

Première Année Master M.A.E.F. 2016 – 2017

Econométrie I

Examen final, janvier 2017

Examen de 3h00. Tout document ou calculatrice est interdit.

Exercice 1 (Sur 20 points) Rappel (Lemme de Slutsky): Si (X_n) et (Y_n) sont deux suites de v.a. sur $(\Omega, \mathcal{A}, \mathbb{P})$ telles que $X_n \xrightarrow[n \rightarrow \infty]{\mathcal{L}} X$ et $Y_n \xrightarrow[n \rightarrow +\infty]{\mathcal{P}} c$, où $c \in \mathbf{R}$, alors $X_n Y_n \xrightarrow[n \rightarrow \infty]{\mathcal{L}} cX$.

On dispose des données donnant l'espérance de vie E en France en fonction du sexe S et de l'année depuis 1900. Plus précisément, on notera $E_{1,i}$ (respectivement $E_{2,i}$) pour $i = 1, \dots, n$, l'espérance de vie moyenne des femmes (resp. des hommes) l'an $1900 + i$. On posera $F = (F_j)_{1 \leq j \leq 2n} = {}^t(E_{1,1}, \dots, E_{1,n}, E_{2,1}, \dots, E_{2,n})$.

1. Dans un premier temps, sans tenir compte du sexe, on veut savoir si l'espérance de vie dépend linéairement de l'année considérée.

(a) Ecrire le modèle vectoriel sous-jacent vérifié par F en détaillant la matrice X .

(b) Montrer que $({}^t X X)^{-1} = \frac{1}{n(n^2-1)} \begin{pmatrix} (n+1)(2n+1) & -3(n+1) \\ -3(n+1) & 6 \end{pmatrix}$ pour tout $n \in \mathbf{N}^*$ (on rappelle que

$\sum_{i=1}^k i^2 = \frac{1}{6} (2k+1)(k+1)k$). En déduire que l'estimateur par moindres carrés $\hat{\theta} = {}^t(\hat{\theta}_0, \hat{\theta}_1)$ des paramètres $\theta = {}^t(\theta_0, \theta_1)$ du modèle vaut:

$$\hat{\theta}_0 = \frac{1}{n(n-1)} \sum_{i=1}^n ((2n+1) - 3i)(E_{1,i} + E_{2,i}) \quad \text{et} \quad \hat{\theta}_1 = \frac{3}{n(n^2-1)} \sum_{i=1}^n (2i - (n+1))(E_{1,i} + E_{2,i}),$$

où θ_0 désigne l'intercept du modèle.

(c) Si on suppose que les erreurs pour ce modèle forment une suite de v.a.i.i.d. centrées de variance σ^2 , l'estimateur $\hat{\theta}$ est-il sans biais? Déterminer sa matrice de variance en fonction de n et de σ^2 . Montrer qu'il est asymptotiquement gaussien et en déduire en particulier que

$$\sqrt{\frac{n^3}{6}} (\hat{\theta}_1 - \theta_1) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2).$$

(d) Donner l'expression de $\hat{\sigma}^2$ estimateur sans biais de σ^2 en fonction des $E_{1,i}$ et $E_{2,i}$, de $\hat{\theta}_0$ et $\hat{\theta}_1$. Est-ce un estimateur convergent?

(e) Déterminer le comportement asymptotique de $\hat{T}_n = \sqrt{\frac{n^3}{6\hat{\sigma}^2}} \hat{\theta}_1$ dans le cas où $\theta_1 = 0$, puis dans le cas où $\theta_1 \neq 0$. Est-ce que \hat{T}_n est la statistique de Student d'un test? Si on prend les données jusqu'à l'an 2000, avec $\hat{\sigma} \simeq 5$, on obtient $\hat{\theta}_1 \simeq 0.37$. Peut-on dire que le modèle est légitime avec un risque de 5%?

2. Dans un deuxième temps, sans tenir compte de l'année, on veut savoir si l'espérance de vie dépend du sexe sous la forme

$$E_{k,i} = \alpha_k + \varepsilon_{k,i} \quad \text{pour } k = 1, 2 \text{ et } i = 1, \dots, n,$$

où les $(\varepsilon_{k,i})$ forment une suite de v.a.i.i.d. centrées de variance σ^2 .

(a) Ecrire le modèle sous forme vectorielle et en déduire les estimateurs par moindres carrés de α_1 et α_2 .

(b) Montrer que $\hat{\alpha}_1$ et $\hat{\alpha}_2$ vérifient des théorèmes de la limite centrale que l'on précisera.

(c) En déduire que:

$$\sqrt{\frac{n}{2\hat{\sigma}^2}} ((\hat{\alpha}_1 - \hat{\alpha}_2) - (\alpha_1 - \alpha_2)) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, 1),$$

où $\hat{\sigma}^2$ est l'estimateur non biaisé de σ^2 que l'on précisera en fonction des $E_{1,i}$, $E_{2,i}$, $\hat{\alpha}_1$ et $\hat{\alpha}_2$.

(d) Donner l'expression de la statistique de Student permettant de tester si le sexe est un facteur explicatif significatif de l'espérance de vie. En l'an 2000 on trouve que $\hat{\alpha}_1 \simeq 67$, $\hat{\alpha}_2 \simeq 60$ et $\hat{\sigma} \simeq 12$. Grâce au théorème précédent, que peut-on conclure avec un risque de 5%?

3. On considère désormais un modèle où l'espérance de vie dépend linéairement à la fois du sexe et de l'année considérée, soit:

$$E_{k,i} = \beta_k + i\gamma_k + \varepsilon_{k,i} \quad \text{pour } k = 1, 2 \text{ et } i = 1, \dots, n,$$

toujours avec les $(\varepsilon_{k,i})$ formant une suite de v.a.i.i.d. centrées de variance σ^2 .

- (a) A partir de ce qui précède donner l'expression des estimateurs par moindres carrés $\hat{\beta}_1, \hat{\beta}_2, \hat{\gamma}_1$ et $\hat{\gamma}_2$.
 (b) Démontrer que:

$$\sqrt{\frac{n^3}{24\hat{\sigma}^2}} ((\hat{\gamma}_1 - \hat{\gamma}_2) - (\gamma_1 - \gamma_2)) \xrightarrow[n \rightarrow \infty]{\mathcal{L}} \mathcal{N}(0, 1),$$

où $\hat{\sigma}^2$ est l'estimateur non biaisé de σ^2 que l'on précisera en fonction des $E_{1,i}, E_{2,i}, \hat{\beta}_1, \hat{\beta}_2, \hat{\gamma}_1$ et $\hat{\gamma}_2$.

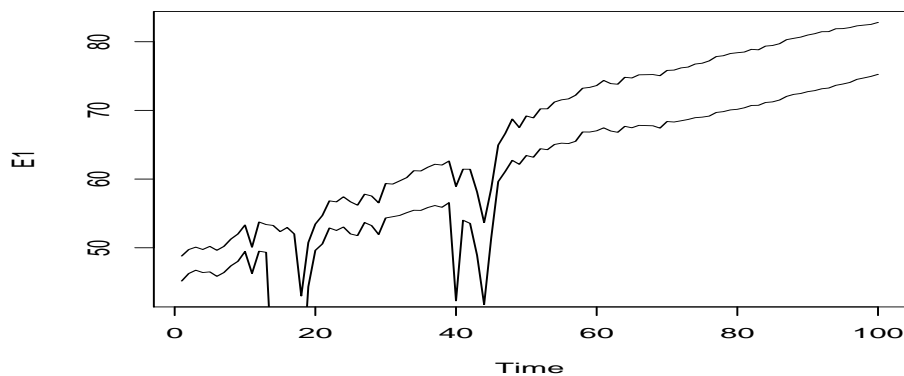
- (c) On trouve numériquement $\hat{\sigma} \simeq 4, \hat{\gamma}_1 \simeq 0.382$ et $\hat{\gamma}_2 \simeq 0.365$. Peut-on alors admettre que $\gamma_1 = \gamma_2$ avec un risque de 5%? Quel autre modèle peut-on proposer?

Exercice 2 (Sur 7 points)

On reprend les données et notations de l'Exercice 1 et on effectue l'étude évoquée précédemment avec le logiciel R. On commence à lire les données à partir d'un fichier (donnant les espérances de vie depuis 1816) puis on les représente graphiquement:

```
Esp=read.table('EsperanceVie.txt',header=TRUE)
E1=Esp$Fem[86:185]
E2=Esp$Ma[86:185]
I=1:100
ts.plot(E1)
lines(I,E2)
```

Soit le graphe:



On reprend le même ordre que les questions précédentes

1. On tape les commandes suivantes:

```
F=c(E1,E2)
J=c(I,I)
Reg1=lm(F~J)
summary(Reg1)
```

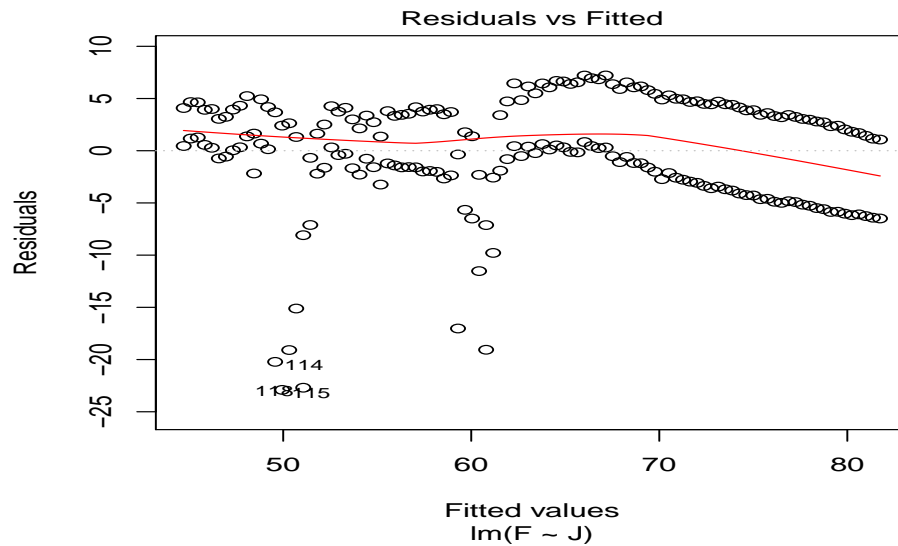
On obtient alors les résultats numériques et le graphe suivants:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	44.33430	0.76904	57.65	<2e-16 ***
J	0.37404	0.01322	28.29	<2e-16 ***

Residual standard error: 5.397 on 198 degrees of freedom
 Multiple R-squared: 0.8017, Adjusted R-squared: 0.8007
 F-statistic: 800.4 on 1 and 198 DF, p-value: < 2.2e-16

Quelle serait la prédiction d'espérance de vie d'une femme française en 2100? Quelles sont vos conclusions sur la régression?



2. On tape les commandes suivantes:

```
S1=c(rep(1,100),rep(0,100))
S2=c(rep(0,100),rep(1,100))
Reg2=lm(F~S1+S2-1)
summary(Reg2)
plot(c(1:100),Reg2$res[1:100])
```

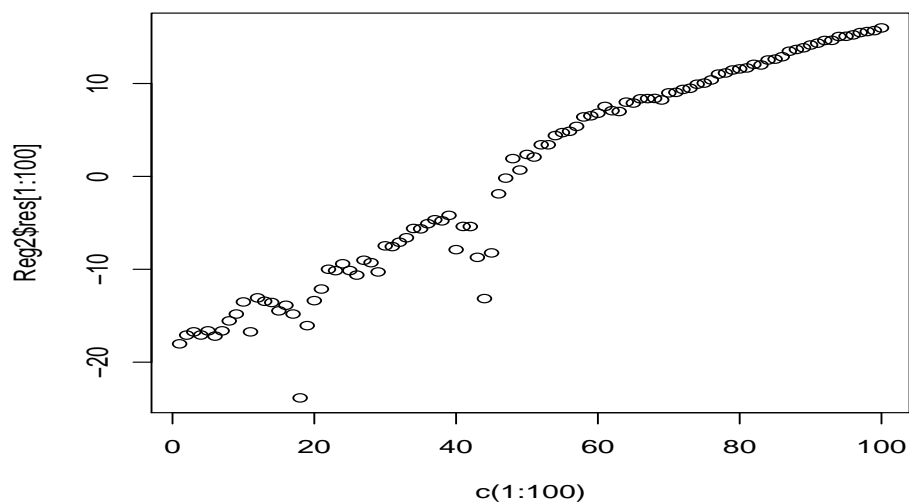
On obtient alors les résultats numériques et le graphe suivants:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
S1	66.823	1.157	57.77	<2e-16 ***
S2	59.624	1.157	51.55	<2e-16 ***

Residual standard error: 11.57 on 198 degrees of freedom
 Multiple R-squared: 0.968, Adjusted R-squared: 0.9677
 F-statistic: 2997 on 2 and 198 DF, p-value: < 2.2e-16

Quelle serait votre prédiction pour l'espérance de vie des femmes françaises en 2100? Quelles sont vos conclusions



quant à la régression?

3. On tape les commandes suivantes:

```
J1=c(I,rep(0,100))
J2=c(rep(0,100),I)
Reg3=lm(F~S1+S2+J1+J2-1)
summary(Reg3)
```

On obtient alors les résultats suivants:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
S1	47.49914	0.80961	58.67	<2e-16	***
S2	41.16945	0.80961	50.85	<2e-16	***
J1	0.38265	0.01392	27.49	<2e-16	***
J2	0.36543	0.01392	26.25	<2e-16	***

Residual standard error: 4.018 on 196 degrees of freedom
 Multiple R-squared: 0.9962, Adjusted R-squared: 0.9961
 F-statistic: 1.278e+04 on 4 and 196 DF, p-value: < 2.2e-16

Expliquer pourquoi les écarts-types empiriques sont les mêmes pour $S1$ et $S2$, ainsi que pour $J1$ et $J2$. Quelle serait votre prédiction pour l'espérance de vie des femmes françaises en 2100? Quelles sont vos conclusions par rapport à ce nouveau modèle?