

Linear projections

- The linear projection is defined as follows

$$L(y|1, x_1, \dots, x_K) = L(y|1, \mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K = \beta_0 + \mathbf{x}\beta$$

and β is defined as

$$\beta = [\text{Var}(\mathbf{x})]^{-1} \text{Cov}(\mathbf{x}, y)$$

$$\beta_0 = E(y) - E(\mathbf{x})\beta$$

- We can always write

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + u$$

and from definition of linear projection $E(u) = 0$, $\text{Cov}(x_j, u) = 0$, $j = 1, \dots, N$

- It can be shown that linear projection is **minimum mean square error linear predictor** of y

$$\min_{b_0, \mathbf{b}} E [y - b_0 - \mathbf{x}\mathbf{b}]^2$$

- Iteration property

$$L(y | 1, \mathbf{x}) = L(L(y | 1, \mathbf{x}, \mathbf{z}) | \mathbf{x})$$

- Also

$$L(y | 1, \mathbf{x}) = L(E(y | 1, \mathbf{x}, \mathbf{z}) | \mathbf{x})$$

Endogeneity problem

- Explanatory variable x_j is said to be endogenous if $E(u|x_j) \neq 0$
- Endogeneity problem arises because:
 - **Omitted variable problem**

$$E(u|x, q) \neq E(u|x)$$

* In such a case model in error form is:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + \beta_q q + u$$

- * if q is omitted from regression part then it has to be included in error term

$$u^* = \beta_q q + u$$

- * **But:** if q and x_j are correlated than x_j and u^* are correlated and x_j is endogenous
 - **Measurement error.** If instead of x_j we approximate value x_j^* , than measurement error will become part of u and u can become correlated with x_j^*
 - **Simultaneity.** Simultaneity arises if y influences (e.g. with feedback) one of the explanatory variables
- Sometimes the distinction between this three sources of endogeneity is not quite sharp

Omitted variable problem

- Model with additive omitted variable

$$E(y | x_1, x_2, \dots, x_K, q) = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + \gamma q$$

where q is omitted variable

- We are interested in partial effects of x_j on y holding all other variables *including* q constant
- Model in error form

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + \gamma q + v$$

$$E(v | x_1, \dots, x_K, q) = 0$$

- v - structural error
- Assume that $E(q) = 0$, this assumption only influences constant
- Remove v from model in error form

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + u$$

$$u = \gamma q + v$$

- From $E(q) = E(v) = 0$ we have $E(u) = 0$
- **But:** u is only uncorrelated with all x_j if q is uncorrelated with all x_j

- If q is correlated with one of x_j than we have endogeneity problem and by *OLS* we cannot estimate consistently *any* of β_j
- The asymptotic bias resulting from omitted variable problem is called omitted variable inconsistency or omitted variable bias
- Linear projection of q onto explanatory variables

$$q = \delta_0 + \delta_1 x_1 + \dots + \delta_K x_K + r$$

- Substituting this equation into our model we obtain:

$$y = (\beta_0 + \gamma\delta_0) + (\beta_1 + \gamma\delta_1) x_1 + \dots + (\beta_K + \gamma\delta_K) x_K + v + \gamma r$$

- By definition of linear projection $E(r) = 0$, $\text{Cov}(x_j, r) = 0$ for $j = 1, \dots, K$

- The error $v + \gamma r$ has zero mean and is uncorrelated with all the regressors
- Therefore *OLS* estimate $\text{plim} \left(\hat{\beta}_j \right) = \beta_j + \gamma \delta_j$
- If we assume that q is correlated with only one variable, say x_K so that $\delta_j = 0$ for $j \neq K$ and from definition of linear projection

$$\delta_K = \frac{\text{Cov}(x_K, q)}{\text{Var}(x_K)}$$

So

$$\text{plim} \hat{\beta}_K = \beta_K + \gamma \frac{\text{Cov}(x_K, q)}{\text{Var}(x_K)}$$

- This means that:
 - $\gamma \text{Cov}(x_K, q) > 0$ than bias of $\hat{\beta}_K$ is positive

– $\gamma \text{Cov}(x_K, q) < 0$ than bias of $\hat{\beta}_K$ is negative

Example 1. *(Wooldridge) Wage equation with unobserved ability*

$$\log(\text{wage}) = \beta_1 \text{exper} + \beta_2 \text{exper}^2 + \beta_3 \text{educ} + \gamma \text{abil} + v$$

if abil is only correlated with educ:

$$\text{abil} = \delta_0 + \delta_3 \text{educ} + r$$

and abil is omitted from the model, then estimated coefficient for educ is equal in the limit $\text{plim}(\hat{\beta}_3) = \beta_3 + \gamma\delta_3$. If $\delta_3 > 0$ then the influence of education is overestimated.

Proxy variable solution for omitted variable

- The proxy variable should be
 - redundant (ignorable)

$$E(y | \mathbf{x}, q, z) = E(y | \mathbf{x}, q)$$

- correlation between q and x_1, \dots, x_K should be zero once the effect of z is removed out

$$L(q | 1, x_1, \dots, x_K, z) = L(q | z)$$

This condition can also be expressed in error form

$$q = \boldsymbol{\theta}_0 + \boldsymbol{\theta}_1 z + r$$

by definition of linear projection $E(r) = 0$, $\text{Cov}(x_j, r) = 0$

- Now we can write

$$y = \beta_0 + \gamma\theta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_kx_K + \gamma\theta_1z + (\gamma r + v)$$

- Under two assumptions made $(\gamma r + v)$ is uncorrelated with x_j for $j = 1, \dots, N$
- Conditions of consistency of *OLS* imply that β can be consistently estimated in this case
- Imperfect proxy

$$q = \theta_0 + \rho_1x_1 + \dots + \rho_Kx_K + \theta_1z + r$$

- This gives $\text{plim} \left(\widehat{\beta}_j \right) = \beta_j + \gamma\rho_j$

- We may hope that the bias is much smaller than β_j
- Models with interactions with unobservables

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K + \gamma_1 q + \gamma_2 x_K q + v$$

- Partial effect of x_K

$$\frac{\partial \mathbb{E}(y | \mathbf{x}, q)}{\partial x_K} = \beta_K + \gamma_2 q$$

- We cannot estimate the partial effect as q is not observable
- Assuming that $\mathbb{E}(q) = 0$ we can however estimate the *average partial effect*

$$\mathbb{E}(\beta_K + \gamma_2 q) = \beta_K$$

or for binary variable

$$E(y | x_1, \dots, x_{K-1}, 1, q) - E(y | x_1, \dots, x_{K-1}, 0, q) = \beta_K$$

- If $E(q | \mathbf{x}) = 0$ than we can estimate this with *OLS*
- In the case of using proxy

$$E(q | \mathbf{x}, z) = \theta_1 z + r$$

- Estimated equation

$$\begin{aligned} E(y | \mathbf{x}, z) &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K \\ &\quad + \gamma_1 \theta_1 z + \gamma_2 x_K \theta_1 z + \gamma_1 r + \gamma_2 x_K r + v \end{aligned}$$

- The result of estimation will be correct if proxy variable has expected value equal to zero
- If proxy variable has mean significantly different from zero then we may demean it to have the zero mean condition fulfilled

Example 2. *(Wooldridge) Using IQ as proxy for ability (Blackburn and Neumark 1992)*

$$\log(wage) = \beta_0 + \beta_1 exper + \beta_2 tenure + \beta_3 married + \beta_4 south + \beta_5 urban + \beta_6 black + \beta_7 educ + \gamma IQ + v$$

Estimated coefficient for education: $\beta_7 = .065$ and with IQ included $\beta_7 = .054$. Indeed it seems that coefficient for education overestimated.

Instrumental Variable Estimator (IV)

- Linear model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_K x_K + u \quad (1)$$

$$E(u) = 0, \quad \text{Cov}(x_j, u) = 0, \quad j = 1, 2, \dots, K - 1$$

- x_1, \dots, x_{K-1} are uncorrelated with u but x_K is correlated with u (x_K is endogenous)
- *OLS* estimator of the parameters is inconsistent for *all* the parameters β_j for $j = 1, \dots, K$

- The method of instrumental variables gives the estimation technique which solves this problem
- **Instrumental variable estimator (IV)**
- We need an observable variable z_1 , not included in regression (redundant) which satisfies:
 1. In uncorrelated with the error term u

$$\text{Cov}(z_1, u) = 0$$

2. Coefficient on z_1 in linear projection of x_K on x_1, \dots, x_{K-1}, z_1 is not equal to zero

$$x_K = \delta_0 + \delta_1 x_1 + \dots + \delta_{K-1} x_{K-1} + \theta_1 z_1 + r_K \quad (2)$$

$$\theta_1 \neq 0$$

- The second condition can also be formulated as the requirement that z_1 is partially correlated with x_K
- z_1 satisfying these two conditions is called *instrumental variable* (instrument) for x_K
- Equation (2) is called *reduced form equation* for endogenous variable x_K
- Equation (1) can be rewritten as

$$y = \mathbf{x}\boldsymbol{\beta} + u$$

- The vector of all exogenous (uncorrelated with u) variables is

$$\mathbf{z} = (1, x_1, \dots, x_{K-1}, z_1)$$

- From assumptions

$$E(z'u) = 0$$

- Multiply (2), by z' and take expectations

$$E(z'y) = E(z'x\beta) + \underbrace{E(z'u)}_0$$

- If $\text{Rank}[E(z'z)] = K$ (this condition is satisfied if $\theta_1 \neq 0$) then

$$\beta = [E(z'x)]^{-1} E(z'y)$$

- Condition that $\theta_1 \neq 0$ can be tested by checking whether z_1 is significant in reduced form equation
- Parameter β is said to be identified if it can be expressed as a function of expectations of the data

- Using analogy principle (replacing expectations with data means) we obtain:

$$\hat{\beta} = \left(N^{-1} \sum_{i=1}^N \mathbf{z}'_i \mathbf{x}_i \right)^{-1} \left(N^{-1} \sum_{i=1}^N \mathbf{x}'_i y_i \right) = (\mathbf{Z}' \mathbf{X})^{-1} \mathbf{Z}' \mathbf{y}$$

Example 3. *(Wooldridge) Instrumental Variables for education in wage equation*

$$\log(\text{wage}) = \beta_0 + \text{exper} + \beta_2 \text{exper}^2 + \beta_3 \text{educ} + u$$

As an instrument for child education mother education can be used: it is correlated with child education but it should not directly influence wages.

Problem: mother education can be correlated with other omitted factors in wage education.

- Angrist and Kruger (1991) instrument for education: quarter of birth
- Card (1995): instrument for education: college proximity
- *Natural experiment instruments*: Angrist (1990) effect of veteran status on wages - instrument: draft lottery number

2 Step Least Squares (2SLS)

- Model:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_{K-1} x_{K-1} + \beta_K x_K + u$$

- For simplicity we assume that only x_K is endogenous
- More than one instrument for endogenous variable x_K
- M instruments, z_1, \dots, z_K such that for any h

$$\text{Cov}(z_h, u) = 0$$

- Vector of exogenous variables contains exogenous explanatory and instrumental variables:

$$z = (x_1, \dots, x_{K-1}, z_1, \dots, z_M)$$

- It is possible to find M instrumental variable (IV) estimators of parameter β which are consistent
- Which estimator we should use?
- The most efficient one is the one using the instrument being the linear combination of z most highly correlated with x_K
- The most highly correlated with x_K combination of z is the linear projection of x_K on z .

- Reduced form for x_K is

$$x_K = \delta_0 + \delta_1 x_1 + \dots + \delta_{K-1} x_{K-1} + \theta_1 z_1 + \dots + \theta_M z_M + r_K$$

- New instrument is:

$$x_K^* = \delta_0 + \delta_1 x_1 + \dots + \delta_{K-1} x_{K-1} + \theta_1 z_1 + \dots + \theta_M z_M$$

- As the linear combination of exogenous variables it is uncorrelated with u
- δ_i and θ_j can easily be estimated as *OLS* gives the consistent estimators of coefficients in linear projection

$$\hat{x}_{iK}^* = \hat{\delta}_0 + \hat{\delta}_1 x_{i1} + \dots + \hat{\delta}_{K-1} x_{iK-1} + \hat{\theta}_1 z_{i1} + \dots + \hat{\theta}_M z_{iM}$$

where $\hat{\theta} = (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'\mathbf{X}$.

- Using \widehat{x}_{iK}^* as instrument in *IV* we get

$$\widehat{\beta} = \left(\sum_{i=1}^N \widehat{x}_i' x_i \right)^{-1} \left(\sum_{i=1}^N \widehat{x}_i' y_i \right) = \left(\widehat{\mathbf{X}}' \mathbf{X} \right)^{-1} \widehat{\mathbf{X}}' \mathbf{Y}$$

- Notice that $\widehat{\mathbf{X}} = \mathbf{Z}' \widehat{\theta} = \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} = \mathbf{P}_X \mathbf{X}$, where projection matrix $\mathbf{P}_X = \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}'$ is symmetric ($\mathbf{P}_X' = \mathbf{P}_X$) and idempotent ($\mathbf{P}_X \mathbf{P}_X = \mathbf{P}_X$)

- Then

$$\begin{aligned}
 (\widehat{\mathbf{X}}' \mathbf{X})^{-1} \widehat{\mathbf{X}}' \mathbf{Y} &= (\mathbf{X}' \mathbf{P}'_X \mathbf{X})^{-1} \widehat{\mathbf{X}}' \mathbf{Y} \\
 &= (\mathbf{X}' \mathbf{P}'_X \mathbf{P}_X \mathbf{X})^{-1} \widehat{\mathbf{X}}' \mathbf{Y} \\
 &= (\widehat{\mathbf{X}}' \widehat{\mathbf{X}})^{-1} \widehat{\mathbf{X}}' \mathbf{Y}
 \end{aligned}$$

- **Two Stage Least Squares:**

1. Obtain fitted values \widehat{x}_K from regression of x_K on x_1, \dots, x_{K-1} and z_1, \dots, z_K
2. Run regression of y_i on $x_1, \dots, x_{K-1}, \widehat{x}_K$

- Instruments have to be correlated with endogenous explanatory variable: at least one of $\theta_j \neq 0$

- This assumption can be tested by testing hypothesis

$$H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2 \dots = \boldsymbol{\theta}_M$$

- General case:
 - G number of endogenous explanatory variables
 - L_X - number of exogenous explanatory variables
 - L_Z - number of instruments which are not explanatory variables
 - $L = L_Z + L_X$ total number of instruments (including L_X exogenous explanatory variables)
 - K - total number of explanatory variables
- Assumptions needed for consistency of $2SLS$:
 1. $E(\mathbf{z}'\mathbf{u}) = 0$
 2. Rank conditions:

(a) $\text{Rank} [\mathbb{E} (\mathbf{z}'\mathbf{z})] = L$

(b) $\text{Rank} [\mathbb{E} (\mathbf{z}'\mathbf{x})] = K$

- Condition (2a) is technical: no linear dependence between instruments
- Condition (2b) is important.
- Total number of explanatory variables $K = G + L_X$
- Necessary condition for condition (2b) is $L \geq K \implies L_Z \geq G$ - number of instruments has to be equal or bigger than the number of endogenous variables in the regression equation.
- Asymptotic distribution of $\sqrt{N} (\hat{\beta} - \beta)$ is normal if $\frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{z}_i u_i$ which can be proven from central limit theorem

- If the homoscedasticity assumption is true: $E(u^2 z' z) = \sigma^2 E(z' z)$, where $\sigma^2 = E(u^2)$ then asymptotic variance of $\sqrt{N}(\hat{\beta} - \beta)$ is

$$Avar(\hat{\beta}) = \sigma^2 \left\{ E(x' z) [E(z' z)]^{-1} E(z' x) \right\}$$

- If we define 2SLS residuals as

$$\hat{u}_i = y_i - x_i \hat{\beta}$$

then the unbiased estimator of σ^2 is

$$\hat{\sigma}^2 = (N - K)^{-1} \sum_{i=1}^N \hat{u}_i^2$$

- The estimator of asymptotic variance is then

$$\hat{\sigma}^2 \left(\sum_{i=1}^N \hat{\mathbf{x}}_i' \hat{\mathbf{x}}_i \right)^{-1} = \hat{\sigma}^2 \left(\widehat{\mathbf{X}}' \widehat{\mathbf{X}} \right)^{-1}$$

Example 4. (Wooldridge) Instruments for *educ* is *motheduc*, *fatheduc*, *huseduc*. Reduced form equation for *educ*

$$educ = \delta_0 + \delta_1 exper + \delta_2 exper^2 + \theta_1 motheduc + \theta_2 fatheduc + \theta_3 huseduc + r$$

t-statistics for all θ_i significant

Structural equation

$$\log(wage) = -\underset{(.285)}{.187} + \underset{(.013)}{.043} exper - \underset{(.00040)}{.00086} exper^2 + \underset{(.022)}{.080} educ$$

OLS coefficient for education .107.

- For testing hypothesis we may use the Wald statistics or use the sum of squares

$$F = \frac{(SSR_R - SSR) / g}{SSR / (N - K)}$$

where SSR is a sum of $2SLS$ residuals in unrestricted model and SSR_R is a sum of $2SLS$ residuals in restricted model

- Beware that $SSR = \sum_{i=1}^N \hat{u}_i^2$ and $\hat{u}_i = y_i - \mathbf{x}_i \hat{\boldsymbol{\beta}}$ (not $\hat{u}_i = y_i - \hat{\mathbf{x}}_i \hat{\boldsymbol{\beta}}$)
- We can also use the LM statistics for testing
- The heteroscedasticity robust variance matrix and test statistics can be computed

Possible Pitfalls using 2SLS

- The weak correlation between endogenous variable on instrument (x and z)
- It can be shown that for $y = \beta_0 + \beta_1 x_1 + u$

$$\text{plim } \hat{\beta} = \beta + (\sigma_u / \sigma_{x_1}) (\rho_{z_1 u} / \rho_{z_1 x_1})$$

- If correlation between z_1 and x_1 is small then even very small correlation between z_1 and u can result in large asymptotic bias
- It is also the case that if instruments are poor (weakly partially correlated with endogenous explanatory variable x_K) that standard deviation of estimator $\hat{\beta}_K$ will be large

Example 5. *(Wooldridge) Bound, Jaeger and Baker (1995) have shown that Angrist and Kruger (1991) 2SLS estimator using instruments for education based on the date of birth behave poorly even for sample size of 500,000. Reason: very weak correlation between the date of birth and the number of years of education.*

IV solution to omitted variable problem

- Model

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_{K-1} x_{K-1} + \beta_K x_K + \gamma q + v$$

- Explanatory variable q unobserved
- **Instrumental variable solution:**
- Find instruments
 1. redundant in structural equation

2. uncorrelated with omitted variable q
3. correlated with endogenous variable (with a variable correlated with omitted variable)

- Use *IV* or *2SLS* for estimating model

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_{K-1} x_{K-1} + \beta_K x_K + u$$

where $u = \gamma q + v$

- **Indicator variables solutionn:**

1. Find 2 variables redundant in structural equation ($\text{Cov}(q_1, v) = \text{Cov}(q_2, v) = 0$) but correlated with q

$$q_1 = \delta_0 + \delta_1 q + a_1$$

$$q_2 = \rho_0 + \rho_1 q + a_2$$

where

$$\text{Cov}(q, a_1) = \text{Cov}(q, a_2) = \text{Cov}(\mathbf{x}, a_1) = \text{Cov}(\mathbf{x}, a_2) = \text{Cov}(a_1, a_2) = 0$$

2. Rearranging we get $q = -\frac{\delta_0}{\delta_1} + \frac{1}{\delta_1}q_1 - \frac{a_1}{\delta_1} = \gamma_0 + \gamma_1q_1 - \gamma_1a_1$ and plugging for q we obtain

$$y = \alpha_0 + \mathbf{x}\beta + \gamma_1q_1 + (\mathbf{v} - \gamma_1a_1)$$

where $\gamma_1 = \frac{\gamma}{\delta_1}$ but still q_1 is correlated with a_1

3. From assumptions q_1 and a_1 are correlated but q_2 is not correlated with a_1 as a_1 is not correlated with q and a_2 . But q_1 and q_2 are correlated with each other. Then, indicator q_2 may be used as instrument in estimation of the last equation.

Example 6. (Wooldridge) *IQ and KKW (Knowledge of the Working World test score) as Indicators of Ability.* $IQ = \delta_0 + \delta_1 \text{abil} + a_1$, $KKW = \rho_0 + \rho_1 \text{abil} + a_2$.

We add IQ to equation and use KKW as instrument

$$\begin{aligned} \log(\textit{wage}) = & 4.59 + .014 \textit{exper} + .010 \textit{tenure} + .201 \textit{married} \\ & (0.33) \quad (.003) \quad (.003) \quad (.042) \\ & + .051 \textit{south} + .177 \textit{urban} + .023 \textit{black} + .025 \textit{educ} + .013 \textit{IQ} \\ & (.031) \quad (.028) \quad (.074) \quad (.017) \quad (.005) \end{aligned}$$

return to education only 2.5% and not significant.