

## Première Année Master M.A.E.F. 2020 – 2021

## Econométrie II

Contrôle continu n°2, avril 2021

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

## 1. Exercice 1 (Sur 19 points)

Soit  $(Y_i)_{1 \leq i \leq n}$  une famille de variables aléatoires définie par:

$$Y_i = \theta_0 + \sum_{k=1}^p \theta_k Z_i^{(k)} + \varepsilon_i \quad \text{pour tout } i \in \{1, \dots, n\}, \quad \text{où:} \quad (1)$$

- $\theta = {}^t(\theta_0, \theta_1, \dots, \theta_p)$  est un vecteur composé de  $p + 1$  réels inconnus.
  - pour  $1 \leq j \leq p$ , les  $(Z_i^{(j)})_{1 \leq i \leq n}$  sont  $p$  familles de réels connues. On note  $X = \begin{pmatrix} 1 & Z_1^{(1)} & \dots & Z_1^{(p)} \\ \vdots & \vdots & & \vdots \\ 1 & Z_n^{(1)} & \dots & Z_n^{(p)} \end{pmatrix}$  et on suppose que son rang est  $p + 1$  avec  $p + 1 \leq n + 1$ .
  - la suite  $(\varepsilon_i)_i$  est une suite de v.a.i.i.d. de loi gaussienne centrée de variance  $\sigma^2 > 0$ .
- (a) On note  $Y = (Y_i)_{1 \leq i \leq n}$  et  $\varepsilon = (\varepsilon_i)_{1 \leq i \leq n}$ . Ecrire le modèle (1) sous une forme matricielle, en précisant la loi du vecteur d'erreur  $\varepsilon$  (**0.5pts**).
- (b) Rappeler l'expression de l'estimateur  $\hat{\theta}$  de  $\theta$  par moindres carrés en fonction de  $X$  et  $Y$  (**0.5pts**). On note  $\hat{Y} = X \hat{\theta}$ . On mesure la qualité de la prédiction par cet estimateur avec le risque quadratique  $R(\hat{Y}) = \mathbb{E}(\|\hat{Y} - X \theta\|^2)$ , où  $\|\cdot\|$  désigne la norme euclidienne classique. Déterminer  $R(\hat{Y})$  en justifiant votre réponse (**1.5pts**).
- (c) A partir du modèle (1), on veut tester l'hypothèse  $H_0: \theta_i = 0$  pour tout  $i = p - p_0 + 1, \dots, p$ , où  $p_0 \in \mathbf{N}^*$ , contre l'hypothèse  $H_1$ , son complément. On note  $\hat{\sigma}^2 = \frac{1}{n - (p_0 + 1)} \|Y - \hat{Y}\|^2$ . Déterminer sous  $H_0$  la loi de  $\hat{\sigma}^2$  (**1.5pts**).
- (d) On note  $X^0$  la matrice extraite de  $X$  contenant uniquement ses  $p - p_0 + 1$  premières colonnes et  $\hat{Y}^0 = X^0 \hat{\theta}^0$ , où  $\hat{\theta}^0$  est obtenu par régression par moindres carrés sur les  $p - p_0$  premières variables. On définit:

$$\hat{F} = \frac{\frac{1}{p_0} \|\hat{Y} - \hat{Y}^0\|^2}{\hat{\sigma}^2}.$$

Montrer que sous  $H_0$ ,  $\|\hat{Y} - \hat{Y}^0\|^2 = \|P_A \varepsilon\|^2$  où  $A$  est un sous-espace vectoriel de  $\mathbf{R}^n$  de dimension  $p_0$  que l'on précisera et  $P_A$  est la matrice de la projection orthogonale sur  $A$  (**3.5pts**). En déduire la loi du numérateur de  $\hat{F}$  (**1pt**). Montrer que  $\hat{F}$  suit une loi de Fisher à  $(p_0, n - p - 1)$  degrés de liberté (**1.5pts**). Quelle règle de décision s'en déduit pour décider de  $H_0$  avec un risque de première espèce  $\alpha \in ]0, 1[$ ? (**1pt**)

- (e) On suppose jusqu'à la fin du problème que  $\theta_i = 0$  pour tout  $i = p - p_0 + 1, \dots, p$ . Déterminer alors  $R(\hat{Y})$  et  $R(\hat{Y}^0)$  (**1.5pts**). Quel estimateur vaut-il mieux choisir entre  $\hat{\theta}$  et  $\hat{\theta}^0$  (**0.5pts**)?
- (f) Pour estimer  $\sigma^2$ , on utilise les estimateurs par moindres carrés non biaisés  $\hat{\sigma}^2$  et  $\hat{\sigma}_0^2$  construits respectivement à partir de  $\hat{\theta}$  et  $\hat{\theta}^0$ . Déterminer en justifiant la loi de  $\hat{\sigma}_0^2$  (**0.5pts**). Montrer que pour  $Z$  une variable de loi  $\mathcal{N}(0, 1)$ ,  $\text{var}(Z^2) = 2$  (**1.5pts**). Déterminer alors les risques quadratiques de  $\hat{\sigma}^2$  et  $\hat{\sigma}_0^2$ , soit  $\mathbb{E}[(\hat{\sigma}^2 - \sigma^2)^2]$  et  $\mathbb{E}[(\hat{\sigma}_0^2 - \sigma^2)^2]$  (**1.5pt**). Quel estimateur de  $\sigma^2$  vaut-il mieux choisir entre les deux? (**0.5pts**)

- (g) On note  $\hat{R}^2$  et  $\hat{R}_0^2$  les coefficients de détermination  $R^2$  respectifs pour les modèles avec  $\hat{\theta}$  et avec  $\hat{\theta}^0$ . Montrer que  $\hat{R}^2 \geq \hat{R}_0^2$  (**1.5pts**). Par rapport à ce critère, quel estimateur choisiriez-vous? (**0.5pts**)

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

On considère un jeu de données réelles mettant en relation une mesure quantitative du degré d'avancement du diabète (variable `prog`) et 10 variables cliniques dans une cohorte de 442 patients. Chacun des 442 patients est décrit par les variables suivantes :

- `prog`: l'indice de progression de la maladie (le plus grand, le plus en progrès)
- `age`: l'âge du patient
- `sex`: le sexe du patient, codé numériquement (1, homme ou 2, femme)
- `bmi`: l'indice de masse corporelle
- `map`: la pression artérielle moyenne
- `ser1` - `ser6`: 6 mesures sérologiques

On a séparé au préalable et aléatoirement le jeu de données en deux sous-ensembles: `diabetes1` et `diabetes2`.

- (a) On cherche à construire un modèle linéaire avec de bonnes propriétés prédictives pour la variable `prog`. On a ainsi tapé les commandes suivantes avec le logiciel R:

```
Diab1=read.table("C:/Donnees/diabetes1.txt",header=TRUE)
Diab1$sex=as.factor(Diab1$sex)
attach(Diab1)
reg1=lm(prog~age+bmi+map+ser1+ser2+ser3+ser4+ser5+ser6+sex)
summary(reg1)
```

Voici le début de `Diab1`:

```
> Diab1
   prog age sex  bmi  map ser1 ser2 ser3 ser4 ser5 ser6
1   151  59  2  32.1 101.00 157  93.2 38.0 4.00 4.8598  87
3   141  72  2  30.5  93.00 156  93.6 41.0 4.00 4.6728  85
6    97  23  1  22.6  89.00 139  64.8 61.0 2.00 4.1897  68
9   110  60  2  32.1  83.00 179 119.4 42.0 4.00 4.4773  94
10  310  29  1  30.0  85.00 180  93.4 43.0 4.00 5.3845  88
:      :   :   :      :      :      :      :      :      :
:      :   :   :      :      :      :      :      :      :
:      :   :   :      :      :      :      :      :      :
```

Une partie des résultats obtenus est présente ci-dessous:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-319.8486	84.3700	-3.791	0.000202 ***
age	0.1468	0.3097	0.474	0.635981
bmi	6.2159	1.0607	5.860	2.02e-08 ***
map	1.1904	0.3172	3.753	0.000232 ***
ser1	-0.3994	0.7728	-0.517	0.605918
ser2	-0.2271	0.7468	-0.304	0.761353
ser3	0.2344	0.9817	0.239	0.811551
ser4	20.6894	8.1506	2.538	0.011943 *
ser5	41.5733	21.7245	1.914	0.057174 .
ser6	0.1199	0.3858	0.311	0.756264
sex2	-22.9688	8.2902	-2.771	0.006154 **

Residual standard error: 52.35 on 189 degrees of freedom

Multiple R-squared: 0.5669, Adjusted R-squared: 0.544

F-statistic: 24.74 on 10 and 189 DF, p-value: < 2.2e-16

Questions I.1: Pourquoi a-t-on tapé la deuxième commande? Montrer mathématiquement que dans le cadre de cette variable à deux modalités il n'était pas obligatoire de la taper. Combien y-a-t-il d'individus dans cette base de données? Ecrire précisément le modèle sous forme matricielle en précisant exactement la première ligne de la matrice  $X$ . En notant les paramètres  $\theta_j$ ,  $j = 0, \dots, p$ , dans l'ordre donné ci-dessus, que vaut numériquement  $\hat{\theta}_4$ ? Ecrire mathématiquement ce qu'est le nombre 0.3097. De ces résultats, pourriez-vous conclure statistiquement que plus on est en surpoids, plus la maladie progresse? Que pensez-vous de cette régression? (6pts)

(b) On tape ensuite les commandes:

```
library(leaps)
Z=matrix(c(age,bmi,map,ser1,ser2,ser3,ser4,ser5,ser6,sex),ncol=10);
colnames(Z)=c("age","bmi","map","ser1","ser2","ser3","ser4","ser5","ser6","sex");
r=leaps(Z,prog)
t=(r$Cp==min(r$Cp))
colnames(Z)[r$whi[t]]
reg2=lm(prog~bmi+map+ser2+ser4+ser5+sex)
summary(reg2); plot(reg2)
```

Voici les résultats obtenus:

```
> colnames(Z)[r$whi[t]]
[1] "bmi" "map" "ser2" "ser4" "ser5" "sex"
```

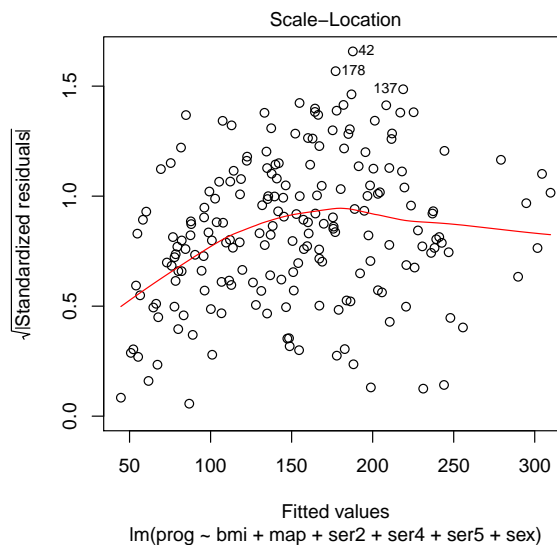
Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-291.5707	39.0335	-7.470	2.71e-12	***
bmi	6.4745	0.9778	6.621	3.44e-10	***
map	1.2302	0.3008	4.089	6.35e-05	***
ser2	-0.6142	0.1653	-3.716	0.000265	***
ser4	22.0946	4.6624	4.739	4.16e-06	***
ser5	31.5945	9.1086	3.469	0.000645	***
sex2	-21.9458	8.1080	-2.707	0.007405	**

Residual standard error: 51.92 on 193 degrees of freedom

Multiple R-squared: 0.5651, Adjusted R-squared: 0.5516

F-statistic: 41.8 on 6 and 193 DF, p-value: < 2.2e-16



Questions I.2: Qu'a-t-on fait en tapant ces commandes? Que concluez vous quant aux résultats et graphe obtenus? (1pt)

On tape maintenant les commandes suivantes:

```
library(MASS)
X=as.data.frame(Z)
reg3=stepAIC(lm(prog~.,data=X),k=log(200))
summary(reg3); plot(reg3)
```

Voici ce que l'on obtient à la fin des résultats:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-311.7870	36.4236	-8.560	3.19e-15	***
bmi	8.1540	0.9918	8.221	2.67e-14	***
ser5	53.7323	8.1519	6.591	3.90e-10	***

Residual standard error: 55.9 on 197 degrees of freedom  
 Multiple R-squared: 0.4854, Adjusted R-squared: 0.4801  
 F-statistic: 92.9 on 2 and 197 DF, p-value: < 2.2e-16

*Questions I.3: Qu'a-t-on fait en tapant ces commandes? Que concluez vous quant aux résultats obtenus? (1pt)*

On effectue maintenant ceci:

```
Diab2=read.table("C:/Donnees/diabetes2.txt",header=TRUE)
Diab2$sex=as.factor(Diab2$sex)
new=data.frame(age=Diab2$age,bmi=Diab2$bmi,map=Diab2$map,ser1=Diab2$ser1,ser2=Diab2$ser2,
ser3=Diab2$ser3,ser4=Diab2$ser4,ser5=Diab2$ser5,ser6=Diab2$ser6,sex=Diab2$sex)
pred1=predict(reg1,new); pred2=predict(reg2,new); pred3=predict(reg3,new)
MSE1=sqrt(mean((Diab2$prog-pred1)^2)); MSE2=sqrt(mean((Diab2$prog-pred2)^2));
MSE3=sqrt(mean((Diab2$prog-pred3)^2)); MSE1; MSE2; MSE3
```

Et on obtient:

```
> MSE1; MSE2; MSE3
[1] 56.98467
[1] 57.47944
[1] 57.73582
```

*Questions I.4: Qu'a-t-on fait en tapant ces commandes? Que concluez vous à partir des résultats obtenus? Est-ce surprenant? (2pts)*