

Économétrie

Cours et exercices corrigés

Régis Bourbonnais

9^e édition

DUNOD

Le pictogramme qui figure ci-contre mérite une explication. Son objet est d'alerter le lecteur sur la menace que représente pour l'avenir de l'écrit, particulièrement dans le domaine de l'édition technique et universitaire, le développement massif du photocopillage.

Le Code de la propriété intellectuelle du 1^{er} juillet 1992 interdit en effet expressément la photocopie à usage collectif sans autori-

sation des ayants droit. Or, cette pratique s'est généralisée dans les établissements

d'enseignement supérieur, provoquant une baisse brutale des achats de livres et de revues, au point que la possibilité même pour

les auteurs de créer des œuvres nouvelles et de les faire éditer correctement est aujourd'hui menacée. Nous rappelons donc que toute reproduction, partielle ou totale, de la présente publication est interdite sans autorisation de l'auteur, de son éditeur ou du Centre français d'exploitation du

droit de copie (CFC, 20, rue des Grands-Augustins, 75006 Paris).



© Dunod, 2015

5 rue Laromiguière, 75005 Paris

www.dunod.com

ISBN 978-2-10-072151-1

Le Code de la propriété intellectuelle n'autorisant, aux termes de l'article L. 122-5, 2° et 3° a), d'une part, que les « copies ou reproductions strictement réservées à l'usage privé du copiste et non destinées à une utilisation collective » et, d'autre part, que les analyses et les courtes citations dans un but d'exemple et d'illustration, « toute représentation ou reproduction intégrale ou partielle faite sans le consentement de l'auteur ou de ses ayants droit ou ayants cause est illicite » (art. L. 122-4).

Cette représentation ou reproduction, par quelque procédé que ce soit, constituerait donc une contrefaçon sanctionnée par les articles L. 335-2 et suivants du Code de la propriété intellectuelle.

Table des matières

Avant-propos	IX
1. Qu'est-ce que l'économétrie ?	1
I. La notion de modèle	1
A. Définition	1
B. La construction des modèles en économétrie	2
II. Le rôle de l'économétrie	5
A. L'économétrie comme validation de la théorie	5
B. L'économétrie comme outil d'investigation	5
III. La théorie de la corrélation	6
A. Présentation générale	6
B. Mesure et limite du coefficient de corrélation	8
2. Le modèle de régression simple	13
I. Présentation du modèle	13
A. Exemple introductif	13
B. Rôle du terme aléatoire	14
C. Conséquences du terme aléatoire	16
II. Estimation des paramètres	17
A. Modèle et hypothèses	17
B. Formulation des estimateurs	18
C. Les différentes écritures du modèle : erreur et résidu	21
D. Propriétés des estimateurs	22
III. Conséquences des hypothèses : construction des tests	24
A. Hypothèse de normalité des erreurs	24
B. Conséquences de l'hypothèse de normalité des erreurs	25
C. Test bilatéral, test unilatéral et probabilité critique d'un test	27
IV. Équation et tableau d'analyse de la variance	33
A. Équation d'analyse de la variance	33
B. Tableau d'analyse de la variance	34
V. La prévision dans le modèle de régression simple	39

3. Le modèle de régression multiple	47
I. Le modèle linéaire général	47
A. Présentation	47
B. Forme matricielle	48
II. Estimation et propriétés des estimateurs	49
A. Estimation des coefficients de régression	49
B. Hypothèses et propriétés des estimateurs	51
C. Équation d'analyse de la variance et qualité d'un ajustement	54
III. Les tests statistiques	59
A. Le rôle des hypothèses	59
B. Construction des tests	60
C. Tests sur les résidus : valeur anormale, effet de levier et point d'influence	62
IV. L'analyse de la variance	67
A. Construction du tableau d'analyse de la variance et test de signification globale d'une régression	67
B. Autres tests à partir du tableau d'analyse de la variance	68
C. Généralisation des tests par analyse de la variance	73
V. L'utilisation de variables indicatrices	75
A. Constitution et finalités des variables indicatrices	75
B. Exemples d'utilisation	76
VI. La prévision à l'aide du modèle linéaire général et la régression récursive	81
A. Prédiction conditionnelle	81
B. Fiabilité de la prévision et intervalle de prévision	82
C. Les tests de stabilité par la régression récursive	84
D. Le test de spécification de Ramsey	86
VII. Exercices récapitulatifs	90
<i>Annexe</i>	102
A) Interprétation géométrique de la méthode des moindres carrés	102
B) Résolution de l'exercice 1 par des logiciels informatiques de régression multiple	103
C) Estimation de la variance de l'erreur	105
4. Multicolinéarité et sélection du modèle optimal	107
I. Corrélations partielles	107
A. Exemple introductif	107
B. Généralisation de la notion de corrélation partielle	108
II. Relation entre coefficients de corrélation simple, partielle et multiple	112

III. Multicolinéarité : conséquences et détection	114
A. Conséquences de la multicolinéarité	114
B. Tests de détection d'une multicolinéarité	115
C. Comment remédier à la multicolinéarité ?	118
IV. Sélection du modèle optimal	119
5. Problèmes particuliers : la violation des hypothèses	125
I. L'autocorrélation des erreurs	125
A. Présentation du problème	125
B. L'estimateur des Moindres Carrés Généralisés (MCG)	126
C. Les causes et la détection de l'autocorrélation des erreurs	127
D. Les procédures d'estimation en cas d'autocorrélation des erreurs	134
II. L'hétéroscédasticité	142
A. Présentation du problème	142
B. Correction de l'hétéroscédasticité	144
C. Tests de détection de l'hétéroscédasticité	147
D. Autre test d'hétéroscédasticité : le test ARCH	153
III. Modèles à erreurs sur les variables	154
A. Conséquences lorsque les variables sont entachées d'erreurs	154
B. La méthode des variables instrumentales	155
C. Le test d'exogénéité d'Hausman	156
D. La méthode des moments généralisée	157
6. Les modèles non linéaires	165
I. Les différents types de modèles non linéaires	165
A. Les fonctions de type exponentiel	165
B. Les modèles de diffusion	168
II. Méthodes d'estimation des modèles non linéaires	170
A. Initiation aux méthodes d'estimation non linéaires	170
B. Exemples d'application	172
7. Les modèles à décalages temporels	177
I. Les modèles linéaires autorégressifs	177
A. Formulation générale	177
B. Test d'autocorrélation et méthodes d'estimation	178
II. Les modèles à retards échelonnés	183
A. Formulation générale	183
B. Détermination du nombre de retards	184

C. Distribution finie des retards	188
D. Distribution infinie des retards	192
III. Deux exemples de modèles dynamiques	198
A. Le modèle d'ajustement partiel	198
B. Le modèle d'anticipations adaptatives	199
8. Introduction aux modèles à équations simultanées	217
I. Équations structurelles et équations réduites	218
A. Exemple introductif	218
B. Le modèle général	220
C. Cas particulier : les modèles récurrents	221
II. Le problème de l'identification	221
A. Restrictions sur les coefficients	221
B. Conditions d'identification	222
III. Les méthodes d'estimation	223
A. Les moindres carrés indirects	223
B. Les doubles moindres carrés	223
C. Autres méthodes d'estimation	224
<i>Annexe</i>	236
<i>Identification : les conditions de rang</i>	236
9. Éléments d'analyse des séries temporelles	239
I. Stationnarité	239
A. Définition et propriétés	239
B. Fonctions d'autocorrélation simple et partielle	240
C. Tests de « bruit blanc » et de stationnarité	241
II. La non-stationnarité et les tests de racine unitaire	245
A. La non-stationnarité : les processus TS et DS	245
B. Les tests de racine unitaire et la stratégie séquentielle de test	248
III. Les modèles ARIMA	256
A. Typologie des modèles AR, MA et ARMA	256
B. L'extension aux processus ARIMA et SARIMA	259
IV. La méthode de Box et Jenkins	260
A. Recherche de la représentation adéquate : l'identification	260
B. Estimation des paramètres	261
C. Tests d'adéquation du modèle et prévision	262
10. La modélisation VAR	275
I. Représentation d'un modèle VAR	276
A. Exemple introductif	276

B. La représentation générale	277
C. La représentation ARMAX	278
II. Estimation des paramètres	279
A. Méthode d'estimation	279
B. Détermination du nombre de retards	279
C. Prévision	280
III. Dynamique d'un modèle VAR	284
A. Représentation VMA d'un processus VAR	284
B. Analyse et orthogonalisation des « chocs »	285
C. Décomposition de la variance	288
D. Choix de l'ordre de décomposition	288
IV. La causalité	292
A. Causalité au sens de Granger	292
B. Causalité au sens de Sims	293
11. La cointégration et le modèle à correction d'erreur	297
I. Exemples introductifs	297
II. Le concept de cointégration	299
A. Propriétés de l'ordre d'intégration d'une série	299
B. Conditions de cointégration	301
C. Le modèle à correction d'erreur (ECM)	301
III. Cointégration entre deux variables	302
A. Test de cointégration entre deux variables	303
B. Estimation du modèle à correction d'erreur	303
IV. Généralisation à k variables	306
A. La cointégration entre k variables	306
B. Estimation du modèle à correction d'erreur	307
C. Le modèle à correction d'erreur vectoriel	308
D. Tests de relation de cointégration	310
E. Test d'exogénéité faible	313
F. Synthèse de la procédure d'estimation	314
12. Introduction à l'économétrie des variables qualitatives	319
I. Les problèmes et les conséquences de la spécification binaire	320
II. Les modèles de choix binaires	322
A. Le modèle linéaire sur variable latente	322
B. Les modèles Probit et Logit	323
C. Interprétation des résultats et tests statistiques	325
III. Les modèles à choix multiples	330
A. Les modèles Probit et Logit ordonnés	331

B. Le modèle de choix multiples non ordonné : le Logit multinomial	335
IV. Les modèles à variable dépendante limitée : le modèle Tobit	337
A. Le modèle Tobit simple : modèle de régression tronqué ou censuré	338
B. Estimation et interprétation des résultats	340
13. Introduction à l'économétrie des données de panel	345
I. Présentation des modèles à données de panel	346
A. Spécificités des données de panel	346
B. La méthode SUR	347
C. Le modèle linéaire simple	348
II. Les tests d'homogénéité	349
A. Procédure séquentielle de tests	349
B. Construction des tests	350
III. Spécifications et estimations des modèles à effets individuels	355
A. Le modèle à effets fixes individuels	355
B. Le modèle à effets aléatoires	357
C. Effets fixes ou effets aléatoires ? Le test d'Hausman	358
Liste des exercices	363
Tables statistiques	367
Bibliographie	375
Index	379

Avant-propos

Cette neuvième édition est enrichie de nouveaux exercices et des développements les plus récents de l'économétrie. Ce livre couvre tous les champs de l'économétrie : régression simple et multiple, violation des hypothèses (hétéroscédasticité, autocorrélation des erreurs, variables explicatives aléatoires), modèle à décalage, analyse des séries temporelles, tests de racine unitaire, équations multiples, VAR, cointégration, VECM, économétrie des variables qualitatives et des données de panel...

Sur l'ensemble de ces thèmes, ce livre vous propose un cours, des exercices corrigés, et une présentation des logiciels d'économétrie les plus répandus. Souhaitons qu'il corresponde à votre attente.

En effet, nous avons voulu, par une alternance systématique de cours et d'exercices, répondre à un besoin pédagogique qui est de mettre rapidement en pratique les connaissances théoriques et ainsi, d'utiliser de manière opérationnelle les acquis du cours ; les exercices sont repérés grâce à un bandeau grisé. De surcroît, le recours à des logiciels¹, lors de la résolution des exercices, permet une découverte de ces outils et donne une dimension pratique que recherchent l'étudiant et le praticien.

Afin que le lecteur puisse lui-même refaire les exercices, les données utilisées (sous format Excel, ASCII, RATS et Eviews) ainsi que les programmes de traitement « Batch » de Eviews ou de RATS sont disponibles gratuitement par téléchargement sur le serveur web :

<http://regisbourbonnais.dauphine.fr>

Pour chaque exercice faisant appel à un fichier de données, le nom du fichier est cité en tête de l'exercice et repéré par l'icône suivante : 

Nous avons voulu faire de ce manuel un livre d'apprentissage facilement accessible ; c'est pourquoi les démonstrations les plus complexes font l'objet de renvois à une bibliographie plus spécialisée. Cependant, il convient de préciser que l'économétrie fait appel à des notions d'algèbre linéaire et d'induction statistique qu'il est souhaitable de connaître.

1. Trois logiciels sont utilisés : EXCEL (© Microsoft), RATS (© Var Econometrics version 3 et Estima version 4), Eviews (© Quantitative Micro Software). Nous recommandons aussi particulièrement le logiciel GRETL (<http://gretl.sourceforge.net>) qui est un logiciel d'économétrie gratuit, complet et très facile d'apprentissage.

Dans le terme « économétrie » figure la racine du mot « économie » car son utilisation est surtout destinée à des fins de traitement de données économiques ; cependant, d'autres domaines tels que la finance, la recherche agronomique, la médecine, etc., font maintenant le plus souvent appel à ces techniques.

Ce livre s'adresse en premier lieu aux étudiants (sciences économiques, gestion, écoles de commerce et d'ingénieurs, etc.) dont la formation requiert une connaissance de l'économétrie. Gageons qu'il sera un support de cours indispensable et un allié précieux pour préparer les séances de travaux dirigés.

N'oublions pas cependant le praticien de l'économétrie (économiste d'entreprise, chercheur, etc.) qui, confronté à des problèmes d'estimation statistique, trouvera dans ce livre les réponses pratiques aux différentes questions qu'il peut se poser.

Enfin, j'exprime toute ma gratitude à toutes les personnes – collègues et étudiants – qui ont eu la gentillesse de me faire des commentaires et dont les conseils et suggestions contribuent à la qualité pédagogique de ce livre. Je reste, bien entendu, le seul responsable des erreurs qui subsisteraient¹.

1. Les lecteurs souhaitant faire des commentaires ou des remarques peuvent me contacter : Régis Bourbonnais, université de Paris-Dauphine, place du Maréchal de Lattre de Tassigny, 75775 Paris Cedex 16, E-mail : regis.bourbonnais@dauphine.fr

1. Qu'est-ce que l'économétrie ?

Ce premier chapitre est consacré à la présentation de l'économétrie et à sa liaison avec la théorie économique. Nous abordons tout d'abord la notion de modèle ainsi que les différentes étapes de la modélisation. L'apport de l'économétrie en tant qu'outil de validation est étudié en II. Enfin, la théorie de la corrélation – fondement de l'économétrie – fait l'objet de la section III.

I. La notion de modèle

A. Définition

Il est délicat de fournir une définition unique de la notion de modèle¹. Dans le cadre de l'économétrie, nous pouvons considérer qu'un modèle consiste en une *présentation formalisée d'un phénomène* sous forme d'équations dont les variables sont des grandeurs économiques. L'objectif du modèle est de représenter les traits les plus marquants d'une réalité qu'il cherche à styliser. Le modèle est donc l'outil que le modélisateur utilise lorsqu'il cherche à comprendre et à expliquer des phénomènes. Pour ce faire, il émet des hypothèses et explicite des relations.

1. La notion de modèle est relative au point de vue auquel nous nous plaçons : la physique, l'épistémologie...

► *Pourquoi des modèles ?*

► *Nombreux sont ceux – sociologues, économistes ou physiciens – qui fondent leurs analyses ou leurs jugements sur des raisonnements construits et élaborés. Ces constructions réfèrent implicitement à des modèles ; alors pourquoi ne pas expliciter clairement les hypothèses et les relations au sein d'un modèle ?*

Le modèle est donc une présentation schématique et partielle d'une réalité naturellement plus complexe. Toute la difficulté de la modélisation consiste à ne retenir que la ou les représentations intéressantes pour le problème que le modélisateur cherche à expliciter. Ce choix dépend de la nature du problème, du type de décision ou de l'étude à effectuer. La même réalité peut ainsi être formalisée de diverses manières en fonction des objectifs.

B. La construction des modèles en économétrie

Dans les sciences sociales, et particulièrement en économie, les phénomènes étudiés concernent le plus souvent des comportements afin de mieux comprendre la nature et le fonctionnement des systèmes économiques. L'objectif du modélisateur est, dans le cadre de l'économétrie et au travers d'une mesure statistique, de permettre aux agents économiques (ménages, entreprises, État...) d'intervenir de manière plus efficace. La construction d'un modèle comporte un certain nombre d'étapes qui sont toutes importantes. En effet, en cas de faiblesse d'un des « maillons », le modèle peut se trouver invalidé pour cause d'hypothèses manquantes, de données non représentatives ou observées avec des erreurs, etc. Examinons les différentes étapes à suivre lors de la construction d'un modèle, ceci à partir de l'exemple du modèle keynésien simplifié.

1) Référence à une théorie

Une théorie s'exprime au travers d'hypothèses auxquelles le modèle fait référence. Dans la théorie keynésienne, quatre propositions sont fondamentales :

1. la consommation et le revenu sont liés ;
2. le niveau d'investissement privé et le taux d'intérêt sont également liés ;
3. il existe un investissement autonome public ;
4. enfin, le produit national est égal à la consommation plus l'investissement privé et public.

2) Formalisation des relations et choix de la forme des fonctions

À partir des propositions précédentes, nous pouvons construire des relations :

1. la consommation est fonction du revenu : $C = f(Y)$ avec $f' > 0$;
2. l'investissement privé dépend du taux d'intérêt : $I = g(r)$ avec $g' < 0$;
3. il existe un investissement autonome public : \bar{I} ;
4. enfin, le produit national (ou le revenu national) est égal à la consommation plus l'investissement : $Y \equiv C + I + \bar{I}$.

À ce stade, nous n'avons postulé aucune forme particulière en ce qui concerne les fonctions f et g . Ainsi, bien que des considérations d'ordre théorique nous renseignent sur le signe des dérivées, il existe une multitude de fonctions de formes très différentes et ayant des signes de dérivées identiques, par exemple $C = a_0 + a_1 Y$ et $C = a_0 Y^{a_1}$. Cependant ces deux relations ne reflètent pas le même comportement ; une augmentation du revenu provoque un accroissement proportionnel pour la première relation, alors que, dans la seconde, l'effet s'estompe avec l'augmentation du revenu (si $0 < a_1 < 1$). Nous appelons « forme fonctionnelle » ce choix (arbitraire ou fondé) de spécification précise du modèle. Dans notre exemple, le modèle explicité s'écrit :

$$\begin{aligned} C &= a_0 + a_1 Y && \text{avec } a_0 > 0 \text{ et } 0 < a_1 < 1 \\ & && a_1 = \text{propension marginale à consommer} \\ & && \text{et } a_0 = \text{consommation incompressible ;} \\ I &= b_0 + b_1 r && \text{avec } b_0 > 0 \text{ et } b_1 < 0 ; \\ Y &\equiv C + I + \bar{I} \end{aligned}$$

Les deux premières équations reflètent des relations de comportements alors que la troisième est une identité (aucun paramètre n'est à estimer).

3) Sélection et mesure des variables

Le modèle étant spécifié, il convient de collecter les variables représentatives des phénomènes économiques. Ce choix n'est pas neutre et peut conduire à des résultats différents, les questions qu'il convient de se poser sont par exemple :

- *Faut-il raisonner en euros constants ou en euros courants ?*
- *Les données sont-elles brutes ou CVS¹ ?*
- *Quel taux d'intérêt faut-il retenir (taux au jour le jour, taux directeur de la Banque Centrale Européenne,...) ? etc.*

1. Corrigées des Variations Saisonnières.

Nous distinguons plusieurs types de données selon que le modèle est spécifié en :

- *série temporelle* : c'est le cas le plus fréquent en économétrie, il s'agit de variables observées à intervalles de temps réguliers (la consommation annuelle, totale France, exprimée en euros courants sur 20 ans) ;
- *coupe instantanée* : les données sont observées au même instant et concernent les valeurs prises par la variable pour un groupe d'individus¹ spécifiques (consommation observée des agriculteurs pour une année donnée) ;
- *panel* : la variable représente les valeurs prises par un échantillon d'individus à intervalles réguliers (la consommation d'un échantillon de ménages de la région parisienne sur 20 ans) ;
- *cohorte* : très proches des données de panel, les données de cohorte se distinguent de la précédente par la constance de l'échantillon, les individus sondés sont les mêmes d'une période sur l'autre.

4) Décalages temporels

Dans le cadre de modèle spécifié en séries temporelles, les relations entre les variables ne sont pas toujours synchrones mais peuvent être décalées dans le temps. Nous pouvons concevoir que la consommation de l'année t est expliquée par le revenu de l'année $t - 1$ et non celui de l'année t . Pour lever cette ambiguïté, il est d'usage d'écrire le modèle en le spécifiant à l'aide d'un indice de temps : $C_t = a_0 + a_1 Y_{t-1}$. La variable Y_{t-1} est appelée « variable endogène retardée ».

- On appelle « variable exogène » une variable dont les valeurs sont prédéterminées, et « variable endogène » une variable dont les valeurs dépendent des variables exogènes.

5) Validation du modèle

La dernière étape est celle de la validation² du modèle :

- Les relations spécifiées sont-elles valides ?
- Peut-on estimer avec suffisamment de précision les coefficients ?
- Le modèle est-il vérifié sur la totalité de la période ?
- Les coefficients sont-ils stables ? Etc.

À toutes ces questions, les techniques économétriques s'efforcent d'apporter des réponses.

1. Le terme d'individu est employé au sens statistique, c'est-à-dire comme un élément d'une population : une personne, une parcelle de terre...

2. Validation, c'est-à-dire en conformité avec les données disponibles.

II. Le rôle de l'économétrie

A. L'économétrie comme validation de la théorie

L'économétrie est un outil à la disposition de l'économiste qui lui permet d'infirmer ou de confirmer les théories qu'il construit. Le théoricien postule des relations ; l'application de méthodes économétriques fournit des estimations sur la valeur des coefficients ainsi que la précision attendue.

Une question se pose alors : pourquoi estimer ces relations, et les tester statistiquement ? Plusieurs raisons incitent à cette démarche : tout d'abord cela force l'individu à établir clairement et à estimer les interrelations sous-jacentes. Ensuite, la confiance aveugle dans l'intuition peut mener à l'ignorance de liaisons importantes ou à leur mauvaise utilisation. De plus, des relations marginales mais néanmoins explicatives, qui ne sont qu'un élément d'un modèle global, doivent être testées et validées afin de les mettre à leur véritable place. Enfin, il est nécessaire de fournir, en même temps que l'estimation des relations, une mesure de la confiance que l'économiste peut avoir en celles-ci, c'est-à-dire la précision que l'on peut en attendre. Là encore, l'utilisation de méthodes purement qualitatives exclut toute mesure quantitative de la fiabilité d'une relation.

B. L'économétrie comme outil d'investigation

L'économétrie n'est pas seulement un système de validation, mais également un outil d'analyse. Nous pouvons citer quelques domaines où l'économétrie apporte une aide à la modélisation, à la réflexion théorique ou à l'action économique par :

- la mise en évidence de relations entre des variables économiques qui n'étaient pas *a priori* évidentes ou pressenties ;
- l'induction statistique ou l'inférence statistique consiste à inférer, à partir des caractéristiques d'un échantillon, les caractéristiques d'une population. Elle permet de déterminer des intervalles de confiance pour des paramètres du modèle ou de tester si un paramètre est significativement¹ inférieur, supérieur ou simplement différent d'une valeur fixée ;

1. Au sens statistique, c'est-à-dire avec un seuil (risque d'erreur à ne pas dépasser, souvent 5 %).

- la simulation qui mesure l'impact de la modification de la valeur d'une variable sur une autre ($\Delta C_t = a_1 \Delta Y_t$) ;
- la prévision¹, par l'utilisation de modèles économétriques, qui est utilisée par les pouvoirs publics ou l'entreprise afin d'anticiper et éventuellement de réagir à l'environnement économique.

Dans cet ouvrage, nous nous efforcerons de montrer, à l'aide d'exemples, les différentes facettes de l'utilisation des techniques économétriques dans des contextes et pour des objectifs différents.

III. La théorie de la corrélation

A. Présentation générale

Lorsque deux phénomènes ont une évolution commune, nous disons qu'ils sont « corrélés ». La corrélation simple mesure le degré de liaison existant entre ces deux phénomènes représentés par des variables. Si nous cherchons une relation entre trois variables ou plus, nous ferons appel alors à la notion de corrélation multiple.

Nous pouvons distinguer la corrélation linéaire, lorsque tous les points du couple de valeurs (x, y) des deux variables semblent alignés sur une droite, de la corrélation non linéaire lorsque le couple de valeurs se trouve sur une même courbe d'allure quelconque.

Deux variables peuvent être :

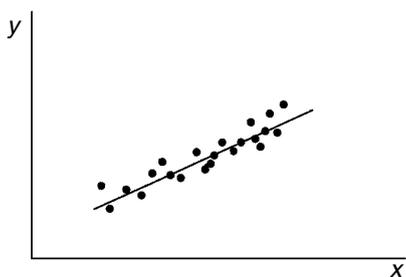
- en corrélation positive ; on constate alors une augmentation (ou diminution, ou constance) simultanée des valeurs des deux variables ;
- en corrélation négative, lorsque les valeurs de l'une augmentent, les valeurs de l'autre diminuent ;
- non corrélées, il n'y a aucune relation entre les variations des valeurs de l'une des variables et les valeurs de l'autre.

Le tableau 1, en croisant les critères de linéarité et de corrélation, renvoie à une représentation graphique.

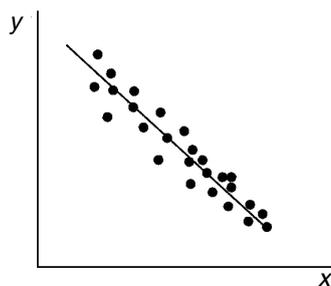
1. Pour découvrir l'utilisation de l'économétrie à des fins de prévision de ventes, voir Bourbonnais R. et Usunier J. C. (2013).

Tableau 1 – Linéarité et corrélation

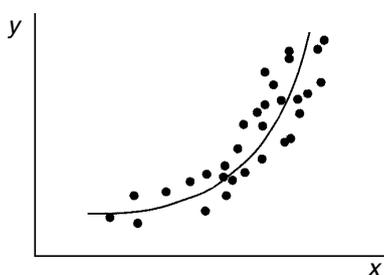
	Corrélation positive	Corrélation négative	Absence de corrélation
Relation linéaire	Graphe 1	Graphe 2	Graphe 5
Relation non linéaire	Graphe 3	Graphe 4	Graphe 5



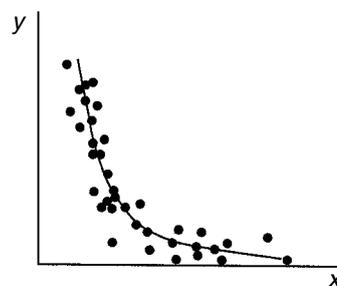
Graphe 1



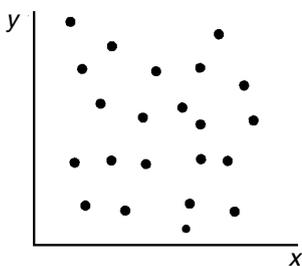
Graphe 2



Graphe 3



Graphe 4



Graphe 5

B. Mesure et limite du coefficient de corrélation

1) Le coefficient de corrélation linéaire

La représentation graphique ne donne qu'une « impression » de la corrélation entre deux variables sans donner une idée précise de l'intensité de la liaison, c'est pourquoi nous calculons une statistique appelée *coefficient de corrélation linéaire simple*, noté $r_{x,y}$. Il est égal à :

$$r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad [1]$$

avec :

$\text{Cov}(x,y)$ = covariance entre x et y ;
 σ_x et σ_y = écart type de x et écart type de y ;
 n = nombre d'observations.

En développant la formule [1], il vient :

$$r_{x,y} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sqrt{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \sqrt{n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i\right)^2}} \quad [2]$$

On peut démontrer que, par construction ce coefficient reste compris entre -1 et 1 :

- proche de 1 , les variables sont corrélées positivement ;
- proche de -1 , les variables sont corrélées négativement ;
- proche de 0 , les variables ne sont pas corrélées.

Dans la pratique, ce coefficient est rarement très proche de l'une de ces trois bornes et il est donc difficile de proposer une interprétation fiable à la simple lecture de ce coefficient. Ceci est surtout vrai en économie où les variables sont toutes plus au moins liées entre elles. De plus, il n'est calculé qu'à partir d'un échantillon d'observations et non pas sur l'ensemble des valeurs. On appelle $\rho_{x,y}$

ce coefficient empirique qui est une estimation du coefficient vrai $r_{x,y}$. La théorie des tests statistiques nous permet de lever cette indétermination.

Soit à tester l'hypothèse $H_0 : r_{x,y} = 0$, contre l'hypothèse $H_1 : r_{x,y} \neq 0$.

Sous l'hypothèse H_0 , nous pouvons démontrer que $\frac{\rho_{x,y}}{\sqrt{\frac{(1 - \rho_{x,y}^2)}{n - 2}}}$ suit une loi

de Student à $n - 2$ degrés de liberté¹. Nous calculons alors une statistique, appelé le t de Student empirique :

$$t^* = \frac{|\rho_{x,y}|}{\sqrt{\frac{(1 - \rho_{x,y}^2)}{n - 2}}} \quad [3]$$

Si $t^* > t_{n-2}^{\alpha/2}$ valeur lue dans une table de Student² au seuil $\alpha = 0,05$ (5 %) à $n - 2$ degrés de liberté³, nous rejetons l'hypothèse H_0 , le coefficient de corrélation est donc significativement différent de 0 ; dans le cas contraire, l'hypothèse d'un coefficient de corrélation nul est acceptée. La loi de Student étant symétrique, nous calculons la valeur absolue du t empirique et nous procédons au test par comparaison avec la valeur lue directement dans la table.

-
1. La notion de degrés de liberté est explicitée au chapitre 2.
 2. Les lois de probabilité sont en fin d'ouvrage.
 3. Si le nombre d'observations n est supérieur à 30, on peut approximer la loi de Student par une loi normale, soit $t^{\alpha/2} \approx 1,96$.

Exercice n° 1

↓ fichier C1EX1

Calcul d'un coefficient de corrélation

Un agronome s'intéresse à la liaison pouvant exister entre le rendement de maïs x (en quintal) d'une parcelle de terre et la quantité d'engrais y (en kilo). Il relève 10 couples de données consignés dans le tableau 2

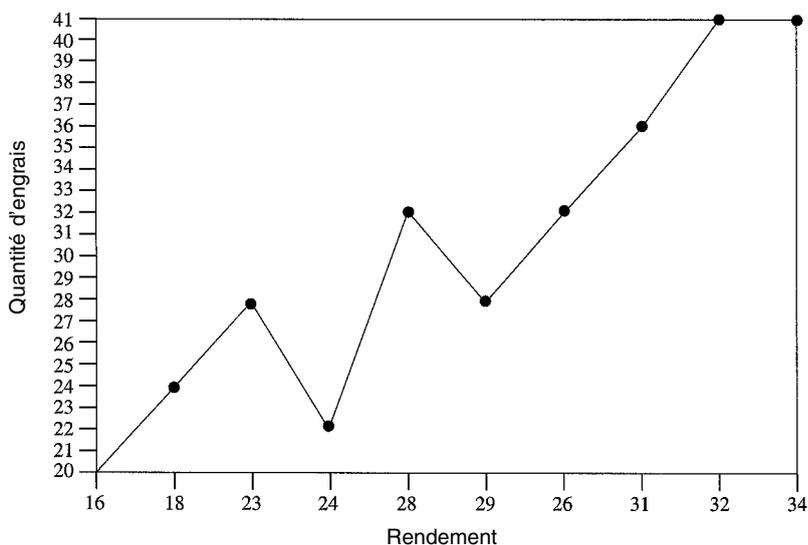
Tableau 2 – Rendement de maïs et quantité d'engrais

Rendement x	16	18	23	24	28	29	26	31	32	34
Engrais y	20	24	28	22	32	28	32	36	41	41

- 1) Tracer le nuage de points et le commenter.
- 2) Calculer le coefficient de corrélation simple et tester sa signification par rapport à 0 pour un seuil $\alpha = 0,05$.

Solution

- 1) Le nuage de points (graphique 6) indique que les couples de valeurs sont approximativement alignés : les deux variables semblent corrélées positivement.



Graphique 6 – Nuage du couple de valeurs :
rendement-quantité d'engrais

2) Afin d'appliquer la formule [2], nous dressons le tableau de calcul 3.

Tableau 3 – Calcul d'un coefficient de corrélation

	x	y	x^2	y^2	xy
	16	20	256	400	320
	18	24	324	576	432
	23	28	529	784	644
	24	22	576	484	528
	28	32	784	1 024	896
	29	28	841	784	812
	26	32	676	1 024	832
	31	36	961	1 296	1 116
	32	41	1 024	1 681	1 312
	34	41	1 156	1 681	1 394
Somme	261	304	7 127	9 734	8 286

$$\rho_{x,y} = \frac{(10)(8\,286) - (261)(304)}{\sqrt{(10)(7\,127) - 261^2} \sqrt{(10)(9\,734) - 304^2}} = \frac{3\,516}{(56,11)(70,17)}$$

soit $\rho_{x,y} = 0,89$ et $\rho_{x,y}^2 = 0,79$

Le t de Student empirique (d'après [3]) est égal à :

$$t^* = \frac{|\rho_{x,y}|}{\sqrt{\frac{1 - \rho_{x,y}^2}{n - 2}}} = \frac{0,89}{0,1620} = 5,49 > t_8^{0,025} = 2,306$$

le coefficient de corrélation entre x et y est significativement différent de 0.

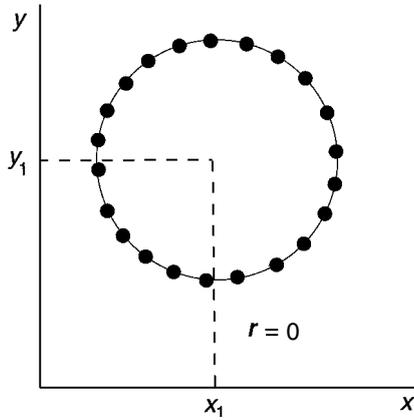
2) Limites de la notion de corrélation

a) La relation testée est linéaire

L'application de la formule [1] ou [2] ne permet de déterminer que des corrélations linéaires entre variables. Un coefficient de corrélation nul indique que la covariance entre la variable x et la variable y est égale à 0. C'est ainsi que deux variables en totale dépendance peuvent avoir un coefficient de corrélation nul, comme l'illustre l'exemple suivant : l'équation d'un cercle nous est donnée par $(x - x_1)^2 + (y - y_1)^2 = R^2$, les variables x et y sont bien liées entre elles fonctionnellement (graphique 7) et pourtant leur covariance est nulle et donc leur coefficient de corrélation égal à 0.

Pour pallier cette limite, il convient éventuellement de transformer les variables, préalablement au calcul du coefficient de corrélation, afin de linéariser

leur relation, par exemple au moyen d'une transformation de type logarithmique.



Graphique 7 – Relation fonctionnelle n'est pas corrélation linéaire

b) Corrélacion n'est pas causalité

Le fait d'avoir un coefficient de corrélation élevé entre deux variables ne signifie pas qu'il existe un autre lien que statistique. En d'autres termes, une covariance significativement différente de 0 n'implique pas une liaison d'ordre économique, physique ou autre. Nous appelons *corrélacion fortuite* ce type de corrélation que rien ne peut expliquer.

L'exemple le plus fameux concerne la forte corrélation existante entre le nombre de taches solaires observées et le taux de criminalité aux États-Unis. Cela ne signifie pas qu'il existe une relation entre les deux variables, mais qu'une troisième variable, l'évolution de long terme (la tendance) ici, explique conjointement les deux phénomènes. La théorie de la cointégration traite de ce problème (cf. chapitre 11).

2. Le modèle de régression simple

Nous commençons notre étude par le modèle le plus simple : une variable endogène est expliquée par une variable exogène. Après avoir étudié les conséquences probabilistes de l'erreur d'observation, nous présentons en I. les formules de base permettant d'estimer les paramètres du modèle. Les hypothèses stochastiques et leurs conséquences sont étudiées au paragraphe II. En III. et IV., la qualité de l'estimation d'un modèle est examinée à l'aide des premiers tests statistiques (Student, Fisher). Enfin, en V., le modèle de régression simple est étudié en tant qu'outil de prévision avec le degré de confiance que nous pouvons en attendre.

I. Présentation du modèle

A. Exemple introductif

Soit la fonction de consommation keynésienne :

$$C = a_0 + a_1 Y$$

où :

C = consommation,

Y = revenu,

a_1 = propension marginale à consommer,

a_0 = consommation autonome ou incompressible.

1) Vocabulaire

- La variable consommation est appelée « variable à expliquer » ou « variable endogène ».
- La variable revenu est appelée « variable explicative » ou « variable exogène » (c'est le revenu qui explique la consommation).
- a_1 et a_0 sont les paramètres du modèle ou encore les coefficients de régression.

2) Spécification

Nous pouvons distinguer deux types de spécifications :

- Les modèles en série temporelle, les variables représentent des phénomènes observés à intervalles de temps réguliers, par exemple la consommation et le revenu annuel sur 20 ans pour un pays donné. Le modèle s'écrit alors :

$$C_t = a_0 + a_1 Y_t \quad t = 1, \dots, 20$$

où :

C_t = consommation au temps t ,

Y_t = revenu au temps t .

- Les modèles en coupe instantanée, les variables représentent des phénomènes observés au même instant mais concernant plusieurs individus, par exemple la consommation et le revenu observés sur un échantillon de 20 pays. Le modèle s'écrit alors :

$$C_i = a_0 + a_1 Y_i \quad i = 1, \dots, 20$$

où :

C_i = consommation du pays i pour une année donnée,

Y_i = revenu du pays i pour une année donnée.

B. Rôle du terme aléatoire

Le modèle tel qu'il vient d'être spécifié n'est qu'une caricature de la réalité. En effet ne retenir que le revenu pour expliquer la consommation est à l'évidence même insuffisant ; il existe une multitude d'autres facteurs susceptibles d'expliquer la consommation. C'est pourquoi nous ajoutons un terme (ε_t) qui synthétise l'ensemble de ces informations non explicitées dans le modèle : $C_t = a_0 + a_1 Y_t + \varepsilon_t$ si le modèle est spécifié en série temporelle ($C_i = a_0 + a_1 Y_i + \varepsilon_i$ si le modèle est spécifié en coupe instantanée), où ε_t représente l'erreur de spécification du modèle, c'est-à-dire l'ensemble des phénomènes explicatifs de la consommation non liés au revenu. Le terme ε_t mesure la

différence entre les valeurs réellement observées de C_t et les valeurs qui auraient été observées si la relation spécifiée avait été rigoureusement exacte. Le terme ε_t regroupe donc trois erreurs :

- une erreur de spécification, c'est-à-dire le fait que la seule variable explicative n'est pas suffisante pour rendre compte de la totalité du phénomène expliqué ;
- une erreur de mesure, les données ne représentent pas exactement le phénomène ;
- une erreur de fluctuation d'échantillonnage, d'un échantillon à l'autre les observations, et donc les estimations, sont légèrement différentes.

Exercice n° 1

↓ fichier C2EX1

Génération d'une consommation aléatoire

Le tableau 1 présente le revenu moyen par habitant sur 10 ans exprimé en dollars pour un pays.

Tableau 1 – Évolution du revenu moyen par habitant en dollars

Année	Revenu
1	8 000
2	9 000
3	9 500
4	9 500
5	9 800
6	11 000
7	12 000
8	13 000
9	15 000
10	16 000

Sachant que la propension marginale à consommer est de 0,8 et que la consommation incompressible est 1 000, on demande :

- 1) de calculer la consommation théorique sur les 10 ans ;
- 2) considérant que notre erreur d'observation suit une loi normale de moyenne 0 et de variance 20 000, de générer cette variable aléatoire et de calculer une consommation observée tenant compte de cette erreur.

Solution

Les calculs des questions 1) et 2) sont présentés dans le tableau 2.

La consommation théorique (colonne 3) est calculée par application directe de la formule : $C_t = 1\,000 + 0,8 Y_t$.