

## Week 9: Simultaneous Equation Models and Miscellaneous Topics

Advanced Econometrics 4EK608

Vysoká škola ekonomická v Praze

- 1 Introduction
- 2 Simultaneity Bias
- 3 Identification problem
- 4 Identification conditions
- 5 Systems with more than two equations
- 6 Miscellaneous topics
  - Simulations, Bootstrap & Monte Carlo studies
  - Alternative approaches to econometric modeling
  - Data mining

## Simultaneity is another important form of endogeneity

Simultaneity occurs if at least two variables are jointly determined. A typical case is when observed outcomes are the result of separate behavioral mechanisms that are coordinated in an equilibrium.

Prototypical case: a system of demand and supply equations:

- $D(p)$  how high *would* demand be if the price was set to  $p$ ?
- $S(p)$  how high *would* supply be if the price was set to  $p$ ?
- Both mechanisms have a ceteris paribus interpretation.
- Observed quantity and price will be determined in equilibrium, where  $D(p) = S(p)$ .

Simultaneous equations systems can be estimated by 2SLS/IVR  
... Identification conditions apply.

## Example 1: Labor supply and demand in agriculture

$$h_s = \alpha_1 w + \beta_1 z_1 + u_1$$

$$h_d = \alpha_2 w + \beta_2 z_2 + u_2$$

- Endogenous variables, exogenous variables, observed and unobserved supply shifter, observed and unobserved demand shifter
- We have  $n$  regions, market sets equilibrium price and quantity in each. We observe the equilibrium values only

$$h_{is} = h_{id} \Rightarrow (h_i, w_i)$$

Example 1: Labor supply and demand in agriculture contnd.

$$h_i = \alpha_1 w_i + \beta_1 z_{i1} + u_{i1}$$

$$h_i = \alpha_2 w_i + \beta_2 z_{i2} + u_{i2}$$

- If we have the same exogenous variables in each equation, we cannot identify (distinguish) equations.
- We assume independence between errors in structural equations & exogenous regressors.

**Example 1:** Labor supply and demand in agriculture contnd.

If we estimate the structural equation with OLS method, estimators will be biased – so called “simultaneity bias”.

$$y_1 = \alpha_1 y_2 + \beta_1 z_1 + u_1$$

$$y_2 = \alpha_2 y_1 + \beta_2 z_2 + u_2$$

$y_2$  is dependent on  $u_1$

(substitute RHS of the 1<sup>st</sup> equation for  $y_1$  in the 2<sup>nd</sup> eq.)

$$\Rightarrow y_2 = \left[ \frac{\alpha_2 \beta_1}{1 - \alpha_2 \alpha_1} \right] z_1 + \left[ \frac{\beta_2}{1 - \alpha_2 \alpha_1} \right] z_2 + \left[ \frac{\alpha_2 u_1 + u_2}{1 - \alpha_2 \alpha_1} \right]$$

# Structural and reduced form equations, 2SLS method

## Structural equations (example)

$$y_1 = \beta_{10} + \beta_{11}y_2 + \beta_{12}z_1 + u_1$$

$$y_2 = \beta_{20} + \beta_{21}y_1 + \beta_{22}z_2 + u_2$$

## Reduced form equations

$$y_1 = \pi_{10} + \pi_{11}z_1 + \pi_{12}z_2 + \varepsilon_1 \quad \Rightarrow \quad \hat{y}_1 \text{ by OLS}$$

$$y_2 = \pi_{20} + \pi_{21}z_1 + \pi_{22}z_2 + \varepsilon_2 \quad \Rightarrow \quad \hat{y}_2 \text{ by OLS}$$

## 2SLS (a special case of IVR)

- 1<sup>st</sup> stage: Estimate reduced forms, get  $\hat{y}_1$  and  $\hat{y}_2$ .
- 2<sup>nd</sup> stage: Replace endogenous regressors in structural equations by fitted values from 1<sup>st</sup> stage, estimate by OLS.

Estimation assumptions and “problems” involved:

- ... Identification of structural equations,
- ... Statistical inference in structural equations (2<sup>nd</sup> stage).

## Example 2: (Structural equations)

Estimation of murder rates

$$murdpc = \beta_{10} + \alpha_1 polpc + \beta_{11} incpc + u_1$$

$$polpc = \beta_{20} + \alpha_2 murdpc + \beta(other\ factors) + u_2$$

- 1<sup>st</sup> equation describes the behaviour of murderers, 2<sup>nd</sup> one the behaviour of municipalities.  
Each one has its ceteris paribus interpretation.
- For the municipality policy, the 1<sup>st</sup> equation is interesting: what is the impact of exogenous increase of police force on the murder rate?
- However, the number of police officers is not exogenous (simultaneity problem).



# Identification problem

## Example 3: (Identification)

Identification problem in a SEM

- Example: Supply and demand for milk

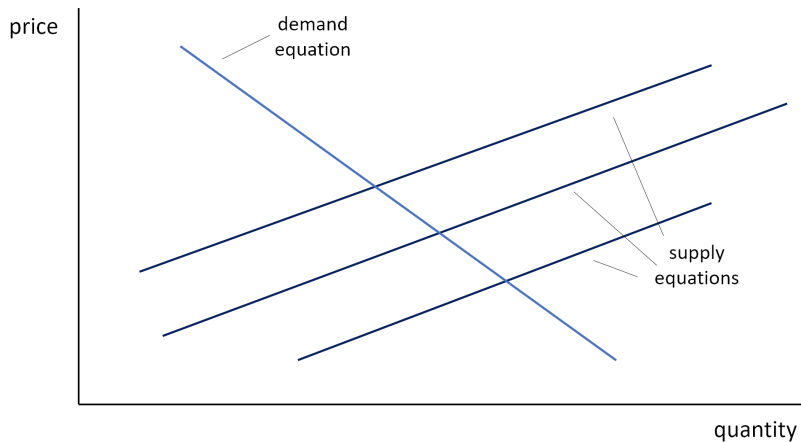
Supply of milk:  $q = \alpha_1 p + \beta_1 z_1 + u_1$

Demand for milk:  $q = \alpha_2 p + u_2$

- Supply of milk cannot be consistently estimated because we do not have (at least) one exogenous variable “available” to be used as instrument for  $p$  in the supply equation.
- Demand for milk can be consistently estimated because we can use exogenous variable  $z_1$  as instrument for  $p$  in the demand equation.

# Identification problem

- Illustration



# Identification conditions

Identification conditions for a sample 2-equation SEM  
(individual  $i$  subscripts omitted)

$$y_1 = \beta_{10} + \alpha_1 y_2 + \beta_{11} z_{11} + \beta_{12} z_{12} + \cdots + \beta_{1k} z_{1k} + u_1$$

$$y_2 = \beta_{20} + \alpha_2 y_1 + \beta_{21} z_{21} + \beta_{22} z_{22} + \cdots + \beta_{2k} z_{2k} + u_2$$

- Order condition (necessary): 1<sup>st</sup> equation is identified if at least one exogenous variable  $z$  is excluded from 1<sup>st</sup> equation (yet in the SEM).
- Rank condition (necessary and sufficient): 1<sup>st</sup> equation is identified if and only if the second equation includes at least one exogenous variable excluded from the first equation with a nonzero coefficient, so that it actually appears in the reduced form.
- For the second equation, the conditions are analogous.
- Some estimation approaches allow for identification through IVs not explicitly included in the SEM.

## Example 4: (Identification)

Labor supply of married working women

Supply (workers):

$$\begin{aligned} \text{hours} = & \alpha_1 \log(\text{wage}) + \beta_{10} + \beta_{11} \text{educ} + \beta_{12} \text{age} + \beta_{13} \text{kidslt6} \\ & + \beta_{14} \text{nwifeinc} + u_1 \end{aligned}$$

Demand (enterprises):

$$\log(\text{wage}) = \alpha_2 \text{hours} + \beta_{20} + \beta_{21} \text{educ} + \beta_{22} \text{exper} + \beta_{23} \text{exper}^2 + u_2$$

Order condition is fulfilled in both equations.

## Example 4: (Identification)

Labor supply of married working women contnd.

- Identification of the first equation (Supply). For the rank condition, either  $\beta_{22}$  or  $\beta_{23}$  non-zero population coefficient (in the second equation) is required – so that *exper*, *exper*<sup>2</sup> (or both) can be used in the reduced form.
- To evaluate the rank condition for supply equation, we estimate the reduced form for  $\log(\textit{wage})$  and test if we can reject the null hypothesis that coefficients for both coefficients for *exper* and *exper*<sup>2</sup> are zero.  
If  $H_0$  is rejected, the rank condition is fulfilled.
- We would do the evaluation of the rank condition for the demand equation analogically.

- We can consistently estimate identified equations with the 2SLS method.
- In the 1<sup>st</sup> stage, we regress each endogenous variable on all exogenous variables (“reduced forms”).
- In the 2<sup>nd</sup> stage we put into the structural equations instead of endogenous variables their predictions from the 1<sup>st</sup> stage and estimate with the OLS method.
- The reduced form can be always estimated (by OLS).
- In the 2<sup>nd</sup> stage, we cannot estimate unidentified structural equations.
- With some additional assumptions, we can use a more efficient estimation method than 2SLS: 3SLS.

# Systems with more than two equations

Example 5: Keynesian macroeconomic model

$$C_t = \beta_0 + \beta_1(Y_t - T_t) + \beta_2 r_t + u_{t1}$$

$$I_t = \gamma_0 + \gamma_1 r_t + u_{t2}$$

$$Y_t \equiv C_t + I_t + G_t$$

Endogenous:  $C_t, I_t, Y_t$

Exogenous:  $T_t, G_t, r_t$

- Order condition for identification is the same as for two equations systems, rank condition is more complicated.
- There exist complicated models based on macroeconomic time series. There is a lot of problems with these models: series are usually not weakly dependent, it is difficult to find enough exogenous variables as instruments. Question is, if any macroeconomic variables are exogenous at all.

# Identification in SEMs with more than two equations

$y_i = X_i\beta + u_i$  is the  $i$ -th equation of a SEM.

$K$  - # of exogenous/predetermined variables in the SEM,

$K_i$  - # of  $K$  in the  $i$ -th equation,

$G_i$  - # of endogenous variables in the  $i$ -th equation.

**Order condition** for the  $i$ -th equation:

necessary, not sufficient condition for identification

$$K - K_i \geq G_i - 1$$

Condition evaluates as:

- = Equation  $i$  is just-identified,
- > Equation  $i$  is over-identified,
- < Equation  $i$  is not identified,  
structural equation  $i$  cannot be estimated by 2SLS/IVR.



# Identification in SEMs with more than two equations

Rank condition: based on matrix algebra & IV estimator

Consider IVR for an identified  $i$ -th equation of SEM

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{u}_i$$

$\mathbf{X}_i$  is a  $(n \times k)$  matrix, includes the intercept column and all endogenous regressors of the  $i$ -th equation,

$\hat{\mathbf{X}}_i$  is a  $(n \times k)$  matrix, includes the intercept column.

Exogenous regressors are repeated from  $\mathbf{X}_i$ , endogenous are projected to the column space of  $\mathbf{Z}$ : a  $(n \times l)$  matrix of all exogenous variables in the SEM.

Single equation (limited information) estimator for each  $i$ -th equation:

- $\hat{\boldsymbol{\beta}}_{IVR} = \hat{\boldsymbol{\beta}}_{2SLS,i} = \left( \hat{\mathbf{X}}_i' \mathbf{X}_i \right)^{-1} \hat{\mathbf{X}}_i' \mathbf{y}$
- $\hat{\mathbf{X}}_i' (\mathbf{y} - \mathbf{X}_i \hat{\boldsymbol{\beta}}) = \mathbf{0}$  (GMM moment conditions)

# Identification in SEMs with more than two equations

Rank condition: based on matrix algebra & IV estimator (cont.)

$$\hat{\beta}_{IVR} = \left( \hat{\mathbf{X}}_i' \mathbf{X}_i \right)^{-1} \hat{\mathbf{X}}_i' \mathbf{y}$$

- **Order condition:** The necessary condition for the  $i$ -th equation to be identified is that the number of columns (exogenous variables of SEM) in  $\mathbf{Z}$  should be no less than the number of columns (explanatory variables) in  $\mathbf{X}_i$ .
- **Rank condition:** The necessary and sufficient condition for identification of the  $i$ -th equation is that  $\hat{\mathbf{X}}_i'$  has full column rank of  $\mathbf{X}_i$ .  
...ensures the existence of  $\left( \hat{\mathbf{X}}_i' \mathbf{X}_i \right)^{-1}$ .

# Identification in SEMs with more than two equations

## Identification: recap & final remarks

- Reduced form equations can always be estimated.
- Structural equations can be estimated (IV/2SLS) only if identified: i.e. if rank condition is met.
- With SW, checking rank condition for  $\left(\hat{\mathbf{X}}_i' \mathbf{X}_i\right)^{-1}$  is easy for finite datasets.
- Asymptotic identification may be “tricky”:  
because some columns in  $\mathbf{X}_i$  are endogenous,  
 $\text{plim } n^{-1} \hat{\mathbf{X}}_i' \mathbf{X}_i$   
depends on the parameters of the DGP.  
...see Davidson-MacKinnon (2009) Econometric theory and methods

## **Miscellaneous topics** - not specifically related to SEMs

- Simulations
- Bootstrap
- Monte Carlo studies
  
- Simple-to-general approach to econometric modeling
- General-to-specific approach to econometric modeling
  
- Data mining

# Simulations & Monte Carlo studies

Simulations and simulation-based methods are used for:

- Inferring characteristics of random variables
- Inferring characteristics of estimators and estimator functions
- Inferring characteristics of tests
- Can be used for construction of numerical (approximate) estimators for highly complex scenarios (see Greene, Chapter 15.6)

Bootstrap and Monte Carlo studies are two common simulation-based methods.

# Bootstrap

Based on repeated draws (with replacement)

“Resampling” of the primary sample

**Sample:** Population (size  $N$ )  $\Rightarrow$  Sample (size  $n$ )

**Bootstrap:** Sample (size  $n$ )  $\Rightarrow$  Bootstrap sample (size  $n$ )

Bootstrap is used for calculating confidence intervals, standard errors and/or bias in estimators.

Bootstrap is based on the assumption that our “primary sample” is a representative sample from the population. By repeatedly sampling (with replacement) from the “primary sample”, we simulate sampling from the (original) population.

While taking repeated samples from the population is the best approach, bootstrap is the second best approach if primary source of sampling is not accessible (cost, time, ...).

# Bootstrap

- ① From a dataset (primary sample with  $n$  observations), take  $B$  (say, 1.000) bootstrap samples (each bootstrap sample has  $n$  observations, chosen randomly with replacement from the original sample).
  - ② For each bootstrap sample, calculate the statistic(s) required: median, variance, regression coefficients in a LRM, etc. Save results for each bootstrap sample.
  - ③ From our saved results ( $B$  bootstrap samples), we obtain the required distributional characteristics.
- Observed data need not follow a specific distribution (e.g. normally distributed errors in LRMs).
  - Useful mainly for CS data (less so for TS analysis)  
Efron, Tibshirani: An Introduction to the Bootstrap (1993)

# Bootstrap: Examples

- Sample mean, sample median and their standard errors

$$\text{s.e.}(\bar{x}) = \frac{\hat{\sigma}}{\sqrt{n}}$$

$$\text{s.e.}(\tilde{x}) = ?$$

- Variance in OLS coefficients - small sample problems:

$\hat{\beta} \sim N(\beta, \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1})$  if and only if  $u_i \sim N(0, \sigma^2), i.i.d.$   
(other G.M. assumptions pending)

First, HC-consistent estimates of variance /  $\text{var}(\hat{\beta})$  / have asymptotic relevance only

Second, the distribution of  $u_i$  may be unknown  
(Jarque-Berra test rejects  $H_0$  of normality in residuals).

Such problems may be solved using the bootstrap method.  
In R, `{boot}` package is available.



# Monte Carlo studies

- Analysis/comparison of properties of estimators.
- For example, in TS analysis, finite-sample properties are often analytically inaccessible.
- We may either compare different estimators (using same dataset) or experiment with DGP (simulate different  $\mathbf{ar}(\mathbf{p})$  properties in error term) and study behavior of the estimators.

Monte Carlo studies (in 4 steps):

- 1 Model the data generating process
- 2 Generate many sets of artificial data
- 3 Use the data and estimator to create repeated estimates
- 4 Use these estimates to gauge the sampling distribution properties (test power, predictive properties ...).

## Simple-to-general approach

- Traditional approach to econometric modeling
- Starts with formulation of the simplest model consistent with the relevant economic theory.
- If this initial model proves unsatisfactory, it is improved in some way – adding or changing variables, using different estimators etc.

## Criticism of the simple-to-general approach

- Revisions to the simple model are carried out arbitrarily and simply reflect investigator's prior beliefs: danger of always finding what you want to find.
- It is open to accusation of data mining: researchers usually presents just the final model (true significance level is problematic).

## General-to-specific approach

- Professor Hendry, London School of Economics started this approach in the 80ies.
- It starts with formulation of a very general and maybe quite complicated model.
- Starting model contains a series of simpler models, nested within it as special cases.
- These simpler models should represent all the alternative economic hypotheses that require consideration.

# Alternative approaches to econometric modeling

## General-to-specific approach

- General model must be able to explain existing data and be able to satisfy various tests of misspecification.
- What follows is simplification search (testing-down procedure). Through parameter restrictions, we test nested models against the containing model. If the nested model does not pass the tests, we can reject the whole branch of sub-nested models.
- If we find more non-nested models satisfying tests, we can compare them using e.g.  $F$ -test.

# Alternative approaches to econometric modeling

## Advantages of the general-to-specific approach

- “Data mining” present in this approach is transparent (for all to see) and it is carried out in a systematic manner that avoids worst data mining problems.
- Researcher usually reports both the initial general model and all steps involved so it is possible to get some idea about the true significance levels.
- Supporters of this approach stress the importance of both testing final models against new data and the ability of the model to provide adequate out-of-sample forecasts.

We use brute-force algorithms to find some statistically significant relationship. It can be completely misleading – as following example shows.

Repetition:

$$t \text{ test: } H_0 : \beta_j = 0 \quad H_1 : \beta_j \neq 0$$

Significance level:

probability of a type I error, i.e. probability of rejecting  $H_0$  when it is in fact true, i.e. finding a regressor significant when -in fact- it does not influence the dependent variable.

Example:

- ❶ Suppose we have 20 “possible” regressors  $x_1, x_2, \dots, x_{20}$ , but all are factually unrelated to the dependent variable  $y$
- ❷ Suppose we have computed 20 simple regressions of the form

$$\hat{y} = \hat{\beta}_{0p} + \hat{\beta}_{1p}x_p$$

- ❸ If we use significance level 0.05, we can expect one of the 20 regressors to appear significant just by chance, even if none of them actually influences  $y$ .



# Data mining

$$\Pr(X_1 \text{ appears significant by chance}) = 0.05$$

$$\Pr(X_1 \text{ does not appear significant}) = 0.95$$

$$\Pr(X_2 \text{ does not appear significant}) = 0.95$$

$$\begin{aligned}\Pr(\text{neither } X_1 \text{ nor } X_2 \text{ appear significant}) &= 0.95 \times 0.95 = \\ &= 0.9025\end{aligned}$$

$$\Pr(\text{at least one of } X_1, X_2 \text{ appear significant}) = 1 - 0.9025 =$$

$$\boxed{\begin{array}{l} \text{true significance level} \\ \alpha^* = (1 - (1 - \alpha)^2) \end{array}} \longrightarrow = 0.0975$$

If we want the true significance level to be 0.05, we must solve the equation  $0.05 = 1 - (1 - \alpha)^2$  and do all  $t$ -tests on significance level  $\alpha^*$  (here,  $\alpha^* \doteq 0.0253$ ).

For  $c$  independent candidates:  $\alpha^* = 1 - (1 - \alpha)^c$ .

Lovell (1983): rule of thumb for finding the true significance level in the case where  $k$  regressors are selected from  $c$  possible candidates.

$$\alpha^* = 1 - (1 - \alpha)^{\frac{c}{k}}$$