

I EXERCICE-1

- Une des hypothèses de la méthode MCO est l'homoscédasticité du terme d'erreur. La variance $V(\varepsilon_i) = E((\varepsilon_i - E(\varepsilon_i))^2)$ est constante pour tout i , soit $V(\varepsilon_i) = \sigma_\varepsilon^2$. Cette hypothèse de variance constante est l'hypothèse d'homoscédasticité ; on parle alors de série homoscédastique ; par opposition si le terme d'erreur n'a pas une variance constante, on parle de terme d'erreur hétéroscédastique.
- On donne l'écriture d'un modèle sous la forme : $Y = \hat{a}_0 + \hat{a}_1 X_1 + \hat{a}_2 X_2 + \hat{a}_3 X_3 + \varepsilon$, Y étant la variable expliquée et les X_i trois variables explicatives.
 - \hat{a}_0 est la valeur estimée de Y pour $X_i = 0$. \hat{a}_1 , \hat{a}_2 et \hat{a}_3 donnent respectivement une estimation de la variation de Y quand une des variables explicatives X_1 , X_2 ou X_3 augmente d'une unité, les autres restant constantes.

b. Rappel de cours :

ε étant le **terme d'erreur** regroupe en fait trois types d'erreur :

Une erreur de spécification :

On ne peut expliquer entièrement la variable Y par les variables explicatives retenues ; de nombreuses autres variables non prises en compte ont une influence sur Y ; elles ont été omises du fait de leur faible influence sur Y . Donc une part de l'erreur est due aux variables omises.

Une erreur de mesure :

les données ne représentent pas pleinement le phénomène,

Une erreur d'échantillonnage :

les observations varient d'un échantillon à l'autre engendrant les fluctuations de Y autour de sa moyenne,

$$3. F = \frac{SCE/k}{SCR/(n-k-1)} = \frac{R^2/k}{(1-R^2)/(n-k-1)}. \text{ Si les résidus sont normalement distribués, sous l'hypothèse}$$

nulle $H_0 : a_1 = a_2 = \dots = a_k = 0$, la variable aléatoire F suit la loi de Fisher avec pour degrés de liberté k et $n - k - 1$.

Le ratio F fournit un test d'hypothèse permettant de tester l'hypothèse nulle H_0 .

Règle de décision :

On choisit un seuil de signification, puis :

Si $F \leq F_{(k;n-k-1)}$, on ne rejette pas H_0 ; aucune variable n'est significative.

Si $F > F_{(k;n-k-1)}$, on rejette H_0 , on accepte l'hypothèse que les paramètres ne sont pas tous nuls et que R^2 est significativement différent de 0, il existe au moins une variable significative.

- Estimateur BLUE : sans biais et de variance minimale et linéaire. Exemple : modèle de régression simple $\hat{a}_1 = \frac{\sum_{i=1}^{i=n} x_i y_i}{\sum_{i=1}^{i=n} x_i^2}$; en posant

$$w_i = \frac{x_i}{\sum_{i=1}^{i=n} x_i^2}, \hat{a}_1 = \sum_{i=1}^{i=n} w_i Y_i.$$

- Théorème de Gauss-Markov. Si les hypothèses de la MCO sont vérifiées, les estimateurs \hat{a}_1 et \hat{a}_0 sont BLUE (Best Linear Unbiased Estimator).

II EXERCICE-2

- $Y_t = Y_0 (1+r)^t \iff \ln Y_t = \ln Y_0 (1+r)^t = \ln Y_0 + t \ln(1+r)$; si on pose : $Z_t = \ln Y_t$, $B = \ln Y_0$ et $A = \ln(1+r)$, on a : $Z_t = At + B$, soit une relation linéaire entre t et $\ln Y_t$.
- La calculatrice :

2-Var Stats	2-Var Stats
$\bar{x}=8.5$	$\uparrow n=16$
$\Sigma x=136$	$\bar{y}=2570.16875$
$\Sigma x^2=1496$	$\Sigma y=41122.7$
$Sx=4.760952286$	$\Sigma y^2=105792705$
$\sigma x=4.609772239$	$Sy=81.82372695$
$\downarrow n=16$	$\downarrow \sigma y=79.22548295$

a.

On obtient : $\bar{t} = 8.5$ et $\bar{Y} = 2570.17$.

L1	L2	L3	L1	L2	L3	LinReg(ax+b) L1, L2, L3	LinReg y=ax+b
1	2445.3	7.8062	1	2445.3	7.8062		y=ax+b
2	2455.9	7.816	2	2455.9	7.816		a=.0066753196
3	2480	7.8218	3	2480	7.8218		b=7.794511944
4	2494.4	7.8284	4	2494.4	7.8284		r²=.9977952349
5	2510.9	7.8365	5	2510.9	7.8365		r=.9988970092
6	2531.4	7.8414	6	2531.4	7.8414		
7	2543.8		7	2543.8			

b. L3 = ln(L2) L3(1)=7.801923093...

On effectue la regression linéaire de $\ln(Y_t)$ en t , ce qui donne : $\widehat{\ln(Y_t)} = \hat{a}_0 + \hat{a}_1 t$, avec : $\hat{a}_0 \simeq 7.7945$ et $\hat{a}_1 \simeq 0.0067$.

- c. \hat{a}_0 donne l'estimation de la valeur de $\widehat{\ln(Y_t)}$ pour $t = 0$, donc au dernier trimestre 1992, soit : $\widehat{Y}_0 \simeq e^{7.7945} \simeq 2427.22$.
- d. Pour interpréter \hat{a}_1 on dérive : $\widehat{\ln(Y_t)} = \hat{a}_0 + \hat{a}_1 t$ donne $\frac{Y'_t}{Y_t} = \hat{a}_1$ soit $\frac{dY_t}{dt} * \frac{1}{Y_t} = \hat{a}_1$ soit $\hat{a}_1 = \left(\frac{dY_t}{Y_t} \right) / dt$; c'est le quotient de la variation relative de Y sur la variation absolue de t ; si $dt = 1$, $\hat{a}_1 = \frac{dY_t}{Y_t}$ donne une estimation de la variation relative de Y_t , soit de Y_t à Y_{t+1} une variation estimée à une augmentation de 0.67%.

Source de variation	ddl	Somme des carrés	Moyenne des carrés (variances)
Régression	1	SCE = $\sum (\hat{Z}_t - \bar{Z})^2$	SCE/1
Résiduelle	n - 2	SCR = $\sum (\hat{Z}_t - Z_t)^2$	SCR/(n - 2)
Totale	n - 1	SCT = $\sum (Z_t - \bar{Z})^2$	

	Degré de liberté	Somme des carrés	Moyenne des carrés
Régression	1	0.01515	0.01515
Résidus	14	0.00003	0.00000
Total	15	0.01518	

La calculatrice donne les paramètres de $Z = \ln Y$:

```
2-Var Stats
n=16
x̄=7.851252161
ȳ=125.6200346
Σx²=986.289752
Σy²=.0318159707
↓σy=.0308056812
```

et on a : $SCT = \sum (Z_t - \bar{Z})^2 = nV(Z_t) = 16 * 0.030806^2 \simeq 0.01518$; par ailleurs :

$\hat{Z}_t - \bar{Z} = \hat{a}_1 (t - \bar{t})$ et $\sum (\hat{Z}_t - \bar{Z})^2 = \hat{a}_1^2 \sum (t - \bar{t})^2 = nV(t) \hat{a}_1^2$ soit $SCE = 16 * 0.00667^2 * 4.60977^2 \simeq 0.01513$ et par différence : $SCR = SCT - SCE \simeq 0.01518 - 0.01513 = 0.00005$

Les degrés de liberté :

Pour SCT , nous devons estimer la moyenne \bar{Y} , donc le degré de liberté est $n - 1$, soit ici 15 ; pour SCR , il faut estimer \hat{a}_0 et \hat{a}_1 , ce qui constitue une perte de deux degrés de liberté donc un degré de liberté de $n - 2$, soit 14 et enfin pour SCE , il suffit d'estimer \hat{a}_1 car $V(\hat{Y} - \bar{Y}) = \hat{a}_1^2 V(X)$, le degré de liberté est 1.

- f. L'erreur type, est donnée par : $S_\varepsilon^2 = \frac{\sum_{i=1}^{i=n} e_i^2}{n-2} \simeq \frac{0.00005}{14} = 3.57 \times 10^{-6}$ et $S_\varepsilon \simeq \sqrt{\frac{0.00005}{14}} \simeq 0.0019$ (Excel donne 0.0015).

- g. $I = [\hat{a}_1 - t_{\alpha/2; \nu} S_{\hat{a}_1} ; \hat{a}_1 + t_{\alpha/2; \nu} S_{\hat{a}_1}]$, avec : $S_{\hat{a}_1}^2 = \frac{S_\varepsilon^2}{\sum_{i=1}^{i=n} x_i^2} = \frac{S_\varepsilon^2}{nV(t)} \simeq \frac{3.57 \times 10^{-6}}{340} = 0.0000001015$ soit

$$S_{\hat{a}_1} \simeq \sqrt{\frac{3.57 \times 10^{-6}}{340}} = 1.024695 \times 10^{-4} \text{ soit } I = [\hat{a}_1 - t_{\alpha/2; \nu} S_{\hat{a}_1} ; \hat{a}_1 + t_{\alpha/2; \nu} S_{\hat{a}_1}]$$

soit $I = [0.006675 - t_{\alpha/2; \nu} S_{\hat{a}_1} ; 0.006675 + t_{\alpha/2; \nu} S_{\hat{a}_1}]$.

$t_{0.025; 14} = 2.145$; on obtient : $0.006675 - 2.145 * 1.024695 \times 10^{-4} = 0.006455$ et

$0.006675 + 2.145 * 1.024695 \times 10^{-4} = 0.006895$, soit $I = [0.006455 ; 0.006895]$

h. Le premier trimestre 97 correspond à $t = 17$; on a pour estimation ponctuelle : $\widehat{Ln}(Y_t) = 7.7945 + 0.0067t$, soit

$$\widehat{Ln}(Y_{17}) = 7.7945 + 0.0067 * 17 = 7.9084 \text{ et } \widehat{Y}_{17} \simeq e^{7.9084} = 2720.03, \text{ de plus } S_f^2 = S_\epsilon^2 \left(1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum_{i=1}^n x_i^2} \right)$$

soit $S_f^2 = 3.57 \times 10^{-6} \left(1 + \frac{1}{16} + \frac{(17-8.5)^2}{16 * 4.60977^2} \right) \simeq 4.55 \times 10^{-6}$, soit $S_f \simeq \sqrt{4.55 \times 10^{-6}} = 2.1331 \times 10^{-3}$.

$\frac{f}{S_f}$ suit une loi de Student avec un ddl de $(n - 2)$, soit 14, ce qui donne pour Z_t , l'intervalle : $7.9084 - 2.145 * 2.1331 \times 10^{-3} = 7.9038$ et $7.9084 + 2.145 * 2.1331 \times 10^{-3} = 7.91298$ soit pour Y_t : $e^{7.9038} = 2707.55$ et $e^{7.91298} = 2732.52$ soit

$$I = [2707.55 ; 2732.52]$$

III EXERCICE-3

Un chercheur a mené une étude sur le salaire horaire Y d'ouvriers, à partir d'un échantillon comportant 250 hommes et 280 femmes.

On considère la variable indicatrice définie par : $M = \begin{cases} 1 & \text{si l'ouvrier est un homme} \\ 0 & \text{si il s'agit d'une femme.} \end{cases}$

1. a. Il a obtenu pour la régression :

$$\begin{cases} \hat{a}_1 = 2.12 & S_{\hat{a}_1} = 0.36 \\ \hat{a}_0 = 12.52 & S_{\hat{a}_0} = 0.23 \end{cases} \text{ que l'on note : } \boxed{\hat{Y} = 2.12 M + 12.52} \text{ avec } S_\epsilon \simeq 4.2.$$

2. $M = 0$, donne une estimation du salaire horaire moyen des femmes, soit $\hat{a}_0 = 12.52$ et $M = 1$, donne une estimation du salaire horaire moyen des hommes, soit $\hat{a}_0 + \hat{a}_1 = 2.12 + 12.52 = 14.64$, \hat{a}_1 représentant donc la différence entre le salaire moyen des hommes et celui des femmes.

3. La différence de salaire dans les deux sous échantillons (hommes et femmes) est \hat{a}_1 . Il nous faut faire un test sur \hat{a}_1 ; $\begin{cases} H_0 : a_1 = 0 \\ H_1 : a_1 \neq 0 \end{cases}$
 . Sous l'hypothèse H_0 , $t_{\hat{a}_1} = \frac{\hat{a}_1}{S_{\hat{a}_1}}$, le ratio de Student, suit une distribution de Student avec $(n - 2)$ degrés de liberté, soit ici 528 ;

ici $\frac{\hat{a}_1}{S_{\hat{a}_1}} = \frac{2.12}{0.36} = 5.89$; il reste à comparer ce quotient avec la valeur lue dans la table de Student, soit ici 2, car $ddl > 20$. La valeur

du quotient $\frac{\hat{a}_1}{S_{\hat{a}_1}}$ est supérieure à 2, on en déduit donc que l'on rejette l'hypothèse H_0 et donc \hat{a}_1 suffisamment différent de zéro pour affirmer que a_1 est significativement différent de zéro. La différence de salaire dans les deux sous échantillons est différente de zéro au niveau de 5%.

4. L'autre chercheur : $\boxed{\hat{Y} = \hat{\beta}_1 M + \hat{\beta}_0}$. On a ici $\hat{\beta}_0$ qui est une estimation du salaire des hommes ($M = 0$, donc $\hat{\beta}_0 = 14.64$ et $\hat{\beta}_1 + \hat{\beta}_0$ une estimation du salaire des femmes donc $\hat{\beta}_1 = 12.52 - 14.64 = -2.12$).