

# Many Instruments and Judges

---

Michal Kolesár

ECO539B, Fall 2022

April 1, 2025

- Consider linear instrumental variables (IV) model when number of instruments  $K$  large.

Arises in 2 ways:

1. Interact low-dim instrument with controls (Angrist and Krueger 1991; Belloni et al. 2012).
2. Use fixed effects (group indicators) as instruments: “judge” or “examiner” designs. Studies of effects of incarceration on economic outcomes (Kling 2006; Aizer and Doyle 2015), studies that exploit random assignment of judges to bankruptcy cases, criminal cases, or patent cases, doctors to shifts...

### Key issues

1. two-stage least squares (TSLS) biased with many instruments. Alternatives?
2. standard errors
3. Interpretation of first-stage  $F$  statistic

Setup and Estimation

Inference

First stage  $F$

Summary and illustration

- Same notation as in previous lecture, with reduced form and first stage

$$Y_i = Z_i' \delta + W_i' \psi_Y + u_{Yi}, \quad (1)$$

$$D_i = Z_i' \pi + W_i' \psi_D + u_{Di}. \quad (2)$$

The parameter of interest is  $\beta := \frac{\delta' Q \pi}{\pi' Q \pi}$ , where  $Q = E[\tilde{Z}_i \tilde{Z}_i']$ , and  $\tilde{Z}_i = Z_i - E[Z_i W_i'] E[W_i W_i']^{-1} W_i$ .

- To measure overall strength of instruments, define measure of “effective sample size”

$$r_n = n \pi' Q \pi = n E[(\tilde{Z}_i' \pi)^2].$$

- $X_i = (Z_i', W_i')'$  collects the right-hand side (RHS) (or “exogenous”) variables.

- For any matrix  $A$ , let  $H_A = A(A'A)^{-1}A'$  denote hat matrix (also called a projection matrix), and let  $\ddot{A} = A - H_W A$  denote residuals after projecting  $A$  onto the covariates.
  - $\tilde{Z}_i$  is population analog of  $\ddot{Z}_i$

## Simple judge design

- Individual  $i$  assigned judge  $Q_i$ , random conditional on  $i$ 's geographic location  $G_i$ .
- Covariates and instruments both indicators:  $W_{i\ell} = \mathbb{1}\{G_i = \ell\}$ , and  $Z_{ik} = \mathbb{1}\{Q_i = k\}$ .
- There are  $L$  locations and  $K + L$  judges; in each location we drop one judge to avoid collinearity—call this judge reference judge.
- $\psi_{D,\ell}$  is average sentencing rate of the reference judge in location  $\ell$ , and  $\pi_k$  is sentencing rate of judge  $k$  relative to reference judge in same location.

# Monotonicity i

- Monotonicity assumption requires that judges **agree on ranking** of defendants, they just disagree on **cutoff** at which they start sentencing them.
- At odds with economic theory as well as statistical evidence:
  1. Chan, Gentzkow, and Yu (2022) point out that decisions of physicians differ due to both skill and preferences.
  2. Mueller-Smith (2015) argues that a judge may not use a common cutoff for each defendant: for instance, the cutoff may vary by crime type.
  3. Kleinberg et al. (2018) find that the increase in crime associated with judges who are more likely to release defendants on bail is about the same as if these more lenient judges randomly picked the extra defendants to release on bail
  4. Statistical tests of monotonicity reject in empirical applications (e.g. Frandsen, Lefgren, and Leslie 2023; Agan, Doleac, and Harvey 2023; Coulibaly et al. 2024)

- 5. Direct tests of monotonicity using judge panels also reject (Sigstad 2025)
- Failure of monotonicity only affects the interpretation of  $\beta$  if we want to allow for treatment effect heterogeneity; to keep statistical issues separate from identification issues, we put these issues aside here.
  - Evidence from Sigstad (2025) suggests problem not first order, but that's probably TBD

Consider asymptotics in which  $K = \dim(Z_i)$  and  $L = \dim(W_i)$  can to grow with sample size

### Lemma

The TSLS estimator suffers from own observation bias towards ordinary least squares (OLS). In particular, suppose  $E[u_i | X_i] = 0$ , with  $\Omega(X_i) = E[u_i u_i' | X_i]$  bounded,  $E[\tilde{Z}_i | W_i] = 0$ , and  $L/n \rightarrow 0$ .

Then consistency of TSLS is in general requires  $K/r_n \rightarrow 0$ . Under homoskedastic errors,

$$\hat{\beta}_{\text{TSLS}} = \beta + \frac{(\Omega_{YD} - \Omega_{DD}\beta)K}{r_n + \Omega_{DD}K} + o_p(1) = (1 - w)\beta + w\beta_{\text{OLS}} + o_p(1), \quad w = \frac{K}{r_n/\Omega_{DD} + K},$$

where  $\beta_{\text{OLS}} = (\Omega_{YD} - \Omega_{DD}\beta)/\Omega_{DD} + \beta$  is the probability limit of OLS.

TSLS bias scales with  $K/r_n$ , rather than  $K/n$ !



- Bias arises because single constructed instrument  $\hat{Z}_{\text{TSLs},i} = Z_i' \hat{\pi} + W_i' \hat{\psi}_D$  used by TSLS puts positive weight on own treatment status  $D_i$ .
- Total net weight, holding overall instrument strength constant, scales linearly with  $K$ , so TSLS bias scales with number of instruments.

### Simple judge design

Without covariates,  $\hat{\beta}_{\text{TSLs}} = \sum_i Y_i \hat{Z}_{\text{TSLs},i} / \sum_i Y_i \hat{Z}_{\text{TSLs},i}$ , and  $\hat{Z}_{\text{TSLs},i}$  is average sentencing rate of judge  $i$  is assigned to—including  $i$ 's own treatment. If # cases per judge same, puts weight  $K/n$  on  $D_i$ .

- Solution 1 (Bekker [1994](#)): use limited information maximum likelihood (LIML), or variants thereof. Only need  $\sqrt{K}/r_n \rightarrow 0$  for consistency, substantially weaker requirement. But LIML-like estimators not robust to heterogeneous treatment effects.
- Solution 2: subtract estimate of bias. Easy to do under homoskedastic errors: leads to variants of the bias-corrected TSLS estimator that dates back to Nagar ([1959](#)). But consistency does depend on homoskedasticity.

- Bias caused by using  $D_i$  in constructing the single instrument  $\hat{Z}_{\text{TSLs},i}$ : obvious solution is to not use it!
- Without covariates, clear we want to use  $\hat{Z}_{\text{JIVE1},i} = Z_i' \hat{\pi}_{-i}$ , where  $\hat{\pi}_{-i}$  based on regression of  $D$  onto  $Z$  with  $i$  **excluded**. How to incorporate covariates?
- Option 1: construct predictors  $\hat{\pi}_{-i}$  and  $\hat{\psi}_{D,-i}$  of  $\pi$  and  $\psi_D$  based on a regression of  $D$  onto  $(Z, W)$ , but with the  $i$ th observation removed, to construct a single instrument
$$\hat{Z}_{\text{JIVE1},i} = Z_i' \hat{\pi}_{-i} + W_i' \hat{\psi}_{D,-i}$$
  - What is this in judges design?

- Then run an IV regression of  $Y_i$  onto  $D_i$  and  $W_i$ , using  $\hat{Z}_{\text{JIVE1},i}$  as an instrument for  $D_i$ .  
Resulting estimator known as jackknife instrumental variables estimator (JIVE1):

$$\hat{\beta}_{\text{JIVE1}} = \frac{\hat{Z}'_{\text{JIVE1}} \ddot{Y}}{\hat{Z}'_{\text{JIVE1}} \ddot{D}}, \quad \hat{Z}_{\text{JIVE1}} = \ddot{H}D, \quad \ddot{H} = (I - \text{diag}(H_X))^{-1}(H_X - \text{diag}(H_X)),$$

- Don't need to run  $n$  first-stage regressions
- Problem: when we project out the covariates from  $\hat{Z}_i$ , we re-introduce the own observation bias
  - In judges design, adjust instrument by average sentencing rate in the location of  $i$

- Generally need  $L/r_n \rightarrow 0$  for consistency (see Evdokimov and Kolesár (2018)). Under homo, bias goes in opposite direction to TSLS bias:

$$\hat{\beta}_{\text{JIVE}_1} = \beta - \frac{(\Omega_{YD} - \Omega_{DD}\beta)L}{r_n - \Omega_{DD}L} + o_p(1) = (1 + \lambda)\beta - \lambda\beta_{OLS} + o_p(1), \quad \lambda = \frac{L}{r_n/\Omega_{DD} - L}.$$

- Leads to option 2: *first* partial out the covariates, and then do the leave-one-out prediction. Called improved improved jackknife instrumental variables estimator (IJIVE<sub>1</sub>) estimator, proposed in Akerberg and Devereux (2009).

$$\hat{\beta}_{\text{IJIVE}_1} = \frac{\ddot{D}'\ddot{H}'\ddot{Y}}{\ddot{D}'\ddot{H}'\ddot{D}}, \quad \ddot{H} = (I - \text{diag}(H_{\ddot{Z}}))^{-1}(H_{\ddot{Z}} - \text{diag}(H_{\ddot{Z}})).$$

- Option 3: Problem is to estimate

$$\beta = \frac{\text{cov}(Y_i, \tilde{Z}_i' \pi)}{\text{cov}(Y_i, \tilde{Z}_i' \pi)} = \frac{\sum_i E[Y_i \tilde{Z}_i' \pi]}{\sum_i E[D_i \tilde{Z}_i' \pi]}.$$

Just use sample analog with  $\hat{\pi}_{-i}$  instead of  $\pi$ ! Kolesár (2013) calls this estimator unbiased jackknife instrumental variables estimator (UJIVE)

- (Actually, version in the paper slightly inferior, wait for revision)

That is, (i) run jackknife the regression of  $D$  onto  $(W, Z)$  to compute  $\hat{\pi}_{-i}$ . (ii) covariate-adjust  $Z$  to get  $\ddot{Z}$ , and then use  $\hat{D}_{\text{UJIVE},i} = \ddot{Z}_i \hat{\pi}_{-i}$  as single instrument,  $\hat{\beta}_{\text{UJIVE}} = \hat{D}'_{\text{UJIVE}} Y / \hat{D}'_{\text{UJIVE}} D$ .

- Both UJIVE doesn't restrict number of covariates while IJIVE1 consistent so long as  $LK/r_n n \rightarrow 0$

- Chao, Swanson, and Woutersen (2023) offer alternative to UJIVE based on Hadamard products (analog to Hadamard variance estimator for high-dimensional OLS), but implementing requires inverting  $n \times n$  matrices...

Setup and Estimation

Inference

First stage  $F$

Summary and illustration



- Define  $\pi_\Delta = \delta - \pi\beta$ ,  $u_{\Delta,i} = u_{Yi} - u_{Di}\beta$ ,  
 $\gamma = E[W_i W_i']^{-1} E[W_i (Y_i - D_i \beta)] = E[W_i W_i']^{-1} E[W_i Z_i] \pi_\Delta + \psi_Y - \psi_D \beta$ , and  $\epsilon_i = \tilde{Z}_i' \pi_\Delta + u_{\Delta,i}$ .
- Write the “structural” equation as

$$Y_i = D_i \beta + W_i' \gamma + \epsilon_i.$$

- saw in the previous lecture that the oracle estimator satisfied

$$\mathcal{V}_{1,n}^{-1/2} (\hat{\beta}^* - \beta) \Rightarrow \mathcal{N}(0, 1), \quad \mathcal{V}_{1,n} = \frac{E[\epsilon_i^2 (\tilde{Z}_i' \pi)^2]}{r_n E[(\tilde{Z}_i' \pi)^2]}.$$

Variance  $\mathcal{V}_{1,n}$  is estimated by the conventional robust standard errors, such as those in Stata. Also correct standard errors for TSLS, UJIVE, or IJIVE<sub>1</sub> if (i) there is no treatment effect heterogeneity, and (ii)  $K/r_n \rightarrow 0$  for UJIVE, or IJIVE<sub>1</sub>, and  $K^2/r_n \rightarrow 0$  for TSLS.

- If (i) fails, then, as discussed last time, don't achieve oracle variance, but instead the correct asymptotic variance is given by

$$\mathcal{V}_{2,n} = \frac{E[((\tilde{Z}'_i \pi_\Delta)u_{D,i} + \epsilon_i(\tilde{Z}'_i \pi))^2]}{r_n E[(\tilde{Z}'_i \pi)^2]}.$$

- What if (ii) fails? What if  $K/r_n \rightarrow 0$ , but  $K^2/r_n \not\rightarrow 0$  and we use TSLS?

## More robust variance ii

- If  $K/r_n \not\rightarrow 0$ , and use UJIVE or IJIVE<sub>1</sub>, we need to account for the presence of many instruments in the asymptotic variance formula (Evdokimov and Kolesár 2018, Theorem 5.4)

$$(\mathcal{V}_{2,n} + \mathcal{V}_{MI,n})^{-1/2}(\hat{\beta}_{\text{UJIVE}} - \beta) \Rightarrow \mathcal{N}(0, 1),$$

$$\mathcal{V}_{MI,n} = \frac{1}{r_n^2} \sum_{i \neq j} (H_{\tilde{Z},ij}^2 u_{\Delta,i}^2 u_{2,j}^2 + H_{\tilde{Z},ij}^2 u_{\Delta,i} u_{2,i} \cdot u_{\Delta,j} u_{2,j}), \quad (3)$$

- Under homoskedasticity,  $\Omega(X_i) = \Omega = E[u_i u_i']$ , additional many instrument term simplifies to

$$\mathcal{V}_{MI,n} = \frac{K}{r_n^2} (E[(u_{i1} - u_{2i}\beta)^2] \cdot E[u_{2i}^2] + E[(u_{i1} - u_{2i}\beta)u_{2i}]^2)(1 + o_p(1)).$$

Setup and Estimation

Inference

First stage  $F$

Summary and illustration

- Some papers calculate the instrument  $\hat{Z}_i$  manually, often as a leave-one-out prediction, effectively computing JIVE<sub>1</sub> by hand. But this does not mean that  $K = 1$ ! That will overstate actual instrument strength.
- When using the (correctly computed) first-stage  $F$  for diagnostics, remember that the  $F > 10$  rule of thumb tests hypothesis that TSLS bias, relative to the bias of OLS exceeds 0.1.
- But small  $F$  statistic not necessarily a concern when UJIVE or IJIVE<sub>1</sub> are used. Under homoskedasticity,

$$E[F] = \frac{E[\hat{\pi}_2 \tilde{Z}' \tilde{Z} \hat{\pi}_2]}{KE[u_{2i}^2]} = \frac{\pi_2 E[\tilde{Z}' \tilde{Z}] \pi_2}{KE[u_{2i}^2]} + 1 \approx \frac{r_n}{KE[u_{2i}^2]} + 1.$$

- If  $r_n/K \rightarrow 0$ , TSLS will be inconsistent; but jackknife estimators will remain consistent so long as  $r_n/\sqrt{K} \rightarrow \infty$

Setup and Estimation

Inference

First stage  $F$

Summary and illustration

- If  $K$  non-negligible relative to effective sample size  $r_n$ , TSLS biased. Instead, use IJIVE<sub>1</sub> or UJIVE.
  - $K/r_n$  may be large even if  $K/n$  small!
  - JIVE<sub>1</sub> not a good solution
- To ensure reliable inference, standard errors need to account for additional many instrument term in the asymptotic variance
  - Even more important is to avoid downward bias that's present in the default standard errors estimator based on TSLS—see notes.
- Will now illustrate in application to Angrist and Krueger (1991)

Estimator	Estimate	$\hat{\mathcal{V}}_1^{1/2}$	$\hat{\mathcal{V}}_2^{1/2}$	$\sqrt{\hat{\mathcal{V}}_2 + \hat{\mathcal{V}}_{MI}}$	$\hat{r}_n/K$
Panel A: OLS					
OLS	0.0670	0.0004			
Panel B: Instrument is QOB. $F = 34.0$					
TSLS	0.1026	0.0195	0.0198		366.0
JIVE1	0.1039	0.0203	0.0206	0.0209	351.6
UJIVE	0.1036	0.0201	0.0204	0.0207	355.2
Panel C: Instrument is $QOB \times YOB$ . $F = 4.9$					
TSLS	0.0891	0.0162	0.0176		52.6
JIVE1	0.0959	0.0224	0.0244	0.0273	38.3
UJIVE	0.0938	0.0204	0.0222	0.0211	41.9
Panel D: Instrument is $QOB \times YOB + QOB \times SOB$ , $F = 2.6$					
TSLS	0.0928	0.0097	0.0112		26.2
JIVE1	0.1211	0.0205	0.0243	0.0273	12.7
UJIVE	0.1096	0.0160	0.0187	0.0211	16.1
Panel E: Instrument is $QOB \times YOB \times SOB$ , $F = 1.1$					
TSLS	0.0721	0.0049	0.0067		11.6
JIVE1	0.0320	0.0307	0.0425	0.0515	-1.9
UJIVE	0.1110	0.0397	0.0548	0.0663	1.4



# References i

- Ackerberg, Daniel A., and Paul J. Devereux. 2009. "Improved JIVE Estimators for Overidentified Linear Models with and without Heteroskedasticity." *Review of Economics and Statistics* 91, no. 2 (May): 351–362. <https://doi.org/10.1162/rest.91.2.351>.
- Agan, Amanda, Jennifer L Doleac, and Anna Harvey. 2023. "Misdemeanor Prosecution." *The Quarterly Journal of Economics* 138, no. 3 (June): 1453–1505. <https://doi.org/10.1093/qje/qjad005>.
- Aizer, Anna, and Joseph J. Doyle Jr. 2015. "Juvenile Incarceration, Human Capital, and Future Crime: Evidence from Randomly Assigned Judges." *The Quarterly Journal of Economics* 130, no. 2 (May): 759–803. <https://doi.org/10.1093/qje/qjv003>.
- Angrist, Joshua D., and Alan B. Krueger. 1991. "Does Compulsory School Attendance Affect Schooling and Earnings?" *The Quarterly Journal of Economics* 106, no. 4 (November): 979–1014. <https://doi.org/10.2307/2937954>.
- Bekker, Paul A. 1994. "Alternative Approximations to the Distributions of Instrumental Variable Estimators." *Econometrica* 62, no. 3 (May): 657–681. <https://doi.org/10.2307/2951662>.
- Belloni, Alexandre, Daniel L. Chen, Victor Chernozhukov, and Christian B. Hansen. 2012. "Sparse Models and Methods for Optimal Instruments with an Application to Eminent Domain." *Econometrica* 80, no. 6 (November): 2369–2429. <https://doi.org/10.3982/ECTA9626>.
- Chan, David C., Matthew Gentzkow, and Chuan Yu. 2022. "Selection with Variation in Diagnostic Skill: Evidence from Radiologists." *The Quarterly Journal of Economics* 137, no. 2 (May): 729–783. <https://doi.org/10.1093/qje/qjab048>.

## References ii

- Chao, John C., Norman R. Swanson, and Tiemen Woutersen. 2023. “Jackknife Estimation of a Cluster-Sample IV Regression Model with Many Weak Instruments.” *Journal of Econometrics* 235, no. 2 (August): 1747–1769. <https://doi.org/10.1016/j.jeconom.2022.12.011>.
- Coulibaly, Mohamed, Yu-Chin Hsu, Ismael Mourifié, and Yuanyuan Wan. 2024. “A Sharp Test for the Judge Leniency Design,” May. arXiv: [2405.06156](https://arxiv.org/abs/2405.06156).
- Evdokimov, Kirill, and Michal Kolesár. 2018. “Inference in Instrumental Variable Regression Analysis with Heterogeneous Treatment Effects.” January. [https://www.princeton.edu/~mkolesar/papers/het\\_iv.pdf](https://www.princeton.edu/~mkolesar/papers/het_iv.pdf).
- Frandsen, Brigham, Lars Lefgren, and Emily Leslie. 2023. “Judging Judge Fixed Effects.” *American Economic Review* 113, no. 1 (January): 253–277. <https://doi.org/10.1257/aer.20201860>.
- Kleinberg, Jon, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. 2018. “Human Decisions and Machine Predictions.” *The Quarterly Journal of Economics* 133, no. 1 (February): 237–293. <https://doi.org/10.1093/qje/qjx032>.
- Kling, Jeffrey R. 2006. “Incarceration Length, Employment, and Earnings.” *American Economic Review* 96, no. 3 (May): 863–876. <https://doi.org/10.1257/aer.96.3.863>.
- Kolesár, Michal. 2013. “Estimation in an Instrumental Variables Model With Treatment Effect Heterogeneity.” Working paper, Princeton University, November. [https://www.princeton.edu/~mkolesar/papers/late\\_estimation.pdf](https://www.princeton.edu/~mkolesar/papers/late_estimation.pdf).

- Mueller-Smith, Michael. 2015. "The Criminal and Labor Market Impacts of Incarceration." Working paper, University of Michigan, August.
- Nagar, Anirudh Lal. 1959. "The Bias and Moment Matrix of the General  $k$ -Class Estimators of the Parameters in Simultaneous Equations." *Econometrica* 27, no. 4 (October): 575–595. <https://doi.org/10.2307/1909352>.
- Sigstad, Henrik. 2025. "Monotonicity among Judges: Evidence from Judicial Panels and Consequences for Judge IV Designs." March. <https://doi.org/10.2139/ssrn.4534809>.