

Lecture Notes on Advanced Econometrics

Lecture 8: Instrumental Variables Estimation

Endogenous Variables

Consider a population model:

$$y_{i1} = \alpha_1 y_{i2} + \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$$

We call y_2 an **endogenous** variable when y_2 is correlated with u . As we have studied earlier, y_2 would be correlated with u if (a) there are omitted variables that are correlated with y_2 and y_1 , (b) y_2 is measured with errors, and (c) y_1 and y_2 are simultaneously determined (we will cover this issue later). All of these problems, we can identify the source of the problems as the correlation between the error term and one or some of the independent variables.

For all of these problems, we can apply **instrumental variables (IV)** estimations because instrumental variables are used to cut correlations between the error term and independent variables. To conduct IV estimations, we need to have instrumental variables (or instruments in short) that are **(R1) uncorrelated with u** but **(R2) partially and sufficiently strongly correlated with y_2 once the other independent variables are controlled for**.

It turns out that it is very difficult to find proper instruments!

In practice, we can test the second requirement (b), but we can not test the first requirement (a) because u is unobservable. To test the second requirement (b), we need to express a **reduced form** equation of y_2 with all of **exogenous** variables. Exogenous variables include all of independent variables that are not correlated with the error term and the instrumental variable, z . The reduced form equation for y_2 is

$$y_2 = \delta_z z + \delta_1 + \delta_2 x_2 + \dots + \delta_{k-1} x_{k-1} + u$$

For the instrumental variable to satisfy the second requirement (R2), the estimated coefficient of z must be significant.

In this case, we have one endogenous variable and one instrumental variable. When we have the same number of endogenous and instrumental variables, we say the endogenous variables are **just identified**. When we have more instrumental variables than endogenous variables, we say the endogenous variables are **over-identified**. In this case,

we need to use “two stage least squares” (2SLS) estimation. We will come back to 2SLS later.

Define $\mathbf{x} = (y_2, 1, x_2, \dots, x_{k-1})$ as a 1-by-k vector, $\mathbf{z} = (z, 1, x_2, \dots, x_{k-1})$ a 1-by-k vector of all exogenous variables, X as a n-by-k matrix that includes one endogenous variable and k-1 independent variables, and Z as a n-by-k matrix that include one instrumental variable and (k-1) independent variables:

$$X = \begin{bmatrix} y_{12} & 1 & x_{12} & x_{13} & \dots & x_{1k-1} \\ y_{22} & 1 & x_{22} & x_{23} & \dots & x_{2k-1} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ y_{n2} & 1 & x_{n2} & x_{n3} & \dots & x_{nk-1} \end{bmatrix}, \quad Z = \begin{bmatrix} z_1 & 1 & x_{12} & x_{13} & \dots & x_{1k-1} \\ z_2 & 1 & x_{22} & x_{23} & \dots & x_{2k-1} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ z_n & 1 & x_{n2} & x_{n3} & \dots & x_{nk-1} \end{bmatrix}.$$

The instrumental variables (IV) estimator is

$$\hat{\beta}_{IV} = (Z'X)^{-1} Z'Y$$

Notice that we can take the inverse of $Z'X$ because both Z and X are n-by-k matrices and $Z'X$ is a k-by-k matrix which has full rank, k. This indicates that there is no perfect collinearity in Z . The condition that $Z'X$ has full rank of k is called the **rank condition**.

Problems with Small property.

The consistency of the IV estimators can be shown by using the two requirements for IVs:

$$\begin{aligned} \hat{\beta}_{IV} &= (Z'X)^{-1} Z'(X\beta + u) \\ &= \beta + (Z'X)^{-1} Z'u \\ &= \beta + (Z'X/n)^{-1} Z'u/n \end{aligned}$$

From the first requirement (R1), $p \lim Z'u/n \rightarrow 0$.

From the second requirement (R2), $p \lim Z'X/n \rightarrow A$, where $A \equiv E(z'x)$.

Therefore, the IV estimator is consistent when IVs satisfy the two requirements.

A Bivariate IV model

Let's consider a simple bivariate model:

$$y_1 = \beta_0 + \beta_1 y_2 + u$$

We suspect that y_2 is an endogenous variable, $\text{cov}(y_2, u) \neq 0$. Now, consider a variable, z , which is correlated y_2 but not correlated with u : $\text{cov}(z, y_2) \neq 0$ but $\text{cov}(z, u) = 0$.

Consider $\text{cov}(z, y_1)$:

$$\text{cov}(z, y_1) = \text{cov}(z, \beta_0 + \beta_1 y_2 + u)$$

$$= \beta_0 \text{cov}(z, 1) + \beta_1 \text{cov}(z, y_2) + \text{cov}(z, u)$$

Because $\text{cov}(z, 1) = 0$ and $\text{cov}(z, u) = 0$, we find that

$$\begin{aligned} \beta_1 &= \frac{\text{cov}(z, y_1)}{\text{cov}(z, y_2)} \\ &= \frac{\sum_{i=1}^n (z_i - \bar{z})(y_{i1} - \bar{y}_1)}{\sum_{i=1}^n (z_i - \bar{z})(y_{i2} - \bar{y}_2)}, \text{ which is} \\ &= (Z'X)^{-1} Z'Y. \end{aligned}$$

Thus, we find the same conclusion as using the matrix form.

The problem in practice is the first requirement, $\text{cov}(z, u) = 0$. We can not empirically confirm this requirement because u cannot be observed. Thus, the validity of this assumption is left to economic theory or economists' common sense.

Recent studies show that even the first requirement can be problematic when the correlation between the endogenous and instrumental variables is weak. We will discuss this later.

Example 1: Card (1995), CARD.dta.

A dummy variable grew up near a 4 year collage as an IV on *educ*.

OLS

```
. reg lwage educ
```

Source	SS	df	MS			
Model	58.5153536	1	58.5153536	Number of obs =	3010	
Residual	534.126258	3008	.17756857	F(1, 3008) =	329.54	
				Prob > F =	0.0000	
				R-squared =	0.0987	
				Adj R-squared =	0.0984	
				Root MSE =	.42139	

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	.0520942	.0028697	18.153	0.000	.0464674	.057721
_cons	5.570883	.0388295	143.470	0.000	5.494748	5.647017

Correlation between *nearc4* (an IV) and *educ*

```
. reg educ nearc4
```

Source	SS	df	MS			
Model	448.604204	1	448.604204	Number of obs =	3010	
Residual	21113.4759	3008	7.01910767	F(1, 3008) =	63.91	
				Prob > F =	0.0000	
				R-squared =	0.0208	
				Adj R-squared =	0.0205	
				Root MSE =	2.6494	

	educ	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
	nearc4	.829019	.1036988	7.994	0.000	.6256913	1.032347
	_cons	12.69801	.0856416	148.269	0.000	12.53009	12.86594

Thus, *nearc4* satisfies the one of the two requirements to be a good candidate as an IV.

IV Estimation:

```
. ivreg lwage (educ=nearc4)
```

Instrumental variables (2SLS) regression

Source	SS	df	MS	Number of obs = 3010		
Model	-340.111443	1	-340.111443	F(1, 3008)	=	51.17
Residual	932.753054	3008	.310090776	Prob > F	=	0.0000
				R-squared	=	.
				Adj R-squared	=	.
Total	592.641611	3009	.196956335	Root MSE	=	.55686

	lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
	educ	.1880626	.0262913	7.153	0.000	.1365118	.2396134
	_cons	3.767472	.3488617	10.799	0.000	3.08344	4.451504

```
Instrumented: educ
Instruments: nearc4
```

Note that you can obtain the same coefficient by estimating OLS on *lwage* with the predicted *educ* (predicted by *nearc4*). However, the standard error would be incorrect. In the above IV Estimation, the standard error is already corrected.

End of Example 1

The Two-Stage Least Squares Estimation

Again, let's consider a population model:

$$y_1 = \alpha_1 y_2 + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u \quad (1)$$

where y_2 is an endogenous variable. Suppose that there are m instrumental variables. Instruments, $\mathbf{z} = (1, x_1, \dots, x_k, z_1, \dots, z_m)$, are correlated with y_2 . From the reduced form equation of y_2 with all exogenous variables (exogenous independent variables plus instruments), we have

$$y_2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \dots + \delta_k x_k + \delta_{k+1} z_1 + \dots + \delta_{k+m} z_m + \varepsilon$$

$$y_2 = \hat{y}_2 + \varepsilon$$

\hat{y}_2 is a linear projection of y_2 with all exogenous variables. Because \hat{y}_2 is projected with all exogenous variables that are not correlated with the error term, u , in (1), \hat{y}_2 is not correlated with u , while ε is correlated with u . Thus, we can say that by estimating y_2 with all exogenous variables, we have divided into two parts: one is correlated with u and the other is not.

The projection of y_2 with Z can be written as

$$\hat{y}_2 = Z\hat{\delta} = Z(Z'Z)^{-1}Z'y_2$$

When we use the two-step procedure (as we discuss later), we use this \hat{y}_2 in the place of y_2 . But now, we treat y_2 as a variable in X and project X itself with Z :

$$\hat{X} = Z\hat{\Pi} = Z(Z'Z)^{-1}Z'X = P_Z X$$

$\hat{\Pi}$ is a $(k+m-1)$ -by- k matrix with coefficients, which should look like:

$$\hat{\Pi} = \begin{bmatrix} \delta_1 & 1 & 0 & 0 \\ \delta_2 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ \delta_{k+m-1} & 0 & 0 & 0 \end{bmatrix}.$$

Thus, y_2 in X should be expressed as a linear projection, and other independent variables in X should be expressed by itself. $P_Z = Z(Z'Z)^{-1}Z'$ is a n -by- n symmetric matrix and idempotent (i.e., $P_Z'P_Z = P_Z$). We use \hat{X} as instruments for X and apply the IV estimation as in

$$\begin{aligned} \hat{\beta}_{2SLS} &= (\hat{X}'X)^{-1}\hat{X}'Y \\ &= (X'P_Z X)^{-1}X'P_Z Y \\ &= (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y \end{aligned} \quad (2)$$

This can be also written as

$$\hat{\beta}_{2SLS} = (\hat{X}'\hat{X})^{-1}\hat{X}'Y$$

This is the 2SLS estimator. It is called as two-stage because it looks like we take two steps by creating projected X to estimate the 2SLS estimators. We do not need to take two steps as we show in (2). We can just estimate 2SLS estimators in one step by using X and Z . (This is what econometrics packages do.)

The Two-Step procedure

It is still a good idea to know how to estimate the 2SLS estimators by a two-step procedure:

Step 1: Obtain \hat{y}_2 by estimating an OLS against all of exogenous variables, including all of instruments (**the first-stage regression**)

Step 2: Use \hat{y}_2 in the place of y_2 to estimate y_1 against \hat{y}_2 and all of exogenous independent variables, not instruments (**the second stage regression**)

The estimated coefficients from the two-step procedure should exactly the same as 2SLS. **However, you must be aware that the standard errors from the two-step procedure are incorrect, usually smaller than the correct ones.** Thus, in practice,

avoid using predicted variables as much as you can!

Econometric packages will provide you 2SLS results based on (2). So you do not need to use the two-step procedure.

We use the first step procedure to test the second requirement for IVs. In the first stage regression, we should conduct a F-test on all instruments to see if instruments are jointly significant in the endogenous variable, y_2 . As we discuss later, instruments should be strongly correlated with y_2 to have reliable 2SLS estimators.

Consistency of 2SLS

Assumption 2SLS.1: For vector \mathbf{z} , $E(\mathbf{z}'u)=0$,

where $\mathbf{z} = (1, x_1, \dots, x_k, z_1, \dots, z_m)$.

Assumption 2SLS.2: (a) $\text{rank } E(\mathbf{z}'\mathbf{z}) = k+m+1$; (b) $\text{rank } E(\mathbf{z}'\mathbf{x}) = k$.

(b) is **the rank condition** for identification that \mathbf{z} is sufficiently linearly related to \mathbf{x} so that $\text{rank } E(\mathbf{z}'\mathbf{x})$ has full column rank.

The order condition is $k-1+m \geq k-1+h$, where h is the number of endogenous variable. Thus, the order condition indicates that we must have at least as many instruments as endogenous variables.

Under assumption 2SLS1 and 2SLS2, the 2SLS estimators in (10-5) are consistent.

Under homoskedasticity,

$$\hat{\sigma}^2 (\hat{X}'\hat{X})^{-1} = (n-k)^{-1} \sum_{i=1}^n \hat{u}_i (\hat{X}'\hat{X})^{-1}$$

is a valid estimator of the asymptotic variance of $\hat{\beta}_{2SLS}$.

Under heteroskedasticity, the heteroskedasticity-robust (White) standard errors is

$$(\hat{X}'\hat{X})^{-1} \hat{X}'\hat{\Sigma}\hat{X}(\hat{X}'\hat{X})^{-1}$$

Here are some tests for IV estimation. See Wooldridge (Introductory Econometrics) for details.

Testing for Endogeneity

- (i) Estimate the reduced form model using the endogenous variable as the dependent variable:

$$y_2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \dots + \delta_k x_k + \delta_{k+1} z_1 + \dots + \delta_{k+m} z_m + \varepsilon$$
- (ii) Obtain the residual, $\hat{\varepsilon}$.
- (iii) Estimate

$$y_1 = \alpha_1 y_2 + \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \delta \hat{\varepsilon} + u$$
- (iv) If $\hat{\delta}$ is significant, then y_2 is endogenous.

Testing the Over-Identification

- (i) Estimate $\hat{\beta}_{2SLS}$ and obtain \hat{u} .
- (ii) Regress \hat{u} on z and x .
- (iii) Get R^2 and get nR^2 , which is chi-squared. If this is significantly different from zero, then at least some of the IVs are not exogenous.

Example 2: Card (1995), card.dta again.

OLS

```
. reg lwage educ exper expersq black smsa south
```

Source	SS	df	MS	
Model	172.165615	6	28.6942691	Number of obs = 3010
Residual	420.475997	3003	.140018647	F(6, 3003) = 204.93
Total	592.641611	3009	.196956335	Prob > F = 0.0000
				R-squared = 0.2905
				Adj R-squared = 0.2891
				Root MSE = .37419

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
educ	.074009	.0035054	21.113	0.000	.0671357 .0808823
exper	.0835958	.0066478	12.575	0.000	.0705612 .0966305
expersq	-.0022409	.0003178	-7.050	0.000	-.0028641 -.0016177
black	-.1896316	.0176266	-10.758	0.000	-.2241929 -.1550702
smsa	.161423	.0155733	10.365	0.000	.1308876 .1919583
south	-.1248615	.0151182	-8.259	0.000	-.1545046 -.0952184
_cons	4.733664	.0676026	70.022	0.000	4.601112 4.866217

IV Estimation: nearc2 nearc4 as IVs

```
. ivreg lwage (educ= nearc2 nearc4) exper expersq black smsa south
```

Instrumental variables (2SLS) regression

Source	SS	df	MS	Number of obs =	3010
Model	86.2368644	6	14.3728107	F(6, 3003) =	110.30
Residual	506.404747	3003	.168632949	Prob > F =	0.0000
				R-squared =	0.1455
				Adj R-squared =	0.1438
Total	592.641611	3009	.196956335	Root MSE =	.41065

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	.1608487	.0486291	3.308	0.001	.065499	.2561983
exper	.1192111	.0211779	5.629	0.000	.0776865	.1607358
expersq	-.0023052	.0003507	-6.574	0.000	-.0029928	-.0016177
black	-.1019727	.0526187	-1.938	0.053	-.205145	.0011996
smsa	.1165736	.0303135	3.846	0.000	.0571363	.1760109
south	-.0951187	.0234721	-4.052	0.000	-.1411418	-.0490956
_cons	3.272103	.8192562	3.994	0.000	1.665743	4.878463

```
Instrumented: educ
Instruments: nearc2 nearc4 + exper expersq ... south
```

```
. ivreg lwage (educ= nearc2 nearc4 fatheduc motheduc) exper expersq black sms
> a south
```

Instrumental variables (2SLS) regression

Source	SS	df	MS	Number of obs =	2220
Model	108.419483	6	18.0699139	F(6, 2213) =	83.66
Residual	320.580001	2213	.144862178	Prob > F =	0.0000
				R-squared =	0.2527
				Adj R-squared =	0.2507
Total	428.999484	2219	.193330097	Root MSE =	.38061

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	.1000713	.01263	7.923	0.000	.0753034	.1248392
exper	.0989441	.009482	10.435	0.000	.0803496	.1175385
expersq	-.002449	.0004013	-6.103	0.000	-.0032359	-.0016621
black	-.1504635	.0259113	-5.807	0.000	-.2012765	-.0996505
smsa	.150854	.0195975	7.698	0.000	.1124226	.1892854
south	-.1072406	.0180661	-5.936	0.000	-.1426688	-.0718123
_cons	4.26178	.216812	19.657	0.000	3.836604	4.686956

```
Instrumented: educ
Instruments: nearc2 nearc4 fatheduc motheduc + exper expersq ... south
```


IV Estimation: *fatheduc* *motheduc* as IVs

```
. ivreg lwage (educ= fatheduc motheduc) exper expersq black smsa south
Instrumental variables (2SLS) regression
```

Source	SS	df	MS	Number of obs = 2220		
Model	108.477154	6	18.0795257	F(6, 2213)	=	83.44
Residual	320.52233	2213	.144836118	Prob > F	=	0.0000
				R-squared	=	0.2529
				Adj R-squared	=	0.2508
				Root MSE	=	.38057

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
educ	.099931	.012756	7.834	0.000	.0749161	.124946
exper	.098884	.0095123	10.395	0.000	.0802299	.117538
expersq	-.0024487	.0004012	-6.103	0.000	-.0032356	-.0016619
black	-.1505902	.0259598	-5.801	0.000	-.2014983	-.0996822
smsa	.1509271	.0196181	7.693	0.000	.1124553	.1893988
south	-.1072797	.0180714	-5.936	0.000	-.1427183	-.071841
_cons	4.26415	.2189075	19.479	0.000	3.834865	4.693436

```
Instrumented: educ
Instruments: fatheduc motheduc + exper expersq ... south
```

Which ones are better?

End of Example 2

Weak Instruments

Remember the two requirements for instrumental variables to be valid: **(R1) uncorrelated with u** but **(R2) partially and sufficiently strongly correlated with y_2 once the other independent variables are controlled for.**

The first requirement is difficult to be confirmed because we cannot observe u . So, we rely on economic theory or natural experiments to find instrumental variables that satisfy the first requirement. The second requirement can be checked by conducting some analyses.

One way of checking the second requirement is to estimate a regression model of endogenous variables on exogenous variables which include instrumental variables and other exogenous variables. Suppose that we have one endogenous variable, y_2 , and m instrumental variables, z_1, \dots, z_m . Then estimate the following model:

$$y_{i2} = \beta_0 + \delta_1 z_{i1} + \dots + \delta_m z_{im} + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i. \quad (1)$$

Then obtain a F-statistics on the estimators of the instrumental variables:

$H_0 : \delta_1 = \dots = \delta_m = 0$. If the F-statistics is small, then we call the instrumental variables *weak*. When the instrumental variables are weak, the IV or 2SLS estimators could be inconsistent or have large standard errors.

Inconsistency

To examine how weak instruments can make IV and 2SLS estimators inconsistent, let us consider a simple bivariate model:

$$y_1 = \beta_0 + \beta_1 y_2 + u$$

In a bivariate model, we write

$$\text{plim } \hat{\beta}_{IV} = \beta + \frac{\text{cov}(z, u)}{\text{cov}(z, y_2)}$$

because $\text{Corr}(z, u) = \text{cov}(z, u) / [\text{sd}(z)\text{sd}(u)]$ (see Wooldridge pp714)

$$\text{plim } \hat{\beta}_{IV} = \beta + \frac{\text{corr}(z, u) / (\text{sd}(z)\text{sd}(u))}{\text{corr}(z, y_2) / (\text{sd}(z)\text{sd}(y_2))}$$

$$\text{plim } \hat{\beta}_{IV} = \beta + \frac{\text{corr}(z, u) \text{sd}(y_2)}{\text{corr}(z, y_2) \text{sd}(u)}$$

Thus if z is only weakly correlated with the endogenous variable, y_2 , i.e., $\text{corr}(z, y_2)$ is very small, the IV estimator could be severely biased even when the correlation between z and u is very small.

Consider now the OLS estimator of the bivariate model ignoring the endogeneity of y_2 :

$$\begin{aligned} \text{plim } \hat{\beta}_{OLS} &= \beta + \frac{\text{cov}(y_2, u)}{\text{var}(y_2)} \\ &= \beta + \frac{\text{corr}(y_2, u) \text{sd}(y_2) \text{sd}(u)}{\text{sd}(y_2)^2} \\ &= \beta + \frac{\text{corr}(y_2, u) \text{sd}(u)}{\text{sd}(y_2)} \end{aligned}$$

Here we have the endogenous bias. But the size of the endogenous bias could be smaller than the bias created by the weak instrument. To compare these two, let us take the ratio of these two:

$$\frac{\text{plim } \hat{\beta}_{IV} - \beta}{\text{plim } \hat{\beta}_{OLS} - \beta} = \frac{\text{corr}(z, u)}{\text{corr}(z, y_2)} \frac{1}{\text{corr}(y_2, u)}$$

When this ratio is larger than 1, the bias in the IV estimator is larger than the OLS estimator. When the instrumental variable is completely uncorrelated with u , the bias in

the IV estimator is zero. When it has even a very small correlation with u , the size of the ratio depends on the size of the denominator, which could be very small when the correlation between the instrumental variable and the endogenous variable, $corr(z, y_2)$, is small. The implication from this simple model could be also applied on a more complicated model where there are more than one endogenous variable and one instrumental variable.

Low Precision

Weak instrumental variables can lead to large standard errors of the IV/2SLS estimators. The variance of the IV estimator is

$$V(\hat{\beta}_{IV}) = \sigma^2 (Z'X)^{-1} Z'Z(Z'X)^{-1}.$$

This could be rearranged as:

$$\begin{aligned} V(\hat{\beta}_{IV}) &= \sigma^2 (X'X)^{-1} X'X (Z'X)^{-1} Z'Z(Z'X)^{-1} \\ &= V(\hat{\beta}_{OLS}) X'X (Z'X)^{-1} Z'Z(Z'X)^{-1} \\ &= V(\hat{\beta}_{OLS}) [(X'X)^{-1} Z'X]^{-1} [(Z'Z)^{-1} Z'X]^{-1} \\ &= V(\hat{\beta}_{OLS}) [\hat{\Pi}_{Z,X}]^{-1} [\hat{\Pi}_{X,Z}]^{-1} \end{aligned}$$

where $\hat{\Pi}_{Z,X}$ and $\hat{\Pi}_{X,Z}$ are projections of Z on X and X on Z , respectively. If the correlations between Z and X are low, then $\hat{\Pi}_{Z,X}$ and $\hat{\Pi}_{X,Z}$ have low values, which would make the variance of the IV estimators large. This could be applied on the 2SLS estimators.

In empirical estimation, we often find large standard errors in the IV/2SLS estimators. This could be caused by weak instruments.

Thus, weak instrumental variables can cause inconsistency and imprecision in the IV/2SLS estimators. But how weak is weak?

A Rule of Thumb to find Weak Instruments

Staiger and Stock (1997) suggest that the F-statistics of instrumental variables in (1), of this lecture note, should be larger than 10 to ensure that the maximum bias in IV estimators to be less than 10 %. If you are willing to accept the maximum bias in IV estimators to be less than 20 %, the threshold is F-stat being larger than 5. If the number of instrumental variables is one, the F-statistics should be replaced by the t-statistics.

In practice, it is quite difficult to find valid instrumental variables that are not weak. I personally find searching for instrumental variables time-consuming and not so rewarding. One can never be sure about validness of IVs. You should look for natural experiments or randomized experiments that could be used as instrumental variables.

