

---

**Sujet 10: tests de comparaison de 2 variances ou plus**

---

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

Soit  $X$  une variable aléatoire à valeurs dans un ensemble  $\mathcal{X}$ , de loi  $P_\theta$  où  $\theta \in \Theta$ . Supposons que l'on veuille effectuer un test statistique d'hypothèses portant sur  $\theta$ . 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  $\Theta$  notés  $\Theta_0$  et de  $\Theta_1$  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  $\alpha$ .

Dans ce cadre, on parle de **fonction de risque de 1ère espèce** définie sur  $\Theta_0$  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**  $\alpha$  si sa taille est égale à  $\alpha$  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  $\Theta_1$  pour

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

On parle de **fonction puissance** définie sur  $\Theta_1$  pour

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

• **Test de comparaison de 2 variances:**

Soient  $X$  et  $Y$  deux variables aléatoires indépendantes. On souhaite tester l'égalité de leurs variances respectives. Formellement, on souhaite tester l'hypothèse nulle  $H_0: \text{Var}(X) = \text{Var}(Y)$  au niveau  $\alpha$ . Soit  $(X_1, \dots, X_{n_X})$  un échantillon i.i.d. distribué comme la variable  $X$ . Soit  $(Y_1, \dots, Y_{n_Y})$  un échantillon i.i.d. distribué comme la variable  $Y$ , indépendant de  $(X_1, \dots, X_{n_X})$ . Soit

$$S'_{n_X}{}^2 = \frac{1}{n_X - 1} \sum_{i=1}^{n_X} (\bar{X}_{n_X} - X_i)^2$$

un estimateur de  $\text{Var}(X)$  avec  $\bar{X}_{n_X} = \frac{1}{n_X} \sum_{i=1}^{n_X} X_i$ . Soit

$$S'_{n_Y}{}^2 = \frac{1}{n_Y - 1} \sum_{i=1}^{n_Y} (\bar{Y}_{n_Y} - Y_i)^2$$

un estimateur de  $\text{Var}(Y)$  avec  $\bar{Y}_{n_Y} = \frac{1}{n_Y} \sum_{i=1}^{n_Y} Y_i$ . Notons également

$$\begin{aligned} S'_{n_X+n_Y}{}^2 &= \frac{(n_X - 1)S'_{n_X}{}^2 + (n_Y - 1)S'_{n_Y}{}^2}{n_X + n_Y - 2} \\ &= \frac{1}{n_X + n_Y - 2} \left( \sum_{i=1}^{n_X} (\bar{X}_{n_X} - X_i)^2 + \sum_{i=1}^{n_Y} (\bar{Y}_{n_Y} - Y_i)^2 \right). \end{aligned}$$

**1er cas:**  $n_X \geq 30$  et  $n_Y \geq 30$

Soit la statistique de test:

$$T_{n_X, n_Y} = \frac{S'_{n_X}{}^2 - S'_{n_Y}{}^2}{S'_{n_X+n_Y}{}^2 \sqrt{\frac{2}{n_X} + \frac{2}{n_Y}}}.$$

Sous  $H_0$ , la statique  $T_{n_X, n_Y}$  converge en loi vers une variable de loi  $\mathcal{N}(0, 1)$ . Lorsque  $H_0$  n'est pas satisfaite, la statistique  $T_{n_X, n_Y}$  diverge presque-sûrement vers  $\infty$ .

- Lorsque  $H_1$ :  $\text{Var}(X) \neq \text{Var}(Y)$ , la région critique est

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_X}) \in \mathbb{R}^{n_X}, (y_1, \dots, y_{n_Y}) \in \mathbb{R}^{n_Y} : |t_{n_X, n_Y}| > F_{\mathcal{N}(0,1)}^{-1}(1 - \alpha/2) \right\}.$$

- Lorsque  $H_1$ :  $\text{Var}(X) > \text{Var}(Y)$ , la région critique est

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_X}) \in \mathbb{R}^{n_X}, (y_1, \dots, y_{n_Y}) \in \mathbb{R}^{n_Y} : t_{n_X, n_Y} > F_{\mathcal{N}(0,1)}^{-1}(1 - \alpha) \right\}.$$

- Lorsque  $H_1$ :  $\text{Var}(X) < \text{Var}(Y)$ , la région critique est

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_X}) \in \mathbb{R}^{n_X}, (y_1, \dots, y_{n_Y}) \in \mathbb{R}^{n_Y} : t_{n_X, n_Y} < F_{\mathcal{N}(0,1)}^{-1}(\alpha) \right\}.$$

**2ème cas** (test de Fisher ou de Fisher-Snedecor):  $X \sim \mathcal{N}(\mathbb{E}[X], \text{Var}(X))$  et  $Y \sim \mathcal{N}(\mathbb{E}[Y], \text{Var}(Y))$ .

Soit la statistique de test:

$$T_{n_X, n_Y} = \frac{S'_{n_X}{}^2}{S'_{n_Y}{}^2}.$$

Sous  $H_0$ , la statique  $T_{n_X, n_Y}$  suit la loi  $F(n_X, n_Y)$  dite de Fisher à  $n_X$  et  $n_Y$  degrés de liberté. Heuristiquement, lorsque  $H_0$  est vraie, le rapport  $T_{n_X, n_Y}$  ne doit pas trop s'éloigner de 1.

- Lorsque  $H_1$ :  $\text{Var}(X) \neq \text{Var}(Y)$ , la région critique est

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_X}) \in \mathbb{R}^{n_X}, (y_1, \dots, y_{n_Y}) \in \mathbb{R}^{n_Y} : t_{n_X, n_Y} > F_{F(n_X, n_Y)}^{-1}(1 - \alpha/2) \text{ ou } t_{n_X, n_Y} < F_{F(n_X, n_Y)}^{-1}(\alpha/2) \right\}$$

- Lorsque  $H_1$ :  $\text{Var}(X) > \text{Var}(Y)$ , la région critique est

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_X}) \in \mathbb{R}^{n_X}, (y_1, \dots, y_{n_Y}) \in \mathbb{R}^{n_Y} : t_{n_X, n_Y} > F_{F(n_X, n_Y)}^{-1}(1 - \alpha) \right\}.$$

- Lorsque  $H_1: \text{Var}(X) < \text{Var}(Y)$ , la région critique est

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_X}) \in \mathbb{R}^{n_X}, (y_1, \dots, y_{n_Y}) \in \mathbb{R}^{n_Y} : t_{n_X, n_Y} < F_{F(n_X, n_Y)}^{-1}(\alpha) \right\}.$$

• **Test de comparaison de  $K$  variances:**

Soient  $X^{(k)}$ , pour  $k = 1, \dots, K$ , des variables aléatoires indépendantes où  $X^{(k)} \sim \mathcal{N}(\mathbb{E}[X^{(k)}], \text{Var}(X^{(k)}))$ , pour  $k = 1, \dots, K$ . On souhaite tester l'égalité de leurs variances respectives. Formellement, on souhaite tester au niveau  $\alpha$  l'hypothèse nulle  $H_0: \text{Var}(X^{(k)}) = \sigma^2$  pour  $k = 1, \dots, K$  contre l'hypothèse alternative  $H_1: \exists k_1 \neq k_2 \in \{1, \dots, K\}$  tels que  $\text{Var}(X^{(k_1)}) \neq \text{Var}(X^{(k_2)})$ . Pour  $k = 1, \dots, K$ , soit  $(X_1^{(k)}, \dots, X_{n_k}^{(k)})$  un échantillon i.i.d. distribué comme la variable  $X^{(k)}$ . Pour  $k = 1, \dots, K$ , soit

$$S'_{k, n_k}{}^2 = \frac{1}{n_k - 1} \sum_{i=1}^{n_k} \left( \bar{X}_{n_k}^{(k)} - X_i^{(k)} \right)^2$$

un estimateur de  $\text{Var}(X^{(k)})$  avec

$$\bar{X}_{n_k}^{(k)} = \frac{1}{n_k} \sum_{i=1}^{n_k} X_i^{(k)}.$$

**1er cas:** Test de Bartlett.

Soit la statistique de test:

$$T_{n_1, \dots, n_K} = \frac{(n - K) \left[ \log \left( \frac{1}{n - K} \sum_{k=1}^K (n_k - 1) S'_{k, n_k}{}^2 \right) - \frac{1}{n - K} \sum_{k=1}^K (n_k - 1) \log(S'_{k, n_k}{}^2) \right]}{1 + \frac{1}{3(K-1)} \left( \sum_{k=1}^K \frac{1}{n_k - 1} - \frac{1}{n - K} \right)}$$

en notant  $n = \sum_{k=1}^K n_k$ . Sous  $H_0$ , la statistique  $T_{n_1, \dots, n_K}$  suit la loi  $\chi^2(K - 1)$ . Notons  $t_{n_1, \dots, n_K}$  la réalisation de  $T_{n_1, \dots, n_K}$ . La région critique est:

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_k})_{k=1, \dots, K} : t_{n_1, \dots, n_K} > F_{\chi^2(K-1)}^{-1}(1 - \alpha) \right\}.$$

La fonction `bartlett.test` du logiciel R implémente ce test.

**2ème cas:** Test de Levene.

Introduisons:

$$\begin{aligned} Z_i^{(k)} &= \left| X_i^{(k)} - \bar{X}_{n_k}^{(k)} \right| \\ \bar{Z}_{n_k}^{(k)} &= \frac{1}{n_k} \sum_{i=1}^{n_k} Z_i^{(k)} \\ \bar{Z}_n &= \frac{1}{n} \sum_{k=1}^K \sum_{i=1}^{n_k} Z_i^{(k)} = \frac{1}{n} \sum_{k=1}^K n_k \bar{Z}_{n_k}^{(k)}. \end{aligned}$$

Soit la statistique de test:

$$T_{n_1, \dots, n_K} = \frac{n - K}{K - 1} \frac{\sum_{k=1}^K \left( \bar{Z}_{n_k}^{(k)} - \bar{Z}_n \right)^2}{\sum_{k=1}^K \sum_{i=1}^{n_k} \left( Z_i^{(k)} - \bar{Z}_{n_k}^{(k)} \right)^2}.$$

Sous  $H_0$ , la statistique  $T_{n_1, \dots, n_K}$  suit la loi de Fisher  $F(K - 1, n - K)$ . Notons  $t_{n_1, \dots, n_K}$  la réalisation de  $T_{n_1, \dots, n_K}$ . La région critique est:

$$\mathcal{R} = \left\{ (x_1, \dots, x_{n_k})_{k=1, \dots, K} : t_{n_1, \dots, n_K} > F_{F(K-1, n-K)}^{-1}(1 - \alpha) \right\}.$$

La fonction `leveneTest` du package `car` du logiciel `R` implémente ce test ainsi que la fonction `levene.test` du package `lawstat`.

### 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.
3. Tests paramétriques: Illustrer de manière empirique à partir de données simulées la robustesse ou au contraire la sensibilité du test aux hypothèses sous-jacentes.
4. Tests non-paramétriques: faire varier la loi servant à générer les données ainsi que la valeur de son/ses paramètre(s).

Vous veillerez à faire varier tour à tour:

- le risque de 1ère espèce  $\alpha$ ,
- la taille de l'échantillon  $n$ ,
- la variabilité du phénomène.

Le test de Levene a la réputation d'être robuste aux écarts à la normalité. Qu'en pensez-vous?