r_{\min}} \boldsymbol{I}_{g_2}$ and by assumptions $\frac{k}{r_{\min}^2} \to 0$, $r_{\max}/k = O(1)$ we have that $\operatorname{Var}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{X}_2\right] \to \boldsymbol{O}$, so $\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{X}_2 \to_p \boldsymbol{I}_{g_2}$, and therefore, $\left(\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{X}_2\right)^{-1} = O_p(1)$. From Lemma A.4 it also follows that

$$\mathbb{E}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] = \boldsymbol{0}, \text{ Var}\left[\boldsymbol{H}_{22}^{-1}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] \leq c_{u}\boldsymbol{H}_{22}^{-1} + c_{u}k\boldsymbol{H}_{22}^{-2}.$$

This variance goes to O for the same reason as above, so $H_{22}^{-1}X_2'C\varepsilon = o_p(1)$. Therefore, $\widetilde{\beta}_2 - \beta_2 = o_p(1)$.

Let now

$$\widehat{V}(\boldsymbol{\beta}) = \frac{2}{k} \boldsymbol{\varepsilon}^{(2)} \boldsymbol{C}^{(2)} \boldsymbol{\varepsilon}^{(2)}, \quad \boldsymbol{\varepsilon} = \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}.$$

Note that Δ in (23) now is

$$\Delta = -\boldsymbol{\varepsilon}' \boldsymbol{C} \boldsymbol{X}_2 \left(\boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{X}_2 \right)^{-1} \boldsymbol{X}_2' \boldsymbol{C} \boldsymbol{\varepsilon}.$$

This can be written as

$$\Delta = -\varepsilon' C X_2 H_{22}^{-1/2} \left(H_{22}^{-1/2} X_2' C X_2 H_{22}^{-1/2} \right)^{-1} H_{22}^{-1/2} X_2' C \varepsilon.$$
(B.2)

From Lemma A.4 we know that

$$\mathbb{E}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\boldsymbol{H}_{22}^{-1/2}\right] = \boldsymbol{I}_{g_{2}},$$

$$\text{Var}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\boldsymbol{H}_{22}^{-1/2}\right] \leq c_{u}\boldsymbol{H}_{22}^{-1} + c_{u}k\boldsymbol{H}_{22}^{-2} + c_{u}r_{\max}\boldsymbol{H}_{22}^{-2}.$$

$$\boldsymbol{H}^{-1} \leq \frac{1}{2}\boldsymbol{I} \quad \text{and by assumptions} \quad \frac{k}{2} \rightarrow 0 \quad r \quad /k = O(1) \text{ we have that}$$

Since $\boldsymbol{H}_{22}^{-1} \leq \frac{1}{r_{\min}} I_{g_2}$ and by assumptions $\frac{\kappa}{r_{\min}^2} \to 0$, $r_{\max}/k = O(1)$ we have that

$$\operatorname{Var}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{X}_{2}\boldsymbol{H}_{22}^{-1/2}\right]\to\boldsymbol{O},$$

so $H_{22}^{-1/2} X_2' C X_2 H_{22}^{-1/2} \rightarrow_p I_{g_2}$. Consequently, $\left(H_{22}^{-1/2} X_2' C X_2 H_{22}^{-1/2} \right)^{-1} = O_p(1)$. By Lemma A.4 we know that

$$\mathbf{E}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] = \boldsymbol{0}.$$
(B.3)

Next we show that under $\frac{\sqrt{k}}{r_{\min}} \to 0$ it holds that

$$\frac{1}{\sqrt{k}}\operatorname{Var}\left(\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right)=o\left(1\right).$$
(B.4)

From Lemma A.4 we know that

$$\frac{1}{\sqrt{k}}\operatorname{Var}\left[\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}\right] \leq \frac{c_{u}}{\sqrt{k}}\boldsymbol{I}_{g_{2}} + c_{u}\sqrt{k}\boldsymbol{H}_{22}^{-1} \leq \left(\frac{c_{u}}{\sqrt{k}} + c_{u}\frac{\sqrt{k}}{r_{\min}}\right)\boldsymbol{I}_{g_{2}},$$

where the last inequality follows from $H_{22}^{-1} \leq \frac{1}{r_{\min}} I_{g_2}$. So (B.4) holds under $k \to \infty$, $\frac{\sqrt{k}}{r_{\min}} \to 0$, and therefore, taking also (B.3) into account we obtain

$$rac{1}{k^{1/4}}oldsymbol{H}_{22}^{-1/2}oldsymbol{X}_{2}^{\prime}oldsymbol{C}oldsymbol{arepsilon}=o_{p}\left(1
ight).$$

Consequently,

$$B_{2} = \frac{1}{\sqrt{k}} \frac{\Delta}{\sqrt{\widehat{V}\left(\widetilde{\boldsymbol{\beta}}\right)}} = \left(\frac{1}{k^{1/4}} \boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_{2}^{\prime} \boldsymbol{C} \boldsymbol{\varepsilon}\right)^{\prime} \frac{1}{k^{1/4}} \boldsymbol{H}_{22}^{-1/2} \boldsymbol{X}_{2}^{\prime} \boldsymbol{C} \boldsymbol{\varepsilon} \cdot O_{p}\left(1\right) = o_{p}\left(1\right).$$

This result is not very different from Theorem 2 in the main text, but it is useful because, on the one hand, the convenient expression of the JIV1 estimator allows us to explain why underrejection of the null hypothesis occurs in most cases.¹ On the other hand, this result allows us to better link the weak instrument assumption $\sqrt{k}/r_{\min} \rightarrow 0$ to the asymptotic distribution of T_2 . Specifically, the proof of this result suggests that the assumption $\sqrt{k}/r_{\min} \rightarrow 0$ appears to be necessary for the asymptotic standard normality of the statistic T_2 .²

Derivation of Example 1. Suppose that there are ℓ groups with group g having n_g observations and

$$oldsymbol{Z} = \left(egin{array}{cccc} oldsymbol{\iota}_{n_1} & \cdots & oldsymbol{0} \ dots & \ddots & dots \ oldsymbol{0} & \cdots & oldsymbol{\iota}_{n_\ell} \end{array}
ight),$$

 $^{^1\}mathrm{We}$ discuss this in more detail in Remark B.1 below. $^2\mathrm{See}$ also Remark B.2.

where $\boldsymbol{\iota}_m$ is an $m\times 1$ vector of ones. In this case

The expression

$$E_2 = \frac{1}{\sqrt{k}} \sum_{i=1}^n \left(P_{ii} - \frac{k}{n} \right) \varepsilon_i^2$$

from (17) can be written as

$$E_2 = \frac{1}{\sqrt{\ell}} \sum_{g=1}^{\ell} \sum_{i \in G_g} \left(\frac{1}{n_g} - \frac{\ell}{n} \right) \varepsilon_i^2,$$

where G_g is the set of observations belonging to group g.

Suppose now that the groups have either 2 or 3 observations. In this case

$$E_2 = \frac{1}{\sqrt{\ell}} \sum_{g:n_g=2} \sum_{i \in G_g} \left(\frac{1}{2} - \frac{\ell}{n}\right) \varepsilon_i^2 + \frac{1}{\sqrt{\ell}} \sum_{g:n_g=3} \sum_{i \in G_g} \left(\frac{1}{3} - \frac{\ell}{n}\right) \varepsilon_i^2$$
$$= \left(\frac{1}{2} - \frac{\ell}{n}\right) \frac{1}{\sqrt{\ell}} \sum_{g:n_g=2} \sum_{i \in G_g} \varepsilon_i^2 + \left(\frac{1}{3} - \frac{\ell}{n}\right) \frac{1}{\sqrt{\ell}} \sum_{g:n_g=3} \sum_{i \in G_g} \varepsilon_i^2.$$

Suppose homoskedasticity with $E[\varepsilon_i^2] = \sigma^2$ and let ℓ_2 and ℓ_3 denote the number of 2-observation and 3-observation groups, respectively. In this case

$$E_2 = \left(\frac{1}{2} - \frac{\ell}{n}\right) \frac{2\ell_2}{\sqrt{\ell}} \frac{\sum_{g:n_g=2} \sum_{i \in G_g} \varepsilon_i^2}{2\ell_2} + \left(\frac{1}{3} - \frac{\ell}{n}\right) \frac{3\ell_3}{\sqrt{\ell}} \frac{\sum_{g:n_g=3} \sum_{i \in G_g} \varepsilon_i^2}{3\ell_3}$$

Note that $\ell = \ell_2 + \ell_3$ and $n = 2\ell_2 + 3\ell_3$, so

$$E_{2} = \frac{\ell_{3}}{2\ell_{2} + 3\ell_{3}} \frac{\ell_{2}}{\sqrt{\ell}} \frac{\sum_{g:n_{g}=2} \sum_{i \in G_{g}} \varepsilon_{i}^{2}}{2\ell_{2}} - \frac{\ell_{2}}{2\ell_{2} + 3\ell_{3}} \frac{\ell_{3}}{\sqrt{\ell}} \frac{\sum_{g:n_{g}=3} \sum_{i \in G_{g}} \varepsilon_{i}^{2}}{3\ell_{3}} \\ = \frac{\ell_{2}\ell_{3}}{\ell_{n}} \sqrt{\ell} \left(\frac{\sum_{g:n_{g}=2} \sum_{i \in G_{g}} \varepsilon_{i}^{2}}{2\ell_{2}} - \sigma^{2} - \left[\frac{\sum_{g:n_{g}=3} \sum_{i \in G_{g}} \varepsilon_{i}^{2}}{3\ell_{3}} - \sigma^{2} \right] \right).$$
(B.5)

By the CLT for IID observations

$$\sqrt{2\ell_2} \left(\frac{\sum_{g:n_g=2} \sum_{i \in G_g} \varepsilon_i^2}{2\ell_2} - \sigma^2 \right) \stackrel{d}{\to} \mathcal{N}(0, v) \quad \text{and}$$
$$\sqrt{3\ell_3} \left(\frac{\sum_{g:n_g=3} \sum_{i \in G_g} \varepsilon_i^2}{3\ell_3} - \sigma^2 \right) \stackrel{d}{\to} \mathcal{N}(0, v) ,$$

where $v = \operatorname{Var}[\varepsilon_i^2]$. The limit $\frac{\ell}{n} \to \lambda \in (0, 1)$ implies that $\frac{\ell}{2\ell_2} \to \frac{\lambda}{6\lambda - 2}$ and $\frac{\ell}{3\ell_3} \to \frac{\lambda}{3 - 6\lambda}$, so we obtain

$$\sqrt{\ell} \left(\frac{\sum_{g:n_g=2} \sum_{i \in G_g} \varepsilon_i^2}{2\ell_2} - \sigma^2 \right) \xrightarrow{d} \mathcal{N} \left(0, \frac{\lambda}{6\lambda - 2} v \right) \quad \text{and}$$
$$\sqrt{\ell} \left(\frac{\sum_{g:n_g=3} \sum_{i \in G_g} \varepsilon_i^2}{3\ell_3} - \sigma^2 \right) \xrightarrow{d} \mathcal{N} \left(0, \frac{\lambda}{3 - 6\lambda} v \right).$$

Therefore, from (B.5) we obtain

$$E_2 \xrightarrow{d} \mathcal{N}\left(0, \frac{(3\lambda - 1)(1 - 2\lambda)}{6\lambda}v\right).$$

Since its variance does not vanish in the limit, E_2 will not converge to 0 in probability. \Box

Example B.1. In this example we consider the indicator instruments discussed in Example 1 under heteroskedasticity when there are only groups of size 2 and 3, and we study whether E_2 defined in B.5 has convergent or divergent mean. That is, we study the limit of

$$\mathbf{E}\left[E_{2}\right] = \frac{1}{\sqrt{\ell}} \frac{\ell_{2}\ell_{3}}{n} \left(\frac{\sum_{g:n_{g}=2} \sum_{i \in G_{g}} \sigma_{i}^{2}}{2\ell_{2}} - \frac{\sum_{g:n_{g}=3} \sum_{i \in G_{g}} \sigma_{i}^{2}}{3\ell_{3}}\right)$$

with respect to the growth rate of ℓ_2 and ℓ_3 . First note that the assumption $\frac{1}{k} \sum_i \left(P_{ii} - \frac{k}{n}\right)^2 \rightarrow 0$ from Proposition 1 is equivalent to

$$\frac{1}{\ell} \left(\sum_{g:n_g=2} \sum_{i \in G_g} \left(\frac{1}{2} - \frac{\ell}{n} \right)^2 + \sum_{g:n_g=3} \sum_{i \in G_g} \left(\frac{1}{3} - \frac{\ell}{n} \right)^2 \right) \to 0.$$

Further, since

$$\sum_{g:n_g=2} \sum_{i \in G_g} \left(\frac{1}{2} - \frac{\ell}{n}\right)^2 + \sum_{g:n_g=3} \sum_{i \in G_g} \left(\frac{1}{3} - \frac{\ell}{n}\right)^2 = \frac{\ell_2 \ell_3}{6n},$$

this is equivalent to

$$\frac{\ell_2\ell_3}{\ell n} \to 0$$

Recalling that $\ell = \ell_2 + \ell_3$ and $n = 2\ell_2 + 3\ell_3$, we conclude that this can only happen if either $\ell_2/\ell_3 \to 0$ or $\ell_3/\ell_2 \to 0$. Suppose $\ell_3/\ell_2 \to 0$, which implies $\ell_2 \to \infty$.

Suppose that the variance averages $\frac{\sum_{g:n_g=2}\sum_{i\in G_g}\sigma_i^2}{2\ell_2}$ and $\frac{\sum_{g:n_g=3}\sum_{i\in G_g}\sigma_i^2}{3\ell_3}$ converge as $n \to \infty$; let

$$\overline{\sigma}_2^2 = \lim_{n \to \infty} \frac{\sum_{g: n_g = 2} \sum_{i \in G_g} \sigma_i^2}{2\ell_2}, \quad \overline{\sigma}_3^2 = \lim_{n \to \infty} \frac{\sum_{g: n_g = 3} \sum_{i \in G_g} \sigma_i^2}{3\ell_3}.$$

Note that

$$\frac{1}{\sqrt{\ell}} \frac{\ell_2 \ell_3}{2\ell_2 + 3\ell_3} = \frac{1}{\sqrt{\ell_2 + \ell_3}} \frac{\ell_2 \ell_3}{2\ell_2 + 3\ell_3} = \frac{1}{\sqrt{1 + \ell_3/\ell_2}} \frac{\ell_3/\sqrt{\ell_2}}{2 + 3\ell_3/\ell_2}$$

Therefore, if $\ell_3/\sqrt{\ell_2} \to 0$ then $\frac{1}{\sqrt{\ell}} \frac{\ell_2 \ell_3}{2\ell_2 + 3\ell_3} \to 0$. In this case

 $\mathrm{E}\left[E_2\right] \to 0.$

If $\ell_3/\sqrt{\ell_2} \to \infty$ then $\frac{1}{\sqrt{\ell}} \frac{\ell_2 \ell_3}{2\ell_2 + 3\ell_3} \to \infty$. In this case $E[E_2]$ can be unbounded; specifically

$$\mathbf{E}\left[E_{2}\right] \rightarrow \begin{cases} \infty & \text{if } \overline{\sigma}_{2}^{2} > \overline{\sigma}_{3}^{2}, \\ -\infty & \text{if } \overline{\sigma}_{2}^{2} < \overline{\sigma}_{3}^{2}, \end{cases}$$

and therefore, E_2 is not bounded in probability. Consequently, the statistic AR_{AG} is not bounded in probability. We summarize our findings in the following. **Proposition B.1.** Suppose that $\ell_3/\ell_2 \to 0$ and that the variance averages $\frac{\sum_{g:n_g=2}\sum_{i\in G_g}\sigma_i^2}{2\ell_2}$ and $\frac{\sum_{g:n_g=3}\sum_{i\in G_g}\sigma_i^2}{3\ell_3}$ converge to $\overline{\sigma}_2^2$ and $\overline{\sigma}_3^2$, respectively. Then, if $\ell_3/\sqrt{\ell_2} \to 0$, $E_2 = o_p(1)$; if $\ell_3/\sqrt{\ell_2} \to \infty$ and $\overline{\sigma}_2^2 \neq \overline{\sigma}_3^2$, E_2 is not bounded in probability.

Remark B.1. When the plug-in is the JIV1 estimator $\widetilde{\beta}_2 = (X'_2 C X_2)^{-1} X'_2 C y_0$ we obtain that

$$T_{2} = \frac{1}{\sqrt{k}} \frac{\varepsilon' C \varepsilon - \varepsilon' C X_{2} (X_{2}' C X_{2})^{-1} X_{2}' C \varepsilon}{\sqrt{\widehat{V}\left(\widetilde{\beta}\right)}}.$$
(B.6)

The formula in (B.6) suggests that T_2 is more likely to be negative than positive in finite samples, which may explain the underrejection results in our Monte Carlo simulations. See, e.g., Figures C.7 and C.8 in Section B. Indeed, we know that $\mathbb{E}\left[\boldsymbol{\varepsilon}'\boldsymbol{C}\boldsymbol{\varepsilon}\right] = 0$ and since $\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{X}_2\boldsymbol{H}_{22}^{-1/2} \rightarrow_p \boldsymbol{I}_{g_2}, \left(\boldsymbol{H}_{22}^{-1/2}\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{X}_2\boldsymbol{H}_{22}^{-1/2}\right)^{-1}$ is likely to be positive definite in sufficiently large finite samples. Therefore, $\boldsymbol{\varepsilon}'\boldsymbol{C}\boldsymbol{X}_2(\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{X}_2)^{-1}\boldsymbol{X}_2'\boldsymbol{C}\boldsymbol{\varepsilon} \geq 0$, so the numerator of T_2 is more likely to take negative values, unless the sign of $\boldsymbol{\varepsilon}'\boldsymbol{C}\boldsymbol{\varepsilon}$ interacts with the magnitude of $\hat{V}\left(\boldsymbol{\beta}\right)$ in a special way. This suggests that the density of T_2 is shifted to the left, which leads to underrejection.

Remark B.2. The assumption $\frac{\sqrt{k}}{r_{\min}} \to 0$ in Theorem B.2 appears to be necessary. Suppose that this assumption is violated while Assumptions 1, 2, 3 hold; for simplicity consider the case when $g_2 = 1$ and denote the single endogenous variable as \mathbf{x}_2 . Moreover, $\mathbf{x}_2 = \mathbf{Z}\mathbf{\pi}_2 + \mathbf{u}_2$. In this case $r = r_{\min} = H_{22}$ and suppose that $\sqrt{k}/r = \tau_n$ with $\tau_n \ge c_{\tau} > 0$ for any n. One important special case is when τ_n converges to a positive number; another special case is when τ_n goes to ∞ .

We note first that in this case the JIV1-type estimator β₂ = β₂ + (**x**'₂C**x**₂)⁻¹**x**'₂Cε is not consistent. Indeed, by (A.4)

$$\frac{1}{\sqrt{k}} \bm{x}_{2}' \bm{C} \bm{x}_{2} = \frac{1}{\sqrt{k}} \left(\bm{Z} \bm{\pi}_{2} \right)' \bm{C} \bm{Z} \bm{\pi}_{2} + \frac{1}{\sqrt{k}} \left(\bm{Z} \bm{\pi}_{2} \right)' \bm{C} \bm{u}_{2} + \frac{1}{\sqrt{k}} \bm{u}_{2}' \bm{C} \bm{Z} \bm{\pi}_{2} + \frac{1}{\sqrt{k}} \bm{u}_{2}' \bm{C} \bm{u}_{2}$$

Note that

$$\mathbf{E}\left[\frac{1}{\sqrt{k}} \left(\boldsymbol{Z}\boldsymbol{\pi}_{2}\right)' \boldsymbol{C}\boldsymbol{u}_{2}\right] = 0, \\ \operatorname{Var}\left[\frac{1}{\sqrt{k}} \left(\boldsymbol{Z}\boldsymbol{\pi}_{2}\right)' \boldsymbol{C}\boldsymbol{u}_{2}\right] = \frac{1}{k} \mathbf{E}\left[\left(\boldsymbol{Z}\boldsymbol{\pi}_{2}\right)' \boldsymbol{C}\boldsymbol{u}_{2}\boldsymbol{u}_{2}' \boldsymbol{C}\boldsymbol{Z}\boldsymbol{\pi}_{2}\right] \leq \frac{c_{u}}{k} r \leq \frac{c_{u}}{c_{\tau}\sqrt{k}},$$

where the first inequality follows from (A.7), so

$$\frac{1}{\sqrt{k}}\boldsymbol{x}_{2}^{\prime}\boldsymbol{C}\boldsymbol{x}_{2}=\frac{1}{\tau_{n}}+\frac{1}{\sqrt{k}}\boldsymbol{u}_{2}^{\prime}\boldsymbol{C}\boldsymbol{u}_{2}+o_{p}\left(1\right).$$

Since by Lemma A.1 $\frac{1}{\sqrt{k}} \boldsymbol{u}_2' \boldsymbol{C} \boldsymbol{u}_2 \rightarrow_d N_1$ with N_1 normal with mean 0, we obtain that

$$\frac{1}{\sqrt{k}} \boldsymbol{x}_{2}^{\prime} \boldsymbol{C} \boldsymbol{x}_{2} = \frac{1}{\tau_{n}} + N_{1} + o_{p} (1) .$$
(B.7)

Further,

$$rac{1}{\sqrt{k}}oldsymbol{x}_2'oldsymbol{C}oldsymbol{arepsilon} = rac{1}{\sqrt{k}}\left(oldsymbol{Z}oldsymbol{\pi}_2
ight)'oldsymbol{C}oldsymbol{arepsilon} + rac{1}{\sqrt{k}}oldsymbol{u}_2'oldsymbol{C}oldsymbol{arepsilon},$$

where the first term is $o_p(1)$ for similar reasons as above and the second term is asymptotically normal with mean 0 (from Lemma A2 in Chao et al., 2012). Therefore, we can write

$$\frac{1}{\sqrt{k}}\boldsymbol{x}_{2}^{\prime}\boldsymbol{C}\boldsymbol{\varepsilon}=N_{2}+o_{p}\left(1\right),\tag{B.8}$$

with N_2 normal with mean 0. So, from (B.7) and (B.8)

$$\widetilde{\beta}_2 - \beta_2 = \frac{N_2 + o_p(1)}{1/\lambda_n + N_1 + o_p(1)},$$

which in general is not $o_p(1)$, so $\tilde{\beta}_2$ is not consistent. Therefore, we cannot prove that $\hat{V}\left(\tilde{\boldsymbol{\beta}}\right) - \hat{V}\left(\boldsymbol{\beta}\right) = o_p(1)$ in the way we do above (Lemmas A.2 and A.3).

2. In order to derive the limit of T_2 in Theorem B.2 we proved that $\frac{1}{\sqrt{k}} \boldsymbol{\varepsilon}' \boldsymbol{C} \boldsymbol{x}_2 (\boldsymbol{x}_2' \boldsymbol{C} \boldsymbol{x}_2)^{-1} \boldsymbol{x}_2' \boldsymbol{C} \boldsymbol{\varepsilon} = o_p(1)$. In this case (B.7) and (B.8) imply

$$\frac{1}{\sqrt{k}}\boldsymbol{\varepsilon}'\boldsymbol{C}\boldsymbol{x}_{2}(\boldsymbol{x}_{2}'\boldsymbol{C}\boldsymbol{x}_{2})^{-1}\boldsymbol{x}_{2}'\boldsymbol{C}\boldsymbol{\varepsilon} = \left(\frac{1}{\sqrt{k}}\boldsymbol{\varepsilon}'\boldsymbol{C}\boldsymbol{x}_{2}\right)^{2}\left(\frac{1}{\sqrt{k}}\boldsymbol{x}_{2}'\boldsymbol{C}\boldsymbol{x}_{2}\right)^{-1} = \frac{\left(N_{2}+o_{p}\left(1\right)\right)^{2}}{1/\lambda_{n}+N_{1}+o_{p}\left(1\right)}$$

which is not $o_p(1)$ in general.

Example B.2. This example is motivated by the fact that, in practice, applied researchers may erroneously choose an inconsistent plug-in estimator. It is reasonable to think that such a choice may affect the behaviour of T_2 . In order to simplify the analysis we assume that the plug-in is consistent but it converges at an arbitrary slow rate to the true value. The assumption of consistency allows us to use Lemma A.3. Let us consider a simple two-regressor model

$$oldsymbol{y} = oldsymbol{x}_1eta_1 + oldsymbol{x}_2eta_2 + oldsymbol{arepsilon}$$

where \mathbf{x}_1 and \mathbf{x}_2 may both be endogenous and suppose that we want to test the following null $H_0: \beta_1 = \beta_{10}$. Let us assume that there exists an estimator for β_2 , say $\tilde{\beta}_2$, such that $\sqrt{a_n}(\tilde{\beta}_2 - \beta_2) = O_p(1)$ where $a_n \to \infty$ as $n \to \infty$. This situation defines a consistent but potentially slowly converging estimator. Let us also define

$$oldsymbol{x}_2 = oldsymbol{Z}oldsymbol{\pi}_2 + oldsymbol{u}_2$$

with $H_{22} = \pi'_2 \mathbf{Z}' \mathbf{Z} \pi_2$, in this case $r = r_{\min} = H_{22}$ and $r \to \infty$ as $n \to \infty$. Let us suppose that a_n and r diverge to infinity possibly at different rates. If we assume that $\widetilde{\boldsymbol{\beta}} = (\beta_{10}, \widetilde{\beta}_2)'$ is consistent we can use Lemmas A.2 and A.3. From Equation (23) in the main text we notice that the fact that T_2 converges to a standard normal would now depend only on the behaviour of Δ (see Equation (24) in the main text), which in this case is

$$\Delta = \left(\widetilde{\beta}_2 - \beta_2\right) \boldsymbol{x}_2' \boldsymbol{C} \boldsymbol{x}_2 \left(\widetilde{\beta}_2 - \beta_2\right) - 2\left(\widetilde{\beta}_2 - \beta_2\right) \boldsymbol{x}_2' \boldsymbol{C} \boldsymbol{\varepsilon}$$

Since
$$\frac{\mathbf{x}_2'\mathbf{C}\mathbf{x}_2}{r} \to_p 1$$
, we get

$$\frac{1}{\sqrt{k}} \left(\widetilde{\beta}_2 - \beta_2\right) \mathbf{x}_2'\mathbf{C}\mathbf{x}_2 \left(\widetilde{\beta}_2 - \beta_2\right) = \frac{1}{\sqrt{k}} \frac{r}{a_n} \sqrt{a_n} \left(\widetilde{\beta}_2 - \beta_2\right) \frac{\mathbf{x}_2'\mathbf{C}\mathbf{x}_2}{r} \sqrt{a_n} \left(\widetilde{\beta}_2 - \beta_2\right) = O_p \left(\frac{r}{a_n\sqrt{k}}\right)$$
(B.9)

Moreover,

$$\frac{1}{\sqrt{k}} \left(\widetilde{\beta}_2 - \beta_2 \right) \boldsymbol{x}_2' \boldsymbol{C} \boldsymbol{\varepsilon} = \frac{\sqrt{r}}{\sqrt{ka_n}} \sqrt{a_n} \left(\widetilde{\beta}_2 - \beta_2 \right) \frac{\boldsymbol{x}_2' \boldsymbol{C} \boldsymbol{\varepsilon}}{\sqrt{r}} = o_p \left(\frac{\sqrt{r}}{\sqrt{a_n} k^{1/4}} \right).$$
(B.10)

This means that if $\frac{r}{a_n\sqrt{k}}$ diverges, Δ does not go to zero and T_2 would not converge to a standard normal. Notice that Δ does not go to zero when a_n grows slower or at the same rate of the boundary condition r/\sqrt{k} . In this case the distribution of T_2 will be shifted to the right causing the test to overreject.

C Monte Carlo Experiments

This Section collects some complementary Monte Carlo results on the finite sample properties of T_1 and T_2 . The simulations consider two DGPs and both the homoskedastic and heteroskedastic case. Apart from the T_1 and T_2 statistics we include the AR_{AG} test of Anatolyev and Gospodinov (2011) and three test statistics due to Bun *et al.* (2018)

$$\widehat{AR} = n\widehat{g}(\boldsymbol{\beta})'\widehat{\boldsymbol{\Omega}}(\boldsymbol{\beta})^{-1}\widehat{g}(\boldsymbol{\beta})$$
$$\widetilde{AR} = n\widehat{g}(\boldsymbol{\beta})'\widetilde{\boldsymbol{\Omega}}(\boldsymbol{\beta})^{-1}\widehat{g}(\boldsymbol{\beta})$$
$$\widetilde{AR}_{df} = n\widehat{g}(\boldsymbol{\beta})'\widetilde{\boldsymbol{\Omega}}_{df}(\boldsymbol{\beta})^{-1}\widehat{g}(\boldsymbol{\beta})$$

where $\widehat{\Omega}(\beta) = \frac{1}{n} \sum_{i=1}^{n} g(\beta) g(\beta)'$, $\widetilde{\Omega}(\beta) = \widehat{\Omega}(\beta) - \widehat{g}(\beta) \widehat{g}(\beta)'$ and $\widetilde{\Omega}_{df}(\beta) = \frac{n}{n-k} \widetilde{\Omega}(\beta)$. In our case the moment condition model is defined as $g_i(\beta) = Z_i(y_i - X'_i\beta)$ and $\widehat{g}(\beta) = \frac{1}{n} \sum_{i=1}^{n} g_i(\beta)$. The evaluation of the performance is made in terms of size and power. Furthermore, the second DGP is also used to assess the quality of the asymptotic approximations as presented in Corollary 1 in the main text.

The first DGP (DGP I) is similar to Bekker and Van der Ploeg (2005) where the instruments are dummies. In this experiment the observations are stratified in k groups where each group contains n_j observations and $n = \sum_{j=1}^k n_j$ and each group contains a different number of observations. Let us define the model

$$egin{aligned} egin{aligned} egi$$

where the true value of β is zero and Z is a $n \times k$ matrix of dummy variables, such that each of its rows is a versor. Moreover, for each group, the disturbances are jointly normally distributed with zero mean and variance covariance matrix equal to

$$\boldsymbol{\Sigma}_{j} = \begin{pmatrix} \sigma_{j}^{2} & \rho \sigma_{j} \sigma_{vj} \\ \rho \sigma_{j} \sigma_{vj} & \sigma_{vj} \end{pmatrix}, \quad j = 1, \dots, k.$$

We choose $\rho = 0.5$ and $(k, n) \in \{(7, 146), (40, 140), (60, 168)\}$. The parameters σ_j and σ_{vj} are sampled independently from a uniform distribution $\mathcal{U}(0.5, 1)$. We consider both the homoskedastic case where Σ_j is the same for any j and the corresponding heteroskedastic case. Furthermore, the elements of π are sampled from $\mathcal{U}(0.05, 0.1)$. The experiment is replicated 5000 times.

The second DGP (DGP II) (Hausman *et al.*, 2012) is given by

$$y = \iota \gamma + x\beta + \varepsilon$$
(C.2)
$$x = z\pi + v$$

where $\gamma = \beta = 1$, while $\pi = 0.1$ in the analysis of size and $\pi \in \{0.1, 1\}$ in the analysis of power. The sample size is n = 800, $\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_n)$ and independently $\boldsymbol{v} \sim \mathcal{N}(\boldsymbol{0}, 0.1^2 \times \boldsymbol{I}_n)$.

The disturbances vector $\boldsymbol{\varepsilon}$ is generated as

$$\boldsymbol{\varepsilon} = \rho \boldsymbol{v} + \sqrt{\frac{1 - \rho^2}{\phi^2 + \psi^4}} (\phi \boldsymbol{w}_1 + \psi \boldsymbol{w}_2), \tag{C.3}$$

where $\rho = 0.3$, $\psi = 0.86$ and conditional on z, independent of v, $w_1 \sim \mathcal{N}(\mathbf{0}, \operatorname{Diag}(z)^2)$ where $\operatorname{Diag}(z)$ is a diagonal matrix where the diagonal elements are the elements of zand $w_2 \sim \mathcal{N}(\mathbf{0}, \psi^2 I_n)$. Moreover, $\phi \in \{0, 1.38072\}$, where $\phi = 0$ is the homoskedastic case. The instrument matrix Z is given by matrices with rows $(1, z_i, z_i^2, z_i^3, z_i^4)$ and $(1, z_i, z_i^2, z_i^3, z_i^4, z_i b_{1i}, \ldots, z_i b_{\ell i}), \ell = 95, 695$, where, independent of other random variables, the elements $b_{1i}, \ldots, b_{\ell i}$ are i.i.d. Bernoulli distributed with p = 1/2.³ We consider also two rather extreme situations: k = 2 and k = 700. We replicate our experiments 5000 times. When using the T_1 test and the T_2 test we consider $H_0: (\gamma, \beta)' = (1, 1)'$ and $H_0: \beta = 1$ respectively.

C.1 Simulation results

We first discuss the quality of the potential approximations for T_1 and T_2 when k = 2, then we provide some interpretation of the simulations by separately analysing the results on size and power. We also discuss the behaviour of T_2 when an inconsistent plug-in is used.

Approximations. In Figure C.1, we explore the behaviour of T_1 for k = 2 and n = 50, 100, 200, 400, 800. It seems clear that, in this case, the chi square approximation for T_1 is more accurate than its Gaussian counterpart. This result is less evident in the case of T_2 , since, as shown in Corollary 1, three alternative chi square approximations are available. Nonetheless, Figure C.2 panel (c) suggests that result (*iii*) in Corollary 1 may cause the test to reject too often. On the other hand, the approximations in (*ii*) and (*iv*) of Corollary 1 deliver more reliable results (Figure C.2 panels (b) and (d)).

Size. In the case of DGP I (Figure C.3 and Figure C.4), heteroskedasticity is rather mild and, as expected, the various statistics perform quite similarly in the homoskedastic

³The same set of instruments is used throughout the various repetitions.

and heteroskedastic case. In addition, we observe that T_1 , AR_{AG} and AR_{df} work well for the three combinations of k and n considered. On the other hand, \widehat{AR} tends to underreject as $\frac{k}{n}$ gets larger, while \widehat{AR} tends to underreject. The case of DGP II (Figure C.5 to Figure C.8) is more complex, as the type of heteroskedasticity introduced in the model may have a non trivial impact on the performance of the tests. In general, we observe that T_1 and T_2 work well in all the considered cases and AR_{AG} performs well in general under homoskedasticity and, as it is expected, it shows some tendency to overreject when $k = 700.^4$ As we introduce heteroskedasticity, the performance of the AR_{AG} test dramatically deteriorates. The tests introduced in Bun *et al.* (2018) work well for most of the cases but tend to either underreject $(\widehat{AR}, \widehat{AR})$ or overreject (\widehat{AR}_{df}) when k is large.

Power. Under homoskedasticity and k small the tests are indistinguishable (Figure C.9 to Figure C.12 panels (a) and (b)). The picture gets more complicated as k increases. In particular, with $\pi = 1$ all the test apart from \widetilde{AR} can control size and have excellent power properties (Figure C.11 and Figure C.12 panel (c)). However, when $\pi = 0.1$, the power properties of all the tests, in particular \widehat{AR} and \widehat{AR}_{df} , deteriorate (Figure C.9 and Figure C.10 panel (c)). In the heteroskedastic case and when k = 2, 5, the T_1 and the T_2 tests along with the tests of Bun *et al.* (2018) are able to discriminate among alternatives (Figure C.13 to Figure C.16 panels (a) and (b)). To some extent the same could be said about the case where k = 100 (Figure C.13 to Figure C.16 panel (c)). When k = 700, $\pi = 0.1$ no test statistic among those considered seems to work well in this case. Only the AR_{AG} test has some power in the homoskedastic case (Figure C.9 and C.10 panel (d)). However, when $\pi = 1$, the T_1 and T_2 tests tend to outperform their competitors (Figure C.9 to Figure C.16 panel (d)).

⁴It is worth noticing that, in general, for the hypothesis $H_0: \beta_1 = \beta_{10}$ all the tests tend to underreject for small values of k.

C.2 Figures



Figure C.1: PP-plots for T_1 under DGP II with heteroskedasticity, k = 2 and n = 50, 100, 200, 400, 800.



Figure C.2: PP-plots for T_2 under DGP II with heteroskedasticity, k = 2 and n = 50, 100, 200, 400, 800.



Figure C.3: PP-plots with homoskedasticity under DGP I, $H_0: \beta = \beta_0$.



Figure C.4: PP-plots with heteroskedasticity under DGP I, $H_0: \beta = \beta_0$.


Figure C.5: PP-plots with homoskedasticity under DGP II, $H_0: \beta = \beta_0$.



Figure C.6: PP-plots with heteroskedasticity under DGP II, $H_0: \beta = \beta_0$.



Figure C.7: PP-plots with homoskedasticity under DGP II, $H_0: \beta_1 = \beta_{10}$.



Figure C.8: PP-plots with heteroskedasticity under DGP II, $H_0: \beta_1 = \beta_{10}$.



Figure C.9: Power curves with homoskedasticity and $\pi = 0.1$, $H_0: \beta = \beta_0$.



Figure C.10: Power curves with homoskedasticity and $\pi = 0.1$, $H_0: \beta_1 = \beta_{10}$.



Figure C.11: Power curves with homoskedasticity and $\pi = 1$, $H_0: \beta = \beta_0$.



Figure C.12: Power curves with homoskedasticity and $\pi = 1$, $H_0: \beta_1 = \beta_{10}$.



Figure C.13: Power curves with heteroskedasticity and $\pi = 0.1$, $H_0: \beta = \beta_0$.



Figure C.14: Power curves with heteroskedasticity and $\pi = 0.1$, $H_0: \beta_1 = \beta_{10}$.



Figure C.15: Power curves with heteroskedasticity and $\pi = 1$, $H_0: \beta = \beta_0$.



Figure C.16: Power curves with heteroskedasticity and $\pi = 1$, $H_0: \beta_1 = \beta_{10}$.

References

- Anatolyev, S. and Gospodinov, N. (2011) Specification Testing in Models with Many Instruments. *Econometric Theory* 27, 427–441.
- Bekker, P.A. and Crudu, F. (2015) Jackknife Instrumental Variable Estimation with Heteroskedasticity. *The Journal of Econometrics* 185, 332–342.
- Bekker, P.A. and Van der Ploeg, J. (2005) Instrumental variable estimation based on grouped data. *Statistica Neerlandica* 59, 239–267.

- Bun, M., Farbmacher, H. and Poldermans, R. (2018) Finite sample properties of the Anderson and Rubin (1949) test. working paper .
- Chao, J.C., Swanson, N.R., Hausman, J.A., Newey, W.K. and Woutersen, T. (2012) Asymptotic distribution of JIVE in a heteroskedastic IV regression with many instruments. *Econometric Theory* 28, 42–86.
- Hausman, J.A., Newey, W.K., Woutersen, T., Chao, J.C. and Swanson, N.R. (2012) Instrumental variable estimation with heteroskedasticity and many instruments. *Quantitative Economics* 3, 211–255.