

Inference on LATEs with covariates*

Tom Boot
University of Groningen
t.boot@rug.nl

Didier Nibbering
Monash University
didier.nibbering@monash.edu

February 21, 2024

Abstract

In theory, two-stage least squares (TSLS) identifies a weighted average of covariate-specific local average treatment effects (LATEs) from a saturated specification without making parametric assumptions on how available covariates enter the model. In practice, TSLS is severely biased when saturation leads to a number of control dummies that is of the same order of magnitude as the sample size, and the use of many, arguably weak, instruments. This paper derives asymptotically valid tests and confidence intervals for an estimand that identifies the weighted average of LATEs targeted by saturated TSLS, even when the number of control dummies and instrument interactions is large. The proposed inference procedure is robust against four key features of saturated economic data: treatment effect heterogeneity, covariates with rich support, weak identification strength, and conditional heteroskedasticity.

JEL codes: C12, C14, C21, C26.

Keywords: two-stage least squares, local average treatment effect, many controls, many instruments.

*Tom Boot acknowledges financial support by the Dutch Research Council (NWO) as part of grant VI.Veni.201E.11.

1 Introduction

With endogenous treatment and a binary instrument, [Imbens and Angrist \(1994\)](#) show that the two-stage least squares (TSLS) estimand has a causal interpretation as a local average treatment effect (LATE). Recently, [Blandhol, Bonney, Mogstad, and Torgovitsky \(2022\)](#) point out that the causal interpretation of the TSLS estimand is lost if we linearly include controls *unless* this is a correct parametric assumption. In the absence of a credible justification to linearly include the controls, [Angrist and Imbens \(1995\)](#) show that TSLS can consistently estimate a weighted average of LATEs provided that the number of possible values of the vector of controls is fixed: a researcher can select a saturated specification that includes (i) a dummy for each unique realized value of the vector of control variables and (ii) interactions of the instrument with these dummies. However, in most empirical settings the vector of controls has rich support, and saturated TSLS breaks down. As a result, applied work continues to use the more parsimonious linear specification at the risk of targeting a non-causal estimand.

In this paper we propose a new method for inference in saturated specifications that do not require parametric assumptions even when the support of the controls is rich. More specifically, this method provides asymptotically valid tests and confidence intervals for an estimand that identifies the weighted average of LATEs targeted by saturated TSLS, allowing the number of control dummies and instrument interactions to be large. As saturated specifications contain only dummy variables, we can provide low-level assumptions for our results. The most important and rather mild assumption is that each covariate group contains at least two individuals for which the instrument is active and two for which the instrument is inactive. Settings with multiple or multivalued instruments also reduce to this setting as in that case dummies for each value of the instrument are interacted with the control dummies ([Angrist and Imbens, 1995](#)).

The proposed inference procedure is robust against four key features of saturated economic data: treatment effect heterogeneity, many control dummies and instrument interactions, weak identification strength, and conditional heteroskedasticity. While methods exist to address each of these features individually, we are the first to tackle their combination by building upon two recent advances in the literature. First, [Chao, Swanson, and Woutersen \(2023\)](#) propose an estimator for a homogeneous slope coefficient under many weak instruments in a panel data setting where fixed effects take the role of control dummies. We show that this estimator can consistently estimate a weighted average of LATEs in the setting of saturated instrumental variable estimation (SIVE). Second, in a setting with a fixed number of control dummies and instrument interactions, [Kleibergen and Zhan \(2021\)](#) propose a variance estimator for the score of the continuous updating objective function that is robust to treatment effect heterogeneity. While the SIVE estimator is very different from continuous updating, we use analogous ideas to formulate a heterogeneity robust variance estimator. We show that in the setting we consider, the assumptions in both papers can be relaxed to allow for the number of control dummies and instrument interactions to be asymptotically non-negligible relative to the sample size.

To highlight our contribution, we discuss the four features of the data we consider in turn. First, the inference method is robust to fully heterogeneous treatment effects. [Imbens and Angrist \(1994\)](#) focused attention on allowing for treatment effect heterogeneity in the estimation stage. In our setting with multiple instrument interactions, this heterogeneity is equally important in the inference stage as it affects the variance of the estimator. As such, we cannot use standard variance estima-

Table 1: Empirical examples from [Blandhol et al. \(2022\)](#): sample size and covariate values

	Sample size	Covariate values	Ratio
Gelbach (2002)	440	186	0.42
Dube and Harish (2020)	107	11	0.10
Card (1995)	1780	238	0.13
Angrist and Krueger (1991)	329,463	659	0.002

Note: sample size is the effective sample size as reported in [Blandhol et al. \(2022\)](#) after accounting for perfect multicollinearity in the first stage. Covariates/Instruments indicates the number of distinct values of the available vector of controls, which equals the number of instrument interactions.

tors for TSLS with multiple instruments, as they rely on the assumption of homogeneous treatment effects. A notable exception is the TSLS variance estimator proposed by [Lee \(2018\)](#), which however is not robust to the large number of instrument interactions in typical saturated specifications. We therefore propose a new variance estimator that is robust to treatment heterogeneity when the number of control dummies and instrument interactions are large.

Second, we allow for the number of control dummies and instrument interactions to be a non-negligible fraction of the sample size. To illustrate that this matters in practice, consider the four empirical examples studied by [Blandhol et al. \(2022\)](#): [Table 1](#) shows the number of distinct covariate values and the sample size. We see that the number of covariate values, and hence the number of control dummies and instrument interactions in a saturated specification, is a substantial fraction of the sample size. The fact that TSLS is biased when the number of instruments grows proportionally with the sample size has been shown by [Bekker \(1994\)](#), and alternatives are provided by e.g. [Hansen, Hausman, and Newey \(2008\)](#); [Akerberg and Devereux \(2009\)](#); [Hausman, Newey, Woutersen, Chao, and Swanson \(2012\)](#); [Bekker and Crudu \(2015\)](#). In addition to many instrument bias, [Kolesár \(2013\)](#) proposes estimators that also remove the bias due to many controls. However, inferential procedures based on these estimators have only been developed under the assumption that the number of control dummies is a negligible fraction of the sample size ([Evdokimov and Kolesár, 2018](#)).

The third feature is that we accommodate weak instrument interactions. In particular, we allow the first stage signal to decrease to zero asymptotically. A saturated specification exacerbates the concern of weak identification, as even interacting a strong instrument with control dummies may result in instrument interactions that are only weakly related to the treatment. The condition that we impose on the identification strength has been shown by [Mikusheva and Sun \(2022\)](#) to be the weakest possible. While there is an extensive literature that combines the notion of many and weak instruments, e.g. [Bekker and Kleibergen \(2003\)](#); [Chao and Swanson \(2005\)](#); [Hausman et al. \(2012\)](#); [Mikusheva and Sun \(2022\)](#); [Crudu et al. \(2021\)](#); [Matsushita and Otsu \(2022\)](#); [Lim et al. \(2024\)](#), the focus has been on the linear IV model with a homogeneous slope coefficient.

The fourth feature is that the reduced form errors can be conditionally heteroskedastic. Inference in the presence of heteroskedasticity is non-trivial under many instruments and many controls as consistency results underlying the usual robust standard errors do not apply, see for instance [Hausman, Newey, Woutersen, Chao, and Swanson \(2012\)](#). Our newly proposed variance estimator employs estimators for the variances and covariances of the first stage and reduced form errors as in [Hartley, Rao, and Kiefer \(1969\)](#) to allow for heteroskedasticity. These variance estimators were also recently used by [Cattaneo, Jansson, and Newey \(2018\)](#). Compared to existing methods that allow for

heteroskedasticity with many instruments, our newly proposed variance estimator is also robust to many control dummies and heterogeneous treatment effects.

Provided that the weighted average of covariate-specific LATEs as derived by Angrist and Imbens (1995) is a parameter of interest, this paper presents researchers with an inference framework that is readily applicable to many empirical instrumental variable estimation problems. Evidently, there are settings in which a researcher has a different causal parameter in mind. Słoczyński (2020) points out that in our parameter of interest a larger weight is placed on covariate groups with large variation in the instrument assignment and a strong first stage. Researchers that are concerned about this weighting, can use our identification robust methods to perform a subgroup-specific analysis that zooms in on groups with little variation in the instrument. Inference on subgroup-specific LATEs that are not weighted by the instrument strength may not provide much useful insights, as any unidentified LATE that receives nonzero weight will trigger the confidence interval to be the entire real line (Evdokimov and Lee, 2013).

We conduct a series of Monte Carlo simulations which set-up mimics key features of the data used in Card (1995), and has also been used by Blandhol et al. (2022). The results illustrate that the estimator we study is median unbiased for a range of values for the instrument strength and the number of covariate groups. Fully saturated TSLS and various jackknife estimators incur a bias that increases with the number of covariate groups and as the strength of the instrument decreases. A t -test using our proposed variance estimator yields close to nominal size control regardless of the instrument strength when the number of instruments is small. When the number of instruments increases, the test becomes progressively more conservative under weak instruments, while maintaining close to nominal size control under strong instruments. The standard t -test based on the fully saturated TSLS estimator with heteroskedasticity-robust standard errors shows large size distortions. The exception is the just-identified case where the control can take on only two values, in which TSLS is known to offer close to nominal size control even under weak instruments (Angrist and Kolesár, 2023). Finally, we verify numerically that not taking into account treatment heterogeneity when estimating the variance indeed leads to an oversized test. This underlines the importance of not only accounting for treatment effect heterogeneity in the estimation stage, which has been the main focus of the extant literature, but also in the inference stage.

To illustrate the estimator we briefly revisit the data used by Card (1995) in the specification selected by Słoczyński (2020). The goal of the study is to estimate the effect of going to college on income, where the endogenous decision of going to college is instrumented by the distance to college. In particular, we consider a specification with five binary controls and binarize the treatment to having some college attendance. We document that unrealistically large estimates are obtained when the covariate dummies are not interacted with the instruments. The estimator we study yields much more reasonable point estimates although the effects statistically cannot be distinguished from zero at conventional significance levels.

The remainder of this article is organized as follows. Section 2 explains the current practice of inferring LATEs with covariates from the data and its challenges. Section 3 introduces our proposed causal estimand and its inference procedure, supported by large sample theoretical results. Section 4 discusses the Monte Carlo simulations, Section 5 the empirical application, and Section 6 concludes.

2 LATEs with covariates

Suppose we are interested in the causal effect of a binary treatment T_i on an outcome Y_i , for individuals $i = 1, \dots, n$. For each individual, define the potential outcomes $Y_i(1)$ and $Y_i(0)$ corresponding to the values of Y_i if individual i is treated or not treated, respectively. Hence, the treatment effect is defined as $Y_i(1) - Y_i(0)$. The treatment is potentially endogenous and a binary instrument Q_i and a vector of covariates X_i is available to help identifying a causal effect. The developed theory applies equally well to the extensions to multivalued treatments and instruments in Angrist and Imbens (1995) as we discuss in Appendix B.

Define $\mathbb{X} = \{x_1, \dots, x_G\}$ as the set of all possible G unique realizations of X_i . Denote the potential treatment statuses $T_i(1)$ and $T_i(0)$ corresponding to the values of T_i if individual i 's treatment assignment is given by $Q_i = 1$ and $Q_i = 0$, respectively. In case the outcome Y_i also directly depends on Q_i , its corresponding potential outcomes are given by $Y_i(Q_i, T_i)$. If we condition on the covariates, the four instrumental variable assumptions in the Imbens and Angrist (1994) framework are the following.

Assumption 1

1. **Independence:** $(Y_i(q, t), T_i(q)) \perp Q_i | X_i$ for $q \in \{0, 1\}$ and $t \in \{0, 1\}$,
2. **Exclusion:** $\mathbb{P}(Y_i(1, t) = Y_i(0, t) | X_i) = 1$ a.s. for $t \in \{0, 1\}$,
3. **Relevance:** $\mathbb{P}[T_i(1) \neq T_i(0) | X_i] > 0$ a.s.,
4. **Monotonicity:** $\mathbb{P}[T_i(1) \geq T_i(0) | X_i] = 1$ a.s., or $\mathbb{P}[T_i(1) \leq T_i(0) | X_i] = 1$ a.s.

These assumptions allow for complete treatment effect heterogeneity across all individuals, and do not impose any parametric assumptions. The monotonicity assumption is referred to in Blandhol et al. (2022) as *weak* monotonicity, because it allows the effect of the instrument on the treatment to have a different direction for each covariate group. Strong monotonicity requires $\mathbb{P}[T_i(1) \geq T_i(0)] = 1$ or $\mathbb{P}[T_i(1) \leq T_i(0)] = 1$, and therefore assumes that the effect of switching on the instrument on potential treatment status is (weakly) in the same direction for all individuals.

Within the LATE framework, causal effects are estimated of the form

$$\tau = \sum_g \omega(x_g) \tau(x_g) \text{ with } \sum_g \omega(x_g) = 1, \omega(x_g) \geq 0 \text{ for } g = 1, \dots, G, \text{ and} \quad (1)$$

$$\tau(x_g) = \mathbb{E}[Y_i(1) - Y_i(0) | T_i(1) \neq T_i(0), X_i = x_g]. \quad (2)$$

The causal effect is then a positively weighted average of covariate-specific LATEs. The following well-known result shows that the covariate group specific LATEs $\tau(x_g)$ are indeed identified.

Lemma 1 Assume that $\mathbb{E}[Y_i | Q_i, X_i]$ is almost surely bounded. Under Assumption 1 it holds that

$$\frac{\mathbb{E}[Y_i | Q_i = 1, X_i = x_g] - \mathbb{E}[Y_i | Q_i = 0, X_i = x_g]}{\mathbb{E}[T_i | Q_i = 1, X_i = x_g] - \mathbb{E}[T_i | Q_i = 0, X_i = x_g]} = \tau(x_g). \quad (3)$$

This result is discussed in Angrist and Pischke (2009), among others. For completeness, we provide a short proof in Appendix C.1.

In practice, the number of observations in each covariate-group is usually small, and the moments in Lemma 1 cannot be accurately estimated. This provides a researcher with two options.

First, we can maintain the focus on the LATE parameters, and use an (estimated) propensity score to aggregate the covariate specific LATEs. However, the propensity score can be difficult to estimate and limited overlap produces substantial statistical challenges. To avoid having to estimate the propensity score, an attractive option is to rely on regression to estimate a weighted average of the covariate group specific LATEs as the parameter of interest, where the weights are automatically selected through the regression model that is specified.

In this paper we focus on the regression approach. There are then several strategies for estimating the causal effect in (1). First, we can make a parametric assumption that restricts how the covariates enter the model. The default option is to include the covariates linearly in the first and second stage. [Blandhol et al. \(2022\)](#) shows that if this parametric assumption is incorrect, then TSLS has no causal interpretation.

To avoid parametric assumptions, a substantial and active literature considers semiparametric estimators ([Abadie, 2003](#)), or non-parametric estimators ([Frölich, 2007](#)). In particular, recent developments highlight the potential of machine learning techniques as non-parametric estimators for the effect of the controls on the outcome, endogenous treatment and instrumental variable, e.g. [Chernozhukov et al. \(2018\)](#). However, these algorithms need to attain a particular convergence rate, and it is not immediately clear whether the required conditions are met under weak identification ([Mikusheva and Sun, 2023](#)).

Finally, we can saturate the specification as suggested by [Angrist and Imbens \(1995\)](#). If we saturate, the two stage least squares estimator will only be consistent if the number of possible values of the covariate vector is fixed. However, say we have 10 binary controls, then this already gives us 1,024 possible values of the covariate vector. To accommodate cases with increasingly rich support of the covariates, we propose a new procedure that allows for reliable inference in saturated specifications. The only material additional assumption that we make is that in each covariate group there are at least two individuals for each value of the instrument.

2.1 Saturating the covariates

[Angrist and Imbens \(1995\)](#) show that τ can be estimated by TSLS in saturated specifications: the first stage includes dummies for each possible value of X_i and a full set of interactions between these dummies and the instrument, and the second stage includes the treatment variable and the control dummies. By including dummies for each possible value of the covariates, no parametric assumptions are required. To be more precise, define the $G \times 1$ vector W_i has elements $W_{ig} = \mathbb{1}[X_i = x_g]$, indicating the covariate group of individual i . The $G \times 1$ vector Z_i contains the instrument interactions $Z_{ig} = Q_i \mathbb{1}[X_i = x_g]$. Note that $\sum_g W_{ig} = 1$ and $\sum_g Z_{ig} = Q_i$. Define $n_g = \sum_i W_{ig}$ as the number of individuals in covariate group g , and $m_g = \sum_i Z_{ig}$ as the number of individuals in covariate group g with an active instrument.

Both the full set of covariate group indicators and the full set of instrument interactions are required for nonparametric estimation of τ by TSLS ([Blandhol et al., 2022](#)). This ensures that the conditional expectation of the instrument given the covariate groups $\mathbb{E}[Q_i|X_i]$ is linear in the covariate group indicators. This allows for correctly partialling out the covariates, which otherwise may induce negative weights into τ . Without the instrument interactions, the first stage does not necessarily reproduce the direction of the monotonicity assumption in all covariate groups. With a binary instrument, omitting the instrument interactions requires the direction of the monotonicity to be invariant to the covariate group.

It is clear that saturation can only work if we observe each covariate value more than once. Moreover, to achieve identification each covariate groups has to include both individuals with an active and an inactive instruments. The setup we consider, with many control dummies and many instrument interactions, requires the number of observations in each group to satisfy the following assumption.

Assumption 2 Group sizes: $m_g \geq 2$ and $n_g - m_g \geq 2$ for all $g = 1, \dots, G$.

This assumption is rather mild. It requires that both the number of individuals with an active instrument and nonactive instrument has to be larger than one in each covariate group.

2.2 Estimation challenges

While saturation results in a causal TSLS estimand if the number of possible covariate values is small, it is not straightforward to find a causal estimand in the empirically more common setting in which the controls have rich support. Since each group requires an indicator, and the instrument is interacted with these indicators, this automatically results into a large set of control dummies and a large set of instruments. In this setting, TSLS is known to be biased, see e.g. [Kolesár \(2013\)](#).

Second, the instrument interactions may weaken the instrument strength. Instrument strength is measured by the first stage signal, which can be written as $FS = \sum_g \mathbb{P}[X_i = x_g] \pi(x_g)^2 \mathbb{V}[Q_i | X_i = x_g]$ with complier shares $\pi(x_g) = \mathbb{P}[T_i(1) \neq T_i(0) | X_i = x_g]$ and treatment assignment variation $\mathbb{V}[Q_i | X_i = x_g]$. If the complier share and the variation in treatment assignment is homogeneous across covariate groups, that is $\pi(x_g) = \pi$ and $\mathbb{V}[Q_i | X_i = x_g] = \mathbb{V}[Q_i]$, the strength of the instrument interactions Z equals the strength of the instrument Q . However, in settings where the proportion of compliers is large in groups with a low number of treated or untreated units, while the proportion of compliers is small in groups with a number of treated units close to half of the number of group members, the first stage is likely weak. It is well known that under a large number of potentially weak instruments TSLS can be severely biased, see e.g. [Bekker \(1994\)](#) and [Chao and Swanson \(2005\)](#).

2.3 Two-stage least squares

Define the n -dimensional vectors $Y = (Y_1, \dots, Y_n)'$ and $T = (T_1, \dots, T_n)'$, and the $n \times G$ -dimensional matrices $W = (W_1', \dots, W_n)'$ and $Z = (Z_1', \dots, Z_n)'$. Define the residual maker matrix $M_W = I_n - W(W'W)^{-1}W'$ with I_n the n -dimensional identity matrix. The TSLS estimand is commonly defined as

$$\beta^{\text{TSLS}} = \frac{\mathbb{E}[T'PY|Q, X]}{\mathbb{E}[T'PT|Q, X]}, \quad (4)$$

where $P = M_W Z(Z'M_W Z)^{-1}Z'M_W$ partials out the controls W from the first stage and the second stage. However, this estimand is problematic when the number of covariate groups is large as the following result makes precise.

Lemma 2 Under [Assumption 1](#) and [2](#) it holds that

$$\beta^{\text{TSLS}} = \frac{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i | X_i = x_g] \tau(x_g) + \frac{1}{n} \sum_i \mathbb{E}[u_i \varepsilon_i | Q_i, X_i] P_{ii}}{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i | X_i = x_g] + \frac{1}{n} \sum_i \mathbb{E}[u_i^2 | Q_i, X_i] P_{ii}}, \quad (5)$$

where $\tilde{\mathbb{P}}[X_i = x_g] = \frac{n_g}{n}$, $\pi(x_g) = \mathbb{P}[T_i(1) \neq T_i(0) | X_i = x_g]$, $\tilde{\mathbb{V}}[Q_i | X_i = x_g] = \frac{m_g}{n_g} (1 - \frac{m_g}{n_g})$, $u_i = T_i - \mathbb{E}[T_i | Q_i, X_i]$, $\varepsilon_i = Y_i - \mathbb{E}[Y_i | Q_i, X_i]$, and $P_{ii} = \sum_g \frac{1}{n_g} W_{ig} \frac{(Z_{ig} - \tilde{\mathbb{P}}[Q_i=1 | X_i=x_g])^2}{\tilde{\mathbb{V}}[Q_i | X_i=x_g]}$.

A proof is deferred to [Appendix C.2](#). Note that $\tilde{\mathbb{P}}[X_i = x_g]$, $\tilde{\mathbb{V}}[Q_i|X_i = x_g]$, and $\tilde{\mathbb{P}}[Q_i = 1|X_i = x_g]$ are the sample analogues of $\mathbb{P}[X_i = x_g]$, $\mathbb{V}[Q_i|X_i = x_g]$, and $\mathbb{P}[Q_i = 1|X_i = x_g]$.

[Lemma 2](#) shows that both the second term in the numerator and the second term in the denominator have to go to zero for β^{TSLs} to identify τ . In this case the weights in (1) equal $\omega^{\text{TSLs}}(x_g) = \frac{\tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i|X_i = x_g]}{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i|X_i = x_g]}$. It is clear that the additional terms go to zero when P_{ii} is small. Some interpretation of this result can be obtained by noting that

$$P_{ii} = \sum_g W_{ig} \left(Z_{ig} \frac{1}{m_g} \frac{n_g - m_g}{n_g} + (1 - Z_{ig}) \frac{1}{n_g - m_g} \frac{m_g}{n_g} \right). \quad (6)$$

Hence, for TSLS to have a causal interpretation, in each covariate group the number of individuals for which the instrument is active (m_g) and the number of individuals for which the instrument is inactive ($n_g - m_g$) need to be large. This requirement becomes more stringent as the strength of the instrument interactions measured by $\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i|X_i = x_g]$ decreases.

As we have seen in [Table 1](#), the number of covariate groups is generally large and nonnegligible relative to the number of individuals. In this case, at least a number of covariate groups has to have a small number of observations, and the additional terms in [Lemma 2](#) do not go to zero. This bias is known as the many instrument bias of TSLS. Under homogenous treatment effects, and assuming homoskedastic errors, one can follow [Bekker \(1994\)](#) in using LIML to avoid this bias. However, [Kolesár \(2013\)](#) points out that with treatment effect heterogeneity, the estimand of LIML is generally not causal.

In [Lemma 2](#) we condition both on the instrument and the covariates. This estimand can generally be more accurately inferred from the data relative to the unconditional counterpart. This point is made by [Crump et al. \(2009\)](#) in a regression context. In the IV context, [Evdokimov and Kolesár \(2018\)](#) show that both the conditional and unconditional estimands are a weighted combination of covariate specific LATEs, where the unconditional estimand integrates out sampling uncertainty in the combination weights. As such, confidence intervals for the unconditional estimand are wider. We focus on the conditional estimand in the subsequent analysis.

2.4 Jackknife instrumental variables estimation

It follows from [Lemma 2](#) that the bias in TSLS is due to the diagonal elements P_{ii} . A frequently used approach to reduce many instrument bias is to employ a jackknife-style correction ([Angrist, Imbens, and Krueger, 1999](#); [Ackerberg and Devereux, 2009](#)). In the current setting, we could remove the diagonal of P , denoted by D_P , to obtain an estimand referred to as the JIVE1,

$$\beta^{\text{JIVE1}} = \frac{\mathbb{E}[T'(P - D_P)Y|Q, X]}{\mathbb{E}[T'(P - D_P)T|Q, X]}. \quad (7)$$

This diagonal removal has been the basis of recent papers in the literature on identification-robust inference under many instrument sequences ([Mikusheva and Sun, 2022](#); [Crudu, Mellace, and Sándor, 2021](#); [Matsushita and Otsu, 2022](#)). It ensures that the many instrument bias present in TSLS disappears. However, in the present case it may nevertheless not be an attractive option. With a potentially large set of control variables, the consequence of removing the diagonal elements P_{ii} is that the controls are no longer projected out. Indeed, the following result shows that removing the diagonal of the projection matrix biases the estimand.

Lemma 3 Under [Assumption 1](#) and [2](#) it holds that

$$\beta^{\text{JIVE1}} = \frac{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i | X_i = x_g] (1 - \tilde{\mathbb{P}}[Q_i = 1, X_i = x_g]) \tau(x_g) - B_Y}{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i | X_i = x_g] (1 - \tilde{\mathbb{P}}[Q_i = 1, X_i = x_g]) - B_T}, \quad (8)$$

where

$$B_Y = \sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g) (\phi_g + \tau(x_g) \psi_g) \tilde{\mathbb{V}}[Q_i | X_i = x_g] \tilde{\mathbb{P}}[Q_i = 1, X_i = x_g] + \frac{1}{n} \psi_g \phi_g, \quad (9)$$

$$B_T = 2 \sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g) \psi_g \tilde{\mathbb{V}}[Q_i | X_i = x_g] \tilde{\mathbb{P}}[Q_i = 1, X_i = x_g] + \frac{1}{n} \psi_g^2, \quad (10)$$

with $\tilde{\mathbb{P}}[Q_i = 1, X_i = x_g] = \frac{1}{m_g}$, $\psi_g = \mathbb{E}[T_i | Z_{ig} = 0, W_{ig} = 1]$, and $\phi_g = \mathbb{E}[Y_i | Z_{ig} = 0, W_{ig} = 1]$.

A proof is deferred to [Appendix C.3](#). While removing the diagonal elements has removed the many instrument bias, this is done at a high cost. When the effect of the controls on the treatment and/or the outcome are large, as measured by ϕ_g and ψ_g , the estimand for JIVE1 can be substantially different from τ and seems difficult to interpret.

An alternative jackknife approach is to first partial out the controls before removing the diagonal. This gives the following estimand labeled JIVE2.

$$\beta^{\text{JIVE2}} = \frac{\mathbb{E}[T' M_W (P - D_P) M_W Y | Q, X]}{\mathbb{E}[T' M_W (P - D_P) M_W T | Q, X]}. \quad (11)$$

The following theorem shows that this in fact reintroduces the many instrument bias, although it is smaller than that in TSLS. Also, while the estimand differs from TSLS, the weights on the covariate specific LATEs continue to be positive.

Lemma 4 Under [Assumption 1](#) and [2](#) it holds that

$$\beta^{\text{JIVE2}} = \frac{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i | X_i = x_g] \tau(x_g) + B_{Y,1} + B_{Y,2}}{\sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i | X_i = x_g] + B_{T,1} + B_{T,2}} \quad (12)$$

where

$$\begin{aligned} B_{Y,1} &= -\frac{1}{n} \sum_g \pi(x_g)^2 [1 - 3\tilde{\mathbb{V}}[Q_i | X_i = x_g]] \tau(x_g), & B_{Y,2} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[u_i \varepsilon_i | X_i, Q_i] \left(\frac{2}{n_g} P_{ii} - \frac{1}{n_g^2} \right), \\ B_{T,1} &= -\frac{1}{n} \sum_g \pi(x_g)^2 [1 - 3\tilde{\mathbb{V}}[Q_i | X_i = x_g]], & B_{T,2} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[u_i^2 | X_i, Q_i] \left(\frac{2}{n_g} P_{ii} - \frac{1}{n_g^2} \right). \end{aligned} \quad (13)$$

The proof is similar as for [Lemma 3](#) and omitted. In this case, the causal estimand still changes relative to TSLS, but under [Assumption 2](#) the weights on $\tau(x_g)$ will lie between 0 and the TSLS weights. When the noise terms are nonzero we see that the many instrument bias returns. To quantify this bias, consider momentarily a homoskedastic setting where $\mathbb{E}[u_i \varepsilon_i | X_i, Q_i] = \sigma_{u\varepsilon}$. We see that $B_{Y,2} = \frac{\sigma_{u\varepsilon}}{n} \sum_{g=1}^G \frac{1}{n_g}$ compared to the bias in TSLS that is $\frac{\sigma_{u\varepsilon}}{n} \sum_{g=1}^G 1$. We conclude that JIVE2 offers a substantial bias reduction, especially when n_g is large. However, under many instruments, n_g is fixed and the bias is of the same order as for TSLS.

3 Saturated instrumental variable estimation

3.1 A causal estimand

We now consider an estimand identical to that of TSLS when the number of possible values of the vector of controls is fixed, but that does not suffer from the many instrument bias when the number

of covariate values grows. Define the matrix $V = [Z, W]$ consisting of the instrument interactions and the covariate group indicators, and define the residual maker matrix $M_{Z,W} = I - V(V'V)^{-1}V'$. The SIVE estimand is specified as

$$\beta^{\text{SIVE}} = \frac{\mathbb{E}[T'(P - M_{Z,W}DM_{Z,W})Y|Q, X]}{\mathbb{E}[T'(P - M_{Z,W}DM_{Z,W})T|Q, X]}, \quad (14)$$

where D is a diagonal matrix with diagonal elements such that $P_{ii} = [M_{Z,W}DM_{Z,W}]_{ii}$. It follows that (14) is a jackknife estimand, which removes the diagonal of P and hence the bias in the TSLS estimand. At the same time, by pre- and post-multiplying D by $M_{Z,W}$, the controls are projected out correctly and the bias in JIVE is prevented. In addition to the controls, $M_{Z,W}$ also projects out the instrument interactions, removing the bias in the TSLS estimand.

The SIVE estimator has been proposed by [Chao, Swanson, and Woutersen \(2023\)](#) as Fixed Effect Jackknife IV (FEJIV). They show that the diagonal elements of D can be obtained by solving a system of linear equations with a unique solution, and derive consistency results in a linear instrumental variable regression with fixed effects using panel data with many weak instruments. Within our saturated setting, we can derive a closed-form expression for D , and show that the estimator identifies a weighted average of covariate-specific LATEs. These results require a different, but arguably weaker set of assumptions. For instance, we allow the number of control dummies and instrument interactions to be asymptotically non-negligible relative to the sample size, and [Assumption 2](#) relaxes [Assumption 6](#) in [Chao et al. \(2023\)](#) that requires $m_g \geq 3$ and $n_g - m_g \geq 3$.

The following result shows that the diagonal matrix D in (14) exists under [Assumption 2](#) and its elements are available in closed form.

Lemma 5 *Under [Assumption 2](#), it holds that if D is a diagonal matrix with elements*

$$D_{ii} = \sum_g \frac{1}{n_g} W_{ig} \left[\frac{n_g - m_g}{m_g - 1} Z_{ig} + \frac{m_g}{n_g - m_g - 1} (1 - Z_{ig}) \right] \quad (15)$$

then $P_{ii} = [M_{Z,W}DM_{Z,W}]_{ii}$ with $P = M_W Z(Z'M_W Z)^{-1}Z'M_W$ and $M_W = I_n - W(W'W)^{-1}W'$.

The proof is deferred to [Appendix C.4](#). The result shows that when the number of covariate values is small and both m_g and $n_g - m_g$ are large, the diagonal elements D_{ii} are small and the estimator reduces to the TSLS estimator.

Because the term that is subtracted in the numerator and denominator of (14) is orthogonal to the instruments and controls, it is straightforward to establish that the SIVE estimand has a causal interpretation that is identical to the unbiased TSLS estimand, without requiring the number of control dummies or instrument interactions to be small.

Theorem 1 *Under [Assumption 1](#) and [2](#) it holds that*

$$\beta^{\text{SIVE}} = \frac{\sum_g \hat{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Z_i | X_i = x_g] \tau(x_g)}{\sum_g \hat{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Z_i | X_i = x_g]} = \sum_g \omega(x_g) \tau(x_g). \quad (16)$$

The proof is deferred to [Appendix C.5](#).

3.2 Inference on the estimand

While having a causal estimand is a crucial first step, we also need to be able to infer the estimand from the data. In this section, we therefore consider the testing problem

$$H_0 : \beta^{\text{SIVE}} = \beta_0 \text{ against } H_1 : \beta^{\text{SIVE}} \neq \beta_0, \quad (17)$$

for a given β_0 . We develop a test statistic that is valid when treatment effects are heterogeneous, the number of values x_g that the control vector can take is non-negligible relative to the sample size, identification is weak, and the errors are heteroskedastic. The testing procedure is standard and based on the fact that under H_0 ,

$$\frac{(\hat{\beta}^{\text{SIVE}} - \beta_0)}{\sqrt{\hat{\mathbb{V}}[\hat{\beta}^{\text{SIVE}}|Q, X]}} \rightarrow_d N(0, 1). \quad (18)$$

Here, $\hat{\beta}^{\text{SIVE}}$ is simply the sample analogue of (14) and given by,

$$\hat{\beta}^{\text{SIVE}} = \frac{T'(P - M_{Z,W}DM_{Z,W})Y}{T'(P - M_{Z,W}DM_{Z,W})T}. \quad (19)$$

The crucial part to make the test operational is to find an appropriate estimator for the variance of $\hat{\beta}^{\text{SIVE}}$. Denote by $u_i = T_i - \mathbb{E}[T_i|Q_i, X_i]$, $\varepsilon_i = Y_i - \mathbb{E}[Y_i|Q_i, X_i]$ and $v_i = \varepsilon_i - u_i\beta^{\text{SIVE}}$. We propose the following variance estimator.

$$\hat{\mathbb{V}}[\hat{\beta}^{\text{SIVE}}|Q, X] = \frac{(Y - T\hat{\beta}^{\text{SIVE}})'AD_{\hat{\sigma}_u^2}A(Y - T\hat{\beta}^{\text{SIVE}}) + T'AD_{\hat{\sigma}_v^2}AT + 2(Y - T\hat{\beta}^{\text{SIVE}})'AD_{\hat{\sigma}_{uv}}AT}{(T'AT)^2}, \quad (20)$$

where $A = P - M_{W,Z}DM_{W,Z}$ with D as defined in Lemma 5, and $D_{\hat{\sigma}_u^2}$, $D_{\hat{\sigma}_v^2}$ and $D_{\hat{\sigma}_{uv}}$ are diagonal matrices with $[D_{\hat{\sigma}_u^2}]_{ii} = \hat{\sigma}_{u,i}^2$, $[D_{\hat{\sigma}_v^2}]_{ii} = \hat{\sigma}_{v,i}^2$, $[D_{\hat{\sigma}_{uv}}]_{ii} = \hat{\sigma}_{uv,i}$ on their respective diagonals. The testing procedure is completed by defining the estimators $[\hat{\sigma}_u^2]_i$, $[\hat{\sigma}_v^2]_i$, and $[\hat{\sigma}_{uv}]_i$, for $\sigma_{u,i}^2 = \mathbb{E}[u_i^2|Q_i, X_i]$, $\sigma_{v,i}^2 = \mathbb{E}[v_i^2|Q_i, X_i]$, and $\sigma_{uv,i} = \mathbb{E}[u_i v_i|Q_i, X_i]$, respectively. In the presence of many instruments, standard heteroskedasticity robust Eicker-Huber-White variance estimators are inconsistent (Cattaneo et al., 2018). We therefore consider the Hartley et al. (1969) variance estimators, also discussed in the previous section:

$$\begin{aligned} \hat{\sigma}_{u,i}^2 &= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}T \odot M_{W,Z}T), \\ \hat{\sigma}_{v,i}^2 &= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}}) \odot M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}})), \\ \hat{\sigma}_{uv,i} &= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}}) \odot M_{W,Z}T). \end{aligned} \quad (21)$$

We show that when $\hat{\beta}^{\text{SIVE}}$ is replaced by its population counterpart β^{SIVE} , the estimators are unbiased conditional on the instrument and covariates. This removes the main driver of the inconsistency of standard heteroskedasticity robust variance estimators. We show below that when using (21) and (20) in the test (18) leads to a conservative test under weak identification.

One issue with the estimators in (21) is that they require a strengthening of Assumption 2 to $m_g \geq 3$ and $n_g - m_g \geq 3$ for the inverse of $M_{W,Z} \odot M_{W,Z}$ to exist. Instead, we can use the following estimators on the individuals with an instrument status that is only shared with one other individual in the same covariate group:

$$\begin{aligned} \hat{\sigma}_{u,i}^2 &= 4e_i'(M_{W,Z}T \odot M_{W,Z}T), \\ \hat{\sigma}_{v,i}^2 &= 4e_i'(M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}}) \odot M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}})), \\ \hat{\sigma}_{uv,i} &= 4e_i'(M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}}) \odot M_{W,Z}T). \end{aligned} \quad (22)$$

These estimators can be shown to generate a positive bias in the variance estimator (20). Hence, the presence of many small groups will make the inference procedure more conservative.

3.3 Assumptions

To study the asymptotic properties of the test statistic in (18) with variance estimator (20), we impose the following assumptions. Throughout, C denotes a generic positive constant that can differ between occurrences.

Assumption 3

The error terms (u_i, ε_i) are independent across i , conditionally on Q and X , and for all i it holds that almost surely

1. $\mathbb{E}[Y_i|Q_i, X_i]$ is almost surely bounded.
2. $\mathbb{E}[u_i|Q, X] = 0$ and $\mathbb{E}[\varepsilon_i|Q, X] = 0$.
3. $\mathbb{E}[u_i^2|Q, X] \geq C > 0$ and $\mathbb{E}[\varepsilon_i^2|Q, X] \geq C > 0$ for some positive constant C , and $|\text{corr}[u_i, \varepsilon_i|Q, X]| \leq C < 1$.
4. $\mathbb{E}[u_i^8|Q, X] \leq C < \infty$ and $\mathbb{E}[\varepsilon_i^8|Q, X] \leq C < \infty$.

Assumption 3 part 1 ensures that the treatment effect is bounded for all individuals. Part 2 is a standard assumption on the residuals in the reduced form model. Part 3 ensures that the distribution of the test statistic is non-degenerate. Part 4 is used to control the behavior of the estimators for the conditional variances of u_i and ε_i .

Assumptions 1 to 3 allow for the inference on the weighted average of LATEs in (1) in a wide range of empirically relevant settings. First, the treatment effects $Y_i(1) - Y_i(0)$ are allowed to be heterogeneous across all i . It follows that treatment effects may vary across the covariate groups defined by the elements of \mathbb{X} , and hence SIVE identifies a weighted average of potentially heterogeneous conditional treatment effects. This is in line with the LATE identification literature. Note that the literature on inference in IV models often assumes that the data satisfies a model along the lines of

$$T_i = \pi' Z_i + \delta_1' W_i + \varepsilon_i, \quad Y_i = T_i \beta + \delta_2' W_i + u_i, \quad (23)$$

in which the treatment effect of T_i on Y_i is modelled with β which is specified to be homogeneous. Since we conduct inference with saturated instruments and covariates, no parametric assumptions or a model specification is required.

Second, the first stage and reduced form errors $u_i = T_i - \mathbb{E}[T_i|Q_i, X_i]$ and $\varepsilon_i = Y_i - \mathbb{E}[Y_i|Q_i, X_i]$ are allowed to be heteroskedastic. That is, $\mathbb{E}[u_i^2|Q, X]$, $\mathbb{E}[\varepsilon_i^2|Q, X]$, and $\text{corr}[u_i, \varepsilon_i|Q, X]$ may differ across i . Although heteroskedasticity has been accounted for in existing methods with many instruments (e.g. Hausman et al. (2012), Mikusheva and Sun (2022), Crudu et al. (2021)), this is usually within IV models similar to (23) that restrict treatment effect heterogeneity.

Third, inference based on our test statistic in (18) is robust to weak identification under a minimal assumption on the identification strength. That is, the test works with sets of instrument interactions that have a strong or a weak signal. The instrument interactions are referred to as strong when the concentration parameter

$$\mu_n = \frac{n}{G} \text{FS} \rightarrow \infty \text{ with FS} = \sum_g \tilde{\mathbb{P}}[X_i = x_g] \pi(x_g)^2 \tilde{\mathbb{V}}[Q_i|X_i = x_g]. \quad (24)$$

Since $|\pi(x_g)| \leq 1$ and $\text{Var}[Q_i|X_i = x_g] \leq \frac{1}{4}$, in the scenario with the strongest identification possible within each covariate group, we require $\frac{n}{4} \gg G$ for the instrument interaction set to be considered

as strong. In this case, we show that the SIVE estimator is consistent, its variance estimator is consistent, and the asymptotic confidence intervals constructed from the test statistic attain nominal coverage. Under weak identification, that is if as $G \rightarrow \infty$,

$$\sqrt{G}\mu_n \rightarrow \infty, \quad (25)$$

the SIVE estimator remains consistent. Mikusheva and Sun (2022) show that this is the weakest identification strength under which a consistent estimator exists. In this case, the fact that we take into account treatment effect heterogeneity leads to a positive asymptotic bias in the variance estimator. This is common also in the case when the number of instruments is small (Kleibergen and Zhan, 2021). The positive bias means that the asymptotic confidence intervals may be conservative. In the case that one is worried that the identification is even weaker so that $\sqrt{G}\mu_n \rightarrow C < \infty$ and $G \rightarrow \infty$, the results we provide allow for the construction of fully identification robust confidence intervals by inverting a score statistic as in Kleibergen (2005).

3.4 Large sample theory

We first provide a consistency result for $\hat{\beta}^{\text{SIVE}}$.

Lemma 6 (Consistency SIVE) *Under Assumption 2 and 3 with $\frac{n}{\sqrt{G}}\text{FS} \rightarrow_p \infty$, it holds that $\hat{\beta}^{\text{SIVE}} \rightarrow_p \beta^{\text{SIVE}}$.*

The proof is deferred to Appendix D.2. The most stringent condition on the identification strength occurs when $G \propto n$, in which case we require that $\sqrt{n}\text{FS} \rightarrow_p \infty$. Note that this allows the first stage FS to decrease to zero asymptotically, but it limits the rate at which it can decrease to zero.

The following result shows the asymptotic validity of a standard t -test that uses the variance estimator from (20). What is particularly important to note is that the theorem does not limit the rate at which G can grow with n . In particular, we allow $G/n \rightarrow \alpha \in (0, 1)$ which are the many-instrument sequences by Bekker (1994). In our setting G can never exceed n because G is governed by the number of unique observations on the vector of controls.

Theorem 2 *Let Assumption 2 and 3 hold. If $\frac{n}{G}\text{FS} \rightarrow_{a.s.} \infty$, or $\frac{n}{\sqrt{G}}\text{FS} \rightarrow_{a.s.} \infty$ and $G \rightarrow \infty$, then,*

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\frac{\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}}{\sqrt{\hat{\mathbb{V}}[\hat{\beta}^{\text{SIVE}}|Q, X]}} \leq \Phi^{-1}(1 - \alpha) \right) \geq 1 - \alpha, \quad (26)$$

with $\hat{\mathbb{V}}[\hat{\beta}^{\text{SIVE}}]$ defined in (20).

The proof is deferred to Appendix D.3. Part of the result follows from a central limit theorem for quadratic forms first derived by Chao et al. (2012). We use a version by Evdokimov and Kolesár (2018) that allows us to efficiently verify the necessary conditions for the central limit theorem to apply. The most challenging result to establish is that the variance estimator converges to a quantity at least as large as the population variance, conditional on the covariates X and instrument Q . The requirement that $\frac{n}{\sqrt{G}}\text{FS} \rightarrow_{a.s.} \infty$ ensures that we can apply Lemma 6 when analyzing the variance estimator in (20).

Theorem 2 shows that tests and confidence intervals based on the proposed procedure will be conservative. This property is due to the fact that the variance estimator allows for treatment effect heterogeneity, and has also been found in the setting with a fixed number of instruments (Kleibergen

and Zhan, 2021). A natural question is then how conservative tests and confidence intervals based on the SIVE estimator actually are. To quantify this, we have the following result.

Corollary 1 *Strengthen Assumption 2 to $m_g \geq 3$ and $n_g - m_g \geq 3$ and let Assumption 3 hold. Then, If $\frac{n}{G}FS \rightarrow_{a.s.} \infty$, or $\frac{n}{\sqrt{G}}FS \rightarrow_{a.s.} \infty$ and $G \rightarrow \infty$,*

$$\frac{\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}}{\sqrt{\hat{\mathbb{V}}[\hat{\beta}^{\text{SIVE}}|Q, X]}} \rightarrow_d N(0, \lambda), \quad (27)$$

with $\hat{\mathbb{V}}[\hat{\beta}^{\text{SIVE}}|Q, X]$ as in (20) and where $\lambda = 1$ when $\frac{n}{G}FS \rightarrow_{a.s.} \infty$ (strong identification) and $\lambda \in [1/4, 1]$ when $\frac{n}{G}FS \rightarrow_{a.s.} \mu \in [0, \infty)$ (weak identification) and λ is increasing in μ .

Corollary 1 follows from the proof of Theorem 2. The strengthening of Assumption 2 shuts down one source of positive bias in the variance estimator (20) that is due to the existence of covariate groups with only two individuals with a particular instrument status for which we use (22) to estimate the error variances. When excluding those groups, we see that as the identification strength increases the rejection rates and coverage probabilities of the t -test attain the nominal values. In terms of the confidence intervals, in a very weakly identified model, the confidence intervals are twice as wide as needed to achieve the nominal size.

A second question regarding Theorem 2 concerns settings in which the identification may be even weaker than that required by the theorem. These concerns can be mitigated by constructing an identification-robust procedure. If we replace in (20) the estimator $\hat{\beta}^{\text{SIVE}}$ by β_0 , then Theorem 2 holds under $H_0: \beta^{\text{SIVE}} = \beta_0$ without the requirement that $\frac{n}{G}FS \rightarrow_{a.s.} \infty$ and only requires $G \rightarrow \infty$. The following results formalizes that our testing procedure is fully identification-robust with the asymptotic rejection rate not exceeding the nominal rate.

Corollary 2 *Let Assumption 2 and 3 hold. When $G \rightarrow \infty$,*

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\frac{\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}}{\sqrt{\hat{\mathbb{V}}[\beta^{\text{SIVE}}|Q, X]}} \leq \Phi^{-1}(1 - \alpha) \right) \geq 1 - \alpha, \quad (28)$$

with $\hat{\mathbb{V}}[\beta^{\text{SIVE}}|Q, X]$ as in (20) with $\hat{\beta}^{\text{SIVE}}$ replaced by β^{SIVE} .

Corollary 1 follows from the proof of Theorem 2. As usual, confidence intervals can now be constructed using test inversion. Finally, the analogous result to Corollary 1 can be established.

4 Monte Carlo Study

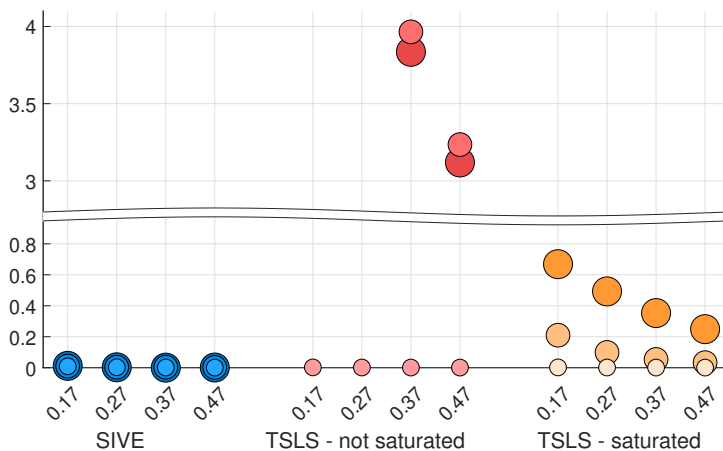
We build on the Monte Carlo set-up considered in Blandhol et al. (2022) that is designed to match some key features of the data used by Card (1995). Consider a sample with $n = 3,000$ observations. We have a single control variable X_i that can take on $L = \{1, 2, 25, 50, 100, 200, 300\}$ values taken from a one-dimensional Halton sequence. The binary instrument satisfies

$$\mathbb{E}[Q_i|X_i] = 0.119 + 1.785X_i - 1.534X_i^2 + 0.597X_i^3. \quad (29)$$

We generate $(u_i, \varepsilon_i)' \sim N(0, \Sigma)$ where $[\Sigma]_{11} = [\Sigma]_{22} = 1$ and $[\Sigma]_{12} = 0.527$. The endogenous treatment and the outcome variable are generated as,

$$\begin{aligned} T_i &= \Phi(u_i) \leq p_0 1[Z_i = 0] + p_1 1[Z_i = 1], \\ Y_i &= \log(129.7 + 1247.7X_i - 2149.0X_i^2 + 1515.7X_i^3) + \beta(\gamma_i T_i) + \varepsilon_i. \end{aligned} \quad (30)$$

Figure 1: Average absolute bias in the estimand



Note: the figure shows the absolute median difference with the causal estimand in a setting without treatment heterogeneity. The size of the circles indicates the number of covariate groups with the small circle corresponding to $L = 1$, the medium circle corresponding to $L = 25$ and the large circle corresponding to $L = 300$. The x -axis is the instrument strength $p(1) - p(0)$, with $p(0) = 0.22$ and $p(1) = \{0.39, 0.49, 0.59, 0.69\}$. Because non-saturated TSLS shows large biases for $L = 25$ and $L = 300$, the y -axis is broken between 0.9 and 3 and limited to $[0, 4]$.

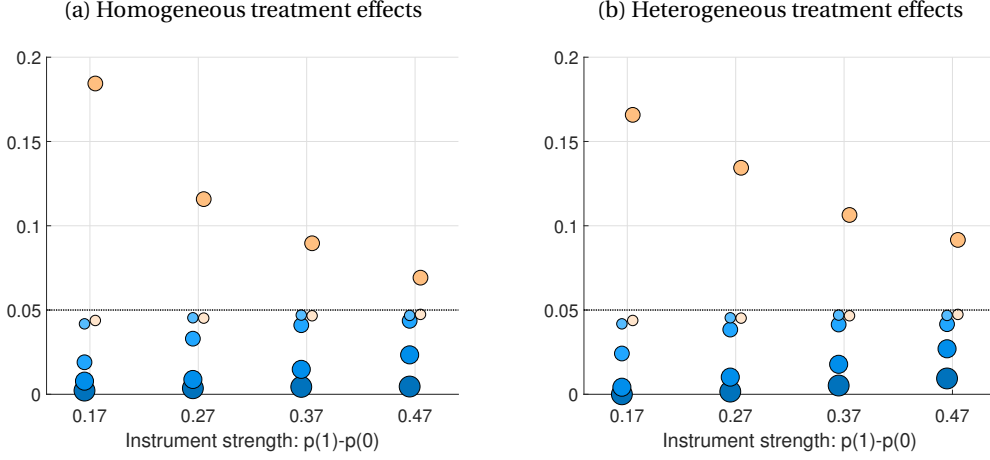
where $\gamma_i = 1$ in the homogeneous treatment effect design and $\beta = 0.2$. We set $p_0 = 0.22$ and vary the identification strength through $p_1 \in \{0.29, 0.39, \dots, 0.69\}$. To simulate heterogeneous treatment effects, we set $\gamma_i = 1 + h$ with $h \in \{2, 4, \dots, 10\}$ for the 900 observations with the smallest value of X_i and $\gamma_i = 1$ for the remaining. Throughout, we only use covariate groups that satisfy [Assumption 2](#).

Bias We first study the bias of the following estimators: SIVE, nonsaturated TSLS and saturated TSLS. Results for the saturated JIVE estimators considered in [Section 2.4](#) are reported in the appendix. [Figure 1](#) shows the absolute median bias of the various estimators as a function of the instrument strength in the absence of treatment heterogeneity. The small circles correspond to $L = 1$ covariate groups, the medium circles to $L = 25$ groups and the large circles to $L = 300$ groups. We analyze the effects of varying the instrument strength by setting $p(2) = \{0.39, 0.49, 0.59, 0.69\}$.

We see that SIVE is median unbiased regardless of the number of covariate groups and the values of $p(2)$ under consideration. For TSLS, we see that it is median unbiased for $L = 1$ covariate group. However, as the number of covariate groups increases, non-saturated TSLS incurs a bias because estimand is non-causal. For saturated TSLS the (many instrument) bias enters. As is well known, this bias is more pronounced in settings where the instruments are weak. In the appendix we find that for JIVE1 removing the diagonal of the projection matrix leads to a large (omitted variable) bias as the controls are no longer correctly projected out. As we have shown in [Section 2.4](#) this effect is mitigated by moving to the JIVE2 estimator. However, this estimator also shows an increasing bias with increasing number of covariate groups and an decreasing instrument strength.

Size In [Figure 2](#) we show that the size of a test of $H_0: \beta_0 = \beta$ with $\beta = 0.2$ in a setting without treatment heterogeneity (left panel) and with treatment heterogeneity (right panel). In the latter case, in [\(30\)](#) we set $\gamma_i = 11$ for the 900 observations with the lowest values of X_i and $\gamma_i = 1$ for the

Figure 2: Size versus instrument strength for an increasing number of instruments



Note: the figure shows the size of testing $H_0: \beta = 0.2$ at a nominal level of 5%. The x -axis is the instrument strength, with $p(0) = 0.22$ and $p(1) = \{0.39, 0.49, 0.59, 0.69\}$. The circles of increasing size correspond to $L = \{1, 25, 100, 300\}$. TSLS is fully saturated and we use heteroskedasticity-robust (HC0) standard errors to construct the t -statistic. SIVE uses standard errors based on (20). The observed size for TSLS when $L = \{100, 300\}$ is above 0.2.

remaining. The size of the circles indicates the number of covariate groups, which we choose as $L = \{1, 25, 100, 300\}$. Given the large bias in non-saturated TSLS observed in Figure 1, we now only consider saturated TSLS. As expected TSLS yields accurate size control when $L = 1$. For $L = 25$, we have seen a substantial bias in Figure 1 and consequently we observe a size distortion that is increasing with decreasing instrument strength. For SIVE, we obtain close to nominal size control for $L = 1$ for all values of the instrument strength. As expected based on the theory we see that for a larger number of covariate groups, the test becomes progressively more conservative. Increasing the instrument strength makes the test less conservative, as formalized in Corollary 1. The results with and without treatment heterogeneity do not show any qualitative differences.

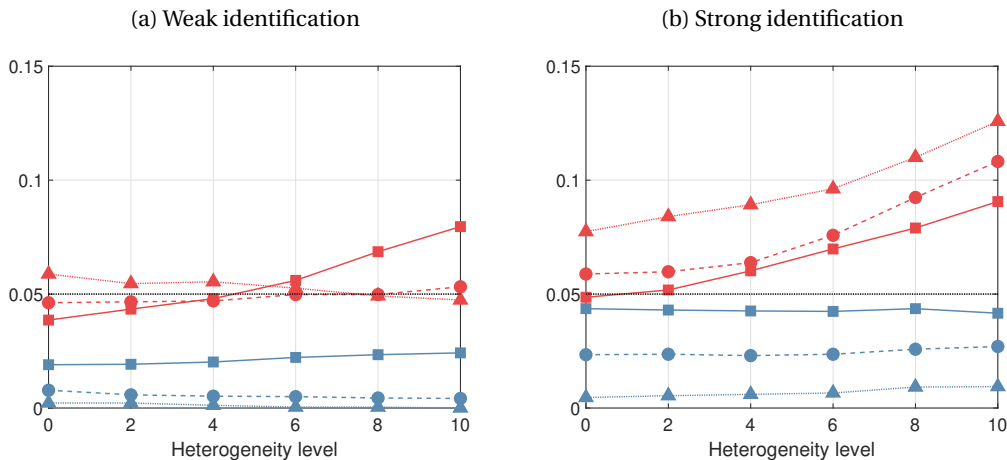
Alternative variance estimators. For SIVE, we use the variance estimator given in (20), which is robust to treatment effect heterogeneity. We first compare this with using the variance estimator proposed by Chao et al. (2023) that is given by

$$V_c = (T'AD_1AT + (\hat{\varepsilon} \odot \hat{u})'J(A \odot A)J(\hat{\varepsilon} \odot \hat{u})) / (T'AT)^2, \quad (31)$$

where $J = (M_W \odot M_W)^{-1}$, $[D_1]_i = [J(\hat{\varepsilon} \odot \hat{\varepsilon})]_i$, $\hat{\varepsilon} = M_{W,Z}(y - T\hat{\beta}^{\text{SIVE}}) = M_{W,Z}(\varepsilon - u\hat{\beta}^{\text{SIVE}})$ and $\hat{u} = M_{W,Z}T = M_{W,Z}u$. The variance estimator (31) is proposed in the context of a linear panel data model with homogeneous slope coefficient. We now assess the effect of treatment effect heterogeneity on tests that rely on (31).

In Figure 3 we show the size of the test of $H_0: \beta_0 = \beta$ with $\beta = 0.2$ in a setting with weak instruments ($p(2) = 0.39$, left panel) and with strong instruments ($p(2) = 0.69$, right panel). On the x -axis we vary the treatment effect heterogeneity through the parameter $\gamma_i = 1 + h$ with $h \in \{2, 4, \dots, 10\}$ for the 900 observations with the smallest value of X_i and $\gamma_i = 1$ for the remaining. In the left panel, we again observe that SIVE offers a conservative test under weak instruments. As expected based on the theory, the level of heterogeneity has no effect on size of the test. For the alternative variance estimator, we see that increasing the level of treatment effect heterogeneity leads to a slightly oversized

Figure 3: Size versus treatment effect heterogeneity: compared to Chao et al. (2023).



Note: the figure shows the size of testing $H_0: \beta = \beta^{\text{SIVE}}$ at a nominal level of 5%. The left panel is for weak instruments, $p(2) - p(1) = 0.17$, the right panel for strong instruments $p(2) - p(1) = 0.47$. The x -axis is the heterogeneity level h . In (30), we set $\gamma_i = 1 + h$ with $h \in \{0, 2, \dots, 10\}$ for the 900 observations with the smallest value of X_i and $\gamma_i = 1$ for the remaining. SIVC uses the SIVE estimator, but the variance estimator (31) proposed by Chao et al. (2023). We consider $L = \{25, 100, 300\}$ for the number of covariate groups, which correspond to the solid, dashed and dotted lines respectively.

test, but no ordering in terms of the number of covariate groups is observed. The fact that treatment effect heterogeneity has only a mild effect in this case is due to the fact that the heterogeneity is flooded by the additional uncertainty introduced by the presence of many weak instruments. When the instruments are strong, as in the right panel of Figure 3, accounting for treatment effect heterogeneity becomes more important. Again, SIVE shows no dependence on the level of treatment effect heterogeneity. The alternative variance estimator now become progressively oversized as the level of heterogeneity increases for all values of the number of covariate groups.

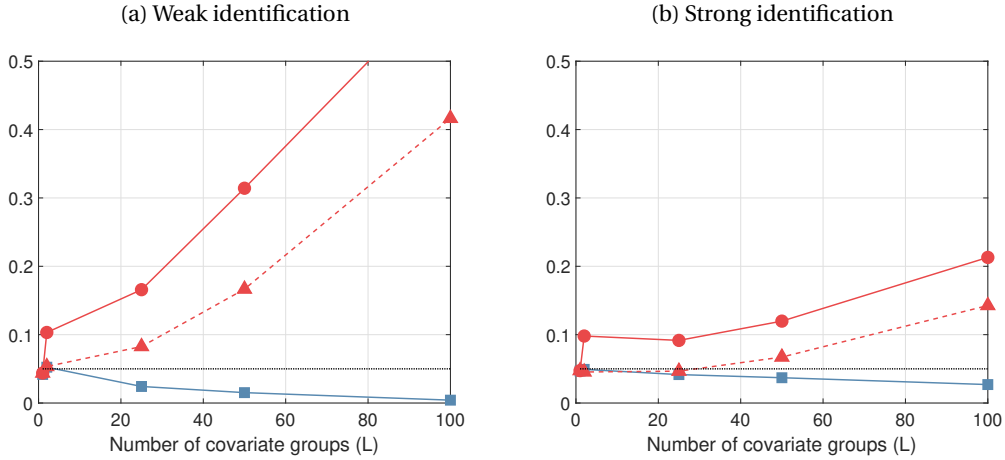
Lee (2018) proposes a variance estimator for TSLS that is valid in overidentified systems where each instrument identifies a different LATE. It therefore allows for treatment effect heterogeneity. However, the analysis in Lee (2018) proceeds under the assumption that the number of instruments is fixed relative to the sample size. We therefore study its performance in a setting where the number of covariate groups L increases.

In Figure 4, the left panel shows a setting with weak instruments $p(2) - p(1) = 0.17$, the right panel a setting with strong instruments $p(2) - p(1)$. The variance estimator in Lee (2018) is designed for the right panel with a number of covariate groups that is small relative to the sample size. Indeed, for $L = \{2, 25\}$ the variance estimator (dashed line, triangle markers) shows excellent size control and a substantial improvement over the standard Eicker-Huber-White variance estimator (solid line, circle markers). However, when the number of covariate groups increases, the many instrument bias manifests itself and the size increases above its nominal value. This effect is even stronger in the weakly identified setting in the left panel.

5 Application: Card (1995)

As an illustration of the proposed estimator and inference procedures we revisit the study by Card (1995) that uses the distance to the nearest college to instrument educational attainment. The data

Figure 4: Size versus covariate groups: compared to Lee (2018).



Note: the figure shows the size of testing $H_0: \beta = \beta^{\text{SIVE}}$ at a nominal level of 5%. The x-axis is the number of covariate groups L . The left panel is for weak instruments, $p(2) - p(1) = 0.17$, the right panel for strong instruments $p(2) - p(1) = 0.47$. In (30), we set $\gamma_i = 11$ for the 900 observations with the smallest value of X_i and $\gamma_i = 1$ for the remaining. The red solid line with circle marker is for TSLS, the red dashed line with triangle marker uses the variance estimator from Lee (2018), the blue solid line with square marker is for SIVE.

considers men aged 14-24 sampled in 1966 from the National Longitudinal Survey of Young Men (NLSYM). These men were followed until 1981. Following Card (1995), we consider individuals that provided education and wage information when they were interviewed in 1976.

The instrument used by Card (1995) is the distance to the nearest four-year college. We consider some adjustments to the original model as proposed by Kitagawa (2015) and Słoczyński (2020). In particular, the specification includes five binary controls (Black, living in a metropolitan area (SMSA) in 1966, living in a metropolitan area (SMSA) in 1976, living in the South in 1966, living in the South in 1976). With these five binary controls, we potentially have 32 covariate groups after saturation. The original sample size is 3,010. We restrict the sample by requiring at least five observations in each covariate group, which brings the sample size to 2,988. In each of the covariate group we have at least two treated individuals and two non-treated individuals. Finally, we follow Kitagawa (2015) and Słoczyński (2020) and redefine the instrument to equal 1 if individuals have some college attendance (defined as having strictly more than 12 years of education) and 0 otherwise.

We consider the TSLS and SIVE estimators under different specifications for the controls and the instruments. First, the standard TSLS estimator that uses the binary instrument and linearly includes the controls. This estimator is inconsistent if the assumption of a linear relation with the controls is violated. Moreover, it supposes strong monotonicity in the instrument. We then saturate the model in the controls. To allow for weak monotonicity interact these controls with the instruments. This estimator is inconsistent due to the many instrument bias. We consider then two restricted versions of the fully saturated TSLS estimator. First, we only saturate in the controls, assuming strong monotonicity such that saturation in the instrument is not necessary. Second, we only saturate in the instrument, assuming a linear relation with the controls so that we do not need to saturate the controls. We follow the same specifications for SIVE with the exception of the model without any saturation.

Table 2 shows the point estimates, standard errors and 95% confidence intervals. For TSLS these

Table 2: Empirical application: estimates, standard errors and confidence intervals

$(m_g, n_g - m_g)$	Estimator	Specification	Estimate	Standard error	95% CI
≥ 2	2SLS	not saturated	0.524	0.296	[-0.056, 1.104]
		fully saturated	0.156	0.138	[-0.116, 0.427]
		saturated instruments	0.209	0.102	[0.009, 0.408]
		saturated controls	0.570	0.298	[-0.014, 1.154]
	SIVE	fully saturated	0.125	0.342	[-0.546, 0.795]
		saturated instruments	0.217	0.171	[-0.119, 0.553]
≥ 3	2SLS	not saturated	0.499	0.278	[0.041, 0.957]
		fully saturated	0.190	0.139	[-0.038, 0.417]
		saturated instruments	0.218	0.106	[0.044, 0.392]
		saturated controls	0.538	0.282	[0.074, 1.001]
	SIVE	fully saturated	0.215	0.273	[-0.234, 0.664]
		saturated instruments	0.233	0.159	[-0.079, 0.545]
		saturated controls	0.599	0.388	[-0.040, 1.237]

are based on the HC0 based estimator for the variance covariance matrix. We reproduce the key findings by [Słoczyński \(2020\)](#). In particular, not interacting the instruments leads to unreasonably large estimates for the effect of schooling both when using TSLS and SIVE. When we saturate the instrument, the point estimates drop from around 0.5 to 0.2, which is much more in line with the recent literature on wage gains resulting from education. In terms of statistical efficiency, we see that the standard errors from SIVE are generally higher. This is to be expected as it takes into account the many instrument effect, as well as treatment effect heterogeneity.

If we have individuals that share the treatment status with only one other individual in the covariate group, we use the estimators from (22) that leads to an upward bias in the variance estimator (20). To analyze whether these individuals drive the results, we remove those individuals from the data. This reduces the sample size to 2,957 individuals. the bottom panel of [Table 2](#) shows the point estimates, standard errors and confidence intervals of the methods. The main finding of interest is that the point estimate from the fully saturated SIVE slightly increases and is almost equal to that of the SIVE estimator that only saturates the instrument. Despite the smaller sample size, the standard errors are also somewhat lower.

6 Conclusion

We show how to conduct inference in a saturated IV model where the number of covariate values is of the same order as the sample size. The estimator is consistent under a minimal assumption on the identification strength. Crucially, and unlike existing procedures we allow arbitrary treatment effect heterogeneity. Our findings are confirmed through numerical experiments that rely on data with similar characteristics to the data analyzed by [Card \(1995\)](#). Applying the proposed estimator to that data yields realistic point estimates.

References

- Abadie, A. (2003). Semiparametric instrumental variable estimation of treatment response models. *Journal of Econometrics* 113(2), 231–263.
- Ackerberg, D. A. and P. J. Devereux (2009). Improved JIVE estimators for overidentified linear models with and without heteroskedasticity. *The Review of Economics and Statistics* 91(2), 351–362.
- Angrist, J. and M. Kolesár (2023). One instrument to rule them all: The bias and coverage of just-identified. *Journal of Econometrics* 0(just-accepted), 1–18.
- Angrist, J. D. and G. W. Imbens (1995). Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American Statistical Association* 90(430), 431–442.
- Angrist, J. D., G. W. Imbens, and A. B. Krueger (1999). Jackknife instrumental variables estimation. *Journal of Applied Econometrics* 14(1), 57–67.
- Angrist, J. D. and A. B. Krueger (1991). Does compulsory school attendance affect schooling and earnings? *The Quarterly Journal of Economics* 106(4), 979–1014.
- Angrist, J. D. and J.-S. Pischke (2009). *Mostly harmless econometrics: An empiricist's companion*. Princeton University Press.
- Bekker, P. A. (1994). Alternative approximations to the distributions of instrumental variable estimators. *Econometrica* 62(3), 657–681.
- Bekker, P. A. and F. Cruadu (2015). Jackknife instrumental variable estimation with heteroskedasticity. *Journal of Econometrics* 185(2), 332–342.
- Bekker, P. A. and F. Kleibergen (2003). Finite-sample instrumental variables inference using an asymptotically pivotal statistic. *Econometric Theory* 19(5), 744–753.
- Blandhol, C., J. Bonney, M. Mogstad, and A. Torgovitsky (2022). When is TSLS actually LATE? Working paper, National Bureau of Economic Research.
- Card, D. (1995). Using geographic variation in college proximity to estimate the return to schooling. In L. Christofides, E. Grant, and R. Swidinsky (Eds.), *Aspects of Labor Market Behaviour: Essays in Honour of John Vanderkamp*. University of Toronto Press.
- Cattaneo, M. D., M. Jansson, and W. K. Newey (2018). Inference in linear regression models with many covariates and heteroscedasticity. *Journal of the American Statistical Association* 113(523), 1350–1361.
- Chao, J. C. and N. R. Swanson (2005). Consistent estimation with a large number of weak instruments. *Econometrica* 73(5), 1673–1692.
- 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.

- Chao, J. C., N. R. Swanson, and T. Woutersen (2023). Jackknife estimation of a cluster-sample IV regression model with many weak instruments. *Journal of Econometrics* 235(2), 1747–1769.
- Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* 21(1), C1–C68.
- Crudu, F., G. Mellace, and Z. Sándor (2021). Inference in instrumental variable models with heteroskedasticity and many instruments. *Econometric Theory* 37(2), 281–310.
- Crump, R. K., V. J. Hotz, G. W. Imbens, and O. A. Mitnik (2009). Dealing with limited overlap in estimation of average treatment effects. *Biometrika* 96(1), 187–199.
- Dube, O. and S. Harish (2020). Queens. *Journal of Political Economy* 128(7), 2579–2652.
- Evdokimov, K. S. and M. Kolesár (2018). Inference in instrumental variables analysis with heterogeneous treatment effects. Working paper, Princeton University.
- Evdokimov, K. S. and D. Lee (2013). Diagnostics for exclusion restrictions in instrumental variables estimation. Working paper, Princeton University.
- Frölich, M. (2007). Nonparametric IV estimation of local average treatment effects with covariates. *Journal of Econometrics* 139(1), 35–75.
- Gelbach, J. B. (2002). Public schooling for young children and maternal labor supply. *American Economic Review* 92(1), 307–322.
- Hansen, C., J. Hausman, and W. Newey (2008). Estimation with many instrumental variables. *Journal of Business & Economic Statistics* 26(4), 398–422.
- Hartley, H., J. Rao, and G. Kiefer (1969). Variance estimation with one unit per stratum. *Journal of the American Statistical Association* 64(327), 841–851.
- Hausman, J. A., W. K. Newey, T. Woutersen, J. C. Chao, and N. R. Swanson (2012). Instrumental variable estimation with heteroskedasticity and many instruments. *Quantitative Economics* 3(2), 211–255.
- Imbens, G. W. and J. D. Angrist (1994). Identification and estimation of local average treatment effects. *Econometrica* 62(2), 467–475.
- Kitagawa, T. (2015). A test for instrument validity. *Econometrica* 83(5), 2043–2063.
- Kleibergen, F. (2005). Testing parameters in GMM without assuming that they are identified. *Econometrica* 73(4), 1103–1123.
- Kleibergen, F. and Z. Zhan (2021). Double robust inference for continuous updating GMM. Working paper, arXiv preprint arXiv:2105.08345.
- Kolesár, M. (2013). Estimation in an instrumental variables model with treatment effect heterogeneity. Working paper, Princeton University.
- Lee, S. (2018). A consistent variance estimator for 2SLS when instruments identify different LATEs. *Journal of Business & Economic Statistics* 36(3), 400–410.

- Lim, D., W. Wang, and Y. Zhang (2024). A conditional linear combination test with many weak instruments. *Journal of Econometrics* 283(2), 105602.
- Matsushita, Y. and T. Otsu (2022). A jackknife Lagrange multiplier test with many weak instruments. *Econometric Theory* 0(just-accepted), 1–24.
- Mikusheva, A. and L. Sun (2022). Inference with many weak instruments. *Review of Economic Studies* 89(5), 2663–2686.
- Mikusheva, A. and L. Sun (2023). Weak identification with many instruments. Working paper, arXiv preprint arXiv:2308.09535.
- Słoczyński, T. (2020). When should we (not) interpret linear IV estimands as LATE? Working paper, arXiv preprint arXiv:2011.06695.

Appendix to “Inference on LATEs with covariates”

A Conventions and notation

Without loss of generality, we assume that the observations are ordered according to the columns of the matrix W that contains the dummies indicating the values of the control variate(s) in the sense that

$$W = \begin{pmatrix} \iota_{n_1} & \mathbf{0}_{n_1} & \cdots & \mathbf{0}_{n_1} \\ \mathbf{0}_{n_2} & \iota_{n_2} & \cdots & \mathbf{0}_{n_2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{n_G} & \mathbf{0}_{n_G} & \cdots & \iota_{n_G} \end{pmatrix}. \quad (32)$$

As a subsequent ordering, we assume again without loss of generality that the matrix with instrument interactions has the following structure

$$Z = \begin{pmatrix} \iota_{m_1} & \mathbf{0}_{m_1} & \cdots & \mathbf{0}_{m_1} \\ \mathbf{0}_{n_1-m_1} & \mathbf{0}_{n_1-m_1} & \cdots & \mathbf{0}_{n_1-m_1} \\ \mathbf{0}_{m_2} & \iota_{m_2} & \cdots & \mathbf{0}_{m_2} \\ \mathbf{0}_{n_2-m_2} & \mathbf{0}_{n_2-m_2} & \cdots & \mathbf{0}_{n_2-m_2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{m_G} & \mathbf{0}_{m_G} & \cdots & \iota_{m_G} \\ \mathbf{0}_{n_G-m_G} & \mathbf{0}_{n_G-m_G} & \cdots & \mathbf{0}_{n_G-m_G} \end{pmatrix}. \quad (33)$$

Throughout we denote by $r_n = G\mu_n$ under strong identification and $r_n = G$ under weak identification as defined in the main paper.

For any vector v , $v_{(g)}$ is the vector of observations from v that have $W_{ig} = 1$. Additionally, $v_{(g,1)}$ is the vector of observations that have $Z_{ig} = 1$ and $v_{(g,2)}$ is the vector of observations that have $W_{ig} = 1$ and $Z_{ig} = 0$.

B Generalizations

This section extends the binary treatment and binary instrument setting to a multivalued treatment and instrument using respectively Theorem 1 and Theorem 2 in (Angrist and Imbens, 1995).

B.1 Multivalued treatment

Suppose the multivalued treatment T_i takes values in the set $\{0, 1, 2, \dots, J\}$, where $T_i = 0$ corresponds to no treatment, and $T_i = 1, \dots, J$ correspond to different treatment levels. Define $D_{ij} = \mathbb{1}[T_i \geq j]$ with $D_{ij} = (1 - Q_i)D_{ij}(0) + Q_i D_{ij}(1)$. Under [Assumption 1.2](#), we have

$$\begin{aligned} Y_i &= \sum_{j=0}^J (D_{i,j} - D_{i,j+1}) Y_i(j) \\ &= (1 - Q_i) \sum_{j=0}^J (D_{i,j}(0) - D_{i,j+1}(0)) Y_i(j) + Q_i \sum_{j=0}^J (D_{i,j}(1) - D_{i,j+1}(1)) Y_i(j), \end{aligned} \quad (34)$$

where we use that $D_{i0} = 1$ and $D_{i,J+1} = 0$. It follows from [Assumption 1.1](#) that

$$\begin{aligned}\theta(x_g) &= \mathbb{E}[Y_i|Q_i = 1, X_i = x_g] - \mathbb{E}[Y_i|Q_i = 0, X_i = x_g] \\ &= \sum_{j=0}^J \mathbb{E}[(D_{i,j}(1) - D_{i,j+1}(1) - D_{i,j}(0) + D_{i,j+1}(0))Y_i(j)|X_i = x_g] \\ &= \sum_{j=1}^J \mathbb{E}[(D_{i,j}(1) - D_{i,j}(0))(Y_i(j) - Y_i(j-1))|X_i = x_g].\end{aligned}\tag{35}$$

From [Assumption 1.4](#) follows that either $T_i(1) \geq T_i(0)$ or $T_i(1) \leq T_i(0)$ for all i with $X_i = x_g$. Hence, we either have $D_{i,j}(1) = D_{i,j}(0)$, $D_{i,j}(1) > D_{i,j}(0)$, or $D_{i,j}(1) < D_{i,j}(0)$, where the latter two cases occur if $\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))$. It follows that

$$\theta(x_g) = \sum_{j=1}^J \mathbb{E}[Y_i(j) - Y_i(j-1)|D_{i,j}(1) - D_{i,j}(0) \neq 0, X_i = x_g] \mathbb{P}[D_{i,j}(1) - D_{i,j}(0) \neq 0|X_i = x_g]\tag{36}$$

$$= \sum_{j=1}^J \mathbb{E}[Y_i(j) - Y_i(j-1)|\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0)), X_i = x_g]\tag{37}$$

$$\times \mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g].\tag{38}$$

Similarly, we have

$$\begin{aligned}\pi(x_g) &= \mathbb{E}[T_i|Q_i = 1, X_i = x_g] - \mathbb{E}[T_i|Q_i = 0, X_i = x_g] \\ &= \sum_{j=1}^J \mathbb{E}[D_{i,j}(1) - D_{i,j}(0)|X_i = x_g] \\ &= \sum_{j=1}^J \mathbb{P}[D_{i,j}(1) - D_{i,j}(0) \neq 0|X_i = x_g] \\ &= \sum_{j=1}^J \mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g].\end{aligned}\tag{39}$$

Hence, with a multivalued treatment we obtain a covariate-specific weighted average of LATEs, known as an average causal response:

$$\tau(x_g) = \frac{\theta(x_g)}{\pi(x_g)} = \sum_{j=1}^J \mathbb{E}[Y_i(j) - Y_i(j-1)|\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0)), X_i = x_g]\tag{40}$$

$$\times \frac{\mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g]}{\sum_{j=1}^J \mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g]}\tag{41}$$

$$= \sum_{j=1}^J \eta_j \mathbb{E}[Y_i(j) - Y_i(j-1)|\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0)), X_i = x_g],\tag{42}$$

where

$$\eta_j = \frac{\mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g]}{\sum_{j=1}^J \mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g]},\tag{43}$$

with $\sum_j \eta_j = 1$. From [Assumption 1.3](#) follows that $\mathbb{P}[\max(T_i(1), T_i(0)) \geq j > \min(T_i(1), T_i(0))|X_i = x_g] > 0$ for at least one treatment level j , and therefore $\eta_j \geq 0$ for all j . It follows that $\tau(x_g)$ identifies an retains its causal interpretation with a multivalued treatment T . Note that with $J = 1$ the result boils down to the covariate-specific LATE with a binary treatment.

Since all results in this paper are written in terms of $\tau(x_g)$ using the reduced form and first stage (in which T is written in terms of Z and W), they directly apply to the multivalued treatment setting.

B.2 Multivalued instrument

Suppose we have L mutually exclusive binary instruments Q_{il} with $l = 1, \dots, L$. This set of instruments may be interpreted as L different instruments, which includes a full set of interactions across an original set of instruments in case they are not mutually exclusive, or the L levels in a single instrument. The instruments in the set are ordered such that $l < k$ implies that $\mathbb{E}[T_i|Q_{il} = 1, X_i = x_g] < \mathbb{E}[T_i|Q_{ik} = 1, X_i = x_g]$. Define $Q_{i,0} = 1 - \sum_{l=1}^L Q_{il}$. The $LG \times 1$ vector Z_i contains the instrument interactions $Z_{igl} = Q_{il}W_{ig}$, and Z is defined as an $n \times LG$ matrix with Z_i as rows.

Define

$$\tau_{l,l-1}(x_g) = \frac{E[Y_i|Q_{il} = 1, X_i = x_g] - E[Y_i|Q_{i,l-1} = 1, X_i = x_g]}{E[T_i|Q_{il} = 1, X_i = x_g] - E[T_i|Q_{i,l-1} = 1, X_i = x_g]}. \quad (44)$$

It follows from [Lemma 1](#) that $\tau_{l,l-1}(x_g) = \mathbb{E}[Y_i(1) - Y_i(0)|T_i(l) \neq T_i(l-1), X_i = x_g]$ identifies a LATE.

We can write

$$\begin{aligned} & E[Y_i|Q_{il} = 1, X_i = x_g] \\ &= \tau_{l,l-1}(x_g) (E[T_i|Q_{il} = 1, X_i = x_g] - E[T_i|Q_{i,l-1} = 1, X_i = x_g]) + E[Y_i|Q_{i,l-1} = 1, X_i = x_g] \\ &= \sum_{k=1}^l \tau_{k,k-1}(x_g) (E[T_i|Q_{ik} = 1, X_i = x_g] - E[T_i|Q_{i,k-1} = 1, X_i = x_g]) + E[Y_i|Q_{i,0} = 1, X_i = x_g], \end{aligned} \quad (45)$$

and therefore

$$\theta_l(x_g) = E[Y_i|Q_{il} = 1, X_i = x_g] - E[Y_i|Q_{i,0} = 1, X_i = x_g] \quad (46)$$

$$= \sum_{k=1}^l \tau_{k,k-1}(x_g) (E[T_i|Q_{ik} = 1, X_i = x_g] - E[T_i|Q_{i,k-1} = 1, X_i = x_g]). \quad (47)$$

Similarly, we have

$$\pi_l(x_g) = E[T_i|Q_{il} = 1, X_i = x_g] - E[T_i|Q_{i,0} = 1, X_i = x_g] \quad (48)$$

$$= \sum_{k=1}^l (E[T_i|Q_{ik} = 1, X_i = x_g] - E[T_i|Q_{i,k-1} = 1, X_i = x_g]). \quad (49)$$

Hence,

$$\tau_l(x_g) = \frac{\theta_l(x_g)}{\pi_l(x_g)} = \sum_{k=1}^l \tau_{k,k-1}(x_g) \frac{(E[T_i|Q_{ik} = 1, X_i = x_g] - E[T_i|Q_{i,k-1} = 1, X_i = x_g])}{\sum_{k=1}^l (E[T_i|Q_{ik} = 1, X_i = x_g] - E[T_i|Q_{i,k-1} = 1, X_i = x_g])}, \quad (50)$$

which is a weighted average of LATEs with weights that sum up to one and are nonnegative due to the ordering of the instruments. Note that with $l = 1$ the result boils down to the covariate-specific LATE with one instrument.

Instead of averaging over G covariate-specific LATEs $\tau(x_g)$, τ will be a weighted average of LG covariate-specific weighted average of LATEs $\tau_l(x_g)$ with multiple instruments. By substituting the following two expressions in the theoretical derivations in the paper, it follows that the weights are nonnegative and sum to one:

$$\pi' Z' M_W Z \theta = \sum_g \sum_l m_{gl} \left(1 - \frac{m_{gl}}{n_g}\right) \pi_{gl}(x_g) \theta_l(x_g), \quad (51)$$

$$\pi' Z' M_W Z \pi = \sum_g \sum_l m_{gl} \left(1 - \frac{m_{gl}}{n_g}\right) \pi_l(x_g)^2, \quad (52)$$

where π and θ are $LG \times 1$ vectors containing respectively $\pi_{gl}(x_g)$ and $\theta_l(x_g)$, and $m_{gl} = \sum_i Z_{igl}$.

C Proofs - Estimands

C.1 Proof Lemma 1

Using the exclusion restriction in [Assumption 1.2](#), the observed outcomes are linked to the potential outcomes as $Y_i = (1 - T_i)Y_i(0) + T_i Y_i(1) = Y_i(0) + T_i \Delta_i$ with $\Delta_i = Y_i(1) - Y_i(0)$. The observed treatment is linked to the potential treatment statuses as $T_i = (1 - Q_i)T_i(0) + Q_i T_i(1) = T_i(0) + (T_i(1) - T_i(0))Q_i$. Hence, $Y_i = Y_i(0) + T_i(0)\Delta_i + (T_i(1) - T_i(0))\Delta_i Q_i$. We can now write

$$\begin{aligned}
& \frac{E[Y_i|Q_i = 1, X_i = x_g] - E[Y_i|Q_i = 0, X_i = x_g]}{E[T_i|Q_i = 1, X_i = x_g] - E[T_i|Q_i = 0, X_i = x_g]} \\
&= \frac{E[Y_i(0) + T_i(0)\Delta_i + (T_i(1) - T_i(0))\Delta_i|Q_i = 1, X_i = x_g] - E[Y_i(0) + T_i(0)\Delta_i|Q_i = 0, X_i = x_g]}{E[T_i(1)|Q_i = 1, X_i = x_g] - E[T_i(0)|Q_i = 0, X_i = x_g]} \quad (53) \\
&= \frac{E[(T_i(1) - T_i(0))\Delta_i|X_i = x_g]}{E[T_i(1) - T_i(0)|X_i = x_g]} = \frac{E[\Delta_i|T_i(1) \neq T_i(0), X_i = x_g]P[T_i(1) \neq T_i(0) = 1|X_i = x_g]}{P[T_i(1) \neq T_i(0) = 1|X_i = x_g]} \\
&= E[\Delta_i|T_i(1) \neq T_i(0), X_i = x_g],
\end{aligned}$$

using subsequently [Assumption 1.1](#) independence and [1.4](#) monotonicity in the third line, and [1.3](#) relevance in the fourth line.

C.2 Proof Lemma 2

In a saturated specification, we can write

$$\begin{aligned}
T_i &= \sum_g (\mathbb{E}[T_i|Z_{ig} = 1, W_{ig} = 1] - \mathbb{E}[T_i|Z_{ig} = 0, W_{ig} = 1])Z_{ig} + \mathbb{E}[T_i|Z_{ig} = 0, W_{ig} = 1]W_{ig} + u_i, \\
Y_i &= \sum_g (\mathbb{E}[Y_i|Z_{ig} = 1, W_{ig} = 1] - \mathbb{E}[Y_i|Z_{ig} = 0, W_{ig} = 1])Z_{ig} + \mathbb{E}[Y_i|Z_{ig} = 0, W_{ig} = 1]W_{ig} + \varepsilon_i,
\end{aligned} \quad (54)$$

where $u_i = T_i - \mathbb{E}[T_i|Q_i, X_i]$ and $\varepsilon_i = Y_i - \mathbb{E}[Y_i|Q_i, X_i]$. It follows from the derivations in [Appendix C.1](#) that this can be written as

$$T_i = \pi' Z_i + \psi' W_i + u_i, \quad (55)$$

$$Y_i = \theta' Z_i + \phi' W_i + \varepsilon_i, \quad (56)$$

where $\pi = (\pi(x_1), \dots, \pi(x_g))'$ with $\pi(x_g) = \mathbb{P}[T_i(1) \neq T_i(0)|X_i = x_g]$, $\psi = (\psi_1, \dots, \psi_g)'$ with $\psi_g = \mathbb{E}[T_i|Z_{ig} = 0, W_{ig} = 1]$, $\theta = (\theta(x_1), \dots, \theta(x_g))'$ with $\theta(x_g) = \tau(x_g)\mathbb{P}[T_i(1) \neq T_i(0)|X_i = x_g]$, and $\phi = (\phi_1, \dots, \phi_g)'$ with $\phi_g = \mathbb{E}[Y_i|Z_{ig} = 0, W_{ig} = 1]$.

For the TSLS estimand, we now obtain the following.

$$\begin{aligned}
\beta^{\text{TSLS}} &= \frac{\mathbb{E}[T'PY|Q, X]}{\mathbb{E}[T'PT|Q, X]} = \frac{\mathbb{E}[(Z\pi + u)'P(Z\theta + \varepsilon)|Q, X]}{\mathbb{E}[(Z\pi + u)'P(Z\pi + u)|Q, X]} \\
&= \frac{\pi'Z'PZ\theta + \pi'Z'P\mathbb{E}[\varepsilon|Q, X] + \mathbb{E}[u|Q, X]'PZ\theta + \mathbb{E}[u'P\varepsilon|Q, X]}{\pi'Z'PZ\pi + 2\pi'Z'P\mathbb{E}[u|Q, X] + \mathbb{E}[u|Q, X]'PZ\theta + \mathbb{E}[u'Pu|Q, X]} \\
&= \frac{\pi'Z'M_W Z\theta + \mathbb{E}[u'P\varepsilon|Q, X]}{\pi'Z'M_W Z\pi + \mathbb{E}[u'Pu|Q, X]} \quad (57) \\
&= \frac{\sum_g \frac{\theta(x_g)}{\pi(x_g)} \pi(x_g)^2 n_g \frac{m_g}{n_g} (1 - \frac{m_g}{n_g}) + \sum_i \mathbb{E}[u_i \varepsilon_i | Q_i, X_i] P_{ii}}{\sum_g \pi(x_g)^2 n_g \frac{m_g}{n_g} (1 - \frac{m_g}{n_g}) + \sum_i \mathbb{E}[u_i^2 | Q_i, X_i] P_{ii}} \\
&= \frac{\sum_g \frac{n_g}{n} \pi(x_g)^2 \frac{m_g}{n_g} (1 - \frac{m_g}{n_g}) \tau(x_g) + \frac{1}{n} \sum_i \mathbb{E}[u_i \varepsilon_i | Q_i, X_i] P_{ii}}{\sum_g \frac{n_g}{n} \pi(x_g)^2 \frac{m_g}{n_g} (1 - \frac{m_g}{n_g}) + \frac{1}{n} \sum_i \mathbb{E}[u_i^2 | Q_i, X_i] P_{ii}},
\end{aligned}$$

where we use that $Z' M_W Z$ is a $G \times G$ diagonal matrix with elements $n_g \frac{m_g}{n_g} (1 - \frac{m_g}{n_g})$.

Note that $e_i' M_W Z$ is a $1 \times G$ vector with as elements the residual of observation i in the regression of Z_{ig} on W_{ig} , so $[e_i' M_W Z]_g = Z_{ig} - W_{ig} \tilde{P}[Z_i = 1 | X_i = x_g]$, where $\tilde{P}[Z_i = 1 | X_i = x_g] = \frac{m_g}{n_g}$. Hence, $P_{ii} = e_i' M_W Z (Z' M_W Z)^{-1} Z' M_W e_i = \sum_g \frac{1}{n_g} W_{ig} \frac{(Z_{ig} - \tilde{P}[Z_i = 1 | X_i = x_g])^2}{\text{Var}[Q_i | X_i = x_g]}$.

C.3 Proof Lemma 3

For the JIVE1 estimator we have

$$\beta^{\text{JIVE1}} = \frac{\mathbb{E}[T'(P - D_P)Y | Q, X]}{\mathbb{E}[T'(P - D_P)T | Q, X]} = \frac{\mathbb{E}[T'PY | Q, X] - \mathbb{E}[T'D_P Y | Q, X]}{\mathbb{E}[T'PT | Q, X] - \mathbb{E}[T'D_P T | Q, X]}, \quad (58)$$

with

$$\begin{aligned} \mathbb{E}[T'D_P Y | Q, X] &= \mathbb{E}[(Z\pi + W\psi + u)' D_P (Z\theta + W\phi + \varepsilon) | Q, X] \\ &= \pi' Z' D_P Z \theta + \pi' Z' D_P W \phi + \pi' Z' D_P \mathbb{E}[\varepsilon | Q, X] \\ &\quad + \psi' W' D_P Z \theta + \psi' W' D_P W \phi + \psi' W' D_P \mathbb{E}[\varepsilon | Q, X] \\ &\quad + \mathbb{E}[u | Q, X]' D_P W \phi + \mathbb{E}[u | Q, X]' D_P Z \theta + \mathbb{E}[u' D_P \varepsilon | Q, X] \\ &= \pi' Z' D_P Z \theta + \pi' Z' D_P W \phi + \psi' W' D_P Z \theta + \psi' W' D_P W \phi + \mathbb{E}[u' D_P \varepsilon | Q, X] \\ &= \sum_g \left(\frac{\theta(x_g)}{\pi(x_g)} \pi(x_g)^2 + \pi(x_g) \phi_g + \theta(x_g) \psi_g \right) \left(1 - \frac{m_g}{n_g} \right) + \psi_g \phi_g + \sum_i \mathbb{E}[u_i \varepsilon_i | Q_i, X_i] P_{ii}, \end{aligned} \quad (59)$$

where $W' D_P W = I_G$ as $\sum_i W_{ig}^2 P_{ii} = 1$, and $Z' D_P Z = W' D_P Z$ is a $G \times G$ diagonal matrix with elements $\sum_i Z_{ig}^2 P_{ii} = (1 - \frac{m_g}{n_g})$. Similarly, we get

$$\mathbb{E}[T'D_P T | Q, X] = \sum_g \left(\pi(x_g)^2 + 2\pi(x_g) \psi_g \right) \left(1 - \frac{m_g}{n_g} \right) + \psi_g^2 + \sum_i \mathbb{E}[u_i \varepsilon_i | Q_i, X_i] P_{ii}. \quad (60)$$

Hence, using the results in [Appendix C.2](#), we have

$$\begin{aligned} \beta^{\text{JIVE1}} &= \frac{\sum_g \frac{n_g}{n} \pi(x_g)^2 \frac{m_g}{n_g} \left(1 - \frac{m_g}{n_g} \right) \tau(x_g) \left(1 - \frac{1}{m_g} \right) - \sum_g \frac{n_g}{n} \left(\pi(x_g) \phi_g + \theta(x_g) \psi_g \right) \frac{m_g}{n_g} \left(1 - \frac{m_g}{n_g} \right) \frac{1}{m_g} - \frac{1}{n} \psi_g \phi_g}{\sum_g \frac{n_g}{n} \pi(x_g)^2 \frac{m_g}{n_g} \left(1 - \frac{m_g}{n_g} \right) \left(1 - \frac{1}{m_g} \right) - 2 \sum_g \frac{n_g}{n} \pi(x_g) \psi_g \frac{m_g}{n_g} \left(1 - \frac{m_g}{n_g} \right) \frac{1}{m_g} - \frac{1}{n} \psi_g^2} \\ &= \frac{\sum_g \mathbb{P}[X_i = x_g] \pi(x_g)^2 \text{Var}[Q_i | X_i = x_g] \left(1 - \mathbb{P}[Q_i = 1, X_i = x_g] \right) \tau(x_g) - B_Y}{\sum_g \mathbb{P}[X_i = x_g] \pi(x_g)^2 \text{Var}[Q_i | X_i = x_g] \left(1 - \mathbb{P}[Q_i = 1, X_i = x_g] \right) - B_T} \end{aligned} \quad (61)$$

where

$$\begin{aligned} B_Y &= \sum_g \mathbb{P}[X_i = x_g] \left(\pi(x_g) \phi_g + \theta(x_g) \psi_g \right) \text{Var}[Q_i | X_i = x_g] \mathbb{P}[Q_i = 1, X_i = x_g] - \frac{1}{n} \psi_g \phi_g, \\ B_T &= 2 \sum_g \mathbb{P}[X_i = x_g] \pi(x_g) \psi_g \text{Var}[Q_i | X_i = x_g] \mathbb{P}[Q_i = 1, X_i = x_g] - \frac{1}{n} \psi_g^2, \end{aligned} \quad (62)$$

with $\tilde{P}[Q_i = 1, X_i = x_g] = \frac{1}{m_g}$, $\psi_g = \mathbb{E}[T_i | Z_{ig} = 0, W_{ig} = 1]$, and $\phi_g = \mathbb{E}[Y_i | Z_{ig} = 0, W_{ig} = 1]$.

C.4 Proof Lemma 5

First we show that $P_{ii} = [M_{W,Z}DM_{W,Z}]_{ii}$ if the elements of the diagonal matrix D are set equal to $D_{ii} = \sum_{k=1}^n [(M_{W,Z} \odot M_{W,Z})^{-1}]_{ik} P_{kk}$. Using that D is diagonal and $M_{W,Z}$ symmetric, we have

$$\begin{aligned} [M_{W,Z}DM_{W,Z}]_{ii} &= \sum_j [M_{W,Z} \odot M_{W,Z}]_{ij} D_{jj} \\ &= \sum_j [M_{W,Z} \odot M_{W,Z}]_{ij} \sum_{k=1}^n P_{kk} [(M_{W,Z} \odot M_{W,Z})^{-1}]_{kj} \\ &= \sum_{k=1}^n P_{kk} \sum_j [M_{W,Z} \odot M_{W,Z}]_{ij} [(M_{W,Z} \odot M_{W,Z})^{-1}]_{kj} \\ &= P_{ii}. \end{aligned} \tag{63}$$

Next, we derive an expression for D_{ii} . Note that $M_{W,Z}$ is a block diagonal matrix with the g th block an $n_g \times n_g$ matrix

$$M_{W,Z,(g)} = I_{n_g} - P_{W,Z,(g)} = I_{n_g} - \begin{pmatrix} m_g^{-1} t_{m_g} t'_{m_g} & O_{m_g, n_g - m_g} \\ O_{n_g - m_g, m_g} & (n_g - m_g)^{-1} t_{n_g - m_g} t'_{n_g - m_g} \end{pmatrix}, \tag{64}$$

where the observations are ordered according to covariate group and within covariate group on active instrument, without loss of generality. It now follows that

$$[M_{W,Z} \odot M_{W,Z}]_{(g)} = \begin{pmatrix} \left(1 - \frac{2}{m_g}\right) I_{m_g} + \frac{1}{m_g^2} t_{m_g} t'_{m_g} & O_{m_g, n_g - m_g} \\ O_{n_g - m_g, m_g} & \left(1 - \frac{2}{n_g - m_g}\right) I_{n_g - m_g} + \frac{1}{(n_g - m_g)^2} t_{n_g - m_g} t'_{n_g - m_g} \end{pmatrix}, \tag{65}$$

and hence

$$[(M_{W,Z} \odot M_{W,Z})^{-1}]_{(g)} = \begin{pmatrix} \frac{m_g}{m_g - 2} \left(I_{m_g} - \frac{1}{m_g(m_g - 1)} t_{m_g} t'_{m_g} \right) & O_{m_g, n_g - m_g} \\ O_{n_g - m_g, m_g} & \frac{n_g - m_g}{n_g - m_g - 2} \left(I_{n_g - m_g} - \frac{1}{(n_g - m_g)(n_g - m_g - 1)} t_{n_g - m_g} t'_{n_g - m_g} \right) \end{pmatrix}. \tag{66}$$

Note that $P_{ii} = \sum_g \frac{1}{n_g} W_{ig} \frac{(Z_{ig} - \frac{m_g}{n_g})^2}{\frac{m_g}{n_g} (1 - \frac{m_g}{n_g})}$, as derived in [Appendix C.2](#), and hence $P_{ii} = \frac{1}{m_g} - \frac{1}{n_g}$ if $W_{ig} = Z_{ig} = 1$ and $P_{ii} = \frac{m_g/n_g}{n_g - m_g}$ if $W_{ig} = 1$ and $Z_{ig} = 0$. Define $[D_P]_{(g)}$ as the elements in the diagonal of P corresponding to the observations in $[M_{W,Z} \odot M_{W,Z}]_{(g)}$:

$$\begin{aligned} D_{ii} &= \sum_g W_{ig} e'_i [(M_{W,Z} \odot M_{W,Z})^{-1}]_{(g)} [D_P]_{(g)} \\ &= \sum_g W_{ig} \left[\frac{m_g}{m_g - 2} \left(1 - \frac{m_g}{m_g(m_g - 1)} \right) \left(\frac{1}{m_g} - \frac{1}{n_g} \right) Z_{ig} \right. \\ &\quad \left. + \frac{n_g - m_g}{n_g - m_g - 2} \left(1 - \frac{n_g - m_g}{(n_g - m_g)(n_g - m_g - 1)} \right) \frac{m_g}{n_g(n_g - m_g)} (1 - Z_{ig}) \right] \\ &= \sum_g \frac{1}{n_g} W_{ig} \left[\frac{n_g - m_g}{m_g - 1} Z_{ig} + \frac{m_g}{n_g - m_g - 1} (1 - Z_{ig}) \right]. \end{aligned} \tag{67}$$

The derivation above implicitly assumes that $m_g \geq 3$ and $n_g - m_g \geq 3$ for all expressions to be well-defined. However, the end result only requires $m_g \geq 2$ and $n_g - m_g \geq 2$. We can briefly verify whether D_{ii} as given on the final line of (67) indeed yields a zero diagonal for $P - M_{W,Z}DM_{W,Z}$.

Consider observation i for which $Z_{ig} = 1$. Then,

$$P_{ii} = \frac{1}{n_g} \frac{n_g - m_g}{m_g}, \quad D_{ii} = \frac{1}{n_g} \frac{n_g - m_g}{m_g - 1}, \quad P_{W,Z,ii} = \frac{1}{m_g}. \quad (68)$$

Since D_{ii} is the same for all observations that have $Z_{ig} = 1$, we have

$$[M_{W,Z} D M_{W,Z}]_{ii} = D_{ii}(1 - 2P_{W,Z,ii}) + m_g \frac{1}{m_g^2} D_{ii} = D_{ii}(1 - P_{W,Z,ii}) = D_{ii} \frac{m_g - 1}{m_g} = P_{ii}. \quad (69)$$

The case where $Z_{ig} = 0$ follows analogously.

C.5 Proof Theorem 1

$$\beta^{\text{SIVE}} = \frac{\mathbb{E}[T'(P - M_{W,Z} D M_{W,Z})Y|Q, X]}{\mathbb{E}[T'(P - M_{W,Z} D M_{W,Z})T|Q, X]} = \frac{\mathbb{E}[T'PY|Q, X] - \mathbb{E}[T'M_{W,Z} D M_{W,Z} Y|Q, X]}{\mathbb{E}[T'PT|Q, X] - \mathbb{E}[T'M_{W,Z} D M_{W,Z} T|Q, X]}, \quad (70)$$

with

$$\mathbb{E}[T'M_{W,Z} D M_{W,Z} Y|Q, X] = \mathbb{E}[u' M_{W,Z} D M_{W,Z} \varepsilon|Q, X] = \sum_i \mathbb{E}[u_i \varepsilon_i | Q_i, X_i] P_{ii}, \quad (71)$$

and similarly we have $\mathbb{E}[T'M_{W,Z} D M_{W,Z} Y|Q, X] = \sum_i \mathbb{E}[u_i^2 | Q_i, X_i] P_{ii}$. Combining this result with the result in [Appendix C.2](#), we have

$$\beta^{\text{SIVE}} = \frac{\sum_g \frac{n_g}{n} \pi(x_g)^2 \frac{m_g}{n_g} (1 - \frac{m_g}{n_g}) \tau(x_g)}{\sum_g \frac{n_g}{n} \pi(x_g)^2 \frac{m_g}{n_g} (1 - \frac{m_g}{n_g})}. \quad (72)$$

D Proofs - Inference

D.1 Preliminary results

D.1.1 The elements of the projection matrix P

Since the columns in W are orthogonal, $P_W = W(W'W)^{-1}W$ is a block diagonal matrix consisting of G blocks of dimension $n_g \times n_g$. The g th block is given by $P_{(g)} = \frac{1}{n_g} l_{n_g} l_{n_g}'$. We then have that

$$M_W Z = \begin{pmatrix} M_{l_{n_1} z(1)} & 0_{n_1} & \dots \\ 0_{n_2} & M_{l_{n_2} z(2)} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (73)$$

The matrix P is again block diagonal, with

$$P = \begin{pmatrix} (z'_{(1)} M_{l_{n_1} z(1)})^{-1} M_{l_{n_1} z(1)} z'_{(1)} M_{l_{n_1}} & 0_{n_1, n_2} & \dots \\ 0_{n_2, n_1} & (z'_{(2)} M_{l_{n_2} z(2)})^{-1} M_{l_{n_2} z(2)} z'_{(2)} M_{l_{n_2}} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (74)$$

The elements of P are given by

$$[P]_{ij} = \begin{cases} \frac{1}{n_g} \frac{1 - m_g/n_g}{m_g/n_g} & W_{ig} = W_{jg} = 1, Z_{ig} = 1, Z_{jg} = 1, \\ -1/n_g & W_{ig} = W_{jg} = 1, Z_{ig} = 1, Z_{jg} = 0, \\ -1/n_g & W_{ig} = W_{jg} = 1, Z_{ig} = 0, Z_{jg} = 1, \\ \frac{1}{n_g} \frac{m_g/n_g}{1 - m_g/n_g} & W_{ig} = W_{jg} = 1, Z_{ig} = 0, Z_{jg} = 0, \\ 0 & \text{otherwise.} \end{cases} \quad (75)$$

D.1.2 Results on the matrix A used for asymptotic normality

Recall that $A = P - M_{W,Z} D M_{W,Z}$. This matrix is block diagonal with blocks

$$A_{(g)} = \begin{pmatrix} \frac{n_g - m_g}{n_g(m_g - 1)}(l_{m_g} l'_{m_g} - I_{m_g}) & -\frac{1}{n_g} l_{m_g} l'_{n_g - m_g} \\ -\frac{1}{n_g} l_{n_g - m_g} l'_{m_g} & \frac{m_g}{n_g(n_g - m_g - 1)}(l_{n_g - m_g} l'_{n_g - m_g} - I_{n_g - m_g}) \end{pmatrix}. \quad (76)$$

The following result is used in the proof of the central limit theorem for the SIVE estimator,

$$\begin{aligned} \text{tr}(A^2) &= G + \text{tr}(M_{W,Z} D M_{W,Z} D) \\ &= G + \sum_{i=1}^n D_{ii}^2 - 2 \sum_{i=1}^n D_{ii}^2 [P_{W,Z}]_{ii} + \text{tr}(P_{W,Z} D P_{W,Z} D) \\ &= G + \sum_{g=1}^G m_g \frac{1}{n_g^2} \frac{(n_g - m_g)^2}{(m_g - 1)^2} + (n_g - m_g) \frac{1}{n_g^2} \frac{m_g^2}{(n_g - m_g - 1)^2} \\ &\quad - \sum_{g=1}^G \frac{1}{n_g^2} \frac{(n_g - m_g)^2}{(m_g - 1)^2} + \frac{1}{n_g^2} \frac{m_g^2}{(n_g - m_g - 1)^2} \\ &= G + \sum_{g=1}^G \frac{1}{n_g^2} \left[\frac{(n_g - m_g)^2}{m_g - 1} + \frac{m_g^2}{n_g - m_g - 1} \right] \\ &\leq G + \sum_{g=1}^G \left[\frac{1}{m_g - 1} + \frac{1}{n_g - m_g - 1} \right] \\ &\leq 3G, \end{aligned} \quad (77)$$

where the last line uses that $m_g \geq 2$ and $n_g - m_g \geq 2$. From the first line we also immediately have that

$$\text{tr}(A^2) \geq G. \quad (78)$$

The eigenvalues of A are 0 with multiplicity 1, 1 with multiplicity 1, $-(n_g - m_g)/(n_g(m_g - 1))$ with multiplicity $m_g - 1$ and $-m_g/(n_g(n_g - m_g - 1))$ with multiplicity $n_g - m_g - 1$. Since $m_g \geq 2$ and $n_g - m_g \geq 2$, we conclude that the eigenvalues are on the $[-1, 1]$ interval.

In the SIVE estimator, we encounter the product $AZ\pi$ and $AZ\zeta$. We require an elementwise bound on these products, which can be established as follows.

$$\begin{aligned} |[AZ\pi]_i| &= |[M_{W,Z}\pi]_i| \\ &= \pi(x_g) \begin{cases} \frac{n_g - m_g}{n_g} & \text{if } Z_{ig} = 1, W_{ig} = 1, \\ \frac{m_g}{n_g} & \text{if } Z_{ig} = 0, W_{ig} = 1, \\ 0 & \text{otherwise.} \end{cases} \\ &\leq C < \infty, \quad a.s. \end{aligned} \quad (79)$$

The same result follows for $AZ\zeta = AZ(\theta - \beta^{\text{SIVE}}\pi)$ by noting that $\theta(x_g)/\pi(x_g) = \tau(x_g) \leq C < \infty$ *a.s.* by [Assumption 3.1](#).

D.1.3 Results on the matrix A used for consistency of the variance estimator

To show consistency of the variance estimator, we need to upper bound terms that are of the form

$$S(A, \bar{P}) = \sum_{i_1, \dots, i_{16}} |A_{i_1 i_2}| |A_{i_3 i_4}| |A_{i_5 i_6}| |A_{i_7 i_8}| |\bar{P}_{i_9 i_{10}}| |\bar{P}_{i_{11} i_{12}}| |\bar{P}_{i_{13} i_{14}}| |\bar{P}_{i_{15} i_{16}}|, \quad (80)$$

with some restrictions on the indices over which we sum. We will first find a matrix \tilde{A} such that $S(A, \bar{P}) \leq S(\tilde{A}, \bar{P})$. To save on notation, define

$$H_g = \begin{pmatrix} \iota_{m_g} & 0_{m_g} \\ 0_{n_g-m_g} & \iota_{n_g-m_g} \end{pmatrix}, \quad E_g = \begin{pmatrix} m_g^{-1} & 0 \\ 0 & (n_g - m_g)^{-1} \end{pmatrix} \quad (81)$$

We then have

$$\begin{aligned} A_{(g)} &= \begin{pmatrix} \frac{n_g-m_g}{n_g(m_g-1)}(\iota_{m_g} \iota'_{m_g} - I_{m_g}) & -\frac{1}{n_g} \iota_{m_g} \iota'_{n_g-m_g} \\ -\frac{1}{n_g} \iota_{n_g-m_g} \iota'_{m_g} & \frac{m_g}{n_g(n_g-m_g-1)}(\iota_{n_g-m_g} \iota'_{n_g-m_g} - I_{n_g-m_g}) \end{pmatrix} \\ &= \frac{m_g(n_g - m_g)}{n_g} \begin{pmatrix} \frac{1}{m_g(m_g-1)}(\iota_{m_g} \iota'_{m_g} - I_{m_g}) & -\frac{1}{m_g(n_g-m_g)} \iota_{m_g} \iota'_{n_g-m_g} \\ -\frac{1}{m_g(n_g-m_g)} \iota_{n_g-m_g} \iota'_{m_g} & \frac{1}{(n_g-m_g)(n_g-m_g-1)}(\iota_{n_g-m_g} \iota'_{n_g-m_g} - I_{n_g-m_g}) \end{pmatrix} \\ &\leq 2 \frac{m_g(n_g - m_g)}{n_g} H_g E_g \iota_2 \iota_2' E_g H_g' \equiv \tilde{A}_{(g)}, \end{aligned} \quad (82)$$

where \leq indicates an elementwise inequality.

We can now establish the following results

$$\begin{aligned} (R1) \quad & \text{tr}(\tilde{A}_{(g)}) \leq C, \\ (R2) \quad & \lambda_{\max}(\tilde{A}_{(g)}) \leq C, \\ (R3) \quad & \tilde{A}_{(g)}^2 = 2\tilde{A}_{(g)}, \\ (R4) \quad & \tilde{A}_{(g)} P_{W,Z,(g)} = \tilde{A}_{(g)}, \\ (R5) \quad & \iota'_{n_g} (\tilde{A}_{(g)} \odot P_{W,Z,(g)}) \iota_{n_g} \leq C, \\ (R6) \quad & \iota'_{n_g} (\tilde{A}_{(g)} \odot P_{W,Z,(g)})^2 \iota_{n_g} \leq C, \\ (R7) \quad & \iota'_{n_g} (\tilde{A}_{(g)} \odot P_{W,Z,(g)}) \tilde{A}_{(g)} \iota_{n_g} \leq C, \\ (R8) \quad & \iota'_{n_g} \tilde{A}_{(g)} (P_{W,Z,(g)} \odot P_{W,Z,(g)}) \tilde{A}_{(g)} \iota_{n_g} \leq C. \end{aligned} \quad (83)$$

Proof: (R1) follows from the fact that $H_g' H_g = E_g^{-1}$, $\iota'_{n_g} H_g = \iota_2' E_g^{-1}$ and $\iota_2' E_g \iota_2 = \frac{n_g}{m_g(n_g-m_g)}$. Since $\text{rank}(\tilde{A}_{(g)}) = 1$, (R2) follows from (R1). (R3) follows from the fact that $H_g' H_g = E_g^{-1}$. For (R4)–(R8), we first note that the g th diagonal block of the projection matrix $P_{W,Z}$ satisfies $P_{W,Z,(g)} = H_g E_g H_g'$. Using this result and again the fact that $H_g' H_g = E_g^{-1}$ yields (R4). For (R5), we note that $\iota'_{n_g} (\tilde{A}_{(g)} \odot P_{W,Z,(g)}) \iota_{n_g} = \text{tr}(\tilde{A}_{(g)} P_{W,Z,(g)}) = \text{tr}(\tilde{A}_{(g)})$ with the last equality by (R4). For (R6) and (R7) we first calculate the elementwise product,

$$\tilde{A}_{(g)} \odot P_{W,Z,(g)} = 2 \frac{m_g(n_g - m_g)}{n_g} H_g E_g^3 H_g'. \quad (84)$$

We can now explicitly calculate bounds for (R6) and (R7),

$$\begin{aligned} \iota'_{n_g} (\tilde{A}_{(g)} \odot P_{W,Z,(g)})^2 \iota_{n_g} &= 4 \frac{m_g^2(n_g - m_g)^2}{n_g^2} (m_g^{-3} + (n_g - m_g)^{-3}) \\ &= 4 \frac{n_g^3 - 3n_g m_g^2 - 3n_g^2 m_g}{n_g^3(n_g - m_g)} \leq 4 \frac{1}{n_g} \leq 2, \\ \iota'_{n_g} (\tilde{A}_{(g)} \odot P_{W,Z,(g)}) \tilde{A}_{(g)} \iota_{n_g} &= 4 \frac{m_g^2(n_g - m_g)^2}{n_g^2} (m_g^{-2} + (n_g - m_g)^{-2}) \leq 8, \end{aligned} \quad (85)$$

For (R8), we have a similar result as,

$$\iota'_{n_g} \tilde{A}_{(g)} (P_{W,Z,(g)} \odot P_{W,Z,(g)}) \tilde{A}_{(g)} \iota_{n_g} = 4 \frac{m_g^2(n_g - m_g)^2}{n_g^2} (m_g^{-2} + (n_g - m_g)^{-2}) \leq 8. \quad (86)$$

Finally, we establish the following results.

$$(R9) \quad \iota'_{n_g}(\tilde{A}_{(g)} \odot \tilde{A}_{(g)})^2 \iota_{n_g} = 16 \frac{m_g^2(n_g - m_g)^2}{n_g^2} \iota'_2 E_g^3 \iota_2 \leq 16,$$

$$(R10) \quad \iota'_{n_g} \tilde{A}_{(g)} (\tilde{A}_{(g)} \odot \tilde{A}_{(g)}) \tilde{A}_{(g)} \iota_{n_g} \leq 256, \quad (87)$$

$$(R11) \quad (\iota'_{n_g} (\tilde{A}_{(g)} \odot \tilde{A}_{(g)})) ((\tilde{A}_{(g)} \iota_{n_g}) \odot (\tilde{A}_{(g)} \iota_{n_g})) \leq 128.$$

From the definition of $\tilde{A}_{(g)}$, we have

$$\tilde{A}_{(g)} \odot \tilde{A}_{(g)} = 4 \frac{m_g^2(n_g - m_g)^2}{n_g^2} H_g E_g^2 \iota_2 \iota'_2 E_g^2 H'_g, \quad (88)$$

$$\tilde{A}_{(g)} \iota_{n_g} = 4 \frac{m_g(n_g - m_g)}{n_g} H_g E_g \iota_2, \quad (89)$$

$$(\tilde{A}_{(g)} \iota_{n_g}) \odot (\tilde{A}_{(g)} \iota_{n_g}) = 16 \frac{m_g^2(n_g - m_g)^2}{n_g^2} H_g E_g^2 \iota_2, \quad (90)$$

$$(\tilde{A}_{(g)} \odot \tilde{A}_{(g)}) \iota_{n_g} = 4 \frac{m_g(n_g - m_g)}{n_g} H_g E_g^2 \iota_2. \quad (91)$$

The first equality in (R9) then follows from (91) and using that $H'_g H_g = E_g^{-1}$. Then (R9) follows from the following.

$$\begin{aligned} \frac{m_g^2(n_g - m_g)^2}{n_g^2} \iota'_2 E_g^3 \iota_2 &= \frac{m_g^2(n_g - m_g)^2}{n_g^2} \frac{m_g^3 + (n_g - m_g)^3}{m_g^3(n_g - m_g)^3} \\ &= \frac{n_g^3 - 3n_g^2 m_g - 3n_g m_g^2}{n_g^2 m_g (n_g - m_g)} \\ &\leq \frac{n_g^3 - 3n_g^2 m_g - 3n_g m_g^2}{n_g^3} \leq 1, \end{aligned} \quad (92)$$

where the first inequality uses that $m_g(n_g - m_g) \geq 2(n_g - 2) \geq 2(n_g - n_g/2) = n_g$ since $n_g \geq 4$ by [Assumption 2](#).

For (R10), using first (88) and then (89),

$$\iota'_{n_g} \tilde{A}_{(g)} (\tilde{A}_{(g)} \odot \tilde{A}_{(g)}) \tilde{A}_{(g)} \iota_{n_g} = 64 \frac{m_g^4(n_g - m_g)^4}{n_g^4} (\iota'_2 E_g^2 \iota_2)^2 = 64 \frac{(m_g^2 + (n_g - m_g)^2)^2}{n_g^4} \leq 256. \quad (93)$$

Finally, for (R11) using (90) and (91), we have

$$(\iota'_{n_g} (\tilde{A}_{(g)} \odot \tilde{A}_{(g)})) ((\tilde{A}_{(g)} \iota_{n_g}) \odot (\tilde{A}_{(g)} \iota_{n_g})) = 64 \frac{m_g^3(n_g - m_g)^3}{n_g^3} \iota'_2 E_g^3 \iota_2 = 64 \frac{m_g^3 + (n_g - m_g)^3}{n_g^3} \leq 128. \quad (94)$$

D.1.4 First stage and reduced form error (co)variance estimators

From the reduced form equation (56) and first stage (55), we see that $M_{W,Z}T = M_{W,Z}u$ and $M_{W,Z}(Y - T\beta^{\text{SIVE}}) = M_{W,Z}v$. We now have the following estimators,

$$\begin{aligned}
\hat{\sigma}_{u,i}^2 &= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}T \odot M_{W,Z}T) \\
&= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}u \odot M_{W,Z}u) \\
\hat{\sigma}_{uv,i} &= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}}) \odot M_{W,Z}T) \\
&= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}(Y - T\beta^{\text{SIVE}}) \odot M_{W,Z}T) - (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})\hat{\sigma}_{u,i}^2 \\
&= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}v \odot M_{W,Z}u) - (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})\hat{\sigma}_{u,i}^2 \\
&= \hat{\sigma}_{uv,i}(\beta^{\text{SIVE}}) - (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})\hat{\sigma}_{u,i}^2, \\
\hat{\sigma}_{v,i}^2 &= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}}) \odot M_{W,Z}(Y - T\hat{\beta}^{\text{SIVE}})) \\
&= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}(Y - T\beta^{\text{SIVE}}) \odot M_{W,Z}(Y - T\beta^{\text{SIVE}})) \\
&\quad - 2(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}}) + (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})^2\hat{\sigma}_{u,i}^2 \\
&= e_i'(M_{W,Z} \odot M_{W,Z})^{-1}(M_{W,Z}v \odot M_{W,Z}v) \\
&\quad - 2(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}}) + (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})^2\hat{\sigma}_{u,i}^2 \\
&= \hat{\sigma}_{v,i}^2(\beta^{\text{SIVE}}) - 2(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}}) + (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})^2\hat{\sigma}_{u,i}^2,
\end{aligned} \tag{95}$$

where we write $\hat{\sigma}_{v,i}^2(\beta^{\text{SIVE}})$ and $\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}})$ as the infeasible analogues of $\hat{\sigma}_{v,i}^2$ and $\hat{\sigma}_{uv,i}$ with $\hat{\beta}^{\text{SIVE}}$ replaced by β^{SIVE} . Since $\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}} = o_p(1)$, we will see that the contribution of the corresponding terms in $\hat{\sigma}_{uv,i}$ and $\hat{\sigma}_{v,i}^2$ to the variance estimator for the score is $o_p(1)$ as well.

We first show that $\hat{\sigma}_{u,i}^2$, $\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}})$ and $\hat{\sigma}_{v,i}^2(\beta^{\text{SIVE}})$ are unbiased for $\sigma_{u,i}^2$, $\sigma_{uv,i}$ and $\sigma_{v,i}^2$, respectively. To save on space, we denote the elements of $P_{W,Z}$ as \bar{P}_{ij} . From the properties of $M_{W,Z}$ and using that $\bar{P}_{ij} = m_g^{-1}$ if $Z_{ig} = Z_{jg} = 1$, $\bar{P}_{ij} = (n_g - m_g)^{-1}$ if $Z_{ig} = Z_{jg} = 0$, and $W_{ig} = W_{jg} = 1$ and $P_{ij} = 0$ elsewhere, we can write

$$\begin{aligned}
\hat{\sigma}_{u,i}^2 &= u_i^2 - u_i \frac{2}{1-2\bar{P}_{ii}} \sum_{j \neq i} \bar{P}_{ij} u_j + \frac{1}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})} \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} u_j u_k, \\
\hat{\sigma}_{v,i}^2(\beta^{\text{SIVE}}) &= v_i^2 - v_i \frac{2}{1-2\bar{P}_{ii}} \sum_{j \neq i} \bar{P}_{ij} v_j + \frac{1}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})} \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} v_j v_k, \\
\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}}) &= u_i v_i - v_i \frac{1}{1-2\bar{P}_{ii}} \sum_{j \neq i} \bar{P}_{ij} u_j - u_i \frac{1}{1-2\bar{P}_{ii}} \sum_{j \neq i} \bar{P}_{ij} v_j + \frac{1}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})} \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} u_j v_k.
\end{aligned} \tag{96}$$

Taking the conditional expectation and using independence across i , we see that $\hat{\sigma}_{u,i}^2$, $\hat{\sigma}_{uv,i}(\beta^{\text{SIVE}})$ and $\hat{\sigma}_{v,i}^2(\beta^{\text{SIVE}})$ are (conditionally) unbiased for $\sigma_{u,i}^2$, $\sigma_{uv,i}$ and $\sigma_{v,i}^2$, respectively. What is important to note is that the estimators only exist if $m_g > 2$ and $n_g - m_g > 2$ for all $g = 1, \dots, G$. We discuss how we handle groups with $m_g = 2$ or $n_g - m_g = 2$ in [Appendix D.3.3](#).

D.2 Proof Lemma 6

Let $\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}} = S_n/D_n$ with

$$\begin{aligned}
S_n &= r_n^{-1/2} T' A (Y - T\beta^{\text{SIVE}}), \\
D_n &= r_n^{-1/2} T' A T.
\end{aligned} \tag{97}$$

We refer to S_n as the score vector. Let $\zeta = \theta - \beta^{\text{SIVE}}\pi$ and $v_i = \varepsilon_i - u_i\beta^{\text{SIVE}}$ with ε_i and u_i defined in the reduced form equations (55) and (56). Denote $\sigma_{u,i}^2 = \mathbb{E}[u_i^2|Q, X]$, $\sigma_{v,i}^2 = \mathbb{E}[v_i^2|Q, X]$, and $\sigma_{uv,i} = \mathbb{E}[u_i v_i|Q, X]$.

We have $\mathbb{E}[S_n|Q, X] = 0$ and

$$\begin{aligned} \mathbb{V}(S_n|Q, X) &= r_n^{-1} \sum_{i=1}^n \sigma_{u,i}^2 [M_W Z \zeta]_i^2 + \sigma_{v,i}^2 [M_W Z \pi]_i^2 + 2\sigma_{uv,i} [M_W Z \pi]_i [M_W Z \zeta]_i \\ &\quad + r_n^{-1} \sum_{i=1}^n \sum_{j \neq i} A_{ij}^2 (\sigma_{uv,i} \sigma_{uv,j} + \sigma_{v,i}^2 \sigma_{u,i}^2) \\ &\leq C [r_n^{-1} \zeta' Z' M_W Z \zeta + r_n^{-1} \pi' Z' M_W Z \pi + r_n^{-1} \text{tr}(A^2)] \leq C, \quad a.s., \end{aligned} \quad (98)$$

where the first inequality follows from [Assumption 3.3](#) and the second inequality from [Assumption 3.1](#), the definition of r_n and (77). For D_n , we have $\mathbb{E}[D_n|Q, X] = r_n^{-1/2} \pi' Z' M_W Z \pi \equiv r_n^{1/2} H_n$ and

$$\begin{aligned} \mathbb{V}(D_n|Q, X) &\leq 4r_n^{-1} \mathbb{E}[(u' M_W Z \pi)^2|Q, X] + r_n^{-1} \mathbb{E}[(u' A u)^2|Q, X] \\ &\leq 4 \max_i \sigma_{u,i}^2 r_n^{-1} \pi' Z' M_W Z \pi + 2r_n^{-1} \max_i \sigma_{u,i}^4 \text{tr}(A^2) \leq C, \quad a.s. \end{aligned} \quad (99)$$

We can now write

$$\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}} = \frac{r_n^{-1/2} T' A (Y - T \beta^{\text{SIVE}}) / (r_n^{1/2} H_n)}{1 + r_n^{1/2} (r_n^{-1} T' A T - H_n) / (r_n^{1/2} H_n)} = \frac{A_n}{B_n}. \quad (100)$$

Assume that $r_n^{-1/2} H_n^{-1} \rightarrow_p 0$. From (98) it follows that $A_n \rightarrow_p 0$. From (99) we have that $r_n^{1/2} (r_n^{-1} T' A T - H_n) / (r_n^{1/2} H_n) \rightarrow_p 0$, and hence $B_n \rightarrow_p 1$. We conclude that when $(r_n^{1/2} H_n)^{-1} = (r_n^{-1/2} \pi' Z' M_W Z \pi)^{-1} \rightarrow_p 0$, the SIVE estimator is consistent: $\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}} \rightarrow_p 0$. Finally, we have under weak identification, $r_n^{1/2} H_n = \sqrt{G} \frac{\pi' Z' M_W Z \pi}{G} = \sqrt{G} \frac{n\text{FS}}{G} = \sqrt{G} \mu_n$ so that requiring that $r_n^{-1/2} H_n^{-1} \rightarrow_p 0$ is the same as the requirement in [Lemma 6](#) that $\sqrt{G} \mu_n \rightarrow_p \infty$. This completes the proof.

D.3 Proof Theorem 2

The proof of Theorem 2 is long, so we start with a brief overview of the steps. First, we establish asymptotic normality of the score vector S_n defined in (97) after normalizing it with the square root of its conditional variance. This step relies on a central limit theorem for quadratic forms of growing rank. We then show that the difference between the estimator for the conditional variance given in (20) and the population conditional variance converges in probability to a quantity larger than zero (under weak identification) or equal to zero (under strong identification). This step uses the properties of the regression error variance estimators (21). The derivation is somewhat tedious as it requires to keep track of higher-order interactions of the first stage and reduced form regression errors. Finally, we analyze the use of the variance estimator (22) when we encounter covariate groups with $m_g = 2$ or $n_g - m_g = 2$.

D.3.1 Asymptotic normality of the self-normalized score

As in the previous subsection, we rewrite

$$r_n^{1/2} (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}) = \frac{r_n^{-1/2} T' A (Y - T \beta^{\text{SIVE}})}{r_n^{-1} T' A T} = \frac{S_n}{r_n^{-1/2} D_n}. \quad (101)$$

The variance of S_n is given in the first line of (98). We now study the self-normalized version of S_n given by

$$S_n^* = \frac{r_n^{-1/2} T' A (Y - T \beta)}{\mathbb{V}(S_n | Q, X)^{1/2}}. \quad (102)$$

To show that $S_n^* \rightarrow_d N(0, 1)$, we use Lemma D.2 in [Evdokimov and Kolesár \(2018\)](#), which in turn is based on Lemma A.2 in [Chao et al. \(2012\)](#).

Lemma 7 [[Evdokimov and Kolesár \(2018\)](#), Lemma D.2] *Define $R = t' u + s' v + u' \bar{A} v$, where $[\bar{A}]_{ii} = 0$. Suppose that the following holds almost surely.*

1. $\mathbb{E}[u_i | Q, X] = \mathbb{E}[v_i | Q, X] = 0$ and $\mathbb{E}[u_i^4 | Q, X] \leq C < \infty$ and $\mathbb{E}[v_i^4 | Q, X] \leq C < \infty$,
2. $\mathbb{V}(R | Q, X)^{-1} \leq C < \infty$,
3. $\sum_{i=1}^n (t_i^4 + s_i^4) \rightarrow 0$,
4. $\text{tr}(\bar{A}^4) \rightarrow 0$.

Then, $\mathbb{V}(R | Q, X)^{-1/2} R \rightarrow_d N(0, 1)$.

Let $R = S_n$, so that $t = r_n^{-1/2} M_W Z \zeta$ and $s = r_n^{-1/2} M_W Z \pi$ and $\bar{A} = r_n^{-1/2} A$. Moreover, u_i is as defined in (55) and $v_i = \varepsilon_i - u_i \beta^{\text{SIVE}}$ with ε_i as in (56).

Verifying Conditions 1–4. Condition 1: holds by [Assumption 3](#). Condition 3: From (79) we see that $|e_i' M_W Z \pi| \leq C < \infty$ a.s. uniformly over i . Similarly, $|e_i' M_W Z \zeta| \leq C < \infty$ a.s. uniformly over i . We then obtain

$$\sum_{i=1}^n (t_i^4 + s_i^4) \leq C r_n^{-2} \sum_{i=1}^n (t_i^2 + s_i^2) = C r_n^{-1} (\pi' Z' M_W Z \pi / r_n + \zeta' Z' M_W Z \zeta / r_n). \quad (103)$$

The term in brackets is almost surely bounded as shown in [Appendix D.2](#). Since $r_n^{-1} \rightarrow_{a.s.} 0$, Condition (3) of [Lemma 7](#) holds. To verify Condition 2 and Condition 4 we separately consider strong and weak/vanishing identification.

Strong identification We start by showing that in [Lemma 7](#) we have $u' \bar{A} v = o_p(1)$ so that we only have to verify Condition 2. We first note that

$$\begin{aligned} \mathbb{V}(G^{-1/2} u' A v | Q, X) &= G^{-1} \sum_{j \neq i} \sum_{k \neq m} A_{ij} A_{mk} \mathbb{E}[u_i u_k v_j v_m | Q, X] \\ &= G^{-1} \sum_{j \neq i} A_{ij}^2 (\mathbb{E}[u_i v_i u_j v_j | Q, X] + \mathbb{E}[u_i^2 v_j^2 | Q, X]) \\ &\leq C G^{-1} \text{tr}(A^2) \leq C \quad a.s., \end{aligned} \quad (104)$$

using that u_i and v_i have bounded fourth conditional moment by [Assumption 3](#) and the bound established in (77). Then, by the definition of r_n in the case of strong identification,

$$\mathbb{V}(u' A v | Q, X) / r_n \leq C G / r_n \rightarrow_{a.s.} 0. \quad (105)$$

We can then conclude that $r_n^{-1/2} u' A v = o_p(1)$. Therefore, it suffices to consider $\tilde{R} = t' u + s' v$ and we only need to verify Condition 2. As in (98),

$$\mathbb{V}(\tilde{R} | Q, X) = r_n^{-1} \sum_{i=1}^n \sigma_{u,i}^2 [M_W Z \zeta]_i^2 + \sigma_{v,i}^2 [M_W Z \pi]_i^2 + 2\sigma_{uv,i} [M_W Z \pi]_i [M_W Z \zeta]_i. \quad (106)$$

Denote $g(i)$ a function such that $g(i) = g$ if $w_{ig} = 1$. By [Assumption 3](#), using [\(79\)](#),

$$\begin{aligned}\mathbb{V}(\tilde{R}|Q, X) &= r_n^{-1} \sum_{i=1}^n c_i^2 \left(\zeta(x_{g(i)})^2 \sigma_{u,i}^2 + \pi(x_{g(i)})^2 \sigma_{v,i}^2 + 2\pi(x_{g(i)})\zeta(x_{g(i)})\sigma_{uv,i} \right) \\ &= r_n^{-1} \sum_{i=1}^n c_i^2 [\zeta(x_{g(i)}), \pi(x_{g(i)})] \begin{pmatrix} 1 & 0 \\ -\beta & 1 \end{pmatrix} \Sigma_i \begin{pmatrix} 1 & -\beta \\ 0 & 1 \end{pmatrix} [\zeta(x_{g(i)}), \pi(x_{g(i)})]' \\ &\geq C r_n^{-1} \sum_{i=1}^n c_i^2 (\zeta(x_{g(i)})^2 + \pi(x_{g(i)})^2),\end{aligned}\tag{107}$$

where $\Sigma_i = \begin{pmatrix} \sigma_{u,i}^2 & \sigma_{u\varepsilon,i} \\ \sigma_{u\varepsilon,i} & \sigma_{\varepsilon_i}^2 \end{pmatrix}$ and $c_i = \sum_g Z_{ig} \frac{n_g - m_g}{n_g} - (W_{ig} - Z_{ig}) \frac{m_g}{n_g}$. We have

$$\begin{aligned}r_n^{-1} \sum_{i=1}^n c_i^2 \pi(x_{g(i)})^2 &= r_n^{-1} \sum_{g=1}^G \pi(x_g)^2 \left[m_g \left(\frac{n_g - m_g}{n_g} \right)^2 + (n_g - m_g) \left(\frac{m_g}{n_g} \right)^2 \right] \\ &= r_n^{-1} \sum_{g=1}^G \pi(x_g)^2 n_g \frac{m_g (n_g - m_g)}{n_g^2} \\ &= 1,\end{aligned}\tag{108}$$

where the last line uses the definition of r_n . We conclude that $\mathbb{V}(\tilde{R}|Q, X) \geq C > 0$ *a.s.* and [Condition 2](#) of [Lemma 7](#) holds.

Weak identification We start with verifying [Condition 2](#) in [Lemma 7](#). By [Assumption 3](#), $\min_{i=1,\dots,n} \sigma_{v,i}^2 \geq C > 0$ and $\min_{i=1,\dots,n} \sigma_{u,i}^2 \geq C > 0$. Then, using [\(78\)](#)

$$\begin{aligned}\mathbb{V}(S_n|Q, X) &\geq G^{-1} \text{tr}(D_{\sigma_v^2} A D_{\sigma_u^2} A) \\ &\geq C G^{-1} \text{tr}(A^2) \\ &\geq C > 0 \quad \textit{a.s.}\end{aligned}\tag{109}$$

For [Condition 4](#), we have that $G^{-1} \text{tr}(A^2) \leq 3$ *a.s.* by [\(77\)](#) and $\lambda_{\max}(A^2) = 1$ since $\lambda_{\max}(A) = 1$. Then,

$$\text{tr}(\bar{A}_n^4) = r_n^{-2} \text{tr}(A^4) \leq \frac{G^2 \text{tr}(A^2) \lambda_{\max}(A^2)}{r_n^2 G} \rightarrow_{\textit{a.s.}} 0.\tag{110}$$

We note that these results are shown without any assumptions on the identification strength, so that we can base identification-robust inference on the asymptotic normality of the score vector.

We have now established asymptotic normality of the renormalized score vector under both weak and strong identification. To prove [Theorem 2](#) the next step is to show that the numerator of the variance estimator in [\(20\)](#) converges to $\mathbb{V}(S_n|Q, X)$.

D.3.2 Consistency of the estimator for the variance of the score

We momentarily assume that $m_g \geq 3$ and $n_g - m_g \geq 3$ so that we can use [\(21\)](#) for the regression error variances and covariances. We relax this assumption in [Appendix D.3.3](#).

We can rewrite the variance given in [\(98\)](#) as

$$\begin{aligned}\mathbb{V}(S_n|Q, X) &= r_n^{-1} \left[\zeta' Z' M_W D_{\sigma_u^2} M_W Z \zeta + \pi' Z' M_W D_{\sigma_v^2} M_W Z \pi + 2\zeta' Z' M_W D_{\sigma_{uv}} M_W Z \pi \right. \\ &\quad \left. + \underbrace{\iota' D_{\sigma_{uv}} (A \odot A) D_{\sigma_{uv}} \iota + \iota' D_{\sigma_v^2} (A \odot A) D_{\sigma_u^2} \iota}_{(B.0)} \right].\end{aligned}\tag{111}$$

In (20) we use the following estimator for $\mathbb{V}(S_n|Q, X)$.

$$\hat{\mathbb{V}}(S|Q, X) = r_n^{-1} \left[(Y - T\hat{\beta}^{\text{SIVE}})' AD_{\hat{\sigma}_u^2} A (Y - T\hat{\beta}^{\text{SIVE}}) + T' AD_{\hat{\sigma}_v^2} AT + 2(Y - T\hat{\beta}^{\text{SIVE}})' AD_{\hat{\sigma}_{uv}} AT \right]. \quad (112)$$

Let D_{u^2} be a diagonal matrix with $[D_{u^2}]_{ii} = u_i^2$ and likewise for D_{v^2} and D_{uv} . Consider the infeasible variance estimator

$$\mathbb{V}_{\text{inf}}(S_n|Q, X) = r_n^{-1} \left[\zeta' Z' M_W D_{u^2} M_W Z \zeta + \pi' Z' M_W D_{v^2} M_W Z \pi + 2\zeta' Z' M_W D_{uv} M_W Z \pi + \iota' D_{uv} (A \odot A) D_{uv} \iota + \iota' D_{v^2} (A \odot A) D_{u^2} \iota \right]. \quad (113)$$

We decompose (112) as

$$\begin{aligned} \hat{\mathbb{V}}(S|Q, X) &= r_n^{-1} \left[(Z\zeta + v - T(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}))' AD_{\hat{\sigma}_u^2} A (Z\zeta + v - T(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})) \right. \\ &\quad \left. + (Z\pi + u)' AD_{\hat{\sigma}_v^2} A (Z\pi + u) + 2(Z\zeta + v - T(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}))' AD_{\hat{\sigma}_{uv}} A (Z\pi + u) \right] \\ &= r_n^{-1} \left[(Z\zeta + v)' AD_{\hat{\sigma}_u^2} A (Z\zeta + v) + (Z\pi + u)' AD_{\hat{\sigma}_v^2(\beta^{\text{SIVE}})} A (Z\pi + u) \right. \\ &\quad \left. + 2(Z\zeta + v)' AD_{\hat{\sigma}_{uv}(\beta^{\text{SIVE}})} A (Z\pi + u) \right] + R_1 \\ &= \mathbb{V}(S_n|Q, X) + \mathbb{V}_{\text{inf}}(S_n|Q, X) - \mathbb{V}(S_n|Q, X) \\ &\quad + \underbrace{r_n^{-1} \iota' D_v^2 (A \odot A) D_u^2 \iota + r_n^{-1} \iota' D_{uv} (A \odot A) D_{uv} \iota}_{(B.1)} + \underbrace{r_n^{-1} v' AD_{\hat{\sigma}_u^2} Av - \iota' D_v^2 (A \odot A) D_u^2 \iota}_{(B.2)} \\ &\quad + \underbrace{r_n^{-1} u' AD_{\hat{\sigma}_v^2(\beta^{\text{SIVE}})} Au - \iota' D_u^2 (A \odot A) D_v^2 \iota + 2(v' AD_{\hat{\sigma}_{uv}(\beta^{\text{SIVE}})} Au - \iota' D_{uv} (A \odot A) D_{uv} \iota)}_{(B.3)} \\ &\quad + \underbrace{r_n^{-1} \zeta' Z' M_W (D_{\hat{\sigma}_u^2} - D_{u^2}) M_W Z \zeta}_{(Z.1)} + \underbrace{2r_n^{-1} v' AD_{\hat{\sigma}_u^2} M_W Z \zeta}_{(Z.2)} \\ &\quad + r_n^{-1} \pi' Z' M_W (D_{\hat{\sigma}_v^2(\beta^{\text{SIVE}})} - D_{v^2}) M_W Z \pi + 2r_n^{-1} \zeta' Z' M_W (D_{\hat{\sigma}_{uv}(\beta^{\text{SIVE}})} - D_{uv}) M_W Z \pi \\ &\quad + 2r_n^{-1} u' AD_{\hat{\sigma}_v^2(\beta^{\text{SIVE}})} M_W Z \pi + 2r_n^{-1} \zeta' Z' M_W D_{\hat{\sigma}_{uv}(\beta^{\text{SIVE}})} Au + 2r_n^{-1} \pi' Z' M_W D_{\hat{\sigma}_{uv}(\beta^{\text{SIVE}})} Av \\ &\quad + R_1, \end{aligned} \quad (114)$$

where we use (95) and we define the remainder term

$$R_1 = 4r_n^{-1} \left[(\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}})^2 (Z\pi + u)' AD_{\hat{\sigma}_u^2} A (Z\pi + u) - (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}) (Z\pi + u)' AD_{\hat{\sigma}_{uv}(\beta^{\text{SIVE}})} A (Z\pi + u) - (\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}) (Z\pi + u)' AD_{\hat{\sigma}_u^2} A (Z\zeta + v) \right]. \quad (115)$$

When we do identification-robust inference, under the null we replace $\hat{\beta}^{\text{SIVE}} = \beta^{\text{SIVE}}$ so that $R_1 = 0$. Therefore, the proof of [Corollary 1](#) follows from the proof presented here for [Theorem 2](#).

We now will prove the following. First, $\mathbb{V}_{\text{inf}}(S_n|Q, X) - \mathbb{V}(S_n|Q, X) \rightarrow_p 0$. Second, we show that under weak identification (B.1), (B.2) and (B.3) from (114) together converge to a three times the value of (B.0) defined in (111). Finally, we show that (Z.1) and (Z.2) converge to zero in probability. The remaining terms converge to zero by the same arguments.

Part 1: $\mathbb{V}_{\text{inf}}(S|Q, X) - \mathbb{V}(S|Q, X) \rightarrow_p 0$. It is sufficient to show that

$$r_n^{-1} \zeta' Z' M_W (D_u^2 - D_{\sigma_u^2}) M_W Z \zeta \rightarrow_p 0, \quad (116)$$

$$r_n^{-1} (\iota' D_v^2 (A \odot A) D_u^2 \iota - \iota' D_{\sigma_v^2} (A \odot A) D_{\sigma_u^2} \iota) \rightarrow_p 0, \quad (117)$$

The other terms follow by the same arguments.

For (116), $\mathbb{E}[r_n^{-1}\zeta'Z'M(D_u^2 - D_{\sigma_u^2})MZ\zeta] = 0$, and the variance can be upper bounded as

$$\begin{aligned} \mathbb{V}(r_n^{-1}\zeta'Z'M_W(D_u^2 - D_{\sigma_u^2})M_WZ\zeta|Q, X) &= r_n^{-2} \sum_{i=1}^n \mathbb{E}[(u_i^2 - \sigma_{u,i}^2)^2|Q, X][MZ\zeta]_i^4 \\ &\leq C r_n^{-2} \sum_{i=1}^n [MZ\zeta]_i^2 \rightarrow_p 0. \end{aligned} \quad (118)$$

For (117), $\mathbb{E}[r_n^{-1}(l'D_v^2(A \odot A)D_u^2 l - l'D_{\sigma_v}^2(A \odot A)D_{\sigma_u}^2 l)|Q, X] = 0$. The variance can be bounded as

$$\begin{aligned} &\mathbb{V}(r_n^{-1}(l'D_v^2(A \odot A)D_u^2 l - l'D_{\sigma_v}^2(A \odot A)D_{\sigma_u}^2 l)|Q, X) \\ &= r_n^{-2} \sum_{i=1}^n \sum_{j \neq i}^n \sum_{k=1}^n \sum_{l \neq k}^n \mathbb{E}[(v_i^2 u_j^2 - \sigma_{v,i}^2 \sigma_{u,i}^2)(v_k^2 u_l^2 - \sigma_{v,k}^2 \sigma_{u,l}^2)|Q, X] A_{ij}^2 A_{kl}^2 \\ &= r_n^{-2} \sum_{i=1}^n \sum_{j \neq i}^n (\mathbb{E}[(v_i^2 u_j^2 - \sigma_{v,i}^2 \sigma_{u,i}^2)^2|Q, X] + \mathbb{E}[(v_i^2 u_j^2 - \sigma_{v,i}^2 \sigma_{u,i}^2)(v_j^2 u_i^2 - \sigma_{v,j}^2 \sigma_{u,i}^2)|Q, X]) A_{ij}^4 \\ &\leq C r_n^{-2} \sum_{i=1}^n \sum_{j \neq i}^n A_{ij}^4 \leq C r_n^{-2} \text{tr}(A^2) \rightarrow_p 0, \end{aligned} \quad (119)$$

where we use [Assumption 3.4](#) and (77).

Part 2. It follows from the results in Part 1 that (B.1)–(B.0) $\rightarrow_p 0$. For (B.2), we decompose it further as

$$(B.2) = \underbrace{r_n^{-1}(v'AD_{\sigma_u^2}Av - v'AD_{u^2}Av)}_{(B.2a)} + \underbrace{r_n^{-1}(v'AD_{u^2}Av - l'D_v^2(A \odot A)D_u^2 l)}_{(B.2b)}. \quad (120)$$

For (B.2b),

$$\mathbb{E}[r_n^{-1}v'AD_u^2Av - l'D_v^2(A \odot A)D_u^2 l|Q, X] = r_n^{-1} \sum_{i=1}^n \sum_{k \neq i} \sum_{j \neq \{i,k\}} \mathbb{E}[v_i v_k u_j^2|Q, X] A_{ij} A_{kj} = 0. \quad (121)$$

The variance can be upper bounded as

$$\begin{aligned} &\mathbb{V}(r_n^{-1}(v'AD_u^2Av - l'D_v^2(A \odot A)D_u^2 l)|Z, W) \\ &= r_n^{-2} \mathbb{E} \left[\left(\sum_{i=1}^n \sum_{k \neq i} \sum_{j \neq \{i,k\}} v_i v_k u_j^2 A_{ij} A_{kj} \right)^2 \middle| Q, X \right] \\ &= r_n^{-2} \sum_{i_1, i_2} \sum_{i_3 \neq i_1} \sum_{i_4 \neq i_2} \sum_{i_5 \neq \{i_3, i_1\}} \sum_{i_6 \neq \{i_2, i_4\}} \mathbb{E}[v_{i_1} v_{i_2} v_{i_3} v_{i_4} (u_{i_5}^2 - \sigma_{u, i_5}^2)(u_{i_6}^2 - \sigma_{u, i_6}^2)|Q, X] A_{i_1 i_3} A_{i_1 i_5} A_{i_2 i_4} A_{i_2 i_6} \\ &\quad + r_n^{-2} \sum_{i_1, i_2} \sum_{i_3 \neq i_1} \sum_{i_4 \neq i_2} \sum_{i_5 \neq \{i_3, i_1\}} \sum_{i_6 \neq \{i_2, i_4\}} \mathbb{E}[v_{i_1} v_{i_2} v_{i_3} v_{i_4} |Q, X] \sigma_{u, i_5}^2 \sigma_{u, i_6}^2 A_{i_1 i_3} A_{i_1 i_5} A_{i_2 i_4} A_{i_2 i_6} \\ &\leq C r_n^{-2} \left[\sum_{i_1 \dots i_3} (\tilde{A}_{i_1 i_2}^2 \tilde{A}_{i_1 i_3}^2 + \tilde{A}_{i_1 i_2}^2 \tilde{A}_{i_1 i_3} \tilde{A}_{i_2 i_3}) + \sum_{i_1 \dots i_4} (\tilde{A}_{i_1 i_2}^2 \tilde{A}_{i_1 i_3} \tilde{A}_{i_1 i_4} + \tilde{A}_{i_1 i_2}^2 \tilde{A}_{i_1 i_3} \tilde{A}_{i_2 i_4}) \right] \\ &= C r_n^{-2} \left[l'(\tilde{A} \odot \tilde{A})^2 l + \text{tr}(\tilde{A}(\tilde{A} \odot \tilde{A})\tilde{A}) + \sum_{i_1} l'(\tilde{A} \odot \tilde{A})e_{i_1} (e_{i_1} \tilde{A} l)^2 + l' \tilde{A}(\tilde{A} \odot \tilde{A})\tilde{A} l \right] \\ &\leq C r_n^{-2} G \rightarrow_p 0, \end{aligned} \quad (122)$$

where the last inequality follows from (R9)–(R11) in (87).

We now turn to (B.2a). Using the expression for $\hat{\sigma}_u^2$ from [Appendix D.1.4](#),

$$\begin{aligned}
(B.2a) &= \underbrace{-r_n^{-1} \sum_{i=1}^n \sum_{j \neq i}^n \frac{2}{1-2\bar{P}_{ii}} (e_i' A v)^2 u_i \bar{P}_{ij} u_j}_{(B.2a.1)} \\
&\quad + \underbrace{r_n^{-1} \sum_{i=1}^n \sum_{j=1}^n \sum_{k \neq j}^n \frac{1}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})} (e_i' A v)^2 \bar{P}_{ij} \bar{P}_{ik} u_j u_k}_{(B.2a.2)}.
\end{aligned} \tag{123}$$

Since $[A]_{ii} = 0$, we have that $\mathbb{E}[(B.2a.1)|Q, X] = 0$. For (B.2a.2), we have

$$\mathbb{E}[(B.2a.2)|Q, X] = r_n^{-1} \sum_{i=1}^n \frac{1}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})} \sum_{j=1}^n \sum_{k \neq j}^n \bar{P}_{ij} \bar{P}_{ik} A_{ij} A_{ik} \sigma_{uv,j} \sigma_{uv,k}. \tag{124}$$

Define $B = (A \odot P) J (A \odot P)$, where J is diagonal with $[J]_{ii} = (1-\bar{P}_{ii})^{-1} (1-2\bar{P}_{ii})^{-1}$. Also define $C = B - D_B$, so the matrix B with diagonal elements set to zero. Then, C is a block diagonal matrix with blocks

$$C_{(g)} = \begin{pmatrix} \frac{(n_g - m_g)^2}{n_g^2 (m_g - 1)^2 (m_g - 2)} (\iota_{m_g} \iota_{m_g}' - I_{m_g}) & O \\ O & \frac{m_g^2}{n_g^2 (n_g - m_g - 1)^2 (n_g - m_g - 2)} (\iota_{n_g - m_g} \iota_{n_g - m_g}' - I_{n_g - m_g}) \end{pmatrix}. \tag{125}$$

We will use below that $C_{ij} \leq A_{ij}^2$. For the bias arising from (B.2a.2), we have

$$\mathbb{E}[(B.2a.2)|Q, X] = r_n^{-1} \sigma_{uv}' C \sigma_{uv}. \tag{126}$$

Similarly, (B.3) will yield a bias totaling up to

$$\begin{aligned}
r_n^{-1} \mathbb{E}[(B.2) + (B.3)|Q, X] &= 2r_n^{-1} \sigma_{uv}' C \sigma_{uv} + 2r_n^{-1} \sigma_u^2 C \sigma_v^2 \\
&= 2r_n^{-1} \sum_{i=1}^n \sum_{j \neq i}^n C_{ij} (\sigma_{uv,i} \sigma_{uv,j} + \sigma_{u,i}^2 \sigma_{v,j}^2) \\
&= 2r_n^{-1} \sum_{i=1}^n \sum_{j > i}^n C_{ij} (2\sigma_{uv,i} \sigma_{uv,j} + \sigma_{u,i}^2 \sigma_{v,j}^2 + \sigma_{v,i}^2 \sigma_{u,j}^2) \\
&= 2r_n^{-1} \sum_{i=1}^n \sum_{j > i}^n C_{ij} \mathbb{E}[(u_i v_j + u_j v_i)^2 | Q, X] \geq 0,
\end{aligned} \tag{127}$$

where the conclusion holds since $C_{ij} \geq 0$ for $i \neq j$. The bias is bounded, as

$$r_n^{-1} \sum_{i=1}^n \sum_{j > i}^n C_{ij} \mathbb{E}[(u_i v_j + u_j v_i)^2 | Q, X] \leq C r_n^{-1} \sum_{i=1}^n \sum_{j > i}^n C_{ij} \leq C r_n^{-1} \sum_{i=1}^n \sum_{j > i}^n A_{ij}^2 \leq C r_n^{-1} G \leq C. \tag{128}$$

We conclude that

$$\begin{aligned}
\mathbb{E}[(B.0) + (B.1) + (B.2) + (B.3)|Q, X] &= 2\iota' D_{\sigma_{uv}} [(A \odot A) + C] D_{\sigma_{uv}} \iota + 2\iota' D_{\sigma_v^2} [(A \odot A) + C] D_{\sigma_u^2} \iota \\
&\leq 4(\iota' D_{\sigma_{uv}} (A \odot A) D_{\sigma_{uv}} \iota + \iota' D_{\sigma_v^2} (A \odot A) D_{\sigma_u^2} \iota),
\end{aligned} \tag{129}$$

where the inequality uses that $C_{ij} \leq A_{ij}^2$. We conclude that under weak identification, we may overestimate the variance term (B.0) by a factor 4 if all other variance terms are negligible and to a lesser degree when the remaining variance terms are not negligible.

It remains to be shown that the terms (B.2) and (B.3) converge to their expectation. We show this for (B.2). The result for (B.3) follows analogously. We first note that the result for (B.2b) is already established in (122). For the variance of (B.2a.1) and (B.2a.2), we have

$$\begin{aligned}
\mathbb{V}(B.2a.1|Q, X) &= 4r_n^{-2} \sum_{i_1=1}^n \sum_{i_3=1}^n \sum_{i_2 \neq i_1} \sum_{i_4 \neq i_3} \frac{1}{(1-2\bar{P}_{i_1 i_1})(1-2\bar{P}_{i_3 i_3})} \mathbb{E}[(e'_{i_1} Av)^2 (e'_{i_3} Av)^2 \bar{P}_{i_1 i_2} \bar{P}_{i_3 i_4} u_{i_1} u_{i_2} u_{i_3} u_{i_4} | Q, X] \\
&\leq C r_n^{-2} \sum_{i_1=1}^n \sum_{i_3=1}^n \sum_{i_2 \neq i_1} \sum_{i_4 \neq i_3} \sum_{i_5 \neq i_1} \sum_{i_6 \neq i_1} \sum_{i_7 \neq i_3} \sum_{i_8 \neq i_3} |A_{i_1 i_5}| |A_{i_1 i_6}| |A_{i_3 i_7}| |A_{i_3 i_8}| \bar{P}_{i_1 i_2} \bar{P}_{i_3 i_4} \\
&\quad \times |\mathbb{E}[u_{i_1} u_{i_2} u_{i_3} u_{i_4} v_{i_5} v_{i_6} v_{i_7} v_{i_8} | Q, X]|, \\
\mathbb{V}(B.2a.2|Q, X) &= r_n^{-2} \sum_{i_1=1}^n \sum_{i_2=1}^n \sum_{i_6=1}^n \sum_{i_7=1}^n \sum_{i_3 \neq i_2} \sum_{i_4 \neq i_1} \sum_{i_5 \neq i_1} \sum_{i_8 \neq i_7} \sum_{i_9 \neq i_6} \sum_{i_{10} \neq i_6} c_{i_1} c_{i_6} A_{i_1 i_4} A_{i_6 i_9} A_{i_1 i_5} A_{i_6 i_{10}} \bar{P}_{i_1 i_2} \bar{P}_{i_1 i_3} \bar{P}_{i_6 i_7} \bar{P}_{i_6 i_8} \\
&\quad \times \mathbb{E}[u_{i_2} u_{i_3} v_{i_4} v_{i_5} u_{i_7} u_{i_8} v_{i_9} v_{i_{10}} | Q, X] \\
&\quad - 4r_n^{-2} \sum_{i_1=1}^n \sum_{i_2 \neq i_1} \sum_{i_3 \neq \{i_2, i_1\}} \sum_{i_4=1}^n \sum_{i_5 \neq i_4} \sum_{i_6 \neq \{i_5, i_4\}} c_{i_1} c_{i_4} A_{i_1 i_2} A_{i_4 i_5} A_{i_1 i_3} A_{i_4 i_6} \bar{P}_{i_1 i_2} \bar{P}_{i_4 i_5} \bar{P}_{i_1 i_3} \bar{P}_{i_4 i_6} \\
&\quad \times \mathbb{E}[v_{i_2} v_{i_3} u_{i_2} u_{i_3} v_{i_5} v_{i_6} u_{i_5} u_{i_6} | Q, X],
\end{aligned} \tag{130}$$

where $c_i = (1 - \bar{P}_{ii})^{-1}(1 - 2\bar{P}_{ii})^{-1}$. To bound these expressions, we use independence of u_i and v_i across i and take into account the restrictions on the indices as indicated under the summations signs. The conditional expectations of products of u_i and v_i in these terms are all almost surely bounded by Assumption 3, so we can take these out of the summations. Finding the nonzero terms in (130) is now a combinatorial exercise that can be executed using symbolic programming. The nonzero terms are listed and bounded in Appendix D.3.4. The results show that both variances in (130) converge to zero in probability.

Part 3. Starting with (Z.1), and using the expressions in (96) we have

$$\begin{aligned}
(Z.1) &= r_n^{-1} \zeta' Z' M_W (D_{\hat{\sigma}_u^2} - D_{u^2}) M_W Z \zeta \\
&= -r_n^{-1} \underbrace{\sum_{i=1}^n \sum_{j \neq i} \frac{2}{1 - \bar{P}_{ii}} [M_W Z \zeta]_i^2 u_i \bar{P}_{ij} u_j}_{(Z.1a)} \\
&\quad + r_n^{-1} \underbrace{\sum_{i=1}^n \sum_{j=1}^n \sum_{k \neq j} \frac{1}{(1 - \bar{P}_{ii})(1 - 2\bar{P}_{ii})} [M_W Z \zeta]_i^2 \bar{P}_{ij} \bar{P}_{ik} u_j u_k}_{(Z.1b)}.
\end{aligned} \tag{131}$$

Using independence across i , we find that (Z.1a) has expectation zero, conditional on Q and X . For the variance of the individual components, we have,

$$\begin{aligned}
\mathbb{V}(Z.1a|Q, X) &= r_n^{-2} \sum_{i=1}^n \sum_{k=1}^n \frac{8}{(1-2\bar{P}_{ii})(1-2\bar{P}_{kk})} \sigma_{u,i}^2 \sigma_{u,k}^2 \bar{P}_{ik}^2 [M_W Z \zeta]_i^2 [M_W Z \zeta]_k^2 \\
&\leq C \max_k [M_W Z \zeta]_k^2 r_n^{-2} \sum_{i=1}^n [M_W Z \zeta]_i^2 \rightarrow_p 0,
\end{aligned} \tag{132}$$

where we use that $[M_W Z \zeta]_i^2 \leq C < \infty$ a.s. and $\sum_{k=1}^n \bar{P}_{ik}^2 = \bar{P}_{ii} < 1/2$.

$$\begin{aligned}
\mathbb{V}(Z.1b|Q, X) &= r_n^{-2} \sum_{i=1}^n \sum_{l=1}^n \frac{2}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})(1-\bar{P}_{ll})(1-2\bar{P}_{ll})} \\
&\quad \times \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} \bar{P}_{lj} \bar{P}_{lk} \sigma_{u,j}^2 \sigma_{u,k}^2 [M_W Z \zeta]_i^2 [M_W Z \zeta]_l^2 \\
&\leq C r_n^{-2} \sum_{i=1}^n \sum_{l=1}^n [M_W Z \zeta]_i^2 [M_W Z \zeta]_l^2 \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} \bar{P}_{lj} \bar{P}_{lk} \\
&\leq C \max_l [M_W Z \zeta]_l^2 r_n^{-2} \sum_{i=1}^n [M_W Z \zeta]_i^2 \rightarrow_p 0,
\end{aligned} \tag{133}$$

using that $\sum_{l=1}^n \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} \bar{P}_{lj} \bar{P}_{lk} = \sum_{j=1}^n \sum_{k \neq j} \bar{P}_{ij} \bar{P}_{ik} \bar{P}_{jk} = \sum_{j=1}^n \bar{P}_{ij}^2 - \sum_{j=1}^n \bar{P}_{ij}^2 \bar{P}_{jj} \leq \bar{P}_{ii} < 1/2$. From (132) and (133) we conclude that

$$r_n^{-1} \zeta' Z' M_W (D_{\hat{\sigma}_u^2} - D_{u^2}) M_W Z \zeta \rightarrow_p 0. \tag{134}$$

By the same arguments,

$$r_n^{-1} \pi' Z' M_W (D_{\hat{\sigma}_v^2} - D_{v^2}) M_W Z \pi \rightarrow_p 0, \quad r_n^{-1} \zeta' Z' M_W (D_{\hat{\sigma}_{uv}} - D_{uv}) M_W Z \pi \rightarrow_p 0. \tag{135}$$

We now turn to the final variance term (Z.2).

$$\begin{aligned}
(Z.2) &= 2 r_n^{-1} \underbrace{\sum_{k=1}^n \sum_{i \neq k} v_k A_{ik} [M_W Z \zeta]_i u_i \frac{2}{1-2\bar{P}_{ii}} \sum_{j \neq i} \bar{P}_{ij} u_j}_{(Z.2a)} \\
&\quad + 2 r_n^{-1} \underbrace{\sum_{k=1}^n \sum_{i \neq k} v_k A_{ik} [M_W Z \zeta]_i \frac{1}{(1-\bar{P}_{ii})(1-2\bar{P}_{ii})} \sum_{j=1}^n \sum_{s \neq j} \bar{P}_{ij} \bar{P}_{is} u_j u_s}_{(Z.2b)}.
\end{aligned} \tag{136}$$

We have $\mathbb{E}[(Z.2)|Q, X] = 0$. For the variance,

$$\begin{aligned}
\mathbb{V}(Z.2a|Q, X) &= r_n^{-2} \sum_{k=1}^n \sum_{l=1}^n \sum_{i \neq k} \sum_{m \neq l} A_{ik} A_{ml} [M_W Z \zeta]_i [M_W Z \zeta]_m \\
&\quad \times \frac{2}{1-2\bar{P}_{ii}} \frac{2}{1-2\bar{P}_{mm}} \sum_{j \neq i} \bar{P}_{ij} \sum_{s \neq m} \bar{P}_{ms} \mathbb{E}[v_k v_l u_i u_m u_j u_s | Q, X] \\
&\leq C r_n^{-2} \sum_{i,j,k} [\tilde{A}_{ij}^2 \bar{P}_{ij} \bar{P}_{jk} + \tilde{A}_{ij} \tilde{A}_{ik} \bar{P}_{ik} \bar{P}_{jk} + A_{ij}^2 \bar{P}_{jk}^2 + \tilde{A}_{ij} \tilde{A}_{ik} \bar{P}_{jk}^2 + \tilde{A}_{ij} \tilde{A}_{jk} \bar{P}_{ik} \bar{P}_{jk} + \tilde{A}_{ik} \tilde{A}_{jk} \bar{P}_{ik} \bar{P}_{jk}].
\end{aligned} \tag{137}$$

We can bound each on the terms on the final line as follows.

$$\begin{aligned}
\sum_{i,j,k} \tilde{A}_{ij}^2 \tilde{P}_{ij} \tilde{P}_{jk} &\leq \sum_{i,j} \tilde{A}_{ij}^2 = \text{tr}(\tilde{A}^2) \leq CG, \\
\sum_{i,j,k} \tilde{A}_{ij} \tilde{A}_{ik} \tilde{P}_{ik} \tilde{P}_{jk} &\leq C \sum_{i,j,k} \tilde{A}_{ij} \tilde{P}_{ik} \tilde{P}_{jk} = C \sum_{i,j} \tilde{A}_{ij} \tilde{P}_{ij} = C \iota'(\tilde{A} \odot \tilde{P}) \iota \leq CG, \\
\sum_{i,j,k} A_{ij}^2 \tilde{P}_{jk}^2 &\leq \sum_{i,j} A_{ij}^2 = \text{tr}(\tilde{A}^2) \leq CG, \\
\sum_{i,j,k} \tilde{A}_{ij} \tilde{A}_{ik} \tilde{P}_{jk}^2 &\leq \text{tr}(\tilde{A} \tilde{P} \tilde{A}) \leq \text{tr}(\tilde{A}^2) \leq CG, \\
\sum_{i,j,k} \tilde{A}_{ij} \tilde{A}_{jk} \tilde{P}_{ik} \tilde{P}_{jk} &\leq C \sum_{i,j,k} \tilde{A}_{ij} \tilde{P}_{ik} \tilde{P}_{jk} = C \iota'(\tilde{A} \odot \tilde{P}) \iota \leq CG, \\
\sum_{i,j,k} \tilde{A}_{ik} \tilde{A}_{jk} \tilde{P}_{ik} \tilde{P}_{jk} &= \iota'(\tilde{A} \odot \tilde{P})^2 \iota \leq CG.
\end{aligned} \tag{138}$$

On the first line, we use that $\sum_k \tilde{P}_{jk} = 1$.

For the variance of (Z.2b), we have

$$\begin{aligned}
\mathbb{V}(Z.2b|Q, X) &\leq Cr_n^{-2} \sum_{k=1}^n \sum_{l=1}^n \sum_{i \neq k} \sum_{m \neq l} \sum_{j=1}^n \sum_{s \neq j} \sum_{r=1}^n \sum_{t \neq r} |A_{ik}| |A_{ml}| \tilde{P}_{ij} \tilde{P}_{is} \tilde{P}_{mr} \tilde{P}_{ms} \mathbb{E}[v_k v_l u_j u_s u_r u_t | Q, X] \\
&\leq Cr_n^{-2} \sum_{i_1, \dots, i_4} \tilde{A}_{i_1 i_3} \tilde{A}_{i_2 i_4} \tilde{P}_{i_1 i_3} \tilde{P}_{i_1 i_4} \tilde{P}_{i_2 i_3} \tilde{P}_{i_2 i_4} + Cr_n^{-2} \sum_{i_1, \dots, i_5} [\tilde{A}_{i_1 i_4} \tilde{A}_{i_1 i_5} \tilde{P}_{i_2 i_4} \tilde{P}_{i_2 i_5} \tilde{P}_{i_3 i_4} \tilde{P}_{i_3 i_5} \\
&\quad + \tilde{A}_{i_1 i_4} \tilde{A}_{i_2 i_5} \tilde{P}_{i_1 i_4} \tilde{P}_{i_2 i_5} \tilde{P}_{i_3 i_4} \tilde{P}_{i_3 i_5} + \tilde{A}_{i_1 i_4} \tilde{A}_{i_2 i_5} \tilde{P}_{i_1 i_5} \tilde{P}_{i_2 i_4} \tilde{P}_{i_3 i_4} \tilde{P}_{i_3 i_5}].
\end{aligned} \tag{139}$$

We bound each term on the final line as follows.

$$\begin{aligned}
\sum_{i_1, \dots, i_5} \tilde{A}_{i_1 i_4} \tilde{A}_{i_1 i_5} \tilde{P}_{i_2 i_4} \tilde{P}_{i_2 i_5} \tilde{P}_{i_3 i_4} \tilde{P}_{i_3 i_5} &= \sum_{i_1, i_4, i_5} \tilde{A}_{i_1 i_4} \tilde{A}_{i_1 i_5} \tilde{P}_{i_4 i_5}^2 \leq \text{tr}(\tilde{A} \tilde{P} \tilde{A}) \leq C \text{tr}(A^2) \leq CG, \\
\sum_{i_1, \dots, i_4} \tilde{A}_{i_1 i_3} \tilde{A}_{i_2 i_4} \tilde{P}_{i_1 i_3} \tilde{P}_{i_1 i_4} \tilde{P}_{i_2 i_3} \tilde{P}_{i_2 i_4} &= \text{tr}((\tilde{A} \odot \tilde{P}) \tilde{P} (\tilde{A} \odot \tilde{P}) \tilde{P}) \leq \iota'(\tilde{A} \odot \tilde{P})^2 \iota \leq CG, \\
\sum_{i_1, \dots, i_5} \tilde{A}_{i_1 i_4} \tilde{A}_{i_2 i_5} \tilde{P}_{i_1 i_4} \tilde{P}_{i_2 i_5} \tilde{P}_{i_3 i_4} \tilde{P}_{i_3 i_5} &\leq \sum_{i_1, i_2, i_4, i_5} \tilde{A}_{i_1 i_4} \tilde{A}_{i_2 i_5} \tilde{P}_{i_1 i_4} \tilde{P}_{i_2 i_5} \tilde{P}_{i_4 i_5} = \iota'(\tilde{A} \odot \tilde{P}) \tilde{P} (\tilde{A} \odot \tilde{P}) \iota \leq \iota'(\tilde{A} \odot \tilde{P})^2 \iota \leq CG, \\
\sum_{i_1, \dots, i_5} \tilde{A}_{i_1 i_4} \tilde{A}_{i_2 i_5} \tilde{P}_{i_1 i_5} \tilde{P}_{i_2 i_4} \tilde{P}_{i_3 i_4} \tilde{P}_{i_3 i_5} &\leq \sum_{i_1, i_2, i_4, i_5} \tilde{A}_{i_1 i_4} \tilde{A}_{i_2 i_5} \tilde{P}_{i_1 i_5} \tilde{P}_{i_2 i_4} = \text{tr}(\tilde{A} \tilde{P} \tilde{A} \tilde{P}) \leq \text{tr}(\tilde{A}^2) \leq CG.
\end{aligned} \tag{140}$$

We conclude that [Theorem 2](#) holds when [Assumption 2](#) would be strengthened to $m_g \geq 3$ and $n_g - m_g \geq 3$. We now discuss how to weaken the result to [Assumption 2](#).

D.3.3 Small groups

Consider a group with $m_g = 2$ or $n_g - m_g = 2$. To avoid the singularity in the variance estimators of [Appendix D.1.4](#), We can use the following rescaled version of the conventional variance estimator for the observations that only share their instrument status with one other observation in the group,

$$\begin{aligned}
\tilde{\sigma}_{u,i}^2 &= 4e_i'(M_{W,Z} u \odot M_{W,Z} u) \\
&= u_i^2 - 4u_i \sum_{j \neq i} \tilde{P}_{ij} u_j + 4 \sum_{j \neq i} \tilde{P}_{ij}^2 u_j^2.
\end{aligned} \tag{141}$$

Similarly, we find

$$\begin{aligned}\tilde{\sigma}_{v,i}^2(\beta^{\text{SIVE}}) &= 4e_i'(M_{W,Z}v \odot M_{W,Z}v) = u_i^2 - 4v_i \sum_{j \neq i} \bar{P}_{ij}v_j + 4 \sum_{j \neq i} \bar{P}_{ij}^2 v_j^2, \\ \tilde{\sigma}_{uv,i}(\beta^{\text{SIVE}}) &= 4e_i'(M_{W,Z}u \odot M_{W,Z}v) = u_i v_i - 2u_i \sum_{j \neq i} \bar{P}_{ij}v_j - 2v_i \sum_{j \neq i} \bar{P}_{ij}u_j + 4 \sum_{j \neq i} \bar{P}_{ij}^2 v_j u_j.\end{aligned}\quad (142)$$

We notice that the first term in the expressions coincides with the first terms in (96). The same holds up to a scaling constant for the second term. Those terms are therefore covered by the proof in the previous section. The only new terms are last terms of the respective expressions that yield a positive bias in the variance estimator. The bias contribution to the variance by observation i , if observation $j(i)$ is the only observation in the same covariate group with the same instrument status, is

$$\underbrace{[e_i'(Y - T\beta^{\text{SIVE}})]^2 u_{j(i)}^2}_{t_i} + [e_i'AT]^2 v_{j(i)}^2 + 2[e_i'A(Y - T\beta^{\text{SIVE}})e_i'AT]u_{j(i)}v_{j(i)} + R(|\hat{\beta}^{\text{SIVE}} - \beta^{\text{SIVE}}|). \quad (143)$$

We consider the first term and sum over all observations that only have one member of the same covariate group with the same instrument status. Call this set \mathcal{S} and note that $|\mathcal{S}| = \{2, 4, \dots, n/2\}$. We have

$$\begin{aligned}\sum_{i \in \mathcal{S}} t_i &= \sum_{i \in \mathcal{S}} \left[\sigma_{u,j(i)}^2 [M_W Z \zeta]_i^2 + 2[M_W Z \zeta]_i A_{ij(i)} \mathbb{E}[u_{j(i)}^2 v_{j(i)} | Q, X] + \sum_{k \neq i} \mathbb{E}[u_{j(i)}^2 v_k^2 | Q, X] A_{ik}^2 \right. \\ &\quad + (u_{j(i)}^2 - \sigma_{u,j(i)}^2) [M_W Z \zeta]_i^2 + 2[M_W Z \zeta]_i A_{ij(i)} (u_{j(i)}^2 v_{j(i)} - \mathbb{E}[u_{j(i)}^2 v_{j(i)} | Q, X]) \\ &\quad \left. + 2[M_W Z \zeta]_i u_{j(i)}^2 \sum_{k \neq i, j(i)} A_{ik} v_k + \sum_{k \neq i} (u_{j(i)}^2 v_k^2 - \mathbb{E}[u_{j(i)}^2 v_k^2 | Q, X]) A_{ik}^2 + u_{j(i)}^2 \sum_{k \neq i} \sum_{l \neq i, k} A_{ik} A_{il} v_k v_l \right].\end{aligned}\quad (144)$$

It is straightforward to show that $\mathbb{E}[r_n^{-1} \sum_{i \in \mathcal{S}} t_i | Q, X] \leq C < \infty$ a.s. and $\mathbb{V}(r_n^{-1} \sum_{i \in \mathcal{S}} t_i | Q, X) \rightarrow_{a.s.} 0$.

D.3.4 Variance bounds

For notational convenience, in this section we replace \tilde{A} by A and \tilde{P} by P . For $\mathbb{V}(B.2a.1 | Q, X)$ in (130), we need to bound terms of type

$$r_n^{-2} \sum_{i_1, \dots, i_{12}} |A_{i_1 i_2}| |A_{i_3 i_4}| |A_{i_5 i_6}| |A_{i_7 i_8}| \bar{P}_{i_9 i_{10}} \bar{P}_{i_{11} i_{12}}. \quad (145)$$

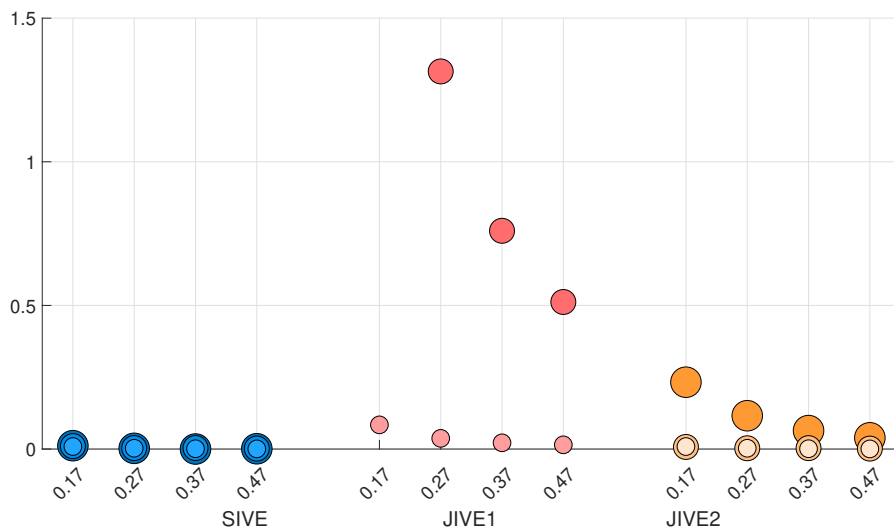
For $\mathbb{V}(B.2a.2 | Q, X)$ in (130), we need to bound terms of type

$$r_n^{-2} \sum_{i_1, \dots, i_{16}} |A_{i_1 i_2}| |A_{i_3 i_4}| |A_{i_5 i_6}| |A_{i_7 i_8}| \bar{P}_{i_9 i_{10}} \bar{P}_{i_{11} i_{12}} \bar{P}_{i_{13} i_{14}} \bar{P}_{i_{15} i_{16}}. \quad (146)$$

We now list the nonzero terms and bound them by $C r_n^{-2} G$ for some constant $C > 0$.

Term	Expression	Bound
1	$\sum_{i_1..i_3} A_{i_1 i_2}^3 A_{i_1 i_3} P_{i_1 i_2} P_{i_1 i_3}$	$\leq C \sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} P_{i_1 i_2} P_{i_1 i_3} \leq C l'(A \odot P)^2 l,$
2	$\sum_{i_1..i_3} A_{i_1 i_2}^2 A_{i_1 i_3}^2 P_{i_1 i_2} P_{i_1 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} P_{i_1 i_2} P_{i_1 i_3} \leq C l'(A \odot P)^2 l,$
3	$\sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3}^3 P_{i_1 i_2}^2$	$\leq C \sum_{i_1..i_3} A_{i_1 i_3}^2 P_{i_1 i_2}^2 \leq C \text{tr}(A^2),$
4	$\sum_{i_1..i_4} A_{i_1 i_2} A_{i_1 i_3} A_{i_1 i_4}^2 P_{i_1 i_2} P_{i_1 i_3}$	$\leq C \sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} e'_i A^2 e_i P_{i_1 i_2} P_{i_1 i_3} \leq C l'(A \odot P)^2 l,$
5	$\sum_{i_1..i_4} A_{i_1 i_3}^2 A_{i_1 i_4}^2 P_{i_1 i_2}^2$	$\leq C \sum_{i_1, i_3, i_4} A_{i_1 i_3}^2 A_{i_1 i_4}^2 \leq C l'(A \odot A)^2 l,$
6	$\sum_{i_1..i_3} A_{i_1 i_2}^3 A_{i_1 i_3} P_{i_1 i_2} P_{i_2 i_3}$	$\leq C \sum_{i_1..i_3} A_{i_1 i_3} P_{i_1 i_2} P_{i_2 i_3} \leq C l'(A \odot P) l,$
7	$\sum_{i_1..i_3} A_{i_1 i_2}^2 A_{i_1 i_3}^2 P_{i_1 i_2} P_{i_2 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} P_{i_1 i_2} P_{i_2 i_3} \leq C l'(A \odot P) l,$
8	$\sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3}^3 P_{i_1 i_2} P_{i_2 i_3}$	$\leq C \sum_{i_1..i_3} A_{i_1 i_3} P_{i_1 i_2} P_{i_2 i_3} \leq C l'(A \odot P) l,$
9	$\sum_{i_1..i_3} A_{i_1 i_3}^4 P_{i_1 i_2} P_{i_2 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} P_{i_1 i_2} P_{i_2 i_3} \leq C l'(A \odot P) l,$
10	$\sum_{i_1..i_3} A_{i_1 i_2}^2 A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2}^2$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} \leq C l'(A^2 \odot P) l.$
11	$\sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3}^2 A_{i_2 i_3} P_{i_1 i_2}^2$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} \leq C l'(A^2 \odot P) l,$
12	$\sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3}^2 A_{i_2 i_3} P_{i_1 i_2} P_{i_1 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} \leq C l'(A^2 \odot P) l,$
13	$\sum_{i_1..i_3} A_{i_1 i_3}^3 A_{i_2 i_3} P_{i_1 i_2} P_{i_1 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} \leq C l'(A^2 \odot P) l,$
14	$\sum_{i_1..i_3} A_{i_1 i_2}^2 A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} P_{i_2 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_2} A_{i_2 i_3} P_{i_1 i_2} P_{i_2 i_3} \leq C l'(A \odot P)^2 l,$
15	$\sum_{i_1..i_4} A_{i_1 i_2} A_{i_1 i_3} A_{i_1 i_4} A_{i_2 i_4} P_{i_1 i_2} P_{i_2 i_3}$	$= \sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} e'_i (A^2 \odot P) e_i P_{i_2 i_3} \leq C l'(A^2 \odot P) l.$
16	$\sum_{i_1..i_4} A_{i_1 i_2} A_{i_1 i_4}^2 A_{i_2 i_3} P_{i_1 i_2} P_{i_2 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_2} e'_i A^2 e_i A_{i_2 i_3} P_{i_1 i_2} P_{i_2 i_3} \leq C l'(A \odot P)^2 l$
17	$\sum_{i_1..i_3} A_{i_1 i_3}^3 A_{i_2 i_3} P_{i_1 i_2} P_{i_2 i_3}$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} \leq C l'(A^2 \odot P) l.$
18	$\sum_{i_1..i_4} A_{i_1 i_3}^2 A_{i_1 i_4} A_{i_2 i_3} P_{i_1 i_2} P_{i_3 i_4}$	$\leq \sum_{i_1..i_3} A_{i_1 i_3} A_{i_2 i_3} P_{i_1 i_2} e'_i (AP) e_i \leq C l'(A^2 \odot P) l.$
19	$\sum_{i_1..i_4} A_{i_1 i_3}^2 A_{i_1 i_4} A_{i_3 i_4} P_{i_1 i_2} P_{i_2 i_3}$	$\leq \sum_{i_1, i_3, i_4} A_{i_1 i_4} A_{i_3 i_4} P_{i_1 i_2} \leq C l'(A^2 \odot P) l.$
20	$\sum_{i_1..i_4} A_{i_1 i_2} A_{i_1 i_3}^2 A_{i_3 i_4} P_{i_1 i_2} P_{i_3 i_4}$	$= l'(A \odot P) A(A \odot P) l \leq C l'(A \odot P)^2 l,$
21	$\sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} A_{i_2 i_3}^2 P_{i_1 i_2}^2$	$\leq \sum_{i_1..i_3} A_{i_2 i_3}^2 P_{i_1 i_2}^2 \leq C \text{tr}(A^2),$
22	$\sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} A_{i_2 i_3}^2 P_{i_1 i_2} P_{i_1 i_3}$	$\leq \sum_{i_1..i_3} A_{i_2 i_3}^2 P_{i_1 i_2} P_{i_1 i_3} \leq C \text{tr}(A^2)$
23	$\sum_{i_1..i_3} A_{i_1 i_3}^2 A_{i_2 i_3}^2 P_{i_1 i_2} P_{i_1 i_3}$	$\leq C \sum_{i_1..i_3} A_{i_2 i_3}^2 P_{i_1 i_2} P_{i_1 i_3} \leq C \text{tr}(A^2),$
24	$\sum_{i_1..i_4} A_{i_1 i_3}^2 A_{i_2 i_4}^2 P_{i_1 i_2}^2$	$= \sum_{i_1, i_2} e'_i A^2 e_i e'_i A^2 e_i P_{i_1 i_2}^2 \leq C \text{tr}(A^2),$
25	$\sum_{i_1..i_4} A_{i_1 i_3} A_{i_1 i_4} A_{i_2 i_3} A_{i_2 i_4} P_{i_1 i_2}^2$	$\leq \sum_{i_1, i_2} (e'_i A^2 e_i)^2 P_{i_1 i_2}^2 \leq C l'(A^2 \odot P) l.$
26	$\sum_{i_1..i_4} A_{i_1 i_3} A_{i_1 i_4} A_{i_2 i_3} A_{i_3 i_4} P_{i_1 i_2} P_{i_1 i_3}$	$\leq \sum_{i_1..i_3} A_{i_2 i_3} P_{i_1 i_2} P_{i_1 i_3} \leq C l'(A \odot P) l,$
27	$\sum_{i_1..i_4} A_{i_1 i_2} A_{i_1 i_3} A_{i_3 i_4}^2 P_{i_1 i_2} P_{i_1 i_3}$	$\leq C \sum_{i_1..i_3} A_{i_1 i_2} A_{i_1 i_3} P_{i_1 i_2} P_{i_1 i_3} \leq C l'(A \odot P)^2 l.$

Figure 5: Average absolute bias in the estimand



Note: the figure shows the absolute median difference with the causal estimand in a setting without treatment heterogeneity. The size of the circles indicates the number of covariate groups with the small circle corresponding to $L = 1$, the medium circle corresponding to $L = 25$ and the large circle corresponding to $L = 300$. The x -axis is the instrument strength $p(1) - p(0)$, with $p(0) = 0.22$ and $p(1) = \{0.39, 0.49, 0.59, 0.69\}$. Because JIVE1 shows large biases for $L = 25$ and $L = 300$, the y -axis is limited to $[0, 0.9]$.