



Notes d'économétrie

Auteur :

Pierre ROUILLARD

Last update : 24 mai 2023

Table des matières

1	Points à savoir	2
1.1	Interprétation du paramètre causal à estimer	2
1.2	Représentations linéaires & OLS	4
1.3	Résumé modèles	5
2	OLS - MCO	6
2.1	Rappels MCO	6

1 Points à savoir

1.1 Interprétation du paramètre causal à estimer

Selon le modèle considéré il est possible ou non d'avoir une interprétation quantitative directe et/ou qualitative du paramètre causal à estimer β_0 .

Définition de l'effet marginal de X_k sur Y : $\frac{\partial E[Y|X=x]}{\partial x_k}$

→ **Modèle linéaire** :

Comme toujours par la suite on considère l'analyse *toutes choses égales par ailleurs*, sur la population considérée...

Interprétation : la variable d'intérêt est Y et en l'absence de puissance ou d'interactions on peut interpréter quantitativement β_0 sur la variable d'intérêt. Le paramètre d'intérêt est l'effet marginal de X_k sur Y qui vaut bien β_{0k} lorsque X_k apparaît simplement dans le modèle. C'est justement pour cela qu'on peut bien interpréter directement quantitativement les coefficients de β_0 !

→ **Modèle binaire** :

$$E[Y|X] = P(Y = 1|X) = F(X'\beta_0)$$

$$E[Y|X] = F(X'\beta_0) \iff Y = \mathbb{1}(Y^* \geq s) : Y^* = X'\beta_0 + \varepsilon \quad \varepsilon \perp\!\!\!\perp X$$

Interprétation quantitative directe : la variable d'intérêt est Y et Y^* n'est qu'une variable latente qui n'a pas forcément de sens quantitatif précis. Le paramètre d'intérêt est l'effet marginal de X_k sur la variable d'intérêt, ici Y . Les coefficients de β_0 concernant Y^* on ne peut donc pas directement interpréter quantitativement ces derniers sur la variable d'intérêt Y . De plus, l'effet marginal de X_k sur Y est différent de β_{0k} : c'est pour cela qu'on ne peut avoir d'interprétation quantitative des coefficients de β_0 !

Interprétation qualitative : en revanche le signe de l'effet marginal de X_k sur Y , i.e. effet positif ou négatif sur $P(Y = 1|X)$, est donné par le signe de β_{0k} .

On peut en revanche comparer quantitativement le ratio des effets marginaux des variables i et j qui vaut $\widehat{\beta}_i/\widehat{\beta}_j$. L'effet sur la proba d'être ... de la variable i est *<quantitativement>* ... que l'effet de la variable $j \iff$ regarder le rapport $\widehat{\beta}_i/\widehat{\beta}_j$.

→ **Modèle de censure / Tobit1** :

1 seul mécanisme détermine la valeur de Y et si on observe la variable d'intérêt ou non. Deux cas sont à distinguer :

⇒ **Données censurées** : la variable d'intérêt est Y^* qui peut ne pas être observée au dessous d'un seuil causant un problème de censure. Le paramètre d'intérêt est l'effet marginal de X_k sur la variable d'intérêt Y^* , qui vaut bien β_{0k} lorsque X_k apparaît simplement dans

le modèle linéaire de Y^* . Ainsi, la variable Y^* ayant un sens quantitatif et malgré la censure liée aux problèmes d'observation on peut bien interpréter quantitativement β_0 sur la variable d'intérêt.

⇒ **Solution en coin** : la variable d'intérêt est bien Y alors que la variable Y^* est une variable latente potentiellement dépourvue de sens quantitatif. Typiquement un pb d'optimisation du consommateur où Y^* mesure l'utilité optimale (en nombre de biens) de consommation d'un bien donné : donc potentiellement négatif. Et Y représente le nombre d'unités effectivement consommées. Les coefficients de β_0 concernant Y^* qui n'as pas de sens quantitatif précis : on ne peut pas interpréter quantitativement les coefficients de β_0 sur la variable d'intérêt Y . Les paramètres d'intérêt sont les effets marginaux : le total $\frac{\partial E[Y|X=x]}{\partial x_k}$ (marge extensive et intensive) et $\frac{\partial E[Y|Y>0, X=x]}{\partial x_k}$ (marge intensive seulement). Ces paramètres sont tous les deux différents de β_{0k} ce qui explique le manque d'interprétation quantitative des coefficients de β_0 .

↪ **Modèle de sélection / Tobit2 :**

Ici on a bien deux processus différents : un qui détermine Y^* **et un autre** qui détermine si on observe cette valeur ou non i.e. modèle sur D .

Interprétation **quantitative directe** : il y a un problème d'observation des données, on observe $Y = D.Y^*$ mais la variable d'intérêt est bien Y^* (variable potentielle qui existe pour tous les *individus*). Par conséquent Y^* suivant un modèle linéaire, les paramètres d'intérêts sont les effets marginaux des variables explicatives sur la variable d'intérêt Y^* et les coefficients de β_0 sont toujours interprétables quantitativement.

1.2 Représentations linéaires & OLS

↪ Représentation non causale - Projection linéaire :

Sous conditions de moments¹, on a toujours par construction/définition de la représentation linéaire théorique (=projection linéaire orthogonale) orthogonalité des *résidus* de cette représentation non causale avec les régresseurs = pas une hypothèse mais une conséquence.

$Y = X' \cdot \tilde{\beta} + \tilde{\varepsilon}$, $E[X\tilde{\varepsilon}] = 0$ toujours définissable sous conditions de moments.

- $\widehat{\beta_{OLS}}$ estime toujours $\tilde{\beta}$: $\widehat{\beta_{OLS}} \xrightarrow{n \rightarrow +\infty} \tilde{\beta}$
- $X' \cdot \tilde{\beta}$ meilleure prédiction linéaire de Y par X : $\tilde{\beta}$ solution MSE.

$$\widehat{\beta_{OLS}} \in \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - X'_i \cdot \beta)^2 \quad \longleftrightarrow \quad \tilde{\beta} \in \underset{\beta}{\operatorname{argmin}} E[(Y - X' \cdot \beta)^2]$$

↪ Représentation causale :

La représentation causale fait intervenir le paramètre causal β_0 qu'on cherche à estimer : dans cette représentation le *terme d'erreur* n'est pas automatiquement orthogonal au régresseur.

$Y = X' \cdot \beta_0 + \varepsilon$, $E[X\varepsilon] \stackrel{?}{=} 0$

- β_0 paramètre causal à estimer.
- ε résidu : agrège les facteurs inobservés qui affectent Y.

Le terme d'erreur ε capte l'hétérogénéité inobservée, i.e capte les déterminants inobservés qui affectent la variable d'intérêt Y : deux individus avec les mêmes variables explicatives auront néanmoins la plupart du temps des variables expliquées différentes.

Avoir orthogonalité (\rightarrow indépendance) entre régresseurs et terme d'erreur est une hypothèse ! C'est **l'hypothèse d'exogénéité**.

↪ Lien :

Sans l'hypothèse d'exogénéité pour la représentation causale, les deux représentations diffèrent et l'estimateur OLS ne permet pas d'identifier le paramètre causal d'intérêt.

En revanche avec hypothèse d'exogénéité les deux représentations coïncident et $\tilde{\beta} = \beta_0$: $\widehat{\beta_{OLS}}$ qui estime toujours $\tilde{\beta}$ est donc un estimateur consistant de β_0 .

En dehors des expériences contrôlées les variables explicatives peuvent parfois être corrélées aux facteurs inobservables et pb d'endogénéité $E[X\varepsilon] \neq 0$.

1. $E[Y^2] < +\infty$, $E[\|X\|^2] < +\infty$ et $E[XX']$ inversible = de rang plein. En particulier : **any level of correlation between covariates except perfect colinearity** : composantes de X linéairement indépendantes mais *n'exclut pas* qu'elles soient corrélées. Si le modèle a une constante et une variable catégorielle il faut exclure une des modalités.

1.3 Résumé modèles

2 OLS - MCO

2.1 Rappels MCO