

Endogeneity. Instrumental variables estimation. Properties of instrumental variables.

Jakub Mućk
SGH Warsaw School of Economics

Least squares estimator

- **Least squares estimator :**

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

where

- ▶ y is the (outcome) dependent variable;
 - ▶ x_1, x_2, \dots, x_K is the set of independent variables;
 - ▶ ε is the error term.
- The dependent variable is explained with the components that vary with the **the dependent variable** and **the error term**.
 - β_0 is the intercept.
 - $\beta_1, \beta_2, \dots, \beta_K$ are the coefficients (slopes) on x_1, x_2, \dots, x_K .

■ Least squares estimator :

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

where

- ▶ y is the (outcome) dependent variable;
 - ▶ x_1, x_2, \dots, x_K is the set of independent variables;
 - ▶ ε is the error term.
- The dependent variable is explained with the components that vary with the **the dependent variable** and **the error term**.
 - β_0 is the intercept.
 - $\beta_1, \beta_2, \dots, \beta_K$ are the coefficients (slopes) on x_1, x_2, \dots, x_K .

$\beta_1, \beta_2, \dots, \beta_K$ measure the effect of change in x_1, x_2, \dots, x_K upon the expected value of y (*ceteris paribus*).

- **Assumption #1:** true DGP (data generating process):

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon. \quad (2)$$

- **Assumption #2:** the expected value of the error term is zero:

$$\mathbb{E}(\varepsilon) = 0, \quad (3)$$

and this implies that $\mathbb{E}(y) = \mathbf{X}\beta$.

- **Assumption #3:** Spherical variance-covariance error matrix.

$$\text{var}(\varepsilon) = \mathbb{E}(\varepsilon\varepsilon') = I\sigma^2 \quad (4)$$

. In particular:

- ▶ the variance of the error term equals σ :

$$\text{var}(\varepsilon) = \sigma^2 = \text{var}(y). \quad (5)$$

- ▶ the covariance between any pair of ε_i and ε_j is zero"

$$\text{cov}(\varepsilon_i, \varepsilon_j) = 0. \quad (6)$$

- **Assumption #4: Exogeneity.** The independent variable are **not random** and therefore they are not correlated with the error term.

$$\mathbb{E}(\mathbf{X}\varepsilon) = 0. \quad (7)$$

- **Assumption #5:** the full rank of matrix of explanatory variables (there is no so-called collinearity):

$$\text{rank}(\mathbf{X}) = K + 1 \leq N. \quad (8)$$

- **Assumption #6 (optional):** the normally distributed error term:

$$\varepsilon \sim \mathcal{N}(0, \sigma^2). \quad (9)$$

Assumptions of the least squares estimators

Under the assumptions A#1-A#5 of the multiple linear regression model, the least squares estimator $\hat{\beta}^{OLS}$ has the smallest variance of all linear and unbiased estimators of β .

$\hat{\beta}^{OLS}$ is the **Best Linear Unbiased Estimators (BLUE)** of β .

- The least squares estimator

$$\hat{\beta}^{OLS} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}. \quad (10)$$

- The variance of the least square estimator

$$\text{Var}(\hat{\beta}^{OLS}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (11)$$

- If the (optional) assumption about normal distribution of the error term is satisfied then

$$\beta \sim \mathcal{N}(\hat{\beta}^{OLS}, \text{Var}(\hat{\beta}^{OLS})). \quad (12)$$

Consequences of endogeneity

Random (stochastic) regressors

Are explanatory variables always predetermined?

Endogenous variables

An explanatory variable is said to be **endogenous** when it is correlated with error term, i.e., $\mathbb{E}(\varepsilon|x) \neq 0$.

1. Bias of the least squares estimator

$$\mathbb{E}(\hat{\beta}^{OLS}) = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}(\varepsilon), \quad (13)$$

2. Inconsistency of the least squares estimator

$$\text{plim}_{N \rightarrow \infty} \hat{\beta}^{OLS} = \beta + \text{plim}_{N \rightarrow \infty} \left(\frac{1}{N} \mathbf{X}'\mathbf{X} \right)^{-1} \frac{1}{N} \mathbf{X}'\mathbb{E}(\varepsilon). \quad (14)$$

Examples of endogeneity problems

Standard cases when explanatory variables are endogenous:

1. Measurement error.
2. Omitted variable bias.
3. Simultaneous causality.

In addition,

1. Dynamic panel data models.
2. Autoregression with the serial correlation of the error term.

- Let's assume that **true** DGP (data generating process) for the consumption (c) is as follows:

$$c = \alpha + \beta inc^* + \varepsilon \quad (15)$$

where inc^* is **the permanent income**.

- Usually, we have data on income inc but not **the permanent income**. If so, we can proxy the permanent income by current income:

$$inc^* = inc + \eta, \quad (16)$$

where η stands for the measurement and $\eta \sim \mathcal{N}(0, \sigma_\eta^2)$.

- The current income (inc) is **proxy variable** for the permanent income (inc^*).
- Substituting the permanent income into (15):

$$c = \alpha + \beta (inc + \eta) + \varepsilon = \alpha + \beta inc + \beta \eta + \varepsilon = \alpha + \beta inc + \nu, \quad (17)$$

where $\nu = \varepsilon + \beta \eta$.

- The covariance between inc and error term (ν):

$$cov(inc, \nu) = \mathbb{E}(inc\nu) = \mathbb{E}((inc^* + \eta)(\varepsilon + \beta\eta)) = \mathbb{E}(\beta\eta^2) = \sigma_\eta^2\beta \neq 0. \quad (18)$$

- **Labor economics:** returns to education.
- Let's assume that **true** DGP (data generating process) for the logged wage (w):

$$w = \alpha + \rho\mathcal{S} + \beta\mathcal{A} + \varepsilon, \quad (19)$$

where \mathcal{S} is the highest grade of schooling completed and \mathcal{A} is a measure of personal ability or (and) motivation.

- **Problem:** data on \mathcal{A} are not unavailable.
- Consider alternative version of (20):

$$w = \alpha + \rho\mathcal{S} + \eta, \quad (20)$$

- where the error term η captures personal abilities \mathcal{A} , i.e., $\eta = \varepsilon + \beta\mathcal{A}$.
- The OLS estimator of ρ can be simplified to:

$$\hat{\rho}^{OLS} = \text{cov}(w, \mathcal{S}) / \text{Var}(\mathcal{S}). \quad (21)$$

- Plugging *true* DGP for wages w :

$$\hat{\rho}^{OLS} = \frac{\text{cov}(\alpha + \rho\mathcal{S} + \beta\mathcal{A} + \varepsilon, \mathcal{S})}{\text{Var}(\mathcal{S})}, \quad (22)$$

- After manipulation we get:

$$\hat{\rho}^{OLS} = \frac{1}{Var(\mathcal{S})} \mathbb{E}[(\alpha + \rho\mathcal{S} + \varepsilon)\mathcal{S} + \beta\mathcal{A}\mathcal{S}] = \rho + \underbrace{\beta \frac{cov(\mathcal{A}, \mathcal{S})}{Var(\mathcal{S})}}_{=bias} \neq \rho. \quad (23)$$

- The OLS coefficient on schooling would be upward biased if the signs of β and $cov(\mathcal{A}, \mathcal{S})/Var(\mathcal{S})$ are the same.

- Simple (Keynesian) model of consumption:

$$c = \alpha + \beta y + \varepsilon \quad (24)$$

$$y = c + i \quad (25)$$

where c is the consumption, y is the aggregate product, i stands for the investment and ε is the error term, i.e., $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$.

- In the above system we have two endogenous variables (c and y) and one exogenous variable (i).
- The reduced form will be defined as model in which endogenous variable(s) is determined by the exogenous variables as well as the stochastic disturbances. In our case:

$$y = c + i$$

$$y = \alpha + \beta y + \varepsilon + i$$

$$(1 - \beta)y = \alpha i + \varepsilon$$

$$y = \frac{\alpha}{(1 - \beta)} + \frac{1}{(1 - \beta)}i + \frac{1}{(1 - \beta)}\varepsilon.$$

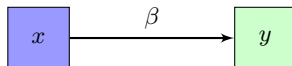
- The general expression of the OLS estimator of the marginal propensity to consume (β) from equation (24):

$$\hat{\beta}^{OLS} = \beta + \frac{\sum (y - \bar{y}) \varepsilon}{\underbrace{\sum (y - \bar{y})^2}} \quad (26)$$

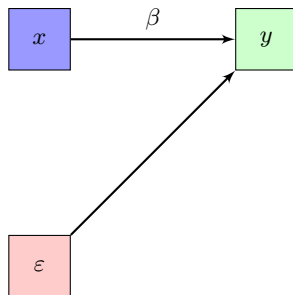
$=0$ if $E(y|\varepsilon)=0$

- But we know that y depends on ε (see the reduced form). If so, then the $\hat{\beta}^{OLS} \neq \beta$ and the OLS estimator is not consistent.

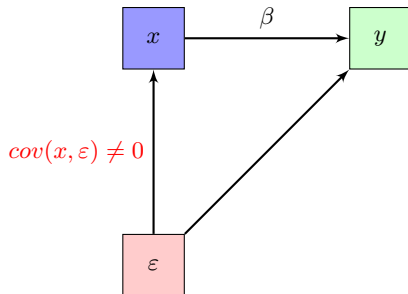
Instrumental variables (IV)



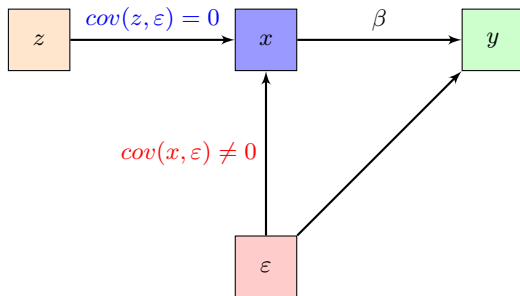
- x – the explanatory variable;
- y – the dependent variable;



- x – the explanatory variable;
- y – the dependent variable;
- ε – the error term;



- x – the explanatory variable;
- y – the dependent variable;
- ε – the error term;



- x – the explanatory variable;
- y – the dependent variable;
- ε – the error term;
- z – the instrumental variable.

- Consider the linear model with single explanatory variable:

$$y = \alpha + \beta x + \varepsilon \quad \text{and} \quad \text{cov}(\varepsilon|x) \neq 0. \quad (27)$$

- The OLS estimates of β will be inconsistent.
- **Instrumental variable regression (IV)** divides variation of the endogenous variable (x) in two parts:
 1. a part that might be **not** correlated with the error term (ε),
 2. a part that might be correlated with the error term (ε).
- It is possible due to using **instrumental variable (instrument, z)** which is not correlated with ε .
- The instrument (z) allows to identify the variation in endogenous variable that is not correlated with ε and, therefore, can be used to estimate β .

- More generally, the IV regression is:

$$y = \alpha + \beta_1 x_1 + \dots + \beta_k x_k + \beta_{k+1} w_1 + \dots + \beta_{k+r} w_r + \varepsilon, \quad (28)$$

where

- y is the dependent variable;
- ε is the error term. In the context of the endogeneity, it might capture omitted factors as well as measurement error;
- x_1, \dots, x_k are k **endogenous** variables that can be correlated with the error term ε ;
- w_1, \dots, w_r are r **exogenous** variables that are potentially not correlated with the error term ε ;
- z_1, \dots, z_m are m **instrumental variables**.

Identification

The coefficients $\beta_1, \dots, \beta_{k+r}$ are said to be:

- exactly identified** if $m = k$;
- underidentified** if $m < k$;
- overidentified** if $m > k$.

The coefficients have to be **exactly identified** or **overidentified** if we want to apply IV regression.

Two conditions for valid instruments

1. Instrument Relevance

A set of instrumental variables (z_1, \dots, z_m) must be related to the endogenous explanatory variables (x_1, \dots, x_k) . Formally,

$$\text{cov}(z_i, x_j) \neq 0.$$

2. Instrument Exogeneity

A set of instrumental variables (z_1, \dots, z_m) cannot be correlated with the error term ε . Formally,

$$\text{cov}(\varepsilon, z_i) = 0.$$

Dependent variable	ENDOGENOUS x	Source of Instrumental variable	Reference
Earnings	Years of schooling	Region and time variation in school construction	Duflo (2001)
Earnings	Years of schooling	Proximity to college	Card (1995)
Earnings	Years of schooling	Quarter of birth	Angrist and Krueger (1991)
Earnings	Veteran status	Cohort dummies	Imbens and van der Klaauw (1995)
Birth weight	Maternal smoking	State cigarette taxes	Evans and Ringel (1999)
Health	Heart attack surgery	Proximity to cardiac care centers	McClellan, McNeil and Newhouse (1994)
College enrollment	Financial aid	Discontinuities in financial aid formula	van der Klaauw (1996)
Crime	Police	Electoral cycles	Levitt (1997)

The Two Stages Least Squares (TSLS) estimator

Two Stage Least Squares (TSLS):

$$y = \beta_0 + \beta_1 x + \varepsilon. \quad (29)$$

1. First-Stage Regression (reduced form):

Regress the endogenous variable (x) on the instrument (z):

$$x = \pi_0 + \pi_1 z + \eta, \quad (30)$$

Based on the OLS estimates calculate predicted values, i.e., \hat{x} .

2. Second -Stage Regression (structural form/equation):

Using OLS regress dependent variable y on the predicted values \hat{x} :

$$y = \beta_0 + \beta_1 \hat{x} + \varepsilon. \quad (31)$$

The TSLS estimator $\hat{\beta}_1^{TSLS}$ stands for the estimates obtained in the second-stage regression.

Two Stage Least Squares (TSLS):

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \beta_{k+1} w_1 + \dots + \beta_{k+r} w_r + \varepsilon. \quad (32)$$

1. First-Stage Regression(s) (reduced form):

Regress each of the endogenous variable (x_i) on the instruments (z_1, \dots, z_m) as well as the exogenous variables (w_1, \dots, w_r):

$$\forall_{i \in 1, \dots, k} \quad x_i = \pi_0 + \pi_1 z_1 + \dots + \pi_m z_m + \pi_{m+1} w_1 + \dots + \pi_{m+r} w_r + \eta, \quad (33)$$

Based on the OLS estimates calculate predicted values, i.e., \hat{x}_i .

2. Second -Stage Regression (structural form/equation):

Using OLS regress dependent variable y on the predicted values $\hat{x}_1, \dots, \hat{x}_k$ as well as the exogenous variables (w_1, \dots, w_r):

$$y = \beta_0 + \beta_1 \hat{x}_1 + \dots + \beta_k \hat{x}_k + \beta_{k+1} w_1 + \dots + \beta_{k+r} w_r + \varepsilon. \quad (34)$$

The TSLS estimator $\hat{\beta}_1^{TSLS}, \dots, \hat{\beta}_k^{TSLS}, \dots, \hat{\beta}_{k+r}^{TSLS}$ stands for the estimates obtained in the second-stage regression.

Specification Tests

- Consider case when there is one endogenous variable (x) and one instrumental variable (z).
- Then the asymptotic properties of the OLS and TSLS estimators:

$$\text{plim}(\hat{\beta}^{OLS}) = \beta + \text{cor}(x, e) \frac{\sigma_e}{\sigma_x}, \quad (35)$$

$$\text{plim}(\hat{\beta}^{IV}) = \beta + \frac{\text{cor}(z, e) \sigma_e}{\text{cor}(z, x) \sigma_z}. \quad (36)$$

- Then, the two stages least squares estimator is consistent if

$$\frac{\text{cor}(z, e)}{\text{cor}(z, x)} = 0, \quad (37)$$

so when the instrumental variable is **strong** and **exogenous**.

- In the IV estimation the TSLS estimator is less efficient than the OLS estimator.

- To test the strength of instruments it is useful to analyze the first step regression:

$$x_i = \pi_0 + \pi_1 z_1 + \dots + \pi_m z_m + \pi_{m+1} w_1 + \dots + \pi_{m+r} w_r + \eta, \quad (38)$$

the intuitive null hypothesis:

$$\mathcal{H}_0 \quad \pi_1 = \pi_2 = \dots = \pi_m = 0, \quad (39)$$

is related to the weak instruments case.

- It can be tested with \mathcal{F} test.
- The rule of the thumb: if the \mathcal{F} statistic is larger than 10 then the null can be rejected.

- The Hausman test allows to investigate endogeneity of explanatory variable.
- **Key assumption:** the IV estimates are unbiased, i.e., instrumental variables are strong and exogenous.
- The null hypothesis:

$$\mathcal{H}_0 : \text{cov}(x, \varepsilon) = 0, \quad (40)$$

refers to exogeneity of explanatory variable while in the alternative hypothesis

$$\mathcal{H}_0 : \text{cov}(x, \varepsilon) \neq 0, \quad (41)$$

the explanatory variable is endogenous and, therefore, the OLS estimates are inconsistent.

- There are several version of the Hausman test.

- The Hausman-Wu test statistic:

$$\mathcal{H} = (\hat{\beta}^{OLS} - \hat{\beta}^{TOLS})' (Var(\hat{\beta}^{TOLS}) - Var(\hat{\beta}^{OLS}))^{-1} (\hat{\beta}^{OLS} - \hat{\beta}^{TOLS}) \quad (42)$$

is χ^2 distributed with K degrees of freedom.

- In general, the above version of Hausman test allows to check whether the IV and OLS estimates are statistically significant.
- It is assumed that the IV estimator is consistent under both null and alternative.
- The null is about a lack of differences between the OLS and TOLS estimators. In other words, the null is about consistency of the least squares.
- The alternative hypothesis postulates that the difference between the OLS and TOLS estimator is systematic/significant. In other words, the alternative states that the OLS estimator is inconsistent.

- The other version of Hausman test bases on including the residuals from the first step regression into the structural equation. Let's assume simply regression model:

$$y = \beta_0 + \beta_1 x + \varepsilon, \quad (43)$$

where x is being tested for endogeneity.

- First step is the regression of x on the instrumental variables (e.g. z_1 and z_2):

$$x = \theta_0 + \theta_1 z_1 + \theta_2 z_2 + \eta, \quad (44)$$

and obtaining the residuals $\hat{\eta}$.

- In the second step, the structural equation is extended by residuals from the first step, i.e., $\hat{\eta}$

$$y = \beta_0 + \beta_1 x + \delta \hat{\eta} + \varepsilon, \quad (45)$$

The null hypothesis is related to exogeneity:

$$\mathcal{H} : \delta = 0, \quad (46)$$

or no correlation between x and ε . It can be tested with a standard t -test.

- When there are more explanatory variables that are tested to be endogenous then:

- ▶ regression in the first step is repeated for all variables. The residuals from each equation are collected,
- ▶ the auxiliary regression in the second step is extended by residuals from each first step regression,
- ▶ the \mathcal{F} statistic is used to test joint significance of the coefficients on the included residuals.

When the number of instruments is larger than number of endogenous variables (overidentification) we can test their validity.

1. Informal strategy

- ▶ Try different combinations of instrumental variables and compare estimates.

2. Sargan test

- ▶ In the first step, we perform IV estimation using all instrumental variables in order to get the residuals $\hat{\varepsilon}$.
- ▶ Then, the residuals $\hat{\varepsilon}$ are regressed on all available instruments.
- ▶ The surplus instruments are tested with the statistics NR^2 , where N is the number of observations and R^2 is the coefficient of determination.
- ▶ NR^2 is χ^2 distributed with $m - k$ degrees of freedom. $m - k$ is the surplus of the instruments.
- ▶ The null refers to validity of instruments.

- **Standards errors** are little bit more complicated than in the OLS estimator, the variation of endogenous variables is exploited together with consistent estimates.
- **Serial correlation and heteroskedasticity** should also be taken into account.
- **Weak instruments** explain little of variation of the endogenous variables. If the instruments are weak then the TSLS estimates are not reliable.
 - ▶ It can be tested with standard \mathcal{F} statistics (testing the hypothesis that the coefficients on the all instruments are zero) in the first stage.
- **Endogeneity of instruments**
 - ▶ There is no formal statistical test allowing for testing whether instruments are correlated with the error term.