

Robust Inference for the Frisch Labor Supply Elasticity

By MICHAEL KEANE AND TIMOTHY NEAL*

January 18, 2022

There is a long standing controversy over the magnitude of the Frisch labor supply elasticity. Macro economists using DSGE models often calibrate it to be large, while many micro data studies find it is small. Several papers attempt to reconcile the micro and macro results. We offer a new and simple explanation: Most micro studies estimate the Frisch using a 2SLS regression of hours changes on wage changes. However, due to a little appreciated power asymmetry property of 2SLS that we clarify, estimates of the Frisch will (spuriously) appear more precise when they are more shifted in the direction of the OLS bias, which is negative. As a result, Frisch elasticity estimates near zero appear (spuriously) precise, while large positive estimates appear (spuriously) imprecise. This pattern makes it difficult for a 2SLS t -test to detect a true positive Frisch elasticity. Fortunately, the Anderson-Rubin (AR) test does not suffer from this power asymmetry problem. The AR test leads us to conclude the Frisch elasticity is large and significant in the NLSY97 data. In contrast, a conventional 2SLS t -test would lead us to conclude it is not significantly different from zero. Our application illustrates a fundamental problem with 2SLS t -tests that arises quite generally. This problem is severe when instruments are weak, but persists even if they are strong. Thus, we argue the AR test should be widely adopted in lieu of the t -test.

Keywords: *Frisch elasticity, labor supply, weak instruments, 2SLS, Anderson-Rubin test*

JEL: *J22, D15, C12, C26*

I. Introduction

The elasticity of labor supply with respect to predictable wage changes – known as the Frisch elasticity – plays a key role in many economic policy debates. The Frisch plays a central role because predictable wage changes have pure substitution effects. As an example of its importance, Conesa, Kitao and Krueger (2009) argue that higher values of the Frisch imply a higher optimal tax rate on capital income. And macro models where real shocks play a key role in business cycles often require the Frisch to be large to match observed fluctuations in work hours over the cycle, see Prescott (2006). Because of its importance, a large literature attempts to estimate the Frisch elasticity, as exemplified by classic papers by MaCurdy (1981) and Altonji (1986) and surveyed in Keane (2011, 2021).

* CEPAR & School of Economics, UNSW. Corresponding author: m.keane@unsw.edu.au

Classic micro data studies in the style of MaCurdy (1981) typically find the Frisch elasticity is small, while macro economists using DSGE models often calibrate it to be large. This led to a long-standing “macro-micro controversy” over the magnitude of the Frisch. Keane and Rogerson (2012, 2015) discuss attempts to resolve the controversy. Here we present a new type of resolution based on a critique of the micro-econometrics itself: We argue the classic studies were inherently biased against finding the Frisch is both large and significant, due to a little appreciated power asymmetry property of 2SLS t -tests that we clarify.

The idea behind most micro studies is as follows: Given panel data on workers, one may run an OLS regression of changes in log hours on changes in log wages. But this fails to deliver the elasticity of hours with respect to *predictable* wage changes – as some changes are surprises. Instead, the approach pioneered by MaCurdy (1981) involves running a 2SLS regression where the instrument for the change in log wages has two properties: First, it predicts wage growth at the individual level. Second, it is known at the start of the time period over which changes in wages are calculated (so it is uncorrelated with surprise wage changes). Then, fitted values from the first-stage of 2SLS give us predictable changes in wages, and the second stage delivers an estimate of the elasticity of hours with respect to these predictable wage changes, which is the Frisch concept.

A little appreciated property of 2SLS is that it generates a strong association between the estimate and the standard error of regression, which is minimized when the estimate is close to $E(\hat{\beta}_{OLS})$, see Phillips (1989). As we show in Keane and Neal (2021), this generates a strong association between 2SLS estimates and their standard errors, which is positive if the OLS bias is negative. Hence, positive Frisch elasticity estimates have artificially inflated standard errors. As a result, a 2SLS t -test has little power to detect a true positive Frisch elasticity.

We further show that the Anderson-Rubin (AR) test does not suffer from the power asymmetry that afflicts the t -test. Hence it is a far more reliable guide to inference in 2SLS applications. Using NLSY97 data, we obtain a Frisch elasticity of 0.60, which is quite large for young men.¹ The AR test indicates this estimate is highly significant ($p=.018$), while a conventional 2SLS t -test indicates it is not significantly different from zero. Thus, application of a superior inferential procedure – the AR test – reveals clear evidence to support a large Frisch elasticity.

It is well-known that the literature on estimating the Frisch elasticity using 2SLS has been hampered by weak instrument problems, as it is hard to find instruments that are strong predictors of wage growth. This issue was explored by Lee (2001), who shows how weak instrument problems in classic papers like MaCurdy (1981) and Altonji (1986) biased their estimates of the Frisch elasticity towards zero (i.e., towards OLS). Those authors used over-identified models, where the instruments used to predict wage growth were primarily age and schooling (including linear, quadratic and interaction terms). As Lee (2001) shows, given PSID samples of

¹There is growing evidence that the Frisch increases substantially with age - see Keane (2021). For example, for men in their 60s French (2005) obtains 1.3, while Imai and Keane (2004) obtain over 2.

the size they had available, and the type of instruments they used, first-stage F -statistics in their models would have been no higher than one. This is far below conventional weak instrument testing thresholds, such as the $F > 10$ rule advocated by Staiger and Stock (1997). Thus, the instruments were very weak in the classic studies. It is well-known that with multiple weak instruments 2SLS is seriously biased towards OLS – see Bound, Jaeger and Baker (1995).

Using more recent data, Lee (2001) constructs a PSID sample five times larger than in the classic studies, and obtains a first-stage F of 26.3 and a Frisch estimate of .503 (se=.092) for 25-60 year old men.² In contrast, in a Monte Carlo exercise where he uses randomly drawn 1/5th sub-samples, he obtains an average Frisch estimate of only .253 (se=.227).³ This is clear evidence of downward bias in the classic studies. This bias towards OLS when instruments are weak is widely appreciated in the IV literature, although Lee (2001)'s result appears to have done little to alter the conventional wisdom that the Frisch elasticity is small.

We emphasize a different issue – completely unrelated to bias – that has been largely overlooked in the prior literature. In contrast to Lee (2001), we focus on the just identified case where instrument strength exceeds conventional weak IV thresholds, so 2SLS is approximately median unbiased (as we explain in the Appendix), so bias towards OLS is not a concern. We demonstrate how the 2SLS t -test exhibits very poor behavior in that supposedly benign context.

Specifically, we show that the association between 2SLS estimates and their standard errors that we document is a very serious problem for inference using 2SLS t -tests even when instruments are strong by conventional standards (such as the $F > 10$ criterion). It is an inherent property of 2SLS that estimates of the Frisch will (spuriously) appear more precise when they are more shifted in the direction of the OLS bias, which is negative. A consequence is that the 2SLS t -test has little power to detect a true positive Frisch elasticity. This will bias studies that rely on 2SLS t -tests against concluding the Frisch is large.

The implications of our results go well beyond the present application: Theorists have often advocated using the AR test when instruments are weak, because it is robust to weak instrument problems. But the association between 2SLS estimates and standard errors renders 2SLS standard errors and t -tests highly misleading even when instruments are well above conventional weak instrument thresholds. Hence, we argue the AR test should replace the t -test in 2SLS applications, not only when instruments are weak but even when they are strong.

The outline of the paper is as follows: Section II discusses our NLSY79 data and estimating equation. Section III presents our 2SLS estimates and t -test results. Section IV presents AR test results. Section V compares the behavior of the AR and t -tests, and Section VI interprets the evidence in light of that comparison. Section VII presents results with multiple instruments. Section VIII concludes.

²Substantively, our estimate of 0.6 is large relative to Lee's estimate of 0.5, as our sample is much younger and there is growing evidence that the Frisch increases substantially with age - see Keane (2021).

³Similarly, if he limits his analysis to the 1/5 of the PSID data available to the classic studies, he obtains a first stage F of only 1.06 and a Frisch estimate of .258 which is insignificant (se=.172).

II. Estimating the Frisch Labor Supply Elasticity

We estimate the Frisch elasticity using data from the National Longitudinal Survey of Youth 1997 (NLSY97). The NLSY97 follows a sample of American youth born in 1980-84. The 8,984 respondents were aged 12-17 when first interviewed in 1997.⁴ We use data from rounds 11 through 15, which contain information on labor income and work hours in 2005 to 2010. The regression we run is:

$$(1) \quad \Delta \ln H_{it} = \alpha + \beta \Delta \ln W_{it} + \gamma \mathbf{C}_{it} + \epsilon_{it}$$

where H_{it} is annual hours worked for respondent i in year t , W_{it} is the wage, and \mathbf{C}_{it} is a vector of control variables which includes year dummies (to capture business cycle effects on hours worked) as well as respondent age and race/ethnicity.

Our hours measure is “Total annual hours worked at all civilian jobs during the year in question” while our income measure is “Annual income from wages, salary, commissions, and tips before tax deductions.” We obtain an annual wage measure by taking the ratio of annual income to annual hours. Regressions that involve percentage changes can be quite sensitive to measurement error and outliers, as these can generate extreme percentage changes. So, as is typical in this literature, we implement a number of sample screens designed to eliminate outliers.⁵

Obviously, OLS estimation of (1) fails to identify the Frisch elasticity, as predictable and unpredictable wage changes have different effects on labor supply. A surprise wage increase has both substitution and income effects. In contrast, a predictable wage increase has no income effect (precisely because it was predictable), so it induces a pure substitution effect that increases labor supply. It is this Frisch substitution effect of predictable wage changes we want to estimate.

Our key task then is to choose an instrument that is known to workers at the start of each year, and that generates predictable wage growth during the year. MaCurdy (1981) and many subsequent papers use education as the primary instrument for wage growth. The motivation is that annual wage growth tends to be faster for more educated workers.⁶ We adopt a closely related approach: The NLSY97 administered an aptitude test called the Armed Services Vocational Aptitude Battery (ASVAB) to respondents when they were 13 to 18 years old.⁷ We find that the ASVAB percentile score is a stronger predictor of wage growth than education, so we use that as our instrument. But the idea is similar: Not surprisingly, wage growth is predictably faster for higher ability workers.

⁴Of that, 6748 is a random sample of the birth cohort while 2236 is an over-sample of minority groups.

⁵Observations were excluded if income was less than \$3,000, the annual wage was less than \$2.70 per hour worked, the total number of hours worked was less than 400 or above 4,160 (roughly 80 hours a week), or if the percentage change in wages from the last year was below -50% or above 70%.

⁶He also used interactions of education and age, to allow the effect of education to differ by age.

⁷The ASVAB measures aptitude in several areas including mathematics, general science, paragraph comprehension, and mechanical skills. It was administered in summer 1997 to spring 1998, when the youth were aged 13 to 18 (those aged 13 to 14 were given an easier version of the test). The NLS grouped respondent’s into three-month age windows and calculated a youth’s percentile rank within his age group.

We did the analysis separately for men and women, as prior literature has shown that their labor supply behavior differs in important ways. Interestingly, the ASVAB score is a much better predictor of wage growth for men than women.⁸ For this reason, we decided to focus only on results for men. Our full data set has 5,931 annual observations on 2,100 young men aged 22 to 30 who we observe over 2 to 6 years (the average being 3.8 years).

III. NLSY97 Estimates of the Frisch Elasticity

Table 1 presents regressions of changes in log hours on changes in log wages, as in equation (1). The first column reports OLS results. The coefficient on the log wage change is -0.42 and highly significant, with a standard error of 0.015.⁹ This implies a 10% wage increase is associated with a 4.2% reduction in work hours. There are two reasons for a negative relationship: Of course surprise wage changes may generate income effects that reduce labor supply. But it is implausible that income effects alone could generate such a large negative effect.

A second key factor driving the OLS estimate negative is “denominator bias” arising because the wage is measured as the ratio of earnings to hours. If hours in the denominator are measured with error, it causes a worker’s measured wage to be too low precisely when his measured hours are too high. This induces an (artificial) negative covariance between measured hours and measured wages that drives the estimated elasticity negative. As a result, the OLS estimate cannot be interpreted causally. A second virtue of instrumenting for wage changes is that it also deals with this measurement error problem - see Altonji (1986).

Next we consider the 2SLS results. The second column of Table 1 reports the first stage, where we regress log wage changes on the ASVAB percentile score to construct predictable wage changes. The coefficient is 0.039 and highly significant (standard error 0.012). The effect size is substantial: A male worker in the 100th percentile of ability is predicted to have annual wage growth 3.9% higher than a male worker in the 1st percentile. The heteroskedasticity robust F-test for significance of ability in the first stage regression is 10.12, which gives a p -value of 0.002. So the ASVAB instrument is significant at well above the 1% level, and passes the Staiger and Stock (1997) $F > 10$ rule of thumb for IV strength.

Notably, however, the R^2 of the first stage regression is only .007, implying a correlation between our predictions and actual wage changes of .084. In fact, the partial R^2 that shows the fraction of wage variation explained by the ASVAB test alone is .002, implying a partial correlation of only .041. This illustrates the point that annual wage growth is very hard to predict. It is important to emphasize,

⁸It is not clear if this is because wages grow relatively faster for high ability men than for high ability women, or because the ASVAB is not as good a proxy for labor market skills of women. It would be interesting to explore this issue in future research.

⁹All standard errors and F-statistics reported in this paper are heteroskedasticity robust or cluster robust. The cluster robust standard errors account for both heteroskedasticity and serial correlation. They are always slightly smaller, because the errors in the hours change regression exhibit negative serial correlation. Hence the heteroskedasticity robust statistics are slightly more conservative.

TABLE 1—FRISCH ELASTICITY ESTIMATES - NLSY97

| | OLS | 2SLS 1 st Stage | 2SLS 2 nd Stage | Reduced Form |
|--|------------------------------|-------------------------------|-------------------------------|-----------------------------|
| Dependent Variable: | ΔH | ΔW | ΔH | ΔH |
| Wage Change | -0.416 (0.015) [0.015] | | 0.597 (0.403) [0.363] | |
| ASVAB Ability Score | | 0.039 (0.012) [0.011] | | 0.024 (0.011) [0.010] |
| F-Stat (Hetero- σ Robust) <i>p-value</i> | | 10.12 0.002 | | 4.47 0.035 |
| F-Stat (Cluster Robust) <i>p-value</i> | | 12.23 0.001 | | 5.64 0.018 |
| R^2 | 0.210 | 0.007 | | 0.009 |

Note: Heteroskedasticity robust standard errors are in parentheses and clustered standard errors (by individual) are in square brackets. All regressions controls for year effects, age, and race/ethnicity. $N = 5,931$

however, that a higher first-stage R^2 would not necessarily be desirable in this context. Measured wage changes contain both unpredictable and measurement error components that we specifically want to filter out, so we expect the R^2 of the first stage regression to be far less than one.

Now consider the second stage 2SLS results, where we regress log hours changes on log predictable wage changes to obtain an estimate of the Frisch elasticity. This is reported in the third column of Table 1. The 2SLS estimate of the Frisch elasticity is 0.597, implying that a 10% predictable wage increase generates a 6% increase in work hours. So the use of 2SLS flips the sign of the coefficient.

This 2SLS estimate is clearly more reasonable: Economic theory predicts a positive Frisch elasticity, as a predictable wage increase should have a positive substitution effect on labor supply. And a Frisch elasticity of 0.6 is well within the range of estimates surveyed in Keane (2011, 2021), although it is clearly towards the high range of estimates for young men.

Notice however, that the (heteroskedasticity robust) standard error on the 2SLS estimate is a substantial .403, giving a t -statistic of only 1.48 and a p -value of 0.138. So, while the estimated Frisch elasticity is a substantial 0.6, it is not even significantly different from zero at the 10% level.¹⁰ This imprecision leaves us in a quandry over what we ought to conclude from the analysis.

¹⁰The cluster robust standard error is slightly smaller, at 0.363, because the serial correlation in the hours change regression is negative. But even then the t -stat is only 1.65 (p -value = 0.099).

The imprecision in our 2SLS estimate is a consequence of the fact that the ASVAB score only explains a small part of the variance of wage changes. Because the partial correlation between the ASVAB score and wage changes is .041, the standard error goes up by a factor of 25 when we go from OLS to 2SLS (i.e., $1/.041 \approx 25$). This imprecision in 2SLS estimates has plagued much of the literature on estimating the Frisch elasticity using the 2SLS approach.

Should we trust the 2SLS results in this case? The first-stage F -statistic for the ASVAB instrument exceeds the commonly used weak IV threshold of 10 (if only marginally), suggesting the results may be viewed as reliable.¹¹ However, weak IV tests of the type developed by Staiger and Stock (1997) and Stock and Yogo (2005) are designed to assess bias in 2SLS estimates and size distortions in t -tests. They say nothing about whether instruments are strong enough for the t -test to have acceptable power properties.¹² In the next section we present AR test results, and in Section V we assess the relative performance of the two tests in this environment. We show that the t -test has very poor power properties in the present application, and that the AR test ought to be relied on instead.

IV. The Anderson-Rubin Approach

Anderson and Rubin (1949) developed an alternative approach to inference that can also be used to test if our estimate of the Frisch elasticity is significant. The Anderson-Rubin (AR) test relies on a reduced form regression of the outcome of interest on the instrument itself, along with the control variables. In our case this is a regression of the change in log hours on the ASVAB score itself, along with the controls (time, age, race). The AR test judges the Frisch elasticity estimate to be significant if the ASVAB score is significant in the reduced form regression.

The logic of the AR test is simple: A fundamental assumption of the IV method is that the instrument only affects the outcome of interest indirectly through its effect on the endogenous variable. Hence, if the instrument is significant in the reduced form, it implies that the endogenous variable has a causal impact on the outcome of interest. In our case, if the ASVAB score is significant in the reduced form, it implies that predictable wage changes influence work hours.

Of course the ASVAB score could appear significant in the reduced form merely because it somehow affects hours growth directly (not indirectly via its effect on wage growth). That is, the ASVAB score may be significant because the exclusion restriction is violated. But in that case the ASVAB score is not a valid instrument, so the 2SLS estimate and t -test results are also invalid. The very assumptions that make the IV approach valid also make the AR test valid.

¹¹For example, Stock and Watson (2015, p.490) say: “One simple rule of thumb is that you do need not to worry about weak instruments if the first stage F -statistic exceeds 10.”

¹²Stock and Yogo (2005) proposed critical values for F based on “worst-case” size distortion in t -tests. For example, in the exactly identified case, a sample $F > 8.96$ gives 95% confidence that a two-tailed 5% t -test will reject $H_0 : \beta = 0$ at a 15% rate or less when the true β is zero. That is, it has a size distortion of no more than 10%. But passing such a test does not imply the t -test will have acceptable power.

The last column of Table 1 reports the reduced form results. Here, the ASVAB score is clearly significant, with a t -stat of 2.18 (p -value 0.035). So we are left with a quandry: The AR test indicates the 2SLS estimate of the Frisch elasticity is significant, while the t -test says it isn't. Which result should we believe?

The AR test is recommended by theory as clearly superior to the t -test when instruments are weak, and no worse when instruments are strong - see Andrews, Stock and Sun (2019). This is because the AR test has three major advantages: First, it is “robust” to weak instrument problems, which means a 5% level AR test rejects a true null hypothesis at the correct 5% rate *regardless* of the strength or weakness of the instruments. In contrast, the t -test suffers from size distortions: If instruments are weak, a 5% t -test may reject a true hypothesis at rates well above or below 5%, depending on details of the situation. Second, the AR test is unbiased, meaning its power is appropriately minimized when the null corresponds to the true β . Third, Moreira (2009) shows that in the case of a single instrument (as we have here) the AR test is the most powerful unbiased test: If the null hypothesis is false, the AR test will reject the null, and conclude the parameter of interest is significant, at least as frequently as any other unbiased test.¹³

Despite its clear advantages, the AR test has been widely neglected by applied researchers. In fact, with the exception of Lee (2001), it has never been adopted in the large literature on estimating the Frisch elasticity. In the next section we present a numerical experiment based on our data that shows the performance of the AR test is *dramatically* superior in practice. We also present analytical results that lead to the same conclusion.

V. Monte Carlo Experiment and Power Analysis

In this section we compare the AR test and the t -test to see which is a more reliable guide to the statistical significance of our Frisch elasticity estimate. To do this, we conduct the following experiment: We start from the NLS sample of $N=5,931$ observations that we used to generate the estimates in Table 1. We can then “bootstrap” a new artificial dataset by sampling 5,931 observations with replacement from the original sample. We do this 10,000 times to form 10,000 artificial datasets. We then repeat the analysis of Table 1, applying OLS and 2SLS to all 10,000 datasets, and summarize the results in Table 2.¹⁴

¹³Andrews, Stock and Sun (2019) argue that the AR test should be widely adopted by applied researchers. They state its advantages more formally: “In just-identified models ... Moreira (2009) shows that the AR test is uniformly most powerful unbiased. ... Thus, the AR test has (weakly) higher power than any other size- α unbiased test no matter the true value of the parameters. In the strongly identified case, the AR test is asymptotically efficient in the usual sense and so does not sacrifice power relative to the conventional t -test. ... Since AR confidence sets are robust to weak identification and are efficient in the just-identified case, there is a strong case for using these procedures in just-identified settings.” Moreira and Moreira (2019) extend this optimality result to models with heteroskedasticity and clustering.

¹⁴By sampling with replacement from the original 5,931 observations we break the panel structure of the data. As a result, the standard errors and F statistics in Table 2 will mimic the heteroskedasticity robust statistics in Table 1, not the cluster robust statistics.

TABLE 2—RESULTS FROM MONTE CARLO BOOTSTRAP SAMPLES

| | OLS | | 2SLS | | First Stage | Reduced Form | |
|-----------|---------------|--------|---------------|----------|---------------|--------------|--------|
| | $\hat{\beta}$ | S.E. | $\hat{\beta}$ | S.E. | F Statistic | $\hat{\pi}$ | S.E. |
| Median | -0.4163 | 0.0146 | 0.5998 | 0.4013 | 10.1314 | 0.0238 | 0.0112 |
| Mean | -0.4164 | 0.0146 | 0.7185 | 4.7202 | 11.0923 | 0.0237 | 0.0112 |
| Std. Dev. | 0.0148 | 0.0004 | 3.8636 | 251.4631 | 6.4577 | 0.0111 | 0.0003 |

Note: $N = 5,931$ for each of the 10,000 samples used to form the results.

A. OLS Estimates and Standard Errors

In Table 2 we see that both the median and mean OLS estimates of β (across all 10,000 datasets) are roughly equal to the (downward biased) value of -0.416 we obtained using the original NLS sample. This is as expected, as our 10,000 “bootstrap” datasets mimic the covariances of the variables in the original NLS sample. The third row of Table 2 reports that the standard deviation of the OLS estimates across the 10,000 artificial samples is 0.015, which equals (to three decimal places) the OLS standard error estimate reported in Table 1. Thus, the estimated OLS standard error is a very good guide to how the OLS estimates actually vary across the different samples.¹⁵

B. 2SLS Estimates and Standard Errors

Now we examine how the 2SLS estimates and standard errors behave. The first thing to note in Table 2 is that the median 2SLS estimate of the Frisch elasticity (across all 10,000 datasets) is 0.600, which is very close to the 2SLS estimate 0.597 we obtained using the original NLS dataset. This is exactly as expected: As our artificial datasets are constructed from our original NLS sample, we can think of the NLS sample as the “population” from which all 10,000 datasets are drawn. In this population, 0.597 is in fact the true value of the Frisch elasticity. We see that the median 2SLS estimate accurately uncovers the true Frisch elasticity value. As Keane and Neal (2021) and Angrist and Kolesár (2021) show, the median bias of 2SLS is negligible at this level of instrument strength.

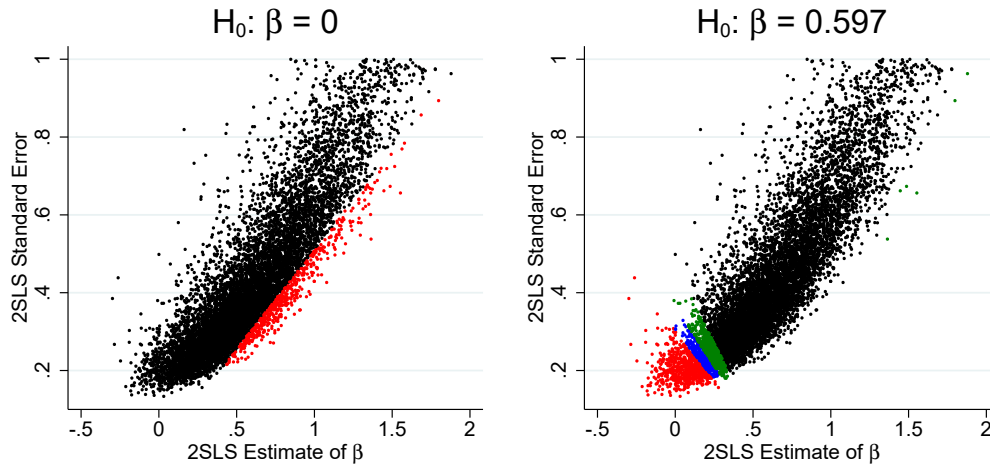
Second, note that the median of the estimated 2SLS standard errors, reported in the first row of Table 2, is 0.401. This agrees closely with the 2SLS standard error estimate of 0.403 in Table 1. However, the actual empirical standard deviation of the 2SLS estimates across the 10,000 data sets is 3.864. In contrast to OLS, this bears no resemblance to the estimated 2SLS standard errors. This is our first indication that the 2SLS standard errors are not a good guide to the actual

¹⁵Table 2 also reports the mean and median of the estimated OLS standard error across the 10,000 artificial datasets. These are again 0.015. And the variation across samples of this standard error estimate is trivially small. So the estimated standard error in each individual sample is a good guide to the actual variability of the OLS estimates across all samples.

variability of the 2SLS estimates across samples.¹⁶ This in turn means that 2SLS t -statistics – which rely on those standard error estimates – will not be a useful guide to significance of 2SLS estimates.

To further explore the behavior of the 2SLS standard error, Figure 1 plots the 2SLS standard errors against the 2SLS estimates of the Frisch elasticity from each of the 10,000 samples. A striking aspect of the figure is the strong positive association between 2SLS estimates and their standard errors: The Spearman correlation is an extraordinarily large 0.905. This means that in samples where the estimated Frisch elasticity is larger, the standard error is also larger. As we will see, this pattern has extremely important empirical implications.

FIGURE 1. STANDARD ERROR OF $\hat{\beta}_{2SLS}$ PLOTTED AGAINST $\hat{\beta}_{2SLS}$ ITSELF



Note: Runs with standard error > 1 are not shown. In the left panel, the red dots indicate $H_0: \beta = 0$ is rejected at the 5% level using the 2SLS t -test, while in the right panel red dots indicate $H_0: \beta = 0.597$ is rejected at the 5% level. Blue and green indicate 10% and 20%.

The association between 2SLS estimates and their standard errors is not specific to this application. It is a generic but little-appreciated property of the 2SLS estimator. Exact finite sample theory can shed light on this phenomenon. Phillips (1989) derives two key properties of 2SLS in the unidentified case. First, the 2SLS estimator converges in distribution to a scale mixture of normals centered on $E(\hat{\beta}_{OLS})$. Second, the 2SLS variance estimator ($\hat{\sigma}^2$) converges in distribution to a quadratic function of $\hat{\beta}_{2SLS}$, with a minimum at $E(\hat{\beta}_{OLS})$. Thus, the standard error of regression ($\hat{\sigma}$) is minimized when $\hat{\beta}_{2SLS}$ is close to $E(\hat{\beta}_{OLS})$. Of course,

¹⁶Of course, in the single instrument case the mean and variance of the 2SLS estimator do not exist, which means that if we did many more than 10,000 runs the mean and variance wouldn't converge. Hence, the standard deviation of the 2SLS standard error cannot be bootstrapped.

the standard error of the regression ($\hat{\sigma}$) is a fundamental driver of the standard error of $\hat{\beta}_{2SLS}$. Thus, in the unidentified case, the standard error of $\hat{\beta}_{2SLS}$ tends to be minimized when the estimate is near $E(\hat{\beta}_{OLS})$.

Importantly, these properties of 2SLS in the unidentified case still influence the behavior of 2SLS estimates and standard errors in strongly identified models. In fact, Phillips (1989) calls this the “leading case” as it provides the leading term of the series expansion of the density of the estimator in the general case. As a result, even in strongly identified models, the standard error of $\hat{\beta}_{2SLS}$ tends to be minimized when the estimate is near $E(\hat{\beta}_{OLS})$, as we see in Figure 1. Intuitively, in finite samples where the exogenous instrument is positively correlated with the structural error, the estimate is shifted toward OLS and appears more precise.¹⁷ In Keane and Neal (2021) we fully explore the implications of this phenomenon. For our present purposes it suffices to note the following: Because of this pattern, large positive 2SLS estimates of the Frisch elasticity will have relatively large standard errors, while estimates near zero will have small standard errors.

This brings us to our key point: The positive association between 2SLS estimates of the Frisch elasticity and their standard errors has important implications for statistical inference. As we now show, this mechanical relationship makes it very difficult for a 2SLS t -test to detect a true positive Frisch elasticity.

Recall that our 10,000 simulated data sets are constructed so the true value of the Frisch elasticity in these data sets is 0.597. Thus, if the 2SLS t -test is reliable it should have two properties: First, if we run 5% t -tests of the hypothesis that the true Frisch elasticity is zero we should reject that false hypothesis at a high rate (indicating the test has good power). Second, if we run 5% t -tests of the true hypothesis that the Frisch is equal to 0.597 (the true value) we should reject that hypothesis approximately 5% of the time (indicating the test has correct size). Furthermore, those rejections should be evenly split between cases where the estimated Frisch elasticity is above and below the true value.

In the left panel of Figure 1 we shade in red the cases where the 2SLS t -test rejects the false null hypothesis that the true Frisch elasticity is equal to zero. These are the cases where the ratio of the estimate to the standard error exceeds the 5% critical level of 1.96 (in absolute value). Notice how the red shaded area

¹⁷The Appendix provides additional mathematical detail. Intuitively, the association between 2SLS estimates and their standard errors arises for the following reason: As we discussed in Section III, in the original NLS sample the partial correlation between the ASVAB score and wage growth is 0.04. When we look across our 10,000 subsamples, the correlation fluctuates around that value due to sampling variation. Two things happen in samples where that correlation is relatively high:

First, the 2SLS standard error estimate is smaller: The stronger is the correlation between the instrument and the endogenous variable, the smaller is the 2SLS standard error.

Second, the 2SLS estimate is more shifted in the direction of the OLS bias (which is negative). This is because, as we discussed in Section III, if the predictable part of wage growth is small, then a high correlation between the instrument and the endogenous variable is not really a good thing. In samples where that correlation rises above 0.04, the instrument is picking up some of the endogenous part of wage growth that arises due to measurement error and surprise wage growth. This in turn means the 2SLS estimate will be shifted in the direction of the OLS bias (negative).

Putting these two facts together, it means that 2SLS estimates that are most shifted in the direction of the OLS bias (negative) appear to be more precise. This is exactly the pattern we see in Figure 1.

is quite small. In fact, the false null hypothesis is only rejected 5.1% of the time. This is an extremely low level of power. It scarcely exceeds the 5% rate at which a well-behaved 5% level test should reject a *true* null hypothesis.

In the right panel of Figure 1 we shade in red the cases where the 2SLS t -test rejects the null hypothesis that the true Frisch elasticity is equal to the true value of 0.597. The test rejects the null hypothesis 6.6% of the time. This is not bad when viewed in isolation, as it is not too far from the correct rate of 5%, so the t -test size distortion is small. But more importantly, the rate of rejecting the true hypothesis that the Frisch equals 0.597 is actually *greater* than the rate of rejecting the false hypothesis that the Frisch equals 0. This is extremely poor behavior for a statistical test: The fact that size exceeds power means the t -test is uninformative about the true parameter value.

These results illustrate our key point: The population F easily passes conventional weak instrument testing thresholds of the type developed by Staiger and Stock (1997) and Stock and Yogo (2005). So, as expected, the median bias of 2SLS and the size distortion in the t -test are trivial. Nevertheless, the 2SLS t -test results are completely uninformative, as size exceeds power.

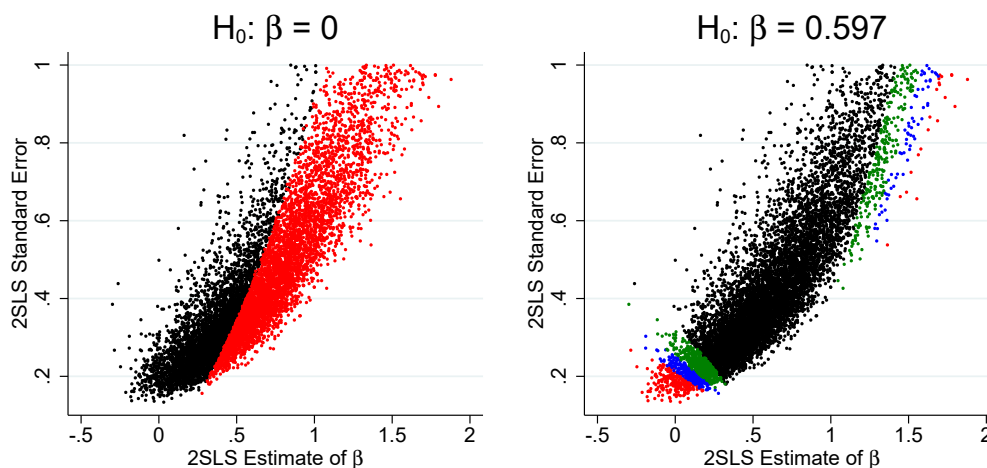
Another notable aspect of the right panel of Figure 1 is that the cases where we reject the null of Frisch = 0.597 are not evenly split between cases where the estimate is above and below the true value. In fact, all the rejections occur when the estimated Frisch elasticity is very small (near zero). This is a direct consequence of the positive association between 2SLS estimates and their standard errors. As large positive estimates of the Frisch elasticity have large standard errors, there is very little chance of concluding a large positive estimate is significant.

C. Anderson-Rubin Test Results

Figure 2 reports the same results for the AR test. The contrast with the t -test is dramatic. We again plot the 2SLS standard errors against the 2SLS estimates, as in Figure 1. But now we plot in red the cases where the AR test rejects the false null hypothesis that $\beta = 0$ at the 5% level. In the left panel we see that the red region is quite large. The AR test rejects the false null that the Frisch is equal to zero 56.5% of the time. This is a good level of power that is more than ten times greater than the 5.1% rate achieved by the t -test.

The right panel of Figure 2 shows the rate of rejecting the true null hypothesis that the Frisch equals 0.597. The AR test rejects 4.9% of the time, which is almost exactly equal to the correct 5% rate. Furthermore, we plot in blue and green the cases where 10% and 20% AR tests reject. These rates are 10% and 19.5%, so again almost perfect. This illustrates how the AR test is “robust” in the sense that it has correct size (rejection rates) regardless of the strength or weakness of the instruments. Thus we see that *the AR test has correct size and ten times the power of the t -test.*

The only limitation of the AR test is that it doesn’t quite generate symmetric rejections when the estimates are above and below the true value. For example,

FIGURE 2. STANDARD ERROR OF $\hat{\beta}_{2SLS}$ PLOTTED AGAINST $\hat{\beta}_{2SLS}$ ITSELF (AR TEST)

Note: Runs with standard error > 1 are not shown. In the left panel, red dots indicate $H_0 : \beta = 0$ is rejected at the 5% level using the AR test. In the right panel red dots indicate $H_0 : \beta = 0.597$ rejected at the 5% level using the AR test. Blue and green indicate rejections at the 10% and 20% levels, respectively.

of the 4.9% rejections in the 5% test, 3.6% occur when the estimate is below 0.597 and 1.3% occur when it is above. This is because, like the t -test, the AR test tends to attribute greater precision to estimates shifted in the (negative) direction of the OLS bias, and less precision to large positive estimates. But this problem is much less severe for the AR test than the t -test.¹⁸

These results make it very obvious that in the data environment of our empirical application the AR test provides a far more reliable guide to the significance of the estimate of the Frisch elasticity than does the t -test. The consequence is that prior work that relied on 2SLS t -tests will have tended to obtain insignificant results even if the true Frisch elasticity is well above zero.

The superiority of the AR test over the t -test is not specific to this example. In Keane and Neal (2021) we show that the superiority of the AR test is evident across a wide range of contexts. 2SLS t -tests perform poorly in general due to the strong association between 2SLS estimates and their standard errors. As a result of this pattern, t -tests have difficulty detecting true negative (positive) effects when the OLS endogeneity bias is positive (negative). This problem is relevant across a wide range of empirical applications, including cases where instruments are much stronger than here. The AR test is much less susceptible to this problem.

¹⁸In Keane and Neal (2021) we show that the power asymmetry in the AR test vanishes quickly as instruments become stronger. But the power asymmetry in the t -test remains substantial even with very strong instruments.

D. Analytical Power Function Comparison: AR vs t -test

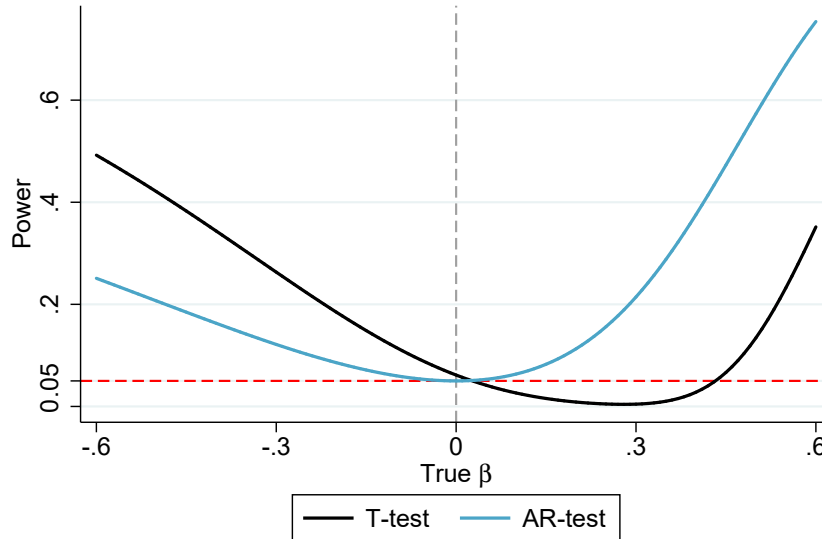
To show the generality of the problem we describe, we now compare the analytical power functions of the AR and T -tests in an exactly identified linear IV model with *iid* normal errors:

$$(2) \quad \begin{aligned} y_i &= \beta x_i + u_i \\ x_i &= \pi z_i + e_i \quad \text{where} \quad e_i = \rho u_i + \sqrt{1 - \rho^2} \eta_i \end{aligned}$$

where $u_i \sim iidN(0, 1)$, $\eta_i \sim iidN(0, 1)$, and $z_i \sim iidN(0, 1)$. Thus the instrument z satisfies $cov(z, u) = 0$ and $cov(z, \eta) = 0$. The parameter $\rho \in (-1, 1)$ determines the severity of the endogeneity problem, while π determines the strength of the instrument. This *iid* normal setup is not as restrictive as it may first appear, as Andrews, Stock and Sun (2019) show that for any heteroskedastic DGP, there exists a homoskedastic DGP yielding equivalent behavior of 2SLS estimates and test statistics. Any exogenous covariates can be partialled out of y and x without changing anything of substance.

Figure 3 compares the power functions of the AR and t -tests, obtained via the procedure described in the Appendix, with population $F=10.12$ and $\rho = -.70$ to mimic our empirical application. The power function is the probability a 5% level test rejects $H_0: \beta = 0$, conditional on each alternative true β listed on the x -axis:

FIGURE 3. POWER OF THE T-TEST VS. AR-TEST WHEN POP. $F = 10$ ($\rho = -0.7$)



Note: Probability a 5% level test rejects $H_0 : \beta = 0$, conditional on each alternative true β listed on the x -axis. In order to most closely match our application, we set $\rho = -0.7$ and the population first-stage F -statistic to 10.12.

Several features of the power functions are notable. First, the power of the 2SLS 5% level t -test to reject $H_0 : \beta = 0$ is close to 5% if the true $\beta = 0$. So as expected the t -test has approximately correct size when population $F = 10$. Both Angrist and Kolesár (2021) and Keane and Neal (2021) argue that size distortions in two-tailed 2SLS t -tests are modest unless instruments are very weak and endogeneity is very severe, and that fact is reflected here.¹⁹

Second, the severe bias of the t -test is evident: Its power dips below 5% for a wide range of positive values of β , dipping to near zero when true β is around 0.25 to 0.30. In the model in (2), β is approximately the standard deviation change in y induced by a one standard deviation change in x . So effect sizes of 0.25 to 0.30 would be substantial in most applications. Yet the 2SLS t -test has essentially no power to detect effect sizes in this range. Third, the unbiasedness of the AR test is also evident, as its power is appropriately minimized at $\beta = 0$.

Fourth, and most importantly, the AR test clearly has far superior power to detect true positive values of β in this environment, in which the OLS bias is negative. The lack of power of the t -test in the positive β range is due to the positive association between 2SLS estimates and their standard errors that we have emphasized. This causes larger positive estimates of β to have spuriously inflated standard errors. Conversely, the t -test *appears* to have better power than AR when the true β is negative, but this is problematic, as it occurs because larger negative estimates of β tend to have spuriously small standard errors.

The Appendix contains additional discussion of analytical power results for different levels of instrument strength. The patterns we discuss here are more pronounced when instruments are weak, but persist even when instruments are far above conventional weak instrument testing thresholds.

VI. Interpreting the Empirical Results in Light of the Experiment

We now return to our empirical results in Table 1, and assess them based on what we have learned from the Monte Carlo and analytical power analysis of Section V. Recall that our 2SLS estimate of the Frisch elasticity based on the ASVAB instrument is 0.597, but the 2SLS t -test indicates this is not significantly different from zero at the 5% level. However, the analysis of Section V clearly indicates that the t -test has little power to detect true positive effects of plausible magnitude in this data environment, characterized by an F of roughly 10 and a correlation between the reduced form residuals of -0.70.

The analysis of Section V also revealed that the AR test is a far more reliable method of inference in this context. The AR test is based on the significance of the instrument (ASVAB) in the reduced form regression of hours changes on wage

¹⁹Stock and Yogo (2005) derived critical values that the sample F must surpass in order for a researcher to have high confidence the “worst-case” size distortion in the two-tailed t -test is modest. But the “worst case” occurs when ρ is near one or minus one, so endogeneity is extremely severe. For smaller values of ρ – that are more typical of applications – much lower levels of F will suffice to render size distortions quite modest. Keane and Neal (2021) and Angrist and Kolesár (2021) also make this point.

changes. It is well-known that AR has desirable properties relative to the t -test: Moreira (2009) shows it is the uniformly most powerful unbiased test in exactly identified linear IV models. It is robust to weak instrument problems, as it is guaranteed to have correct size, yet it is guaranteed to be no less powerful than the t -test when instruments are strong. For this reason, the AR test is widely recommended by theorists for use when instruments are weak.

The main contribution of our analysis in Section V is to show that the AR test has greatly superior power properties to the t -test even when instruments are quite strong. It appears that the heavy focus of the weak IV literature on estimation bias and hypothesis test size has diverted attention from the extremely poor power properties of the t -test in environments where instruments are strong enough to pass conventional weak IV testing thresholds. For this reason we argue the AR test should replace the t -test quite generally for inference in linear IV models, not only when instruments are weak but even when they are strong.

The AR test indicates that our Frisch elasticity estimate of 0.597 is significant at the 3.5% or 1.8% level, depending on whether we rely on the heteroskedasticity robust or cluster robust standard error. We can also invert the AR test to obtain a weak instrument robust confidence interval, as discussed in Anderson and Rubin (1949).²⁰ Using cluster robust statistics we obtain a 95% confidence interval for the Frisch elasticity of 0.082 to 2.03, which is clearly bounded above zero, and covers most of the range often used to calibrate macro models.

VII. Results Based on Multiple Instruments

As we discussed in Section II, much of the prior work on estimating the Frisch elasticity used education as the instrument for wage growth, but we rely on the ASVAB score as we find it is a stronger instrument in the first stage of 2SLS. In this section we consider using both education and the ASVAB score as instruments. In order to keep the sample identical to that in Table 1, we code education as zero if it is missing, and introduce a dummy for missing education as an additional instrument. As we see in the first column of Table 3 both the ASVAB score and education are significant in the first stage of 2SLS, suggesting they capture somewhat different dimensions of ability.²¹

We also report two versions of the partial F -statistic for joint significance of the instruments in the first stage (heteroskedasticity and cluster robust), as well as the Oleva-Pfleuger effective F -test for weak instruments in a non-*iid* setting. These statistics range from 4.3 to 5.1, so they are well below conventional weak

²⁰The basic idea of AR test inversion is to run regressions of $y - xb$ on the instrument and control variables, and find the lower and upper cutoffs for b where the AR test p-value is exactly .05.

²¹To be precise, the p -values for education are .047 or .071 based on the cluster robust or heteroskedasticity robust standard error, respectively.

TABLE 3—FRISCH ELASTICITY - OVER-IDENTIFIED MODELS

| Dependent Variable | 2SLS 1 st Stage | 2SLS 2 nd Stage | Reduced Form | GMM-2S 2 nd Stage | GMM-CU 2 nd Stage |
|----------------------------------|-------------------------------|-------------------------------|------------------------------|---------------------------------|---------------------------------|
| | ΔW | ΔH | ΔH | ΔH | ΔH |
| Wage Change | | 1.017 (0.481) [0.442] | | 0.896 (0.474) [0.433] | 1.310 (0.548) [0.487] |
| ASVAB Ability Score | 0.028 (0.014) [0.012] | | -0.007 (0.017) [0.016] | | |
| Education | 0.002 (0.001) [0.001] | | 0.006 (0.003) [0.002] | | |
| Education Missing | 0.033 (0.034) [0.035] | | 0.033 (0.044) [0.043] | | |
| F-Stat (Hetero- σ Robust) | 4.31 | | 4.21 | | |
| <i>p-value</i> | 0.005 | | 0.006 | | |
| F-Stat (Cluster Robust) | 5.14 | | 4.75 | | |
| <i>p-value</i> | 0.002 | | 0.003 | | |
| Olea-Pfleuger Effective F | 4.57 | | | | |
| Exogeneity Test (AR or J) | | 2.77 | | 3.19 | 3.00 |
| <i>p-value</i> | | 0.428 | | 0.203 | 0.224 |
| R^2 | 0.008 | | 0.013 | | |

Note: ‘GMM-2S’ refers to the 2-step GMM, while ‘GMM-CU’ refers to continuously updated GMM. Heteroskedasticity robust standard errors are in parentheses and clustered standard errors are in square brackets. All regressions controls for year effects, age, and race/ethnicity. $N = 5,931$.

instrument testing thresholds.²² Thus weak instruments are clearly a concern and the 2SLS t -test cannot be viewed as reliable.

The 2SLS estimate of the Frisch elasticity is 1.017, which is much larger than the estimate of 0.597 we obtained in Table 1. Notably, the heteroskedasticity robust standard error increases from 0.403 to 0.481, so $t=2.12$ ($p=.034$) and a 5% t -test judges our estimate significant.²³ It may seem surprising that the 2SLS standard error increases despite the efficiency gain from adding an additional relevant instrument in the first-stage. But we have discussed how the increase

²²Andrews, Stock and Sun (2019) point out that *in general* it is inappropriate to use either a heteroskedasticity-robust or conventional F -test to assess instrument strength in non-homoskedastic settings, and suggest using the Olea and Pflueger (2013) effective first-stage F -statistic. However, as they point out, in the single instrument just-identified case that we considered back in Sections III to VI, this reduces to the conventional heteroskedasticity-robust F .

²³The cluster robust standard error increases from 0.363 to 0.442, giving $t=2.30$ ($p=.021$).

in the Frisch estimate from .597 to 1.017, which moves us further from the OLS bias, will mechanically cause the 2SLS standard error of regression to increase. This tends to inflate the standard error of the 2SLS estimate.

Now consider the AR test, which in the over-identified case is simply the F -test for joint significance of the three instruments in the reduced form. The cluster robust version of the AR test gives a p value of .0026. Moreover, the AR test is not the most powerful test in the over-identified case: the weak instrument robust conditional likelihood ratio (CLR) test of Moreira (2003) is more efficient. The cluster-robust CLR test has a p -value of .0012, so the evidence for a positive Frisch elasticity based on the robust statistics is very strong.

So here we see a milder version of the pattern in Table 1: The 2SLS t -test implies the Frisch elasticity estimate is (just) significant at the 5% level, while the weak instrument robust statistics (the AR and CLR tests) imply much higher levels of confidence. The relative weakness of the 2SLS t -test result is again attributable to the positive covariance between 2SLS estimates and standard errors, which makes it difficult for 2SLS t -tests to detect a positive Frisch elasticity.

A useful feature of the AR test is that we can evaluate it at $\hat{\beta}_{2SLS}$ rather than 0 to obtain a test of the 2SLS over-identifying restrictions. This is just an F -test for joint significance of the excluded instruments in a regression of the 2SLS residuals on all the instruments. Henceforth, we refer to these as the AR(0) and AR($\hat{\beta}_{2SLS}$) tests to distinguish the two. As we see in Table 3, the AR($\hat{\beta}_{2SLS}$) test statistic is 2.77. The test is distributed $\chi^2(3)$ so the p -value is .428. Thus we cannot reject the exogeneity of the instruments. This is important, as a failure of the over-identification test would invalidate the AR test, as it would suggest the instruments may be significant in the reduced form merely because they affect hours changes directly (rather than only indirectly via wages as 2SLS assumes). So the AR(0) and AR($\hat{\beta}_{2SLS}$) statistics should be evaluated in conjunction.

To assess the relative performance of the AR and t -test in the three instrument case, we ran a Monte Carlo analysis like that of Section V, but using the 2SLS estimated model in Table 3 as the data generating process.²⁴ In terms of power, we find that a 5% t -test rejects the false null $H_0 : \beta=0$ at a 60.2% rate, compared to 88.4% for the AR test, and 94.7% for CLR. So the ranking is as expected.

If we invert the AR test (cluster robust F version) we obtain a 95% confidence interval for the Frisch elasticity of 0.241 to 4.336, while inverting the CLR test gives 0.269 to 4.461.²⁵ These intervals sit comfortably above zero, and cover the range of values typically used to calibrate macro models.

Finally, the last two columns of Table 3 report the two-step and continuously updated GMM results. These GMM estimates of the Frisch elasticity are 0.896

²⁴In the one instrument case the instrument is uncorrelated with the 2SLS residuals. So when we treat the full sample as the “population,” the instrument has zero population covariance with the structural error by construction. But in the over-identified case the instruments do have small correlations with the 2SLS residuals. We need to partial out those correlations to set up the experiment.

²⁵We use the Stata command developed by Finlay and Magnusson (2009) to implement the cluster robust version of the CLR test and to do the inversion.

and 1.310 respectively. Notice how the increase in the point estimate to 1.310, moving it even further from the OLS bias, coincides with a further increase in the GMM-CU standard error to 0.548. The GMM estimates and standard errors have the same positive covariance as the 2SLS estimates and standard errors. Thus the GMM standard errors are also unreliable in this context.

However, Stock and Wright (2000) develop a weak instrument robust test that generalizes the AR test to the GMM case. This “S-statistic” is the GMM objective function evaluated at $\hat{\beta}=0$. For GMM-CU we find $S=20.47$. The test is distributed $\chi^2(3)$ so the p -value is .0001 and the Frisch estimate is highly significant. Finally, we consider Hansen’s test of over-identifying restrictions. As we see in Table 3 the J-test has $p > 0.20$, indicating we cannot reject the exogeneity of the instruments. This is important, as a failure of the J-test would invalidate the S test.

In summary, in the over-identified case the weak instrument robust AR and S tests indicate that the 2SLS and GMM estimates of the Frisch elasticity are highly significant. We caution that both tests may reject $H_0:\beta=0$ either because the null is false or because the instruments are endogenous. Hence, before relying on the AR(0) and S test results, it is important to verify, as we have here, that the AR($\hat{\beta}_{2SLS}$) and Hansen J-tests do not reject exogeneity of the instruments.²⁶ However, we emphasize that failure of the exogeneity tests would invalidate 2SLS t -test results as well, so the reliance of the AR(0) and S test results on validity of the instruments is not a disadvantage of these robust tests relative to the t -test.

VIII. Conclusion

The magnitude of the Frisch labor supply elasticity – how work hours respond to predictable wage changes – lies at the center of many economic policy debates, because the pure substitution effect measured by the Frisch is a vital input into tax policy. For example, higher values of the Frisch imply lower optimal tax rates on labor income. Because of its importance, there is a large literature estimating the Frisch elasticity using instrumental variable methods. Most of this literature has been plagued by weak instrument problems, as it is hard to find instruments that strongly predict wage growth. This is one reason the value of the Frisch elasticity remains a topic of intense debate.

Here we revisit that debate. Using the ASVAB ability test as an instrument for wage growth, we estimate a large Frisch elasticity of 0.597 for young men using data from the NLSY97. But, as is typical of this literature, the 2SLS standard error is 0.403, implying our estimate of is very imprecise. Based on this, we can’t reject the hypothesis that the Frisch is zero at conventional levels – a result that is typical of many prior papers. In contrast, the Anderson-Rubin (AR) test indicates that our estimate of the Frisch elasticity is significant at the 3.5% level.

²⁶In the single endogenous variable, K instrument case, the AR($\hat{\beta}_{2SLS}$) and J-tests have power to detect if at least one instrument is endogenous, provided the model is over-identified, which means at least two instruments must be relevant. But power of these tests will be low if $K-1$ instruments are weak.

Importantly, the first-stage F -statistic for our ASVAB instrument is 10.12. This exceeds the Staiger-Stock rule of thumb value of 10, suggesting that bias in 2SLS estimates and size distortions in two-tailed 2SLS t -tests are unlikely to be major concerns in this context. However, conventional weak instrument testing thresholds are completely uninformative about the power properties of test statistics. Thus, we conducted both a Monte Carlo experiment and an analytical analysis to compare the power of the t -test vs. the AR test.

In our data environment we find the AR test has correct size and ten times the power of the t -test. In fact, the power of the t -test is so poor that a 5% level test is more likely to reject a hypothesis that the Frisch elasticity equals its true value than a false hypothesis that it equals zero. In other words, the t -test conveys no information because its size exceeds its power.

The poor power properties of the t -test arise due to a strong positive association between 2SLS estimates and their standard errors that exists whenever the OLS bias is negative, as is the case in the Frisch application. This causes large positive estimates of the Frisch elasticity to have artificially inflated standard errors. As a consequence, a 2SLS t -test has little power to detect a true positive Frisch elasticity. The AR test does not suffer from this power asymmetry problem.

The desirable statistical properties of the AR test are well-known. In exactly identified cases Moreira (2009) shows it is the uniformly most powerful unbiased test. Unlike the t -test, the AR test is robust to weak instrument problems, meaning it has correct size regardless of the strength or weakness of the instruments. For this reason, theorists commonly recommend using the AR test rather than the t -test when instruments are weak – see, e.g., Andrews, Stock and Sun (2019). But we show that the AR test has far superior power properties compared to the t -test even in contexts where instrument strength is well above conventional weak IV testing thresholds. Hence, we argue that the AR test should be preferred over the t -test even in environments where instruments are quite strong.

Given the clear theoretical guidance, along with empirical, Monte Carlo and analytical results like we present here, it is difficult to understand why applied researchers have not widely adopted the AR test in preference to 2SLS t -tests. It appears that the intense focus of the weak instrument literature on bias and size has diverted attention away from the extremely poor power properties of 2SLS t -tests even when instrument strength is well above conventional weak IV thresholds. It also appears likely that applied researchers are simply not aware of the severity of the power problems with 2SLS t -tests that we document here, or how well the AR test deals with these problems.^{27,28}

²⁷For instance, Lee (2001) state “... applied research, with rare exceptions, relies on t -ratio-based inference ... arguably based on ... computational convenience and the presumption that for practical purposes, the distortions in inference are small or negligible.”

²⁸In addition, applied researchers may simply be unfamiliar with AR tests and think they are difficult to implement. That is obviously not true, but the econometric theory literature on AR tests presents them at such a high level of generality that it is indeed difficult for applied researchers to penetrate. And of course, applied researchers may be wedded to t -tests simply because they are so familiar. But we hope that inertia may be overcome so that empirical practice can be improved.

Our estimated Frisch elasticity for young men (.597) may still appear small compared to the values of 1.0 or more often used to calibrate macro models. However, there is accumulating evidence that the Frisch elasticity increases substantially with age (see, e.g., Borella, De Nardi and Yang 2019, Erosa, Fuster and Kambourov 2016, French 2005 and Keane 2021), and clear evidence that it is greater for women than men (see Keane 2011). So a value of .597 for young men is quite consistent with a value of 1.0 or more in the aggregate.

We also consider an over-identified model using both education and ASVAB as instruments for wage growth. Then our estimate of the Frisch increases to 1.02, and the 2SLS t -test indicates it is significant at the 3.4% level. However, the 2SLS standard error is again inflated due to its mechanical positive covariance with the estimate. An additional Monte Carlo experiment shows the AR test has correct size and substantially better power than the t -test in this environment. The AR test gives a much higher significance level of 0.3%. If we invert the AR test we obtain a 95% confidence interval for the Frisch elasticity of 0.241 to 4.336, which covers the range of values typically used to calibrate macro models.

Of course, there have been several previous attempts to reconcile the small and sometimes insignificant 2SLS estimates of the Frisch elasticity that have often been obtained using micro data with the large values often used in macro calibrations. The various approaches are detailed in Keane and Rogerson (2012, 2015). These reconciliations fall into two broad categories: One set of explanations, exemplified by Imai and Keane (2004) and Domeij and Floden (2006) takes issue with the specification of equation (1), arguing that more general models of labor supply (e.g., models that account for human capital or liquidity constraints) imply that estimation of this equation will give downward biased estimates of the Frisch elasticity. The other set of explanations, exemplified by Chang and Kim (2006) and Rogerson and Wallenius (2009), argue that, once one accounts for the participation margin of labor supply and aggregation issues, it is possible for the macro level Frisch elasticity to be large even if the micro level elasticity is small. More recently, Gottlieb, Onken and Valladares-Esteban (2021) have shown how a large macro level Frisch elasticity can be reconciled with modest reactions to tax holidays due to a combination of income and equilibrium effects. These arguments are complementary to our argument here.

Our argument is new in that we criticise the micro-econometric literature on its own terms: Suppose the assumptions necessary for a 2SLS regression of hours changes on wage changes to deliver consistent estimates of the Frisch elasticity do in fact hold. Even then, we show that the econometric methods that have been used to draw inferences from those estimates are inherently biased against finding the Frisch is both large and significant.²⁹ We hope our straightforward econometric argument will prove convincing to economists and econometricians

²⁹Lee (2001) also criticized the econometric literature on its own terms, but from a different angle: He argued that Frisch estimates were downward biased due to weak instrument problems. We argue that even if instruments are strong enough so bias is not a concern, standard errors on positive Frisch estimates are spuriously inflated, making it difficult for a 2SLS t -test to detect a positive Frisch elasticity.

who have not been convinced by the more subtle theoretical arguments based on more complex labor supply models or aggregation issues.

We conclude by restating our key econometric point: The main reason robust statistics lead to different conclusions about the Frisch elasticity from conventional 2SLS t -tests is a basic property of 2SLS that has been generally neglected: Specifically, when the OLS bias is negative, as is the case here, 2SLS estimates and their standard errors have a positive association (as they vary across random samples from the population). This mechanical positive association causes large positive estimates of the Frisch elasticity to appear spuriously imprecise, making it difficult for 2SLS t -tests to detect a true positive. Robust statistics like the Anderson and Rubin (1949) test are much less affected by this problem, so they are better able to detect a true positive Frisch elasticity.

In a different context, where the OLS bias is positive, this pattern would be reversed, and 2SLS standard errors on *positive* estimates would be spuriously precise. That would make it difficult for 2SLS t -tests to detect a true *negative*. In the classic application of instrumental variables to estimate a treatment effect given positive selection into treatment, 2SLS t -tests have difficulty detecting true negative effects, violating a “first do no harm” principle in policy evaluation.

In Keane and Neal (2021) we explore the implications of the association between 2SLS estimates and their standard errors in more detail. There we show that the serious problems with 2SLS t -tests that we have documented here persist even when instruments are very strong, because the association between 2SLS estimates and their standard errors does not vanish as instrument strength increases. Thus it is advisable to use the AR test (or other robust tests) in lieu of the t -test even when instruments are strong.

ACKNOWLEDGEMENTS

We thank Peter Phillips, Josh Angrist, Michal Kolesar and Isaiah Andrews for helpful comments. This research was supported by Australian Research Council (ARC) grants DP210103319 and CE170100005. We have no conflicts of interest to declare.

REFERENCES

- Altonji, J.G.** 1986. “Intertemporal substitution in labor supply: Evidence from micro data.” *Journal of Political Economy*, 94(3, Part 2): S176–S215.
- Anderson, T.W., and H. Rubin.** 1949. “Estimation of the parameters of a single equation in a complete system of stochastic equations.” *Annals of Mathematical statistics*, 20(1): 46–63.
- Andrews, I., J. Stock, and L. Sun.** 2019. “Weak instruments in instrumental variables regression: Theory and practice.” *Annual Review of Economics*, 11: 727–753.

- Angrist, Joshua, and Michal Kolesár.** 2021. “One instrument to rule them all: The bias and coverage of just-id iv.” National Bureau of Economic Research.
- Borella, M., M. De Nardi, and F. Yang.** 2019. “Are marriage-related taxes and Social Security benefits holding back female labor supply?” National Bureau of Economic Research.
- Bound, J., D. Jaeger, and R. Baker.** 1995. “Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak.” *Journal of the American Statistical Association*, 90(430): 443–450.
- Chang, Y., and S. Kim.** 2006. “From individual to aggregate labor supply: A quantitative analysis based on a heterogeneous agent macroeconomy.” *International Economic Review*, 47(1): 1–27.
- Conesa, J.C., S. Kitao, and D. Krueger.** 2009. “Taxing capital? Not a bad idea after all!” *American Economic Review*, 99(1): 25–48.
- Domeij, D., and M. Floden.** 2006. “The labor-supply elasticity and borrowing constraints: Why estimates are biased.” *Review of Economic Dynamics*, 9(2): 242–262.
- Erosa, A., L. Fuster, and G. Kambourov.** 2016. “Towards a micro-founded theory of aggregate labour supply.” *The Review of Economic Studies*, 83(3): 1001–1039.
- Finlay, K., and L.M. Magnusson.** 2009. “Implementing weak-instrument robust tests for a general class of instrumental-variables models.” *The Stata Journal*, 9(3): 398–421.
- French, E.** 2005. “The effects of health, wealth, and wages on labour supply and retirement behaviour.” *The Review of Economic Studies*, 72(2): 395–427.
- Gottlieb, C., J. Onken, and A. Valladares-Esteban.** 2021. “On the Measurement of the Elasticity of Labour.” *European Economic Review*, forthcoming.
- Imai, S., and M.P. Keane.** 2004. “Intertemporal labor supply and human capital accumulation.” *International Economic Review*, 45(2): 601–641.
- Keane, M.P.** 2011. “Labor supply and taxes: A survey.” *Journal of Economic Literature*, 49(4): 961–1075.
- Keane, M.P.** 2021. “Recent Research on Labor Supply: Implications for Tax and Transfer Policy.” *Labour Economics*, 102026.

- Keane, M.P., and R. Rogerson.** 2012. “Micro and macro labor supply elasticities: A reassessment of conventional wisdom.” *Journal of Economic Literature*, 50(2): 464–76.
- Keane, M.P., and R. Rogerson.** 2015. “Reconciling micro and macro labor supply elasticities: A structural perspective.” *Annu. Rev. Econ.*, 7(1): 89–117.
- Keane, M.P., and T. Neal.** 2021. “A Practical Guide to Weak Instruments.” UNSW Economics Working Paper No. 2021-05a. Available at SSRN: <https://ssrn.com/abstract=3846841>.
- Lee, Chul-In.** 2001. “Finite sample bias in IV estimation of intertemporal labor supply models: is the intertemporal substitution elasticity really small?” *Review of Economics and Statistics*, 83(4): 638–646.
- Lee, D., J. McCrary, M. Moreira, and J. Porter.** 2020. “Valid t-ratio Inference for IV.” *arXiv preprint arXiv:2010.05058*.
- MaCurdy, T.E.** 1981. “An empirical model of labor supply in a life-cycle setting.” *Journal of political Economy*, 89(6): 1059–1085.
- Moreira, Humberto, and Marcelo J Moreira.** 2019. “Optimal two-sided tests for instrumental variables regression with heteroskedastic and autocorrelated errors.” *Journal of Econometrics*, 213(2): 398–433.
- Moreira, M.J.** 2003. “A conditional likelihood ratio test for structural models.” *Econometrica*, 71(4): 1027–1048.
- Moreira, M.J.** 2009. “Tests with correct size when instruments can be arbitrarily weak.” *Journal of Econometrics*, 152(2): 131–140.
- Olea, J.L.M., and C. Pflueger.** 2013. “A robust test for weak instruments.” *Journal of Business & Economic Statistics*, 31(3): 358–369.
- Phillips, Peter CB.** 1989. “Partially identified econometric models.” *Econometric Theory*, 5(2): 181–240.
- Prescott, E.C.** 2006. “Nobel lecture: The transformation of macroeconomic policy and research.” *Journal of Political Economy*, 114(2): 203–235.
- Rogerson, R., and J. Wallenius.** 2009. “Micro and macro elasticities in a life cycle model with taxes.” *Journal of Economic theory*, 144(6): 2277–2292.
- Staiger, D., and J. Stock.** 1997. “Instrumental variables regression with weak instruments.” *Econometrica*, 65(3): 557–586.
- Stock, J., and M. Watson.** 2015. *Introduction to econometrics (3rd global ed.)*. Pearson Education.

- Stock, J., and M. Wright.** 2000. "GMM with Weak Identification." *Econometrica*, 68(5): 1055–96.
- Stock, J., and M. Yogo.** 2005. "Testing for weak instruments in linear IV regression." *Identification and inference for econometric models: Essays in honor of Thomas Rothenberg*, 80–108.

Appendix: Analytical Power Calculations for the AR and t -tests

Consider the following just-identified *iid*-normal linear IV model:

$$(A1) \quad \begin{aligned} y_i &= \beta x_i + u_i \\ x_i &= \pi z_i + e_i \quad \text{where } e_i = \rho u_i + \sqrt{1 - \rho^2} \eta_i \\ u_i &\sim iidN(0, 1), \eta_i \sim iidN(0, 1), z_i \sim iidN(0, 1) \end{aligned}$$

The power of both the AR and t -tests depends on three parameters: the true β , the degree of endogeneity ρ , and the population t -statistic on z in the first-stage regression, which we denote λ (= square root of population F). The power of the AR test (i.e., rate of rejecting $H_0:\beta=0$ as a function of the true β) is simply:

$$(A2) \quad Power_{AR} = \Phi(\lambda D - z_{1-\alpha/2}) + \Phi(-z_{1-\alpha/2} - \lambda D)$$

where Φ is the standard normal cdf, $D = \beta/\sqrt{1 + 2\rho\beta + \beta^2}$, and $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the standard normal distribution. We set $\alpha = 0.05$.

To obtain the power function of the of the t -test we follow the analysis in Stock and Yogo (2005), Lee et al. (2020) and Angrist and Kolesár (2021). The power of the two-tailed 2SLS t -test is given by the integral:

$$(A3) \quad Power_t = \int_{-\infty}^{\infty} \left(\mathbb{I}\{t^2 \geq (1 - \rho_0^2)z_{1-\alpha/2}^2\} f(t, D, \lambda, \rho_0) + \mathbb{I}\{t^2 \geq z_{1-\alpha/2}^2\} \right) \phi(t - \lambda) dt$$

where ϕ is the standard normal density, $\rho_0 = (\rho + \beta)/\sqrt{1 + 2\rho\beta + \beta^2}$, and:

$$(A4) \quad f(t, D, \lambda, \rho_0) = \Phi\left(\frac{a_2 - \lambda D - \rho_0(t - \lambda)}{\sqrt{1 - \rho_0^2}}\right) - \Phi\left(\frac{a_1 - \lambda D - \rho_0(t - \lambda)}{\sqrt{1 - \rho_0^2}}\right),$$

$$a_1 = \frac{\rho_0 z_{1-\alpha/2}^2 t - |t| z_{1-\alpha/2} \sqrt{t^2 - (1 - \rho_0^2) z_{1-\alpha/2}^2}}{z_{1-\alpha/2}^2 - t^2},$$

$$a_2 = \frac{\rho_0 z_{1-\alpha/2}^2 t + |t| z_{1-\alpha/2} \sqrt{t^2 - (1 - \rho_0^2) z_{1-\alpha/2}^2}}{z_{1-\alpha/2}^2 - t^2}.$$

The integral in (A3) must be evaluated numerically.

To construct Figure 3 in the text, set $\lambda = 3.186$ and $\rho = -0.7$, which correspond to the empirical example in Section V, as $\lambda = 3.186$ corresponds to a population F of 10.12. The power asymmetry of the t -test is evident in Figure 3, as it has little power to detect a wide range of true positive β values due to the fact that

the OLS bias is negative ($\rho = -0.70$). The AR test is unbiased – its power is appropriately minimized when the true β is 0 – in contrast to the t -test whose power is minimized when true β is in the 0.25 to 0.30 range.

A1. Explaining the Power Asymmetry of the T-test

Next we give a simple mathematical explanation of the power asymmetry in the t -test, to compliment the intuitive explanation in Section V.B, particularly footnote 17. Recall that the 2SLS estimator of β is given by:

$$(A5) \quad \hat{\beta}_{2SLS} = \frac{\sum_{i=1}^n z_i y_i}{\sum_{i=1}^n z_i x_i} = \beta + \frac{\sum_{i=1}^n z_i u_i}{\sum_{i=1}^n z_i x_i} = \beta + \frac{\widehat{cov}(z, u)}{\widehat{cov}(z, x)}$$

In a finite sample a zero covariance between instrument and the structural error is a measure zero event. Hence $\widehat{cov}(z, u) \neq 0$ even though $cov(z, u) = 0$. We assume without loss of generality that the population covariance between the instrument and the endogenous variable is positive, $cov(z, x) > 0$. To clarify the key idea, we further assume that $\widehat{cov}(z, x) > 0$, so the sign of the coefficient on z in the first-stage regression is correct. Violation of this condition is extremely rare if the instrument is reasonably strong. For instance, in our Monte Carlo experiment in Section V, where $F = 10.12$, first-stage sign is correct in all 10,000 replications.

Under these conditions, (A5) indicates that the sign of the sample covariance between the instrument and the structural error $\widehat{cov}(z, u)$ completely determines whether $\hat{\beta}_{2SLS}$ lies above or below the true β .³⁰ Thus, if $\widehat{cov}(z, u)$ is the same (opposite) sign as the OLS bias, the 2SLS estimate is shifted towards (away from) the OLS estimate. The power asymmetry of the 2SLS t -tests arises because it is also true that when $\widehat{cov}(z, u)$ is the same (opposite) sign as the OLS bias the 2SLS standard error is smaller (larger).

To see this clearly, we write $\hat{\beta}_{2SLS} - \beta$ in the instructive form:

$$(A6) \quad \hat{\beta}_{2SLS} - \beta = \frac{\widehat{cov}(z, u)}{\widehat{cov}(z, x)} = \frac{\widehat{cov}(z, u)}{\pi \widehat{var}(z) + \rho \widehat{cov}(z, u) + \sqrt{1 - \rho^2} \widehat{cov}(z, \eta)}$$

The fact that $\widehat{cov}(z, u)$ appears in both numerator and denominator creates the association between 2SLS estimates and standard errors. A positive sample realization of $\rho \widehat{cov}(z, u)$ generates an estimate shifted towards OLS. It also generates a low standard error because a large $\rho \widehat{cov}(z, u)$ leads to a large $\widehat{cov}(z, x)$. (Note: since $Var(\hat{\beta}_{2SLS}) = Var(\hat{\beta}_{OLS})/R_{z,x}^2$, where $R_{z,x}^2$ is first-stage R^2 , the larger is $\widehat{cov}(z, x)$ the smaller is the standard error.) Hence, 2SLS will appear spuriously precise in samples where the estimated coefficient is most shifted towards OLS. Conversely, estimates shifted away from OLS appear spuriously imprecise. This

³⁰It follows that 2SLS is approximately median unbiased provided the instrument is strong enough that an incorrect first-stage sign a rare event, as only such events impart median bias.

generates the power asymmetry in the 2SLS t -test.

For instance, in our Frisch example, the OLS bias is negative, as $cov(z, x) < 0$. Hence a negative sample realization of $\widehat{cov}(z, u)$ leads to a 2SLS estimate shifted towards OLS. It also inflates the sample covariance of the instrument with the endogenous variable, $\widehat{cov}(z, x)$, leading to a spuriously small standard error. Conversely, in samples where $\widehat{cov}(z, u)$ is positive, we obtain a relatively large Frisch elasticity estimate, and $\widehat{cov}(z, x)$ is reduced, leading to an inflated standard error.

A2. A Brief Digression on Weak Instrument Tests

In the remainder of this Appendix we provide additional power comparisons between the AR test and the t -test, looking at different levels of instrument strength as measured by the population $F (= \lambda^2)$. In order to explain why we focus on particular levels of F , it is useful to give an explanation of weak instrument tests. Recall that the population F is given by:

$$(A7) \quad F = N \frac{Var(z\pi)}{\sigma_e^2} = N \frac{\pi^2 \sigma_z^2}{\sigma_e^2} = N \frac{R_{z,x}^2}{1 - R_{z,x}^2}$$

A key insight of the weak IV literature is that the properties of 2SLS do not depend on N or first-stage R^2 *per se*, but only how they combine to form F . Hence, a large sample size alone is not sufficient to ensure 2SLS has an approximately normal sampling distribution. The papers by Staiger and Stock (1997) and Stock and Yogo (2005) derive properties of 2SLS at different levels of $F (= \lambda^2)$.

In any given sample we can only observe the realization $\hat{F} = N\hat{\pi}^2\hat{\sigma}_z^2/\hat{\sigma}_e^2$. The sample \hat{F} is a draw from the non-central F -distribution with non-centrality parameter F . Stock and Yogo (2005) weak IV tests give \hat{F} thresholds that provide high confidence 2SLS will have desirable properties in terms of size and bias. For instance, in the exactly identified case, a sample $\hat{F} > 8.96$ gives 95% confidence that F is at least 1.82, which guarantees, in turn, that in a worst-case scenario a two-tailed 5% t -test will reject $H_0 : \beta = 0$ at a 15% rate or less when the true β is zero. In other words, it has a worst-case size distortion of no more than 10%. The worst case corresponds to very severe endogeneity (i.e., $\rho = -1$ or 1).

In this calculation, the value of 1.82 is obtained numerically by setting $\beta=0$ and $\rho = \pm 1$ in the power expression in equation (A3), and doing a grid search for the level of λ^2 that sets power approximately equal to 15%. (Of course size is simply power against $H_0 : \beta=0$ evaluated at a true β equal to zero).³¹ Both Keane and Neal (2021) and Angrist and Kolesár (2021) have criticized the focus on the worst case scenario of $\rho=\pm 1$, arguing that for most plausible levels of endogeneity the size distortion is much less. But in addition, we also criticize the exclusive focus on size (or bias) to the neglect of power.

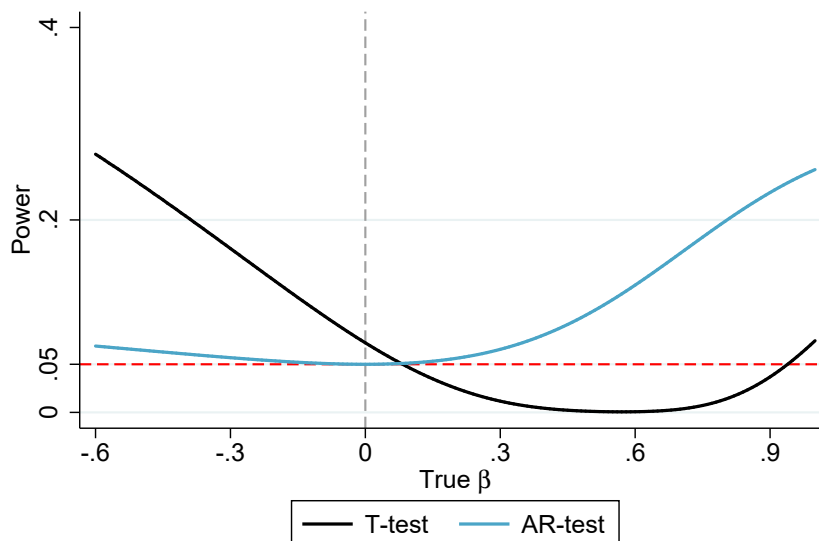
³¹Thus, the figure of 1.82 is subject to numerical error, and should not be viewed as exact. The same is true of other weak IV thresholds in this literature.

A3. Power of the AR vs t -test at Different Levels of Instrument Strength

A critical point we have emphasized is that weak IV tests do not reveal if the estimator has acceptable power properties. We now present some comparisons of the power of the AR and t -tests at different levels of F .

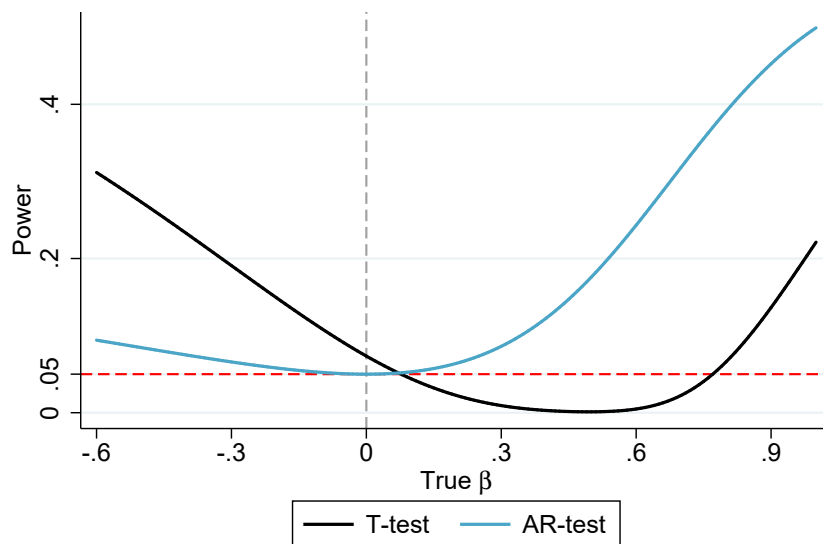
Figure A1 presents the case of $F=1$, which is indicative of the poor instrument strength in classic studies of the Frisch elasticity. We see the 2SLS t -test has essentially no power to detect positive Frisch elasticities in the plausible range of 0.1 to 0.9 (power is less than the 5% size of the test throughout this range). The AR test performs better: It does have power greater than size at all levels of β except $\beta = 0$, reflecting that it is an unbiased test. But its power is still very low: It doesn't pass 20% until the elasticity exceeds 0.8. This reflects the fact that the data is simply not very informative at this low level of instrument strength.

FIGURE A1. POWER OF THE T-TEST VS. AR-TEST WHEN POP. $F = 1$ ($\rho = -0.7$)



Another notable feature of Figure A1 is that the t -test appears to have much better power than the AR test for negative values of true β . This reveals the flip side of the power asymmetry problem: In samples where the 2SLS estimate is shifted in the direction of OLS, which in this case means it is shifted in the negative direction, the 2SLS standard error is spuriously small, which inflates the power of the t -test. This is not a desirable property, as the standard error exaggerates the precision of the estimate in such cases.

Figure A2 considers the case where $F=2.3$. This case is particularly interesting, as a sample \hat{F} of 10 is required to give 95% confidence that F is at least 2.3. Thus, this case corresponds to the widely used Staiger-Stock rule of thumb, that

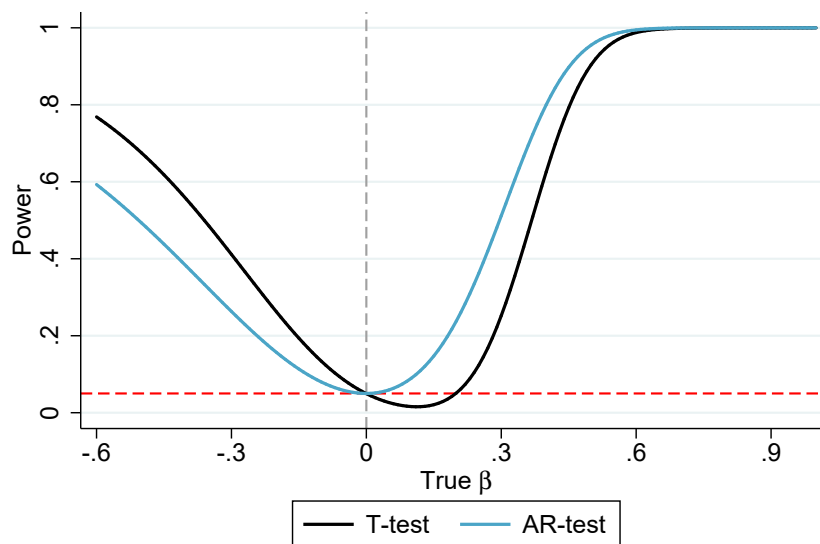
FIGURE A2. POWER OF THE T-TEST VS. AR-TEST WHEN POP. $F = 2.3$ ($\rho = -0.7$)

a first-stage \hat{F} of at least 10 indicates an acceptable level of instrument strength. However, Figure A2 reveals that the power of the 2SLS t -test is still very poor in this case. It has essentially no power to detect positive Frisch elasticities in the plausible range of 0.1 to 0.8, as power is less than the 5% size of the test throughout this range. Again the AR test has much better power to detect a true positive Frisch elasticity, but its power is still rather low (e.g., it doesn't pass 20% until the elasticity exceeds 0.5). So while this level of instrument strength is deemed acceptable by common practice, the data is not very informative.

The severe bias of the t -test is also evident in Figure A2. Power is minimized in the vicinity of $\beta=0.50$ rather than at $\beta=0$. This again reflects the power asymmetry of the t -test, and the fact that it has little power to detect a wide range of plausible positive elasticities because the OLS bias is negative.

Figure 3 in the main text considers the case of $F=10.12$, which is the value we used in the Monte Carlo exercise in Section V. A first-stage \hat{F} of at least 23.2 is required to have 95% that population F is at least 10.12. In this case the t -test has power less than size for true effects in the 0.05 to 0.50 range, so, as we note in the text, it is uninformative over that range. The power of the AR test is far superior, reaching about 60% when true elasticity is 0.50.

An obvious question is how large F must be for the t -test to begin to exhibit acceptable power for plausible elasticity values. Figure A3 reports results for a population F of 29.44. A first-stage \hat{F} of at least 50 is required to have 95% confidence that population F is at least this large. At this level of instrument strength the power of both tests approaches one when the true elasticity approaches 0.6

FIGURE A3. POWER OF THE T-TEST VS. AR-TEST WHEN POP. $F = 29.44$ ($\rho = -0.7$)

range. However, the t -test still has power less than size for elasticities in the 0.0 to 0.2 range, and very poor power compared to the t -test for elasticities in the 0.0 to 0.4 range. The size of the t -test (i.e., power at $\beta = 0$) is 4.96% in this case, so it is very close to the correct 5%. But bias is still evident as power is minimized at an elasticity of roughly $\beta = 0.1$.

Finally, Figure A4 considers the case of true $F=73.75$, which is quite a high level of instrument strength. A first-stage \hat{F} of at least 104.7 is required to have 95% confidence that population F is at least this large. We choose to examine this case because Lee et al. (2020) show that a first-stage sample \hat{F} of at least 104.7 is required for the worst-case size distortion in the t -test to be no more than 5%. Their analysis is subtly different from Stock and Yogo (2005), in that their “worst case” refers to the maximum size distortion over all possible values of endogeneity ρ and all possible values of the true F . The worst case scenario for ρ is again ± 1 , while the worst case for F is $[\hat{F}/(\sqrt{\hat{F}} + 1.96)]^2$.

At this high level of instrument strength the power curves of the two tests are much more similar, and power of both tests approaches 1 for β around 0.40. But the power advantage of the AR test is still evident in the $\beta \in (0.0, 0.30)$ range. For instance, for $\beta=0.15$ the AR test power is 40% vs. 25% for the t -test.

In summary, our results clearly show that the power advantage of the AR test over the t -test is substantial at empirically relevant elasticity values. The power asymmetry of the t -test (i.e., its low power to detect plausible positive elasticities) is dramatic when instruments are weak but persists even when instruments are very strong.

FIGURE A4. POWER OF THE T-TEST VS. AR-TEST WHEN POP. $F = 73.75$ ($\rho = -0.7$)