

**ÉCONOMÉTRIE  
(UGA, S2)  
CHAPITRE 3 :  
ENDOGENÉITÉ ET VARIABLES INSTRUMENTALES(1)**

Michal W. Urdanivia \*

\* UGA, Faculté d'Économie, GAEL,  
e-mail : [michal.wong-urdanivia@univ-grenoble-alpes.fr](mailto:michal.wong-urdanivia@univ-grenoble-alpes.fr)

12 avril 2022

1. Introduction à la notion de variable instrumentale

2. L'estimateur de VIs

## 1. Introduction à la notion de variable instrumentale

## Endogénéité et (non-)identification dans un modèle linéaire simple

- On considère ici le modèle linéaire le plus simple :

$$y_i = \alpha_0 + b_0 x_i + u_i \text{ avec } E[u_i] := 0, \quad (1)$$

mais on considère que l'analyse du PGD indique que  $x_i$  est endogène dans ce modèle, i.e. que :

$$E[u_i | x_i] \neq 0 \Rightarrow \text{Cov}[x_i; u_i] \neq 0.$$

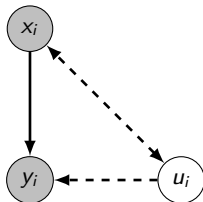
- Rappelons que lorsque  $x_i$  est exogène et que  $\text{Var}[x_i] \neq 0$ , on a par l'exogénéité de  $x_i$  :

$$\text{Cov}[x_i; u_i] = 0 \Leftrightarrow \text{Cov}[x_i; y_i - \alpha_0 + b_0 x_i] = 0 \Rightarrow b_0 = \frac{\text{Cov}[x_i; y_i]}{\text{Var}[x_i]}.$$

- Autrement dit l'exogénéité de  $x_i$  permet d'identifier  $b_0$  (et aussi de  $\alpha_0$ ) comme une fonction de la distribution des variables observées  $y_i$  et  $x_i$ .
- Inversement dans la situation que nous considérons dans ce chapitre  $\text{Cov}[x_i; u_i] \neq 0$ , rend impossible l'identification  $b_0$  (et aussi celle de  $\alpha_0$ ), et la construction d'un estimateur convergent.

## Endogénéité et (non-)identification dans un modèle linéaire simple

- On peut représenter ce problème en utilisant un **graphe causal** :



**Figure 1** – Graphe causal du modèle :  $y_i = \alpha_0 + b_0 x_i + u_i$ , avec  $\text{Cov}(x_i; u_i) \neq 0$ . Les variables  $(x_i, y_i, u_i)$  sont les nœuds du graph et les nœuds foncés correspondent aux variables observées. Les arêtes représentent les relations entre les variables. Les relations observées sont en trait plein.

## Identification avec une variable instrumentale

- L'intuition sous-jacente à la méthode des VIs consiste à répondre à la question de savoir si avec une variable, que nous notons  $z_i$ , il est possible d'obtenir une mesure de la relation causale entre  $x_i$  et  $y_i$  qui ne dépende pas de  $u_i$ .
- Autrement dit,  $z_i$  doit être exogène par rapport à  $u_i$  :

$$E[u_i|z_i] = 0 \Rightarrow \text{Cov}[z_i; u_i] = 0, \quad (2)$$

ce qui nous permet d'écrire :

$$\begin{aligned} \text{Cov}[z_i; u_i] = 0 &\Leftrightarrow \text{Cov}[z_i; y_i - \alpha_0 + b_0 x_i] = 0 \\ &\Leftrightarrow \text{Cov}[z_i; y_i - \alpha_0 + b_0 \text{Cov}[z_i; x_i]] = 0 \end{aligned}$$

- Ceci indique que pour identifier  $b_0$  on doit aussi supposer aussi que,

$$\text{Cov}[z_i; x_i] \neq 0, \quad (3)$$

et  $b_0$  est identifié par :

$$b_0 = \frac{\text{Cov}[z_i; y_i]}{\text{Cov}[z_i; x_i]}, \quad (4)$$

## Identification avec une variable instrumentale

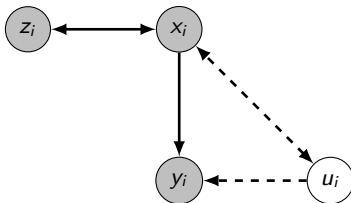
- On peut résumer les conditions (2)-(3) ainsi :

### Définition 1 (Conditions de validité de VIs dans un modèle simple)

Dans le modèle  $y_i = \alpha + b_0 x_i + u_i$  avec  $E[u_i] := 0$ , une variable  $z_i$  est un instrument(de  $x_i$ ) ssi :

- (i)  $\text{Cov}[z_i; u_i] = 0$ , i.e.  $z_i$  est exogène par rapport à  $u_i$  et :
- (ii)  $\text{Cov}[z_i; x_i] \neq 0$ , i.e.  $z_i$  et  $x_i$  sont liées.

- Dans la représentation en termes de graphe causal cela donne :



**Figure 2** – Graphe causal du modèle :  $y_i = \alpha_0 + b_0 x_i + u_i$ , avec  $\text{Cov}(x_i; u_i) \neq 0$ ,  $\text{Cov}(z_i; u_i) = 0$ . Les variables  $(z_i, x_i, y_i, u_i)$  sont les nœuds du graph et les nœuds forcés correspondent aux variables observées. Les arêtes représentent les relations entre les variables. Les relations observées sont en trait plein.

## Estimateur des VIs dans le modèle simple

- L'identification de  $b_0$  par (4) suggère l'estimateur :

$$\hat{b}_N^{VI} = \frac{N^{-1} \sum_{i=1}^N (z_i - \bar{z}_N)(y_i - \bar{y}_N)}{N^{-1} \sum_{i=1}^N (z_i - \bar{z}_N)(x_i - \bar{x}_N)} = \frac{N^{-1} \sum_{i=1}^N (z_i - \bar{z}_N)y_i}{N^{-1} \sum_{i=1}^N (z_i - \bar{z}_N)x_i}, \quad (5)$$

où  $\bar{z}_N$ ,  $\bar{x}_N$ , et  $\bar{y}_N$  sont les moyennes empiriques respectives de  $z_i$ ,  $x_i$ , et  $y_i$ .

- De plus,  $\hat{b}_N^{VI}$  est convergent. Nous avons en effet :

$$\begin{aligned} \text{plim}_{N \rightarrow +\infty} N^{-1} \sum_{i=1}^N (z_i - \bar{z}_N)y_i &\rightarrow \text{Cov}[z_i; y_i], \\ \text{plim}_{N \rightarrow +\infty} N^{-1} \sum_{i=1}^N (z_i - \bar{z}_N)x_i &\rightarrow \text{Cov}[z_i; x_i], \end{aligned}$$

d'où :

$$\begin{aligned} \hat{b}_N^{VI} &\xrightarrow[N \rightarrow +\infty]{p} \frac{\text{Cov}[z_i; y_i]}{\text{Cov}[z_i; x_i]} = \frac{\text{Cov}[z_i; \alpha_0 + b_0 x_i + u_i]}{\text{Cov}[z_i; x_i]}, \\ &= b_0 + \frac{\text{Cov}[z_i; u_i]}{\text{Cov}[z_i; x_i]}, \\ &= b_0. \end{aligned}$$



## Estimateur des VIs dans le modèle simple

### Remarque 1

- ★ On dit des variations de  $z_i$  qu'elles sont des variations exogènes : elles ne sont pas liées à  $u_i$  puisque  $\text{Cov}[z_i; u_i] = 0$ .
- ★ Ce sont les effets de ces variations exogènes sur  $x_i$  qui sont exploitées pour l'identification de  $b_0$  grâce à  $\text{Cov}[z_i; x_i] \neq 0$ .
- ★ Noter qu'il n'est aucunement nécessaire que l'effet de  $z_i$  sur  $x_i$  soit causal.
- ★ L'effet de  $z_i$  sur  $y_i$  ne « transite » que via  $x_i$ . La variable instrumentale  $z_i$  n'est pas une variable explicative dans le modèle de  $y_i$ . On parle alors de relation d'exclusion (de la VI  $z_i$  vis-à-vis du modèle de  $y_i$ ).
- ★ L'estimateur des VIs est parfois appelé estimateur des moindres carrés indirects. Cela provient de ce que  $b_0$  dans (4) peut s'écrire :

$$b_0 = \frac{\text{Cov}[z_i; y_i] / \text{Var}[z_i]}{\text{Cov}[z_i; x_i] / \text{Var}[z_i]},$$

qui est le rapport entre le coefficient de  $z_i$  dans la projection de  $y_i$  sur  $z_i$ , et le coefficient de  $z_i$  dans la projection de  $x_i$  sur  $z_i$ .

## 2. L'estimateur de VIs

## Variables endogènes, exogènes, instruments

- L'objectif est maintenant de généraliser l'approche présentée dans le cas simple précédent au modèle linéaire général :

$$y_i = \mathbf{x}_i' \mathbf{a}_0 + u_i, \text{ avec } E[u_i] := 0. \quad (6)$$

- Plusieurs éléments du vecteur  $\mathbf{x}_i$  peuvent être endogènes de sorte que dans l'estimateur des MCO de  $\mathbf{a}_0$  plusieurs éléments sont potentiellement biaisés (c.f. cours précédent sur les VIs).
- Notons :

$$\mathbf{x}_i = \begin{bmatrix} 1 \\ \tilde{\mathbf{x}}_i^x \\ \tilde{\mathbf{x}}_i^e \end{bmatrix} = \begin{bmatrix} \mathbf{x}_i^x \\ \tilde{\mathbf{x}}_i^e \end{bmatrix} \begin{array}{ll} \{\text{variables explicatives exogènes} & : E[u_i | \mathbf{x}_{k,i}^x] = 0 (k = 1, \dots, M) \\ \{\text{variables explicatives endogènes} & : E[u_i | \mathbf{x}_{k,i}^x] \neq 0 (k = M + 1, \dots, K) \end{array}$$

### Remarque 2

- ★ Il est clair que la variable constante 1 est «exogène» :  $E[1 \times u_i] = E[u_i] = 0$ .
- ★ Comme pour l'estimateur des MCO nous utiliserons la Méthode des Moments pour construire un estimateur convergent de  $\mathbf{a}_0$ , l'estimateur des VI du modèle (6).
- ★ On considère ici que chaque élément  $\mathbf{x}_i^e$  a une variable instrumentale.

## Définition 2 (Variable instrumentale)

$z_{k,i}$  est une variable instrumentale de  $x_{k,i}$  dans le modèle linéaire (6) si :

- (i)  $\text{Cov}[z_{k,i}; u_i] = 0$  i.e.,  $z_{k,i}$  est exogène par rapport à  $u_i$ ,
- (ii)  $z_{k,i}$  « suffisamment » liée à  $x_{k,i}$ .

## Remarque 3

- ★ On verra dans la suite (analyse des conditions de rang) que la condition (ii) doit en fait être définie comme :

$$\text{Cov}[z_{k,i}; e_{k,i}] \neq 0 \text{ pour } k > 1,$$

où  $e_{k,i}$  est la partie spécifique de  $x_{k,i}$  dans  $x_i$ , i.e. le résidu de la projection linéaire de  $x_{k,i}$  sur les autres explicatives  $\mathbf{x}_{-k,i}$ .

$$e_{k,i} = x_{k,i} - \mathcal{EL}[x_{k,i} | \mathbf{x}_{-k,i}].$$

- ★ Dans la définition d'un VI précédente, on voit que lorsqu'une variable explicative  $x_{k,i}$  est exogène alors c'est aussi une variable instrumentale d'elle même. En ce sens que non seulement elle vérifie (i) mais elle vérifie forcément (ii) (car ayant une corrélation de 1 avec elle même)

## Variables endogènes, exogènes, instruments

- On construit le vecteur des variables instrumentales  $\mathbf{z}_i$  avec :

$$\tilde{\mathbf{z}}_i^e = \begin{bmatrix} \tilde{z}_{M+1,i} \\ \tilde{z}_{M+2,i} \\ \vdots \\ \tilde{z}_{K,i} \end{bmatrix} \text{ et } \mathbf{z}_i = \begin{bmatrix} 1 \\ \tilde{\mathbf{x}}_i^x \\ \tilde{\mathbf{z}}_i^e \end{bmatrix} = \begin{bmatrix} \mathbf{x}_i^x \\ \tilde{\mathbf{z}}_i^e \end{bmatrix} \begin{array}{ll} \{ \text{variables exogènes de } \mathbf{x}_i & : E[u_i | \mathbf{x}_{k,i}^x] = 0 (k = 1, \dots, M) \\ \{ \text{variables instrumentales} & : E[u_i | \mathbf{z}_{k,i}] = 0 (k = M + 1, \dots, K) \end{array}$$

- Ce vecteur contient en fait toutes les variables exogènes du modèle. Ce sont ces variables qui assurent l'identification des paramètres du modèle.
- $\mathbf{z}_i$  est parfois nommé ensemble d'information du modèle.

### *Définition 3 (Modèle linéaire à variables instrumentales)*

Le modèle défini par :

$$y_i = \mathbf{x}_i' \mathbf{a}_0 + u_i, \text{ avec } E[u_i | \mathbf{z}_i] = E[u_i] := 0,$$

est un modèle linéaire à variables instrumentales. La condition d'identification de  $\mathbf{a}_0$  dans ce modèle est donnée par :

$$\text{Rang} (E[\mathbf{z}\mathbf{x}']) = K = \dim(\mathbf{x}).$$

### *Remarque 4*

La condition d'exogénéité de  $\mathbf{z}_i$  est définie par  $E[u_i | \mathbf{z}_i] = 0$ , et non par  $\text{Cov}[\mathbf{z}_i; u_i] = \mathbf{0}$ . Ce n'est pas nécessaire pour un modèle linéaire où  $\text{Cov}[\mathbf{z}_i; u_i] = \mathbf{0}$  suffit mais c'est standard et cela simplifie la présentation des hypothèses d'homoscédasticité.

## Variables endogènes, exogènes, instruments

- Comme dans le cas où on a construit l'estimateur des MCO de  $a$  on part de la condition d'exogénéité des  $z_i$  (et non des  $x_i$  comme dans le cas des MCO), i.e. la condition d'orthogonalité donnée par :

$$E[u_i z_i] = 0 \Rightarrow E[z_i u_i] = 0 \Leftrightarrow E[z_i(y_i - x_i' a_0)] = 0.$$

- On a ici la condition de moment estimante pour  $a_0$  est  $E[z_i(y_i - x_i' a_0)] = 0$ . Et on a alors :

$$E[z_i(y_i - x_i' a_0)] = 0 \Leftrightarrow a = a_0.$$

- On suppose ici que  $a_0$  est l'unique solution en  $a$  de  $E[z_i(y_i - x_i' a)] = 0$ .
- Le principe d'analogie définit l'estimateur de la MM de  $a_0$  par :

$$N^{-1} \sum_{i=1}^N z_i(y_i - x_i' a) = 0_{K \times 1} \Leftrightarrow a = \hat{a}_N^{MM}.$$

- L'équation dont  $\hat{a}_N^{MM}$  est définie comme la solution en  $a$  est en fait un système de  $K$  équations linéaires à  $K$  inconnues (les éléments de  $\hat{a}_N^{MM}$ ). Il a solution sous forme explicite. On a :

$$N^{-1} \sum_{i=1}^N z_i(y_i - x_i' \hat{a}_N^{MM}) = 0_{K \times 1}.$$

## Variables endogènes, exogènes, instruments

- Il est aisé de définir la forme de  $\hat{\mathbf{a}}_N^{MM}$ ,

$$N^{-1} \sum_{i=1}^N \mathbf{z}_i y_i - \left[ N^{-1} \sum_{i=1}^N \mathbf{z}_i \mathbf{x}_i' \right] \hat{\mathbf{a}}_N^{MM} = \mathbf{0}_{K \times 1},$$

qui donne finalement :

$$\hat{\mathbf{a}}_N^{MM} = \left[ N^{-1} \sum_{i=1}^N \mathbf{z}_i \mathbf{x}_i' \right]^{-1} N^{-1} \sum_{i=1}^N \mathbf{z}_i y_i,$$

qui définit ce qu'on appelle l'estimateur des VI.



## Références