
Sujet 14: Test sur le coefficient directeur d'une droite de régression

• **Quelques rappels sur la définition mathématique des tests d'hypothèses:**

Soit X une variable aléatoire à valeurs dans un ensemble \mathcal{X} , de loi P_θ où $\theta \in \Theta$. Supposons que l'on veuille effectuer un test statistique d'hypothèses portant sur θ . On formule deux hypothèses contradictoires notées H_0 et H_1 dont on suppose que l'une et seulement l'une est vraie. Mathématiquement, formuler H_0 et H_1 revient à choisir deux sous-ensembles disjoints de Θ notés Θ_0 et de Θ_1 de sorte que l'hypothèse H_0 s'écrit alors $\{\theta \in \Theta_0\}$ tandis que l'hypothèse H_1 s'écrit $\{\theta \in \Theta_1\}$. Soit (X_1, \dots, X_n) un échantillon i.i.d. issu d'une variable parente X et soit $(x_1, \dots, x_n) \in \mathcal{X}^n$ une réalisation de (X_1, \dots, X_n) . Construire un test de H_0 contre H_1 revient à construire une région critique \mathcal{R} de telle sorte que l'on rejette H_0 lorsque $(x_1, \dots, x_n) \in \mathcal{R}$ et que l'on s'assure que le risque de 1ère espèce est inférieur ou égal à α .

Dans ce cadre, on parle de **fonction de risque de 1ère espèce** définie sur Θ_0 pour

$$\theta \rightarrow \alpha(\theta) = \mathbb{P}_\theta((X_1, \dots, X_n) \in \mathcal{R}).$$

On appelle **taille du test** la probabilité maximale de rejeter H_0 alors que H_0 est vraie:

$$\sup_{\theta \in \Theta_0} \{\mathbb{P}_\theta((X_1, \dots, X_n) \in \mathcal{R})\}.$$

On dit qu'un test est de **niveau** α si sa taille est égale à α ie si

$$\sup_{\theta \in \Theta_0} \{\mathbb{P}_\theta((X_1, \dots, X_n) \in \mathcal{R})\} = \alpha.$$

On parle de **fonction de risque de 2nde espèce** définie sur Θ_1 pour

$$\theta \rightarrow \beta(\theta) = \mathbb{P}_\theta((X_1, \dots, X_n) \notin \mathcal{R}).$$

On parle de **fonction puissance** définie sur Θ_1 pour

$$\theta \rightarrow 1 - \beta(\theta) = \mathbb{P}_\theta((X_1, \dots, X_n) \in \mathcal{R}).$$

• **Introduction à la régression:**

Le but de tout modèle de régression (d'une manière générale) est de spécifier une relation (de nature stochastique) entre une variable Y appelée variable à expliquer (ou variable de réponse ou variable endogène ou variable dépendante) et un certain nombre de variables explicatives $X^{(1)}, \dots, X^{(p)}$ (ou variables de contrôle ou variables exogènes ou régresseurs ou covariables). Ici, on supposera que les covariables sont de loi continue.

Notons $(Y_i, X_i^{(1)}, \dots, X_i^{(p)})$ pour $i = 1, \dots, n$ les observations recueillies sur n individus que l'on suppose distribuées comme $(Y, X^{(1)}, \dots, X^{(p)})$. Le modèle de régression linéaire gaussien standard à n observations indépendantes s'écrit sous la forme générique suivante:

$$Y_i = \beta_0 + \sum_{k=1}^p \beta_k X_i^{(k)} + \varepsilon_i, \quad \text{où } \varepsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2). \quad (1)$$

Le modèle de régression linéaire gaussien standard repose donc sur les hypothèses suivantes:

(M₁) **normalité**: conditionnellement à $X_i^{(1)}, \dots, X_i^{(p)}$, la variable de réponse Y_i suit la loi normale.

(M₂) **homoscédasticité**: les lois conditionnelles de Y_i sachant $X_i^{(1)}, \dots, X_i^{(p)}$ ont même variance donc pour $i = 1, \dots, n$, on a $\sigma_i^2 = \sigma^2$.

(M₃) **linéarité**: le prédicteur linéaire $\eta_i = \beta_0 + \sum_{k=1}^p \beta_k X_i^{(k)}$ est une combinaison affine des covariables, la linéarité s'apprécie relativement à la linéarité de η_i en les paramètres. Attention, la linéarité du modèle ne doit pas prêter à confusion: un modèle dit linéaire est linéaire en les paramètres $\beta_0, \beta_1, \dots, \beta_p$. Ainsi, le modèle de régression linéaire simple

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

est, comme son nom l'indique, linéaire, tout comme l'est le modèle de régression linéaire multiple ($p > 1$):

$$Y = \beta_0 + \sum_{k=1}^p \beta_k X^{(k)} + \varepsilon.$$

Le modèle de régression polynomiale

$$Y = \beta_0 + \beta_1 X^{(1)} + \beta_2 (X^{(1)})^2 + \varepsilon$$

est un modèle linéaire bien que la relation entre Y et $X^{(1)}$ soit quadratique. En revanche, un modèle de la forme

$$Y = \beta_0 + \beta_1 \exp(\beta_2 X^{(1)} + \beta_3 X^{(2)}) + \varepsilon_i$$

ou bien de la forme

$$Y = \frac{\exp(\beta_0 + \beta_1 X)}{1 + \exp(\beta_0 + \beta_1 X)} + \varepsilon$$

n'est pas linéaire!

(M₄) **centrage des erreurs**: $\mathbb{E}[\varepsilon_i] = 0$ pour $i = 1, \dots, n$.

Remarquons que cela équivaut à l'existence d'une relation "identité" entre le prédicteur linéaire η_i d'une part et l'espérance conditionnelle de la variable réponse d'autre part:

$$\mu_i := \mathbb{E}[Y_i | X_i^{(1)}, \dots, X_i^{(p)}] = \eta_i := \beta_0 + \sum_{k=1}^p \beta_k X_i^{(k)}.$$

(M₅) **non-corrélation**: les vecteurs $(Y_i, X_i^{(1)}, \dots, X_i^{(p)})$ sont non-corrélés pour $i = 1, \dots, n$. Si la normalité est effectivement respectée, cela revient alors à dire que les vecteurs $(Y_i, X_i^{(1)}, \dots, X_i^{(p)})$ sont mutuellement indépendants.

(M₆) **exogénéité = non-endogénéité**: $(X_i^{(1)}, \dots, X_i^{(p)})$ est indépendant de ε_i , pour $i = 1, \dots, n$.

(M₇) **non-colinéarité des covariables**: les covariables sont non-corrélées entre elles, ce qui se traduit en pratique par le fait que les vecteurs $(X_1^{(k)}, \dots, X_n^{(k)})$ sont non-colinéaires pour $k = 1, \dots, p$. On écartera également le cas trivial où l'un de ces vecteurs serait colinéaire au vecteur de longueur n dont toutes les composantes sont égales à 1.

• **Test de Student sur l'effet des covariables:**

Dans le cadre du modèle linéaire gaussien standard présenté ci-dessus, on se demande si les covariables introduites dans le modèle ont un effet statistiquement significatif ou non. Autrement dit, on souhaite tester au niveau de risque α l'hypothèse nulle $H_0: \beta_j = 0$ contre l'hypothèse alternative $H_1: \beta_j \neq 0$, ceci pour un $j \in \{1, \dots, p\}$ fixé quelconque. Pour cela, on s'intéressera ici au test de Student qui est fondé sur l'estimateur des moindres carrés ordinaires (qui coïncide ici avec l'EMV) de β_j que l'on notera $\hat{\beta}_j$. Notons $\widehat{\text{Var}}(\hat{\beta}_j)$ un estimateur de la variance de $\hat{\beta}_j$. La statistique de test de Student pour le test de $H_0: \beta_j = 0$ s'écrit:

$$T_n^{(j)} = \frac{\hat{\beta}_j}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_j)}}.$$

Sous H_0 , la loi de T_n est une loi de Student à $(n - (p + 1))$ degrés de liberté.

Le logiciel R permet d'effectuer simplement ces calculs. Pour ajuster un modèle linéaire sur un échantillon $(Y_i, X_i^{(1)}, \dots, X_i^{(p)})_{i=1, \dots, n}$, on stockera les n valeurs des Y_i dans un vecteur \mathbf{y} , puis, pour $k = 1, \dots, p$ on stockera les n valeurs des $X_i^{(k)}$ dans un vecteur \mathbf{xk} . Le modèle linéaire est alors ajusté au moyen de l'instruction suivante (où $p = 3$)

```
mylm <- lm(y ~ x1 + x2 + x3)
```

et le résultat est stocké dans l'objet `mylm`. Pour accéder à l'ensemble des quantités calculées, on peut exécuter l'instruction suivante:

```
summary(mylm)
```

Sont notamment calculées, les estimations des β_j pour $j = 0, \dots, p$, les estimations de leurs écarts-types, la statistique de test de Student correspondante et la p-valeur du test associé. L'objet `summary(mylm)` est une liste dont on peut récupérer chacune des composantes qui nous intéresse. Pour cela, l'instruction `str(summary(myfit))` permet de voir le nom et la structure des différentes composantes de cette liste.

Exercice 1.

1. Illustrer de manière empirique à partir de données simulées le comportement de la taille du test.
2. Illustrer de manière empirique à partir de données simulées le comportement de la fonction puissance.

Vous veillerez à faire varier tour à tour:

- le risque de 1ère espèce α ,
- la taille de l'échantillon n ,
- le cas échéant l'ampleur de l'écart à H_0 .

Que pensez-vous de la robustesse du test de Student à la bonne spécification du modèle (cf les hypothèses (\mathbf{M}_1) - (\mathbf{M}_7))?

NB: on dit que le modèle est bien spécifié lorsqu'il correspond effectivement au mécanisme ayant servi à générer les données.