

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

**Econométrie I**

Correction du Contrôle continu n°1, novembre 2016

*Examen de 1h30. Tout document ou calculatrice est interdit.***Exercice 1 (Sur 21 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  $\sigma^2 > 0$ .

1. 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**).
2. 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**).
3. A partir du modèle (1), on veut tester l'hypothèse  $H_0: \theta_i = 0$  pour tout  $i = p - p_0, \dots, p$ , où  $p_0 \in \mathbf{N}^*$ , contre l'hypothèse  $H_1$ , son complément. On note  $\hat{\sigma}^2 = \frac{1}{n - (p+1)} \|Y - \hat{Y}\|^2$ . Déterminer sous  $H_0$  la loi de  $\hat{\sigma}^2$  (**1.5pts**).
4. 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**)

5. On suppose jusqu'à la fin du problème que  $\theta_i = 0$  pour tout  $i = p - p_0, \dots, p$ . 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**)?
6. 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  $\mathbb{E}[(\hat{\sigma}_0^2 - \sigma^2)^2]$  (**1.5pts**). Quel estimateur de  $\sigma^2$  vaut-il mieux choisir entre les deux? (**0.5pts**)
7. Pour  $A$  et  $B$  deux sous-espaces vectoriels de  $\mathbf{R}^n$  tels que  $A \subset B$ , montrer que pour tout  $x \in \mathbf{R}^n$ ,  $\|P_A x\|^2 \leq \|P_B x\|^2$  (**2pts**). 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  $\hat{R}^2 \geq \hat{R}_0^2$  presque sûrement (**1.5pts**). Par rapport à ce critère, quel estimateur choisiriez-vous? (**0.5pts**)

*Proof.* 1. On a  $Y = X\theta + \varepsilon$ , où  $\varepsilon \stackrel{\mathcal{L}}{\sim} \mathcal{N}_n(0, \sigma^2 I_n)$ ,  $I_n$  étant la matrice identité.

2. On a  $\hat{\theta} = ({}^t X X)^{-1} {}^t X Y$ .  
 On a  $\hat{Y} = X \theta + P_{[X]} \varepsilon$ , où  $P_{[X]}$  désigne la matrice de la projection orthogonale sur  $[X]$ . Ainsi  $R(\hat{Y}) = \mathbb{E}(\|\hat{Y} - X \theta\|^2) = \mathbb{E}(\|P_{[X]} \varepsilon\|^2)$ .  
 En utilisant le Théorème de Cochran  $\|P_{[X]} \varepsilon\|^2 \stackrel{\mathcal{L}}{\sim} \sigma^2 \chi^2(\dim([X]))$  et on en déduit donc que  $R(\hat{Y}) = \sigma^2 \dim([X]) = \sigma^2 (p + 1)$ .
3. Après calculs déjà faits et le Théorème de Cochran,  $\widehat{\sigma}^2 = \frac{1}{n-(p+1)} \|P_{[X]^\perp} \varepsilon\|^2 \stackrel{\mathcal{L}}{\sim} \frac{\sigma^2}{n-(p+1)} \chi^2(\dim([X]^\perp)) \stackrel{\mathcal{L}}{\sim} \frac{\sigma^2}{n-(p+1)} \chi^2(n - (p + 1))$ .
4. Sous  $H_0$ , on a  $\hat{Y} = P_{[X]} X^0 \theta^0 + P_{[X]} \varepsilon = X^0 \theta^0 + P_{[X]} \varepsilon$  car  $[X^0] \subset [X]$ . De même,  $\hat{Y}^0 = P_{[X^0]} X^0 \theta^0 + P_{[X^0]} \varepsilon = X^0 \theta^0 + P_{[X^0]} \varepsilon$ . Ainsi  $\|\hat{Y} - \hat{Y}^0\|^2 = \|(P_{[X]} - P_{[X^0]}) \varepsilon\|^2$ . Or  $[X^0] \subset [X]$ , donc  $P_{[X^0]} \varepsilon = P_{[X^0]} P_{[X]} \varepsilon$  et par suite  $\|\hat{Y} - \hat{Y}^0\|^2 = \|(I_n - P_{[X^0]}) P_{[X]} \varepsilon\|^2 = \|P_{[X^0]^\perp} P_{[X]} \varepsilon\|^2 = \|P_A \varepsilon\|^2$  avec  $A = [X^0]^\perp \cap [X]$ . Or  $[X^0] \subset [X]$  donc  $[X]^\perp \subset [X^0]^\perp$ . D'où  $\dim(A) = \dim([X^0]^\perp) - \dim([X]^\perp) = p_0$ .  
 On en déduit ainsi que le numérateur a pour loi  $\frac{\sigma^2}{p_0} \chi^2(p_0)$  d'après le Théorème de Cochran.  
 On a vu à la question 3. que le dénominateur a pour loi  $\frac{\sigma^2}{n-(p+1)} \chi^2(n - (p + 1))$ . Il reste à montrer que numérateur et dénominateur sont indépendants. Mais le numérateur est une projection de  $\varepsilon$  sur  $[X^0]^\perp \cap [X]$  et le dénominateur sur  $[X]^\perp$ , donc deux sous-espaces vectoriels orthogonaux. Le Théorème de Cochran permet donc d'en déduire l'indépendance et la loi de Fisher pour  $\hat{F}$ . La règle de décision est la suivante: si  $\hat{F} \leq q_F(1 - \alpha)$ , où  $q_F(1 - \alpha)$  est le quantile d'ordre  $1 - \alpha$  de la loi de Fisher à  $(p_0, n - (p + 1))$  degrés de liberté, alors on choisit  $H_0$ , et si  $\hat{F} > q_F(1 - \alpha)$ , on rejette  $H_0$ .
5. Comme  $X \theta = X^0 \theta^0$ , le calcul de  $R(\hat{Y})$  est le même et on obtient  $R(\hat{Y}) = \sigma^2 (p + 1)$ . On peut reprendre le point 2. pour obtenir également que  $R(\hat{Y}^0) = \sigma^2 (p - p_0 + 1)$  car  $\dim([X^0]) = p - p_0 + 1$ .  
 Par rapport à ce critère, il vaut mieux utiliser  $\hat{\theta}^0$  que  $\hat{\theta}$ .
6. En utilisant le même raisonnement que 2., la loi de  $\widehat{\sigma}_0^2$  est  $\frac{\sigma^2}{n-(p-p_0+1)} \chi^2(n - (p - p_0 + 1))$ .  
 On a  $\text{var}(Z^2) = \int_{\mathbf{R}} x^4 f_Z(x) dx - (\mathbb{E}(Z^2))^2 = \frac{1}{\sqrt{2\pi}} \int_{\mathbf{R}} x^4 e^{-x^2/2} dx - 1$ . On utilise alors une intégration par partie en découpant  $x^4 = (-x^3) \times (-x)$  et en utilisant le fait qu'une primitive de  $-x e^{-x^2/2}$  est  $e^{-x^2/2}$ . D'où  $\frac{1}{\sqrt{2\pi}} \int_{\mathbf{R}} x^4 e^{-x^2/2} dx = \frac{1}{\sqrt{2\pi}} \left( [-x^3 e^{-x^2/2}]_{\mathbf{R}} + 3 \int_{\mathbf{R}} x^2 e^{-x^2/2} dx \right) = 3$ . Ainsi  $\text{var}(Z^2) = 2$ .  
 Comme les estimateurs sont non biaisés, il revient de calculer la variance d'un  $\chi^2$  à  $(n - (p + 1))$  et  $(n - (p - p_0 + 1))$  degrés de liberté. Avec le calcul précédent, comme un  $\chi^2(q)$  s'écrit comme la somme des carrés de  $q$  variables  $\mathcal{N}(0, 1)$  indépendantes, la variance d'un  $\chi^2(q)$  est donc  $2q$ . Par suite,  $\mathbb{E}[(\widehat{\sigma}^2 - \sigma^2)^2] = \frac{2\sigma^4}{(n-(p+1))}$  et  $\mathbb{E}[(\widehat{\sigma}_0^2 - \sigma^2)^2] = \frac{2\sigma^4}{(n-(p-p_0+1))}$ .  
 On préférera donc utiliser  $\widehat{\sigma}_0^2$  dont le risque quadratique est plus petit.
7. On a  $B = A \oplus (A^\perp \cap B)$ , les deux sous-espaces étant orthogonaux. Donc d'après Pythagore,  $\|P_B x\|^2 = \|P_A x\|^2 + \|P_{A^\perp \cap B} x\|^2 \geq \|P_A x\|^2$ .  
 On a  $\widehat{R}^2 = 1 - \frac{\|Y - \hat{Y}\|^2}{\|Y - \bar{Y}\|^2} = 1 - \frac{\|P_{[X]^\perp} \varepsilon\|^2}{\|Y - \bar{Y}\|^2}$ , alors que  $\widehat{R}_0^2 = 1 - \frac{\|P_{[X^0]^\perp} \varepsilon\|^2}{\|Y - \bar{Y}\|^2}$ . Comme  $[X]^\perp \subset [X^0]^\perp$ , d'après l'inégalité ci-dessus,  $\|P_{[X^0]^\perp} \varepsilon\|^2 \geq \|P_{[X]^\perp} \varepsilon\|^2$  presque sûrement et ainsi  $\widehat{R}_0^2 \leq \widehat{R}^2$  presque sûrement.  
 On en déduit donc que l'adéquation est meilleure avec des variables supplémentaires qui ne sont pas dans le vrai modèle et que suivant ce critère il vaut mieux choisir  $\hat{\theta}$  que  $\hat{\theta}^0$ .

□

## Exercice 2 (Sur 4 points)

On dispose de la consommation hebdomadaire de gaz et la température moyenne externe d'une maison test au sud de l'Angleterre pendant une saison. Une régression pour expliquer la consommation de gaz en fonction de la température est réalisée avec **R**. Les résultats numériques sont les suivants:

```

Residuals:
      Min       1Q   Median       3Q      Max
-0.97802 -0.11082  0.02672  0.25294  0.63803

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.72385     0.12974      ?    < 2e-16 ***
Temp        -0.27793         ?      -11.04 1.05e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3548 on 28 degrees of freedom
Multiple R-Squared:  0.8131,    Adjusted R-squared:  0.8064
F-statistic: 121.8 on 1 and 28 DF,  p-value: 1.046e-11

```

1. Donner le modèle et les hypothèses de la régression (**0.5pts**).
2. Donner une estimation de la variance du terme d'erreur dans le modèle de régression simple (**0.5pts**).
3. Compléter le tableau (**1pts**).

4. Soit  $A$  une variable aléatoire de loi de Student de degré de liberté 28. Quelle est la probabilité que  $A$  soit inférieure à 11.04 **(0.5pts)**?
5. Préciser les éléments du test correspondant à la ligne "Temp" du tableau ( $H_0$ ,  $H_1$ , la statistique de test, sa loi sous  $H_0$  et la règle de décision) **(0.5pt)**.
6. Rappeler la définition et interpréter le nombre "Multiple R-squared" du tableau **(0.5pts)**.
7. Pensez-vous que la température extérieure a un effet sur la consommation de gaz? Justifiez **(0.5pts)**.

*Proof.* 1. **(0.5pts)** On note  $y_i$  les consommations et  $z_i$  les températures, le modèle linéaire suppose:  $y_i = \beta z_i + \mu_i + \varepsilon_i$ ,  $i = 1, \dots, 30$  et les  $\varepsilon_i$  sont des variables aléatoires gaussiennes centrées et de même variance  $\sigma^2$ . On sait que  $n = 30$  car on voit que l'estimateur de la variance des résidus est de 28 degrés de libertés dans le tableau.

2. **(0.5pts)** Un estimateur de la variance  $\sigma^2$  du terme d'erreur est donné par le carré du terme 'Residual standard error' dans le tableau:  $\hat{\sigma}^2 = 0.3548^2 \approx 0.126$ .
3. **(1pts)** Dans le tableau il manque  $\hat{T}_\mu$  et  $\hat{\sigma}_\beta$ . On trouve d'après le tableau:

$$\hat{T}_\mu = \frac{\hat{\mu}}{\hat{\sigma}_\mu} \approx \frac{4.73}{0.13} = \frac{73}{13} \approx 36.3.$$

Puis comme on a  $\hat{T}_\beta = \frac{\hat{\beta}}{\hat{\sigma}_\beta} = -11.04$  d'après le tableau, et que  $\hat{\beta} = -0.27793$  toujours d'après le tableau, on en déduit que:

$$\hat{\sigma}_\beta = \frac{-0.27793}{-11.04} \approx \frac{0.28}{11.04} \approx \frac{28}{11.04} \times \frac{1}{100} \approx \frac{2.5}{100} = 0.025.$$

4. **(0.5pts)** D'après le tableau:  $P(|A| > 11.04) = 1.05 \times 10^{-11}$ .
5. **(0.5pt)** Pour la ligne "Temp", l'hypothèse est  $H_0: \beta = 0$  contre  $\beta \neq 0$ . Sout  $H_0$  la statistique de test  $\hat{T}_\beta = \frac{\hat{\beta}}{\hat{\sigma}_\beta}$  suit une loi de Student à  $n - 2 = 28$  ddl. On rejette  $H_0$  si la p-valeur est très petite (plus petite que 0.05 en pratique le plus souvent). Ici  $|\hat{T}_\beta| = 11.04$  d'après le tableau et la p-valeur:  $P(|A| > 11.01) = 1.05 \times 10^{-11}$  est très faible, hautement significative donc on rejette  $H_0$ .
6. **(0.5pts)** Le nombre "Multiple R-squared" 0.8131 donné par le tableau correspond au coefficient de détermination  $R^2$  du modèle. Il signifie qu'environ 81 % de la variation des données de consommation est expliquée par ce modèle de régression linéaire simple.
7. **(0.5pts)** Au vu du résultat de test de Student (ou du test de Fischer équivalent à il est clair que la température a un impact sur la consommation de gaz (+ il fait froid + on chauffe!).

□

### Exercice 3 (Sur 2 points)

Nous étudions le nombre d'enfants  $y_i$  de 10 ménages en fonction du nombre de fois  $z_i$  où le ménage a mangé dans au restaurant en un trimestre. Nous disposons de 10 couples de mesures  $(z_i, y_i)$  et nous savons que:  $\bar{z} = 15$ ,  $\bar{y} = 2.4$  et

$$\frac{1}{10} \sum_{i=1}^{10} (z_i - \bar{z})^2 = 34.6, \quad \frac{1}{10} \sum_{i=1}^{10} (y_i - \bar{y})^2 = 4.24, \quad \frac{1}{10} \sum_{i=1}^{10} (z_i - \bar{z})(y_i - \bar{y}) = -10.7$$

1. On note  $y = \hat{\mu} + \hat{\beta}z$  la droite de régression. Calculer à  $10^{-2}$  près,  $\hat{\mu}$  et  $\hat{\beta}$  **(1pt)**.
2. Donner la formule du coefficient de détermination et montrer que l'application numérique donne 78%. Commentez. **(1pt)**.

*Proof.* 1. **(1pt)**

$$\hat{\beta} = \frac{-10.7}{34.6} \approx -0.31$$

$$\hat{\mu} = \bar{y} - \hat{\beta}\bar{z} = 2.4 - (-0.31) \times 15 = 2.4 + 4.65 = 7.05$$

2. **(1pt)**

$$R^2 = \frac{(-10.7)^2}{34.6 \times 4.24} = \frac{114.49}{146.704} \approx \frac{114.5}{146.7} \approx 0.78$$

le modèle de régression linéaire simple explique 78% de la variance des données.

□