

Introduction to Econometrics: Recitation 8

Seung-hun Lee*

November 3rd, 2022

1 Maximum Likelihood Estimation Method

Notice that both logit and probit regressions are nonlinear in the sense that the β_0, β_1 parameters are no longer in linear relationship with the X_i 's and subsequently Y_i 's. One of the assumptions used in using OLS is that the linear regression assumption. This is no longer a valid option anymore, which requires a different approach. This is where **maximum likelihood estimation** comes in. To understand the maximum likelihood estimators, you must understand what the likelihood function is. A **likelihood function** is the conditional density of Y_1, \dots, Y_n given X_1, \dots, X_n that is treated as the function of the unknown parameters (β_0, β_1 in our case). In other words, since we have the observations for Y_i 's and X_i 's, but do not know the values of the parameter β_i 's, what we are trying to do here is to find the values of β_i 's that best matches the values of X_i 's and Y_i 's. As a result, the maximum likelihood estimators is the value of β_i 's that best describes the data and maximizes the value of the likelihood function

To nail this home, let's not worry about regression equation for the moment and consider a single variable - Y_i . For example, Let's assume that Y_i 's are IID normal with mean μ and standard error σ , both of which are unknown. The joint probability of Y_i 's are (our likelihood

*Contact me at sl4436@columbia.edu if you spot any errors or have suggestions on improving this note.

function)

$$\begin{aligned}
 \Pr(Y_1 = y_1, \dots, Y_n = y_n | \mu, \sigma) &= \Pr(Y_1 = y_1 | \mu, \sigma) \times \dots \times \Pr(Y_n = y_n | \mu, \sigma) \\
 &= \prod_{i=1}^n f(y_i | \mu, \sigma) \\
 &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y_i - \mu)^2}{2\sigma^2}} \\
 &= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{-\frac{n}{2}} e^{-\sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2}}
 \end{aligned}$$

A convenient way to solve this class of problem is to use the *log-likelihood function*. Taking logs to above equation gets us

$$-\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^2} \quad (*)$$

To find the maximum likelihood estimator of μ , we differentiate the above with respect to μ . This gets us

$$\begin{aligned}
 2 \sum_{i=1}^n \frac{(Y_i - \mu)}{2\sigma^2} &= 0 \implies \sum_{i=1}^n \frac{(Y_i - \mu)}{\sigma^2} = 0 \\
 (\because \sigma, \text{ though unknown, will be a constant}) &\implies \sum_{i=1}^n (Y_i - \mu) = 0 \\
 &\implies n\mu = \sum_{i=1}^n Y_i \implies \mu_{MLE} = \frac{1}{n} \sum_{i=1}^n Y_i
 \end{aligned}$$

We can do the similar for σ^2 by differentiating $(*)$ with respect to σ^2 . This leads to

$$\begin{aligned}
 -\frac{n}{2\sigma^2} + \sum_{i=1}^n \frac{2(Y_i - \mu)^2}{(2\sigma^2)^2} &= 0 \implies \frac{n}{2\sigma^2} = \sum_{i=1}^n \frac{(Y_i - \mu)^2}{2\sigma^4} \\
 &\implies n = \sum_{i=1}^n \frac{(Y_i - \mu)^2}{\sigma^2} \\
 &\implies n = \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \mu)^2
 \end{aligned}$$

By imputing μ_{MLE} in place of μ , we get

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \mu_{MLE})^2$$

2 Instrumental Variable Regression

2.1 Motivation

Recall that when there is a possibility of three of the internal validity threat factors - namely omitted variable bias, measurement error, and simultaneity bias - then we run into the case where the OLS estimation gives us biased estimates. Specifically, these cases bias our result because the error term u is no longer expected to be 0 conditional on X ($E[u|x] \neq 0$).

This is where **instrumental variables regression** comes in handy. Instrumental variables allow us to eliminate bias in the following sense: When you have an independent variable X , there are parts of this variable that are correlated with u and the other parts that are independent of u . If we are able to find Z that is correlated with X but not with u , using this Z variable would allow you to sort variable X into what is correlated with u and what is not. Then, using the part that is not correlated with u , variable Z allows us to get unbiased estimates.

This approach can be very useful in so many settings. Most of our independent variables are **endogenous** in the sense that X is not a given, but determined as a result of *choice* of individuals, along with dependent variable Y . For instance, If Y variable is test score and X is time spent studying, it is possible that the estimates of the effect of time spent studying can be biased. Think about the case where you see your friend (or rival, whichever you prefer) get a really high score. After you know that this person spends a lot of time studying, you decide to study more. Effectively, you see your friend's Y and that affects your selection of X - simultaneity bias. Alternatively, there might be factors uncontrolled by this regression affecting your choice of study hours - maybe hours spent commuting to schools. Lastly, the measure of X might be inaccurate - for instance, there might be a case where you think you are studying but in reality, it is not. These cases can be addressed by finding some variables that are related with time spent studying but is not correlated with the error term.

2.2 Recapping associations vs causality

We use the term **association** to refer to any type of statistical dependence across variables. While this association includes causation, that is a rare case. Most association we talk about are correlations, which is

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)}\sqrt{\text{var}(y)}}$$

We now demonstrate that obtaining correlation between y and x is practically equivalent to running an OLS. Consider a population OLS where we minimize the *expected value* of the squared sum of errors. Thus, we solve

$$\min E[(Y - \beta_0 - \beta_1 X)^2]$$

whose first order conditions are

$$\begin{aligned} E[Y - \beta_0 - \beta_1 X] &= 0 \\ E[X(Y - \beta_0 - \beta_1 X)] &= 0 \end{aligned}$$

After solving for this, we get

$$\beta_1 = \frac{E[X - E[X]]E[Y - E[Y]]}{(E[X - E[X]])^2} = \text{corr}(X, Y) \sqrt{\frac{\text{var}(Y)}{\text{var}(X)}}$$

so β_1 is obtained by multiplying the correlation coefficient with a constant. Since variances are nonnegative, OLS coefficients and the correlations follow the same direction. While this statement is nice on its own, you end up getting the same result (at least the signs) when you switch X and Y . This is not nice for identifying causal effect of X on Y .

2.3 IV conditions

Unfortunately, not every IV is found easily (in fact, you might be able to rebut my examples two subsections ago). Finding variables that qualify as IV is challenging. The conditions they should satisfy are

Relevance : Variable Z satisfies relevancy condition if $\text{cov}(X, Z) \neq 0$

Exogeneity : Variable Z satisfies exogeneity condition if $\text{cov}(Z, u) = 0$

In words, variable Z should be somewhat correlated with the original independent variable X . Additionally, variable Z should not be correlated with u . The last part seems a bit challenging to conceptualize. So I will introduce another way of thinking about this condition.

Exclusion : Variable Z satisfies exclusion condition if it affects Y only through X . In other words, when X is controlled for, Z alone should not affect Y .

So how does this condition come in? Suppose that your model is (I skip subscript i for convenience) $Y = \beta_0 + \beta_1 X + u$. Then $cov(Z, u) = 0$ can be re-written as

$$\begin{aligned} cov(Z, u) &= cov(Z, Y - \beta_0 - \beta_1 X) \\ &= cov(Z, Y) - cov(Z, \beta_0) - cov(Z, \beta_1 X) \end{aligned}$$

By imposing $cov(Z, u) = 0$ and the fact that the covariance of a random variable with a fixed constant is 0, I can get that

$$cov(Z, Y) = cov(Z, \beta_1 X)$$

This condition means that Z is correlated with Y *only through* X . If there is something else other than $cov(Z, \beta_1 X)$ in the right hand side, then this is a violation of the exclusion condition because there is a covariance between Z and Y that is not explained by X . If this happens, the IV estimation fails to be accurate.

2.4 IV estimation

There are many ways to estimate IV. I will be introducing **Two stage least squares**, **covariance method**, and **derivations from the reduced form**.

2.4.1 Two stage least squares

From the discussion in the beginning, IV estimation allows us to separate X into the part that is correlated with u and uncorrelated with u . This method explicitly shows how that is done. The steps for carrying out this approach is as follows:

1. **Regress with X as dependent, Z as independent variable.** You will get an equation that looks like

$$X = \delta_0 + \delta_1 Z + v$$

From this regression, obtain the predicted values of X , denoted as $\hat{X} = \hat{\delta}_0 + \hat{\delta}_1 Z$. This \hat{X} is the part that is related with Z but is uncorrelated with u .

2. **Regress with Y as dependent, \hat{X} as independent variable.** Regressing with this \hat{X} will satisfy the $E[u|\hat{X}] = 0$ condition, as the \hat{X} is uncorrelated with u , as obtained in the previous step. Thus, your regression equation looks like this:

$$Y = \beta_0 + \beta_1 \hat{X} + u$$

Then you run a OLS regression on the above equation and get the two stage least squares estimator $\hat{\beta}_{\text{TSLs}}$.

There is one caveat to this method. In particular, you should be careful about what standard errors you are using. If you take the second stage standard error as is, then you end up with wrong standard errors. This is because the second stage standard error does not take into account that you obtained \hat{X} using a separate (first stage) regression. The following example should make things clear. Going back to the regression on cigarette demand in 1995, we are interested in how demand for cigarettes is affected by income and price of cigarettes. In other words,

$$\ln Q_i^d = \beta_0 + \beta_1 \ln P_i^d + u_i$$

Because quantity demanded can affect price of cigarettes, there is a simultaneity bias. To get around this, we use general sales tax per pack. The following figures demonstrate two different standard errors.

```
. regress lpackpc xhat if year==1995, vce(robust)
```

Linear regression	Number of obs	=	48
	F(1, 46)	=	10.54
	Prob > F	=	0.0022
	R-squared	=	0.1525
	Root MSE	=	.22645

lpackpc	Robust		t	P> t	[95% Conf. Interval]	
	Coef.	Std. Err.				
xhat	-1.083586	.3336948	-3.25	0.002	-1.755279	-.4118936
_cons	9.298537	1.597118	5.82	0.000	6.083705	12.51337

```
. ** ivregress
. ivregress 2sls lpackpc (lavgprs=rstax) if year == 1995, vce(robust)
```

Instrumental variables (2SLS) regression	Number of obs	=	48
	Wald chi2(1)	=	12.05
	Prob > chi2	=	0.0005
	R-squared	=	0.4011
	Root MSE	=	.18635

lpackpc	Robust		z	P> z	[95% Conf. Interval]	
	Coef.	Std. Err.				
lavgprs	-1.083587	.3122037	-3.47	0.001	-1.695495	-.471679
_cons	9.29854	1.496144	6.22	0.000	6.366152	12.23093

Instrumented: lavgprs
Instruments: rstax

Notice that the standard errors are different. Therefore, if you want to work with correct standard errors, you must use `ivregress 2sls y x1 (x2 = z)` command.

2.4.2 Covariance method

Go back to the part where we discussed the exclusion condition. From there, get

$$\text{cov}(Z, Y) = \text{cov}(Z, \beta_1 X) \implies \text{cov}(Z, Y) = \beta_1 \text{cov}(Z, X)$$

From this, we can get

$$\beta_1 = \frac{\text{cov}(Z, Y)}{\text{cov}(Z, X)}$$

where division is possible because we require a relevancy condition ($\text{cov}(Z, X) \neq 0$). If you replace the population covariances with sample covariances, we get the IV estimator

$$\hat{\beta}_{TSLS} = \frac{s_{ZY}}{s_{ZX}}$$

2.4.3 Derivations from the reduced form

This requires some amount of algebra. Denote the two equations as

$$X = \pi_0 + \pi_1 Z + v \text{ (where } \text{cov}(Z, v) = 0 \text{)}$$

$$Y = \gamma_0 + \gamma_1 Z + w \text{ (where } \text{cov}(Z, w) = 0 \text{)}$$

When you rewrite the first equation in terms of Z , you get

$$Z = \frac{X}{\pi_1} - \frac{\pi_0}{\pi_1} - \frac{v}{\pi_1}$$

Then plug this into the second equation. Reorganizing this equation, you should get

$$Y = \left(\gamma_0 - \frac{\pi_0 \gamma_1}{\pi_1} \right) + \left(\frac{\gamma_1}{\pi_1} \right) X + \left(w - \frac{\gamma_1}{\pi_1} v \right)$$

What this gets you is that β_1 from the equation where we had X as independent variable is $\beta_1 = \frac{\gamma_1}{\pi_1}$. This also implies that when Z rises by 1 unit, X exogenously changes by π_1 units and Y changes by γ_1 as a result. In other words, 1 unit of exogenous change in X translates to $\beta_1 = \frac{\gamma_1}{\pi_1}$ units of change in Y .

2.5 Properties of IV estimators

Note that the IV estimator can be written in terms of the ratios of the two covariances. In sample analogs, it is written as

$$\hat{\beta}_{1,\text{TSLs}} = \frac{s_{zy}}{s_{zx}} \simeq \frac{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})(X_i - \bar{X})}$$

Note that $\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})(Y_i - \bar{Y})$ is not an unbiased estimator for the covariance of Z, Y (and similar for $\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})(X_i - \bar{X})$). If you have large samples, however, then $\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})(Y_i - \bar{Y})$ is a **consistent** for covariance of Z, Y ¹. In mathematical terms, $s_{YZ} \xrightarrow{p} \text{cov}(Y, Z), s_{XZ} \xrightarrow{p} \text{cov}(X, Z)$. With these properties, we can say that the IV estimator of β_1 , which is $\hat{\beta}_{1,\text{TSLs}}$ is a consistent estimator for the true β_1

The notion of IV estimation does not change much even in multiple regression settings. Suppose that we have

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \delta_1 W_{1i} + \dots + \delta_l W_{li} + u_i$$

where X variables are endogenous and W variables are exogenous. Assume that we have found a total of m (not necessarily, equal. We'll save the discussion to overidentification test) variables that could qualify as IVs. Additionally, IV regression requires these assumptions

IV1 $E[u_i | W_{1i}, \dots, W_{li}] = 0$ (At least for exogenous variables, this is satisfied)

IV2 $(Y_i, X_{1i}, \dots, X_{ki}, W_{1i}, \dots, W_{li}, Z_{1i}, \dots, Z_{mi})$ are IID

IV3 The Y, X, W, Z variables all have nonzero finite 4th moments

IV4 The instruments are valid. That is $\text{cov}(Z_{ji}, u_i) = 0$ for all $j = 1, \dots, m$ and relevancy conditions are satisfied for all Z 's.

2.6 Identification issues

In a general sense, a parameter is **identified** if different values of the parameter produce different distributions of the data. In other words, there is a one-to-one matching of the parameters and the distributions. For instance, A normal distribution with mean μ and

¹However, $\frac{1}{n-1} \sum_{i=1}^n (Z_i - \bar{Z})(Y_i - \bar{Y})$ is unbiased. For more information refer to <https://math.stackexchange.com/questions/2019122/unbiased-estimate-of-the-covariance>

variance σ^2 is identified because different values of μ, σ^2 produce different distributions. If it is the case that the same distribution can be obtained from different parameter values, we say that the parameters are not identified.

In the context of IV's, the parameters we are interested are coefficients for the independent variables. This depends on how much instrument variables (m in the previous section) we have relative to the endogenous variables (k). There are three possible cases

Just-identified : When $m = k$. There are just enough instruments to identify k endogenous variables

Overidentified When $m > k$. There are more than enough instruments. We later test whether the instruments are valid using **overidentification test**.

Underidentified When $m < k$. There are not enough instruments. The coefficients for X 's will not be identified

So the key takeaway here so far is that when you are using IV, you need at least as much instrumental variables as the number of endogenous regressors you have. Now the issue is having too much IV bad? That depends on the answer we get from the overidentification test, which we cover in the next section.

2.7 Implementing tests to verify that IV conditions hold

2.7.1 Underlying distribution of the IV estimators

As with other estimators, we need to understand the the distribution of IV estimators in order to conduct any sort of tests. Fortunately, IV estimators, at least in the large numbers, have normal distribution. This means there is no drastic changes in the critical values and confidence intervals we select when we conduct hypothesis tests. So why is that the case?

We can find this out by breaking down the 2SLS estimator into the following

$$\begin{aligned}
 \hat{\beta}_{\text{TSL}} &= \frac{s_{zy}}{s_{zx}} = \frac{\sum_{i=1}^n (Z_i - \bar{Z})(Y_i - \bar{Y})}{\sum_{i=1}^n (Z_i - \bar{Z})(X_i - \bar{X})} \\
 &= \frac{\sum_{i=1}^n (Z_i - \bar{Z})Y_i}{\sum_{i=1}^n (Z_i - \bar{Z})X_i} \\
 &= \frac{\sum_{i=1}^n (Z_i - \bar{Z})(\beta_0 + \beta_1 X_i + u_i)}{\sum_{i=1}^n (Z_i - \bar{Z})X_i} \\
 &= \beta_1 + \frac{\sum_{i=1}^n (Z_i - \bar{Z})u_i}{\sum_{i=1}^n (Z_i - \bar{Z})X_i} = \beta_1 + \frac{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})u_i}{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})X_i}
 \end{aligned}$$

After subtracting both sides by β_1 and multiplying both sides by \sqrt{n} , we get

$$\sqrt{n}(\hat{\beta}_{\text{TSL}} - \beta_1) = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n (Z_i - \bar{Z})u_i}{\frac{1}{n} \sum_{i=1}^n (Z_i - \bar{Z})X_i}$$

Note that the denominator converges to $\text{cov}(Z, X)$ by (weak) law of large numbers and the numerator converges to $N(0, \text{var}[(Z - \mu_Z)u])$ by central limit theorem. Thus, $\hat{\beta}_{\text{TSL}}$ has a normal distribution.

2.7.2 Confirming relevance of the instruments

We introduce a test to confirm that the instruments are really relevant. Assume we have

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 W_{1i} + \dots + \beta_{1+r} W_{ri} + u_i$$

where only X_i is endogenous. Suppose we have m instruments. Then our first stage equation looks like

$$X_i = \pi_0 + \pi_1 Z_{1i} + \dots + \pi_m Z_{mi} + \pi_{(m+1)i} W_{1i} + \dots + \pi_{(m+r)i} W_{ri} + v_i$$

We say our instrument is relevant if at least one of π_1, \dots, π_m is statistically nonzero. To confirm that this is the case, we run the F -test with these null and alternative hypothesis

$$H_0 : \pi_1 = \dots = \pi_m = 0 \text{ vs. } H_1 : \neg H_0$$

We calculate the F -statistics from this test. By rule of thumb, we want to have an F -statistics larger than 10. If otherwise, all these coefficients are zero or very close to it. In such case, we

have a **weak instrument** problem - our choice of instruments does little to explain variation in X_i . This affects our inference as the IV estimator may not have a normal distribution even with large numbers of observations.

2.7.3 Tests for exogeneity of the instruments

Here, we partially check for exogeneity of our instruments. Recall that instruments are exogenous if all of them are uncorrelated with the error term. This is how we pick up the part of X that is uncorrelated with u (the \hat{X}). If otherwise, we cannot pick up this part of X . We can *partially* test for this (more on this later) when we have more instrumental variables than the endogenous variables ($m > k$ from above section). We call this the test of overidentifying restrictions.

Consider the following case: We want to regress

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \delta_1 W_{1i} + \dots + \delta_l W_{li} + u_i$$

and we found m possible candidates for instrumental variables, Z_1, \dots, Z_m . Overidentifying restriction test we conduct here is a J -test that uses TSLS estimators and not the hypothesized β values. The steps are as follows

1. Regress the above equation using 2SLS.
2. Obtain the residuals $\hat{u}_i = Y_i - \hat{Y}_i$
3. Regress residuals onto instruments and other controls, namely

$$\hat{u}_i = \alpha_0 + \alpha_1 Z_{1i} + \dots + \alpha_m Z_{mi} + \alpha_{m+1} W_{1i} + \dots + \alpha_{m+l} W_{li} + v_i$$

4. Compute the F -statistics from the following hypothesis

$$H_0 : \alpha_1 = \dots = \alpha_m = 0 \text{ vs. } H_1 : \neg H_0$$

5. Derive the J statistics as follows $J = m \times F$. J follows χ^2_{m-k} distribution
6. If there are endogenous instruments, the α coefficient for that instrument will not be 0. This will inflate the J statistic to the extent that you reject the null of exogenous instruments. (Note that rejecting the null in this case is not something that you want!)

7. If the test is rejected, you need to make a guess on which instrument is violating the exogeneity condition, drop those, and redo the above procedure. (Good luck!)

There are some other ways to understand this. Suppose that we only have one independent variable that is also endogenous - X_i . Also suppose that we have two instrumental variables - Z_{1i}, Z_{2i} . Now, we run two two-stage least squares regression with each one of the instrumental variables. The estimates for the coefficient for X are

$$\hat{\beta}_{Z_1} = \frac{\text{cov}(Z_1, Y)}{\text{cov}(Z_1, X)}, \hat{\beta}_{Z_2} = \frac{\text{cov}(Z_2, Y)}{\text{cov}(Z_2, X)}$$

The numbers look different, although both $\hat{\beta}$'s are suppose to be the same coefficient estimating impact of X_1 on Y . The overidentification test checks whether the differences between $\hat{\beta}_{Z_1}, \hat{\beta}_{Z_2}$ are large. If the two estimates are similar in the sense that the sample analogs of both $\frac{\text{cov}(Z_1, Y)}{\text{cov}(Z_1, X)}$ and $\frac{\text{cov}(Z_2, Y)}{\text{cov}(Z_2, X)}$ converge to the same thing in probability, we don't really have to worry. If they converge to something different, then we may have problem with either one or potentially both instrumental variables. Therefore, the takeaway from overidentification test is that it checks whether there exists a faulty IV variable in the case where we found too many candidates for instrumental variables.

So what do we do in just identified case? To be very honest, there is not much we can do in terms of rigorous testing approach. This is usually the area of judgement call. One way to do this is to use the approach of exclusion restriction - you argue that the instrumental variable Z affects Y only through X . If Z affects Y directly without X , then the exclusion restriction fails. You try to persuade others based on logic, previous practices, or intuition (usually a combination of all three). This is what I meant in the beginning when I said we can partially test for exogeneity conditions.