

## Lecture 4 - agenda

1. FW for the population regression (review)
2. FW for OLS
3. an example
4. Omitted variable formula

*Frisch-Waugh theorem*

Population regression:  $y_i = x_i\beta^* + u_i$

Auxilliary regression of  $x_{ji}$  on all the other  $x's$ :

$$x_{ji} = x'_{(\sim j)i}\pi + \xi_i.$$

The  $j^{th}$  row of  $\beta^*$  is:

$$\beta_j^* = E[\xi_i^2]^{-1}E[\xi_i y_i]$$

$\beta_j^*$  = coefficient from univariate regression of  $y_i$  on  $x_{ji}$ , after “partialling out” other  $x's$ .

How does the proof work? Use FOC for  $\beta^*$  and  $\pi$ !

$$\begin{aligned} E[\xi_i y_i] &= E[\xi_i(\beta_1^* x_{1i} + \beta_2^* x_{2i} + \dots + \beta_j^* x_{ji} + \dots + \beta_K^* x_{Ki} + u_i)] \\ &= \beta_1^* E[\xi_i x_{1i}] + \beta_2^* E[\xi_i x_{2i}] + \dots + \beta_j^* E[\xi_i x_{ji}] + \dots + \beta_K^* E[\xi_i x_{Ki}] \\ &\quad + E[\xi_i u_i] \end{aligned}$$

$$\text{FOC for } \pi \Rightarrow E[x_{(\sim j)i} \xi_i] = 0 \Rightarrow E[\xi_i x_{ni}] = 0 \text{ unless } n = j$$

$$\text{FOC for } \beta^* \Rightarrow E[x_i u_i] = 0 \Rightarrow E[\xi_i u_i] = E[(x_{ji} - x'_{(\sim j)i} \pi) u_i] = 0$$

And:  $E[\xi_i x_{ji}] = E[\xi_i (x'_{(\sim j)i} \pi + \xi_i)] = E[\xi_i^2]$  using the FOC for  $\pi$  (again). So

$$E[\xi_i y_i] = \beta_j^* E[\xi_i^2] \Rightarrow \beta_j^* = E[\xi_i^2]^{-1} E[\xi_i y_i]$$

Suppose we have a constant and one other  $x$  variable:  $x'_i = (1, x_{2i})$ :

$$y_i = \beta_1^* + \beta_2^* x_{2i} + u_i$$

In this case, auxiliary regression is

$$x_{i2} = 1 \bullet \pi + \xi_i$$

And we know  $\pi = E[x_{i2}]$  (population regression = CEF). So in this case,  $\xi_i = x_{i2} - E[x_{i2}]$ . Therefore:

$$\begin{aligned}\beta_2^* &= E[(x_i - E[x_i])^2]^{-1} E[(x_i - E[x_i])y_i] \\ &= Var[x_i]^{-1} Cov[x_i, y_i]\end{aligned}$$

Now let's move from the population regression to the OLS regression. Recall the OLS objective is

$$\min_{\beta} \sum_{i=1}^N (y_i - x_i' \beta)^2$$

The FOC (which defines  $\hat{\beta}$ ) is:

$$\begin{aligned} \sum_{i=1}^N x_i (y_i - x_i' \hat{\beta}) &= 0 \quad \Rightarrow \quad \frac{1}{N} \sum_{i=1}^N x_i (y_i - x_i' \hat{\beta}) \\ \Rightarrow \frac{1}{N} \sum_{i=1}^N x_i y_i &= \left( \frac{1}{N} \sum_{i=1}^N x_i x_i' \right) \hat{\beta} \\ \Rightarrow \hat{\beta} &= \left( \frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \frac{1}{N} \sum_{i=1}^N x_i y_i \end{aligned}$$

$$\hat{\beta} = \left( \frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1} \frac{1}{N} \sum_{i=1}^N x_i y_i$$

c/w population regression:

$$\beta^* = E[x_i x_i']^{-1} E[x_i y_i]$$

repl.  $E[x_i x_i']$  with  $S_{xx} = \frac{1}{N} \sum_{i=1}^N x_i x_i'$

repl.  $E[x_i y_i]$  with  $S_{xy} = \frac{1}{N} \sum_{i=1}^N x_i y_i$

The 3 properties of the population regression are also true of the OLS regression. For the pop. regression, these come from FOC:  $E[x_i(y_i - x_i'\beta^*)] = 0$ .

For the OLS regression, these come from FOC:

$$\sum_{i=1}^N x_i(y_i - x_i'\hat{\beta}) = 0$$

a. if  $x_i$  contains a constant, then  $\bar{y} = \bar{x}'\hat{\beta}$ : the regression model “fits the mean of  $y$ ”

b. if  $x_i$  contains a dummy variable for membership in group  $g$  then  $\bar{y}_g = \bar{x}_g'\hat{\beta}$ : the regression model “fits the mean of  $y$  for subgroup  $g$ ”

c. Frisch-Waugh (FW) for OLS: The  $j^{th}$  row of  $\hat{\beta}$  is:

$$\hat{\beta}_j = [\frac{1}{N} \sum_{i=1}^N \hat{\xi}_i^2]^{-1} [\frac{1}{N} \sum_{i=1}^N \hat{\xi}_i y_i]$$

where  $\hat{\xi}_i$  is the *estimated residual* from an OLS regression of  $x_{ji}$  on all the other  $x's$ :

$$x_{ji} = x'_{(\sim j)i} \hat{\pi} + \hat{\xi}_i.$$



How are we going to prove FW for OLS?

(i) define  $\hat{u}_i = y_i - x'_i \hat{\beta}$ . We know  $\frac{1}{N} \sum_{i=1}^N x_i \hat{u}_i = 0$

(ii) define  $\hat{\xi}_i = x_{ji} - x'_{(\sim j)i} \hat{\pi}$ . We know  $\frac{1}{N} \sum_{i=1}^N x_{(\sim j)i} \hat{\xi}_i = 0$

(iii) write:  $y_i = \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots + \hat{\beta}_j x_{ji} + \dots + \hat{\beta}_K x_{Ki} + \hat{u}_i$

and form

$$\frac{1}{N} \sum_{i=1}^N \hat{\xi}_i y_i = \frac{1}{N} \sum_{i=1}^N \hat{\xi}_i (\hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots + \hat{\beta}_j x_{ji} + \dots + \hat{\beta}_K x_{Ki} + \hat{u}_i)$$

What terms are equal to 0 from the 2 FOC?

An example: calculating the relationship between gender and self-reported health, controlling for wage

Sample: people age 30-60 in March 2014 Current Pop Survey

Questions: age, education, race, gender, earnings last year, hours worked last year, self-reported health (1-5 scale)

# Simple Statistics

| Variable | N     | Mean     | Std Dev  | Sum     | Minimum  | Maximum   |
|----------|-------|----------|----------|---------|----------|-----------|
| age      | 41754 | 44.23804 | 8.63941  | 1847115 | 30.00000 | 60.00000  |
| educ     | 41754 | 13.99313 | 2.91045  | 584269  | 0        | 20.00000  |
| twage    | 41754 | 25.93441 | 19.65316 | 1082865 | 6.00000  | 200.00000 |
| logwage  | 41754 | 3.05289  | 0.61457  | 127470  | 1.79176  | 5.29832   |
| healthr  | 41754 | 3.90564  | 0.92125  | 163076  | 1.00000  | 5.00000   |
| female   | 41754 | 0.47653  | 0.49945  | 19897   | 0        | 1.00000   |

## Pearson Correlation Coefficients, N = 41754

|         | age      | educ     | twage    | logwage  | healthr  | female   |
|---------|----------|----------|----------|----------|----------|----------|
| age     | 1.00000  | -0.02521 | 0.07538  | 0.07627  | -0.13714 | 0.00752  |
| educ    | -0.02521 | 1.00000  | 0.39451  | 0.45478  | 0.19646  | 0.07658  |
| twage   | 0.07538  | 0.39451  | 1.00000  | 0.89902  | 0.13500  | -0.14397 |
| logwage | 0.07627  | 0.45478  | 0.89902  | 1.00000  | 0.16880  | -0.16728 |
| healthr | -0.13714 | 0.19646  | 0.13500  | 0.16880  | 1.00000  | -0.02178 |
| female  | 0.00752  | 0.07658  | -0.14397 | -0.16728 | -0.02178 | 1.00000  |

Self reported Health and Gender  
Age 30-60 and worked last year, March 2014)

|                           | males | females | all  |
|---------------------------|-------|---------|------|
| Health = 1<br>(poor)      | 0.9   | 1.1     | 1.0  |
| Health = 2<br>(fair)      | 4.8   | 5.6     | 5.2  |
| Health = 3<br>(good)      | 26.0  | 26.3    | 26.2 |
| Health = 4<br>(very good) | 37.6  | 37.7    | 37.7 |
| Health = 5<br>(excellent) | 30.7  | 29.3    | 30.0 |

| Dependent variable: |                     |                    |                        |                    |
|---------------------|---------------------|--------------------|------------------------|--------------------|
|                     | Health              | Health             | Female<br>(Aux. Regr.) | Health             |
| Constant            | 3.9248<br>(0.0062)  | 3.1220<br>(0.0239) | 0.8916<br>(0.0122)     | 3.9056<br>(0.0045) |
| Female              | -0.0402<br>(0.0091) | 0.0123<br>(0.0090) | --                     | --                 |
| log(wage)           | --                  | 0.2547<br>(0.0073) | -0.1359<br>(0.0039)    | --                 |
| resid-Female        | --                  | --                 | --                     | 0.0123<br>(0.0090) |

## *Omitted Variables Formula*

FW: if you add a regressor, the coefficient is “as if” you added only the part of that regressor that is *unexplained by the other regressors*:  $x_{ji} - x'_{(\sim j)i}\pi$  (or  $\hat{\pi}$ )

What about the opposite direction? What happens if you forget a regressor?

$$y_i = \beta_1^* x_{1i} + \beta_2^* x_{2i} + \beta_3^* x_{ji} + u_i$$

Suppose we don't include  $x_{3i}$ ?

a. *Population version*

Aux. model for the *omitted variable*:  $x_{3i} = \pi_1 x_{1i} + \pi_2 x_{2i} + \xi_i$

Then:

$$\begin{aligned} y_i &= \beta_1^* x_{1i} + \beta_2^* x_{2i} + \beta_3^* (\pi_1 x_{1i} + \pi_2 x_{2i} + \xi_i) + u_i \\ &= (\beta_1^* + \beta_3^* \pi_1) x_{1i} + (\beta_2^* + \beta_3^* \pi_2) x_{2i} + \beta_3^* \xi_i + u_i \\ &= \beta_1^0 x_{1i} + \beta_2^0 x_{2i} + \eta_i \end{aligned}$$

Notice that  $E[(x_{1i}, x_{2i})' \eta_i] = E[(x_{1i}, x_{2i})' (\beta_3^* \xi_i + u_i)] = (0, 0)'$ .

So  $(\beta_1^0, \beta_2^0)$  satisfy FOC for population regression of  $y_i$  on  $(x_{1i}, x_{2i})$

Conclusion: If

$$y_i = \beta_1^* x_{1i} + \beta_2^* x_{2i} + \beta_3^* x_{ji} + u_i$$

and we don't include  $x_{3i}$ , the coefficient on  $x_{2i}$  is:

$$\beta_2^0 = \beta_2^* + \beta_3^* \pi_2$$

where  $\pi_2$  is the coefficient on  $x_{2i}$  from the regression of the omitted variable on the remaining  $x$ 's:

$$x_{ji} = \pi_1 x_{1i} + \pi_2 x_{2i} + \xi_i$$

Intuition: you forgot  $x_{3i}$  so the house elf is doing the best he can to predict  $y$  given what he has to work with. The best he can do is use the other  $x$ 's to predict  $x_{3i}$ .



*b. OLS (sample) version*

Aux. model for the omitted variable:  $x_{3i} = \hat{\pi}_1 x_{1i} + \hat{\pi}_2 x_{2i} + \hat{\xi}_i$

Then:

$$\begin{aligned} y_i &= \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 (\hat{\pi}_1 x_{1i} + \hat{\pi}_2 x_{2i} + \hat{\xi}_i) + \hat{u}_i \\ &= (\hat{\beta}_1 + \hat{\beta}_3 \hat{\pi}_1) x_{1i} + (\hat{\beta}_2 + \hat{\beta}_3 \hat{\pi}_2) x_{2i} + \hat{\beta}_3 \hat{\xi}_i + \hat{u}_i \\ &= \hat{\beta}_1^0 x_{1i} + \hat{\beta}_2^0 x_{2i} + \hat{\eta}_i \end{aligned}$$

Notice  $\frac{1}{N} \sum_{i=1}^N (x_{1i}, x_{2i})' \hat{\eta}_i = \frac{1}{N} \sum_{i=1}^N (x_{1i}, x_{2i})' (\hat{\beta}_3 \hat{\xi}_i + \hat{u}_i) = (0, 0)'$ .

So  $(\hat{\beta}_1^0, \hat{\beta}_2^0)$  satisfy FOC for OLS regression of  $y_i$  on  $(x_{1i}, x_{2i})$

Summary of OLS version: OLS if you used  $x_{3i}$ :

$$y_i = \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + \hat{u}_i$$

If we don't include  $x_{3i}$ , the OLS coefficient on  $x_{2i}$  is:

$$\hat{\beta}_2^0 = \hat{\beta}_2 + \hat{\beta}_3 \hat{\pi}_2$$

where  $\hat{\pi}_2$  is the coefficient on  $x_{2i}$  from the OLS regression of the omitted variable on the remaining  $x$ 's:

$$x_{3i} = \hat{\pi}_1 x_{1i} + \hat{\pi}_2 x_{2i} + \hat{\xi}_i$$

Same intuition as for population version.

An example: Suppose we are interested in comparing wages of immigrants and natives. We will use a sample of women age 35-44 who were surveyed in the March 2012 CPS about earnings last year. We consider two models:

$$\log(wage) = \beta_1 + \beta_2 Immigrant + \beta_3 Education \quad (1)$$

and a simpler model:

$$\log(wage) = \beta_1^0 + \beta_2^0 Immigrant \quad (2)$$

Using the OVF, we can relate  $\beta_2^0$  to  $\beta_2$ ,  $\beta_3$ , and the coefficient from an auxiliary regression:

$$Education = \pi_1 + \pi_2 Immigrant$$

We know that

$$\beta_2^0 = \beta_2 + \beta_3\pi_2$$

This holds both for the population regression and for the OLS estimates. What is  $\pi_2$ ? In this regression we are including a dummy for immigrant status. So in the population version

$$\pi_2 = E[\textit{education}|\textit{immigrant}] - E[\textit{education}|\textit{native}]$$

and in the sample version  $\pi_2$  will equal the difference in mean education between immigrants and natives. This will be a pretty big negative number! And since  $\beta_3$  is a number like 0.11 we can conclude that if you “leave out” education, you will tend to find that immigrants earn less.

Table 1: Relationships Between Wages, Education and Immigrant Status for Working Women Age 35-44 in March 2012 Current Population Survey

|                    | Log Wage<br>(1)     | Log Wage<br>(2)    | Immigrant<br>Status<br>(3) | Years of<br>Education<br>(4) | Log Wage<br>(5)     |
|--------------------|---------------------|--------------------|----------------------------|------------------------------|---------------------|
| Immigrant Status   | -0.1800<br>(0.0165) | --                 | --                         | -1.4920<br>(0.0674)          | -0.0101<br>(0.0129) |
| Years of Education | --                  | 0.1141<br>(0.0021) | -0.0297<br>(0.0013)        | --                           | 0.1139<br>(0.0021)  |
| R-squared of model | 0.0111              | 0.2239             | 0.0442                     | 0.0442                       | 0.2239              |

Notes: Each column reports separate regression of dependent variable in column heading on regressors shown in rows, plus a constant (coefficient not reported). Sample is females age 35-44 in March 2012 CPS who reported earnings for the last year. "Wage" refers to average hourly earnings last year, n=10,601. Means and standard deviations of dependent variables are: for log wage 2.8527 and 0.6677; for immigrant status 0.1872 and 0.3901; for education 14.1724 and 2.7676.

NOTE:  $-0.0101 + 0.1139 \times (-1.4920) = -0.1800$   
 $0.1139 - 0.0101 \times (-0.0297) = 0.1141$

The real importance of the OVF is that we can often think about how the omission of a variable affects the estimated coefficient of variables we include, even if we can't estimate the auxiliary regression. Classic example: suppose the true model is:

$$\log(wage) = \beta_1 + \beta_2 Education + \beta_3 Ability$$

But we don't know "ability", and estimate the simpler model:

$$\log(wage) = \beta_1^0 + \beta_2^0 Education$$

We know that

$$\beta_2^0 = \beta_2 + \beta_3 \pi_2$$

where in this case,  $\pi_2$  is the coefficient from a regression of ability on education. Many people (especially those with high education) think that  $\beta_3 > 0$  and  $\pi_2 > 0$ , which leads them to believe that a model that does not control for ability "overstates" the effect of education.