

Chapitre 10

Régression multiple

Le modèle étudié dans ce chapitre— la régression multiple— fait appel à des notions théoriques qui dépassent le niveau de cet ouvrage; et son application exige des calculs qui peuvent difficilement se faire sans l'aide d'un logiciel. Donc la théorie ne sera développée que dans ses grandes lignes: les démonstrations de théorèmes, ainsi que certaines formules, seront omises. Et les formules de calcul ne seront pas toutes explicitées: pour celles qui ne le seront pas, nous nous en remettons à Excel.

Ces compromis n'affecteront pas la compréhension. Nous continuerons à énoncer les théorèmes formellement, dans le langage mathématique familier des chapitres précédents: il sera encore question d'estimateurs, de lois d'échantillonnage, d'intervalles de confiance, et de tests d'hypothèses.

Pour quelques calculs relatifs à la régression multiple, Excel ne suffit pas. Il faudra alors recourir à un logiciel spécialisé. Il y en a plusieurs sur le marché et il importe peu lequel est choisi: le format de base des sorties de logiciel est à peu de choses près le même pour tous. Pour les lecteurs qui s'initient à l'emploi d'un logiciel statistique, nous recommandons **MegaStat**, un logiciel associé à Excel. Quelques indications sur l'utilisation de **MegaStat** sont présentées en annexe de ce chapitre.

10.1 Le modèle

La régression linéaire simple développée au chapitre 8 permet de prédire la valeur d'une variable endogène Y à partir d'une variable exogène x . Il est clair, cependant que la valeur d'une variable est généralement affectée par plusieurs facteurs. Considérons les maisons vendues dans une ville au cours d'une période donnée, disons un an, et qu'on s'intéresse au prix de vente. On peut prédire le prix d'une maison donnée à partir d'une variable exogène comme, par exemple, le nombre x de chambres à coucher. Mais le prix d'une maison, même pour un x fixe, peut varier beaucoup à cause des nombreux autres facteurs qui contribuent au prix: le quartier, l'âge, le nombre de salles de bain, etc. Une prédiction basée sur ces informations promet d'être plus précise qu'une prédiction basée sur le seul nombre de chambres à coucher.

Un modèle pour ce faire est une extension toute naturelle du modèle de régression simple, le modèle de *Régression multiple*. Supposons qu'on dispose de k variables exogènes, x_1, x_2, \dots, x_k , et une variable endogène Y . L'échantillon est constitué de n observations sur chacune des variables. On désigne par Y_i la i^e observation sur Y et $x_{1i}; x_{2i}; \dots; x_{ki}$ les i^e valeurs des variables x_1, x_2, \dots, x_k , respectivement. On désigne ces valeurs par le vecteur

$$\mathbf{x}_i = [x_{1i}; x_{2i}; \dots; x_{ki}].$$

Le modèle stipule que, étant donné une valeur $\mathbf{x} = [x_1; x_2; \dots; x_k]$ des variables exogènes, l'espérance de Y est une fonction linéaire de \mathbf{x} :

$$\mu_{y,\mathbf{x}} = E(Y | \mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k,$$

alors que la variance est indépendante de \mathbf{x} :

$$\text{Var}(Y | \mathbf{x}) = \sigma^2.$$

Plusieurs des hypothèses du modèle de régression simple sont retenues:

- Les n variables Y_i sont supposées indépendantes;
- La variance de Y_i pour une valeur donnée de \mathbf{x} est indépendante de \mathbf{x} (l'hypothèse d'homoscédasticité);
- Les valeurs $x_{1i}, x_{2i}, \dots, x_{ki}$ sont des nombres fixes et non des variables aléatoires;
- Les intervalles de confiance et les tests d'hypothèses reposent sur la supposition que les Y_i sont de loi normale:

$$Y_i \sim \mathcal{N}(\mu_{y,x}; \sigma^2).$$

10.2 Estimation des paramètres, intervalles de confiance et tests d'hypothèses

Estimateurs des paramètres

Nous avons maintenant $k + 2$ paramètres à estimer:

les $k + 1$ composantes du vecteur $\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$, et σ^2 .

Plusieurs approches différentes mènent aux mêmes estimateurs $\hat{\beta}_j$ des β_j . Parmi elles se trouve le principe des moindres carrés (tout comme pour une régression simple):

Les estimateurs de $\beta_1, \beta_2, \dots, \beta_k$ sont les valeurs $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ qui minimisent la somme des carrés des erreurs

$$Q = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki})]^2.$$

Les « formules » des $\hat{\beta}_j$ s'expriment en un langage matriciel que nous préférons éviter ici. Nous les omettons donc et laissons les logiciels effectuer le calcul.

L'estimation de la variance σ^2 est basée sur les écarts quadratiques $(y_i - \hat{y}_i)^2$, où

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \dots + \hat{\beta}_k x_{ki}.$$

Un estimateur sans biais de σ^2 est donné par

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - (k + 1)}.$$

Remarque L'estimation de la variance est ici conforme au principe utilisé jusqu'ici: $\sigma^2 = E[(Y_i - \mu_i)^2]$, où

$$\mu_i = E(Y_i / \mathbf{x}_i) = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki}.$$

Il est donc naturel de l'estimer par une moyenne des carrés des écarts $y_i - \mu_i$. Cette moyenne n'étant pas connue, on la remplace par son estimation \hat{y}_i . Cependant, nous entendons par « moyenne » la somme $(y_i - \hat{y}_i)^2$ divisée non pas par le nombre de termes n mais plutôt par $n - (k + 1)$. C'est le nombre de degrés de liberté de la distribution de $\hat{\sigma}^2$, et c'est ce qui fait de $\hat{\sigma}^2$ un estimateur sans biais.

Il existe un estimateur sans biais $\hat{\sigma}_{\beta_j}^2$ de la variance $\sigma_{\beta_j}^2$ de $\hat{\beta}_j$; nous n'en donnons pas la formule, mais nous énonçons ses propriétés.

Distribution des estimateurs

Sous l'hypothèse de normalité,

$$\hat{\beta}_j \sim \mathcal{N}(\beta_j; \sigma_{\beta_j}^2)$$

Donc $\hat{\beta}_j$ est sans biais.

Il existe une formule pour sa variance $\sigma_{\hat{\beta}_j}^2$, formule qui sera elle aussi omise.

Comme on le remarque souvent, l'estimateur de la variance $\hat{\sigma}^2$ suit, à une constante multiplicative près, une loi khi-deux:

Sous l'hypothèse que les y_i sont de loi normale,

$$\frac{[n - (k + 1)]\hat{\sigma}^2}{\sigma^2} = \frac{\sum_{i=1}^n (Y_i - \hat{y}_i)^2}{\sigma^2} \sim \chi_{n-(k+1)}^2$$

Remarque

Ce résultat explique pourquoi le diviseur de $\hat{\sigma}^2$ doit être $n - (k + 1)$ pour que $\hat{\sigma}^2$ soit sans biais. L'espérance d'une variable de loi χ^2 est égale à son nombre de degrés de liberté.

$$E\left[\frac{[n - (k + 1)]\hat{\sigma}^2}{\sigma^2}\right] = n - k + 1 \Rightarrow \frac{n - (k + 1)}{\sigma^2} E(\hat{\sigma}^2) = n - k + 1 \Rightarrow E(\hat{\sigma}^2) = \sigma^2: \hat{\sigma}^2 \text{ est donc sans biais.}$$

Il existe aussi des estimateurs $\hat{\sigma}_{\hat{\beta}_j}^2$ sans biais des variances $\sigma_{\hat{\beta}_j}^2$.

Finalement, les statistiques à la base de la construction des intervalles de confiance et des tests d'hypothèses concernant les β_j suivent la loi de *Student*:

Sous l'hypothèse que les Y sont de loi normale,

$$t_j = \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_{\hat{\beta}_j}} \sim t_{n-(k+1)}$$

Intervalles de confiance

Il s'ensuit du dernier résultat qu'un intervalle de confiance de niveau $100(1-\alpha)\%$ est donné par

$$\hat{\beta}_j - t_{n-k-1;\alpha/2} \hat{\sigma}_{\hat{\beta}_j} \leq \beta_j \leq \hat{\beta}_j + t_{n-k-1;\alpha/2} \hat{\sigma}_{\hat{\beta}_j}$$

Test d'hypothèse concernant β_j

Si β_{j_0} est une valeur donnée, un test bilatéral de l'hypothèse

$$H_0 : \beta_j = \beta_{j_0} \text{ vs } H_1 : \beta_j \neq \beta_{j_0}$$

est donné par:

$$\text{On rejette } H_0 \text{ si } \left| \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_{\hat{\beta}_j}} \right| > t_{n-k-1;\alpha/2}$$

Les hypothèses que nous voudrions tester généralement sont

$$H_0 : \beta_j = 0 \text{ vs } H_1 : \beta_j \neq 0.$$

La statistique de test est $t_j = \frac{\hat{\beta}_j}{\hat{\sigma}_{\hat{\beta}_j}}$ et la région critique est

$$\text{On rejette } H_0 \text{ si } |t_j| > t_{n-k-1;\alpha/2}$$

Un critère analogue est le suivant:

On rejette H_0 si

$$P(|T_{n-(k+1)}| > t_j | H_0) < \alpha.$$

où $T_{n-(k+1)}$ désigne une variable aléatoire de loi de *Student* à $n - (k+1)$ degrés de liberté.

La probabilité

$$p = P(|T_{n-(k+1)}| > t_j | H_0)$$

est appelée *valeur p*.

Exemple 10.2.1 *Le prix d'une maison en fonction de son âge et du nombre de chambres à coucher*

Le tableau 10.1 présente des données sur la vente de $n = 43$ maisons. On se propose de développer une formule permettant de prédire le prix d'une maison à partir de son âge et du nombre de chambres à coucher qu'elle a. Les variables observées sont:

Prix: Le prix auquel la maison s'est vendue (en milliers de dollars)
 age: L'âge de la maison
 cc: Le nombre de chambres à coucher

Le modèle de régression linéaire est

$$\text{Prix} = \beta_0 + \beta_1(\text{age}) + \beta_2(\text{cc})$$

Nous effectuons les calculs à l'aide de **MegaStat**. Le logiciel fournit un grand nombre de données, mais les principaux résultats sont contenus dans le tableau suivant:

Regression output				confidence interval		
variables	coefficients	std. error	t (df=40)	p-value	95% lower	95% upper
Intercept	32.4148	10.8474	2.988	.0048	10.4914	54.3382
age	-0.4295	0.0870	-4.936	1.45E-05	-0.6053	-0.2536
cc	12.1198	3.4052	3.559	.0010	5.2376	19.0020

a) *Estimation des coefficients* β_0 , β_1 , et β_2 .

La réponse découle de la colonne *coefficients*

Sont indiquées dans cette colonne les estimations des coefficients β_0 , β_1 , et β_2 .

Ainsi donc: $\hat{\beta}_0 = 32,4147$; $\hat{\beta}_1 = -0,4295$; $\hat{\beta}_2 = 12,1198$. On estime le prix moyen des maisons d'âge x_1 et de x_2 chambres à coucher par $32,4147 - 0,4295x_1 + 12,1198x_2$. Pour une maison de 15 ans avec 2 chambres à coucher, cela donne $32,4147 - 0,4295(15) + 12,1198(2) = 50\ 212\ \$$

b) *Estimation des écarts-types* σ_{β_0} , σ_{β_1} et σ_{β_2}

La réponse découle de la colonne *std. error*

Cette colonne présente les estimations $\hat{\sigma}_{\beta_j}$ des écarts-types σ_{β_j} des $\hat{\beta}_j$.

Donc $\hat{\sigma}_{\beta_0} = 10,8474$; $\hat{\sigma}_{\beta_1} = 0,0870$; $\hat{\sigma}_{\beta_2} = 3,4052$.

Ces écarts-types serviront au calcul des statistiques de *Student* pour tester des hypothèses et déterminer des intervalles de confiance.

c) *Tests des hypothèses* $\beta_j = 0$.

La réponse découle de la colonne *t (df=40)*

Cette colonne présente les statistiques $t_j = \frac{\hat{\beta}_j}{\hat{\sigma}_{\beta_j}}$. Sous l'hypothèse que $\beta_j = 0$, t_j suit une loi de

Student à $n-k-1 = 40$ degrés de liberté. Avec $\alpha = 0,05$, on rejette l'hypothèse si $|t_j| > t_{40;0,025} = 2,02$. On a $|t_1| = 0,0048 > 2,02$; $|t_2| = 4,936 > 2,02$; et $|t_3| = 3,559 > 2,02$. On rejette les 3 hypothèses, ce qui permet de conclure que $\beta_0 > 0$, $\beta_1 < 0$, et $\beta_2 > 0$.

La colonne *p-value* permet d'éviter ces calculs, puisqu'elle donne les *p-valeurs* $P(|T_{40}| > |t_j|)$ sous chacune des hypothèses $\beta_j = 0$. On rejette l'hypothèse si $p < \alpha$. On constate qu'avec $\alpha = 0,05$ on peut rejeter les trois hypothèses. Bien plus encore: on sait qu'on aurait pu les rejeter avec un α beaucoup plus petit.

d) En termes concrets, comment s'expriment les conclusions en c)?

L'hypothèse que $\beta_0 = 0$ est testée automatiquement par le logiciel, mais elle est peu pertinente. En général, on ne s'attend pas à ce qu'elle soit vraie.

Les hypothèses $\beta_1 = 0$ et $\beta_2 = 0$, par contre, sont fondamentales: ce sont des tests sur la pertinence des variables exogènes.

L'hypothèse que $\beta_1 = 0$ affirme que l'âge de la maison est sans effet sur le prix. Le test conclut le contraire: le prix de la maison décroît avec l'âge. L'inclusion de cette variable exogène améliore la prédiction du prix.

Même chose pour l'hypothèse que $\beta_2 = 0$. Conclure que $\beta_2 > 0$, c'est conclure que le prix de la maison croît avec le nombre de chambres à coucher (ce qui coule de source).

e) *Intervalle de confiance à 95 % pour les β_j .*

Les intervalles de confiance sont présentés dans les colonnes *confidence interval*.

L'intervalle pour β_0 est sans intérêt.

L'intervalle pour β_1 : La valeur $\hat{\beta}_1 = -0,4295$ signifie que le prix de la maison décroît avec l'âge à un taux estimé de 429,5 \$ par année. L'intervalle de confiance entoure cette estimation d'une marge d'erreur de sorte qu'on peut affirmer avec 95 % de confiance que le taux de décroissance se situe entre 253 \$ et 605 \$ (en arrondissant au dollar le plus proche).

L'intervalle pour β_2 : La valeur estimée est $\hat{\beta}_2 = 12,1198$. On estime donc qu'une chambre à coucher de plus ajoute (en moyenne) 12 120 \$ à la valeur d'une maison. Pour être sûr (à 95 %), on dira que cette « valeur ajoutée » est située quelque part entre 5 238 \$ et 19 000 \$.

Intervalles de confiance pour $\mu_{y,x}$ et limites de prédiction

Pour une valeur donnée $\mathbf{x}_0 = [x_{10}, x_{20}, \dots, x_{k0}]$ des variables exogènes, l'espérance

$$\mu_{y,x_0} = E(Y | \mathbf{x}_0) = \beta_0 + \beta_1 x_{10} + \dots + \beta_k x_{k0}$$

est naturellement estimée par

$$\hat{\mu}_{y,x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_{10} + \hat{\beta}_2 x_{20} + \dots + \hat{\beta}_k x_{k0}.$$

Sous la supposition de normalité, nous pouvons déterminer un intervalle de confiance pour μ_{y,x_0} . Quand on détermine un intervalle de confiance pour μ_{y,x_0} , on affirme avec confiance que la *moyenne* des Y se situe dans les limites de l'intervalle.

La régression sert aussi à faire des *prédictions*. Une valeur future de Y non observée, peut être prédite à partir des valeurs des variables exogènes $\mathbf{x}_0 = [x_{10}, x_{20}, \dots, x_{k0}]$. La prédiction \hat{y}_0 est identique à $\hat{\mu}_{y,x_0}$:

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{10} + \hat{\beta}_2 x_{20} + \dots + \hat{\beta}_k x_{k0}$$

Cette prédiction est bien sûr sujette à erreur. Les *limites de prédiction* sont les limites dans lesquelles devraient se trouver la prochaine *observation* sur Y . Ces limites de prédiction sont plus larges que celles de l'intervalle de confiance pour μ_{y,x_0} .

Exemple 10.2.2 (Suite de l'exemple 10.2.1)

[Nous devons utiliser ici le logiciel **MegaStat**.]

Supposons qu'on considère une classe de maisons, celles de 2 chambres à coucher âgées de 20 ans. Voici un intervalle de confiance et des limites de prédiction à 95 % données par MegaStats pour les maisons de cette classe:

Predicted values for: Prix						
			95% Confidence Interval		95% Prediction Interval	
cc	age	Predicted	lower	upper	lower	upper
2	20	48.0651	39.0761	57.0541	15.8470	80.2832

La valeur Predicted signifie deux choses: c'est l'estimation d'une moyenne et c'est aussi la prédiction du prix d'une maison donnée appartenant à la classe considérée.

Estimation d'une moyenne On estime que le prix moyen $\mu_{y,x}$ des maisons de la classe considérée est de 48 065 \$. L'intervalle de confiance pour $\mu_{y,x}$ (*95% Confidence Interval*) permet d'affirmer avec 95 % de confiance que le prix moyen des maisons de 2 chambres à coucher de 20 ans d'âge se situe entre 39 076 \$ et 57 054 \$.

Limites de prédiction Considérons une maison de 2 chambres à coucher âgée de 20 ans à vendre. On prédit qu'elle se vendra à 48 065 \$. Les limites de prédiction (*95% Prediction Interval*) permettent d'affirmer que le prix auquel cette maison se vendra se situe quelque part entre 15 847 \$ et 80 283 \$.

Remarque

Les limites de prédiction dans le dernier exemple sont à toutes fins pratiques inutiles, voire ridicules. Cela est dû à plusieurs facteurs:

1. Bien que la corrélation ($R = 0,68$) ne soit pas négligeable, elle n'est pas assez forte pour fins de prévision. Les prix varient considérablement, même lorsqu'on se limite aux maisons de même âge et ayant le même nombre de chambres à coucher. Avec cette restriction, l'écart-type des prix ($\hat{\sigma}$) est de 15 308 \$—inférieur, forcément, à l'écart-type (20 295 \$) des prix dans la population entière, mais assez grand quand même.
2. L'échantillon n'est pas très grand, ce qui fait que l'estimation des paramètres β_j et donc de $\mu_{y,x}$ est peu fiable, comme le montre l'intervalle $39\,076 \$ \leq \mu_{y,x} \leq 57\,054 \$$.
3. Si la *moyenne* se situe entre 39 076 \$ et 57 054 \$, alors le prix d'une maison donnée est d'autant plus imprévisible que les prix se dispersent par rapport à leur moyenne.
4. Finalement, il faut reconnaître que le niveau de confiance demandé (95 %) est un peu trop ambitieux. Faute de mieux, on pourrait se contenter d'un niveau de confiance inférieur. Que peut-on affirmer, par exemple, à un niveau de confiance de 50 %? Dans ce cas, les limites de prédiction sont 37 215 \$ et 58 917 \$.
5. Quiconque a déjà eu l'occasion de s'intéresser au prix d'une maison serait surpris du peu de succès de modèle, comme on le voit au dernier paragraphe. Il est certain qu'une agente d'immeuble expérimentée pourrait faire une prédiction plus précise avec plus de confiance. Comment ça se fait? La raison est que l'agente tiendra compte de beaucoup plus de facteurs (donc de variables exogènes) que les deux que nous avons considérés ici. Les maisons peuvent être très différentes les unes des autres, même si elles sont du même âge et ont le même nombre de chambres à coucher: elles peuvent se situer dans des rues plus ou moins prisées; elles peuvent être grandes ou petites; être bien entretenues ou pas... autant de variables exogènes qui, si elles sont introduites dans la régression, feraient réduire la variance conditionnelle (estimée à 15 308 \$) et améliorer la précision des prédictions.
6. C'est un fait incontournable que la prédiction d'une *valeur donnée* (par opposition à l'estimation d'une *moyenne*) est généralement peu fiable, et que la régression ne produit pas des prédictions très précises quand on a affaire à une population naturellement très dispersée. Mais sans elle, les prédictions seraient encore moins précises.

10.3 Évaluation globale

Les tests présentés ci-dessus portent chacun sur un seul des coefficients β_j . En général, on voudra également tester la pertinence des variables exogènes prises dans leur *ensemble* ainsi que d'une mesure *globale* de la qualité de la régression (analogue au coefficient de corrélation dans le cas d'une régression simple).

L'outil principal est la décomposition présentée à la section 8.6:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{SCT} = \text{SCE} + \text{SCR}$$

où SCT est la somme des carrés *totale*; SCE est la somme des carrés *expliquée*, et SCR est la somme des carrés *résiduelle*. SCT et SCR peuvent être interprétées comme des mesures d'erreur de prédiction.

- SCT mesure les erreurs commises si on préditait Y sans l'apport des variables exogènes. Sans les variables exogènes, la valeur de Y est normalement estimée par la moyenne \bar{y} . L'erreur commise lorsqu'on prédit y_j par \bar{y} est $(y_j - \bar{y})$ et SCT est la somme des carrés de ces erreurs.
- SCR mesure les erreurs commises lorsque y_j est prédit par \hat{y}_j , la prédiction basée sur le modèle de régression multiple.

On s'attend, bien sûr, à ce que $\text{SCT} > \text{SCR}$. La différence $\text{SCE} = \text{SCT} - \text{SCR}$ mesure la réduction de l'erreur due à l'utilisation des variables exogènes. Si SCE est grande, c'est que la régression a permis de réaliser une bonne réduction de l'erreur.

Coefficient de corrélation multiple et coefficient de corrélation multiple ajusté

La réduction d'erreur SCE est difficile à interpréter, étant donné que sa valeur peut être grande ou petite, dépendamment de l'unité de mesure de Y . Cette difficulté est réglée en exprimant cette différence en termes relatifs, c'est-à-dire, comme fraction de SCT. Ce quotient est appelé *coefficient de détermination* ou *coefficient de corrélation multiple* et désigné par R^2 :

$$R^2 = \text{Coefficient de détermination} = \frac{\text{SCE}}{\text{SCT}} = 1 - \frac{\text{SCR}}{\text{SCT}}$$

Il est clair que $0 \leq R^2 \leq 1$: $R^2 = 0$ si $\text{SCR} = \text{SCT}$ (aucune réduction des erreurs); et $R^2 = 1$ si $\text{SCR} = 0$ (aucune erreur de prédiction dans le modèle de régression).

Le coefficient de corrélation multiple est simplement la racine carrée positive de R^2 :

$$R = \text{Coefficient de corrélation multiple} = \sqrt{R^2} = \sqrt{\frac{\text{SCE}}{\text{SCT}}} = \sqrt{1 - \frac{\text{SCR}}{\text{SCT}}}$$

Autre interprétation de R Une autre interprétation de R le réduit à la notion de coefficient de corrélation simple déjà connue. Pour mesurer la dépendance entre la variable endogène Y et *plusieurs* variables exogènes, x_1, x_2, \dots, x_k , on remplace les k variables exogènes par une fonction affine de celles-ci, $\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$. Comment choisir les β_j ? Eh bien, on prend $\beta_j = \hat{\beta}_j$, ce qui fait que les variables exogènes sont remplacées par une seule, $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_k x_k$, la valeur prédite de y . Alors

R est le coefficient de corrélation entre y et \hat{y} .

Ceci donne un deuxième sens intuitif à R .

Moyennes des carrés

Une somme de carrés augmente avec le nombre de termes et donc ne peut servir que si elle est ajustée pour tenir compte de ce fait, c'est-à-dire, à moins qu'on remplace la somme par une *moyenne* de carrés. Mais ce ne sera pas tout-à-fait une moyenne car on divisera non pas par le nombre de termes mais plutôt par la *nombre de degrés de liberté* (voir plus bas).

Voici le nombre de degrés de liberté des trois sommes de carrés:

Somme des carrés	Degrés de liberté	Moyennes des carrés
SCE	k	$MCE = \frac{SCE}{k}$
SCR	$n-(k+1)$	$MCR = \frac{SCR}{n-(k+1)}$
SCT	$n-1$	$MCT = \frac{SCT}{n-1}$

Remarque Noter que MCR est l'estimateur $\hat{\sigma}^2$ de σ^2 alors que SCT est la variance échantillonnale S_y^2 qui estime la variance des y_i sous un autre modèle, un modèle dans lequel $E(y_i) = \mu$ pour tout i .

On peut alors redéfinir R en remplaçant les sommes de carrés par les moyennes des carrés:

$$R \text{ ajusté} = \sqrt{R^2 \text{ ajusté}} = \sqrt{1 - \frac{MCR}{MCT}}$$

Remarque Il est rare que cet ajustement change la valeur de R de façon radicale. On la préfère comme mesure descriptive parce qu'elle fait payer un prix pour l'ajout de variables exogènes. L'introduction de nouvelles variables exogènes—pertinentes ou pas—ne peut que faire croître R . Le R ajusté réduit la possibilité qu'un accroissement qui ne serait dû qu'à l'ajout d'une nouvelle variable exogène non significative.

Exemple 10.3.2 (Suite de l'exemple 10.2.1)

MegaStat fournit R^2 , R^2 ajusté ainsi que $\hat{\sigma}$ (Std. Error).

Regression Analysis			
	R^2	0.458	
	Adjusted R^2	0.431	n 43
	R	0.677	k 2
	Std. Error	15.308	Dep. Var. Prix

Noter que $\hat{\sigma}$ n'estime pas l'écart-type de toutes les maisons; il estime l'écart-type des prix des maisons ayant toutes les mêmes valeurs des variables exogènes. L'écart-type des prix de toutes les maisons de la population peut être estimé (si les maisons ont été tirées au hasard) par $S_y = 20\,295$ \$. C'est le fait que $\hat{\sigma} = 15\,308$ \$ est plus petit que S_y qui rend la régression utile.

Test global

En général, avant de s'attarder sur la significativité des coefficients β_j , il est bon d'effectuer un test global, soit un test de l'hypothèse

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0.$$

L'hypothèse affirme simplement qu'aucune des variables exogènes ne contribue à la prédiction de

Y. Nous décrivons la procédure d'abord. Nous énoncerons ensuite les théorèmes qui l'appuient et nous justifierons la région critique en termes intuitifs.

Le test est basé sur la statistique $F = \frac{\text{MCE}}{\text{MCR}}$ et la région critique est

$$F > F_{k;n-(k+1); \alpha}$$

où $F_{k;n-(k+1); \alpha}$ est le point critique correspondant à une loi de Fisher à k et $n-(k+1)$ degrés de liberté.

Exemple 10.3.1 Une analyse globale — décomposition d'une somme de carrés

La commande Excel utilisée à l'exemple 10.2.1 fournit aussi le tableau suivant, une *table d'analyse de variance*:

ANOVA table					
Source	SS	df	MS	F	p-value
Regression	5 182.7293	2	2 591.3647	8.56	.0008
Residual	12 417.9950	41	302.8779		
Total	17 600.7243	43			

Les colonnes SS, df et MS donnent, respectivement, les sommes de carrés SCE, SCR et SCT, le nombre de degrés de liberté et les moyennes des carrés.

La colonne F calcule le quotient $F = \text{MCE}/\text{MCR}$.

Si $\alpha = 0,05$, l'hypothèse H_0 est rejetée si $F > F_{2;41;0,05} = 3,22$. La valeur observée $F = 8,56 > 3,22$ entraîne donc le rejet de H_0 .

La colonne *p-value* nous dispense de ces calculs. La valeur 0,0008 est la probabilité $P(F_{2;41} > 8,56)$, où $F_{2;41}$ est une variable de loi F à 2 et 42 degrés de liberté.

La conclusion est que les coefficients β_1 et β_2 ne sont pas tous nuls et donc qu'au moins l'un des deux est non nul. On peut dès lors étudier chacun des coefficients séparément.

Remarque

Il n'est pas impossible que les tests mènent au rejet de l'une des hypothèses $\beta_j = 0$ sans pour autant rejeter l'hypothèse H_0 que tous les coefficients sont nuls. On ne peut pas facilement concilier ces conclusions contradictoires. En général, l'approche conservatrice est recommandée: si H_0 est acceptée, on s'abstient de procéder à des tests individuels. On accepte le verdict du test global : l'utilité des variables exogènes n'est pas démontrée.

On ne peut pas, cependant, se montrer toujours aussi rigide. Il faut se rappeler lorsqu'on « accepte » H_0 , on n'affirme pas qu'elle est vraie; on dit simplement qu'on ne peut pas la rejeter avec confiance. Donc si H_0 est « presque » rejetée (une valeur p proche du seuil), il est défendable de poursuivre l'analyse avec des tests individuels sur les β_j .

Justification théorique

La région critique $F > F_{k;n-(k+1); \alpha}$ est basée sur les propriétés suivantes:

1. $\text{SCR} \sim \chi^2_{n-(k+1)}$
2. Sous H_0 , $\text{SCE} \sim \chi^2_k$
3. SCR et SCE sont indépendantes

Par conséquent, sous H_0 , $\frac{\text{SCE}/k}{\text{SCR}/[n-(k+1)]} = \frac{\text{MCE}}{\text{MCR}} \sim F_{k;n-(k+1)}$.

10.4 Variables dichotomiques

En régression linéaire, la variable endogène est nécessairement quantitative. On pourrait croire que les variables exogènes sont elles aussi nécessairement quantitatives, ce qui exclurait, par exemple, qu'on tienne compte de la variable SEXE pour prédire la valeur Y d'une mesure d'endurance physique. Or ce n'est pas le cas, car une variable qualitative peut toujours s'exprimer par une série de variables dichotomiques, c'est-à-dire, des variables qui ne prennent que deux valeurs—comme la variable SEXE, justement. Ses valeurs (Femme/Homme) peuvent toujours être remplacées par les nombres 0 et 1.

Considérons, pour commencer, le cas d'une variable exogène dichotomique. La procédure est simple: on remplace ses deux valeurs par 0 et 1, respectivement, puis on procède comme avec une variable quantitative ordinaire. Ce qui reste à préciser, c'est l'interprétation des résultats.

Exemple 10.4.1 Une variable dichotomique

Considérons les données du tableau A.3, qui portent sur un échantillon de 28 personnes, dont 16 femmes et 12 hommes. On s'intéresse aux variables suivantes:

Poids	Le poids de la personne (en livres)
Taille	La taille de la personne (en pouces)
Sexe	0 = Féminin; 1 = Masculin

Nous cherchons à développer une équation permettant de prédire le poids d'une personne à partir de son sexe et sa taille. Le modèle est:

$$E(\text{Poids}) = \beta_0 + \beta_1(\text{Sexe}) + \beta_2(\text{Taille})$$

Le logiciel MegaStats donne les estimations suivantes:

$$\hat{\beta}_0 = -124; \quad \hat{\beta}_1 = 6,7356; \quad \hat{\beta}_2 = 3,9477.$$

L'équation de prédiction serait donc

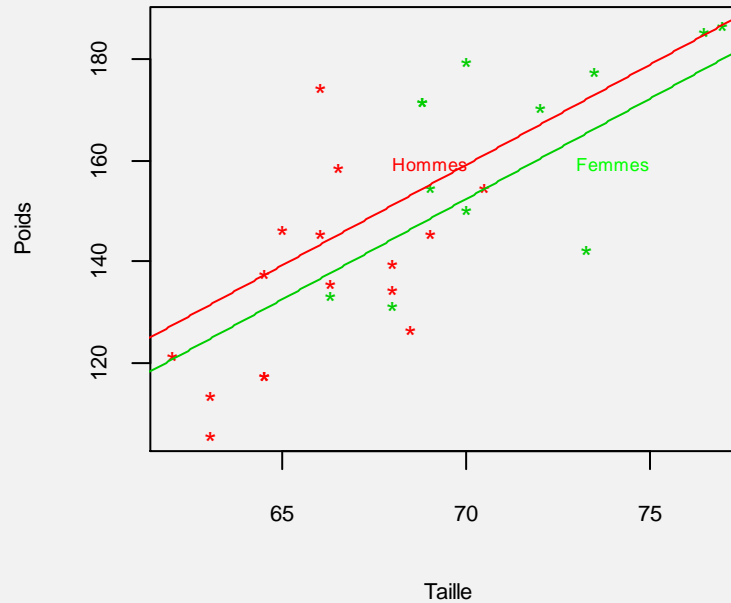
$$\text{Poids} = -124 + 6,7357(\text{Sexe}) + 3,9477(\text{Taille}).$$

Le coefficient de Sexe est 6,7357. Que représente ce nombre? Ce nombre est la différence de poids, pour une taille donnée, entre les hommes et les femmes. Pour le voir, on énonce les résultats d'une autre façon: On présente deux équations, l'une pour les femmes et l'autre pour les hommes:

Femmes:	$\text{Poids} = -124 + 3,9477(\text{Taille})$
Hommes:	$\text{Poids} = -124 + 6,7357 + 3,9477(\text{Taille})$ $= -117 + 3,9477(\text{Taille})$

Graphiquement, le lien entre la taille et le poids s'exprime par deux droites, l'une pour les femmes, l'autre pour les hommes:

Figure 10.1
Relation entre la taille et le poids
Femmes et hommes



Tests d'hypothèses

Nous n'avons pas encore abordé la question de savoir si les coefficients sont significativement différents de 0. Le tableau suivant, produit par **MegaStat**, fournit les tests pertinents.

Regression output					confidence interval	
variables	coefficients	std. error	t (df=25)	p-value	95% lower	95% upper
Intercept	-124.0006	70.5925	-1.757	.0912	-269.3886	21.3875
Sexe	6.7357	8.0845	0.833	.4126	-9.9147	23.3860
Taille	3.9477	1.0687	3.694	.0011	1.7467	6.1487

On constate que la *valeur p* pour l'hypothèse $\beta_1 = 0$ est 0,4126. Donc l'hypothèse ne peut être rejetée à un niveau raisonnable—on ne peut pas affirmer que le sexe d'une personne, une fois sa taille connue, est utile à la prédiction du poids de la personne. Plus concrètement, le poids moyen d'une femme est le même que celui d'un homme de même taille.

On détermine alors une régression avec seule la taille comme variable exogène. Le modèle est:

$$E(\text{Poids}) = \gamma_0 + \gamma_2(\text{Taille})$$

Les sorties suivantes donnent les estimations des paramètres:

Regression output					confidence interval	
variables	coefficients	std. error	t (df=26)	p-value	95% lower	95% upper
Intercept	-162.3787	53.1772	-3.054	.0052	-271.6860	-53.0713
Taille	4.5531	0.7790	5.845	3.68E-06	2.9518	6.1544

Il est intéressant de noter que si seule la variable SEXE est utilisée comme variable exogène, elle se révèle hautement significative:

Regression output					confidence interval	
variables	coefficients	std. error	t (df=26)	p-value	95% lower	95% upper
Intercept	136.3750	4.7314	28.823	2.91E-21	126.6495	146.1005
Sexe	27.0417	7.2273	3.742	.0009	12.1856	41.8977

Question: En fin de compte, la variable SEXE est-elle, oui ou non, utile à la prédiction du poids?
 Les analyses ci-dessus donnent à croire que si les femmes ont un poids inférieur à celui des hommes, c'est dû uniquement au fait qu'elles sont en moyenne plus petites de taille. (Ceci contredit ce qu'on sait des différences anatomiques entre femmes et hommes, mais rappelons que « accepter » une hypothèse n'est pas équivalent à l'affirmer.)

Dans le dernier exemple, nous avons considéré un modèle qui exprime la relation entre la taille et le poids par deux droites, l'une pour les femmes, l'autre pour les hommes. Mais la construction du modèle a contraint les deux droites à être parallèles. L'hypothèse implicite est qu'un accroissement d'une unité de taille entraîne un accroissement de poids qui est le même chez les femmes et les hommes. Il serait préférable de mettre cette hypothèse à l'épreuve car ce n'est peut-être pas le cas. La chose à faire est donc de commencer par un modèle plus général dans lequel les droites ne sont pas nécessairement parallèles et ensuite tester l'hypothèse de parallélisme. La procédure est illustrée dans le prochain exemple.

Exemple 10.4.2 [Suite de l'exemple 10.4.1]

Dans le modèle traité à l'exemple 10.4.1, l'équation $E(\text{Poids}) = \beta_0 + \beta_1(\text{Sexe}) + \beta_2(\text{Taille})$ implique une même pente β_2 indépendamment du sexe (alors que le terme $\beta_1(\text{Sexe})$ qui s'ajoute seulement aux sujets masculins permet que les droites ne soient pas confondues). S'il faut distinguer hommes et femmes quant à la pente, il faudrait, selon le même principe, ajouter un terme, disons $\beta_3(\text{Taille})$ seulement pour les sujets masculins. Pour ce faire, on utilise un artifice suivant: on crée une nouvelle variable *SexeTaille* qui représente la taille pour les hommes mais qui prend la valeur 0 pour les femmes.

$$\text{SexeTaille} = \begin{cases} \text{Taille} & \text{si le sujet est un homme} \\ 0 & \text{sinon} \end{cases}$$

Le modèle est donc

$$E(\text{Poids}) = \beta_0 + \beta_1(\text{Sexe}) + \beta_2(\text{Taille}) + \beta_3(\text{TailleSexe})$$

Regression output				
variables	coefficients	std. error	t (df=24)	p-value
Intercept	-143.2724	112.9266	-1.269	.2167
Sexe	40.1774	151.2235	0.266	.7928
Taille	4.2399	1.7111	2.478	.0206
SexeTaille	-0.4915	2.2192	-0.221	.8266

On résume

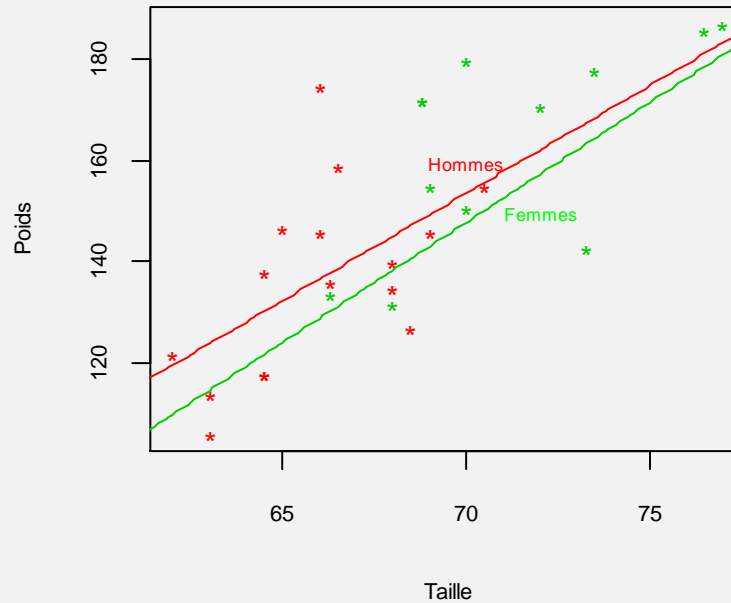
$$\text{Femmes: } \text{Poids} = -143,2724 + 4,2399(\text{Taille})$$

$$\begin{aligned} \text{Hommes: } \text{Poids} &= (-143,2724 + 40,1774) + (4,2399 - 0,4915)(\text{Taille}) \\ &= -103,0950 + 3,7484(\text{Taille}) \end{aligned}$$

La première hypothèse à tester est $\beta_3 = 0$: c'est l'hypothèse que les deux droites sont parallèles. On est loin de pouvoir la rejeter, et on l'exclut du modèle pour revenir au modèle considéré dans le dernier exemple.

La figure 10.2 illustre la relation.

Figure 10.2
Relation entre la taille et le poids
Femmes et hommes



Remarque

Nous n'avons pas rejeté l'hypothèse $\beta_3 = 0$ et nous l'avons donc exclu du modèle, c'est-à-dire que le terme en β_3 n'entrera pas dans notre description de la relation entre la taille et le poids, et ne servira pas à la prédiction poids. Mais cette décision n'est pas automatique. Si on exclut le terme ici, c'est d'abord parce que la valeur p est très élevée (0,8266), mais aussi parce le signe (négatif) de β_3 est a priori peu crédible.

Dans d'autres circonstances, on pourrait très bien conserver un terme qui n'est pas significativement différent de 0—c'est ce qu'on aurait fait ici, par exemple, si $\hat{\beta}_3$ avait été positif et si la valeur p avait été proche du seuil.

10.5 Régression polynomiale

Évidemment, la relation entre deux variables n'est souvent pas linéaire. La figure 10.3, qui illustre la relation entre la grosseur d'un diamant et son prix en est un exemple. Le nuage de points présente une courbure qui montre que le modèle $E(Y) = \beta_0 + \beta_1 x$ qui impose une structure linéaire s'ajuste mal aux données. Il faut alors songer à une fonction non linéaire. Parmi celles-là se trouvent les fonctions polynomiales. Par exemple, $E(Y) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$. Ce modèle est un modèle de régression multiple, avec x , x^2 et x^3 comme variables exogènes.

Exemple 10.5.1 Régression polynomiale

Le tableau 10.3 présente des données sur le prix et la grosseur d'un échantillon de diamants. La figure 10.2 montre bien que la relation n'est pas linéaire. Considérons donc le modèle $E(\text{Prix}) = \beta_0 + \beta_1(\text{Carat}) + \beta_2(\text{Carat}^2)$, où $\text{Carat}^2 = \text{Carat}^2$.

Voici les résultats fournis par MegaStats:

Regression output

variables	coefficients	std. error	t (df=97)	p-value	confidence interval	
					95% lower	95% upper
Intercept	259.2403	466.9891	0.555	.5801	-667.6038	1 186.0843
Carats	1 676.8083	1 694.8114	0.989	.3249	-1 686.9232	5 040.5399
Carats2	7 647.7129	1 346.6735	5.679	1.41E-07	4 974.9388	10 320.4870

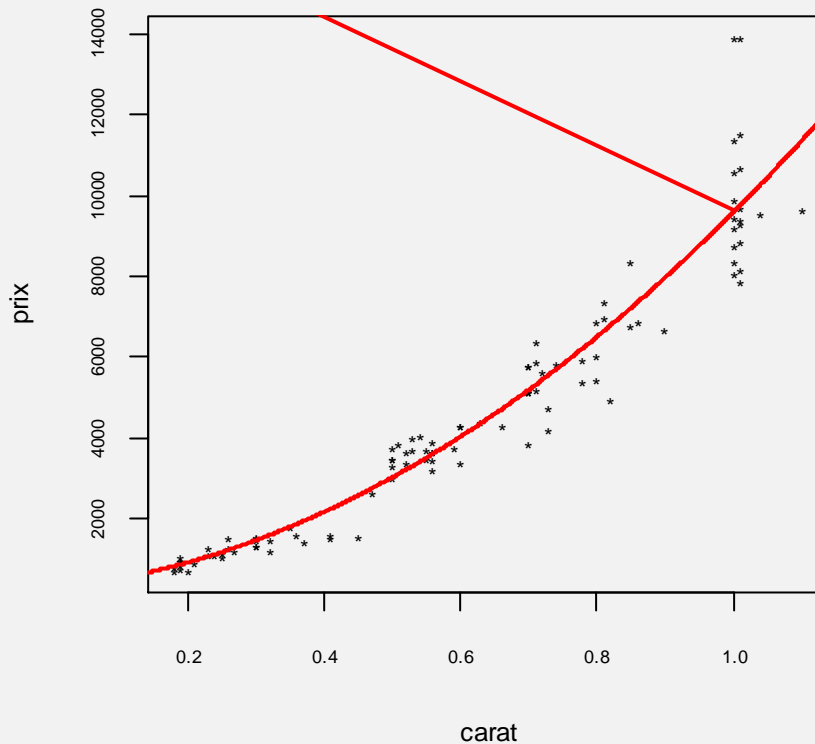
On estime donc que $E(\text{Prix}) = 259 + 1\,677(\text{Carats}) + 7\,648(\text{Carats}^2)$. On constate que le terme Carats2 est fortement significatif, ce qui confirme qu'une droite n'est pas suffisante pour exprimer la relation entre le prix et la grosseur (un fait bien connu dans le domaine, d'ailleurs).

Le terme linéaire $\hat{\beta}_1$ (le coefficient de Carats) est tout à fait non significatif. Peut-on l'éliminer du modèle? Normalement, le terme linéaire est conservé, car même si $\hat{\beta}_1$ est non significatif, il n'est pas prouvé que $\beta_1 = 0$. Dans le cas présent on pourrait l'omettre il ne contribue presque rien à R^2 : il le fait passer de 0,921 à 0,922. On peut alors proposer comme modèle définitif un modèle avec Carats2 comme seul variable exogène.

variables	coefficients	std. error	t (df=98)	p-value	confidence interval	
					95% lower	95% upper
Intercept	696.4480	150.9855	4.613	1.20E-05	396.8222	996.0738
Carats2	8 954.0316	264.9640	33.793	8.23E-56	8 428.2192	9 479.8441

L'équation définitive serait alors $E(\text{Prix}) = 696 + 8954(\text{Carats}^2)$

Figure 10.3
Prix de 100 diamants
en fonction de leur grosseur



10.6 Analyse de variance par la régression

L'analyse de variance à un facteur, traitée au chapitre 9 peut également être considéré comme un cas particulier d'une régression multiple, car il est possible d'exprimer les valeurs d'une variable qualitative à l'aide d'un certain nombre de variables dichotomiques. L'exemple suivant illustre la façon de faire

Exemple 10.6.1 Une analyse de variance effectuée vis une régression multiple

Les données suivantes sont tirées de l'exercice 9.10. Ce sont les prix (en milliers) ses maisons d'une ou deux chambres à coucher vendues dans trois secteurs d'une ville:

Secteur Sud: 249 184 239 314 69 69 80 85 85 269 60 89 89 89 89 142 142
 Secteur Nord: 97 157 170 184 385 145
 Secteur Centre: 98 249 289 299

L'objectif de l'analyse est de savoir s'il y a une différence entre les trois secteurs quant au prix moyen des maisons. Pour effectuer une régression avec *Prix* pour variable endogène et *Secteur* pour variable exogène, nous devons exprimer cette dernière par des variables quantitatives. On commence par prendre l'un des secteurs comme base de comparaison. Ce pourrait être n'importe lequel. Choisissons le secteur sud. Nous définirons alors deux variables exogènes:

$$Nord = \begin{cases} 1 & \text{pour une maison dans le secteur nord,} \\ 0 & \text{pour une maison ailleurs} \end{cases}, \quad Centre = \begin{cases} 1 & \text{pour une maison dans le secteur centre} \\ 0 & \text{pour une maison ailleurs} \end{cases}$$

Les valeurs des variables *Prix*, *Nord* et *Centre* sont présentées au tableau 10.5.1.

Le modèle est

$$E(Prix) = \beta_0 + \beta_1(Nord) + \beta_2(Centre)$$

Le secteur étant le secteur sud, les paramètres β_1 et β_2 représentent l'écart moyen des prix des maisons des secteurs nord et centre, respectivement, par rapport à ceux du secteur sud.

Le logiciel MegsStats produit les estimations des paramètres:

Regression output					confidence interval	
variables	coefficients	std. error	t (df=24)	p-value	95% lower	95% upper
Intercept	137.8235	21.4144	6.436	1.18E-06	93.6263	182.0207
Nord	60.7765	44.9193	1.353	.1887	-31.9324	153.4853
Centre	78.1765	44.9193	1.740	.0946	-14.5324	170.8853

On a $\hat{\beta}_1 = 60,7765$ et $\hat{\beta}_2 = 78,1765$. On estime donc que les maisons des secteurs nord et centre valent en moyenne 60 777 \$ et 78 176 \$ de plus que celles du secteur sud.

Les valeurs $p = 0,1887$ et $p = 0,0946$ testent les hypothèses $\beta_1 = 0$ et $\beta_2 = 0$, respectivement. On ne peut pas raisonnablement rejeter la première hypothèse; la deuxième pourrait l'être, si on adopte un niveau $\alpha = 0,1$.

On peut réunir les deux hypothèses en une seule hypothèse H_0 qui combine les deux, soit:

$$H_0 : \beta_1 = 0 \text{ et } \beta_2 = 0,$$

ce qui revient à dire que le prix d'une maison ne dépend pas du secteur.

C'est précisément ce que teste l'analyse de variance (ANOVA) fournie par MegaStats:

ANOVA table						
Source	SS	df	MS	F	p-value	
Regression	31 148.9961	2	15 574.4980	2.00	.1576	
Residual	187 099.6706	24	7 795.8196			
Total	218 248.6667	26				

Puisque $p = 0,1576$, on ne rejette pas H_0 .

Cette analyse est tout à fait équivalente à celle qu'on aurait faite au chapitre 9.

10.7 Exemples artificiels

Cette section a pour but d'illustrer comment le lien entre deux variables peut être influencé par une troisième. Dans le premier exemple, une dépendance entre deux variables est estompée sous l'effet d'une troisième. Dans le deuxième exemple, au contraire, une dépendance est créée artificiellement par une troisième variable. Les deux exemples qui suivent—fictifs et caricaturaux—sont conçus pour mettre le phénomène en évidence.

Exemple 10.7.1 Une corrélation éliminée sous l'effet d'une troisième variable

Dans la figure 10.4 présente les données du tableau 10.5, l'axe vertical est le score dans un test de vocabulaire et l'axe horizontal est l'âge. On ne décèle aucune dépendance entre les deux.

Un modèle qui ne tient compte que de l'âge, soit $E(\text{voc}) = \beta_0 + \beta_1(\text{age})$, on trouve $\hat{\beta}_0 = 49,65$ et $\hat{\beta}_1 = 0,0037$, ce qui donne l'équation suivante:

$$E(\text{voc}) = 49,65 + 0,0037(\text{age}).$$

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	49.648938	0.549067	90.424	<2e-16	***
age	0.003732	0.011286	0.331	0.742	Coefficients:
Multiple R-squared: 0.01196, Adjusted R-squared: -0.008625					
F-statistic: 0.581 on 1 and 48 DF, p-value: 0.4496					

Mais un r^2 de 0,012 et une valeur p de 0,742 ne permet pas de conclure que $\beta_1 \neq 0$. Conclusion provisoire: pas de relation entre l'âge et le vocabulaire.

Pourtant, le vocabulaire normalement augmente avec l'âge. Ce qui détruit en apparence cette dépendance, c'est une autre relation, celle entre l'âge et la scolarité: il se trouve que, dans la population, les personnes âgées sont moins scolarisées que les plus jeunes, et de ce fait ont un vocabulaire plus faible—plus faible que ne laisserait prévoir leur âge. L'effet positif de l'âge est annulé par l'effet négatif d'une faible scolarisation.

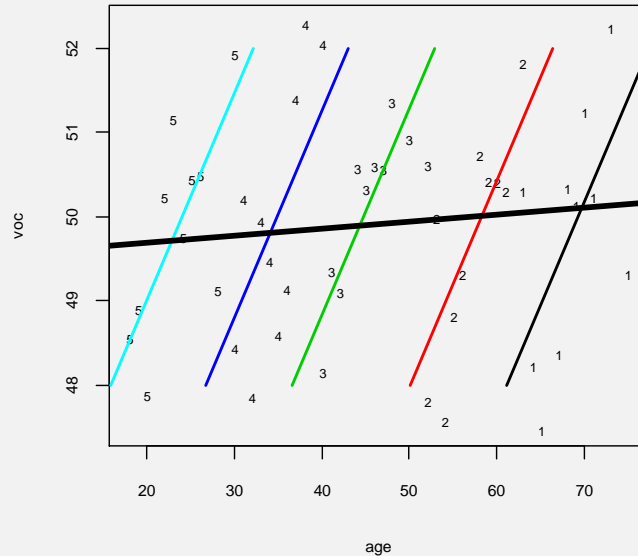
Dans le graphique, les points du nuage sont remplacés par des chiffres de 1 à 5 représentant des niveaux de scolarité: 1 pour le plus bas niveau, 5 plus le plus haut. On constate que le plus haut niveau de scolarité (5) se trouve parmi les plus jeunes.

On remarque que si on se limite à une scolarité donnée (par exemple, aux seuls points identifiés par «5», on décèle une relation très nette entre le vocabulaire et la scolarité—et ce pour tout niveau de scolarité.

On définit alors la variable scol, le niveau de scolarité, et on l'introduit dans le modèle

$$E(\text{voc}) = \gamma_0 + \gamma_1(\text{age}) + \gamma_2(\text{scol}).$$

Figure 10.4
Relation entre le vocabulaire et l'âge
selon le niveau de scolarité



Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	30.32895	3.10354	9.772	6.70e-13 ***
scol	2.80688	0.44696	6.280	1.01e-07 ***
age	0.24150	0.03878	6.227	1.22e-07 ***

Residual standard error: 0.9692 on 47 degrees of freedom
Multiple R-squared: 0.4575, Adjusted R-squared: 0.4344
F-statistic: 19.82 on 2 and 47 DF, p-value: 5.735e-07

On obtient les estimations $\hat{\gamma}_0 = 30,32$; $\hat{\gamma}_1 = 0,2415$; et $\hat{\gamma}_2 = 2,807$.

L'équation est donc la suivante:

$$E(\text{voc}) = 30,33 + 0,2415(\text{age}) + 2,807(\text{scol}).$$

Les tests sur γ_1 et γ_2 permettent de conclure que $\gamma_1 \neq 0$ ($p = 1,22 \times 10^{-7}$) et que $\gamma_2 \neq 0$ ($p = 1,01 \times 10^{-7}$). Le coefficient $\hat{\gamma}_1 = 0,2415$ représente l'accroissement du score de vocabulaire par année d'âge pour des personnes de même scolarité.

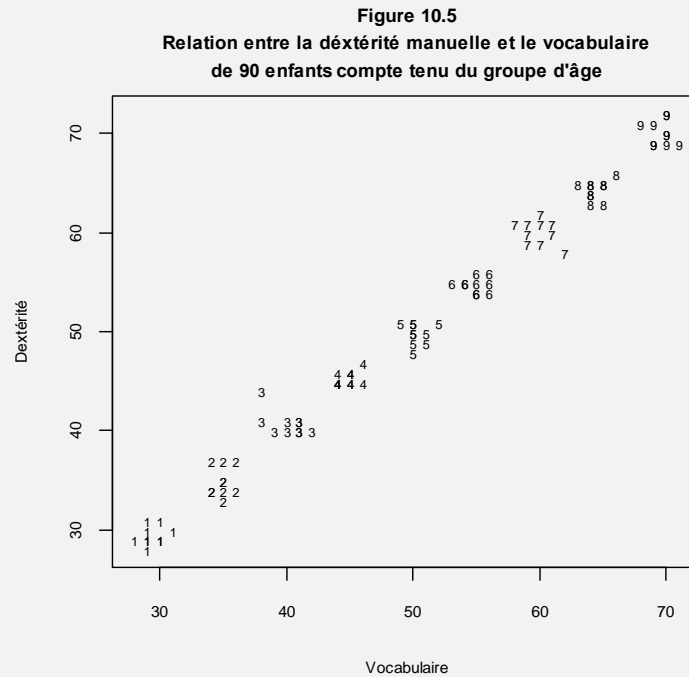
Remarque

Dans le dernier exemple, une dépendance réelle et positive (entre le vocabulaire et l'âge) est effacée sous l'effet d'une troisième variable, la scolarité. L'exemple 10.7.1 illustre la différence entre une corrélation ordinaire et ce qu'on appelle une *corrélation partielle*. La corrélation entre le vocabulaire et l'âge est sujette à l'action d'autres variables, dont la scolarité. Dans le cas de l'exemple, cette action a pour effet de réduire la corrélation à presque rien. La corrélation partielle *étant donné la scolarité* tente d'éliminer l'effet de cette troisième variable. C'est la corrélation qu'on aurait calculée si on pouvait se limiter à des individus ayant tous le même niveau de scolarité. On ne peut pas entrer ici dans le détail de ce calcul, mais le concept de corrélation partielle est révélé dans l'exemple: on voit nettement que si on se limite à un niveau de scolarité donné, la corrélation entre le vocabulaire et l'âge est forte. C'est le coefficient de corrélation partielle.

L'inverse de ce qu'illustre l'exemple 10.7.1 est possible aussi: une corrélation entre deux variables qui n'est due qu'à l'action d'une troisième variable. Ce phénomène est illustré dans l'exemple suivant.

Exemple 10.7.2 Une corrélation créée sous l'effet d'une troisième variable

La figure 10.5 présente la relation entre la dextérité manuelle (en ordonnée) et le score en un test de vocabulaire (en abscisse) (Tableau 10.6). On ne décèle aucune dépendance entre les deux.



Les chiffres dans le graphique représentent les groupes d'âge, « 1 » étant les plus jeunes, « 9 » les plus âgés. Globalement, la relation entre la dextérité et le vocabulaire est positive et prononcée, alors que, dans un groupe donné elle, de toute évidence, inexistante.

Considérons un modèle qui ne tient pas compte de l'âge, soit le modèle $E(Dext) = \beta_0 + \beta_1(voc)$, où la dextérité et le vocabulaire sont désignés par $Dext$ et voc , respectivement. On obtient les estimations $\hat{\beta}_0 = 0,4506$ et $\hat{\beta}_1 = 0,9961$. La relation est donc $E(Dext) = 0,4506 + 0,9961(voc)$.

Pour tenir compte de l'âge, on l'introduit la variable *groupe* dans la régression, le modèle étant $E(Dext) = \gamma_0 + \gamma_1(voc) + \gamma_2(groupe)$. On trouve $\hat{\gamma}_0 = 25,70308$; $\hat{\gamma}_1 = -0,02389$; $\hat{\gamma}_2 = 5,11072$. La relation est donc $E(Dext) = 25,70308 - 0,02389(voc) + 5,11072(groupe)$. On trouve que γ_1 est possiblement nul ($p = 0,678$) alors que γ_2 est certainement positif ($p = 2,2 \times 10^{-16}$).

Exercices

Dans ce qui suit, nous nous permettrons certains raccourcis afin d'alléger le langage.

- Nous dirons « x est significative » pour signifier que l'hypothèse que le coefficient de la variable exogène x est nul est rejetée à un certain niveau α .
- Par défaut, le niveau exigé d'un test d'hypothèse est présumé égal à 5 % (et le niveau d'un intervalle de confiance égal à 95 %).

- 10.1 [Données du tableau 10.4] Le tableau 10.4 présente des données sur un échantillon de 34 logements recueillies afin de déterminer une formule de prédiction du montant de la facture d'électricité. Les variables sont le montant de la facture (*facture*); le revenu du ménage (*revenu*); le nombre de personnes (*personnes*); et de la superficie (*surface*) du plancher du logement.
- Déterminer la matrice de corrélation des variables
 - Analyser le modèle $E(\text{facture}) = \beta_0 + \beta_1(\text{revenu}) + \beta_2(\text{personnes}) + \beta_3(\text{surface})$. Vérifier que la variable *revenu* n'est pas significative.
 - Analyser le modèle $E(\text{facture}) = \gamma_0 + \gamma_1(\text{personnes}) + \gamma_2(\text{surface})$. Tester les hypothèses usuelles $\gamma_1 = 0$ et $\gamma_2 = 0$.
 - Vérifier que la variable *revenu* comme seule variable exogène est néanmoins significative.
 - Vérifier que la variable *revenu* demeure significative en présence de *personnes*.
 - Vérifier que la variable *revenu* en présence de *surface* est à nouveau non significative
 - Il semblerait que *revenu* est significative dans tout modèle, à condition que *surface* n'en fasse pas partie. Peut-on expliquer ceci?
- 10.2 [Données du tableau A.3] Dans le modèle $E(\text{irm}) = \beta_0 + \beta_1g + \beta_2v + \beta_3p$, le coefficient de corrélation multiple est $R = 0,5366586$. Vérifier numériquement que R est le coefficient de corrélation entre $y = \text{irm}$ et $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1g + \hat{\beta}_2v + \hat{\beta}_3p$.
- 10.3 [Données du tableau A.2] Les données du tableau A.2 portent sur un échantillon de maisons vendues. On tente de déterminer lesquelles des variables disponibles (le nombre de salles de bains, le nombre de chambres à coucher, l'âge) pourraient servir à estimer le prix d'une maison.
- Déterminer la matrice de corrélation des variables.
 - S'il fallait utiliser une seule des variables exogènes pour prédire le prix, laquelle choisirait-on? Pourquoi?
 - Introduire les variables exogènes successivement, dans l'ordre décroissant de leur corrélation avec le prix. Examiner à chaque étape
 - l'estimation de l'écart-type conditionnel $\hat{\sigma}$ (il devrait décroître à chaque ajout);
 - le coefficient R ajusté. Commentez les gains apportés par l'ajout de chaque nouvelle variable.
 - À chaque étape, déterminer si les variables exogènes dans le modèle sont significatives.
 - On compare les coefficients dans les trois modèles suivants: $E(\text{prix}) = \beta_0 + \beta_1(\text{bains})$; $E(\text{prix}) = \gamma_0 + \gamma_1(\text{cac})$; et $E(\text{prix}) = \delta_0 + \delta_1(\text{bains}) + \delta_2(\text{cac})$. Expliquer les inégalités suivantes (examiner les signes des corrélations entre les variables concernées): i) $\hat{\beta}_1 > \hat{\delta}_1$; ii) $\hat{\gamma}_1 > \hat{\delta}_2$.
 - On compare les coefficients dans les trois modèles suivants: $E(\text{prix}) = \beta_0 + \beta_1(\text{bains})$; $E(\text{prix}) = \gamma_0 + \gamma_1(\text{age})$; et $E(\text{prix}) = \delta_0 + \delta_1(\text{bains}) + \delta_2(\text{age})$. Expliquer les inégalités suivantes (examiner les signes des corrélations entre les variables concernées): i) $\hat{\beta}_1 > \hat{\delta}_1$; ii) $\hat{\gamma}_1 > \hat{\delta}_2$.

- g) Revenons au modèle qui ne comprend que le nombre de salles de bains comme variable exogène. Estimer ce qu'une salle de bains ajoute en moyenne au prix d'une maison et déterminer un intervalle de confiance pour ce paramètre.

10.4 [Données du tableau A.2] Considérer maintenant un modèle qui comprend les trois variables exogènes énumérées.

- Estimer ce que vaut sur le marché une salle de bains supplémentaire (ce qu'elle ajoute en moyenne au prix d'une maison). Expliquer en quoi le sens de ce paramètre diffère de celui dans la question 10.2-g).
- Déterminer un intervalle de confiance pour la moyenne estimée en a).
- Prédire la valeur d'une maison vieille de 25 ans, ayant deux salles de bains et 3 chambres à coucher.
- Déterminer les limites de la prédiction que vous avez faite en iii) à 95 % et puis à 60 %.

10.5 [Données du Tableau 10.2] Le tableau 10.2 présente des données sur un groupe de professeurs afin d'identifier les facteurs qui contribuent au salaire en 2012 (*sal12*). Les variables exogènes possibles sont l'ancienneté (*anc*), le salaire à l'entrée (*sal0*), le sexe (*sexe*) et l'expérience préalable à l'engagement (*exp*).

- S'il fallait utiliser une seule variable comme variable exogène, laquelle devrait-on choisir? Pourquoi?
- Peut-on expliquer pourquoi le coefficient de corrélation entre *sal12* et *sal0* est négatif?
- Revenons à *anc* comme première variable exogène, et considérons l'ajout de *sal0* comme deuxième variable exogène. On a donc deux modèles:

$$\text{Modèle A: } E(\text{sal12}) = \beta_0 + \beta_1(\text{anc}); \text{ et Modèle B: } E(\text{sal12}) = \gamma_0 + \gamma_1(\text{anc}) + \gamma_2(\text{sal0})$$

- Comparer le coefficient de corrélation du modèle A au coefficient de corrélation multiple R du modèle B. Que peut-on conclure du fait que la différence est minuscule?
 - Vérifier que si R augmente à l'ajout de *sal0*, R ajusté baisse. Est-ce normal?
 - Pouvez-vous expliquer pourquoi dans le modèle B le coefficient de *sal0* est positif (alors qu'il est corrélé négativement avec *sal12*) et pourquoi il n'est plus significatif dans le modèle B?
- d) Considérer deux modèles pour prédire le salaire à l'entrée *sal0*:

$$\text{Modèle A: } E(\text{sal0}) = \beta_0 + \beta_1(\text{exp}) \text{ et Modèle B: } E(\text{sal0}) = \gamma_0 + \gamma_1(\text{exp}) + \gamma_2(\text{anc})$$

On constate que dans le modèle A, on ne peut pas affirmer que $\beta_1 \neq 0$ alors que dans le modèle B, on peut conclure avec confiance que $\gamma_1 > 0$. Comment explique-t-on ce paradoxe? Imaginer un graphique du style des figures 10.4 ou 10.5 pour illustrer ce phénomène.

10.6 [Données du tableau 10.2]

- Comparer les salaires moyens en 2012 des femmes et des hommes: estimer la différence de salaire et montrer qu'elle est significative.
- Vérifier, cependant, que les hommes ont plus d'ancienneté que les femmes.
- Est-ce possible que la différence de salaires ne soit due qu'au fait que les hommes ont plus d'ancienneté? Estimer la différence des salaires en tenant compte de l'ancienneté (déterminer une régression multiple avec l'ancienneté et le sexe comme variables exogènes).
- La différence est maintenant réduite par rapport à celle calculée en a). Est-elle significativement différente de 0?
- Résumer les conclusions en termes concrets.
- Le modèle développé en c) postule que la relation entre *sal12* et *anc* s'exprime par deux droites parallèles, l'une pour les femmes l'autre pour les hommes. Formellement, cela équivaut au

modèle suivant:

$$\text{Femmes: } E(\text{sal}12) = \beta_0 + \beta_1(\text{anc}); \quad \text{Hommes: } E(\text{sal}12) = \gamma_0 + \beta_1(\text{anc})$$

Remarquez que β_1 est le même pour les femmes et les hommes, ce qui assure que les droites sont parallèles. Estimer β_0 , β_1 , et γ_0 et tester l'hypothèse que $\beta_0 = \gamma_0$.

- g) Maintenant analyser un modèle qui n'impose pas de parallélisme entre les deux droites.

Formuler et tester l'hypothèse que chaque année d'ancienneté rapporte (en salaire) le même gain aux femmes qu'aux hommes.

- 10.7 [Données du tableau A.3] Le tableau A.3 présente des données visant à déterminer si un lien peut être établi entre la grosseur du cerveau (*irm*) et certains traits physiques et psychologiques. Dans ce numéro nous comparons la grosseur du cerveau des femmes et des hommes en tentant d'éliminer l'effet de la grandeur du corps.

- Montrer que la grosseur du cerveau (*irm*) est liée à la taille (tester pour montrer que la dépendance est significative)
- Montrer également que la grosseur du cerveau est liée au poids
- Montrer que dans le modèle $E(\text{irm}) = \beta_0 + \beta_1(\text{taille}) + \beta_2(\text{poids})$ on ne peut pas conclure que $\beta_2 \neq 0$, ce qui semble contredire la conclusion en b). Qu'est-ce qui pourrait expliquer la contradiction?
- Tester (indépendamment de toute autre variable) l'hypothèse que le cerveau des femmes est en moyenne égal à celui des hommes.
- La conclusion en d) (que les hommes ont un plus gros cerveau) serait-elle due uniquement au fait que les hommes sont plus grands? Il faudrait, pour que la comparaison soit juste, que la grosseur du cerveau soit relativisée par rapport à la taille (ou du poids). Une façon de le faire est de mesurer la grosseur du cerveau par $\text{irm}T = \text{irm}/\text{taille}$. Montrer qu'avec cette mesure, on ne peut pas conclure à une différence entre hommes et femmes quant à la grosseur du cerveau.
- Refaire l'analyse en e) en prenant pour mesure de la grosseur du cerveau le rapport $\text{irm}p = \text{irm}/\text{poids}$.

- 10.8 [Données du tableau A.3; voir l'exercice 10.6] Le tableau A.3 présente entre autres des données sur le poids (*poids*) et la taille (*taille*) d'un groupe d'hommes et de femmes. Considérer un modèle liant la taille et le poids par deux droites (femmes et hommes) de même ordonnée à l'origine mais de pentes différentes :

$$\text{Femmes: } E(\text{poids}) = \beta_0 + \beta_1(\text{taille})$$

$$\text{Hommes: } E(\text{poids}) = \beta_0 + \gamma_1(\text{taille})$$

- Estimer les deux équations $E(\text{poids}) = \hat{\beta}_0 + \hat{\beta}_1(\text{taille})$ et $E(\text{poids}) = \hat{\beta}_0 + \hat{\gamma}_1(\text{taille})$.
- Tester l'hypothèse $\beta_1 = \gamma_1$.
- Déterminer un intervalle de confiance pour chacun des paramètres β_1 .
- Déterminer un intervalle de confiance pour chacun des paramètres γ_1 .

- 10.9 [Données du tableau A.4]

Le tableau A.4 présente, entre autres, les scores B_1 , B_2 et B_3 de trois posttests composés par 66 sujets à la fin d'une période de formation. On se concentre ici sur B_1 , le posttest de compréhension et les scores A_1 et A_2 à deux pré-tests composés avant la période de formation.

- Montrer que chacune des variables A_1 et A_2 est séparément utile à la prédiction de B_1 .
- Montrer cependant que A_2 n'est plus significative (à 5 %) en présence de A_1 . Comment explique-t-on ce phénomène?

Les données du tableau ont été prélevées afin de comparer trois méthodes d'enseignement (trois *traitements*). L'échantillon est constitué de trois groupes, 1, 2 et 3, identifiés dans le tableau par la variable T . On veut donc savoir si B_1 dépend du traitement. À cette fin on définit les variables dichotomiques t_1 , t_2 et t_3 qui indiquent l'appartenance aux groupes 1, 2 et 3, respectivement.

- c) Montrer par une régression multiple qu'on peut conclure que B_1 dépend en effet du traitement (sans tenir compte des pré-tests).
- d) On tente maintenant, en comparant les traitements, de tenir compte des aptitudes initiales.
 - i) Vérifier pour commencer que les trois groupes diffèrent quant à la moyenne de A_1 .
 - ii) Il y a donc le risque que les différences dans la moyenne de B_1 soient attribuées faussement aux traitements alors qu'elles ne seraient dues qu'aux différences initiales entre les groupes. Considérer le modèle suivant: $E(B_1) = \gamma_j + \beta(A_1)$, où $\gamma_j = \gamma_1, \gamma_2$ ou γ_3 , selon le traitement. Estimer les paramètres γ_j et β .
- e) Tester séparément (mais dans le cadre du même modèle) les trois hypothèses $\gamma_1 = \gamma_2$; $\gamma_1 = \gamma_3$; et $\gamma_2 = \gamma_3$.
- f) Montrer que, selon l'analyse en c), on ne peut pas conclure à une différence entre les traitements 1 et 3 alors que selon l'analyse en e) on peut affirmer avec confiance que le traitement 3 est plus efficace que le traitement 1. Expliquer la contradiction dans un langage aussi concret que possible.
- g) Dans le modèle développé en b)
 - i) Déterminer un intervalle de confiance pour la différence (dans la moyenne de B_1) entre les groupes 1 et 2 pour une même aptitude initiale A_1 .
 - ii) Déterminer un intervalle de confiance pour la différence (dans la moyenne de B_1) entre les groupes 1 et 2 pour une même aptitude initiale A_1 .
 - iii) Déterminer un intervalle de confiance pour la moyenne de B_1 pour des personnes dans le groupe 1 pour lesquelles $A_1 = 9,8$.

10.10 [Données du tableau 10.4] Le tableau 10.4 présente les données suivantes sur un échantillon de 34 logements:

- a) Déterminer la matrice de corrélation des variables
- b) Déterminer une régression multiple pour estimer le montant de la facture à partir du revenu du ménage, du nombre de personnes et de la superficie du plancher du logement.
- c) Vérifier que la variable revenu n'est pas significative dans le modèle en b), et refaire la régression sans elle.
- d) Le tableau suivant présente des estimations dans quatre modèles pour la prédiction de facture, toutes comprenant la variable revenu comme variable exogène.

Modèle	Variables exogènes	Coefficient de revenu	Valeur p
A	revenu	0,2589917	6,973506e-10
B	revenu, personnes	0,2421046	3,887293e-12
C	revenu, surface	-0,1349820	0,1110451
D	revenu, personnes, surface	0,0751369 6	0,5849729

Il semblerait que le *revenu* est significatif dans tout modèle, tant que *surface* n'en fasse pas partie. Comment expliquer ceci? Résumer en une phrase et en termes concrets la conclusion que suggèrent ces résultats.

10.11 [Données du tableau 10.5] Le tableau 10.5 présente des données (fictives) sur l'âge (*age*), le score en un test de vocabulaire (*voc*) et le niveau de scolarité (*scol*).

- a) Considérer les deux modèles suivants:

$$E(voc) = \beta_0 + \beta_1(age) \quad \text{et} \quad E(voc) = \gamma_0 + \gamma_1(age) + \gamma_2(scol)$$

Qu'est-ce qui explique la différence entre $\hat{\beta}_1 = 0,0037$ et $\hat{\gamma}_1 = 0,2415$ (pourquoi $\hat{\beta}_1$ est-il tellement plus petit que $\hat{\gamma}_1$)?

- b) Considérer maintenant la scolarité comme une variable catégorielle : analyser le modèle

$$E(voc) = \gamma_j + \beta_1(age), j = 1, \dots, 5,$$

où j est le niveau de scolarité. Déterminer les 5 équations liant voc à age .

- c) À partir du modèle en b) tester l'hypothèse que $\gamma_1 = \gamma_2$.
 d) À partir du modèle en b) tester l'hypothèse que $\gamma_2 = \gamma_3$.

10.12 [Données du tableau 10.6; suite de l'exemple 10.7.2]

- a) Qu'est-ce qui explique la différence entre $\hat{\beta}_1 = 0,9961$ et $\hat{\gamma}_1 = -0,02389$? Pourquoi $\hat{\beta}_1$ est-il significativement différent de 0 alors que $\hat{\gamma}_1$ ne l'est pas?
 b) Considérer maintenant la scolarité comme une variable catégorielle : analyser le modèle

$$E(Dext) = \gamma_j + \beta(voc), j = 1, \dots, 9, \text{ où } j \text{ est le groupe d'âge.}$$

[Voici les estimations des γ : $\hat{\gamma}_1 = 32,618$; $\hat{\gamma}_2 = 38,701$; $\hat{\gamma}_3 = 45,052$; $\hat{\gamma}_4 = 50,251$; $\hat{\gamma}_5 = 55,334$; $\hat{\gamma}_6 = 60,722$; $\hat{\gamma}_7 = 66,552$; $\hat{\gamma}_8 = 71,329$; $\hat{\gamma}_9 = 77,580$ et $\hat{\beta} = -0,1060$.

10.13 [Données du tableau 10.7] Le tableau 10.7 présente des données médicales sur un échantillon de 332 sujets d'origine indienne Pima d'Arizona. L'objectif dans ce numéro est de tenter d'identifier les facteurs qui contribuent à la tension artérielle.

- a) Analyser le modèle $E(tension) = \beta_0 + \beta_1(imc) + \beta_2(age) + \beta_3(peau) + \beta_4(glu) + \beta_4(gros)$.
 b) À partir du modèle en a), éliminer les variables exogènes non significative s'il y lieu. Procéder par étapes: refaire une régression après avoir éliminé la variable exogène la moins significative (dont la *valeur p* est la plus grande). Recommencer avec le modèle réduit, ainsi de suite jusqu'à ce que toutes les variables exogènes sont significatives. Comparer le R du modèle final au R du modèle initial, question de s'assurer que l'élimination des variables exogènes ne cause pas d'importantes pertes de précision.
 c) Vérifier la relation entre *tension* et *imc*, ainsi que la relation entre *tension* et *peau* sont toutes deux significatives.
 d) Vérifier, cependant, que dans le modèle $E(tension) = \beta_0 + \beta_1(imc) + \beta_2(peau)$, on ne peut pas conclure que $\beta_2 \neq 0$. Comment s'expliquerait cette apparente contradiction?
 e) Vérifier que dans le modèle $E(tension) = \gamma_0 + \gamma_1(age) + \gamma_2(peau)$, on peut conclure avec confiance que $\gamma_2 > 0$. Expliquer pourquoi cette conclusion ne contredit pas nécessairement celles énoncées en a) et en b).
 f) Le tableau suivant présente quelques résultats d'analyse de quatre modèles. La variable endogène est *tension* et une des variables exogènes est *gros* dans tous les cas. Chaque modèle comprend une autre variable exogène. Comment se fait-il que *gros* est significatif dans trois cas et non significatif dans le quatrième?

Modèle	Variables exogènes	Coefficient de <i>gros</i>	Valeur <i>p</i>
A	<i>imc</i> et <i>gros</i>	0,7207	0,0003
B	<i>peau</i> et <i>gros</i>	0,6362	0,0024
C	<i>glu</i> et <i>gros</i>	0,6329	0,0026
D	<i>age</i> et <i>gros</i>	-0,2624	0,3363

10.14 [Données du tableau A.3] Le tableau A.3 présente des données visant à déterminer si un lien peut être établi entre la grosseur du cerveau et certains traits physiques et psychologiques. On désigne par

irm la grosseur du cerveau, par v et p les résultats aux tests V et P de Wechsler. (Nous délaissions la variable g car les données la concernant semblent entachées d'erreurs).

- a) i) Vérifier que la relation entre irm et p est significative.
- ii) Noter, cependant, que dans le modèle $E(irm) = \beta_0 + \beta_1v + \beta_2p$, on ne peut rejeter aucune des deux hypothèses $\beta_1 = 0$ et $\beta_2 = 0$.
- b) Comment s'explique la contradiction entre les conclusions b-i) et b-ii)?

10.15 [Données du tableau A.3] On désigne par *sexe* une variable dichotomique désignant le sexe (1 = homme; 0 = femme).

- a) Montrer que dans le modèle $E(irm) = \beta_0 + \beta_1(sexe) + \beta_2p$, *sexe* et p sont tous deux significatifs.
- b) Supposons qu'on prenne pour mesure de la grosseur du cerveau le rapport $irmt = irm/taille$. Montrer que dans le modèle $E(irmt) = \gamma_0 + \gamma_1(sexe) + \gamma_2p$ *sexe* n'est plus significatif mais p l'est encore.
- c) Résumer les résultats en a) et b), en évitant autant que possible les termes techniques.

10.16 [Données du tableau A.3] Le modèle de régression n'est pas uniquement un outil de prédiction. Parfois, le but est simplement d'établir qu'une corrélation existe. Dans ce cas, il n'y a pas lieu de distinguer les variables endogènes des variables exogènes. On vous demande ici de vérifier empiriquement que le choix qu'on fait n'a pas d'importance.

- a) Considérer les deux modèles suivants: $E(irm) = \beta_0 + \beta_1p$ et $E(p) = \gamma_0 + \gamma_1(irm)$. Vérifier que la valeur p correspondant à l'hypothèse $\beta_1 = 0$ est identique à la valeur p correspondant à l'hypothèse $\gamma_1 = 0$.
- b) Considérer les deux modèles suivants:

Modèle A: $E(irm) = \beta_0 + \beta_1p + \beta_2(sexe)$ et Modèle B: $E(p) = \gamma_0 + \gamma_1(irm) + \gamma_2(sexe)$.

Soit p_A la valeur p obtenue dans le modèle A et p_B la valeur p obtenue dans le modèle B. Avant de calculer, dites quelles seraient vos interprétations sous les hypothèses suivantes? (α est le niveau du test): i) $p_A < \alpha$ et $p_B \geq \alpha$; ii) $p_A \geq \alpha$ et $p_B < \alpha$.

- c) Maintenant vérifier que la valeur p correspondant à l'hypothèse $\beta_1 = 0$ dans le modèle A est identique à la valeur p correspondant à l'hypothèse $\gamma_1 = 0$ dans le modèle B.
- d) Comment concilier les valeurs p correspondant à *sexe* dans les deux modèles? (Concrètement, qu'est-ce qu'on affirme dans un cas et qu'est-ce qu'on affirme dans l'autre?)