

## TD de Régression

### Analyse bivariable

**Ex 1.** *Lien entre 2 variables quantitatives : le coefficient de corrélation linéaire*

On considère deux échantillons de  $n$  variables  $(X_1, \dots, X_n)$  et  $(Y_1, \dots, Y_n)$ .

On rappelle les résultats suivants :

- la densité d'une loi de Student à  $k$  degrés de liberté  $T(k)$  vaut

$$f(x) = \frac{1}{\sqrt{k\pi}} \frac{\Gamma(\frac{k+1}{2})}{\Gamma(\frac{k}{2})} \left(1 + \frac{x^2}{k}\right)^{-\frac{k+1}{2}}, \quad x \in \mathbb{R}.$$

- La fonction bêta  $B(x, y)$  est définie pour tous  $x, y$  par  $B(x, y) = \int_0^1 t^{x-1}(1-t)^{y-1}dt$ .  
On a  $B(x, y) = B(y, x)$  et  $B(x, y+1) = \frac{y}{x+y}B(x, y)$ .

1) Rappeler la définition de la corrélation linéaire empirique  $R$  entre les deux échantillons précédents.

On suppose dans la suite de l'exercice que les couples  $(X_i, Y_i)$ ,  $i = 1, \dots, n$ , sont i.i.d suivant une loi normale, et on note  $\rho$  la corrélation théorique entre  $X_1$  et  $Y_1$ . On admet que dans ce cas, et sous l'hypothèse  $H_0 : \rho = 0$ ,  $\sqrt{n-2}R/\sqrt{1-R^2} \sim T(n-2)$ .

2) En déduire que sous  $H_0$  la densité de  $R$  est  $f(r) = B(\frac{1}{2}, \frac{n-2}{2})^{-1}(1-r^2)^{\frac{n-4}{2}}$ ,  $r \in [-1, 1]$ .

3) Montrer que sous  $H_0$ ,  $E(R) = 0$  et  $V(R) = 1/(n-1)$ .

4) En admettant que pour  $n$  suffisamment grand, la loi de  $R$  peut être approchée par une loi gaussienne, en déduire une région critique pour tester  $H_0 : \rho = 0$  au niveau  $\alpha \in ]0, 1[$ .

**Ex 2.** *Lien entre 2 variables quantitatives : le coefficient de corrélation de Spearman*

On considère deux échantillons de  $n$  variables  $(X_1, \dots, X_n)$  et  $(Y_1, \dots, Y_n)$ . On note  $(r_1, \dots, r_n)$  (resp.  $(s_1, \dots, s_n)$ ) les rangs des variables  $X_i$  (resp.  $Y_i$ ) dans chaque échantillon. On suppose qu'il n'y a pas d'ex-aequo, de telle sorte que les rangs vont de 1 à  $n$ . On rappelle que la corrélation de Spearman  $R_S$  entre les échantillons  $(X_1, \dots, X_n)$  et  $(Y_1, \dots, Y_n)$  correspond à la corrélation linéaire entre leurs rangs.

1) Donner la formule définissant  $R_S$ .

2) Montrer que la moyenne empirique de l'échantillon  $(r_1, \dots, r_n)$  vaut  $(n+1)/2$  et que sa variance empirique vaut  $(n^2-1)/12$ .

3) En déduire que  $R_S = \frac{n^{-1} \sum_{i=1}^n r_i s_i - (n+1)^2/4}{(n^2-1)/12}$ .

4) Soit  $d_i = r_i - s_i$ . Montrer que  $\sum_{i=1}^n r_i s_i = n(n+1)(2n+1)/6 - 1/2 \sum_{i=1}^n d_i^2$ .

5) En déduire que

$$R_S = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2-1)}.$$

**Ex 3.** *Lien variable quantitative/variable qualitative : comparaison de  $k$  moyennes*

On considère un échantillon de  $n$  variables indépendantes  $(X_1, \dots, X_n)$ . On suppose que cet échantillon est composé de  $k$  sous-populations (correspondant aux  $k$  classes d'un facteur).

Pour tout  $i = 1, \dots, k$ , on note  $n_i$  le nombre d'individus dans la sous-population  $i$ , et  $(X_1^i, \dots, X_{n_i}^i)$  les éléments de  $(X_1, \dots, X_n)$  associés à la population  $i$ . Ainsi  $n = \sum_{i=1}^k n_i$  et  $(X_1, \dots, X_n) = \cup_{i=1}^k (X_1^i, \dots, X_{n_i}^i)$ .

On suppose que pour  $i = 1, \dots, k$ , les variables  $(X_1^i, \dots, X_{n_i}^i)$  sont identiquement distribuées selon une loi normale d'espérance  $\mu_i$  et de variance  $\sigma^2$  inconnues (on suppose donc que la variance est la même quelle que soit la sous-population  $i$ ).

1) Pour  $k = 2$ , proposer une procédure pour tester  $H_0 : \mu_1 = \mu_2$  contre  $H_1 : \mu_1 \neq \mu_2$  au niveau  $\alpha \in ]0, 1[$ .

On suppose à présent  $k$  quelconque et on pourra utiliser le résultat suivant.

*Rappel :*

- Si  $Y_1 \sim \chi^2(p)$ ,  $Y_2 \sim \chi^2(q)$ , et  $Y_1$  et  $Y_2$  sont indépendantes, alors  $Y_1 + Y_2 \sim \chi^2(p+q)$
- Réciproquement, si  $Y = Y_1 + Y_2$  et  $Y \sim \chi^2(r)$ ,  $Y_1 \sim \chi^2(p)$ ,  $p < r$ , alors  $Y_2$  est indépendante de  $Y_1$  et  $Y_2 \sim \chi^2(r-p)$ .

2) Montrer que

$$S_T^2 = S_{inter}^2 + S_{intra}^2,$$

où  $S_T^2 = \frac{1}{n} \sum_{l=1}^n (X_l - \bar{X})^2$  représente la variance totale de l'échantillon ;

$S_{inter}^2 = \frac{1}{n} \sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2$  la variance interclasses (entre les sous-populations) ;

$S_{intra}^2 = \frac{1}{n} \sum_{i=1}^k n_i S_i^2$ , où  $S_i^2 = \frac{1}{n_i} \sum_{j=1}^{n_i} (X_j^i - \bar{X}_i)^2$ , la variance intraclasse (moyenne des variances dans chaque sous-population).

3) Quelle est la loi de  $n_i S_i^2 / \sigma^2$  ? En déduire que  $n S_{intra}^2 / \sigma^2 \sim \chi^2(n-k)$ .

4) Sous l'hypothèse  $H_0 : \mu_1 = \dots = \mu_k$ , quelle est la loi de  $n S_T^2 / \sigma^2$  ? En déduire la loi de  $n S_{inter}^2 / \sigma^2$  sous  $H_0$ .

5) En déduire la loi de  $\frac{S_{inter}^2 / (k-1)}{S_{intra}^2 / (n-k)}$  sous  $H_0$  et une région critique pour tester  $H_0$  au niveau  $\alpha \in ]0, 1[$ . Comparer avec la première question pour  $k = 2$ .

Pour quantifier le lien entre une variable quantitative  $X$  et une variable qualitative  $Q$  à  $k$  modalités, on calcule parfois  $\hat{\eta}^2 = \frac{S_{inter}^2}{S_T^2}$ , qui estime  $\eta^2 = V(E(X|Q))/V(X)$ .

6) Montrer que  $0 \leq \eta^2 \leq 1$  et que  $0 \leq \hat{\eta}^2 \leq 1$ . A quoi correspond les cas  $\eta^2 = 0$  et  $\eta^2 = 1$  ? Formuler l'hypothèse  $H_0$  ci-dessus en fonction de  $\eta^2$  et exprimer la région critique du test en fonction de  $\hat{\eta}^2$ .

**Ex 4.** *Lien entre 2 variables qualitatives*

On considère deux variables qualitatives  $Q_1$  et  $Q_2$  ayant respectivement  $I$  et  $J$  modalités. On relève sur un échantillon de  $n$  individus le nombre d'individus appartenant à chacune des modalités croisées de  $(Q_1, Q_2)$ , ce que l'on résume dans un tableau de contingence. Comment tester l'indépendance de  $Q_1$  et  $Q_2$  à l'aide de ce tableau au niveau  $\alpha \in ]0, 1[$  ?

## TD de Régression

### Ex 5. Moyenne empirique

Soit  $z_1, \dots, z_n$  des observations d'une variable  $Z$ .

- 1) Déterminer la valeur de  $\hat{m}$  qui minimise la distance quadratique  $S(m) = \sum_{i=1}^n (z_i - m)^2$ .
- 2) La quantité  $\hat{m}$  correspond à l'estimation par moindres carrés ordinaires dans un modèle de régression linéaire :  $Y = X\beta + \epsilon$ . Préciser ce que valent  $Y$ ,  $X$ ,  $\beta$  et  $\epsilon$ .
- 3) Retrouver le résultat de la première question à partir de la formule générale de l'estimateur des moindres carrés :  $\hat{\beta} = [X'X]^{-1}X'Y$ .

### Ex 6. Reconnaître un modèle de régression linéaire

Les modèles suivants sont-ils des modèles de régression linéaire ? Si non, peut-on appliquer une transformation pour s'y ramener ? Pour chaque modèle de régression linéaire du type  $Y = X\beta + \epsilon$ , on précisera ce que valent  $Y$ ,  $X$ ,  $\beta$  et  $\epsilon$ .

- 1) On observe  $(x_i, y_i), i = 1, \dots, n$  liés théoriquement par la relation  $y_i = a_0 + a_1x_i + \epsilon_i, i = 1, \dots, n$ . où les variables  $\epsilon_i$  sont centrées, de variance  $\sigma^2$  et non-corrélées. On désire estimer  $a_0$  et  $a_1$ .
- 2) On observe  $(x_i, y_i), i = 1, \dots, n$  liés théoriquement par la relation  $y_i = a_1x_i + a_2x_i^2 + \epsilon_i, i = 1, \dots, n$ . où les variables  $\epsilon_i$  sont centrées, de variance  $\sigma^2$  et non-corrélées. On désire estimer  $a_1$  et  $a_2$ .
- 3) On relève pour différents pays ( $i = 1, \dots, n$ ) leur production  $P_i$ , leur capital  $K_i$ , leur facteur travail  $T_i$  qui sont théoriquement liées par la relation de Cobb-Douglas  $P = \alpha_1 K^{\alpha_2} T^{\alpha_3}$ . On désire vérifier cette relation et estimer  $\alpha_1, \alpha_2$  et  $\alpha_3$ .
- 4) Le taux de produit actif  $y$  dans un médicament est supposé évoluer au cours du temps  $t$  selon la relation  $y = \beta_1 e^{-\beta_2 t}$ . On dispose des mesures de  $n$  taux  $y_i$  effectués à  $n$  instants  $t_i$ . On désire vérifier cette relation et estimer  $\beta_1$  et  $\beta_2$ .
- 5) Même problème que précédemment mais le modèle théorique entre les observations s'écrit  $y_i = \beta_1 e^{-\beta_2 t_i} + u_i, i = 1, \dots, n$ , où les variables  $u_i$  sont centrées, de variance  $\sigma^2$  et non-corrélées.

### Ex 7. Régression simple

On considère le modèle de régression linéaire simple où l'on observe  $n$  réalisations  $(x_i, y_i), i = 1, \dots, n$  d'un couple de variables aléatoires liées par la relation

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, i = 1, \dots, n.$$

On suppose les variables  $\epsilon_i$  centrées, de variance  $\sigma^2$  et non-corrélées.

- 1) Ecrire le modèle sous forme matricielle.

- 2) De quel problème de minimisation l'estimateur des moindres carrés  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$  est-il la solution ?
- 3) Calculer  $\hat{\beta}$  en résolvant directement le problème de minimisation précédent et vérifier la formule  $\hat{\beta} = [X'X]^{-1}X'y$ .
- 4) Calculer  $\text{var}(\hat{\beta})$  en utilisant la formule précédente et en déduire  $\text{var}(\hat{\beta}_0)$ ,  $\text{var}(\hat{\beta}_1)$  et  $\text{cov}(\hat{\beta}_0, \hat{\beta}_1)$ .
- 5) Montrer que la moyenne empirique des résidus  $\hat{\epsilon}_i$  est nulle.
- 6) Calculer la matrice de variance-covariance du vecteur de résidus  $\hat{\epsilon}$  et du vecteur des valeurs estimées  $\hat{y}$ .

**Ex 8.** *Régression simple avec bruit gaussien*

On se place dans le même modèle que pour l'exercice précédent, mais on suppose maintenant que les  $\epsilon_i$  sont iid de loi normale  $\mathcal{N}(0, \sigma^2)$  où  $\sigma^2$  est inconnue.

- 1) Quelle est la vraisemblance de l'échantillon  $(y_1, \dots, y_n)$  ?
- 2) Calculer les estimateurs du maximum de vraisemblance de  $\beta_0, \beta_1$  et  $\sigma^2$ . Comparer ces estimateurs avec les estimateurs des moindres carrés.
- 3) Dans cette question, on suppose pour simplifier que  $\beta_0$  et  $\sigma^2$  sont connus. Montrer que l'estimateur de  $\beta_1$  précédent est efficace.
- 4) Déterminer la loi des estimateurs  $\hat{\beta}_0, \hat{\beta}_1$ .
- 5) Justifier l'écriture :  $\hat{\sigma}^2 = \frac{1}{n} \|(I - X[X'X]^{-1}X')\epsilon\|^2$ . En déduire la loi de  $\hat{\sigma}^2$ .
- 6) En déduire des intervalles de confiance pour  $\beta_0$  et  $\beta_1$  de niveau  $\alpha \in ]0; 1[$ .
- 7) On observe un nouveau point  $x_{n+1}$  et on cherche à prédire la valeur  $y_{n+1} = \beta_0 + \beta_1 x_{n+1} + \epsilon_{n+1}$ . Donner la loi de l'erreur de prédiction  $e = y_{n+1} - \hat{y}_{n+1}$  basée sur le modèle de régression. En déduire un intervalle de confiance pour  $y_{n+1}$ .

**Ex 9.** *Convergence des estimateurs*

On sait que si  $X = \begin{pmatrix} 1, \dots, 1 \\ x_1, \dots, x_n \end{pmatrix}'$ , alors en notant  $\bar{x}_n = n^{-1} \sum_{i=1}^n x_i$ ,

$$(X'X)^{-1} = \frac{1}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n x_i^2 & -\bar{x}_n \\ -\bar{x}_n & 1 \end{pmatrix}$$

- 1) On a observé un échantillon i.i.d de couples  $(x_i, y_i)$ ,  $i = 1, \dots, n$ . On suppose le lien suivant :  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ , où les variables  $\epsilon_i$  sont i.i.d, centrées, de variance  $\sigma^2$ . Les régresseurs  $x_i$  sont supposés ici aléatoires et de carré intégrable. On note  $\epsilon = (\epsilon_1, \dots, \epsilon_n)'$ ,  $\beta = (\beta_0, \beta_1)'$  et  $\hat{\beta}$  l'estimateur de  $\beta$  par MCO. On suppose que  $X$  et  $\epsilon$  sont indépendants.
  - a) Exprimer  $\hat{\beta} - \beta$  en fonction des vecteurs  $X$  et  $\epsilon$ .
  - b) En déduire que  $\hat{\beta}$  converge presque sûrement vers  $\beta$  lorsque  $n \rightarrow \infty$ .
- 2) Lors d'une expérience chimique, on observe la teneur d'un certain produit à différents instants réguliers allant de 1 à  $n$ . Le résultat à l'instant  $i$  est noté  $y_i$ . On suppose le lien temporel suivant :  $y_i = \beta_0 + \beta_1 i + \epsilon_i$ ,  $i = 1, \dots, n$ , où les variables  $\epsilon_i$  représentent les erreurs de mesures. Elles sont supposées aléatoires, centrées, de variance  $\sigma^2$  et non corrélées. Soit  $\hat{\beta}$  l'estimateur de  $\beta$  par MCO.
  - a) Calculer  $\text{Var}(\hat{\beta})$  et donner sa limite lorsque  $n \rightarrow \infty$ .
  - b) En déduire le comportement asymptotique en moyenne quadratique de  $\hat{\beta}_0$  et  $\hat{\beta}_1$ .

3) On se place sous les mêmes hypothèses que la question précédente mais on suppose cette fois-ci que le lien temporel est :  $y_i = \beta_0 + \beta_1/i + \epsilon_i$ ,  $i = 1, \dots, n$ .

a) Calculer  $Var(\hat{\beta})$  et donner sa limite lorsque  $n \rightarrow \infty$ .

b) En déduire le comportement asymptotique en moyenne quadratique de  $\hat{\beta}_0$  et  $\hat{\beta}_1$ .

**Ex 10.**

Nous souhaitons exprimer la hauteur  $y$  d'un arbre en fonction de son diamètre  $x$  à 1m30 du sol. Pour cela, nous avons mesuré 20 couples diamètre-hauteur et les résultats ci-dessous sont disponibles :

$$\bar{x} = 34.9, \bar{y} = 18.34, \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})^2 = 28.29,$$

$$\frac{1}{20} \sum_{i=1}^{20} (y_i - \bar{y})^2 = 2.85, \frac{1}{20} \sum_{i=1}^{20} (x_i - \bar{x})(y_i - \bar{y}) = 6.26$$

1) On note  $\hat{y} = \hat{\beta}_0 + X\hat{\beta}$  l'estimation de la droite de régression. Donner l'expression de  $\hat{\beta}_0$  et  $\hat{\beta}_1$  en fonction des statistiques élémentaires ci-dessus. Calculer  $\hat{\beta}_0$  et  $\hat{\beta}_1$ .

2) Donner une mesure de qualité d'ajustement des données au modèle. Exprimer cette mesure à l'aide des statistiques élémentaires. Calculer et commenter.

3) Tester  $H_0 : \beta_j = 0$  contre  $H_1 : \beta_j \neq 0$  pour  $j = 0, 1$ . Commenter.

**Ex 11.** *Théorème de Gauss-Markov*

On considère un modèle de régression linéaire multiple

$$y = X\beta + \epsilon,$$

où  $\beta \in \mathbb{R}^k$ ,  $X$  est une matrice de taille  $n \times k$  et  $\epsilon$  est un vecteur aléatoire de taille  $n$ , centré et de matrice de covariance  $\sigma^2 I$  ( $I$  est la matrice identité). On veut montrer que l'estimateur des moindres carrés  $\hat{\beta}$  est l'estimateur linéaire sans biais de variance minimale, c'est-à-dire, pour tout estimateur linéaire sans biais  $\tilde{\beta}$ ,  $\text{var}(\tilde{\beta}) - \text{var}(\hat{\beta})$  est semi définie positive.

1) On pose  $\tilde{\beta} = My$  où  $M$  est une matrice de taille  $k \times n$ . Montrer que si  $\tilde{\beta}$  est sans biais, alors  $XM = I$ .

2) Montrer que  $\hat{\beta}$  et  $(\tilde{\beta} - \hat{\beta})$  sont non-corrélés.

3) Calculer la variance de  $\tilde{\beta}$  en écrivant  $\tilde{\beta} = \hat{\beta} + (\tilde{\beta} - \hat{\beta})$ . Conclure.

**Ex 12.** *Estimateur de la variance résiduelle*

On se place dans le même cadre que dans l'exercice précédent. On cherche à montrer que l'estimateur de  $\sigma^2$

$$\hat{\sigma}^2 = \frac{1}{n-k} \|y - X\hat{\beta}\|^2$$

est sans biais.

1) Montrer que  $(n-k)\hat{\sigma}^2 = \text{Tr}(\epsilon'(I - \Pi_X)\epsilon)$  où  $\text{Tr}$  désigne la trace et  $\Pi_X = X[X'X]^{-1}X'$  est la matrice de projection orthogonale sur l'espace engendré par les colonnes de  $X$ .

2) En utilisant la linéarité de la trace et le fait que  $\text{Tr}(AB) = \text{Tr}(BA)$ , montrer que

$$(n-k)\mathbb{E}(\hat{\sigma}^2) = \text{Tr}(\mathbb{E}((I - \Pi_X)\epsilon'\epsilon)).$$

3) Conclure.

**Ex 13.** *Le coefficient de corrélation multiple*

On considère le modèle de régression

$$y_i = \beta_0 + \beta_1 x_{1,i} + \cdots + \beta_p x_{p,i} + \epsilon_i, \quad i = 1, \dots, n, \quad (*)$$

où les variables  $\epsilon_i$  sont centrées, de variance  $\sigma^2$  et non-corrélées. On pose  $Y = (y_1, \dots, y_n)^T$ ,  $X_k = (x_{k,1}, \dots, x_{k,n})^T$  et  $\mathbf{1} = (1, \dots, 1)^T$ . On note  $\bar{y}$  la moyenne empirique de  $y$  et  $\hat{y} = \hat{\beta}_0 \mathbf{1} + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p$  où les estimateurs sont ceux obtenus par les moindres carrés ordinaires.

1) Que représente géométriquement  $\hat{y}$ ? Représenter sur un schéma les vecteurs  $y$ ,  $\hat{y}$ ,  $\bar{y}\mathbf{1}$ ,  $y - \bar{y}\mathbf{1}$ ,  $\hat{y} - \bar{y}\mathbf{1}$  et  $\hat{\epsilon}$ .

2) En déduire les égalités suivantes :

$$\begin{aligned} 1. \quad \sum_{i=1}^n y_i^2 &= \sum_{i=1}^n \hat{\epsilon}_i^2 + \sum_{i=1}^n \hat{y}_i^2 \\ 2. \quad \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n \hat{\epsilon}_i^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \end{aligned}$$

3) On considère les ratios :

$$R_1^2 = \frac{\sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2} \quad R_2^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Justifier (géométriquement) que  $R_1^2 \geq R_2^2$ . Dans quel cas a-t-on égalité?

4) Quelle est la définition du coefficient de corrélation multiple pour le modèle (\*)?

5) On considère à présent un modèle de régression sans constante, c'est à dire que l'on fixe  $\beta_0 = 0$  dans (\*). Les égalités montrées en 2) restent-elles valables? Quelle est dans ce cas la définition du coefficient de corrélation multiple?

6) Après estimation du modèle (\*) avec ou sans la constante, on obtient  $R^2 = 0.72$  avec la constante et  $R^2 = 0.96$  sans la constante. L'introduction de la constante est-elle pertinente?

**Ex 14.** *Le test de Fisher en pratique*

Dans un modèle de régression linéaire multiple comprenant 5 variables explicatives (dont éventuellement la constante), estimé sur  $n$  individus, on considère le test de Fisher de  $q \leq 5$  contraintes linéaires sur les coefficients :  $H_0 : R\beta = 0$  contre  $H_1 : R\beta \neq 0$ , où  $R$  est une matrice de taille  $(q, 5)$ .

1) Rappeler la statistique utilisée pour le test précédent.

2) Dans les cas suivants, donner l'expression de la matrice  $R$ , la loi suivie par la statistique de test sous  $H_0$  et la démarche pratique pour mettre en oeuvre le test :

i)  $H_0 : \beta_1 = 0$ ; ii)  $H_0 : \beta_1 = \beta_2 = \beta_4 = 0$ ; iii)  $H_0 : \beta_2 = \beta_3$ ; iv)  $H_0 : \beta_1 = \beta_2$  et  $\beta_2 = 2\beta_3$ .

**Ex 15.** *Le test de Fisher et le  $R^2$* 

On considère un modèle de régression linéaire multiple  $y = X\beta + \epsilon$  où  $\beta \in \mathbb{R}^p$ ,  $X$  est une matrice de taille  $n \times p$  et  $\epsilon$  est un vecteur aléatoire de taille  $n$ , centré et de matrice de covariance  $\sigma^2 I$  ( $I$  est la matrice identité).

On désire tester  $q$  contraintes linéaires sur le paramètre  $\beta$ , c'est à dire tester  $H_0 : R\beta = 0$  contre  $H_1 : R\beta \neq 0$ , où  $R$  est une matrice de taille  $(q, p)$ .

On note  $SCR$  la somme des carrés résiduelle du modèle initial, et  $SCR_c$  la somme des carrés résiduelle du modèle contraint (c'est à dire pour lequel l'hypothèse  $H_0$  est vérifiée).

- 1) Rappeler la statistique utilisée pour effectuer ce test. On la notera  $F$  et on donnera son expression en fonction de  $SCR$  et  $SCR_c$ .
- 2) Quelle loi suit cette statistique sous  $H_0$  lorsque  $\epsilon$  suit une loi normale ? Que peut-on dire de sa loi si aucune hypothèse de normalité n'est faite sur  $\epsilon$  ?
- 3) Montrer que si une constante est présente dans le modèle contraint,

$$F = \frac{R^2 - R_c^2}{1 - R^2} \frac{n - p}{q},$$

où  $R^2$  (respectivement  $R_c^2$ ) désigne le coefficient de détermination du modèle initial (respectivement du modèle contraint).

**Ex 16.** (issu du livre "Régression, Théorie et Applications")

Nous voulons expliquer la concentration de l'ozone sur Rennes en fonction des variables T9, T12, Ne9, Ne12 et Vx. Suite à l'estimation du modèle de régression linéaire, les sorties données par R sont (aux points d'interrogation près) :

Coefficients :

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	62	10	?	0
T9	-4	?	-5	0
T12	5	0.75	?	0
Ne9	-1.5	1	?	0.14
Ne12	-0.5	0.5	?	0.32
Vx	0.8	0.15	5.3	0

--

Multiple R-Squared: 0.6666, Adjusted R-squared: 0.6532

Residual standard error: 16 on 124 degrees of freedom

F-statistic: ? on ? and ? DF, p-value: 0

- 1) Retrouver les valeurs manquantes dans la sortie ci-dessus.
- 2) Rappeler la statistique de test et tester la nullité des paramètres séparément au seuil de 5%.
- 3) Rappeler la statistique de test et tester la nullité simultanée des paramètres autres que la constante au seuil de 5%.
- 4) Les variables Ne9 et Ne12 ne semblent pas influentes et nous souhaitons tester la nullité simultanée de  $\beta_{Ne9}$  et  $\beta_{Ne12}$ . Proposer un test permettant de tester ces contraintes et l'effectuer en vous aidant de la sortie R du modèle sans Ne9 et Ne12 suivante :

Coefficients :

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	66	11	6	0
T9	-5	1	-5	0
T12	6	0.75	8	0
Vx	1	0.2	5	0

--

Multiple R-Squared: 0.6525, Adjusted R-squared: 0.6442  
Residual standard error: 16.2 on 126 degrees of freedom

**Ex 17.** *Effet de la multicollinéarité*

On considère un modèle à deux variables explicatives. De l'estimation sur  $n$  individus, on a obtenu les matrices  $X'X$  et  $X'Y$  suivantes :

$$X'X = \begin{pmatrix} 200 & 150 \\ 150 & 113 \end{pmatrix} \quad X'Y = \begin{pmatrix} 350 \\ 263 \end{pmatrix}.$$

L'ajout d'une observation a modifié ces matrices de la façon suivante :

$$X'X = \begin{pmatrix} 199 & 149 \\ 149 & 112 \end{pmatrix} \quad X'Y = \begin{pmatrix} 347.5 \\ 261.5 \end{pmatrix}.$$

- 1) Calculer les coefficients estimés de la régression dans les deux cas.
- 2) Calculer le coefficient de corrélation linéaire entre les deux variables explicatives.
- 3) Commenter.

**Ex 18.** *Comparaison des critères de sélection d'un modèle*

On considère un modèle de régression linéaire visant à expliquer  $Y$  en fonction des variables  $X_1, \dots, X_p$ . On désire choisir entre le modèle avec  $X_p$  et le modèle sans  $X_p$  (les autres variables étant incluses dans les deux cas), sur la base d'un échantillon de  $n$  individus.

On note  $F$  la statistique :

$$F = (n - p) \frac{SCR_c - SCR}{SCR},$$

où  $SCR$  désigne la somme des carrés résiduelle dans le modèle avec  $X_p$ , et  $SCR_c$  désigne la somme des carrés résiduelle dans le modèle sans  $X_p$ .

1) En appliquant un test de Fisher de modèles emboîtés, selon quelle règle de décision, basée sur  $F$ , choisira-t-on d'inclure la variable  $X_p$  dans le modèle ?

2) On rappelle que le  $R^2$  ajusté dans un modèle à  $k$  variables et  $n$  individus est défini par

$$R_a^2 = 1 - \frac{n - 1}{n - k} \frac{SCR_k}{SCT},$$

où  $SCR_k$  désigne la somme des carrés résiduelles dans le modèle, et  $SCT$  la somme des carrés totaux.

Montrer qu'on décidera d'inclure  $X_p$  selon le critère du  $R^2$  ajusté si  $F > 1$ .

3) On rappelle que le  $C_p$  de Mallows dans un modèle à  $k$  variables et  $n$  individus est défini par

$$C_p = \frac{SCR_k}{\hat{\sigma}^2} - n + 2(k + 1),$$

où  $SCR_k$  désigne la somme des carrés résiduelles dans le modèle, et  $\hat{\sigma}^2$  est un estimateur de  $\sigma^2$  basé sur le plus gros modèle possible. On prendra ici  $\hat{\sigma}^2 = SCR/(n - p)$ , où  $SCR$  désigne la somme des carrés résiduelle dans le modèle avec  $X_p$ .

Montrer qu'on décidera d'inclure  $X_p$  selon le critère du  $C_p$  de Mallows si  $F > 2$ .



4) On rappelle que le critère  $AIC$  dans un modèle à  $k$  variables, à  $n$  individus, avec des résidus gaussiens, est défini par

$$AIC = n(1 + \log(2\pi)) + n \log \frac{SCR_k}{n} + 2(k + 1),$$

où  $SCR_k$  désigne la somme des carrés résiduelles dans le modèle.

Montrer qu'on décidera d'inclure  $X_p$  selon le critère  $AIC$  si  $F > (n - p)(e^{2/n} - 1)$ .

5) On rappelle que le critère  $BIC$  (aussi parfois appelé  $SBC$ ) dans un modèle à  $k$  variables, à  $n$  individus, avec des résidus gaussiens, est défini par

$$BIC = n(1 + \log(2\pi)) + n \log \frac{SCR_k}{n} + \log(n)(k + 1),$$

où  $SCR_k$  désigne la somme des carrés résiduelles dans le modèle.

Montrer qu'on décidera d'inclure  $X_p$  selon le critère  $BIC$  si  $F > (n - p)(e^{\log(n)/n} - 1)$ .

6) En admettant que les quantiles à 95% d'une loi de Fisher de degré de liberté  $(1, \nu)$  prennent leurs valeurs dans l'intervalle  $[3.8, 5]$  dès que  $\nu > 10$ , classer les critères précédents du plus conservatif (i.e. ayant tendance à refuser plus facilement l'introduction de  $X_p$ ) au moins conservatif (i.e. ayant tendance à accepter plus facilement l'introduction de  $X_p$ ). On pourra utiliser un développement limité pour l'étude des critères  $AIC$  et  $BIC$ , en supposant que  $n$  est suffisamment grand.

**Ex 19.** *Probabilité de sur-ajustement des critères de sélection*

On se place dans le cadre de l'exercice précédent, mais on suppose de plus que la variable  $X_p$  n'est pas significative dans le modèle (i.e. son coefficient est nul dans la régression) et que les résidus sont i.i.d. gaussiens. On admet les résultats énoncés dans les questions de l'exercice précédent.

1) Quelle loi suit la statistique  $F$ ? Montrer que lorsque  $n \rightarrow \infty$ , cette loi est équivalente à une loi  $\chi^2(1)$ .

2) Lors de la mise en oeuvre du test de Fisher des modèles emboîtés au niveau  $\alpha \in [0, 1]$ , quelle est la probabilité de décider (à tort) d'inclure la variable  $X_p$  dans le modèle?

3) Vers quoi tend la probabilité précédente si on base la décision sur le  $R_a^2$ ?

4) Même question si la décision est basée sur le  $C_p$  de Mallows.

5) Même question si la décision est basée sur le critère  $AIC$ .

6) Même question si la décision est basée sur le critère  $BIC$ .

7) Quel critère est-il préférable de choisir si l'on souhaite minimiser le risque d'inclure une variable en trop dans le modèle?

*Complément :* Dans la situation inverse où  $X_p$  est significative dans le modèle et qu'il est donc préférable de l'inclure, on peut montrer (mais c'est plus difficile) qu'en se fiant à n'importe lequel des critères ci-dessus, la probabilité de décider (à tort) ne pas inclure  $X_p$  tend vers 0 lorsque  $n \rightarrow \infty$ .

**Ex 20.** *Estimation dans un modèle ANOVA*

On considère une variable quantitative  $Y$  et un facteur  $A$  ayant  $I$  modalités. On suppose disposer d'un échantillon de  $n$  individus répartis en  $n_i$  individus dans chaque modalité  $A_i$  de  $A$ , pour  $i = 1, \dots, I$ . On note  $y_{i,k}$  la valeur de  $Y$  pour l'individu  $k$  appartenant à la

modalité  $A_i$  de  $A$ , et on note  $\bar{y}_i$  la moyenne de  $Y$  dans  $A_i$ .  
On considère le modèle ANOVA liant  $Y$  à  $A$  :

$$y_{i,k} = m + \alpha_i + \epsilon_{i,k}, \quad (1)$$

pour tout  $i = 1, \dots, I$ ,  $k = 1, \dots, n_i$ .

1) Le modèle précédent peut s'écrire sous la forme  $y = X\beta + \epsilon$  où  $y = (y_1, \dots, y_n)$  et  $\epsilon = (\epsilon_1, \dots, \epsilon_n)$ . On suppose que dans cette écriture les individus ont été rangés par modalités de  $A$ , i.e.  $y = (y_{1,1}, \dots, y_{1,n_1}, y_{2,1}, \dots, y_{2,n_2}, \dots, y_{I,1}, \dots, y_{I,n_I})'$ . Donner la forme de la matrice  $X$  et du vecteur  $\beta$ . Pourquoi l'estimation de  $\beta$  dans ce modèle n'est pas possible par les MCO ?

2) On considère la contrainte  $m = 0$ . Ecrire le modèle (1) avec cette contrainte sous la forme  $Y = X\beta + \epsilon$  en précisant la nouvelle matrice  $X$  et le nouveau vecteur  $\beta$ . Calculer l'estimateur  $\hat{\beta}$  issu des MCO.

3) On considère la contrainte  $\alpha_1 = 0$ . Ecrire le modèle (1) avec cette contrainte sous la forme  $Y = X\beta + \epsilon$  en précisant la nouvelle matrice  $X$  et le nouveau vecteur  $\beta$ . Calculer l'estimateur  $\hat{\beta}$  issu des MCO.

4) Montrer que quelle que soit la contrainte choisie précédemment,  $\hat{y}_{i,k} = \bar{y}_i$ , pour tout  $i = 1, \dots, I$  et  $k = 1, \dots, n_i$ .

5) On note  $\mu_i = E(Y|A_i)$  et on désire tester l'effet du facteur  $A$  sur  $Y$ , c'est à dire  $H_0 : \mu_1 = \dots = \mu_I$ . Comment se traduit cette hypothèse nulle sur les paramètres  $\beta$  du modèle (selon chaque contrainte précédemment choisie) ?

6) Montrer que la statistique de Fisher permettant d'effectuer le test précédent sur  $\beta$ , s'écrit (quelle que soit la contrainte sur  $\beta$  choisie initialement) :

$$F = \frac{S_A^2 / (I - 1)}{S_R^2 / (n - I)},$$

où  $S_A^2 = \sum_{i=1}^I n_i (\bar{y}_i - \bar{y})^2$  et  $S_R^2 = \sum_{i=1}^I \sum_{k=1}^{n_i} (y_{i,k} - \bar{y}_i)^2$ . Quelle est la région critique du test au niveau  $\alpha$  ? Sous quelle(s) hypothèse(s) ce test est-il valable ?

**Ex 21.** *Application de l'ANOVA à 1 facteur*

On veut étudier l'impact d'une ancienne mine d'arsenic sur les composantes hydrochimiques et hydrobiologiques d'un réseau hydrographique de Corse. Les mesures ont été faites sur 4 stations : B2, B3 (sur la Bravona) et P2 (sur un affluent la Presa) où est située la mine d'arsenic. Le tableau suivant résume la bioaccumulation de l'arsenic (en  $\mu g/g$ ) sur les branchies des truites capturées pour chaque station.

Station	effectif	moyenne	variance
P2	22	4.83	1.58
B2	21	0.66	0.07
B3	24	0.24	0.02

1) On désire tester l'effet des stations sur la bioaccumulation de l'arsenic. Quelle statistique peut-on utiliser pour mettre en oeuvre le test ? Les hypothèses d'application du test sont-elles réunies ?

2) On propose de transformer les données d'arsenic à l'aide de la fonction  $x \mapsto \sqrt{x}$ . Les résultats sont fournis dans le tableau suivant. Pourquoi cette transformation rend plus raisonnable un test d'effet des stations sur la teneur en arsenic ? Mettre en oeuvre le test et conclure.

Station	effectif	moyenne	variance
P2	22	2.18	0.08
B2	21	0.8	0.02
B3	24	0.48	0.01

**Ex 22.** *Estimation dans un modèle ANCOVA*

On observe chez  $n$  individus une variable quantitative  $Y$ , une autre  $X$  et un facteur  $A$  composé de  $I$  modalités. Pour  $i = 1, \dots, I$ , on note  $n_i$  le nombre d'individus dans la modalité  $A_i$  de  $A$ . Pour  $i = 1, \dots, I$ , et  $k = 1, \dots, n_i$ , on note  $y_{i,k}$  (resp.  $x_{i,k}$ ) la valeur de  $Y$  (resp. de  $X$ ) pour l'individu  $k$  appartenant à la modalité  $A_i$ .

On considère le modèle ANCOVA, pour  $i = 1, \dots, I$  et  $k = 1, \dots, n_i$ ,

$$y_{i,k} = \alpha_i + \beta_i x_{i,k} + \epsilon_{i,k}, \quad (2)$$

où les variables  $\epsilon_{i,k}$  sont supposées centrées, non-corrélées et de même variance  $\sigma^2$ .

1) Quelle contrainte a été implicitement imposée dans le modèle précédent pour le rendre identifiable ?

2) On définit les vecteurs

$$\begin{aligned} y &= (y_{1,1}, \dots, y_{1,n_1}, \dots, y_{I,1}, \dots, y_{I,n_I})', \\ \epsilon &= (\epsilon_{1,1}, \dots, \epsilon_{1,n_1}, \dots, \epsilon_{I,1}, \dots, \epsilon_{I,n_I})', \\ \gamma &= (\alpha_1, \beta_1, \dots, \alpha_I, \beta_I)'. \end{aligned}$$

Ecrire le modèle (2) sous la forme  $y = X\gamma + \epsilon$  en précisant la forme de la matrice  $X$ .

3) Montrer que les estimateurs par MCO des paramètres  $\alpha_i$  et  $\beta_i$  dans le modèle (2) (faisant intervenir toutes les observations) coïncident avec les estimateurs par MCO dans les  $I$  modèles estimés indépendamment dans chaque modalité  $A_i : y_{i,k} = \alpha_i + \beta_i x_{i,k} + \epsilon_{i,k}$ ,  $k = 1, \dots, n_i$ , où pour  $i$  fixé les variables  $\epsilon_{i,k}$  sont supposées centrées, non-corrélées et de même variance  $\sigma_i^2$  (mais aucune hypothèse n'est faite sur le lien entre  $\epsilon_{i,k}$  et  $\epsilon_{j,k}$  pour  $i \neq j$ ).

4) Quel est l'intérêt de considérer le modèle (2) plutôt que de considérer les  $I$  régressions séparément sur chaque modalité  $A_i$  ?

**Ex 23.** *Moindres carrés généralisés*

Soit un modèle de régression linéaire multiple

$$y = X\beta + \epsilon,$$

où  $\beta \in \mathbb{R}^k$ ,  $X$  est une matrice de taille  $n \times k$  et  $\epsilon$  est un vecteur aléatoire de taille  $n$ , centré. On considère ici la situation où les variables  $\epsilon_i$  ne sont plus homoscédastiques et non-corrélés, mais que  $\text{var}(\epsilon) = \Sigma$  où  $\Sigma$  est une matrice de rang  $n$ . On suppose dans cet exercice que  $\Sigma$  est connue (il conviendra dans la pratique de l'estimer).

1) Préciser la matrice  $\Sigma$  lorsque les variables  $\epsilon_i$  sont non-corrélés mais hétéroscédastiques de variance  $\sigma_i^2$ .

2) Déterminer l'espérance et la variance de l'estimateur  $\hat{\beta}$  des moindres carrés ordinaires.

3) On définit pour  $T \in \mathbb{R}^n$ ,  $\|T\|_{\Sigma}^2 = T'\Sigma^{-1}T$ . Donner la forme explicite de l'estimateur  $\hat{\beta}_G$  des moindres carrés généralisés défini comme le minimiseur de  $\|Y - X\beta\|_{\Sigma}$ . Calculer alors son espérance et sa variance.

- 4) En déduire que  $\hat{\beta}_G$  est plus efficace que  $\hat{\beta}$  (au sens du coût quadratique).
- 5) On suppose maintenant  $\epsilon \sim \mathcal{N}(0, \Sigma)$ . Montrer que  $\hat{\beta}_G$  est l'estimateur du maximum de vraisemblance.

**Ex 24.** *Régression biaisée*

Soit un modèle de régression linéaire multiple

$$y = X\beta + \epsilon,$$

où  $\beta \in \mathbb{R}^k$ ,  $X$  est une matrice de taille  $n \times k$  et où les variables  $\epsilon_i$ ,  $i = 1, \dots, p$ , sont centrées, homoscédastiques et non-corrélées.

Pour  $\kappa \geq 0$ , on considère l'estimateur  $\hat{\beta}_R = (X'X + \kappa I_p)^{-1}X'Y$  où  $I_p$  désigne la matrice identité de taille  $p$ .

- 1) Comment s'appelle l'estimateur  $\hat{\beta}_R$ ? Dans quel cas est-il égal à l'estimateur des MCO?

On suppose dans la suite que les variables explicatives ont été centrées, réduites et qu'elles sont non-corrélées entre elles.

- 2) Que vaut la matrice  $X'X$  sous les hypothèses précédentes?
- 3) Exprimer  $\hat{\beta}_R$  en fonction de  $\hat{\beta}$ .
- 4) En déduire l'espérance et la matrice de variance de  $\hat{\beta}_R$ .
- 5) Soit  $\hat{\beta}_{R,i}$ ,  $\hat{\beta}_i$  et  $\beta_i$  la  $i$ -ème composante des vecteurs  $\hat{\beta}_R$ ,  $\hat{\beta}$  et  $\beta$  ( $i = 1, \dots, p$ ). Montrer que si  $\kappa \leq 2\sigma^2/\beta_i^2$ , alors  $EQM(\hat{\beta}_{R,i}) \leq EQM(\hat{\beta}_i)$ .

- 6) En déduire que si  $\kappa \leq 2\sigma^2/||\beta||^2$ , alors  $\hat{\beta}_R$  est meilleur que  $\hat{\beta}$  au sens du coût quadratique.

- 7) La condition précédente est-elle utilisable en pratique? Quelle procédure peut-on mettre en place en pratique pour choisir  $\kappa$ ?