

## Première Année Master M.A.E.F. 2021 – 2022

**Econométrie II**

Examen final, mai 2022

*Examen de 2h00. Tout document ou calculatrice est interdit.*

1. **(18 points)** Soit  $Y_i$ ,  $i = 1, \dots, N$  le chiffre d'affaire de  $N$  magasins de même taille appartenant au même groupe. On s'intéresse à l'impact sur ce chiffre d'affaire de 2 variables:  $Z^{(1)}$  qui vaut 1 si le magasin est en centre-ville et 0 sinon et  $Z^{(2)}$  qui vaut 1 si le magasin est ouvert le dimanche et 0 sinon.

(a) Lors d'une première étude, on choisit les magasins de telle manière que pour  $i = 1, \dots, n_1$ , les magasins sont tous au centre-ville et fermés le dimanche, et pour  $i = n_1 + 1, \dots, N$ , les magasins sont tous ouverts le dimanche sans qu'aucun ne soit au centre-ville. On pose le modèle:

$$Y_i = \theta_1 Z_i^{(1)} + \theta_2 Z_i^{(2)} + \varepsilon_i \quad \text{pour } i = 1, \dots, N,$$

où  $(\varepsilon_i)$  est une suite de variables aléatoires indépendantes gaussiennes centrées de même variance  $\sigma^2 > 0$  inconnue.

- i. Déterminer explicitement et non vectoriellement les estimateurs  $\hat{\theta}_1$  et  $\hat{\theta}_2$  de  $\theta_1$  et  $\theta_2$  par moindres carrés ordinaires et un estimateur  $\hat{\sigma}^2$  de  $\sigma^2$  **(2pts)**.
- ii. On veut tester  $H_0: \theta_1 = 0$ . Donner explicitement la statistique du test de Student permettant de tester cette hypothèse **(1pt)**. Quelle est sa loi **(0.5pts)**? Lorsque  $N$  est grand, quelle est asymptotiquement la loi de ce test **(0.5pts)**?
- iii. On veut tester si l'effet des 2 variables  $Z^{(1)}$  et  $Z^{(2)}$  est le même, c'est-à-dire si  $\theta_1 = \theta_2$  ( $H'_0$ ), contre l'hypothèse  $\theta_1 > \theta_2$  ( $H'_1$ ). Pour cela on considère la statistique:

$$\hat{T}_N = \frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{\frac{N}{n_1(N-n_1)} \hat{\sigma}^2}}.$$

Déterminer, en justifiant, la loi de  $\hat{T}_N$  **(1.5pts)**. Pour quel seuil franchi par  $\hat{T}_N$  allez-vous choisir  $H'_0$  plutôt que  $H'_1$  avec un risque de 5% **(0.5pts)**?

- iv. Si on suppose que pour  $i = 1, \dots, n_1$ , la variance de  $\varepsilon_i$  est  $\sigma_1^2$  et pour  $i = n_1 + 1, \dots, N$ , la variance de  $\varepsilon_i$  est  $\sigma_2^2$ , déterminer explicitement les estimateurs par moindres carrés généralisés de  $\theta_1$  et de  $\theta_2$ . Que remarquez-vous? **(2pts)** Quels sont les estimateurs de  $\sigma_1^2$  et  $\sigma_2^2$  par maximum de vraisemblance **(2pts)**? A partir de la loi de  $\hat{\theta}_1 - \hat{\theta}_2$ , en déduire une autre statistique permettant de tester  $H'_0$  contre  $H'_1$  **(1pt)**.
- (b) Pour chaque magasin  $i$  de l'étude décrite en (a), on considère  $X_i$  la variable telle que  $X_i = 1$  si le magasin a un chiffre d'affaire croissant et  $X_i = 0$  sinon. On veut expliquer  $X_i$  à l'aide de  $Z_i^{(1)}$  et  $Z_i^{(2)}$  en utilisant une régression logistique, avec la fonction logit. Présenter le modèle et donner explicitement les estimateurs des paramètres de cette régression **(2pts)**. Si un nouveau magasin numéroté  $N + 1$  est tel que  $Z_{N+1}^{(1)} = 0$   $Z_{N+1}^{(2)} = 1$ , comment prédire  $X_{N+1}$  **(2pts)**?
- (c) Désormais on a choisi  $N$  magasins pour lesquels il est possible d'avoir  $Z^{(1)} = p$  et  $Z^{(2)} = q$  pour tout  $(p, q) \in \{0, 1\}^2$ . Ainsi, pour  $(p, q) \in \{0, 1\}^2$ , on note maintenant  $Z^{(pq)}$  la variable telle que  $Z_i^{(pq)} = 1$  si  $Z_i^{(1)} = p$  et  $Z_i^{(2)} = q$ ,  $Z_i^{(pq)} = 0$  sinon. On suppose que:

- pour  $i = 1, \dots, n_{00}$ ,  $Z_i^{(00)} = 1$  et 0 sinon;
- pour  $i = n_{00} + 1, \dots, n_{00} + n_{01}$ ,  $Z_i^{(01)} = 1$  et 0 sinon;
- pour  $i = n_{00} + n_{01} + 1, \dots, n_{00} + n_{01} + n_{10}$ ,  $Z_i^{(10)} = 1$  et 0 sinon;
- pour  $i = n_{00} + n_{01} + n_{10} + 1, \dots, N = n_{00} + n_{01} + n_{10} + n_{11}$  alors  $Z_i^{(11)} = 1$  et 0 sinon.

On pose le modèle:

$$Y_i = \theta_{00} Z_i^{(00)} + \theta_{01} Z_i^{(01)} + \theta_{10} Z_i^{(10)} + \theta_{11} Z_i^{(11)} + \varepsilon_i, \quad \text{pour } i = 1, \dots, N = n_{00} + n_{01} + n_{10} + n_{11},$$

où  $(\varepsilon_i)$  est une suite de variables aléatoires indépendantes gaussiennes centrées de variance  $\sigma^2 > 0$ . Déterminer les estimateurs par moindres carrés ordinaires des  $\theta_{pq}$  **(1pt)**. Proposer un test permettant de tester si  $\theta_{00} = \theta_{01}$  **(1pt)**. Si on rejette cette hypothèse, que conclure sur l'influence mutuelle des variables  $Z^{(1)}$  et  $Z^{(2)}$  **(1pt)**?

*Proof.* (a) i. On trouve facilement que vectoriellement  $Y = Z\theta + \varepsilon$  et  ${}^tZZ = \begin{pmatrix} n_1 & 0 \\ 0 & N - n_1 \end{pmatrix}$ , d'où  $\hat{\theta}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} Y_i$  et  $\hat{\theta}_2 = \frac{1}{N-n_1} \sum_{i=n_1+1}^N Y_i$ .

Par ailleurs, l'estimateur non biaisé de la variance est  $\hat{\sigma}^2 = \frac{1}{N-2} \left( \sum_{i=1}^{n_1} (Y_i - \hat{\theta}_1)^2 + \sum_{i=n_1+1}^N (Y_i - \hat{\theta}_2)^2 \right)$ .

ii. D'après le cours, la statistique de test de Student pour ce test est, avec  $C = {}^t(1, 0)$ :

$$\hat{T} = \frac{\hat{\theta}_1}{\sqrt{\hat{\sigma}^2 {}^tC({}^tZZ)^{-1}C}} = \sqrt{n_1} \frac{\hat{\theta}_1}{\hat{\sigma}}.$$

D'après le cours, on sait que  $\hat{T} \stackrel{\mathcal{L}}{\sim} t(N-2)$  sous l'hypothèse  $H_0$ .

Quand  $N \rightarrow \infty$ , alors  $\hat{T} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$  sous l'hypothèse  $H_0$ .

iii. Une nouvelle fois, d'après le cours, la statistique de test de Student pour ce test est, avec  $C = {}^t(1, -1)$ :

$$\hat{T} = \frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{\hat{\sigma}^2 {}^tC({}^tZZ)^{-1}C}} = \frac{\hat{\theta}_1 - \hat{\theta}_2}{\hat{\sigma} \times \sqrt{(1, -1) \begin{pmatrix} 1/n_1 & 0 \\ 0 & 1/(N-n_1) \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix}}} = \frac{\hat{\theta}_1 - \hat{\theta}_2}{\hat{\sigma} \times \sqrt{\frac{N}{n_1(N-n_1)}}}.$$

Ainsi, on sait d'après le cours que sous  $H'_0$ ,  $\hat{T}_N \stackrel{\mathcal{L}}{\sim} t(N-2)$ .

De ceci, on déduit que si  $\hat{T}_N \leq q_{95}$  alors on choisit avec un risque de 5% plutôt  $H'_0$  que  $H'_1$ , avec  $q_{95}$  le quantile de niveau 0.95 de la loi  $t(N-2)$ .

iv. Soit  $\Sigma$  la matrice de covariance de  $\varepsilon$ . Alors  $\Sigma^{-1}$  est une matrice diagonale avec  $n_1$ -fois  $1/\sigma_1^2$  puis  $N - n_1$  fois  $1/\sigma_2^2$  sur la diagonale. Par suite  ${}^tZ\Sigma^{-1}Z = \begin{pmatrix} n_1/\sigma_1^2 & 0 \\ 0 & (N-n_1)/\sigma_2^2 \end{pmatrix}$ . Si on note  $\hat{\theta}_{MCG}$  l'estimateur de  $\theta = {}^t(\theta_1, \theta_2)$  par MCG, alors

$$\hat{\theta}_{MCG} = \begin{pmatrix} \sigma_1^2/n_1 & 0 \\ 0 & \sigma_2^2/(N-n_1) \end{pmatrix} {}^tZ\Sigma^{-1}Y = \hat{\theta}.$$

Par maximum de vraisemblance gaussien, ou plutôt de log-vraisemblance gaussienne

$$(\hat{\sigma}_1^2, \hat{\sigma}_2^2) = \text{Arg} \max_{\sigma_1^2, \sigma_2^2} \frac{1}{2} \left( -N \log(2\pi) - n_1 \log(\sigma_1^2) - (N-n_1) \log(\sigma_2^2) - \frac{1}{\sigma_1^2} \sum_{i=1}^{n_1} (Y_i - \hat{\theta}_1)^2 - \frac{1}{\sigma_2^2} \sum_{i=n_1+1}^N (Y_i - \hat{\theta}_2)^2 \right).$$

Les deux maximisations peuvent être séparées, et on trouve  $\hat{\sigma}_1^2 = \frac{1}{n_1} \sum_{i=1}^{n_1} (Y_i - \hat{\theta}_1)^2$  et  $\hat{\sigma}_2^2 = \frac{1}{N-n_1} \sum_{i=n_1+1}^N (Y_i - \hat{\theta}_2)^2$ .

On a  $\hat{\theta}_1 - \hat{\theta}_2 \stackrel{\mathcal{L}}{\sim} \mathcal{N}\left(0, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{N-n_1}\right)$  sous  $H'_0$ . On définira donc pour statistique de test  $\hat{T}'_N = \frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{\frac{\hat{\sigma}_1^2}{n_1} + \frac{\hat{\sigma}_2^2}{N-n_1}}}$  et

$\hat{T}'_N \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$  quand  $N \rightarrow \infty$  sous  $H'_0$ .

(b) Soit  $g(x) = e^x/(1+e^x)$ . Alors le modèle est tel que  $X_i \stackrel{\mathcal{L}}{\sim} \mathcal{B}\left(g(\theta_1 Z_i^{(1)} + \theta_2 Z_i^{(2)})\right)$ , les  $X_i$  étant indépendantes.

La log-vraisemblance du modèle est:

$$L_{(X_1, \dots, X_N)}(\theta_1, \theta_2) = \sum_{i=1}^{n_1} \left( X_i \log(g(\theta_1)) + (1 - X_i) \log(1 - g(\theta_1)) \right) + \sum_{i=n_1+1}^N \left( X_i \log(g(\theta_2)) + (1 - X_i) \log(1 - g(\theta_2)) \right).$$

On peut noter  $p_1 = g(\theta_1)$  et  $p_2 = g(\theta_2)$ ,  $p_1$  et  $p_2$  étant 2 réels de  $[0, 1]$ . Si on maximise  $L_{(X_1, \dots, X_N)}(p_1, p_2)$  en  $p_1$  et  $p_2$ , la maximisation peut être séparée et revient aux cas classiques du maximum de vraisemblance pour des variables de Bernoulli de même paramètre. On obtient ainsi:

$$\hat{p}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_i \quad \text{et} \quad \hat{p}_2 = \frac{1}{N-n_1} \sum_{i=n_1+1}^N X_i.$$

Pour le magasin  $N+1$ , il est clair qu'il est régi par le même comportement que les individus  $n_1+1, \dots, N$ . Donc si  $\hat{p}_2 \geq 1/2$ , alors on choisit  $\hat{X}_{N+1} = 1$ , et si  $\hat{p}_2 < 1/2$ , alors  $\hat{X}_{N+1} = 0$ .

(c) En reprenant la démarche du (a), la matrice  ${}^tZZ$  étant encore diagonale, on obtient:

$$\hat{\theta}_{00} = \frac{1}{n_{00}} \sum_{i=1}^{n_{00}} Y_i, \quad \hat{\theta}_{01} = \frac{1}{n_{01}} \sum_{i=n_{00}+1}^{n_{00}+n_{01}} Y_i, \quad \hat{\theta}_{10} = \frac{1}{n_{10}} \sum_{i=n_{00}+n_{01}+1}^{n_{00}+n_{01}+n_{10}} Y_i \quad \text{et} \quad \hat{\theta}_{11} = \frac{1}{n_{11}} \sum_{i=n_{00}+n_{01}+n_{10}+1}^N Y_i$$

Pour le test, on utilise une version adaptée de  $\hat{T}$  soit:

$$\hat{T}_1 = \frac{\hat{\theta}_{00} - \hat{\theta}_{01}}{\hat{\sigma} \times \sqrt{\frac{1}{n_{00}} + \frac{1}{n_{01}}}} \quad \text{où} \quad \hat{\sigma}^2 = \frac{1}{N-4} \left( \sum_{i=1}^{n_{00}} (Y_i - \hat{\theta}_{00})^2 + \sum_{i=n_{00}+1}^{n_{00}+n_{01}} (Y_i - \hat{\theta}_{01})^2 + \sum_{i=n_{00}+n_{01}+1}^{n_{00}+n_{01}+n_{10}} (Y_i - \hat{\theta}_{10})^2 + \sum_{i=n_{00}+n_{01}+n_{10}+1}^N (Y_i - \hat{\theta}_{11})^2 \right).$$

De ceci, on en déduit que  $\hat{T}_1 \stackrel{\mathcal{L}}{\sim} t(N-4)$  sous l'hypothèse  $\theta_{00} = \theta_{01}$ .

Ce test permet de savoir si  $Z^{(2)}$  influe sur  $Z^{(1)}$ , l'hypothèse nulle montrant l'absence d'interaction.

□

## 2. (10 points) Exercice de TP utilisant le logiciel R

(a) On a tapé les commandes suivantes:

```
t=c(1:100)
Z1=t
Z2=sqrt(t)*(2+sin(5*t))
Z3=1/t
epsilon=4*rnorm(100)
X=6+0.05*Z1-3*Z2+4*Z3+epsilon
reg=lm(X~Z1+Z2+Z3)
summary(reg)
```

Voici les résultats:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6.93171	1.30383	5.316	6.88e-07 ***
Z1	0.03387	0.02218	1.527	0.130
Z2	-2.98914	0.09127	-32.750	< 2e-16 ***
Z3	-1.30067	4.55615	-0.285	0.776

Residual standard error: 4.615 on 96 degrees of freedom  
 Multiple R-squared: 0.9492, Adjusted R-squared: 0.9476  
 F-statistic: 598 on 3 and 96 DF, p-value: < 2.2e-16

*Question 1: Qu'a-t-on fait et est-ce satisfaisant? Formellement, expliquer ce que sont les valeurs 0.03387 0.02218 1.527 0.130. Pouvait-on s'attendre à la valeur -2.98914? Comment expliquer la valeur 0.776? (3pts)*

(b) On tape ensuite les commandes:

```
M=cbind(1+0*t,Z1,Z2,Z3)
solve(t(M)%*%M)
```

Voici les résultats:

	Z1	Z2	Z3
	0.0798087578	-5.180381e-04	-2.637467e-03
Z1	-0.0005180381	2.310538e-05	-5.419362e-05
Z2	-0.0026374665	-5.419362e-05	3.910900e-04
Z3	-0.1630346962	1.432305e-03	3.009123e-03

*Question 2: Quel intérêt d'avoir ainsi effectué ces calculs? A quelle conclusion arrive-t-on? (2pts)*

(c) On a ensuite tapé les commandes:

```
Z4=t^2
Z5=runif(100,5,9)
Y=6+0.05*Z1-3*Z2+epsilon
Z=matrix(c(Z1,Z2,Z3,Z4,Z5),ncol=5);
colnames(Z)=c("Z1","Z2","Z3","Z4","Z5");
library(MASS)
ZZ=as.data.frame(Z);
y.lm=lm(Y~.,data=ZZ);
y.bic=stepAIC(y.lm,k=log(100))
```

Voici les résultats:

```
Start: AIC=292.72
Y ~ Z1 + Z2 + Z3 + Z4 + Z5
```

```
Step: AIC=288.13
Y ~ Z1 + Z2 + Z4 + Z5
```

```
Step: AIC=283.88
Y ~ Z1 + Z2 + Z5
```

```
Step: AIC=282.39
Y ~ Z1 + Z2
```

	Df	Sum of Sq	RSS	AIC
<none>			1466.9	282.39
- Z1	1	99.4	1566.3	284.34
- Z2	1	22283.2	23750.1	556.23

