

Rappel:

modèle de régression multiple

\Rightarrow Hypothèses:

* H_1 : La matrice X est non aléatoire

$$\Rightarrow \text{cov}(X, \varepsilon) = 0$$

* H_2 : X est de plein rang

\Rightarrow les variables explicatives sont linéairement indépendantes ou orthogonales

$$\Rightarrow (X'X)^{-1} \text{ existe}$$

* H_3 : $E(\varepsilon) = 0$

* H_4 : $\text{Var}(\varepsilon) = \sigma^2 \cdot I$

Hypothèse d'absence d'autocorrelation des erreurs.

$\left. \begin{array}{l} \text{Var}(\varepsilon_t) = \sigma^2 \text{ (constante)} \\ \text{cov}(\varepsilon_t, \varepsilon_s) = 0 \quad (\forall t \neq s) \end{array} \right\} \text{Hypothèse d'homoscedasticité}$

* H_5 : $\varepsilon \sim N(0, \sigma^2 \cdot I)$

\Rightarrow Hypothèse de normalité des erreurs.

H_1 : Problèmes liés aux variables explicatives

I La multicollinéarité:

1) Problèmes
Il y a multicollinéarité lorsque l'hypothèse de l'orthogonalité des var explicatives ou encore de leur indépendance linéaire est violée c'est-à-dire lorsque la matrice X n'est plus de plein rang. Ceci implique que la matrice $X'X$ devient singulière et par conséquent son inverse $(X'X)^{-1}$ n'existe pas.

ce qui rend la méthode MCO complètement défectueuse, il n'est donc pas possible de voir une telle situation d'estimer les paramètres du modèle.

On distingue 2 types de multicollinéarité: la multicollinéarité exacte parfaite et la multicollinéarité imparfaite ou quasi-multicollinéarité.

Dans la pratique c'est plutôt le cas de quasi-multicollinéarité est fréquent. En effet, la multicollinéarité imparfaite correspond au cas où la matrice $X'X$ est non singulière mais son det est proche de 0. La conséquence directe est que l'on aura des valeurs très grandes dans la matrice inverse $(X'X)^{-1}$ et ceci implique que l'on aura aussi des valeurs assez grande pour la matrice des variances-covariances $\hat{\sigma}^2 (X'X)^{-1}$ estimée de l'estimateur.

\Rightarrow Les estimateurs deviennent moins précis

$$\Rightarrow \text{test de Student } (t_c = \left| \frac{\hat{\beta}_j}{\hat{\sigma}_{\hat{\beta}_j}} \right|)$$

deviennent de plus en plus faibles donc les paramètres deviennent statistiquement non significatifs.

$\Rightarrow F$ (Fisher) très élevée

$\Rightarrow R^2$ élevée

2) Tests de détection

Les tests de détection de la multicollinéarité

Les plus fréquemment utilisés sont:

- le test de Klein
- le test de Fomaret Glauken
- a) Test de Klein:

Le test de Klein n'est pas un test-stat au sens strict d'hypothèse mais il s'agit simplement d'un critère de présomption de multicollinéarité

soit le test modèle suivant:

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \dots + \beta_k x_{kt} + \varepsilon_t$$

$\forall t = 1, \dots, T$

Le test de Klein se fait en 3 étapes qui se présentent comme suit:

1. Estimer le modèle de régression et calculer son coeff de détermination R^2

2. calculer la matrice des coefficients de corrélation linéaire, prise 2 à 2 entre les variables explicatives telle que:

$$\begin{pmatrix} 1 & r_{x_1 x_2} & \dots & r_{x_1 x_k} \\ r_{x_2 x_1} & 1 & \dots & r_{x_2 x_k} \\ \vdots & \vdots & \ddots & \vdots \\ r_{x_k x_1} & r_{x_k x_2} & \dots & 1 \end{pmatrix}$$

3. Comparer le R^2 avec différents coefficients de corrélation.

* Règle de décision

si au moins 1 des coeff de corrélation au carré

$$r_{x_i x_j}^2 > R^2 \text{ alors il y a présomption}$$

de collinéarité.

b) Test de Fomaret et Glauken.

Ce test repose sur les hypothèses suivantes

H_0 : les variables explicatives sont orthogonales

contre

H_1 : les variables explicatives sont dépendantes

La première étape consiste à calculer le det de la matrice des coeff de corrélation entre les variables explicatives lorsque la valeur du det $D \rightarrow 0$ le risque de multicollinéarité est très grand par exemple pour un modèle à 2 variables explicatives si les deux séries sont parfaitement corrélées

$$D = \begin{vmatrix} 1 & r_{x_1 x_2} \\ r_{x_2 x_1} & 1 \end{vmatrix} = \begin{vmatrix} 1 & 1 \\ 1 & 1 \end{vmatrix} = 0$$

\Rightarrow Variables explicatives parfaitement corrélées

Dans le cas opposé si les variables explicatives sont orthogonales

$$D = \begin{vmatrix} 1 & r_{x_1 x_2} \\ r_{x_2 x_1} & 1 \end{vmatrix} = \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} = 1$$

2^{ème} étape:

La 2^{ème} étape consiste à effectuer un test de Kruskal en posant les hypothèses suivantes

H_0 : les variables sont orthogonales

$\Rightarrow D = 1$

contre

H_1 : les séries sont dépendantes

la statistique utilisée et sa loi de probabilité.

$$\chi^2_c = - \left[T - 1 - \frac{1}{6}(2K + 5) \right] \chi^2_{\frac{1}{2}(K/K)} \text{ avec } x \text{ by } D \text{ D.D.L}$$

T : taille de l'échantillon

$K = p + 1$ le nombre de paramètres à estimer du modèle.

* Règle de décision :

• si $\chi^2_c < \chi^2_\alpha$ alors H_0 est vraie
 \Rightarrow les séries sont orthogonales (indép)

• si $\chi^2_c > \chi^2_\alpha$ alors H_1 est vraie
 \Rightarrow présence de multicollinéarité

3) Les solutions au problème de multicollinéarité :

La solution la plus efficace reste, pas de la spécification du modèle, l'élimination des variables explicatives susceptibles de représenter les mêmes phénomènes et donc d'être corrélées avec elles. Sinon, il existe d'autres techniques permettant d'apporter des remèdes à la multicollinéarité. Parmi ces techniques on peut citer l'augmentation de la taille de l'échantillon, cette technique n'est efficace que si l'ajout d'observations diffère significativement de celle figurant déjà dans le modèle sinon il y aura

reconstruction de la multicollinéarité.
 * 3^{ème} Solution : la "Ridge Regression".
 c'est une technique proposée par Hoerl et Kennard (1970) qui consiste à traiter la multicollinéarité mécaniquement d'une façon numérique. Il s'agit de transformer la matrice $(X'X)$ en une matrice $(X'X + CI)$ où C est une constante choisie arbitrairement. Le problème vient du fait que la multicollinéarité se traduit par l'existence dans la matrice X d'une colonne représentant une combinaison linéaire d'autres colonnes ce qui affecte aussi la matrice $(X'X)$, on propose alors de détruire cette combinaison linéaire en ajoutant une constante aux éléments de la diagonale de la matrice $(X'X)$. Le principe de la "Ridge Regression" consiste ainsi à définir l'estimateur Ridge suivant :

$$\hat{\beta}_{RCCO} = (X'X)^{-1} X'Y$$

\Downarrow si multicollinéarité

Estimateur "Ridge"

$$\Rightarrow \hat{\beta}_{Ridge} = (X'X + CI)^{-1} X'Y$$

Il est à noter que l'estimateur "Ridge" peut être estimé en fct de l'estimateur $\hat{\beta}_{RCCO}$ telle que :

$$\hat{\beta}_{Ridge} = [I + C(X'X)^{-1}]^{-1} \hat{\beta}_{RCCO}$$

Par ailleurs, on peut montrer que $\hat{\beta}_{Ridge}$

est un estimateur biaisé mais convergent (déviance minimale par rapport à l'estimateur des MCO)
Dem Greene (2007)

II - Variables aléatoires :

explicative: problème d'endogénéité:

Lorsque l'hypothèse H_0 des MCO est rejetée, cela implique que les variables explicatives sont dépendantes du terme d'erreur dans ces conditions les estimateurs des MCO ne sont plus convergents

et il est nécessaire de recourir à un autre estimateur appelé **estimateur des variables instrumentales**.

1) Les causes de l'endogénéité:

Le problème d'endogénéité peut être dû à plusieurs causes parmi lesquelles on peut citer:

- existence d'un double sens de causalité entre la variable à expliquer (y) et la/ou les variable(s) explicative(s)

Exemple: croissance économique et développement financier
 - Omission de variables explicatives dans le modèle

$$y_t = \alpha + \beta_1 x_{1t} + \beta_2 x_{2t} + \varepsilon_t$$

$$\Rightarrow y_t = \alpha + \beta_1 x_{1t} + u_t$$

- Erreurs de mesure, évoqué dans le cas où on utilise les données d'enquête

2) La méthode des variables instrumentales (VI):

Soit X la matrice des variables explicatives et Z la matrice des variables instrumentales ou instruments.

Un instrument Z d'une variable X est une variable qui permet de mesurer X et telle que la covariance entre Z et le terme d'erreur est égale à 0 et Z doit être corrélée avec la variable X

$$\text{cov}(X, Z) \neq 0$$

(Le seul lien) Par ailleurs il faut que le seul lien qui peut exister entre Z et y est la variable X .

Autrement dit l'instrument Z ne peut pas être une variable explicative en elle-même.

Remarque:

On peut parfois utiliser plusieurs instruments pour une même variable endogène ceci permet d'avoir l'estimation de β_1 plus précise (explicative endo)

Dans tous les cas, il faut que le nombre de variables instrumentales soit supérieur ou égal au nombre de variables explicatives

~~l'estimateur des~~

(9) estimation
 - estimateur des variables instrumentales
 - V.I. donné par la formule suivante:

$$\hat{\beta}_{VI} = (Z'X)^{-1} Z'Y$$

avec :

X = matrice initiale des variables explicatives

Z = la matrice des variables instrumentales.

2) L'estimateur de la matrice de variances-covariances des paramètres estimés et défini tel que :

$$\hat{\sigma}_{\hat{\beta}_{VI}}^2 = \hat{\sigma}_\varepsilon^2 (Z'X)^{-1} Z'Z(X'Z)^{-1}$$

Une méthode particulière des variables instrumentales est les doublets moindres carrés ordinaires (DTCO) (2SLS)

qui vient à appliquer la MCO en deux étapes

soit le modèle suivant

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \beta_3 z_{3t} + \varepsilon_t$$

et soit (z_{3t}) un instrument de la variable (x_{3t}) . \Rightarrow On suppose que (x_{3t}) est endogène

Comme son nom l'indique la méthode des doublets moindres carrés ordinaires se fait en deux étapes :

La 1^{ère} étape consiste à régresser x_{3t} sur z_{3t} et les autres variables

explicatives du modèle :

$$x_{3t} = \alpha_0 + \alpha_1 z_{1t} + \alpha_2 x_{1t} + \alpha_3 x_{2t} + u_t$$

$$\Rightarrow \hat{x}_{3t} = \hat{\alpha}_0 + \hat{\alpha}_1 z_{1t} + \hat{\alpha}_2 x_{1t} + \hat{\alpha}_3 x_{2t}$$

$\Rightarrow \hat{x}_{3t}$ valeurs ajustées

2^{ème} étape :

Régression de (y_t) sur (\hat{x}_{3t}) et le reste des variables explicatives

$$y_t = \alpha_0 + \alpha_1 x_{1t} + \alpha_2 x_{2t} + \alpha_3 \hat{x}_{3t} + u_t$$

3) Test d'endogénéité d'Hausman

Le test d'endogénéité d'Hausman permet de détecter une éventuelle corrélation entre le terme d'erreur ε_t et une ou plusieurs variables explicatives x_{it} . Dans ce cas on ne peut plus utiliser l'estimateur des MCO qui devient non convergent.

Les hypothèses du test sont les suivantes :

$H_0 : \text{cov}(x_{it}, \varepsilon_t) = 0 \Rightarrow$ Pas de pb d'endogénéité

contre

$H_1 : \text{cov}(x_{it}, \varepsilon_t) \neq 0$

$\Rightarrow x_{it}$ est endogène

On retient les estimateurs VI

* Statistique du test et loi de probabilité :

$$H = (\hat{\beta}_{VI} - \hat{\beta}_{MCO})' \begin{bmatrix} \hat{\sigma}_{\beta_{VI}}^2 & \hat{\sigma}_{\beta_{VI}\beta_{MCO}} \\ \hat{\sigma}_{\beta_{VI}\beta_{MCO}} & \hat{\sigma}_{\beta_{MCO}}^2 \end{bmatrix}^{-1} (\hat{\beta}_{VI} - \hat{\beta}_{MCO})$$

individuelle pour le paramètre α_1 associé à la VI (α_1, α_2)

$$(\hat{\beta}_{VI} - \hat{\beta}_{MCO}) \sim \chi^2(K)$$

\Rightarrow Test de Student

K = nombre de variables explicatives + constante ($g+1$)

Hypothèses :

+ constante ($g+1$)

$H_0: \alpha_1 = 0$ contre $H_1: \alpha_1 \neq 0$

= rang de la matrice $(\hat{\sigma}_{\beta_{VI}}^2 - \hat{\sigma}_{\beta_{MCO}}^2)$

Statistique du test et loi :

sous H_0 vraie : $\frac{\alpha_1}{\hat{\sigma}_{\alpha_1}} \sim St(T-K)$

* Règle de décision :

1) si $H < \chi_{\alpha}^2$ alors on accepte H_0

* Règle de décision :

\Rightarrow pas de problème d'endogénéité
 \Rightarrow on retient l'estimation par MCO.

2) si $t_c = \left| \frac{\alpha_1}{\hat{\sigma}_{\alpha_1}} \right| < t_{\frac{\alpha}{2}}^{T-K}$ alors

H_0 est vraie $\Rightarrow \alpha_1$ est statistiquement non significatif $\Rightarrow g_1$ n'est corrélé avec la variable instrumentée

3) si $H > \chi_{\alpha}^2$ alors on accepte H_1
 \Rightarrow il y a un problème d'endogénéité
 \Rightarrow on doit retenir l'estimation par VI

4) si $t_c > t_{\frac{\alpha}{2}}^{T-K}$ donc H_1 est vraie

$\Rightarrow \alpha_1$ est statistiquement significatif
 \Rightarrow l'instrument g_1 est corrélé avec x_{1t}

4) Test sur l'instrument :

* Test de Student :

Ce test permet de vérifier que l'instrument est bien corrélé avec la variable instrumentée, il est appliqué au niveau de la première étape de la méthode des doubles moindres carrés ordinaires

$$x_{3t} = \alpha_0 + \alpha_1 g_t + \alpha_2 g_{1t} + \alpha_3 x_{1t} + u_t$$

* On estime les paramètres

* On effectue le test significativité

* Test de Sargan :

Ce test permet de vérifier que les instruments utilisés sont exogènes et donc valide. Ce test n'est valable que lorsque utilise au moins deux instruments ou une variable endogène.

* Hypothèses :

$H_0: \text{cov}(Z, \varepsilon) = 0 \Rightarrow Z$ est exogène
 contre

$H_1: \text{cov}(Z, \varepsilon) \neq 0 \Rightarrow Z$ n'est pas

Statistique du test et for:

$$S = T \cdot R^2$$

nb observations totales α

r = degré de liberté

= nombre d'instruments -

nombre de variables explicatives
endogènes.

R^2 = coefficient de détermination
relatif à la régression suivante (1^{ère} étape)

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \beta_3 x_{3t} + \varepsilon_t$$

1^{ère} étape:

$$y_t = \beta_0 + \beta_1 x_{1t} + \beta_2 x_{2t} + \beta_3 x_{3t} + \varepsilon_t$$

$$\Rightarrow \hat{y}_t \Rightarrow \hat{u}_{1t}$$

2^{ème} étape:

$$\hat{u}_{1t} = \alpha_0 + \alpha_1 z_{1t} + \alpha_2 z_{2t} + v_t$$

$$\Rightarrow R^2$$

* Règle de décision :

- si $S < \chi^2_{\alpha}$ alors H_0 est vraie

$\Rightarrow z$ est endogène.

- si $S > \chi^2_{\alpha}$ alors H_1 est vraie

$\Rightarrow z$ n'est pas exogène