

11.1. Les modèles de régression à variables qualitatives indépendantes ("dummy variable reg. models")

A) Régression sur une seule variable qualitative : ANOVA

Exemple:

$$Y_i = \alpha + \beta D_i + u_i$$

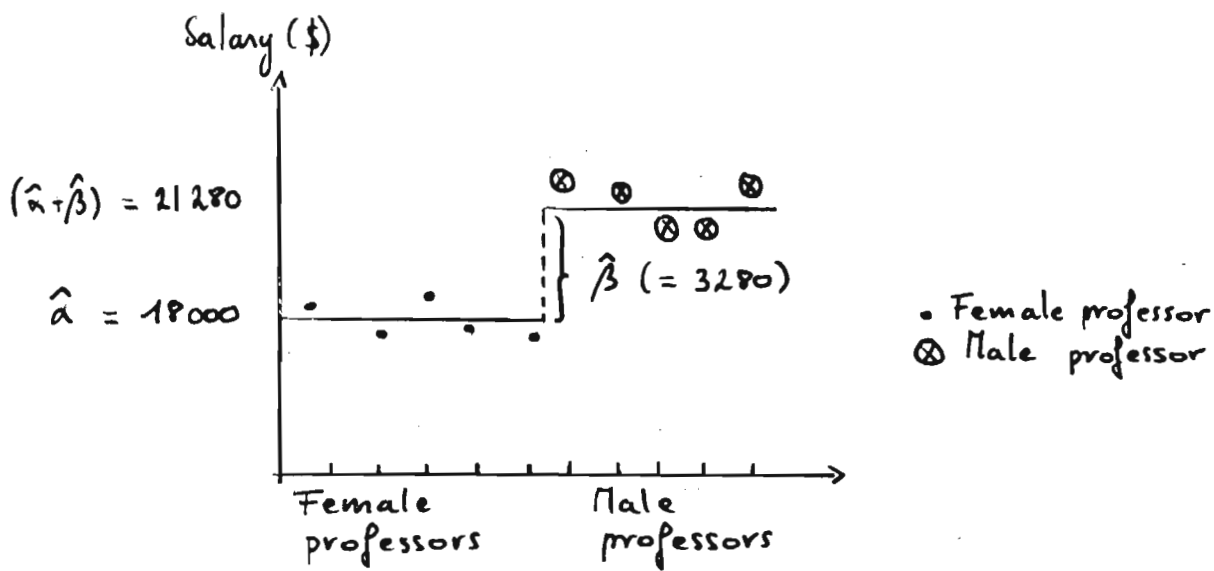
$$\text{où } \begin{cases} Y = \text{salaire annuel d'un prof. d'université.} \\ D_i = 1 \text{ si professeur est un homme.} \\ \quad = 0 \text{ sinon} \\ u_i = \text{terme d'erreur.} \end{cases}$$

Hypothèse : u_i satisfait aux hyp. du modèle de rég. linéaire classique \Rightarrow $E(Y_i / D_i = 0) = \alpha \equiv$ salaire moyen ♀
 $E(Y_i / D_i = 1) = \alpha + \beta \equiv$ salaire moyen ♂

$$\left[\begin{array}{l} \hat{Y}_i = 18,00 + 3,28 D_i \\ \text{s.e.} = (0,32) \quad (0,44) \\ t = (57,74) \quad (7,439) \\ R^2 = 0,8737, \quad n = 10, \quad Y \text{ en milliers de \$} \end{array} \right.$$

$$\hat{\alpha} = 18000 \$$$

$$\hat{\alpha} + \hat{\beta} = \pm 21000 \$ \quad (21280 \$)$$



B) Régression sur une variable qualitative et une variable quantitative : ANCOVA

Variables quantitatives de modèles ANCOVA = "covariates"

$$Y_i = \alpha_1 + \alpha_2 D_i + \beta X_i + u_i$$

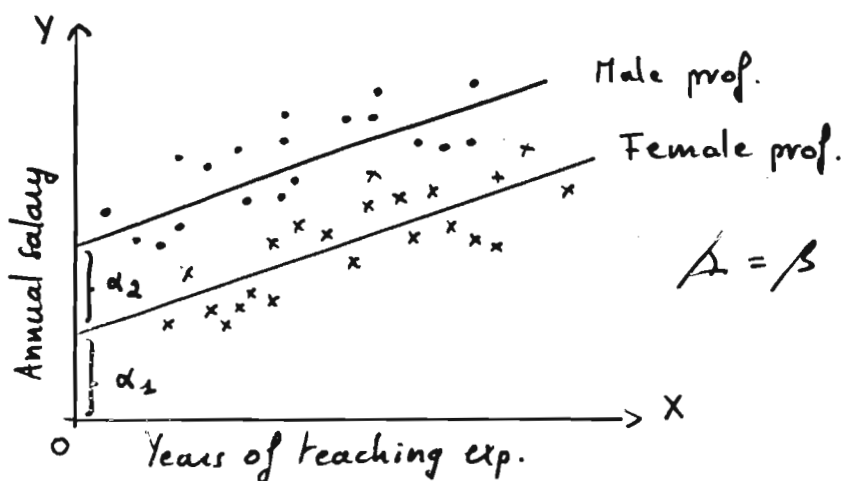
où

- Y = salaire annuel d'un prof.
- X_i = nombre d'années d'expérience.
- $D_i = 1$, si le prof. est un homme
 $= 0$ sinon
- u_i = terme d'erreur

Si $E(u_i) = 0$:

$$E(Y_i | X_i, D_i = 0) = \alpha_1 + \beta X_i \quad \equiv \text{sal. moyen } \text{♀}$$

$$E(Y_i | X_i, D_i = 1) = (\alpha_1 + \alpha_2) + \beta X_i \quad \equiv \text{sal. moyen } \text{♂}$$



• Caractéristiques :

•) Une seule dummy (D_i) pour distinguer les ♀ des ♂.
Quid si une dummy pour chaque sexe?

$$Y_i = \alpha_1 + \alpha_2 D_{2i} + \alpha_3 D_{3i} + \beta X_i + u_i$$

où $\left\{ \begin{array}{l} Y = \text{salaire annuel d'un prof.} \\ X_i = \text{nbre d'années d'exp.} \\ D_{2i} = 1 \text{ si prof. est un homme} \\ \quad = 0 \text{ sinon} \\ D_{3i} = 1 \text{ si prof. est une femme} \\ \quad = 0 \text{ sinon} \\ u_i = \text{terme d'erreur} \end{array} \right.$

Colinéarité pftre les D_{2i} et D_{3i} .

Illustration : éch. de 3 prof ♂ et 2 prof. ♀

$\exists \lambda_1, \lambda_2, \lambda_3$
 non tous nuls
 tq: $\lambda_1 X_1 + \lambda_2 X_2 + \lambda_3 X_3 = 0$
 pq?
 soit $\lambda_1 = 1$
 $\lambda_2 = -1$
 $\lambda_3 = -1$

		X_1	X_2	X_3	X
		Intercepte	D_2	D_3	
Homme	Y_1	1	1	0	X_1
Homme	Y_2	1	1	0	X_2
Femme	Y_3	1	0	1	X_3
Homme	Y_4	1	1	0	X_4
Femme	Y_5	1	0	1	X_5

$\Rightarrow X_1 - X_2 - X_3 = 0 \quad D_2 = 1 - D_3 \quad \text{ou} \quad D_3 = 1 - D_2 \quad (1 \equiv \text{intercepte})$

\Rightarrow multicolin. Estimation par MCO n'est pas possible.

Règle générale : si une variable qualitative comporte m catégories, $(m-1)$ variables dummy doivent être incluses dans la régression

(si règle non respectée : "dummy variable trap")

2) L'attribution des chiffres 0 et 1 est peut-être arbitraire.

$$\left. \begin{array}{l} \text{Si } D = 1 \text{ si prof } \varphi \\ = 0 \text{ sinon} \end{array} \right\} \Rightarrow \begin{array}{l} E(Y_i | X_i, D_i = 1) = (\alpha_1 + \alpha_2) + \beta X_i \\ E(Y_i | X_i, D_i = 0) = \alpha_1 + \beta X_i \end{array}$$

Catégorie de référence / de base !

3) Intercepte mesure valeur moyenne de la catégorie de référence.

4) Coeff. régress. associés aux variables dummy
≡ différentiels d'intercepte
("differential intercept coefficients")

c) Régression sur une variable quantitative et une variable qualitative à plus de 2 catégories.

$$Y_i = \alpha_1 + \alpha_2 D_{2i} + \alpha_3 D_{3i} + \beta X_i + u_i$$

où

$$\left\{ \begin{array}{l} Y_i = \text{dépenses annuelles de santé de l'individu } i \\ X_i = \text{revenu annuel de l'individu } i \\ D_{2i} = 1 \text{ si l'individu } i \text{ a un diplôme du second} \\ = 0 \text{ sinon} \\ D_{3i} = 1 \text{ si l'individu } i \text{ a un diplôme du supérieur} \\ \text{ou un diplôme universitaire} \\ = 0 \text{ sinon} \\ u_i = \text{terme d'erreur stochastique} \end{array} \right.$$

Catégorie de référence = niveau d'éducation primaire ou moins.

α_1 = intercepte catég. de base

α_2 et α_3 indiquent de combien les interceptes des deux autres catég. varient // intercepte catég. de base.

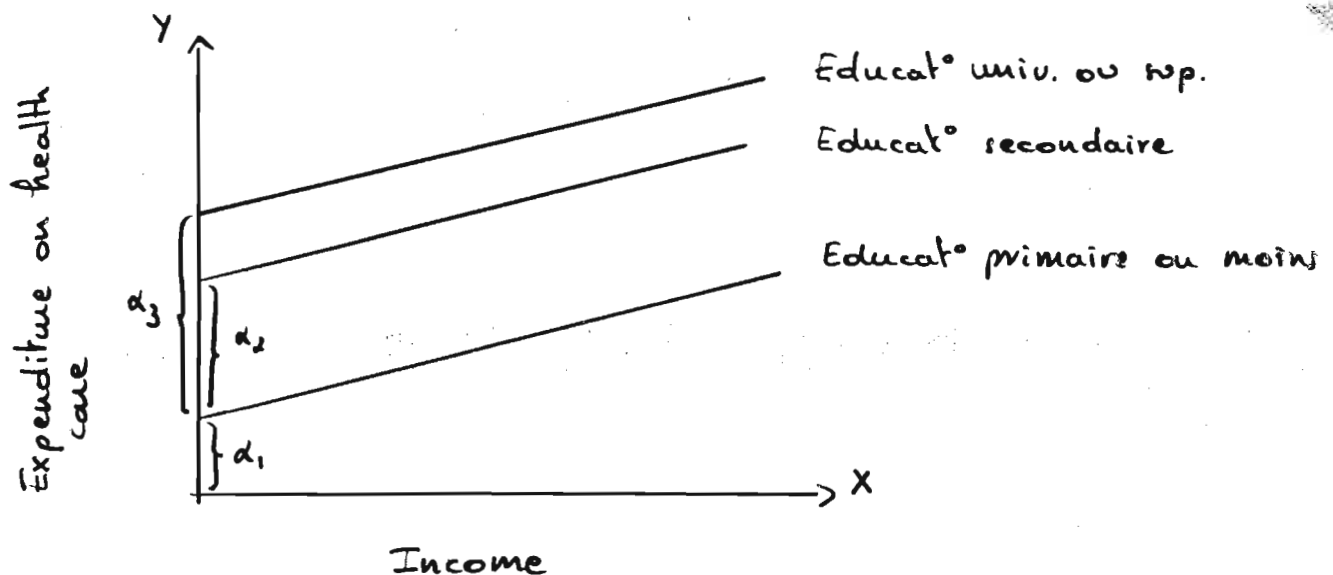
$$\text{Si } E(u_i) = 0$$

$$E(Y_i | D_2 = 0, D_3 = 0, X_i) = \alpha_1 + \beta X_i$$

$$E(Y_i | D_2 = 1, D_3 = 0, X_i) = (\alpha_1 + \alpha_2) + \beta X_i$$

$$E(Y_i | D_2 = 0, D_3 = 1, X_i) = (\alpha_1 + \alpha_3) + \beta X_i$$

$$\text{si } \alpha_3 > \alpha_2 > 0$$



D) Régression sur ~~une~~ 2 variables qualitatives et ~~deux~~ 1 variable quantitative

$$Y_i = \alpha_1 + \alpha_2 D_{2i} + \alpha_3 D_{3i} + \beta X_i + u_i$$

où

- Y_i = salaire annuel de l'individu i
- X_i = nombre d'années d'expérience de l'individu i
- $D_{2i} = 1$ si l'individu i est un homme
= 0 sinon
- $D_{3i} = 1$ si l'individu i est blanc
= 0 sinon
- u_i = terme d'erreur

individu = professeur.

Catégorie de référence : professeurs féminins noirs

Si $E(u_i) = 0$:

$E(Y_i / D_2 = 0, D_3 = 0, X_i) = \alpha_1 + \beta X_i \equiv \text{sal. moyen prof } \text{♀} \text{ noirs}$

$E(Y_i / D_2 = 1, D_3 = 0, X_i) = (\alpha_1 + \alpha_2) + \beta X_i \equiv \text{sal. moyen prof. } \text{♂} \text{ noirs}$

$E(Y_i / D_2 = 0, D_3 = 1, X_i) = (\alpha_1 + \alpha_3) + \beta X_i \equiv \text{sal. moyen prof } \text{♀} \text{ blancs}$

$E(Y_i / D_2 = 1, D_3 = 1, X_i) = (\alpha_1 + \alpha_2 + \alpha_3) + \beta X_i \equiv \text{sal. moyen prof } \text{♂} \text{ blancs}$

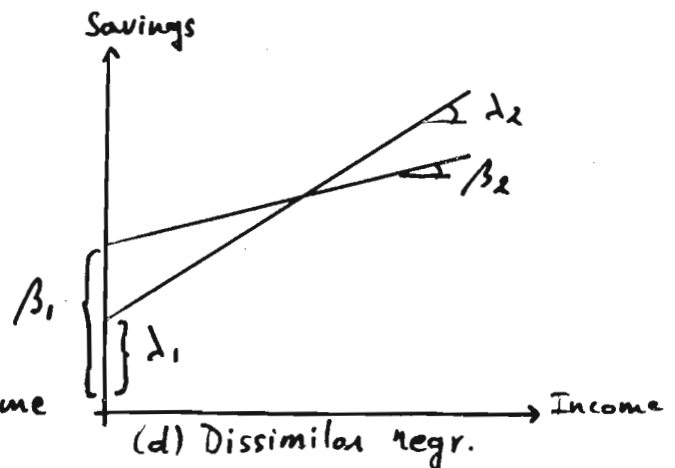
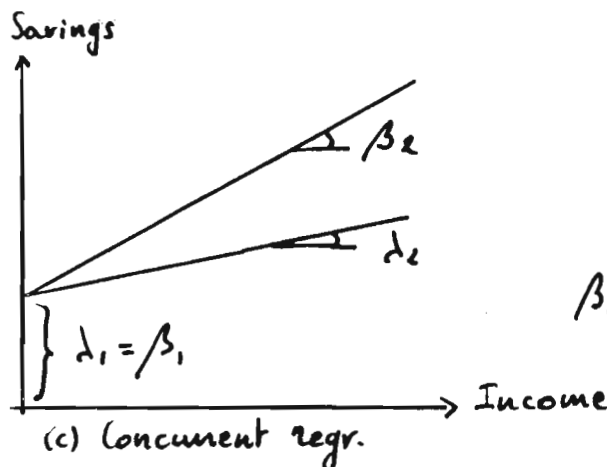
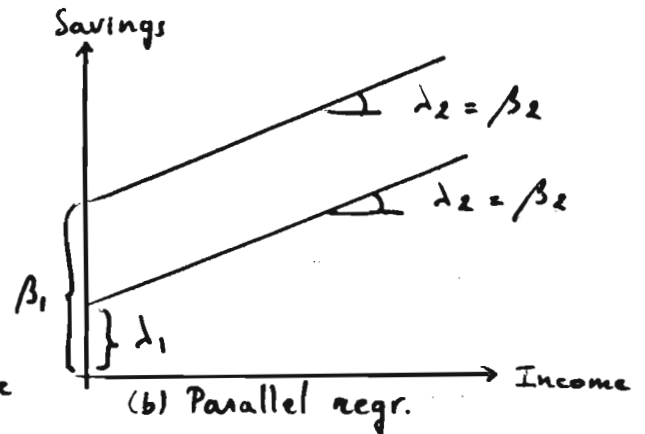
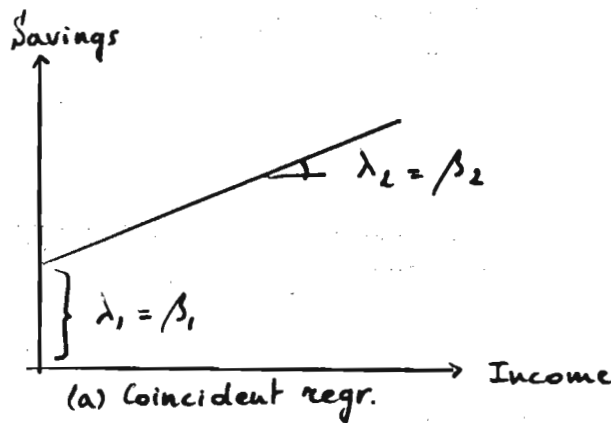
E) Les variables dummy comme alternative au test de Chow

- Rappel : test de Chow pour tester stabilité des paramètres dans relation tel épargne et revenu des ménages sur la période 1970 - 1995 aux USA

Période 1970 - 1981 : $Y_t = \lambda_1 + \lambda_2 X_t + u_{1t} \quad (n_1 = 12) \quad (1)$

Période 1982 - 1995 : $Y_t = \beta_1 + \beta_2 X_t + u_{2t} \quad (n_2 = 14) \quad (2)$

Période 1970 - 1995 : $Y_t = \alpha_1 + \alpha_2 X_t + u_{3t} \quad (n = 26) \quad (3)$



- En pratique, pour déterminer la source de la \neq les 2 régressions, on regroupe ens. des dos. et on estime modèle suivant:

$$Y_t = \alpha_1 + \alpha_2 D_t + \delta_1 X_t + \delta_2 (D_t \cdot X_t) + u_t \quad (4)$$

où

$$\left\{ \begin{array}{l} Y = \text{épargne} \\ X = \text{revenu} \\ t = \text{temps} \\ D = 1 \text{ pour obs. de la période } 1982-1995 \\ \quad = 0 \text{ sinon} \\ u = \text{terme d'erreur} \end{array} \right.$$

- ▷ Si $E(u_i) = 0$:

$$E(Y_t / D_t = 0, X_t) = \alpha_1 + \delta_1 X_t \quad \equiv \text{Épargne moyenne période } 70-81$$

$$E(Y_t / D_t = 1, X_t) = (\alpha_1 + \alpha_2) + (\delta_1 + \delta_2) X_t \quad \equiv \text{Épargne moyenne période } 82-95$$

- ▷ Si $(\lambda_1 = \alpha_1)$ et $(\lambda_2 = \delta_1)$
 $(\beta_1 = \alpha_1 + \alpha_2)$ et $(\beta_2 = \delta_1 + \delta_2)$ } on obtient:

$$\Rightarrow Y_t = \lambda_1 + \lambda_2 X_t + u_{1t} \quad (\text{période } 1970 - 1981) \quad (1)$$

$$Y_t = \beta_1 + \beta_2 X_t + u_{2t} \quad (\text{période } 1982 - 1995) \quad (2)$$

- ▷ Dans Equation (4):

$\alpha_2 \equiv$ différentiel d'intercepte les 2 périodes

$\delta_2 \equiv$ différentiel de pente les 2 périodes

- ▷ Inclusion d'une dummy sous la forme:

- interactive / multiplicative permet de tester des \neq de pentes

- additive permet de tester des \neq d'interceptes

- $$\hat{Y}_t = 1,0161 + 152,4786 D_t + 0,0803 X_t - 0,0655 (D_t X_t)$$

$$se = (20,1648) \quad (33,0824) \quad (0,0144) \quad (0,0159)$$

$$t = (5,5043) \quad (4,6090) \quad (5,5413) \quad (-4,0963)$$

où

$$\left\{ \begin{array}{l} Y = \text{Épargne} \\ X = \text{revenu} \\ t = \text{temps} \\ D = 1 \text{ pour obs. 1e1 1982 et 1995} \\ \quad = 0 \text{ sinon} \end{array} \right.$$

- Fonctions de régression de chaque période :

$$\hat{Y}_t = 1,0161 + 0,0803 X_t \quad \equiv \text{régr. Épargne - revenu pour 1970 - 1981}$$

$$\begin{aligned} \hat{Y}_t &= (1,0161 + 152,4786) \\ &\quad + (0,0803 - 0,0655) X_t \\ &= 153,4947 + 0,0148 X_t \quad \equiv \text{régr. Épargne - revenu pour 1982 - 1995} \end{aligned}$$

- Avantages "dummy" par rapport test de Chow ?

1) On estime 1 seule régression.

2) Estimation + précise car + de df.

3) Permet de tester de nombreuses hyp :

* tester si $\alpha_2 = \delta_2 = 0$ avec F-test

permet de vérifier l'∃ d'un choc structural

* tester à l'aide d'un test en t si ce sont

les interceptes et/ou les pentes qui sont ≠ :

($\alpha_2 = 0$ et/ou $\delta_2 = 0$)

5) Les effets d'interactions en utilisant des dummy

Impact sur la variable dépendante de l'interaction entre plusieurs variables explicatives.

Illustration:

$$Y_i = \alpha_1 + \alpha_2 D_{2i} + \alpha_3 D_{3i} + \beta X_i + u_i \quad (1)$$

où

- Y = salaire horaire en \$
- X = nombre d'années d'éducation
- $D_{2i} = 1$ si individu i est une femme
 $= 0$ sinon
- $D_{3i} = 1$ si individu i est ni blanc, ni hispanique
 $= 0$ sinon
- u_i = terme d'erreur

Hypothèses implicites :

- i) α_2 est identique pour individus ayant couleurs de peau \neq
- ii) α_3 est identique pour φ et σ

Pas très réalistes*. Il est possible que lorsque les 2 variables indicatrices valent 1 qu'il y ait un effet additionnel sur le salaire moyen. (interaction tel D_2 et D_3)

\Rightarrow effet des variables D_2 et D_3 sur salaire moyen n'est pas nécess. uniq. additif comme dans Equation (1) il peut aussi être multiplicatif.

* Il est possible que l'effet du sexe sur la rémunération soit + sup. pour les blancs -1.49- ou les hisp que pour les noirs.

- Prise en compte de l'effet multiplicatif :

$$Y_i = \alpha_1 + \alpha_2 D_{2i} + \alpha_3 D_{3i} + \alpha_4 (D_{2i} D_{3i}) + \beta X_i + u_i \quad (2)$$

$$\Rightarrow E(Y_i / D_{2i} = 1, D_{3i} = 1, X_i) = (\alpha_1 + \alpha_2 + \alpha_3 + \underline{\alpha_4}) + \beta X_i$$

Salaires horaires moyens des ♀ ni blanches, ni hisp. =

$\alpha_1 \equiv$ sal. horaire moyen de la catég. de réf.

(hommes blancs ou hisp.)

$\alpha_2 \equiv$ effet différentiel d'être une ♀

$\alpha_3 \equiv$ effet différentiel d'être ni blanc, ni hisp.

$\alpha_4 \equiv$ effet différentiel d'être une ♀ ni blanche, ni hisp.

- En l'absence de variables d'interaction (uniq. effet additif) :

$$E(Y_i / D_{2i} = 1, D_{3i} = 1, X_i) = (\alpha_1 + \alpha_2 + \alpha_3) + \beta X_i$$

\Rightarrow Introduction des variables dummy sous une forme multiplicative modifie l'effet combiné des variables qualitatives sur la variable dépendante.

Exemple :

$$\begin{aligned} \hat{Y}_i &= -0,2610 - 2,3606 D_{2i} - 1,7327 D_{3i} + 0,8028 X_i \\ (1) \quad t &= (-5,367) \quad (-5,4873) \quad (-2,1803) \quad (9,9094) \\ R^2 &= 0,2032 \quad , \quad n = 528 \end{aligned}$$

$$\begin{aligned} \hat{Y}_i &= -0,2610 - 2,3606 D_{2i} - 1,7327 D_{3i} + 2,1289 \underline{D_{2i} D_{3i}} + 0,8028 X_i \\ (2) \quad t &= (-5,367) \quad (-5,4873) \quad (-2,1803) \quad (1,9802) \quad (9,9095) \\ R^2 &= 0,2532 \quad , \quad n = 528 \end{aligned}$$

$$\Rightarrow \hat{\alpha}_2 + \hat{\alpha}_3 + \hat{\alpha}_4 = -1,964 \quad \text{compris } |e| - 2,3606 \text{ \& } -1,7327$$

Interaction 1el dummy et variable quantitative

Rendement d'une année d'éducation identique pour les ♀ et le ♂ ?

$$Y_i = \alpha_1 + \alpha_2 D_{2i} + \alpha_3 D_{3i} + \beta X_i + \delta (D_{2i} X_i) + u_i$$

Si $E(u_i) = 0$:

$$E(Y_i | D_{2i} = 1, D_{3i} = 0, X_i) = \alpha_1 + \alpha_2 + (\beta + \underline{\delta}) X_i \quad \text{♀ BH}$$

$$E(Y_i | D_{2i} = 0, D_{3i} = 0, X_i) = \alpha_1 + \beta X_i \quad \text{♂ BH}$$

$$E(Y_i | D_{2i} = 1, D_{3i} = 1, X_i) = \alpha_1 + \alpha_2 + \alpha_3 + (\beta + \underline{\delta}) X_i \quad \text{♀ A}$$

$$E(Y_i | D_{2i} = 0, D_{3i} = 1, X_i) = \alpha_1 + \alpha_3 + \beta X_i \quad \text{♂ A}$$

δ = différentiel de pente (= de rendement) pour les ♀

6) L'utilisation des variables dummy ds l'analyse saisonnière

Nbreuses séries éco. (p.ex. mensuelles, trimestrielles) présentent des motifs oscillatoires réguliers (structure saisonnière).

Ajustement saisonnier = processus qui consiste à enlever comp. saisonnière d'une série temporelle.

Méthode des variables dummy pour désaisonnaliser une série temporelle.

- Exemple : Ventes de réfrigérateurs aux Etats-Unis entre 1978 et 1995

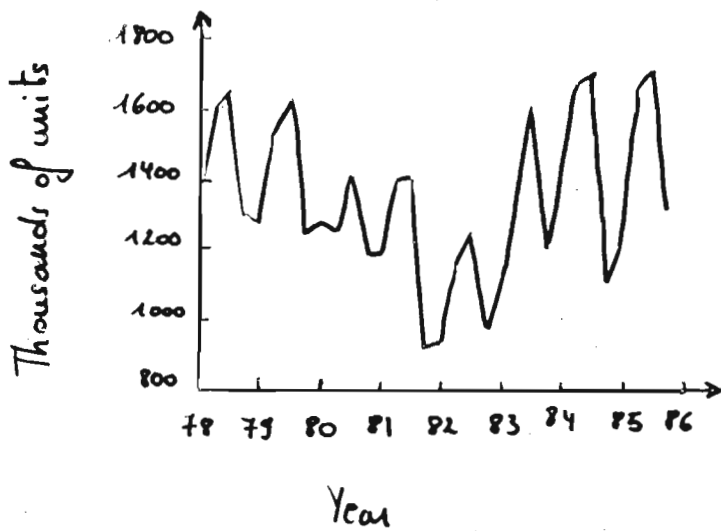


Fig. : Sales of refrigerators 1978-1986 (quarterly)

$$Y_t = \alpha_1 + \alpha_2 D_{2t} + \alpha_3 D_{3t} + \alpha_4 D_{4t} + u_t$$

où $\left\{ \begin{array}{l} Y = \text{ventes de réfrigérateurs (milliers d'unités)} \\ D_{2,3,4,t} = \text{variables dummy qui prennent valeur 1 (ou 0 sinon)} \\ \text{si donnée relative resp. au 2ème, 3ème ou} \\ \text{4ème trimestre} \end{array} \right.$

$$\Rightarrow \hat{Y}_t = 1222,13 + 245,37 D_{2t} + 347,63 D_{3t} - 62,13 D_{4t}$$

$$t = (20,37) \quad (2,89) \quad (4,10) \quad (-0,73)$$

$$R^2 = 0,53$$

- Obtention d'une série désaisonnalisée ?

$$\hat{u}_t = Y_t - \hat{Y}_t$$

Série désaisonnalisée = résidus \hat{u}_t de la régression des ventes de réfrigérateurs sur une cst et 3 dummy.

Série de résidus = série initiale - composante saisonnière
 (il reste : composante cyclique, tendance et comp. aléatoire)

- Introduction d'une variable expl. quantitative (covariée) :

$$\hat{Y}_t = 456,24 + 242,50 D_{2t} + 325,26 D_{3t} - 86,08 D_{4t} + 2,77 X_t$$

$$t = (2,56) \quad (3,70) \quad (4,94) \quad (-1,31) \quad (4,45)$$

$$R^2 = 0,729$$

où $\begin{cases} Y = \text{ventes de réfrig. (milliers d'unités)} \\ D_t = \text{dummy pour trimestres} \\ X_t = \text{dépenses en biens durables des ménages (milliers de \$)} \\ t = \text{temps} \end{cases}$

X_t se caractérise églmt par une comp. saisonnière ?

- Théorème de Frish - Waugh montre que l'inclusion des variables dummy D_2, D_3, D_4 (qui contrôlent pour la saisonnalité de la variable Y) suppriment la composante saisonnière de la variable X .

Vérification :

- 1) $Y \sim \text{cst}, D_2, D_3, D_4 \rightarrow S_1$ (variable Y désais.)
- 2) $X \sim \text{cst}, D_2, D_3, D_4 \rightarrow S_2$ (var. X désais.)
- 3) $S_1 \sim S_2$ (cst) (coeff. ass. à S_2 nû que coeff. ass. à X)
 (ds régr. $Y \sim \text{cst}, D_2, D_3, D_4, X$)

H) La régression linéaire par morceau ("piecewise linear regr.")

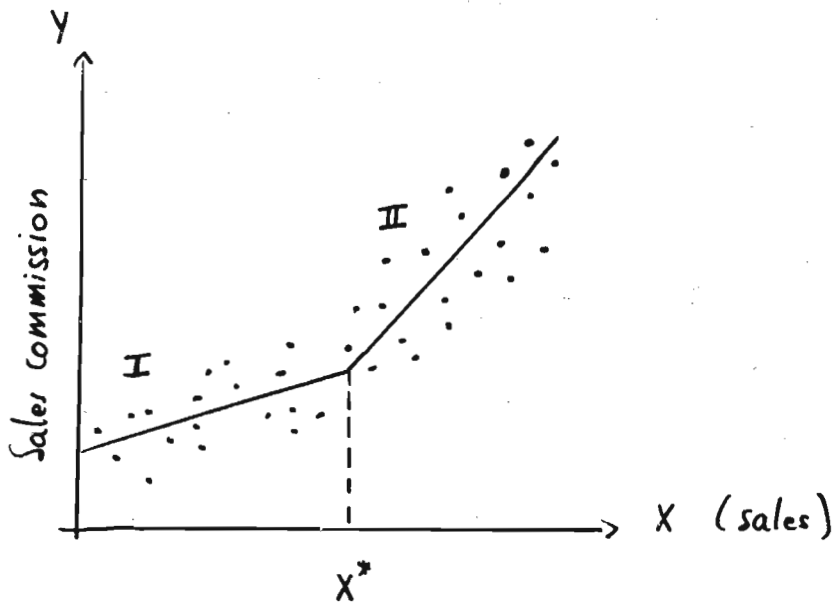


Fig.: Hyp. relationship b/w sales commission and sales volume

Note: The intercept on the Y axis denotes min. guaranteed commission.

• Estimation des segments linéaires (I) et (II) ?

$$Y_i = \alpha_i + \beta_1 X_i + \beta_2 (X_i - X^*) D_i + u_i$$

où

$$\begin{cases} Y_i = \text{commission de l'individu "i"} \\ X_i = \text{volume de vente de l'individu "i"} \\ X^* = \text{valeur critique pour les ventes (connue a priori)} \\ D = \begin{cases} 1 & \text{si } X_i > X^* \\ 0 & \text{si } X_i < X^* \end{cases} \end{cases}$$

Remarque: X^* pas tjs connu a priori

i) Analyse graphique

ii) "Switching regression models"

- Si $E(u_i) = 0$:

$$E(Y_i | D_i = 0, X_i, X^*) = \alpha_1 + \beta_2 X_i \quad (X_i < X^*)$$

$$E(Y_i | D_i = 1, X_i, X^*) = \alpha_1 - \beta_2 X^* + (\beta_1 + \beta_2) X_i \quad (X_i > X^*)$$

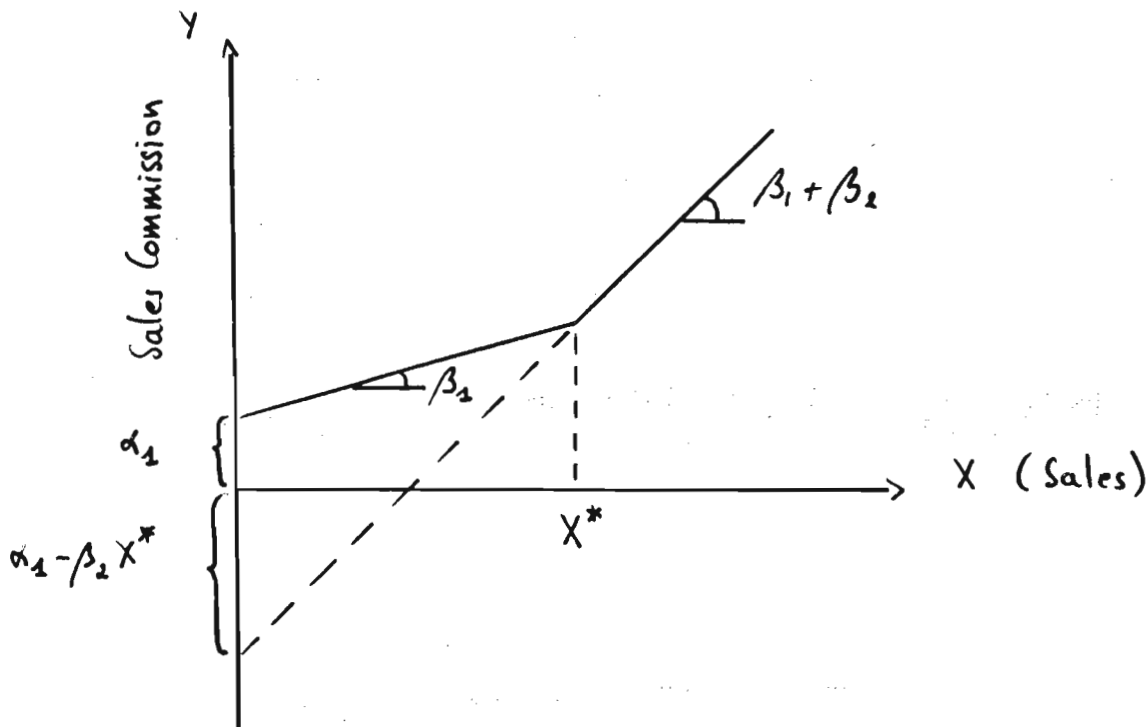


Fig. : Parameters of piecewise linear regression

Exemple : Relation les coûts totaux & niveau de production

Hyp. : coûts marg. de prod. \uparrow lorsque prod. > 5500 unités

$$\hat{Y}_i = -145,72 + 0,28 X_i + 0,09 (X_i - X^*) D_i$$

$$t = (-0,82) \quad (6,07) \quad (1,14)$$

où

$$\begin{cases} Y_i = \text{coûts totaux de l'entreprise } i \text{ (en \$)} \\ X_i = \text{niveau de production de l'entrep. } i \text{ (en unités)} \\ X^* = 5500 \text{ unités} \\ D_i = 1 \text{ si } X_i > X^* \\ \quad = 0 \text{ si } X_i < X^* \end{cases}$$

(1 unité = 1 kg)

I) Interprétation des variables dummy de régr. semi-logarithmique

Coeff. de régr. du modèle sem-log. \rightarrow semi-élasticités

(= mesure variation en % de la variable dép. suite à un chgmt unitaire de la variable explicative).

Interprétation s'applique uniq. aux variables expl. quant.

Quid si var. expl. qualitative ?

$$\ln Y_i = \beta_1 + \beta_2 D_i + u_i$$

où $\left\{ \begin{array}{l} Y_i = \text{salaire horaire de l'individu } i \text{ (en \$)} \\ D_i = 1 \text{ si individu } i \text{ est une } \text{♀} \\ \quad = 0 \text{ sinon} \\ u_i = \text{terme d'erreur} \end{array} \right.$

Si $E(u_i) = 0$:

Fctiou de salaire des ♂ : $E(\ln Y_i | D_i = 0) = \beta_1$

Fctiou de salaire des ♀ : $E(\ln Y_i | D_i = 1) = \beta_1 + \beta_2$

β_1 = moyenne du log. des salaires horaires de ♂

β_2 = différence tel moyenne du log. des salaires horaire ♂ et ♀

$\exp(\beta_1)$ = salaire horaire "médian" des ♂

$\exp(\beta_1 + \beta_2)$ = salaire horaire "médian" des ♀

Rem. : $\exp(\text{moyenne}(\ln)) = \text{médiane}$

$\exp(\ln(\text{moyenne})) = \text{moyenne}$