

Linear Regression with Weak Exogeneity^{*}

ANNA MIKUSHEVA[†], MIKKEL SØLVSTEN[‡]

August 2023

Abstract

This paper studies linear time series regressions with many regressors. Weak exogeneity is the most used identifying assumption in time series. Weak exogeneity requires the structural error to have zero conditional expectation given the present and past regressor values, allowing errors to correlate with future regressor realizations. We show that weak exogeneity in time series regressions with many controls may produce substantial biases and even render the least squares (OLS) estimator inconsistent. The bias arises in settings with many regressors because the normalized OLS design matrix remains asymptotically random and correlates with the regression error when only weak (but not strict) exogeneity holds. This bias's magnitude increases with the number of regressors and their average autocorrelation. To address this issue, we propose an innovative approach to bias correction that yields a new estimator with improved properties relative to OLS. We establish consistency and conditional asymptotic Gaussianity of this new estimator and provide a method for inference.

KEYWORDS: time series regression, weak exogeneity, many controls, feedback bias, OLS inconsistency, bias correction, valid inference

JEL CODES: C13, C22

^{*}We are grateful to Isaiah Andrews, Morten Ø. Nielsen, Jack Porter, and Jim Stock for useful discussions. Harvey Barnhard, Bas Sanders, and Chris Walker provided excellent research assistance.

[†]Department of Economics, M.I.T., 50 Memorial Drive, E52-526, Cambridge, MA, 02142, United States. E-mail: amikushe@mit.edu.

[‡]Department of Economics and Business Economics, Aarhus University, Fuglsangs Allé 4, Building 2621, B 13a, 8210 Aarhus V, Denmark. E-mail: miso@econ.au.dk. This research was supported by grants from the Danish National and Aarhus University Research Foundations (DNRF Chair #DNRF154 and AUFF Grant #AUFF-E-2022-7-3)

1 Introduction

Structural estimation in macroeconomics, finance and other economic fields studying dynamic models often employs time series data. The most used identifying assumption for structural estimation in time series settings is weak exogeneity. Weak exogeneity postulates that the structural shock has zero conditional expectation given the present and past regressor values. It is a less restrictive assumption than strict exogeneity, which additionally requires that the shocks have zero conditional expectation given future values of regressors. Strict exogeneity is implausible in most settings due to *feedback*, i.e., the outcome variable in one period affects the values of the regressors in future periods.¹ Specifically, if the lagged outcome variable is among the regressors, strict exogeneity cannot hold.

Another common feature of modern structural regressions is the presence of many regressors, all of which may be autocorrelated. Various motivations for the use of many regressors in time series are that the economic system generating the data is partially observed (Zellner and Palm, 1974; Wallis, 1977), additional controls in local projections may ensure uniformity (Jordà, 2005; Montiel Olea and Plagborg-Møller, 2021), and long memory may be arising from an underlying high-dimensional model (Schennach, 2018; Chevillon et al., 2018). See also Bauwens et al. (2023) for two applications.

We show that these two features – weak exogeneity and many autocorrelated regressors – can produce substantial biases and even lead to inconsistency of the ordinary least squares (OLS) estimator. The large bias in OLS may arise even when all variables are stationary (i.e., no unit roots or strong persistence is needed) and when the feedback effect violating strict exogeneity is limited to just one period. Finite sample unbiasedness of OLS relies heavily on strict exogeneity. It is well understood that OLS is biased in most time series regressions, but there is also a common belief that the biases are relatively small and of second order (see, e.g., Bao and Ullah, 2007). Our results show that such beliefs are unwarranted and that the bias in OLS can be a first-order issue.

This paper contains several results. First, we derive a formula for the asymptotically non-negligible part of the OLS bias, explain which features of the data are responsible for it, and provide a tool to assess the potential for OLS bias in a given time series application. Second,

¹The formalization of feedback is typically ascribed to Granger (1969), while Engle et al. (1983) provide an early rigorous distinction between weak and strict exogeneity. See also Sims (1972); Chamberlain (1982) for further discussions and an empirical example.

we propose a new estimator which is asymptotically unbiased if the data contains a one-period violation of strict exogeneity. Third, we derive the asymptotic distribution for this new estimator and discuss how to conduct inferences using our new approach. Surprisingly, bias correction does not necessarily trade off with increased variance. Based on our simulations, the standard deviations of OLS and our new estimator are very close to each other with no clear ordering. Finally, we show how our new estimator generalizes to settings where strict exogeneity is violated by multiple periods of feedback effects.

The bias of OLS arises in a setting with many regressors because the properly normalized design matrix, $\frac{1}{T}X'X$, remains asymptotically random even in large samples. Weak (but not strict) exogeneity allows for correlation between the randomness in the design matrix and the numerator of the OLS estimator, which leads to a bias. Specifying the feedback structure allows us to derive a formula for the leading term of the bias. Specifically, once we assume that strict exogeneity is violated by one-period feedback from the outcome variable to the next period's regressors, we show that the OLS bias is aligned with the feedback direction. The size of the bias increases with the number of regressors and their one-period ahead linear predictability.

We propose a new estimator which eliminates the bias asymptotically and is consistent under the same assumptions that may lead to inconsistency of OLS. Our proposal mimics an instrumental variables (IV) estimator with an intentionally endogenous 'technical' instrument: a linear combination of the regressors and their leads (future values). The main insight is that future values of the regressors in the instrument induce an endogeneity bias along the feedback direction only, the same direction along which the OLS is biased. It is therefore possible to pick the weights in the linear combination to ensure that the bias stemming from the endogenous instrument offsets the bias originating from weak exogeneity.

An important feature of our bias correction is that it is constructed based on knowledge about the regressors only. The correction is the same for any outcome variable and does not require knowledge about or estimation of the direction of the feedback mechanism. We show that the new estimator is consistent and, after proper normalization, asymptotically Gaussian when there is a one-period violation of strict exogeneity. The main results can be generalized to the case where the violation of strict exogeneity is for a finite number of periods, that is when the outcome variable has feedback effects on the regressors for L periods. In such a case, the bias of OLS contains L terms, corresponding to the L directions

of feedback.

We conduct a simulation study aimed at assessing how common and how large the OLS bias is in typical macroeconomic regressions. We take a large collection of US macroeconomic indexes observed at quarterly frequency for 200 periods, extract its business cycle part, and randomly draw a regression from this data set. We show that the time series dependence and magnitude of feedback typical for these macroeconomic data produce empirically important OLS biases. For example, in a typical regression with 25 regressors we find a bias in the feedback direction equal to half of the standard deviation, while in regression with 50 regressors this bias approximately equals one standard deviation. Depending on the number of regressors, approximately 6-21percent of coefficients display a statistically significant difference between the OLS estimator and our proposed estimator. Our proposed IV-type estimator provides full correction of the bias.

Our results are related to three distinct strands of literature. First, there is a classical literature on linear equations in time series which shows the inappropriateness of the Generalized Least Squares (GLS) estimator in linear models with only weak exogeneity. For example, [Hansen and West \(2002\)](#) points out that GLS mixes up the timing of observations and – in the case of weakly exogenous regressors – leads to significant biases. The authors argue against using GLS in macroeconomic data. A concern that weak exogeneity may lead to large biases in OLS was raised by [Stambaugh \(1999\)](#), who considered a regression model with a very persistent (near-unit-root) regressor. Also in the absence of persistence, we show that even larger biases arise when the number of regressors is large. The second literature is that on estimation of dynamic effects in panel data, where the presence of fixed effects (many regressors) produces a sizable bias in the coefficient on the lagged outcome variable (the weakly exogenous regressor) ([Nickell, 1981](#)). Unlike the solutions proposed in that literature (e.g., [Arellano and Bond, 1991](#)), our solution for the time series context does not rely on the knowledge that only a single known regressor fails to be strictly exogenous nor does it require that the many regressors in the model are fixed effects for mutually exclusive groups. Finally, the underlying algebraic source of the bias issues, as well as some asymptotic statements related to Gaussianity of quadratic forms, are connected to the issues arising in linear models with many instruments and/or many regressors – see [Hansen et al. \(2008\)](#); [Chao et al. \(2012\)](#); [Kline et al. \(2020\)](#).

The rest of the paper is organized as follows. Section 2 derives the formula for the leading

term of the OLS bias, provides intuition for the findings, and discusses features of the data responsible for the bias. Section 3 introduces a new asymptotically unbiased estimator under the assumption of one-period feedback and establishes its consistency. Section 4 establishes asymptotic Gaussianity of the newly proposed estimator and suggests a valid inference procedure. Section 5 extends some results to settings with multi-period feedback. Section 6 contains simulation studies assessing the empirical relevance of the discussed issues in typical macroeconomic data sets. All proofs are in Appendix A.

Notation: For any vector x , $\|x\| = \sqrt{x^\top x}$ gives the L_2 norm of x . For any matrix A (not necessarily square) $\text{rk}(A)$ is the rank of A , $\|A\| = \sup_x \frac{\|Ax\|}{\|x\|}$ denotes the operator norm (the positive square root of the largest eigenvalue of $A^\top A$), and $\|A\|_F = \sqrt{\text{tr}(A^\top A)}$ denotes the Frobenius norm. Throughout, we let c and C be positive and finite (generic) constants, that may have different values in different equations, but do not depend on the sample size T . The $q \times q$ identity matrix is I_q while $I = I_T$.

2 Inconsistency of OLS

2.1 Model and assumptions

Consider a linear time series regression

$$y_t = x_t^\top \beta + \varepsilon_t, \quad t = 1, \dots, T,$$

where the regressors $x_t \in \mathbb{R}^K$ are weakly exogenous and the number of regressors K is large. We model this by assuming that K may diverge proportionally to T or slower. All features of the data generating process are implicitly indexed by T , but we drop this index for compactness of notation. The object of interest is the linear contrast $\theta = r^\top \beta$ for a known (non-random) K vector r . A leading case is $r^\top \beta = \beta_1$.

The assumption of weak exogeneity imposes:

$$E[\varepsilon_t | x_t, x_{t-1}, \dots] = 0.$$

Weak exogeneity is considerably less restrictive than strict exogeneity, which assumes that

$E[\varepsilon_t/X] = 0$, with X denoting the complete $T \times K$ set of regressors (including past, present, and future). In the context of time series economic data, strict exogeneity is seldom plausible. In contrast, weak exogeneity allows the structural shock to influence future regressor values through feedback. The likelihood of such an effect is substantial, as variables employed in macro estimation often evolve simultaneously within dynamic processes. When outcome variable lags are incorporated as regressors, strict exogeneity is certainly violated. For a textbook discussion regarding the plausibility of weak and strict exogeneity in various applications, see [Stock and Watson \(2019, ch. 16\)](#).

The standard argument for unbiasedness of the OLS estimator $r^\theta \hat{\beta}^{\text{OLS}} = r^\theta (X^\theta X)^{-1} X^\theta y$ heavily relies on strict exogeneity. Although it is well-known that OLS exhibits bias under weak exogeneity (see, e.g., [Hamilton, 1994](#), chapter 8.2), a common belief persists that this bias is small and asymptotically negligible. Here we claim that this widespread belief is both incorrect and misleading: the bias of OLS can be substantial, potentially leading to inconsistency.

Standard statements regarding OLS consistency in a time series context hinge on the assumption that the normalized design matrix concentrates around its expectation: $k \frac{1}{T} X^\theta X^\theta Q k \neq 0$, where Q is non-singular (see, e.g., [Hamilton, 1994](#), Assumption 8.6). For example, [Gupta and Seo \(2023\)](#) assumes that $K^3/T \neq 0$ to invoke a Law of Large Numbers for the normalized design matrix. However, the normalized design matrix remains asymptotically random when the number of regressors (and thus design matrix dimensionality) grows fast enough. Under strict exogeneity the randomness of $\frac{1}{T} X^\theta X^\theta$ does not cause a problem as one can condition on X in the analysis, thus treating the design matrix as non-random. This approach is infeasible under weak exogeneity. In such instances, $\frac{1}{T} X^\theta X^\theta$ is not just a random matrix, but is also correlated with $X^\theta y$. This correlation produces a bias that may become the leading term in the asymptotic analysis of OLS.

The size and form of the OLS bias depend on the feedback mechanism, or how past errors in the outcome variable affect future regressor values. We start by assuming one-period linear feedback, i.e., the present error term ε_t only affects the regressors in the subsequent period, namely x_{t+1} . [Section 5](#) extends our results to feedback lasting a finite number of periods.

Assumption 1. (i) *The observed regressors x_t can be decomposed as $x_t = \tilde{x}_t + \alpha \varepsilon_{t-1}$, where the $T \times K$ matrix $\tilde{X} = [\tilde{x}_1, \dots, \tilde{x}_T]^\theta$ has full rank.*

- (ii) The errors $\tilde{f}_{\varepsilon_t} g_{t=0}^T$ are i.i.d. conditionally on \tilde{X} with $E[\varepsilon_t j \tilde{X}] = 0$, $c < \sigma^2 := E[\varepsilon_t^2 j \tilde{X}] < C$, and $E[\varepsilon_t^4 j \tilde{X}] < C$ almost surely.
- (iii) The size of the non-random vectors $\alpha, r \in \mathbb{R}^K$ relative to the strictly exogenous design matrix $\frac{1}{T} \tilde{X}^\theta \tilde{X}$ is bounded: $\alpha^\theta (\frac{1}{T} \tilde{X}^\theta \tilde{X})^{-1} \alpha = O_p(1)$ and $r^\theta (\frac{1}{T} \tilde{X}^\theta \tilde{X})^{-1} r = O_p(1)$.
- (iv) The number of regressors K may diverge with sample size T such that $\frac{K}{T} < 1 - c$.

Part (i) describes a specific violation of strict exogeneity where the present error term only affects the regressors in the subsequent period. If $\alpha = 0$, then all regressors $x_t = \tilde{x}_t$ are strictly exogenous. The amount of assumptions imposed on the strictly exogenous part of the regressors $\tilde{f}_{\tilde{x}_t} g_{t=1}^T$ is kept to a minimum, the main part of which is a full rank condition. Specifically, we do not assume stationarity or impose any moment conditions on \tilde{X} . Such generality is possible mostly because our analysis is done conditionally on \tilde{X} . Part (ii) is a standard set of assumptions on the error terms in homoskedastic regression models. Part (iii) places a very loose bound on the magnitudes of α and r relative to the scaled design matrix $\frac{1}{T} \tilde{X}^\theta \tilde{X}$. This condition is, for example, satisfied if $\|\alpha\|$ and $\|r\|$ are bounded and the smallest eigenvalue of $\frac{1}{T} \tilde{X}^\theta \tilde{X}$ is separated from zero. Part (iv) allows a wide range for the number of regressors, including a small fixed number. It requires that the degrees of freedom $T - K$ diverges to infinity as the sample size increases.

The violation of strict exogeneity in Assumption 1 seems minimal as it is solely due to feedback from the dependent variable y_t to the one-period-ahead regressors x_{t+1} and the magnitude of the feedback is bounded. However, this violation is enough to produce inconsistency of the OLS estimator.

2.2 Parameter estimator

In order to provide some intuition for how and why the OLS bias arises, we consider a special case. Suppose, therefore, that only the first regressor fails to be strictly exogenous. Then, the first regressor experiences a one-period feedback: $x_{1t} = \tilde{x}_{1t} + a \varepsilon_{t-1}$, while the $T - (K - 1)$ matrix of all other regressors, $X_{-1} = \tilde{X}_{-1}$, is strictly exogenous. Let X_1 be the $T - 1$ vector with elements x_{1t} and \tilde{X}_1 its strictly exogenous part. Furthermore, we take \tilde{X} as fixed and suppose that $\frac{1}{T} \tilde{X}^\theta \tilde{X} = I_K$.

In this special case, we now derive the bias of the OLS estimator for the first coefficient β_1 . According to the Frisch-Waugh-Lovell theorem, we have

$$\hat{\beta}_1^{\text{OLS}} = \frac{X_1^\theta M_{-1} y}{X_1^\theta M_{-1} X_1} \quad \text{and} \quad \hat{\beta}_1^{\text{OLS}} - \beta_1 = \frac{X_1^\theta M_{-1} \varepsilon}{X_1^\theta M_{-1} X_1},$$

where the partialling out operator $M_{-1} = I - \tilde{X}_{-1}(\tilde{X}_{-1}^\theta \tilde{X}_{-1})^{-1} \tilde{X}_{-1}^\theta$ is the projection on the space orthogonal to \tilde{X}_{-1} . Notice that $M_{-1} = (M_{st})$ is fixed (or can be conditioned on) so that all randomness comes from ε . The denominator of $\hat{\beta}_1^{\text{OLS}} - \beta_1$ has a standard asymptotic behavior and once normalized by $\frac{1}{T}$ converges to a non-random and non-zero limit.

If all regressors were strictly exogenous ($X_1 = \tilde{X}_1$), then the OLS estimator would have been properly centered. However, X_1 contains randomness that is correlated with ε . This leads to a non-zero expectation of the numerator in $\hat{\beta}_1^{\text{OLS}} - \beta_1$:

$$\mathbb{E}[X_1^\theta M_{-1} \varepsilon j \tilde{X}] = \mathbb{E}[\tilde{X}_1^\theta M_{-1} \varepsilon j \tilde{X}] + a \mathbb{E}\left[\sum_{s,t} M_{st} \varepsilon_s - 1 \varepsilon_t j \tilde{X}\right] = a \sigma^2 \sum_t M_{tt} - 1.$$

While the weakly exogenous regressor x_{1t} is not correlated with the contemporaneous error ε_t , the process of partialling out the remaining regressors mixes up the timing of the observations and makes the partialled out regressor correlated with the contemporaneous error. This derivation suggest a form for the leading term of the bias of $\hat{\beta}_1^{\text{OLS}}$. We may notice that the bias depends on a which governs the magnitude of feedback in the regressors and on the trace of the lower diagonal of the projection matrix M_{-1} ($\sum_t M_{tt} - 1$). As discussed below, this lower trace may be on the order of the number of regressors and increases with their average autocorrelation. It is also possible to show that the remaining OLS coefficients (e.g., $\hat{\beta}_2^{\text{OLS}}$) have an asymptotically negligible bias (although they have a finite sample bias). Thus for any linear combination $\theta = r^\theta \beta$, the bias depends on the weight placed on β_1 .

The insight from the preceding special case can be extended. In fact, the OLS estimator of a linear contrast remains invariant under linear transformations of the regressors. For instance, if we perform a regression of Y on XA , where A is a $K \times K$ matrix with full rank that linearly transforms the regressors, then the OLS serves as an estimator of $A^{-1} \beta$. To achieve the contrast θ , we should apply $A^\theta r$ as the weighting.

Any OLS scenario simplifies to the one discussed earlier when we choose $A = (\frac{1}{T} \tilde{X}^\theta \tilde{X})^{-1/2}$, with the square root selected to ensure that $A^\theta \alpha$ is proportional to the first basis vector. These

insights, coupled with certain technical derivations, lead to the following characterization of the OLS bias.

Theorem 1 (Inconsistency of OLS estimator). *Suppose Assumption 1 holds. Then,*

$$r^\ell \hat{\beta}^{\text{OLS}} - r^\ell \beta = \sigma^2 r^\ell \bar{S}^{-1} \alpha \sum_{t=2}^T \tilde{M}_{tt-1} + o_p(1),$$

where $\bar{S} = \tilde{X}^\ell \tilde{X} + \alpha \alpha^\ell \sigma^2 (T - K)$ and $\tilde{M} = I - \tilde{X}(\tilde{X}^\ell \tilde{X})^{-1} \tilde{X}^\ell$.

Theorem 1 presents the leading term of the OLS bias under Assumption 1. Note that this formula uses the lower diagonal trace $\sum_t \tilde{M}_{tt-1}$ of the projection matrix orthogonal to \tilde{X} instead of $\sum_t M_{tt-1}$. Theorem 1 establish asymptotic negligibility of this difference.

It is worth discussing the size of the lower trace of the projection matrix as well as the size of the biases we may see in applications. Consider the term $\sum_t \tilde{M}_{tt-1} = \sum_t \tilde{P}_{tt-1}$, where $\tilde{P} = \tilde{X}(\tilde{X}^\ell \tilde{X})^{-1} \tilde{X}^\ell$. Note that this quantity is unchanged by any full rank rotation of the regressors. For simplicity, suppose that the regressors are generated by a stationary and ergodic process with $\frac{1}{T} \tilde{X}^\ell \tilde{X} = I_K$. Then

$$\sum_t \tilde{M}_{tt-1} = \frac{1}{T} \sum_t \tilde{X}_t^\ell \tilde{X}_{t-1} \quad \mathbb{E} \tilde{X}_t^\ell \tilde{X}_{t-1}.$$

Thus, the lower diagonal trace of \tilde{M} measures a linear connection between \tilde{X}_t and \tilde{X}_{t-1} . We can also notice that $\frac{1}{T} \sum_t \tilde{X}_t^\ell \tilde{X}_t = \mathbb{E} \tilde{X}_t^\ell \tilde{X}_t = K$, since we are dealing with K -dimensional vectors. In time series settings we tend to work with data where the regressors are autocorrelated. This way we should expect $\sum_t \tilde{M}_{tt-1} \approx \rho K$, where ρ is the average (over different regressors) scalar measure of regressor predictability.

We can notice that $\bar{S} = \tilde{X}^\ell \tilde{X}$ which leads to Assumption 1(iii) implying $r^\ell \bar{S}^{-1} \alpha = O_p(1/T)$. Let us momentarily assume $Tr^\ell \bar{S}^{-1} \alpha$ converges to a constant r_α . When K grows proportionally to the sample size while ρ remains bounded away from zero, the leading bias term becomes approximately $\sigma^2 r_\alpha \frac{\rho K}{T}$. This renders the OLS estimator for $r^\ell \beta$ inconsistent. Notably, the potential bias term is more pronounced for regressors with higher first-order autocorrelation and a larger number of regressors.

Although $\frac{1}{T} \sum_{t=2}^T \tilde{M}_{tt-1}$ (the lower diagonal trace of \tilde{M} over the sample size) is typically unavailable in practice as the strictly exogenous part of the regressors is unobservable, the

asymptotically equivalent lower diagonal trace of $M = I - X(X^\theta X)^{-1}X^\theta$ over the sample size can be calculated from available data. If this quantity is non-negligible in an application, even a small violation of strict exogeneity may result in substantial biases.

Lastly, the extent of bias depends on the alignment between the contrast r and the feedback direction α , making the feedback direction the most affected contrast direction. Contrasts in all directions orthogonal to α (using the scalar product weighted by \bar{S}^{-1}) experience only negligible bias. In our special case where only the first regressor is weakly exogenous, this observation corresponds to the OLS estimator of β_1 being the sole estimator with significant bias. Thus, for any linear contrast $r^\theta\beta$, the bias depends on the weight placed on β_1 . In real applications, the feedback direction is unknown and challenging to estimate empirically, given that α is a $K - 1$ vector.

Simulations We demonstrate the potential for OLS bias through a small scale simulation that varies the number of regressors and their short-term dependence. Data is generated following Assumption 1. The outcome vector is generated as $y = X\beta + \varepsilon$ with $\varepsilon \sim N(0, I)$ and $\beta = 0$. The design matrix is generated as $x_{1t} = \tilde{x}_{1t} + a\varepsilon_{t-1}$ and $X_{-1} = \tilde{X}_{-1}$, where \tilde{X} is generated as a rotated VAR(1) process with $\tilde{X}\tilde{X}^\theta/T = I_K$, independent from ε . Specifically, we generate $V_t = \rho V_{t-1} + u_t$ with $u_t \sim i.i.d. N(0, I_K)$ and define $\tilde{X} = V(V^\theta V/T)^{-1/2}$, where the square root comes from Cholesky decomposition. Across simulations, we fix the sample size at $T = 200$ and the coefficient on the feedback mechanism at $a = 1.5$. Simulation results are summarized in Figure 1 with the left panel showing results for the number of regressors K between 4 and 100 (fixing ρ at 0.8).² The right panel reports the results for the auto-correlation in regressors ρ between 0 and 0.98 (fixing K at 50). We report simulated values of absolute bias and standard deviation for the first coordinate of OLS together with the mean absolute value of the ratio of the lower trace of M to the sample size. Additionally, we report the same summary statistics for the new estimator proposed in Section 3.

The results presented in Figure 1 suggest that the bias of the OLS estimator can easily surpass its standard deviation, leading to highly unreliable statistical inferences. For example, even a regression with just 20 regressors may exhibit an OLS bias comparable in magnitude to the standard deviation in a sample of size 200 — a very common setting in macroeconomic applications. The observable lower trace of M divided by the sample size

²All results are presented as sixth order polynomial fits to the actual results across K .

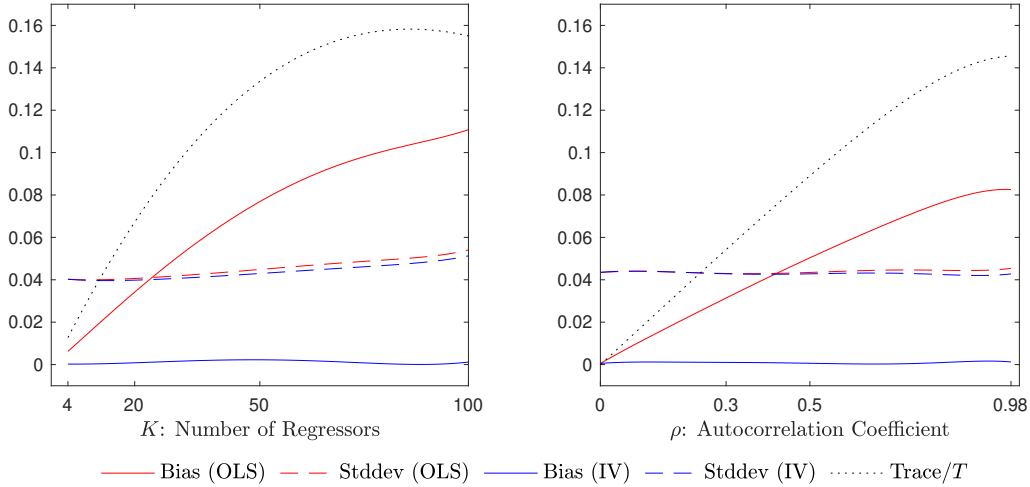


Figure 1: Absolute Bias and Standard Deviation of OLS and IV with $T = 200$

provides a highly predictive measure of the magnitude of the bias in the most affected direction. Notably, both the bias in the most affected direction and the lower diagonal trace tend to increase with the number of regressors and the first auto-correlation of the regressors.³

Figure 2 presents results for the same simulation design but with a sample size of $T = 800$, while the number of regressors varies from 16 to 400. Here, we observe that the bias reaches the same level as in Figure 1 when the number of regressors is the same fraction of the sample size, while the standard deviations drops two-fold. This demonstrates the inconsistency of the OLS for the worst direction when the number of regressors K grows proportionally to T . In essence, the estimator concentrates around an incorrect value as the sample size increases.

2.3 Variance estimator

In order to make reliable statistical inferences in a linear regression we usually need an estimator for the variance of the error term σ^2 . It is well known that in a regression with many regressors and strict exogeneity one has to adjust properly for the degrees of freedom in order to correct for over-fitting. The most commonly applied OLS estimator of the error variance uses this adjustment $\hat{\sigma}^2 = \frac{y'My}{T-K}$. It is not obvious ex ante whether this estimator

³One may be worried that the observed biases are due to persistence in the regressors as the coefficient ρ in an AR(1) process measures not only short-term dependence but also the long-run persistence. We re-run simulations generating $V_t = \rho u_{t-1} + u_t$ as an MA(1) process. The results are presented in Appendix B and are essentially identical to those reported in Figure 1.

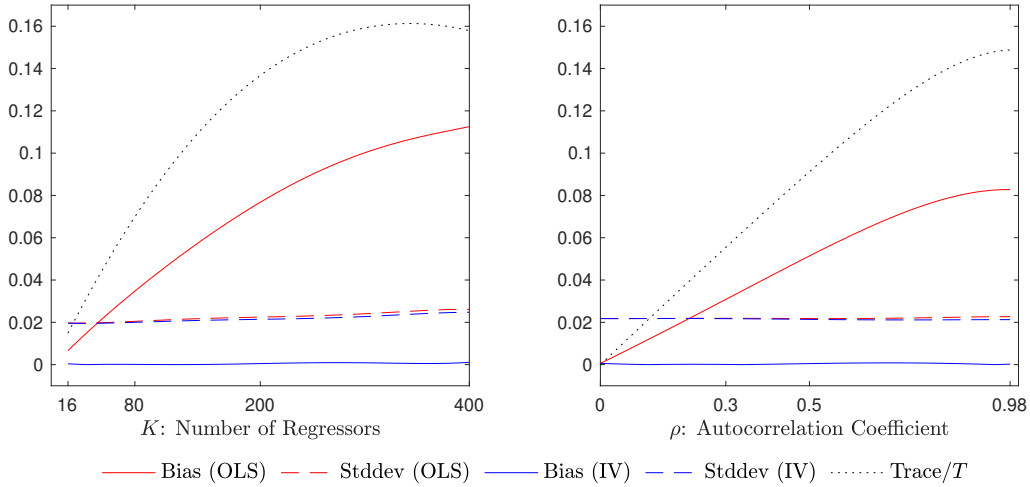


Figure 2: Absolute Bias and Standard Deviation of OLS and IV with $T = 800$

retains consistency with many regressors and only weak exogeneity. On the one hand, a large bias of the OLS estimator for coefficients raises concerns about the consistency of the variance estimator. On the other hand, the bias arises only in the direction of the feedback but not in all directions orthogonal to the feedback, which are numerous. Theorem 2 answers the question of the consistency of the OLS variance estimator.

Theorem 2 (Inconsistency of OLS variance). *Suppose Assumption 1 holds. Then,*

$$\frac{\hat{\sigma}^2}{\sigma^2} = 1 - \frac{\sigma^2 \alpha' \bar{S}^{-1} \alpha}{T K} \left(\sum_{t=2}^T \tilde{M}_{tt-1} \right)^2 + o_p(1).$$

Theorem 2 states that the OLS variance of the error estimate tends to be biased downward and reports an overly optimistic measure of fit in a setting with one-period feedback and many regressors. The size of the bias can be judged by the trace of the lower diagonal of the projection matrix M . When the number and predictability of the regressors are high enough to imply that the OLS estimator of the linear contrast in the worst direction (the feedback direction) is inconsistent, the OLS estimator of variance is inconsistent as well. Despite the inconsistency result, the biases for the OLS variance we observe in simulations tend to be relatively minor.

3 A consistent IV estimator

3.1 The idea of the proposed estimator

Let us introduce a $T \times T$ shift matrix D (or lag operator matrix) that shifts the time series index back by one; its only non-zero elements are $D_{t,t-1} = 1$ for all t . The transpose D^θ is the lead operator, moving the time-series index forward by one. Notice that the lower trace appearing in the bias of OLS can be written as $\sum_{t=2}^T \tilde{M}_{tt-1} = \text{tr}(D^\theta \tilde{M})$. Consider also a $T \times T$ matrix Γ measurable with respect to \tilde{X} such that $k\Gamma k = 1 - c$.

We propose an IV-inspired estimator that relies on an endogenous instrument $Z = (I - \Gamma^\theta)X$:

$$\hat{\beta}^{\text{IV}}(\Gamma) = (Z^\theta X)^{-1} Z^\theta y = (X^\theta (I - \Gamma)X)^{-1} X^\theta (I - \Gamma)Y.$$

The key insight behind the estimator is that we use a deliberately invalid instrument created in a way so the ‘invalid instrument’ bias offsets the bias in OLS. In order to eliminate the asymptotic bias we require that Γ solves the following (non-linear) one-dimensional equation:

$$\text{tr}[D^\theta (I - \Gamma) \tilde{M}_\Gamma] = 0 \tag{1}$$

where $\tilde{M}_\Gamma = I - \tilde{X}(\tilde{X}^\theta (I - \Gamma)\tilde{X})^{-1} \tilde{X}^\theta (I - \Gamma) = I - \tilde{X}(\tilde{Z}^\theta \tilde{X})^{-1} \tilde{Z}^\theta$ is an oblique projection off \tilde{X} in the direction of $\tilde{Z} = (I - \Gamma^\theta)\tilde{X}$. Further below, we establish that the new estimator $\hat{\beta}^{\text{IV}}(\Gamma)$ is a properly centered $\rho_{\bar{T}}$ asymptotically Gaussian (conditionally on \tilde{X}) estimator under Assumption 1. The main theoretical result about bias is obtained for a general Γ with the leading example of $\Gamma = \gamma D$ for $j\gamma j < 1 - c$.

Let us give some intuition of why this approach would work. Consider again the special case when all regressors but the first are strictly exogenous. Namely $x_{1t} = \tilde{x}_{1t} + a\varepsilon_{t-1}$, while $X_{-1} = \tilde{X}_{-1}$ and consider the first coefficient β_1 only, as it is the most biased direction. Consider also the case of $\Gamma = \gamma D$, thus $z_t = x_t - \gamma x_{t+1}$. In this setting all but the first instruments are strictly exogenous, while the first instrument $z_{1t} = \tilde{z}_{1t} + \varepsilon_{t-1} - \gamma\varepsilon_t$ is endogenous as it correlates with the contemporaneous regression error. The direction of the instrument endogeneity in general coincides with the feedback direction, as transformation Γ mixes up timing of the observations but preserves the feedback direction.

Oblique projections underlie the geometry of IV estimation similarly to how orthogonal projections explain the geometry of OLS. An oblique projection is defined as $M_{Z,X} = I - X(Z^\ell X)^{-1}Z^\ell$, where X and Z are of the same dimension and $Z^\ell X$ is invertible. Oblique projections satisfy idempotency, $M_{Z,X}^2 = M_{Z,X}$, but not symmetry, $M_{Z,X}^\ell \notin M_{Z,X}$. One can easily show that the Frisch-Waugh-Lowell theorem holds for oblique projections as well. Specifically, let $Z = [Z_1; Z_{-1}]$ where the dimensions of the corresponding Z 's and X 's coincide and all proper matrices are invertible. Then

$$\hat{\beta}_1^{\text{IV}}(\Gamma) = \frac{Z_1^\ell M_{Z_{-1}, X_{-1}} Y}{Z_1^\ell M_{Z_{-1}, X_{-1}} X_1} \quad \text{and} \quad \hat{\beta}_1^{\text{IV}}(\Gamma) - \beta_1 = \frac{Z_1^\ell M_{Z_{-1}, X_{-1}} \varepsilon}{Z_1^\ell M_{Z_{-1}, X_{-1}} X_1}.$$

Notice that since X_{-1} is strictly exogenous, Z_{-1} is strictly exogenous as well. If we condition on \tilde{X} we may then treat $M_{Z_{-1}, X_{-1}}$ as fixed. We look at the numerator:

$$\mathbb{E}[Z_1^\ell M_{Z_{-1}, X_{-1}} \varepsilon] = \mathbb{E}[\tilde{Z}_1^\ell M_{Z_{-1}, X_{-1}} \varepsilon] + a \mathbb{E}[\sum_{s,t} M_{st} (\varepsilon_{s-1} - \gamma \varepsilon_s) \varepsilon_t] = a \sigma^2 \sum_t (M_{tt-1} - \gamma M_{tt}).$$

Here we used $M_{Z_{-1}, X_{-1}} = (M_{st})$ for shortness of notation. Our goal is to choose γ in a manner that renders the last sum equal to zero. This should generally be feasible, given that the diagonal elements of projection matrices tend to dominate those on the lower diagonal. In the proof, we show that changing $M_{Z_{-1}, X_{-1}}$ to \tilde{M}_Γ in the bias expression introduces only an asymptotically negligible difference. Thus, the expectation of the numerator is asymptotically equivalent to $a \sigma^2 \text{tr}[D^\ell (I - \Gamma) \tilde{M}_\Gamma]$. By selecting Γ to solve equation (1), we achieve an asymptotically unbiased estimator.

3.2 Consistency of estimator

Let γ_0 denote the solution to (1) among matrices $\Gamma_0 = \gamma_0 D$. Since we only observe regressors X , knowing \tilde{X} is equivalent to knowing the feedback direction α . This makes solving equation (1) infeasible in practice. Let $\hat{\gamma}$ and $\hat{\Gamma} = \hat{\gamma} D$ be the solution to the empirically feasible equation:

$$\text{tr}[D^\ell (I - \hat{\Gamma}) M_{\hat{\Gamma}}] = 0 \quad \text{where} \quad M_{\hat{\Gamma}} = I - X(X^\ell (I - \hat{\Gamma}) X)^{-1} X^\ell (I - \hat{\Gamma}). \quad (2)$$

The following Lemma gives sufficient conditions for the existence and uniqueness of a solution to equation (1).

Lemma 1. *Suppose that \tilde{X} is a $T \times K$ matrix of rank K and $\Gamma = \gamma D$.*

(i) *If $K < T/5$, then there exists a unique $\gamma \in [1/2, 1/2]$ solving equation (1).*

(ii) *If $j\text{tr}(D^\theta \tilde{M})j < \mu^2 K$ and $K < T/[1 + (1 + \mu)^2]$ for some $\mu \in [0, 1]$, then there exists a unique γ solving equation (1) such that $j\gamma j < \mu/(1 + \mu)$.*

One may state an analog of Lemma 1 for equation (2) as well due to their analog structure. According to Lemma 1, equation (1) typically can be solved in a one-dimensional family of transformations $\Gamma = \gamma D$. It is possible to search for a solution in other classes of matrices; we leave the question of finding an optimal class of transformations Γ to future research. The proof of Lemma 1 reveals that though equation (1) is non-linear, its solution can be found relatively fast as a fixed point of a contraction.

One may think that $K < T/5$ in Part (i) of Lemma 1 is a restrictive condition. It arises because we do not put any restriction on \tilde{X} other than full rank. This accommodates an extensive range of idiosyncrasies in the generation of regressors, surpassing those encountered in typical macroeconomic time series data sets. Placing some restrictions on regressors may weaken the restriction on the sample size significantly. For example, putting a bound on the average auto-correlation of the regressors, μ , eases this restriction considerably as shown in Part (ii) of Lemma 1. We also notice that if the original potential for bias is small, then the γ that removes the bias is small as well. Specifically, if $\text{tr}[D^\theta \tilde{M}] = 0$ so the original OLS has no (asymptotic) bias, our estimator defaults back to OLS.

The following theorem establishes the asymptotic bias of IV for a generic choice of Γ and shows consistency for the specific choice of $\hat{\Gamma} = \hat{\gamma} D$.

Theorem 3. *Suppose Assumption 1 holds.*

(i) *If Γ is \tilde{X} -measurable and $k\Gamma k < 1 - c$, then*

$$r^\theta \hat{\beta}^{\text{IV}}(\Gamma) - r^\theta \beta = \sigma^2 r^\theta \bar{S}_\Gamma^{-1} \alpha \text{tr}[D^\theta (I - \Gamma) \tilde{M}_\Gamma] + o_p(1),$$

where $\bar{S}_\Gamma = \tilde{X}^\theta (I - \Gamma) \tilde{X} + \sigma^2 \alpha \alpha^\theta \text{tr}[(I - \Gamma) \tilde{M}_\Gamma]$. Specifically, it follows that $r^\theta \hat{\beta}^{\text{IV}}(\Gamma)$ is consistent for $r^\theta \beta$ when Γ solves equation (1).

(ii) If $K < T/5$, then $\hat{\gamma} - \gamma_0 = O_p(\frac{1}{T})$, and $r^\theta \hat{\beta}^{\text{IV}}(\hat{\Gamma}) - r^\theta \hat{\beta}^{\text{IV}}(\Gamma_0) = o_p(\frac{1}{T})$.

The result of Theorem 1 is a special case of Part (i) of Theorem 3 for $\Gamma = 0$. Using Γ that solves equation (1) produces a consistent estimator for any reasonable contrast, and the proposed solution does not depend on the contrast of interest. The ideal Γ_0 solving equation (1) can be random but is strictly exogenous, as it depends only on the strictly exogenous part of the regressors \tilde{X} . Solution $\hat{\Gamma}$ to equation (2) is random and depends on error term ε . However, as stated in part (ii) of Theorem 3 the feasible estimator $r^\theta \hat{\beta}^{\text{IV}}(\hat{\Gamma})$ is consistent and has the same asymptotic distribution as the infeasible one using the ideal Γ_0 .

An appealing feature of our proposed solution is that our estimator is linear in the outcome variable. Our proposal eliminates bias by working only with the regressors. The same bias correction would work for any outcome variables that satisfy the weak exogeneity assumption with one-period feedback to regressors. Our solution does not require one to estimate the direction of feedback α .

Figure 1 shows that in the simulations the proposed IV estimator fixes the bias of OLS with essentially no increase in the standard deviation of the estimator.

3.3 Consistency of variance estimator

Let us introduce $T - K_\Gamma = \text{tr}[(I - \Gamma)\tilde{M}_\Gamma]$ and define an estimator for variance $\hat{\sigma}^2(\Gamma) = \frac{y^\theta(I - \Gamma)M_\Gamma y}{T - K_\Gamma}$ for some matrix Γ . Using equation (4) in the Appendix, we find that the value of $\hat{\sigma}^2(\Gamma)$ is invariant to the value of β and $\hat{\sigma}^2(\Gamma) = \frac{\varepsilon^\theta(I - \Gamma)M_\Gamma \varepsilon}{T - K_\Gamma}$. Part (ii) of Lemma 3 from the Appendix implies that $\hat{\sigma}^2(\Gamma)$ possesses a very desirable property for a variance estimator, namely, it is non-negative for any realization of the data. If Γ solves equation (1) then $\hat{\sigma}^2(\Gamma)$ is consistent for σ^2 .

Theorem 4. *Suppose Assumption 1 holds, Γ is \tilde{X} -measurable, and $k|\Gamma|k < 1 - c$. Then*

$$\frac{\hat{\sigma}^2(\Gamma)}{\sigma^2} = 1 + \frac{\sigma^2 \alpha^\theta \bar{S}_\Gamma^{-1} \alpha}{T - K_\Gamma} \text{tr}[D^\theta(I - \Gamma)\tilde{M}_\Gamma] \text{tr}[D(I - \Gamma)\tilde{M}_\Gamma] + o_p(1).$$

Speci cally, it follows that $\hat{\sigma}^2(\Gamma)$ is consistent for σ^2 when Γ is such that equation (1) holds.

As with Theorem 3, the theoretically desirable Γ that solves (1) is not known, but one can search for an empirically feasible value $\hat{\Gamma}$ that solves (2). One can easily use Part (ii) of Theorem 3 to extend the argument of Theorem 4 and also show that $\hat{\sigma}^2(\hat{\Gamma})$ is consistent.

4 Inference

In this section we show that under some additional assumptions the proposed IV estimator is asymptotically Gaussian conditionally on \tilde{X} . We also suggest standard errors that can be paired with our estimator in order to achieve asymptotically valid inferences, that is, to construct confidence sets and/or to produce reliable t-tests.

There are two theoretical and practical challenges we encounter. The first is a need to correctly account for the asymptotic importance of a quadratic form. Specifically, the bias of OLS arises due to presence of a quadratic form in errors, which has a non-trivial mean. We construct our estimator to guarantee a zero mean of its corresponding quadratic form and thus eliminate the bias asymptotically. The quadratic form in error terms with zero mean is asymptotically Gaussian when its rank is growing to infinity under a condition of eigenvalue negligibility. We refer the reader to [Anatolyev \(2019\)](#); [Chao et al. \(2012\)](#); [Sølvsten \(2020\)](#); [Kline et al. \(2020\)](#) for examples of the Central Limit Theorems for quadratic forms. The naive standard errors tend to improperly account for asymptotic uncertainty of a quadratic form. The asymptotic importance of such quadratic forms has appeared previously in the literature on linear models with many instruments and/or many regressors. For example, [Hansen et al. \(2008\)](#) show the importance to adjust standard errors for presence of a quadratic form in many instrument settings. See also [Anatolyev \(2019\)](#) for a comprehensive survey of the issue. The second challenge is that the quadratic form we end up with has non-zero diagonal elements that are on average zero. This makes the asymptotic variance depend on skewness and kurtosis of errors that are hard to estimate. This is similar to the issue of proper inference for LIML-type estimators with many instruments (as in [Hansen et al. \(2008\)](#)). Below we consider two instances when we can handle the first challenge with relative ease and ignore the second challenge.

4.1 Inference when $K/T \rightarrow 0$

Theorem 5. *Suppose Assumption 1 holds, $\max_t \left\| (\tilde{X}^{\top} \tilde{X})^{-1/2} \tilde{x}_t \right\| = o_p(1)$. Then, as $T \rightarrow \infty$,*

$$\frac{r^{\top} \hat{\beta}^{\text{IV}}(\hat{\Gamma}) - r^{\top} \beta}{\sqrt{\hat{\Sigma}_T}} \rightarrow N(0, 1)$$

where $\hat{\Sigma}_T = \hat{\sigma}^2(\hat{\Gamma})kr^\theta(X^\theta(I - \hat{\Gamma})X)^{-1}X^\theta(I - \hat{\Gamma})k^2$.

Condition $\max_t \left\| (\tilde{X}^\theta \tilde{X})^{-1/2} \tilde{x}_t \right\| = o_p(1)$ is a relatively standard negligibility condition often invoked in order to obtain asymptotic Gaussianity of OLS using Lindeberg CLT (see, e.g., [Koenker and Machado, 1999](#), page 334). It implies among other things that the maximal diagonal element \tilde{P}_{tt} of the projection matrix \tilde{P} is asymptotically negligible. At the same time, the average of these diagonal elements is equal to K/T , and thus, this condition can hold only when we have a moderately large number of regressors ($K/T \rightarrow 0$). When the number of regressors is moderately large, both challenges described at the beginning of this section are asymptotically negligible. The standard errors suggested by Theorem 5 look like the usual IV-type standard errors with one important change, they use the newly proposed estimator of variance $\hat{\sigma}^2(\hat{\Gamma})$.

4.2 Inference with Gaussian errors

Theorem 6. *Suppose Assumption 1 holds and ε_1 is Gaussian conditionally on \tilde{X} . Assume that Γ solves equation (1) and $k\Gamma k < 1 - c$. Then, as $T \rightarrow \infty$,*

$$\frac{r^\theta \hat{\beta}^{\text{IV}}(\Gamma) - r^\theta \beta}{\sqrt{\hat{\Sigma}_T}} \rightarrow N(0, 1)$$

where $\hat{\Sigma}_T$ is measurable with respect to \tilde{X} . With probability asymptotically approaching one, $\hat{\Sigma}_T = (1 + \psi)\hat{\sigma}^2(\Gamma)kr^\theta(X^\theta(I - \Gamma)X)^{-1}X^\theta(I - \Gamma)k^2$, where $\psi = \frac{\text{tr}(B^2)}{\text{tr}(B^\theta B)}$ and $B = D^\theta(I - \Gamma)\tilde{M}_\Gamma$.

Theorem 6 allows the number of regressors to grow proportionally with the sample size. The assumption that errors are Gaussian is used in several ways. First, it allows us to avoid imposing the Lindeberg-type negligibility condition, as any linear combination of errors with weights depending on \tilde{X} is conditionally Gaussian in finite-samples. Secondly, it removes the need to estimate skewness and kurtosis.

However, in the case when K grows proportionally to the sample size, the standard errors stated in Theorem 5 do not correctly account for the uncertainty induced by the quadratic form. The missing term in the asymptotic variance formula depends on the importance of the feedback and is hard to estimate. We instead provide an upper bound on the asymptotic variance that makes the confidence sets asymptotically valid but conservative. In our simu-

lation we noticed that the correction ψ tends to be tiny and does not change the standard errors much. We advise the calculation of it as a robustness check.

5 Extension to multiple periods

In the previous section we maintained the assumption that the violation of strict exogeneity happens for one period only. Some results can be extended to feedback lasting a fixed finite number of periods.

Assumption 2. (i) *The observed regressors x_t can be decomposed as $x_t = \tilde{x}_t + \sum_{\ell=1}^L \alpha_\ell \varepsilon_{t-\ell}$ where the $T \times K$ matrix $\tilde{X} = [\tilde{x}_1, \dots, \tilde{x}_T]'$ has full rank.*

(ii) *The errors ε_t are i.i.d. conditional on \tilde{X} , $E[\varepsilon_t | \tilde{X}] = 0$, $c < \sigma^2 := E[\varepsilon_t^2] < C$, and $E[\varepsilon_t^4] < C$ almost surely.*

(iii) *The non-random vectors $\alpha_1, \dots, \alpha_L, r \in \mathbb{R}^K$ satisfy $\alpha_\ell' (\frac{1}{T} \tilde{X}' \tilde{X})^{-1} \alpha_\ell = O_p(1)$ for $\ell = 1, \dots, L$ and $r' (\frac{1}{T} \tilde{X}' \tilde{X})^{-1} r = O_p(1)$.*

(iv) *L is fixed and $K/T < 1 - c$.*

Vectors α_ℓ describe how the shock to the outcome variable affects the regressors ℓ periods later. The direction of the feedback may vary freely with the lag as the regressors differ in the speed of reaction/adjustment. However, we probably should expect that the size of the ℓ -period after feedback measured as $\alpha_\ell' (\frac{1}{T} \tilde{X}' \tilde{X})^{-1} \alpha_\ell$ should become negligible for large enough ℓ in typical macroeconomic settings.

Theorem 7. *Suppose Assumption 2 holds, Γ is \tilde{X} -measurable, and $k\Gamma k < 1 - c$. Then,*

$$r' \hat{\beta}^{\text{IV}}(\Gamma) - r' \beta = \sigma^2 \sum_{\ell=1}^L r' \bar{S}_\Gamma^{-1} \alpha_\ell \text{tr}[(D^\ell)^\ell (I - \Gamma) \tilde{M}_\Gamma] + o_p(1), \quad (3)$$

where $\bar{S}_\Gamma = \tilde{X}' (I - \Gamma) \tilde{X} + \sigma^2 \sum_{j,\ell=1}^L \alpha_j \alpha_\ell' \text{tr}[(D^\ell)^j (I - \Gamma) \tilde{M}_\Gamma D^\ell]$.

Theorem 7 is a direct generalization of Part (i) of Theorem 3. A special case of $\Gamma = 0$ shows that the bias of the OLS is a linear combination of L terms connected to the lower diagonal traces of the projection matrix \tilde{M} . Lower diagonal traces, $\text{tr}[(D^\ell)^\ell \tilde{M}]$, correspond

to average measures of the regressors’ auto-correlations of order ℓ , and are expected to decay for stationary regressors. Combined with the expected decrease in the size of α_ℓ , the first few terms in the bias formula should capture most of the bias in stationary applications.

The bias formula (3) is derived for any IV-type estimator for an \tilde{X} -measurable $T \times T$ matrix Γ with $\|\Gamma\| < 1 - c$. Theorem 7 suggests that if Γ solves a system of L equations $\text{tr}((D^\theta)^\ell(I - \Gamma)\tilde{M}_\Gamma) = 0$ for $\ell = 1, \dots, L$ then $r^\theta \hat{\beta}^{\text{IV}}(\Gamma)$ is a consistent estimator. A natural suggestion is to search for Γ in the class of matrices $\Gamma = \sum_{\ell=1}^L \gamma_\ell D^\ell$ choosing L parameters γ_ℓ in such a way as to solve the system of equations. We leave the question of considering other classes of matrices as well as getting some guarantees for finding a solution to future research.

6 Simulations

The goal of this section is to assess the size of the OLS bias in a ‘typical’ regression using US macroeconomic data. We use the data set from [Stock and Watson \(2016\)](#) containing quarterly observations from 1964 to 2013 ($T = 200$) on 108 US macro indicators. This data set is largely similar to the [McCracken and Ng \(2020\)](#) FRED-QD data set. The data set includes a broad class of variables with diverse time series properties.

Many macro and financial indicators tend to be very persistent and may be integrated up to the second order (have one or two unit roots). A prevailing (but not uniformly accepted) practice is to transform all variables to stationary before running a regression in order to avoid issues of co-integration and near co-integration or biases related to persistent regressors (see [Stambaugh \(1999\)](#) on such biases). The way applied researchers transform variables to stationary or make decisions about variable stationarity varies widely across the literature with many such decisions based both on statistical tests and expert judgements. Given that a large fraction of regressions is aimed at business-cycle parameters and in order to not take a stand and to unify the pre-treatment of variables we apply [Hamilton \(2018\)](#) transformation to all variables in the data set. Specifically, for each variable we define its cyclical component to be a two-year-ahead forecast error to this variable based on a univariate AR(4) regression. According to [Hamilton \(2018\)](#), this filtering transforms all types of stationary and up to second order integrated variables into stationary ones and extracts their business cycle component.

For each K in the set $\{5, 15, 25, \dots, 85, 100\}$ we perform 100 experiments where we randomly draw K distinct variables from the transformed data set, denote them X_r and an additional variable y_r . We calculate through simulations, the biases and standard deviations of the OLS and of our proposed estimator (referred to as IV) for the linear contrast in the feedback direction in the regression of y_r on X_r under the assumption of one-period violation of strict exogeneity. For this, we simulate $N = 1000$ samples from a data generating process satisfying Assumption 1 that preserves the time series behavior of regressors and the feedback size/direction of the observed (y_r, X_r) . Specifically, we use as true parameters the empirical OLS values $\beta = (X_r^\ell X_r)^\ell^{-1} X_r^\ell y_r$, $\sigma^2 = e^\ell e / (T - K)$ for $e = y_r - X_r \beta$, and $\alpha = X_r^\ell D^\ell e / (e^\ell e)$. We simulate samples as $X = X_r + D^\ell \varepsilon \alpha^\ell$ and $y = X \beta + \varepsilon$ where $\varepsilon \sim N(0, \sigma^2 I)$. We calculate the bias and standard deviation for the OLS and IV estimators of the linear contrast with $r = \alpha$.

The left panel of Figure 3 depicts the results of the experiments (for different K) that fall into the 10th percentile of the OLS bias. For those experiments, we report the OLS bias and standard deviation and the IV bias and standard deviation in the feedback direction alone along with the normalized lower trace of $M_r = I - X_r (X_r^\ell X_r)^\ell^{-1} X_r^\ell$, that is $\text{tr}(D^\ell M_r) / T$. The right panel contains the results of the experiments that fall into the 90th percentile of the bias.

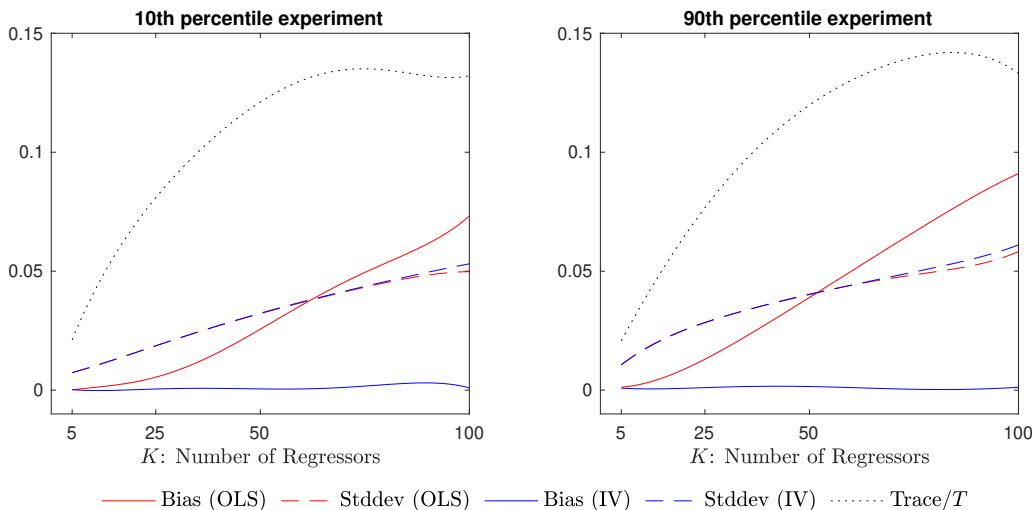


Figure 3: Absolute Bias and Standard Deviation of OLS and IV with $T = 200$

Figure 3 shows that a typical macroeconomic data set demonstrates enough time series

dependence or short-term linear predictability that it creates a potential for significant OLS biases. A typical size/direction of one-period feedback for a randomly picked-up regression using macro indicators is such that for a sample with 200 time periods in a regression with 25 regressors the OLS bias in the feedback direction equals about half of the standard deviation and is approximately equal to the standard deviation when the number of regressors is 50. This leads to invalid statistical inferences when relying on the OLS. Figure 4 reports the size distortions in the experiments described above for the 5% tests about the linear contrast in the direction of the feedback. Our new proposed estimator completely corrects the bias without any significant change to the standard deviation and restores the correct size for statistical tests.

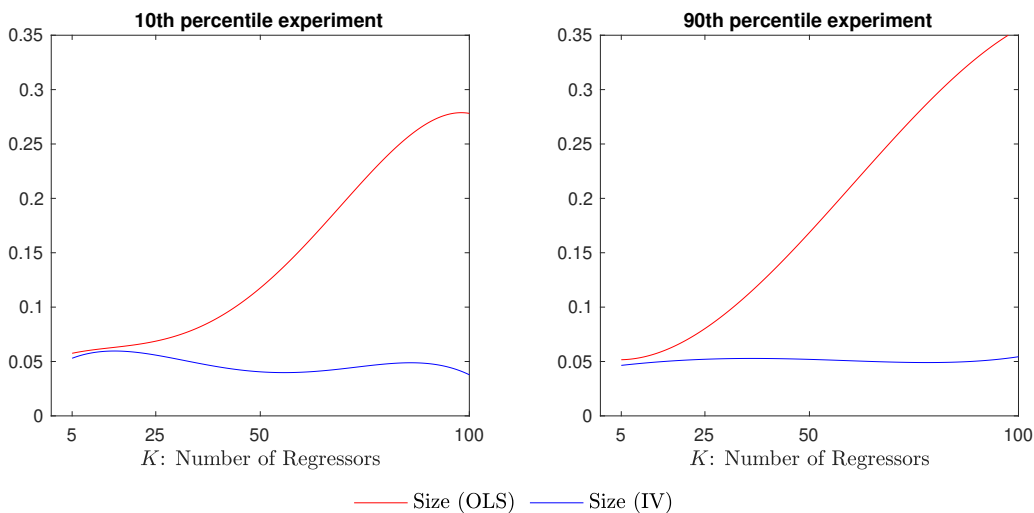


Figure 4: Size of Nominal 5% two-sided tests using OLS and IV with $T = 200$

Another observation from Figure 3 is that the ratio of the lower trace of the regressor projection matrix to the sample size is highly indicative of the size of the worst bias. Applied researchers should be worried when this indicator exceeds 5–10%. It is worth pointing out that in all of our experiments and simulations we encountered no problems of finding the solution to equation (2), which supports our assertion that the sample size requirement imposed in Lemma 1 is sufficient but not necessary for the existence of the solution.

Finally, we note that the results presented in Figure 3 are both qualitatively and quantitatively similar to the left panels of Figures 1 and 6. The difference between figures is that in Figures 1 and 6 we simulated artificial regressors following a simple AR(1) or MA(1) process.

Our theoretical results are quite agnostic about the time series properties of the regressors, specifically, we make no assumptions about stationarity or origin of the regressors. The data generating processes underlying Figure 3 mimics the time series behavior and feedback size/direction of a 'typical' macroeconomic application to the best of our ability.

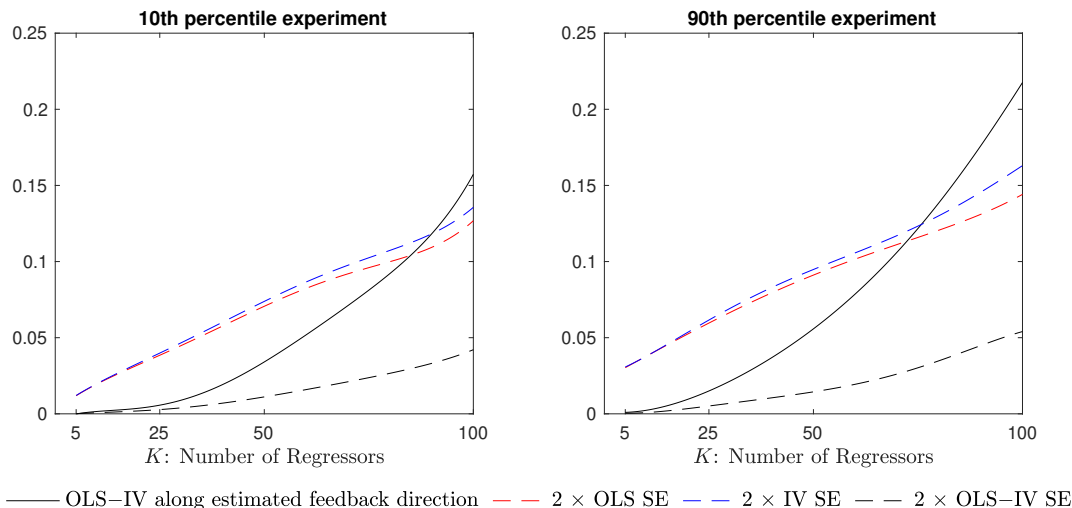


Figure 5: Difference between IV and OLS along estimated feedback direction, $T = 200$

Figure 5 answers the question of how different the results of the OLS and our new proposed estimators are in a randomly selected regression based on typical macroeconomic data. As in the experiments described above we randomly select (y_t, X_t) but rather than evaluate the theoretical bias in simulations we calculate the realization of the difference between two estimators for the contrast in the estimated feedback direction. We report the experiment corresponding to 10th percentile of the absolute difference on the left panel and 90th percentile on the right. For the described experiments we report the absolute value of the difference between the OLS and the IV estimators, double of the OLS and IV standard errors as well as the double standard deviations of the difference. We report doubled standard errors to directly relate the results to the corresponding t -statistics.

As we see from Figure 5, the difference between the OLS and IV estimators is statistically significant in all experiments we report. This indicates that the observed difference between estimators are mostly due to the biases and cannot be explained by randomness in realizations. An alternative explanation is that a hypothesis of no feedback is rejected in almost all regressions we considered. This is a sign of prevalence of feedback mechanisms in

most macroeconomic applications. Comparing the difference between the two estimators to the doubled standard errors shows that when the number of regressors is above 70 the IV estimate falls outside of the standard OLS confidence set and the OLS estimate falls outside the IV confidence set. This demonstrates the potential for drastic disagreement between the two estimation methods. Finally, we notice that the standard deviation of the OLS versus IV difference is much smaller than the standard error of either the OLS or the IV estimators, in some cases more than five times smaller. This shows that stochastic deviations in two estimators are greatly aligned, and that most of the difference between the two estimators comes from the bias, not the variance.

Table 1: Statistical significance of the differences in OLS and IV coefficients

K		5	15	25	35	45	55	65	75	85	100
100	ave($t_{\Delta} > 1.96$)	21.20	13.93	11.72	11.03	10.47	8.93	8.12	7.64	6.91	6.41

NOTE: $t_{\Delta} = j\hat{\beta}^{\text{LS}} - \hat{\beta}^{\text{IV}} / se(j\hat{\beta}^{\text{LS}} - \hat{\beta}^{\text{IV}})$. The average is taken over all coefficients in 100 randomly chosen models for each value of K .

While Figure 5 reports the difference between the two estimators in the most affected direction, one may also ask how different coefficients on individual regressors are. The average bias of individual coefficients is a counter-play of two forces. On one side, a larger number of regressors leads to a larger lower diagonal trace $\text{tr}(D^{\ell}M_r)/T$ and thus to a larger bias in the most affected direction. At the same time, when the dimensionality of regressors is large, (a randomly selected) feedback direction is on average less aligned with any coordinate direction. Thus, the same size of the worst bias results in a smaller average individual coefficient bias when the number of regressors increases, as it spreads out among many individual coefficients.

In Table 1, we report the average fraction of coefficients that display a statistically significant difference between the OLS and the IV estimators. The average is taken over the K coefficients in the regression and over the 100 random regressions with K regressors that we draw from the macroeconomic data base described above. While the fraction of regressors with a statistically significant difference between the OLS and the IV estimators is declining with larger K , the absolute number of such regressors increases. We conclude that while most directions/coefficients are immune to the biases, a non-trivial fraction of coefficients are significantly affected.

References

- Anatolyev, S. (2019). Many instruments and/or regressors: A friendly guide. *Journal of Economic Surveys* 33(2), 689–726.
- Arellano, M. and S. Bond (1991). Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations. *The Review of Economic Studies* 58(2), 277–297.
- Bao, Y. and A. Ullah (2007). The second-order bias and mean squared error of estimators in time-series models. *Journal of Econometrics* 140(2), 650–669.
- Bauwens, L., G. Chevillon, and S. Laurent (2023). We modeled long memory with just one lag! *Journal of Econometrics* 236(1), 105467.
- Chamberlain, G. (1982). The general equivalence of Granger and Sims causality. *Econometrica: Journal of the Econometric Society* 50(3), 569–581.
- Chao, J. C., N. R. Swanson, J. A. Hausman, W. K. Newey, and T. Woutersen (2012). Asymptotic distribution of JIVE in a heteroskedastic IV regression with many instruments. *Econometric Theory* 28(1), 42–86.
- Chevillon, G., A. Hecq, and S. Laurent (2018). Generating univariate fractional integration within a large VAR(1). *Journal of Econometrics* 204(1), 54–65.
- Engle, R. F., D. F. Hendry, and J.-F. Richard (1983). Exogeneity. *Econometrica: Journal of the Econometric Society* 51(2), 277–304.
- Granger, C. W. (1969). Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: Journal of the Econometric Society* 37(3), 424–438.
- Gupta, A. and M. H. Seo (2023). Robust inference on infinite and growing dimensional time series regression. *Econometrica: Journal of the Econometric Society*, 1333–1361.
- Hamilton, J. D. (1994). *Time series analysis*. Princeton University Press.
- Hamilton, J. D. (2018). Why you should never use the Hodrick-Prescott filter. *Review of Economics and Statistics* 100(5), 831–843.
- Hansen, B. E. and K. D. West (2002). Generalized method of moments and macroeconomics. *Journal of Business & Economic Statistics* 20(4), 460–469.
- Hansen, C., J. Hausman, and W. Newey (2008). Estimation with many instrumental variables. *Journal of Business & Economic Statistics* 26(4), 398–422.
- Jordà, Ò. (2005). Estimation and inference of impulse responses by local projections. *Amer-*

- ican Economic Review* 95(1), 161–182.
- Kline, P., R. Saggio, and M. S¸olvsten (2020). Leave-out estimation of variance components. *Econometrica: Journal of the Econometric Society* 88(5), 1859–1898.
- Koenker, R. and J. A. Machado (1999). GMM inference when the number of moment conditions is large. *Journal of Econometrics* 93(2), 327–344.
- McCracken, M. W. and S. Ng (2020). FRED-QD: A quarterly database for macroeconomic research. *Federal Reserve Bank of St. Louis* (Working Paper 2020-005).
- Montiel Olea, J. L. and M. Plagborg-M¸oller (2021). Local projection inference is simpler and more robust than you think. *Econometrica: Journal of the Econometric Society* 89(4), 1789–1823.
- Nickell, S. (1981). Biases in dynamic models with fixed effects. *Econometrica: Journal of the Econometric Society* 49(6), 1417–1426.
- Schennach, S. M. (2018). Long memory via networking. *Econometrica* 86(6), 2221–2248.
- Sims, C. A. (1972). Money, income, and causality. *The American Economic Review* 62(4), 540–552.
- S¸olvsten, M. (2020). Robust estimation with many instruments. *Journal of Econometrics* 214(2), 495–512.
- Stambaugh, R. F. (1999). Predictive regressions. *Journal of Financial Economics* 54(3), 375–421.
- Stock, J. H. and M. W. Watson (2016). Dynamic factor models, factor-augmented vector autoregressions, and structural vector autoregressions in macroeconomics. In *Handbook of Macroeconomics*, Volume 2, pp. 415–525. Elsevier.
- Stock, J. H. and M. W. Watson (2019). *Introduction to econometrics, 4th edition*. Pearson.
- Wallis, K. F. (1977). Multiple time series analysis and the final form of econometric models. *Econometrica: Journal of the Econometric Society* 45(6), 1481–1497.
- Zellner, A. and F. Palm (1974). Time series analysis and simultaneous equation econometric models. *Journal of Econometrics* 2(1), 17–54.

Appendix A Proofs

Notation For brevity, the expectation $E[\cdot]$ is used to denote the conditional expectation $E[\cdot|\tilde{X}]$. The proofs use some well known identities involving matrix traces. For matrices of compatible dimensions, $\text{tr}(ABD) = \text{tr}(BDA)$, $\text{tr}(A) = \text{tr}(A^\theta) = \text{tr}(A + A^\theta)/2$. As in the main paper, let $P = I - M = X(X^\theta X)^{-1}X^\theta$ and $P_\Gamma = I - M_\Gamma = X(X^\theta(I - \Gamma)X)^{-1}X^\theta(I - \Gamma)$.

A.1 Auxilliary lemmas

Lemma 2. (i) For a symmetric matrix A and positive semi-definite (psd) matrix B , we have the bounds $\lambda_{\min}(A) \text{tr}(B) \leq \text{tr}(AB) \leq \lambda_{\max}(A) \text{tr}(B)$.

(ii) For any square matrix A , let λ be the smallest eigenvalue of $\frac{A+A^\theta}{2}$. If $\lambda > 0$, then $\|A^{-1}\| \leq 1/\lambda$.

(iii) For any compatible matrices A and B , we have $|\text{tr}(AB)| \leq \|A\|_F \|B\|_F$.

Proof of Lemma 2. (i) As A is symmetric, there exists U with $U^\theta U = I$ such that $D = UAU^\theta$ is a diagonal matrix with eigenvalues of A along its diagonal. Let $F = UBU^\theta$. Note that $F_{ii} \geq 0$ since B is psd and $\text{tr}(B) = \text{tr}(F) = \sum_i F_{ii}$. Now,

$$\text{tr}(AB) = \text{tr}(UAU^\theta UBU^\theta) = \sum_i D_{ii} F_{ii} \leq \max_i D_{ii} \sum_i F_{ii} = \lambda_{\max}(A) \text{tr}(B),$$

and $\sum_i D_{ii} F_{ii} \geq \min_i D_{ii} \sum_i F_{ii} = \lambda_{\min}(A) \text{tr}(B)$. (ii) Now, $\lambda \|x\|^2 = x^\theta(A + A^\theta)x/2 = x^\theta Ax$. As $\lambda > 0$, it follows that $Ax \neq 0$ when $x \neq 0$, so A is invertible. For $x \neq 0$, we then have $\lambda = \|Ax\|/\|x\| = \|y\|/\|A^{-1}y\|$ where $y = Ax$. Thus, $\|A^{-1}\| = \sup_{y \neq 0} \|A^{-1}y\|/\|y\| = 1/\lambda$. (iii) $|\text{tr}(AB)| = \left| \sum_{ij} A_{ij} B_{ij} \right| \leq \sqrt{\sum_{ij} A_{ij}^2} \sqrt{\sum_{ij} B_{ij}^2} = \|A\|_F \|B\|_F$. \square

Lemma 3. Suppose $X \in \mathbb{R}^{T \times K}$ has full rank and $\Gamma \in \mathbb{R}^{T \times T}$ has $\|\Gamma\| < 1$. Then (i) the matrices $I - \Gamma$, $I - P\Gamma$, $X^\theta(I - \Gamma)X$ and $I + A_\Gamma M$ are invertible where $A_\Gamma = (I - \Gamma)^{-1}\Gamma$;

(ii) the following identities hold:

$$(I - \Gamma)M_\Gamma = M(I + A_\Gamma M)^{-1}, \quad (4)$$

$$P_\Gamma = (I - P\Gamma)^{-1}P(I - \Gamma), \quad (5)$$

$$M_\Gamma = (I - P\Gamma)^{-1}M, \quad (6)$$

$$\hat{\beta}^{\text{IV}}(\Gamma) = (X^\ell X)^{-1}X^\ell(I + A_\Gamma M)^{-1}y. \quad (7)$$

(iii) the matrix $(I - \Gamma)M_\Gamma + M_\Gamma^\ell(I - \Gamma^\ell)$ is positive semi-definite and $c < \frac{T}{T - K_\Gamma} < C$ where $T - K_\Gamma = \text{tr}[(I - \Gamma)M_\Gamma]$;

Proof of Lemma 3. (i) Lemma 2(ii), the triangle inequality, and $k\Gamma k < 1$ yields invertibility of $I - \Gamma$, $I - P\Gamma$, $X^\ell(I - \Gamma)X$, and $I - \Gamma P$. $I + A_\Gamma M$ is invertible since

$$I + A_\Gamma M = (I - \Gamma)^{-1}(I - \Gamma + \Gamma M) = (I - \Gamma)^{-1}(I - \Gamma P).$$

(ii) Next, we note that

$$\begin{aligned} (X^\ell X)^{-1}X^\ell &= (X^\ell X)^{-1}[X^\ell(I - \Gamma)X(X^\ell(I - \Gamma)X)^{-1}]X^\ell \\ &= (X^\ell X)^{-1}X^\ell(I - \Gamma)P_\Gamma(I - \Gamma)^{-1}. \end{aligned} \quad (8)$$

Pre-multiplying (8) by X and using $P_\Gamma = PP_\Gamma$ gives us $P = (I - P\Gamma)P_\Gamma(I - \Gamma)^{-1}$ and hence (5). Therefore, we also have (6):

$$M_\Gamma = I - P_\Gamma = (I - P\Gamma)^{-1}(I - P\Gamma - P(I - \Gamma)) = (I - P\Gamma)^{-1}M.$$

Using the ‘‘push-through’’ identity $(I - P\Gamma)^{-1}P = P(I - \Gamma P)^{-1}$ on (5) similarly yields (4):

$$M_\Gamma = I - P(I + A_\Gamma M)^{-1} = (I + A_\Gamma M)(I + A_\Gamma M)^{-1} = (I - \Gamma)^{-1}M(I + A_\Gamma M)^{-1}.$$

Reusing that $P_\Gamma = P(I + A_\Gamma M)^{-1}$ we get (7) from $\hat{\beta}^{\text{IV}}(\Gamma) = (X^\ell X)^{-1}X^\ell P_\Gamma y$ and $XP = X$.

(iii) Positive semi-definiteness comes from (4):

$$(I - \Gamma)M_\Gamma + M_\Gamma^\ell(I - \Gamma^\ell) = (I + MA_\Gamma^\ell)^{-1}M \{(I - \Gamma^\ell)^{-1} + (I - \Gamma)^{-1}\} M(I + A_\Gamma M)^{-1}$$

and observing that $(I - \Gamma^\theta)^{-1} + (I - \Gamma)^{-1}$ is psd. The rate condition follows from $\text{tr}(M) = T - K$, Lemma 2(i),

$$T - K_\Gamma = \text{tr}[M(I + A_\Gamma M)^{-1}] = \text{tr}\left[\frac{(I + A_\Gamma M)^{-1} + (I + M A_\Gamma^\theta)^{-1}}{2} M\right],$$

and that the eigenvalues of $\frac{(I + A_\Gamma M)^{-1} + (I + M A_\Gamma^\theta)^{-1}}{2}$ are in $\frac{1 - k\Gamma k}{1 + k\Gamma k}$ to $\frac{1 + k\Gamma k}{1 - k\Gamma k}$ with $k\Gamma k < 1 - c$. \square

Lemma 4. Suppose $y_t = x_t^\theta \beta + \varepsilon_t$, $x_t = \tilde{x}_t + \sum_{\ell=1}^L \alpha_\ell \varepsilon_{t-\ell}$ where the $T \times K$ matrix $\tilde{X} = [\tilde{x}_1, \dots, \tilde{x}_T]^\theta$ has full rank, and $\alpha_1, \dots, \alpha_L, r$ are $K - 1$ vectors. Then, there exists an invertible $K \times K$ matrix Θ mapping $f\tilde{X}, r, \beta, f\alpha_\ell g_{\ell=1}^L g$ to $f\tilde{X}\Theta, \Theta^\theta r, \Theta^{-1}\beta, f\Theta^\theta \alpha_\ell g_{\ell=1}^L g$, that satisfies:

- (i) $\Theta^\theta r$ and $f\Theta^\theta \alpha_\ell g_{\ell=1}^L$ are spanned by the first $L+1$ basis vectors so that $r^\theta \Theta = (r^\theta, \mathbf{0}_{K-L-1}^\theta)$ and $\alpha_\ell^\theta \Theta = (\alpha_{\ell, L+1}^\theta, \mathbf{0}_{K-L-1}^\theta)$ with $r, \alpha_{1, L+1}, \dots, \alpha_{L, L+1} \in \mathbb{R}^{L+1}$;
- (ii) $\Theta^\theta \tilde{X}_1^\theta \tilde{X}_1 \Theta / T = I_{L+1}$ and $\Theta^\theta \tilde{Z}_2^\theta \tilde{X}_1 \Theta = 0$ for $\tilde{X} = [\tilde{X}_1, \tilde{X}_2]$, and $\tilde{Z} = [\tilde{Z}_1, \tilde{Z}_2] = (I - \Gamma^\theta) \tilde{X} \Theta$, where \tilde{X}_1 and \tilde{Z}_1 each have $L+1$ columns.

Proof of Lemma 4. Let $\Theta = \Theta_0^{-1} \Theta_1$, where $K \times K$ matrix Θ_0 is the symmetric square root of $[\tilde{X}_1, \tilde{Z}_2]^\theta [\tilde{X}_1, \tilde{Z}_2] / T$ and $\Theta_1 R_1$ is the QR decomposition of $\Theta_0^{-1} [r, \alpha_1, \dots, \alpha_L]$. Specifically, $\Theta_1^\theta = \Theta_1^{-1}$ is a $K \times K$ matrix and R_1 is spanned by the first $L+1$ basis vectors. We then have that $\Theta^\theta [r, \alpha_1, \dots, \alpha_L] = \Theta_1^\theta \Theta_0^{-1} [r, \alpha_1, \dots, \alpha_L] = \Theta_1^\theta \Theta_1 R_1 = R_1$ while we also have $\Theta^\theta [\tilde{X}_1, \tilde{Z}_2]^\theta [\tilde{X}_1, \tilde{Z}_2] \Theta / T = \Theta_1^\theta \Theta_0^{-1} \Theta_0^2 \Theta_0^{-1} \Theta_1 = I_K$. \square

Lemma 5. Suppose X and Z are $T \times K$ matrices with $Z^\theta X$ invertible. Let $X = [X_1, X_2]$ and $Z = [Z_1, Z_2]$, where X_ℓ and Z_ℓ are $T \times K_\ell$ with $K_1 + K_2 = K$. Define the oblique projections $P_Z = X(Z^\theta X)^{-1} Z^\theta = I - M_Z$ and $P_2 = X_2(Z_2^\theta X_2)^{-1} Z_2^\theta = I - M_2$.

- (i) (Generalized Frisch-Waugh-Lowell) If $r = [r_1^\theta; \mathbf{0}_{K_2}^\theta]^\theta$ with $r_1 \in \mathbb{R}^{K_1}$, then

$$r^\theta (Z^\theta X)^{-1} Z^\theta = r_1^\theta (Z_1^\theta M_2 X_1)^{-1} Z_1^\theta M_2.$$

- (ii) $M_Z = M_2 - M_2 X_1 (Z_1^\theta M_2 X_1)^{-1} Z_1^\theta M_2$.

- (iii) If $Z = (I - \Gamma^\theta) X$ for $T \times T$ matrix Γ with $k\Gamma k < 1 - c$ and A is a $T \times T$ matrix, then

$$j\text{tr}(A(M_Z - M_2))j \leq C K_1 k A k.$$

Proof of Lemma 5. Letting $\Delta = (Z_1^\theta M_2 X_1)^{-1}$, the matrix block inversion formula gives

$$(Z^\theta X)^{-1} = \begin{bmatrix} \Delta & \Delta Z_1^\theta X_2 (Z_2^\theta X_2)^{-1} \\ (Z_2^\theta X_2)^{-1} Z_2^\theta X_1 \Delta & (Z_2^\theta X_2)^{-1} I + Z_2^\theta X_1 \Delta Z_1^\theta X_2 g \end{bmatrix} \quad (9)$$

(i) From (9) we have

$$r^\theta (Z^\theta X)^{-1} Z^\theta = r_1^\theta \Delta Z_1^\theta \quad r_1^\theta \Delta Z_1^\theta X_2 (Z_2^\theta X_2)^{-1} Z_2^\theta = r_1^\theta \Delta Z_1^\theta M_2.$$

(ii) Denote $\delta = X_1 \Delta Z_1^\theta$. Using (9) above:

$$P_Z = X (Z^\theta X)^{-1} Z^\theta = \delta \quad \delta P_2 \quad P_2 \delta + P_2 + P_2 \delta P_2 = I \quad M_2 + M_2 \delta M_2.$$

(iii) We impose without loss of generality that $X_2^\theta X_1 = 0$ and $X_1^\theta X_1 = I_{K_1}$. This entails no loss since (4) and (6) yields $Z_2^\theta M_2 = 0$ and $M_2 X_2 = 0$, which in turn implies that (M_Z, M_2) is invariant under the transformation $[X_1, X_2] \mathcal{V} [M \ X_1 (X_1^\theta M \ X_1)^{-1/2}, X_2]$ where $M = I \ X_2 (X_2^\theta X_2)^{-1} X_2^\theta = I \ P$. From (ii), Lemma 2(iii), and $k\Psi k_F = k\Psi k \sqrt{\text{rk}(\Psi)}$:

$$\begin{aligned} j \text{tr}(A(M_\Gamma \ M_2)) j &= \left| \text{tr} \left(Z_1^\theta M_2 A M_2 X_1 (Z_1^\theta M_2 X_1)^{-1} \right) \right| \\ &= K_1 k Z_1^\theta M_2 A M_2 X_1 k \ k (Z_1^\theta M_2 X_1)^{-1} k \\ &= K_1 k I \ \Gamma k \ k M_2 k^2 \ k A k \ k (Z_1^\theta M_2 X_1)^{-1} k. \end{aligned}$$

Equation (6) gives $M_2 = (I \ P \ \Gamma)^{-1} M$ and therefore $kM_2 k = \frac{1}{1 - k\Gamma k} < C$. Together with $X_2^\theta X_1 = 0$, (6) also yields

$$Z_1^\theta M_2 X_1 = X_1^\theta (I \ \Gamma) (I \ P \ \Gamma)^{-1} M \ X_1 = X_1^\theta (I \ \Gamma) (I \ P \ \Gamma)^{-1} X_1,$$

and therefore

$$\frac{1}{2} (Z_1^\theta M_2 X_1 + X_1^\theta M_2^\theta Z_1) = \frac{1}{2} X_1^\theta \{ (I \ \Gamma) (I \ P \ \Gamma)^{-1} + (I \ \Gamma^\theta P)^{-1} (I \ \Gamma^\theta) \} X_1.$$

As $X_1^\theta X_1 = I_{K_2}$, the eigenvalues of the last matrix are larger than $\frac{1 - k\Gamma k}{1 + k\Gamma k}$. Therefore, Lemma 2(ii) and $k\Gamma k < 1 - c$ implies that $k(Z_1^\theta M_2 X_1)^{-1} k < C$. \square

Lemma 6. Suppose \tilde{X} is a $T \times k$ matrix with $\tilde{X}^\theta \tilde{X} / T = I_k$, A is a $T \times T$ matrix that is \tilde{X} -measurable, $\alpha_0, \dots, \alpha_k$ are non-random $k \times 1$ vectors with $k\alpha_\ell k < C$ for all ℓ , and $\tilde{f}_{\varepsilon_t} g_{t=1}^T \ k$

are i.i.d. conditionally on \tilde{X} with $E[\varepsilon_t j \tilde{X}] = 0$, $0 < \sigma^2 = E[\varepsilon_t^2 j \tilde{X}] < C$ and $E[\varepsilon_t^4 j \tilde{X}] < C$. Let $u_\ell = (\varepsilon_{1-\ell}, \dots, \varepsilon_{T-\ell})^\top$. Then, as $T \rightarrow \infty$ while k is fixed: (i) $kT^{-1/2} \tilde{X}^\top A \sum_{\ell=0}^k u_\ell \alpha_\ell k_F = O_p(kAK)$; (ii) $\varepsilon^\top A \varepsilon - E[\varepsilon^\top A \varepsilon j \tilde{X}] = O_p(kAK_F)$.

Proof of Lemma 6. (i) Note that:

$$\begin{aligned} E \left\| \frac{\tilde{X}^\top A}{T} \sum_{\ell=0}^k u_\ell \alpha_\ell \right\|_F^2 &= \text{tr} \left[\frac{\tilde{X}^\top \tilde{X}}{T} A \sum_{\ell,j=0}^k E[u_\ell u_j^\top] \alpha_\ell \alpha_j^\top A^\top \right] = \text{tr} \left[\frac{\tilde{X}^\top \tilde{X}}{T} \right] \left\| A \sum_{\ell,j=0}^k E[u_\ell u_j^\top] \alpha_\ell \alpha_j^\top A^\top \right\| \\ &= k k_A k^2 \left(\sum_{\ell,j=0}^k j \alpha_\ell \alpha_j^\top \left\| E[u_\ell u_j^\top] \right\| \right) = C k^3 k_A k^2. \end{aligned}$$

The last inequality uses that u_ℓ and u_j has $E[u_\ell u_j^\top] = \sigma^2 D^{\ell-j}$ when $\ell > j$, $E[u_\ell u_j^\top] = \sigma^2 (D^\ell)^{j-\ell}$ when $\ell < j$, and $E[u_\ell u_\ell^\top] = \sigma^2 I$. In all cases, $k E[u_\ell u_j^\top] k \leq \sigma^2$. (ii) We have

$$\varepsilon^\top A \varepsilon - E[\varepsilon^\top A \varepsilon] = \sum_t \sum_{s \neq t} A_{ts} \varepsilon_t \varepsilon_s + \sum_t A_{tt} (\varepsilon_t^2 - \sigma^2).$$

In the summations above summands are correlated only when $ft, sg = ft^\ell, s^\ell g$. Therefore

$$\begin{aligned} E \left[\left(\varepsilon^\top A \varepsilon - E[\varepsilon^\top A \varepsilon] \right)^2 \right] &= \sigma^4 \sum_t \sum_{s \neq t} (A_{ts} A_{st} + A_{ts}^2) + \sum_t A_{tt}^2 E[(\varepsilon_t^2 - \sigma^2)^2] \\ &= C \left(\sum_t \sum_{s \neq t} (A_{ts} A_{st} + A_{ts}^2) + \sum_t A_{tt}^2 \right) = 2C \sum_{t,s} A_{ts} A_{st} = 2C k_A k^2, \end{aligned}$$

for $C = \max \{ \sigma^4, E[(\varepsilon_t^2 - \sigma^2)^2] \}$. We used that $j \sum_t \sum_{s \neq t} A_{ts} A_{st} j = \sum_t \sum_{s \neq t} A_{ts}^2$. \square

Lemma 7. Suppose $B = D^\theta (I - \Gamma) M_\Gamma$ where $\Gamma = \gamma D$, $j \gamma_j j < 1 - c$ and $\text{tr}(B) = O(K)$. Then, (i) $\sum_t B_{tt}^2 = O(K)$; (ii) $\text{tr}(B^2) = \sum_{t,s} B_{ts} B_{st} = O(K)$; (iii) $\frac{\text{tr}(B^2)}{\text{tr}(B^2)} = \frac{\text{tr}(B^2)}{\sum_{s,t} B_{s,t}^2} = O(1)$.

Proof of Lemma 7. Due to Lemma 2(i) and equation (5) we have, for any matrix A ,

$$j \text{tr}(AP) j = \left| \text{tr} \left(\frac{A + A^\top}{2} P \right) \right| \left\| \frac{A + A^\top}{2} \right\| \text{tr}(P) = k_A k K; \quad (10)$$

$$j \text{tr}(AP_\Gamma) j = j \text{tr}((I - \Gamma) A (I - P_\Gamma)^\top P) j = \frac{1 + k_\Gamma k}{1 - k_\Gamma k} k_A k K. \quad (11)$$

(i) $B = D^\theta (I - \Gamma) = D^\theta (I - \Gamma) P_\Gamma$, thus $B_{tt} = \gamma F_{tt}$, where $F = D^\theta (I - \Gamma) P_\Gamma$. The

condition $\text{tr}(B) = O(K)$ implies $\gamma T = \sum_t F_{tt} + O(K)$ and $\sum_t B_{tt}^2 = \sum_t F_{tt}^2 - T\gamma^2 + O(K)$
 $\sum_t F_{tt}^2 + O(K)$. Consider diagonal elements of the matrix $F = AP_\Gamma$ with $A = D^\theta(I - \Gamma)$ and
 $kAK = 1 + j\gamma j$:

$$\begin{aligned} \sum_t F_{tt}^2 &= \sum_t (AP_\Gamma)_{tt}^2 = \sum_t \left(\sum_s A_{ts} P_{\Gamma, st} \right)^2 = \sum_t \left\{ \sum_s A_{ts}^2 \sum_s P_{\Gamma, st}^2 \right\} \\ &= \sum_t (AA^\theta)_{tt} (P_\Gamma P_\Gamma^\theta)_{tt} = kAK^2 \text{tr}(P_\Gamma P_\Gamma^\theta) = CK. \end{aligned}$$

(ii) Next we have

$$\begin{aligned} \text{tr}(B^2) &= \text{tr}(D^\theta(I - \Gamma)(I - P_\Gamma)D^\theta(I - \Gamma)(I - P_\Gamma)) \\ &= \text{tr}[D^\theta(I - \Gamma)D^\theta(I - \Gamma)] + \text{tr}[D^\theta(I - \Gamma) \{ 2D^\theta(I - \Gamma) + P_\Gamma D^\theta(I - \Gamma) \} P_\Gamma]. \end{aligned} \quad (12)$$

Since $\sum_t B_{tt}^2 = 0$ reasoning above gives us that

$$\text{tr}(D^\theta(I - \Gamma)D^\theta(I - \Gamma)) = T\gamma^2 - \sum_t F_{tt}^2 = O(K).$$

The second term in (12) is $O(K)$ due to (11). (iii) We use that for any matrix A , $\text{tr}(A^\theta A) = \sum_{t,s} A_{ts}^2 = j \sum_{t,s} A_{ts} A_{st} j = j \text{tr}(A^2) j$:

$$\begin{aligned} \text{tr}(B^\theta B) &= \text{tr}(M_\Gamma^\theta (I - \Gamma^\theta) (I - \Gamma) M_\Gamma) = (1 - j\gamma j)^2 \text{tr}(M_\Gamma^\theta M_\Gamma) \\ &= (1 - j\gamma j)^2 \text{tr}(M_\Gamma^2) = (1 - j\gamma j)^2 \text{tr}(M_\Gamma) = (1 - j\gamma j)^2 (T - K). \quad \square \end{aligned}$$

Lemma 8. *Suppose that \tilde{X} is a $T \times K$ matrix of rank K and $\Gamma = \gamma D$.*

(i) *Equation (1) holds if and only if γ is a fixed point of the transformation f given by*

$$f(\gamma) = \frac{\text{tr}(D^\theta \tilde{M}_\Gamma)}{T - K}. \quad (13)$$

(ii) *If $j \text{tr}(D^\theta \tilde{M}) j = \mu^2 K$ and $K < T/(1 + (1 + \mu)^2)$ for some $\mu \geq [0, 1]$, then f is a contraction on $[\mu, \mu]/(1 + \mu)$ with Lipschitz constant strictly less than μ .*

Proof. (i) Since $\Gamma = \gamma D$ we have $D^\theta(I - \gamma D)\tilde{M}_\Gamma = D^\theta \tilde{M}_\Gamma - D^\theta D \gamma \tilde{M}_\Gamma$. We can therefore

re-write (1) as: $\text{tr}(D^\theta \tilde{M}_\Gamma) - \gamma(T - K) = 0$. This equation is solved if (and only if) $\gamma = f(\gamma)$. (ii) Equation (6) yields $\tilde{M}_\Gamma = \tilde{M} + \gamma \tilde{P} D \tilde{M}_\Gamma$ and $\|k \tilde{M}_\Gamma k\| = \frac{1}{1 - \gamma j}$. Equation (10) gives $\|j \text{tr}(D^\theta \tilde{P} D \tilde{M}_\Gamma) j\| = \|K k D \tilde{M}_\Gamma D^\theta k\|$. Therefore,

$$\|j f(\gamma) j\| = \frac{\|j \text{tr}(D^\theta \tilde{M}) j\|}{T - K} + \gamma j j \frac{\|j \text{tr}(D^\theta \tilde{P} D \tilde{M}_\Gamma) j\|}{T - K} = \frac{\mu^2 K}{T - K} + \gamma j j \frac{K}{(T - K)(1 - j \gamma j)} < \frac{\mu}{1 + \mu}$$

and using equation (7)

$$\frac{\|j f(\gamma_1) j\| - \|f(\gamma_2) j\|}{\|j \gamma_1 j\| - \|j \gamma_2 j\|} = \frac{\|j \text{tr}[D^\theta (\tilde{M}_{\Gamma_1} - \tilde{M}_{\Gamma_2})] j\|}{\|j \gamma_1 j\| - \|j \gamma_2 j\| (T - K)} = \frac{\|j \text{tr}[D^\theta \tilde{P} D (I - \tilde{P} \Gamma_2)^{-1} \tilde{M}_{\Gamma_1}] j\|}{T - K} \frac{K}{T - K} \frac{1}{1 - \|j \gamma_1 j\|} \frac{1}{1 - \|j \gamma_2 j\|} < \mu$$

where the strict inequalities use $K < T/(1 + (1 + \mu)^2)$ and $\|j \gamma j\|, \|j \gamma_1 j\|, \|j \gamma_2 j\| < \frac{\mu}{1 + \mu}$. \square

A.2 Proofs for results stated in the main text

Proof of Theorem 1. Special case of Theorem 3(i) with $\Gamma = 0$. \square

Proof of Theorem 2. Special case of Theorem 4 with $\Gamma = 0$. \square

Proof of Lemma 1. (i) Special case of (ii) with $\mu = 1$ since $\|j \text{tr}(D^\theta \tilde{M}) j\| = \|j \text{tr}(D^\theta \tilde{P}) j\| = K$ follows from (10). (ii) Follows from Lemma 8 and the Banach fixed point theorem. \square

Proof of Theorem 3. (i) Special case of Theorem 7 with $L = 1$. (ii) Let $f(\gamma)$ be as in equation (13) define its empirical analog $\hat{f}(\gamma) = \frac{\text{tr}(D^\theta M)}{T - K} + \gamma \frac{\text{tr}(D^\theta P D M_\Gamma)}{T - K}$ where $\Gamma = \gamma D$. Since $K < T/5$, Lemma 8(ii) yields that both f and \hat{f} are contractions on $[-1/2, 1/2]$ with contraction speed bounded by $\frac{1}{2}$ and therefore has unique fixed points γ_0 and $\hat{\gamma}$ by the Banach fixed point theorem. Furthermore, $\|j \hat{f}(\hat{\gamma}) j\| = \|\hat{f}(\gamma_0) j\| = \frac{1}{2} \|\hat{\gamma} - \gamma_0\|$. For any γ we have:

$$\left| \hat{f}(\gamma) - f(\gamma) \right| \leq \frac{1}{T - K} \left(\left| \text{tr}[D^\theta (P - \tilde{P})] \right| + \gamma \left| \text{tr}[D^\theta P D M_\Gamma - D^\theta \tilde{P} D \tilde{M}_\Gamma] \right| \right).$$

Consider the transformation Θ of Lemma 4. Since the projections $P, \tilde{P}, M_\Gamma, \tilde{M}_\Gamma$ are invariant to linear transformations we may assume that Lemma 4(i) and (ii) hold with $L = 1$ and $\Theta = I_K$. This implies that $M_2 = \tilde{M}_2$, where M_2 is defined as in Lemma 5 using $X, (I - \Gamma^\theta)X$, and $K_1 = 2$, while \tilde{M}_2 is an analogously defined starting from $\tilde{X}, (I - \Gamma^\theta)\tilde{X}$, and $K_1 = 2$.

Therefore, Lemma 5(iii) implies for any compatible matrix A that

$$\left| \operatorname{tr}[A(M_\Gamma \quad \tilde{M}_\Gamma)] \right| = \left| \operatorname{tr}[A(M_\Gamma \quad M_2)] \right| + \left| \operatorname{tr}[A(\tilde{M}_\Gamma \quad M_2)] \right| \leq CkAk$$

A similar statement holds with (P, \tilde{P}) replacing $(M_\Gamma, \tilde{M}_\Gamma)$. Thus $j\hat{f}(\gamma) - f(\gamma)j \leq \frac{C}{T-K}$. As a result:

$$j\hat{\gamma} - \gamma_0j = \left| \hat{f}(\hat{\gamma}) - f(\gamma_0) \right| = \left| \hat{f}(\hat{\gamma}) - \hat{f}(\gamma_0) \right| + \left| \hat{f}(\gamma_0) - f(\gamma_0) \right| \leq \frac{1}{2}j\hat{\gamma} - \gamma_0j + \frac{C}{T-K}.$$

This implies $j\hat{\gamma} - \gamma_0j \leq \frac{1}{2} \frac{C}{T-K} = O_p(1/T)$. Equation (7) yields

$$\left| r^\theta(\hat{\beta}^{\text{IV}}(\hat{\Gamma}) - \hat{\beta}^{\text{IV}}(\Gamma_0)) \right| \leq \left\| (X^\theta X)^{-1/2} r \right\| \left\| (I + A_{\Gamma_0} M)^{-1} (A_{\hat{\Gamma}} - A_{\Gamma_0}) M (I + A_{\hat{\Gamma}} M)^{-1} \varepsilon \right\| \\ \leq \left\| (X^\theta X)^{-1/2} r \right\| \left\| (I + A_{\Gamma_0} M)^{-1} \right\| \|A_{\hat{\Gamma}} - A_{\Gamma_0}\| \left\| (I + A_{\hat{\Gamma}} M)^{-1} \right\| k\varepsilon k,$$

where $A_{\Gamma_0} = (I - \Gamma_0)^{-1} \Gamma_0$ and $A_{\hat{\Gamma}} = (I - \hat{\Gamma})^{-1} \hat{\Gamma}$. We have $k\varepsilon k = O_p(\sqrt{\frac{p}{T}})$ and $kA_{\hat{\Gamma}} - A_{\Gamma_0}k = O_p(\hat{\gamma} - \gamma_0) = O_p(1/T)$. By Assumption 1(iii) $r^\theta(X^\theta X)^{-1} r = O(1/T)$, while $\left\| (I + A_{\Gamma_0} M)^{-1} \right\|$ is uniformly bounded. Thus, $r^\theta \hat{\beta}^{\text{IV}}(\hat{\Gamma}) - r^\theta \hat{\beta}^{\text{IV}}(\Gamma_0) = O_p(T^{-1}) = o_p(1/\sqrt{\frac{p}{T}})$. \square

The proofs of Theorems 4-6 follow after the proof of Theorem 7 as they rely on results established therein.

Proof of Theorem 7. Define $R_{\Gamma, \ell} = \operatorname{tr}[(D^\theta)^\ell (I - \Gamma) \tilde{M}_\Gamma]$. Note that $r^\theta \bar{S}^{-1} \alpha_\ell$, $R_{\Gamma, \ell} \mathcal{G}_{\ell=1}^L$ and $r^\theta \hat{\beta}^{\text{IV}}(\Gamma) - r^\theta \beta$ are invariant under the transformation Θ of Lemma 4. Thus, we may assume without loss of generality that Lemma 4(i) and (ii) hold with $\Theta = I_K$ and $Z = (I - \Gamma^\theta) X$. Since $(\alpha_1, \dots, \alpha_L)$ is spanned by the first $L+1$ basis vectors, we have $M_2 = \tilde{M}_2$, $X_2 = \tilde{X}_2$, and $Z_2 = \tilde{Z}_2$, where M_2 , X_2 , and Z_2 are defined as in Lemma 5 using X , $Z = (I - \Gamma^\theta) X$, and $K_1 = L+1$, while \tilde{M}_2 , \tilde{X}_2 , \tilde{Z}_2 are analogously defined starting from \tilde{X} , $\tilde{Z} = (I - \Gamma^\theta) \tilde{X}$, and $K_1 = L+1$. This also implies that X_2 , Z_2 , and M_2 are non-random conditionally on \tilde{X} . Lemma 5(i) now yields:

$$r^\theta \hat{\beta}^{\text{IV}}(\Gamma) - r^\theta \beta = r^\theta (Z_1^\theta M_2 X_1)^{-1} Z_1^\theta M_2 \varepsilon. \quad (14)$$

Defining $\bar{S}_2 = \mathbb{E}[Z_1^\theta M_2 X_1 j \tilde{X}] = \tilde{Z}_1^\theta \tilde{X}_1 + \sigma^2 \sum_{j, \ell=1}^L \alpha_{\cdot, j} \alpha_{\cdot, \ell}^\theta \operatorname{tr}[(D^\theta)^j (I - \Gamma) M_2 D^\ell]$ and $R_{2, \ell} =$

$\text{tr}[(D^\theta)^\ell(I - \Gamma)M_2]$, we show as a **first** step that

$$r^\theta \hat{\beta}^{\text{IV}}(\Gamma) - r^\theta \beta = \sigma^2 \sum_{\ell=1}^L r^\theta \bar{S}_2^{-1} \alpha_{\cdot, \ell} R_{2, \ell} + o_p(1). \quad (15)$$

Equation (15) follows from equation (14) and the following statements proven below:

$$\|(Z_1^\theta M_2 X_1 - \bar{S}_2)/T\|_F = O_p(\sqrt{1/T}), \quad (16)$$

$$\|(\bar{S}_2/T)^{-1}\| = \frac{1}{(1 - k\Gamma k)}, \quad (17)$$

$$\|(Z_1^\theta M_2 \varepsilon - \sigma^2 \sum_{\ell=1}^L \alpha_{\cdot, \ell} R_{2, \ell})/T\| = O_p(\sqrt{1/T}), \quad (18)$$

and $\|R_{2, \ell}\| \leq T$ and thus $k \sum_{\ell=1}^L \alpha_{\cdot, \ell} R_{2, \ell} k = O(T)$. Let $u_\ell = (\varepsilon_{1-\ell}, \dots, \varepsilon_{T-\ell})^\theta$ as in Lemma 6. Equation (16) considers a $(L+1) \times (L+1)$ matrix with mean of zero:

$$\begin{aligned} Z_1^\theta M_2 X_1 - \bar{S}_2 &= \left(\tilde{X}_1 + \sum_{j=1}^L u_j \alpha_{\cdot, j}^\theta \right)^\theta (I - \Gamma) M_2 \left(\tilde{X}_1 + \sum_{\ell=1}^L u_\ell \alpha_{\cdot, \ell}^\theta \right) - \bar{S}_2 \\ &= \sum_{\ell=1}^L \alpha_{\cdot, \ell} u_\ell^\theta (I - \Gamma) M_2 \tilde{X}_1 + \tilde{X}_1^\theta (I - \Gamma) M_2 \sum_{\ell=1}^L u_\ell \alpha_{\cdot, \ell}^\theta \\ &\quad + \sum_{j, \ell=1}^L \alpha_{\cdot, \ell} \alpha_{\cdot, j}^\theta [u_\ell^\theta (I - \Gamma) M_2 u_j - \mathbb{E}[u_\ell^\theta (I - \Gamma) M_2 u_j]]. \end{aligned}$$

Notice that for $A = (I - \Gamma)M_2$, we have $kAk = \frac{1+k\Gamma k}{1-k\Gamma k}$ and $kAk_F = O(\sqrt{\frac{1}{T}})$. Thus, applying Lemma 6(i) and (ii) we obtain (16). As Lemma 3(iii) yields that $B = (I - \Gamma)M_2 + M_2^\theta (I - \Gamma)^\theta$ is a non-negative definite matrix, we have

$$\begin{aligned} x^\theta (\bar{S}_2 + \bar{S}_2^\theta) x &= x^\theta [\tilde{Z}_1^\theta \tilde{X}_1 + \tilde{X}_1^\theta \tilde{Z}_1] x + \sum_{j, \ell=1}^L x^\theta \alpha_{\cdot, j} \mathbb{E} [u_j^\theta B u_\ell j \tilde{X}] \alpha_{\cdot, \ell}^\theta x \\ &\quad x^\theta \tilde{X}_1^\theta [(I - \Gamma + \Gamma^\theta)/2] \tilde{X}_1 x - (1 - k\Gamma k) k \tilde{X}_1 x k^2 = (1 - k\Gamma k) T k x k^2 \end{aligned}$$

Applying Lemma 2(ii) now implies (17). To prove (18), we note that

$$Z_1^\theta M_2 \varepsilon = \tilde{X}_1^\theta (I - \Gamma) M_2 \varepsilon + \sum_{\ell=1}^L \alpha_{\cdot, \ell} u_\ell^\theta (I - \Gamma) M_2 \varepsilon.$$

The first term is $O_p(\rho \bar{T})$ due to Lemma 6(i) applied with $A = (I - \Gamma) M_2$ and $\alpha_\ell = 0$ for $\ell > 0$. As $\sigma^2 R_{2, \ell} = E[u_\ell^\theta (I - \Gamma) M_2 \varepsilon]$, Lemma 6(ii) yields $\sum_{\ell=1}^L \alpha_{\cdot, \ell} [u_\ell^\theta (I - \Gamma) M_2 \varepsilon - R_{2, \ell}] = O_p(\rho \bar{T})$. This leads to (18).

As the **second** and final step of the proof, we show that (15) implies (3). The first difference is that (15) depends on r^θ , $\alpha_{\cdot, \ell}$, and the $(L+1) \times (L+1)$ matrix \bar{S}_2 , while (3) is described using $r^\theta = (r^\theta, \mathbf{0}_{K-L-1}^\theta)$, $\alpha_\ell^\theta = (\alpha_{\cdot, \ell}^\theta, \mathbf{0}_{K-L-1}^\theta)$, and the $K \times K$ matrix \bar{S}_Γ . The second difference between these two statements is that (15) employs the oblique projection \tilde{M}_2 , which projects off the $T \times (K-L-1)$ matrix \tilde{X}_2 , while (3) employs \tilde{M}_Γ , which projects off the full $T \times K$ matrix of regressors $\tilde{X} = [\tilde{X}_1, \tilde{X}_2]$.

For the first of these discrepancies, from Lemma 4(i) and (ii) we have that the lower left $(K-L-1) \times (L+1)$ blocks of $\alpha_\ell \alpha_j^\theta$ and $\tilde{Z}^\theta \tilde{X}$ are zero. Thus the upper left $(L+1) \times (L+1)$ block of \bar{S}_Γ^{-1} is equal to the inverse of the upper left $(L+1) \times (L+1)$ block of \bar{S}_Γ . Thus

$$r^\theta \bar{S}_\Gamma^{-1} \alpha_\ell = r^\theta \left(\bar{S}_2 + \sigma^2 \sum_{j, \ell=1}^L \alpha_{\cdot, j} \alpha_{\cdot, \ell}^\theta \Delta_{j\ell} \right)^{-1} \alpha_{\cdot, \ell},$$

where $\Delta_{j\ell} = \text{tr}[(D^\theta)^j (I - \Gamma)(\tilde{M}_\Gamma - \tilde{M}_2) D^\theta]$. Now, Lemma 5(iii) gives $j \Delta_{j\ell} j = C k D^j (D^\theta)^i (I - \Gamma) k = O(1)$ so that $\sum_{\ell=1}^L (r^\theta \bar{S}_\Gamma^{-1} \alpha_\ell - r^\theta \bar{S}_2^{-1} \alpha_{\cdot, \ell}) R_{\Gamma, \ell} = o_p(1)$. For the second discrepancy, we have $R_{\Gamma, \ell} - R_{2, \ell} = \text{tr}[(D^\theta)^\ell (I - \Gamma)(\tilde{M}_\Gamma - \tilde{M}_2)]$ and Lemma 5(iii) similarly yielding $j R_{\Gamma, \ell} - R_{2, \ell} j = O(1)$. Thus $\sum_{\ell=1}^L (R_{\Gamma, \ell} - R_{2, \ell}) r^\theta \bar{S}_2^{-1} \alpha_{\cdot, \ell} = o_p(1)$. In conclusion, we have $\sum_{\ell=1}^L r^\theta \bar{S}_\Gamma^{-1} \alpha_\ell R_{\Gamma, \ell} - r^\theta \bar{S}_2^{-1} \alpha_{\cdot, \ell} R_{2, \ell} = o_p(1)$, so (15) implies (3) and Theorem 7 follows. \square

Proof of Theorem 4. From equation (4) we have $\hat{\sigma}^2(\Gamma) = \frac{\varepsilon^\theta (I - \Gamma) M_\Gamma \varepsilon}{T - K_\Gamma}$. Note that σ^2 and $\hat{\sigma}^2(\Gamma)$ are invariant under the transformation Θ of Lemma 4. Thus, we may assume without loss of generality that Lemma 4(i) and (ii) hold with $\Theta = I_K$. As in the proof of Theorem 7, we therefore have $M_2 = \tilde{M}_2$. Applying Lemma 5(ii) yields:

$$\hat{\sigma}^2(\Gamma) = \frac{\varepsilon^\theta (I - \Gamma) \tilde{M}_2 \varepsilon}{T - K_\Gamma} = \frac{1}{T - K_\Gamma} \varepsilon^\theta (I - \Gamma) \tilde{M}_2 X_1 (Z_1^\theta \tilde{M}_2 X_1)^{-1} Z_1^\theta \tilde{M}_2 \varepsilon. \quad (19)$$

Lemma 6(ii) applied with $A = \frac{(I \ \Gamma)M_2}{T \ K_\Gamma}$ implies that $\frac{\varepsilon^\ell(I \ \Gamma)M_2\varepsilon}{\sigma^2(T \ K_\Gamma)} = 1 + o_p(1)$. For the second term in (19), we use statements (16)–(18) and the second part of the proof of Theorem 7 to obtain the conclusion of Theorem 4. \square

Proof of Theorem 5. Let γ_0 be as in Theorem 3. Note that $r^\ell \hat{\beta}^{\text{IV}}(\Gamma_0), r^\ell \beta$ and $\hat{\Sigma}_T(\Gamma_0) = \hat{\sigma}^2(\Gamma_0)k r^\ell (X^\ell(I \ \Gamma_0)X)^\ell X^\ell(I \ \Gamma_0)k^2$ are invariant under the transformation Θ of Lemma 4. Thus, we may assume without loss of generality that Lemma 4(i) and (ii) hold with $\Theta = I_K$ and $Z = (I \ \Gamma_0^\ell)X$. Following the proof of Theorem 7, we have formula (14), where the denominator satisfies (16) and (17). This implies that

$$r^\ell \hat{\beta}^{\text{IV}}(\Gamma_0) - r^\ell \beta = (1 + o_p(1)) \left(\sum_t w_t \varepsilon_t + r^\ell \bar{S}_2^{-1} \alpha^\ell \varepsilon^\ell B \varepsilon \right), \quad (20)$$

where $w^\ell = (w_1, \dots, w_T) = r^\ell \bar{S}_2^{-1} \tilde{Z}_1^\ell M_2$, $B = D^\ell(I \ \Gamma_0)M_2$, and (\bar{S}_2, M_2) is defined as in the proof of Theorem 7 with $\Gamma = \Gamma_0$. The weights $\tilde{r} w_t g$ and the matrix B are measurable with respect to \tilde{X} . Below, we show that

$$\frac{r^\ell \hat{\beta}^{\text{IV}}(\Gamma_0) - r^\ell \beta}{\sqrt{\Sigma_T}} = \frac{\sum_t w_t \varepsilon_t + r^\ell \bar{S}_2^{-1} \alpha^\ell \varepsilon^\ell B \varepsilon}{\sqrt{\Sigma_T}} \Big) \ N(0, 1), \quad (21)$$

where $\Sigma_T = \sigma^2 \sum_t w_t^2 + \sigma^4 (r^\ell \bar{S}_2^{-1} \alpha^\ell)^2 \text{tr}(B^2 + B^\ell B)$. We have

$$\varepsilon^\ell B \varepsilon = \sum_t B_{tt} \varepsilon_t^2 + \sum_t \sum_{s \notin t} \frac{B_{st} + B_{ts}}{2} \varepsilon_t \varepsilon_s.$$

Lemma 7(i) and (iii) together with $K/T \neq 0$ imply that $(r^\ell \bar{S}_2^{-1} \alpha^\ell)^2 (\sum_t B_{tt} \varepsilon_t^2) / \sqrt{\Sigma_T} \xrightarrow{P} 0$ and $\frac{\text{tr}(B^2)}{\text{tr}(B^\ell B)} \neq 0$.

We obtain statement (21) by establishing the four conditions, (i)–(iv), of Sölvsten (2020), Corollary A2.8, located in the Supplementary Appendix of that article. Condition (i) is automatically satisfied if we define $w_{t,T} = \frac{\tilde{r} w_t}{\Sigma_T}$, $M_{st} = \frac{r^\ell \bar{S}_2^{-1} \alpha^\ell}{2 \sqrt{\Sigma_T}} (B_{st} + B_{ts})$ for $s \notin t$, and $M_{tt} = 0$. Condition (iv) is implied by Assumption 1(ii). To establish condition (ii) we note that by Lemma 5(i), we have for any \tilde{r} of the form $\tilde{r}^\ell = (\tilde{r}^\ell, \mathbf{0}_{K-L-1}^\ell)$ that

$$\tilde{r}^\ell (\tilde{Z}^\ell \tilde{X})^\ell \tilde{Z}^\ell = \tilde{r}^\ell (\tilde{Z}_1^\ell M_2 \tilde{X}_1)^\ell \tilde{Z}_1^\ell M_2$$

where we will use that $M_2 = \tilde{M}_2$ since $X_2 = \tilde{X}_2$ is strictly exogenous. Thus for the specific choice of \tilde{r} where $\tilde{r}^\ell = r^\ell \bar{S}_2^{-1} (\tilde{Z}_1^\ell M_2 \tilde{X}_1)$ we have

$$w^\ell = r^\ell \bar{S}_2^{-1} \tilde{Z}_1^\ell M_2 = \tilde{r}^\ell (\tilde{Z}_1^\ell M_2 \tilde{X}_1)^{-1} \tilde{Z}_1^\ell M_2 = \tilde{r}^\ell (\tilde{Z}^\ell \tilde{X})^{-1} \tilde{Z}^\ell$$

so that $w_t = \tilde{r}^\ell (\tilde{Z}^\ell \tilde{X})^{-1} (\tilde{X}_t - \gamma_0 \tilde{X}_{t+1})$. Note, that

$$\begin{aligned} \max_t j w_t j &= \left\| (\tilde{X}^\ell \tilde{X})^{-1/2} (\tilde{X}^\ell \tilde{Z})^{-1} \tilde{r} \right\| (1 + j \gamma_0 j) \max_t \left\| (\tilde{X}^\ell \tilde{X})^{-1/2} \tilde{X}_t \right\|, \\ \sum_t w_t^2 &= \tilde{r}^\ell (\tilde{Z}^\ell \tilde{X})^{-1} \tilde{X}^\ell (I - \Gamma_0) (I - \Gamma_0^\ell) \tilde{X} (\tilde{X}^\ell \tilde{Z})^{-1} \tilde{r} (1 - j \gamma_0 j)^2 \left\| (\tilde{X}^\ell \tilde{X})^{-1/2} (\tilde{X}^\ell \tilde{Z})^{-1} \tilde{r} \right\|^2. \end{aligned}$$

Thus,

$$\max_t j w_{t,T} j = \frac{\max_t j w_t j}{\sqrt{\sum_t w_t^2}} = \frac{1 + j \gamma_0 j}{1 - j \gamma_0 j} \max_t k (\tilde{X}^\ell \tilde{X})^{-1/2} \tilde{X}_t k \neq 0.$$

For condition (iii), we note that Lemma 7 and $K/T \neq 0$ yields

$$\sum_s \sum_{t \neq s} \left[\frac{B_{st} + B_{ts}}{2} \right]^2 = \frac{1}{2} \text{tr}(B^2 + B^\ell B) (1 + o(1)).$$

This yields $\sum_s \sum_{t \neq s} M_{st}^2 \neq 1$. Furthermore,

$$\left\| \frac{B + B^\ell}{2} - \text{diag}(B) \right\| = kBk + \max_t j B_{tt} j = O(1).$$

Therefore we have $k(M_{st})_{s,t} k \neq 0$ and have therefore established (21).

Finally, we prove that $\frac{\hat{\Sigma}_T}{\Sigma_T} \xrightarrow{P} 1$. Reusing the argument in the proof of Theorem 3, we first have that $(\hat{\Sigma}_T - \hat{\Sigma}_T(\Gamma_0))/\Sigma_T = o(1)$. By Lemma 5(ii), we have for $u = (\varepsilon_0, \dots, \varepsilon_{T-1})^\ell$ that

$$\begin{aligned} \hat{\Sigma}_T(\Gamma_0) &= \hat{\sigma}^2(\Gamma_0) k r^\ell (Z_1^\ell M_2 X_1)^{-1} Z_1^\ell M_2 k^2 = [1 + o_p(1)] \sigma^2 r^\ell \bar{S}_2^{-1} Z_1^\ell M_2 M_2^\ell Z_1 (\bar{S}_2^\ell)^{-1} r \\ &= [1 + o_p(1)] \sigma^2 r^\ell \bar{S}_2^{-1} (\tilde{Z}_1^\ell + \alpha u (I - \Gamma_0^\ell)) M_2 M_2^\ell (\tilde{Z}_1 + (I - \Gamma_0^\ell) u \alpha^\ell) (\bar{S}_2^\ell)^{-1} r, \end{aligned}$$

where we also used Theorem 4 and (16). From Lemma 7(i) and (iii), we have

$$\mathbb{E}[u^\ell (I - \Gamma_0^\ell) M_2 M_2^\ell (I - \Gamma_0) u] = [1 + o(1)] \mathbb{E}[\varepsilon^\ell B B^\ell \varepsilon] = [1 + o(1)] \sigma^2 \text{tr}[B^\ell B].$$

From $K/T \neq 0$ we have $\frac{\text{tr}(B^2+B^{\circ}B)}{\text{tr}(B^{\circ}B)} \neq 1$ and therefore

$$\frac{\hat{\Sigma}_T - \Sigma_T}{\Sigma_T} = \frac{2\sigma^2 r^{\circ} \bar{S}_2^{-1} \alpha u(I - \Gamma_0) M_2 M_2^{\circ} \tilde{Z}_1 + \sigma^2 (r^{\circ} \bar{S}_2^{-1} \alpha)^2 (\varepsilon^{\circ} B B^{\circ} \varepsilon - E \varepsilon^{\circ} B B^{\circ} \varepsilon)}{\Sigma_T} + o_p(1). \quad (22)$$

Define $R = M_2^{\circ} \tilde{Z}_1$, $\xi_1 = r^{\circ} \bar{S}_2^{-1} \alpha u(I - \Gamma_0) M_2 M_2^{\circ} \tilde{Z}_1$, $\xi_2 = (r^{\circ} \bar{S}_2^{-1} \alpha)^2 [\varepsilon^{\circ} B B^{\circ} \varepsilon - E \varepsilon^{\circ} B B^{\circ} \varepsilon]$.
Now,

$$\begin{aligned} E[\xi_1^2] &= C(r^{\circ} \bar{S}_2^{-1} \alpha)^2 \text{tr}(R^{\circ} B^{\circ} B R) - C(r^{\circ} \bar{S}_2^{-1} \alpha)^2 \text{tr}(R^{\circ} R) k B^{\circ} B k; \\ \frac{\xi_1}{\Sigma_T} &= O_p \left(\frac{r^{\circ} \bar{S}_2^{-1} \alpha \sqrt{\text{tr}(R^{\circ} R) k B^{\circ} B k}}{\Sigma_T} \right) \\ &= O_p \left(\frac{\sqrt{k B^{\circ} B k} (r^{\circ} \bar{S}_2^{-1} \alpha)^2 \text{tr}(B^{\circ} B) + \text{tr}(R^{\circ} R)}{\Sigma_T} \right) = O_p \left(\sqrt{\frac{k B^{\circ} B k}{\text{tr}(B^{\circ} B)}} \right). \end{aligned}$$

Lemma 6(ii) yields $\varepsilon^{\circ} B B^{\circ} \varepsilon - E \varepsilon^{\circ} B B^{\circ} \varepsilon = O_p(k B^{\circ} B k_F)$. Thus

$$\frac{\xi_2}{\Sigma_T} = O_p \left(\frac{k B^{\circ} B k_F}{\text{tr}(B^{\circ} B)} \right) = O_p \left(\frac{\sqrt{\text{tr}(B^{\circ} B) k B^{\circ} B k}}{\text{tr}(B^{\circ} B)} \right) = O_p \left(\frac{k B k}{\sqrt{\text{tr}(B^{\circ} B)}} \right).$$

Thus, by Lemma 7(iii), both terms in (22) are $O_p(1/\sqrt{T})$. □

Proof of Theorem 6. Note that $r^{\circ} \hat{\beta}^{\text{IV}}(\Gamma)$, $r^{\circ} \beta$ and $\hat{\Sigma}_T(\Gamma) = \hat{\sigma}^2(\Gamma) k r^{\circ} (X^{\circ}(I - \Gamma) X)^{-1} X^{\circ}(I - \Gamma) k^2$ are invariant under the transformation Θ of Lemma 4. Thus, we may assume without loss of generality that Lemma 4(i) and (ii) hold with $\Theta = I_K$ and $Z = (I - \Gamma) X$. Proceeding as in the proof of Theorem 5 we arrive at equation (20) with $B = D^{\circ}(I - \Gamma) M_2$. Due to Gaussianity of the errors, $r^{\circ} \bar{S}_2^{-1} \tilde{Z}_1^{\circ} M_2 \varepsilon = w^{\circ} \varepsilon$ has a Gaussian distribution conditionally on \tilde{X} with conditional variance $\sigma^2 k w k^2 = \sigma^2 k r^{\circ} \bar{S}_2^{-1} \tilde{Z}_1^{\circ} M_2 k^2$. Below we show that

$$\frac{\varepsilon^{\circ} B \varepsilon}{\sigma^2 \sqrt{\text{tr}(B^2) + \text{tr}(B^{\circ} B)}} \mid N(0, 1) \quad (23)$$

and that this term is asymptotically independent from the conditionally Gaussian term $w^{\circ} \varepsilon$.

Define $P_w = \frac{w w^{\circ}}{w^{\circ} w}$ and $M_w = I - P_w$ then $\varepsilon^{\circ} B \varepsilon = 2 \varepsilon^{\circ} B P_w \varepsilon - \varepsilon^{\circ} P_w B P_w \varepsilon + \varepsilon^{\circ} M_w B M_w \varepsilon$. We

notice that

$$\begin{aligned} |\varepsilon^\ell P_w B P_w \varepsilon| &= \left(\frac{w^\ell \varepsilon}{k w k} \right)^2 \left| \frac{w^\ell B w}{w^\ell w} \right| \|B\| \chi_1^2 = O_p(1); \\ \varepsilon^\ell B P_w \varepsilon &= \frac{w^\ell \varepsilon}{k w k} \frac{w^\ell B \varepsilon}{k w k} = \frac{w^\ell \varepsilon}{k w k} N\left(0, \frac{w^\ell B B^\ell w}{w^\ell w}\right) = O_p(1). \end{aligned}$$

Due to Lemma 7 we have $\frac{T}{\text{tr}(B^2) + \text{tr}(B^\ell B)} = O(1)$. Thus

$$\frac{\varepsilon^\ell B \varepsilon}{\sqrt{\text{tr}(B^2) + \text{tr}(B^\ell B)}} = \frac{\varepsilon^\ell M_w B M_w \varepsilon}{\sqrt{\text{tr}(B^2) + \text{tr}(B^\ell B)}} + o_p(1)$$

But $\varepsilon^\ell M_w B M_w \varepsilon$ is independent from $\frac{w^\ell \varepsilon}{k w k} N(0, \sigma^2)$. This implies that $\frac{\varepsilon^\ell B \varepsilon}{\sqrt{\text{tr}(B^2) + \text{tr}(B^\ell B)}}$ is asymptotically independent from the first term. Since $\text{tr}(B) = \sum_t B_{tt} = O_p(1)$:

$$\varepsilon^\ell B \varepsilon = \sum_t \sum_{s \neq t} \frac{B_{ts} + B_{st}}{2} \varepsilon_t \varepsilon_s + \sum_t B_{tt} (\varepsilon_t^2 - \sigma^2) + O_p(1).$$

Using that $E \varepsilon_t^4 = 3\sigma^4$, one can show that

$$\text{Var}(\varepsilon^\ell B \varepsilon) = 2\sigma^4 \sum_t \sum_{s \neq t} \left(\frac{B_{ts} + B_{st}}{2} \right)^2 + 2\sigma^4 \sum_t B_{tt}^2 = \text{tr}(B^2) + \text{tr}(B^\ell B),$$

and due to Lemma 7 the right-hand-side grows no slower than of order T . Since $\max_t \|B_{t,t}\| \|B\| = O(1)$ and the operator norm of $B + B^\ell$ is bounded, conditions (i)–(iv) of Corollary A2.8 in Sølvesten (2020) hold. Therefore (23) holds and we have asymptotic Gaussianity.

By the same argument as in the proof of Theorem 5, we can show that

$$\frac{\hat{\sigma}^2(\Gamma) k r^\ell (X^\ell(I - \Gamma)X)^\top X^\ell (I - \Gamma) k^2}{\sigma^2 k r^\ell \bar{S}_2^{-1} \tilde{Z}_1^\ell M_2 k^2 + \sigma^4 (r^\ell \bar{S}_2^{-1} \alpha)^\top \text{tr}(B^\ell B)} \xrightarrow{p} 1.$$

Pre-multiplying this estimator by $1 + \psi = \frac{j \text{tr}(B^2) + \text{tr}(B^\ell B)}{\text{tr}(B^\ell B)}$ guarantees that the resulting quantity asymptotically weakly exceeds $\Sigma_T = \sigma^2 \sum_t w_t^2 + \sigma^4 (r^\ell \bar{S}_2^{-1} \alpha)^\top \text{tr}(B^2 + B^\ell B)$. \square

Appendix B Additional Simulations

The outcome vector is generated as $y = X\beta + \varepsilon$ with $\varepsilon \sim N(0, I)$ and $\beta = 0$. The design matrix is generated as $x_{1,t} = \tilde{x}_{1,t} + a\varepsilon_{t-1}$ and $X_{-1} = \tilde{X}_{-1}$, where \tilde{X} is generated as a rotated MA(1) process with $\tilde{X}\tilde{X}^\theta/T = I_K$, independent from ε . Specifically, we generate $v_t = \rho u_{t-1} + u_t$ with $u_t, g_{t=1}^T$ *i.i.d.* $N(0, I_K)$ and define $\tilde{X} = V(V^\theta V/T)^{-1/2}$, where the square root comes from Cholesky decomposition. Across simulations, we fix the sample size at $T = 200$ and the coefficient on the feedback mechanism at $a = 1.5$. Simulation results are summarized in Figure 6 with the left panel showing results for number of regressors K between 4 and 150 (fixing ρ at 0.8). The right panel reports the results for the autocorrelation in regressors ρ between 0 and 0.98 (fixing K at 50). We report simulated values of absolute bias and standard deviation for the first coordinate of OLS and IV together with the mean absolute value of the ratio of the lower trace of M to the sample size. The results are extremely similar to the results reported in Section 2 both in term of the size of the bias/standard deviations as well as dependence on the number of regressors and their one-period predictability.

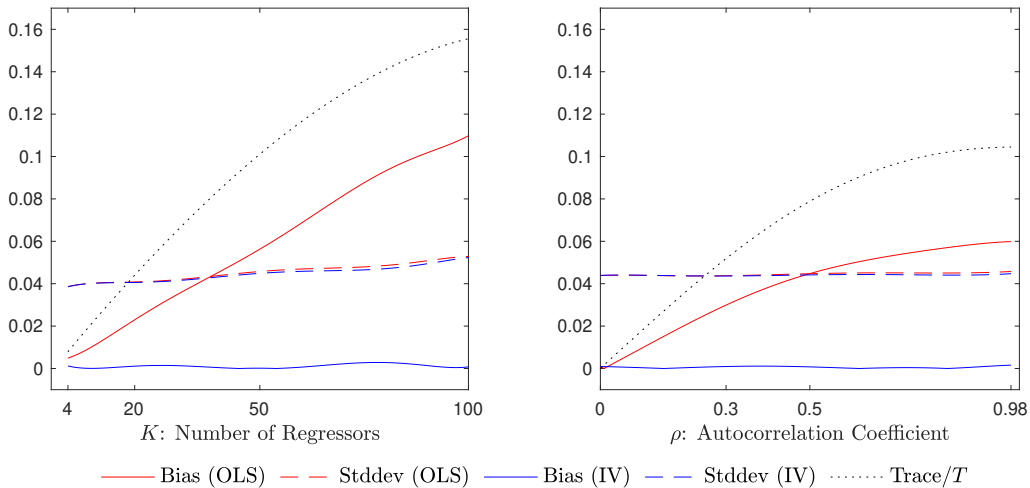


Figure 6: Absolute Bias and Standard Deviation of OLS and IV with T=200