

Quelques tests non paramétriques.¹

1. Les tests non paramétriques sur un échantillon

Dans cette section nous nous intéressons à deux tests non paramétriques :

- le test du signe et
- le test des rangs signés.

Nous utiliserons de préférence le test des rangs signés dès que les conditions de son utilisation sont remplies, sa puissance étant alors supérieure à celle du test du signe.

1.1. Test du signe

Soit un échantillon indépendant et identiquement distribué X_1, \dots, X_n d'une loi continue F dont la valeur médiane est notée m_e et la moyenne μ . Le test du signe permet de tester les hypothèses suivantes.

Hypothèses :

$$\mathcal{H}_0 : m_e = 0 \quad \text{ou de façon équivalente} \quad \mathbb{P}[X_i > 0] = \frac{1}{2}$$

contre

$$\mathcal{H}_1 : m_e \neq 0 \quad \text{ou de façon équivalente} \quad \mathbb{P}[X_i > 0] \neq \frac{1}{2}.$$

Remarque 1.1. La formulation de ce test est bien sûr la formulation d'un test bilatéral. Nous pouvons envisager les deux tests unilatéraux correspondants. À ce moment là, la formulation de l'hypothèse alternative \mathcal{H}_1 est différente et s'écrit soit :

$$\mathcal{H}'_1 : \mathbb{P}[X_i > 0] < \frac{1}{2}$$

soit

$$\mathcal{H}''_1 : \mathbb{P}[X_i > 0] > \frac{1}{2}.$$

Remarque 1.2. Plus généralement ce test permet de tester l'hypothèse nulle

$$\mathcal{H}_0 : m_e = m_0 \quad \text{ou de façon équivalente} \quad \mathbb{P}[X_i > 0] = p$$

contre

$$\mathcal{H}_1 : m_e \neq m_0 \quad \text{ou de façon équivalente} \quad \mathbb{P}[X_i > 0] \neq p$$

où m_0 est un nombre réel et p est une constante comprise entre 0 et 1, ou encore, dans la version unilatérale, contre l'hypothèse alternative

¹Les références [?], [?] et [?] ayant servi à l'élaboration de ce document sont mentionnées dans la bibliographie.

$$\mathcal{H}'_1 : m_e < m_0$$

ou encore, dans la version unilatérale, contre l'hypothèse alternative

$$\mathcal{H}''_1 : m_e > m_0.$$

Pour cela il suffit de considérer l'échantillon Z_1, \dots, Z_n avec $Z_i = X_i - m_0$ et de lui appliquer le test décrit ci-dessous.

Statistique : S_n désigne le nombre de variables X_i , $1 \leq i \leq n$, qui prennent une valeur positive.

Propriétés 1.1. *Lorsque l'hypothèse nulle \mathcal{H}_0 est vraie, la variable aléatoire S_n suit exactement une loi binomiale $\mathcal{B}(n, p)$ de paramètres n et p .*

Concrètement cette hypothèse nulle \mathcal{H}_0 signifie que l'effectif de l'échantillon considéré est faible devant celui de la population dont il est issu.

Remarque 1.3. Nous pourrions prendre comme taille limite des échantillons dont les effectifs sont inférieurs à une fraction de 1/10 de la population. Dans ce cas nous pouvons assimiler les tirages réalisés ici à des tirages avec remise.

Cas le plus souvent utilisé : $p = 1/2$. Nous nous proposons de tester :

Hypothèses :

$$\mathcal{H}_0 : \mathbb{P}[X_i > 0] = \frac{1}{2}$$

contre

$$\mathcal{H}_1 : \mathbb{P}[X_i > 0] \neq \frac{1}{2}.$$

Statistique : S_n désigne le nombre de variables X_i , $1 \leq i \leq n$, qui prennent une valeur positive.

Propriétés 1.2. *Lorsque l'hypothèse nulle \mathcal{H}_0 est vraie, la variable aléatoire S_n a les trois propriétés suivantes :*

1. *La variable aléatoire S_n suit une loi binomiale $\mathcal{B}(n, 1/2)$ de paramètres n et $1/2$. De ce fait, découle les deux propriétés suivantes :*
2. $\mathbb{E}[S_n] = n/2$.
3. $\text{Var}[S_n] = n/4$.

Cette distribution binomiale est symétrique. Pour n grand ($n \geq 40$), nous pouvons utiliser l'approximation normale avec correction de continuité :

$$\mathbb{P}_{\mathcal{H}_0}[S_n \leq h] = \mathbb{P}_{\mathcal{N}(n/2, n/4)}[S_n \geq n - h] = \frac{\Phi(2h + 1 - n)}{\sqrt{n}}$$

où Φ est la fonction de répartition de la loi normale centrée réduite.

Décision 1.1. Pour un seuil donné α ($= 5\%$ en général), nous cherchons le plus grand entier s_α^* tel que $\mathbb{P}[Y \leq s_\alpha^*] \leq \alpha/2$ où Y suit une loi binomiale $\mathcal{B}(n, 1/2)$ de paramètres n et $1/2$. Alors nous décidons :

$$\begin{cases} \mathcal{H}_1 \text{ est vraie si } & S_{n,obs} \notin]s_\alpha^*, n - s_\alpha^*[\\ \mathcal{H}_0 \text{ est vraie si } & S_{n,obs} \in]s_\alpha^*, n - s_\alpha^*[. \end{cases}$$

Remarque 1.4. Le niveau de signification réel du test est alors égal à $2\mathbb{P}[Y \leq s_\alpha^*]$ qui est généralement différent de α .

1.2. Test des rangs signés de Wilcoxon

Soit un échantillon indépendant et identiquement distribué X_1, \dots, X_n d'une loi continue F dont la valeur médiane est notée m_e et la moyenne μ . Le test des rangs signés permet de tester les hypothèses suivantes.

Hypothèses :

\mathcal{H}_0 : La loi continue F est symétrique en 0

contre

\mathcal{H}_1 : La loi continue F n'est pas symétrique en 0.

De plus, si nous savons que la loi continue F est symétrique, alors le test des rangs signés de Wilcoxon devient

\mathcal{H}_0 : $\mu = \mu_0$

contre

\mathcal{H}_1 : $\mu \neq \mu_0$.

Ici μ_0 est un nombre réel et ce jeu d'hypothèses permet alors de s'intéresser à la moyenne de la loi continue F .

1.2.1. Cas où il n'y a pas d'ex æquo.

Soit x_1, \dots, x_n réalisations de l'échantillon précédent. À chaque x_i nous attribuons le rang r_i^a qui correspond au rang de $|x_i|$ lorsque que les n réalisations sont classées par ordre croissant de leurs valeurs absolues.

Statistique : Nous déterminons alors la somme w des rangs r_i^a des seules observations positives. La statistique W_n^+ des rangs signés de Wilcoxon est la variable aléatoire qui prend pour valeur la somme w . Par conséquent, la statistique W_n^+ des rangs signés de Wilcoxon s'écrit

$$W_n^+ = \sum_{1 \leq i \leq n} R_i^a, \quad X_i > 0$$

Propriétés 1.3. Lorsque l'hypothèse nulle \mathcal{H}_0 est vraie, la variable aléatoire W_n^+ a les trois propriétés suivantes :

1. W_n^+ est symétrique autour de sa valeur moyenne $\mathbb{E}[W_n^+] = n(n+1)/4$.
2. $\text{Var}[W_n^+] = n(n+1)(2n+1)/24$.
3. La variable aléatoire W_n^+ est tabulée pour de faibles valeurs de n . Pour $n \geq 15$, nous avons l'approximation normale avec correction de continuité :

$$\mathbb{P}[W_n^+ \leq w] = \Phi \left(\frac{w + 0,5 - n(n+1)/4}{\sqrt{n(n+1)(2n+1)/24}} \right)$$

où Φ est la fonction de répartition de la loi normale centrée réduite.

Décision 1.2.

– **Premier cas :** Pour tester l'hypothèse nulle « \mathcal{H}_0 : La loi continue F est symétrique en 0 » contre l'hypothèse alternative « \mathcal{H}_1 : La loi continue F n'est pas symétrique en 0 » pour un seuil donné α , nous cherchons l'entier w_α tel que $\mathbb{P}[W_n^+ \leq w_\alpha] \approx \alpha/2$. Alors nous décidons :

$$\begin{cases} \mathcal{H}_1 \text{ est vraie si } & W_{n,obs}^+ \notin]w_\alpha + 1, n(n+1)/2 - w_\alpha - 1[\\ \mathcal{H}_0 \text{ est vraie si } & W_{n,obs}^+ \in]w_\alpha + 1, n(n+1)/2 - w_\alpha - 1[. \end{cases}$$

– **Second cas :** Pour tester l'hypothèse nulle « \mathcal{H}_0 : $\mu = \mu_0$ », nous introduisons l'échantillon Z_1, \dots, Z_n avec $Z_i = X_i - \mu$, $1 \leq i \leq n$.

1.2.2. Cas où il y a des ex æquo.

Les observations x_1, \dots, x_n peuvent présenter des ex æquo et a fortiori leurs valeurs absolues. Il s'agit en particulier du cas où la loi F est discrète. Deux procédures sont alors employées.

- *Méthode de départition des ex æquo*

Nous départageons les ex æquo à l'aide d'une table de nombres aléatoires. À chacune des valeurs égales nous associons un entier au hasard puis nous affectons, par ordre croissant de ces entiers, un rang différent à chaque observation. Ainsi chacun des rangs des observations est différent et nous pouvons directement appliquer les résultats du paragraphe précédent.

- *Méthode des rangs moyens*

En associant à la variable X_i son rang moyen R_i^{a*} dans le classement des valeurs absolues et en sommant tous les rangs pour lesquels $X_i > 0$ nous obtenons la statistique :

$$W_n^{+*} = \sum_{1 \leq i \leq n} R_i^{a*}, \quad X_i > 0$$

Les valeurs absolues observées $|x_1|, \dots, |x_n|$ étant ordonnées puis regroupées en classes d'ex æquo, C_0 pour la première classe qui est constituée des nombres $|x_i|$

nuls, s'il en existe, et C_j , $1 \leq j \leq h$ pour les autres nombres, certaines classes C_j pouvant comporter un seul élément, si cet élément n'a pas d'ex æquo, notons d_j le nombre d'ex æquo de la classe C_j . Nous avons

$$d_0 + \sum_{j=1}^h d_j = n.$$

Sous l'hypothèse nulle \mathcal{H}_0 et si $n > 15$, il est d'usage d'utiliser l'approximation normale

$$\frac{W_n^{+*} - m^*}{\sigma^*} \approx \mathcal{N}(0, 1)$$

où

$$m^* = \frac{1}{4}(n(n+1) - d_0(d_0+1))$$

et

$$(\sigma^*)^2 = \frac{1}{24}(n(n+1)(2n+1) - d_0(d_0+1)(2d_0+1)) - \frac{1}{48} \sum_{j=1}^h (d_j^3 - d_j).$$

Dans le cas où nous utilisons cette méthode des rangs moyens, nous ne pouvons pas utiliser les tables statistiques usuelles qui concernent la distribution de la variable aléatoire W_n^{+*} .

2. Les tests non paramétriques sur deux échantillons

2.1. Les échantillons sont indépendants : Test de Mann-Whitney

Le test de Mann-Whitney a été introduit en 1947 indépendamment du test de Wilcoxon de la somme des rangs qui a été élaboré en 1945. Ces deux tests, d'une formulation différente, sont en fait équivalents. En fonction de l'outil informatique que vous utiliserez, la dénomination du test pourra être l'une des suivantes : Test de Mann-Whitney, Test de Wilcoxon de la somme des rangs ou encore Test de Mann-Whitney-Wilcoxon. L'approche de Mann et Whitney paraît souvent plus facile à mettre en pratique. Si vous devez utiliser une table, il vous faudra déterminer quelle a été l'approche utilisée par le logiciel et vous servir de la table appropriée.

Nous observons, de manière indépendante, une variable Y , continue, sur deux populations, ou sur une population divisée en deux sous-populations. Nous notons \mathcal{L}_1 la loi de Y sur la (sous-)population d'ordre i . Nous allons présenter le test des hypothèses suivantes.

Hypothèses :

\mathcal{H}_0 : Les deux lois \mathcal{L}_1 sont égales ou encore de façon équivalente : $\mathcal{L}_1 = \mathcal{L}_2$

contre

\mathcal{H}_1 : Les deux lois \mathcal{L}_i ne sont pas égales ou encore de façon équivalente : $\mathcal{L}_1 \neq \mathcal{L}_2$.

2.1.1. Cas où il n'y a pas d'ex æquo.

Statistique : Pour obtenir la statistique du test notée U_{n_1, n_2} en général, nous devons procéder à des étapes successives :

1. En nous plaçant sous l'hypothèse nulle \mathcal{H}_0 , nous classons par ordre croissant l'ensemble des observations des deux échantillons (x_1, \dots, x_{n_1}) et (y_1, \dots, y_{n_2}) de taille respective n_1 et n_2 .
2. Nous affectons le rang correspondant.
3. Nous effectuons la somme des rangs pour chacun des deux échantillons, notés R_1 et R_2 .
4. Nous en déduisons les quantités U_1 et U_2 qui se calculent ainsi :

$$(2.1) \quad U_1 = n_1 \times n_2 + \frac{n_1(n_1+1)}{2} - R_1$$

et

$$(2.2) \quad U_2 = n_1 \times n_2 + \frac{n_2(n_2+1)}{2} - R_2 = n_1 \times n_2 - U_1.$$

La plus petite des deux valeurs U_1 et U_2 , notée U_{n_1, n_2} , est utilisée pour tester l'hypothèse nulle \mathcal{H}_0 .

Propriétés 2.1. Lorsque l'hypothèse nulle \mathcal{H}_0 est vraie, la variable aléatoire U_{n_1, n_2} a les trois propriétés suivantes :

1. $\mathbb{E}[U_{n_1, n_2}] = (n_1 \times n_2)/2$.
2. $\text{Var}[U_{n_1, n_2}] = (n_1 \times n_2)(n_1 + n_2 + 1)/12$.
3. La variable aléatoire U_{n_1, n_2} est tabulée pour de faibles valeurs de n . Pour $n \geq 20$, nous avons l'approximation normale :

$$\mathbb{P}[U_{n_1, n_2} \leq u] = \Phi \left(\frac{u - (n_1 \times n_2)/2}{\sqrt{(n_1 \times n_2)(n_1 + n_2 + 1)/12}} \right)$$

où Φ est la fonction de répartition de la loi normale centrée réduite.

Décision 2.1.

– **Premier cas :** Si les tailles n_1 ou n_2 sont inférieures à 20, alors, pour un seuil donné α ($= 5\% = 0, 05$ en général), la table de Mann-Whitney nous fournit une valeur critique c . Alors nous décidons :

$$\begin{cases} \mathcal{H}_1 \text{ est vraie si } & U_{n_1, n_2, \text{obs}} \leq c, \\ \mathcal{H}_0 \text{ est vraie si } & U_{n_1, n_2, \text{obs}} > c. \end{cases}$$

– **Second cas :** Si les tailles n_1 et n_2 sont supérieures à 20, alors la quantité est décrite approximativement par une loi normale et nous utilisons alors le test de l'écart réduit :

$$Z_{n_1, n_2} = \frac{U_{n_1, n_2} - (n_1 \times n_2)/2}{\sqrt{(n_1 \times n_2)(n_1 + n_2 + 1)/12}}$$

Pour un seuil donné α ($= 5\% = 0, 05$ en général), la table de la loi normale centrée réduite nous fournit une valeur critique c . Alors nous décidons :

$$\begin{cases} \mathcal{H}_1 \text{ est vraie si } & Z_{n_1, n_2, \text{obs}} \geq c, \\ \mathcal{H}_0 \text{ est vraie si } & Z_{n_1, n_2, \text{obs}} < c. \end{cases}$$

2.1.2. Cas où il y a des ex æquo.

Les observations $x_1, \dots, x_{n_1}, y_1, \dots, y_{n_2}$ peuvent présenter des ex æquo. Il s'agit en particulier du cas où les lois F et G dont sont issus les deux échantillons sont discrètes. Deux procédures sont alors employées.

- *Méthode de répartition des ex æquo*

Nous départageons les ex æquo à l'aide d'une table de nombres aléatoires. À chacune des valeurs égales nous associons un entier au hasard puis nous affectons, par ordre croissant de ces entiers, un rang différent à chaque observation. Ainsi chacun des rangs des observations est différent et nous pouvons directement appliquer les résultats du paragraphe précédent.

- *Méthode des rangs moyens*

Les valeurs absolues observées $x_1, \dots, x_{n_1}, y_1, \dots, y_{n_2}$ étant ordonnées puis regroupées en h classes d'ex æquo C_j , $1 \leq j \leq h$, certaines classes C_j pouvant comporter un seul élément, si cet élément n'a pas d'ex æquo, notons d_j le nombre d'ex æquo de la classe C_j . Nous avons

$$\sum_{j=1}^h d_j = n_1 + n_2.$$

En associant à l'observation X_i son rang moyen R_i^* dans ce classement et en sommant tous les rangs de tous les X_i , nous obtenons la statistique :

$$U_{n_1, n_2}^* = \sum_{i=1}^{n_2} R_i^*.$$

Sous l'hypothèse nulle \mathcal{H}_0 : « X et Y ont la même distribution » et pour $n_1 > 15$ et $n_2 > 15$, il est d'usage d'utiliser l'approximation normale

$$\frac{U_{n_1, n_2}^* - m^*}{\sigma^*} \approx \mathcal{N}(0, 1)$$

où

$$m^* = \frac{1}{2} (n_1(n_1 + n_2 + 1))$$

et

$$(\sigma^*)^2 = \frac{1}{12} (n_1 n_2 (n_1 + n_2 + 1)) - \frac{1}{12} \frac{n_1 n_2}{(n_1 + n_2)(n_1 + n_2 - 1)} \sum_{j=1}^h (d_j^3 - d_j).$$

Dans le cas où nous utilisons cette méthode des rangs moyens nous ne pouvons pas utiliser les tables statistiques usuelles qui concernent la distribution de la variable aléatoire U_{n_1, n_2} .

2.2. Les échantillons sont indépendants : Test de la médiane de Mood

Nous considérons deux échantillons indépendants (X_1, \dots, X_{n_1}) et (Y_1, \dots, Y_{n_2}) . (X_1, \dots, X_{n_1}) est un échantillon indépendant et identiquement distribué d'une loi continue

F et (Y_1, \dots, Y_{n_2}) est un échantillon indépendant et identiquement distribué d'une loi continue G .

Nous allons tester les hypothèses suivantes.

Hypothèses :

\mathcal{H}_0 : Les deux lois F et G sont égales ou encore de façon équivalente : $F = G$

contre

\mathcal{H}_1 : Les deux lois F et G ne sont pas égales ou encore de façon équivalente : $F \neq G$.

Après regroupement des $n_1 + n_2$ valeurs des deux échantillons, $n_1 \times M_N$ est le nombre d'observations X_i qui sont supérieures à la médiane des $N = n_1 + n_2$ observations.

Sous l'hypothèse nulle \mathcal{H}_0 : « Les variables X et Y suivent la même loi continue c'est-à-dire $G = F$ », la variable $n_1 \times M_N$ peut prendre les valeurs $0, 1, \dots, n_1$ selon la distribution hypergéométrique suivante :

$$\mathbb{P} [n_1 \times M_N = k] = \frac{C_{n_1}^k C_{n_2}^{N/2-k}}{C_N^{N/2}}.$$

Ainsi nous avons :

$$\mathbb{E} [n_1 \times M_N] = \frac{n_1(n_1 + n_2 - \epsilon_N)}{2N}$$

$$\text{Var} [n_1 \times M_N] = \frac{n_1 n_2 (n_1 + n_2 + 1)}{4(n_1 + n_2 - 1 + \epsilon_N)(n_1 + n_2 + 1 - \epsilon_N)},$$

où $\epsilon_N = 0$ si N est pair et $\epsilon_N = 1$ si N est impair.

Lorsque les tailles n_1 et n_2 sont grandes, c'est-à-dire $n_1 \geq 25$ et $n_2 \geq 25$, nous utilisons l'approximation normale :

$$\frac{n_1 \times M_N - \mathbb{E} [n_1 \times M_N]}{\sqrt{\text{Var} [n_1 \times M_N]}} \approx \mathcal{N}(0, 1)$$

avec correction de continuité.

La distribution est symétrique lorsque N est pair.

Pour tester l'hypothèse nulle \mathcal{H}_0 : « $G = F$ » contre \mathcal{H}_1 : « $G \neq F$ » avec un niveau de signification égal à α , nous cherchons les entiers k_α et k'_α tels que $\mathbb{P} [n_1 \times M_N \leq k_\alpha] \approx \alpha/2$ et $\mathbb{P} [n_1 \times M_N \geq n_1 - k'_\alpha] \approx \alpha/2$, puis nous rejetons l'hypothèse nulle \mathcal{H}_0 si la réalisation de la statistique du test calculée à l'aide de l'échantillon n'est pas dans l'intervalle $[k_\alpha, k'_\alpha]$. Cette statistique permet également de réaliser des tests unilatéraux.

2.3. Les échantillons sont dépendants : Test de Wilcoxon

Nous considérons deux variables aléatoires X et Y , de même nature, observées toutes les deux sur les mêmes unités d'un n -échantillon. Les observations se présentent alors sous la forme d'une suite de couples $(x_1, y_1), \dots, (x_n, y_n)$.

Ce test concerne les lois des deux variables. Pour ce faire nous testons les hypothèses suivantes.

Hypothèses :

\mathcal{H}_0 : Les deux lois sont égales ou encore de façon équivalente $\mathcal{L}(X) = \mathcal{L}(Y)$

contre

\mathcal{H}_1 : Les deux lois ne sont pas égales ou encore de façon équivalente $\mathcal{L}(X) \neq \mathcal{L}(Y)$.

2.3.1. Cas où il n'y a pas d'ex æquo.

Statistique : Pour obtenir la statistique du test notée S^+ en général, nous devons procéder à des étapes successives :

1. Ce test suppose que la loi de la différence entre les deux variables étudiées est symétrique par rapport à 0.
2. Après avoir calculé les différences d_i , nous classons par ordre croissant les $|d_i|$ non nulles, c'est-à-dire les d_i sans tenir compte des signes.
3. Nous attribuons à chaque $|d_i|$ le rang correspondant.
4. Nous restituons ensuite à chaque rang le signe de la différence correspondante.
5. Enfin, nous calculons la somme S^+ des rangs positifs (P) et la somme S^- des rangs négatifs (M).

La somme S^+ des rangs positifs (P) permet de tester l'hypothèse nulle \mathcal{H}_0 .

Décision 2.2.

- **Premier cas :** Si $n < 15$, nous utilisons une table et nous comparons la valeur de (S^+) à la valeur critique c associée au seuil α du test.
- **Second cas :** Si $n \geq 15$, nous utilisons l'approximation normale avec correction de continuité :

$$\mathbb{P}_{\mathcal{H}_0} [S^+ \leq h] \approx \Phi \left(\frac{h + 0,5 - n(n+1)/4}{\sqrt{(n(n+1)(2n+1))/24}} \right)$$

où Φ est la fonction de répartition de la loi normale centrée réduite.

2.3.2. Cas où il y a des ex æquo.

Il se traite de la même manière que pour la statistique de Wilcoxon pour un échantillon, voir le paragraphe 1.2.

3. Les tests non paramétriques sur $k \geq 3$ échantillons : 1 facteur

3.1. Les échantillons sont indépendants : Test de Kruskal-Wallis

Nous supposons que nous disposons de k échantillons **indépendants** et identiquement distribués $(X_{1,1}, \dots, X_{1,n_1}), \dots, (X_{k,1}, \dots, X_{k,n_k})$. La distribution du i -ème échantillon est notée F_i . Nous admettons **a priori** que $F_i(x) = G(x - \alpha_i)$ où G est une fonction de répartition inconnue mais continue de moyenne μ et les α_i sont des nombres réels. Ainsi nous supposons que le seul paramètre qui diffère d'une distribution F_i à l'autre est un paramètre de position α_i . C'est pourquoi même lorsque vous effectuez un test de Kruskal-Wallis vous devez vous assurer que vous pouvez au moins supposer que les variances des variables sont égales d'un échantillon à l'autre à l'aide d'un test non paramétrique de Levene d'égalité des variances.

Les hypothèses ci-dessus impliquent que nous pouvons écrire, pour tout $1 \leq i \leq k$ la décomposition suivante :

$$X_{i,j} = \mu + \alpha_i + \epsilon_{i,j}, \quad 1 \leq j \leq n_i,$$

les $N = \sum_{i=1}^k n_i$ variables aléatoires $\epsilon_{i,j}$ étant indépendantes et ayant une même distribution inconnue et de moyenne nulle.

La variable KW_N de Kruskal-Wallis est utilisée pour tester les hypothèses suivantes.

Hypothèses :

$\mathcal{H}_0 : \alpha_1 = \dots = \alpha_k = 0$

contre

$\mathcal{H}_1 : \text{Il existe au moins un } i_0 \text{ tel que } \alpha_{i_0} \neq 0.$

3.1.1. Cas où il n'y a pas d'ex æquo.

Nous commençons par calculer le rang $R_{i,j}$ de $X_{i,j}$ parmi les N valeurs, puis la somme des rangs associée à chaque échantillon : $R_{i,\bullet} = \sum_{j=1}^{n_i} R_{i,j}$ et enfin la moyenne des rangs de

chaque échantillon : $\overline{R_{i,\bullet}} = \frac{R_{i,\bullet}}{n_i}$.

La statistique de Kruskal-Wallis KW_N prend en compte l'écart entre la moyenne des rangs de chaque échantillon et la moyenne de tous les rangs, qui vaut $(N+1)/2$:

$$\begin{aligned} KW_N &= \frac{12}{N(N+1)} \sum_{i=1}^k n_i \left(\overline{R_{i,\bullet}} - \frac{N+1}{2} \right)^2 \\ &= \frac{12}{N(N+1)} \sum_{i=1}^k \frac{\overline{R_{i,\bullet}}^2}{n_i} - 3(N+1). \end{aligned}$$

Sous l'hypothèse nulle \mathcal{H}_0 : « X_1, \dots, X_k ont la même distribution continue », qui dans notre cas est équivalente à \mathcal{H}_0 : « $\alpha_1 = \dots = \alpha_k = 0$ », il est possible de déterminer la distribution de KW_N bien que le calcul soit complexe.

– Pour $i = 1, \dots, k$, $W_i = n_i \overline{R_{i,\bullet}}$ est la statistique de Wilcoxon qui compare le i -ème traitement aux $k - 1$ autres traitements.

Sous l'hypothèse nulle \mathcal{H}_0 , nous en déduisons que :

$$\mathbb{E}[W_i] = \frac{n_i(N+1)}{2},$$

$$\text{Var}[W_i] = \frac{n_i(N-n_i)(N+1)}{12}.$$

Ainsi par conséquent, nous avons :

$$KW_N = \frac{1}{N} \sum_{i=1}^k (N - n_i) \frac{(W_i - \mathbb{E}[W_i])^2}{\text{Var}[W_i]}.$$

Nous calculons alors l'espérance et la variance de KW_N sous l'hypothèse nulle \mathcal{H}_0 :

$$\mathbb{E}[KW_N] = k - 1,$$

$$\text{Var}[KW_N] = 2(k-1) - \frac{2[3k^2 - 6k + N(2k^2 - 6k + 1)]}{5N(N+1)} - \frac{6}{5} \sum_{i=1}^k \frac{1}{n_i}.$$

– Si l'un des effectifs n_i , $1 \leq i \leq k$, est inférieur ou égal à 4, nous utilisons une table spécifique.

– Si $n_i \geq 5$, pour tout $1 \leq i \leq k$ nous utilisons l'approximation $KW_N \approx \chi_{k-1}^2$.

Pour un seuil de signification de α , nous déterminons c_α tel que $\mathbb{P}[KW_N \geq c_\alpha] \cong \alpha$ et nous rejetons l'hypothèse nulle \mathcal{H}_0 lorsque la valeur prise par KW_N est supérieure à c_α .

3.1.2. Cas où il y a des ex æquo.

- *Méthode de répartition des ex æquo*

Nous répartissons les ex æquo à l'aide d'une table de nombres aléatoires. À chacune des valeurs égales nous associons un entier au hasard puis nous affectons, par ordre croissant de ces entiers, un rang différent à chaque observation. Ainsi chacun des rangs des observations est différent et nous pouvons directement appliquer les résultats du paragraphe précédent.

- *Méthode des rangs moyens*

À chaque nombre appartenant à un groupe d'ex æquo nous attribuons le rang moyen du groupe auquel il appartient puis nous déterminons la somme $T = \sum_{l=1}^h (t_l^2 - t_l)$ où t_l désigne le nombre d'éléments du l -ème groupe d'ex æquo. Il est d'usage de substituer à KW_N la variable KW_N^* définie par :

$$KW_N^* = \frac{KW_N}{1 - N^3 - N}.$$

Comparaisons multiples

Si nous rejetons l'hypothèse nulle \mathcal{H}_0 : « $\alpha_1 = \dots = \alpha_k = 0$ » d'absence de différence entre les distributions F_i des k échantillons, nous pouvons être amenés à nous demander quelles sont les distributions qui sont différentes.

Nous décidons que **deux distributions** F_i et $F_{i'}$ sont significativement différentes au seuil α si :

$$|\overline{R_{i,\bullet}} - \overline{R_{i',\bullet}}| \geq \sqrt{\chi^2(k-1, 1-\alpha)} \sqrt{\frac{N(N+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_{i'}} \right)},$$

où $\chi^2(k-1, 1-\alpha)$ est le $100(1-\alpha)$ quantile de la loi du χ^2 à $k-1$ degrés de liberté.

Nous décidons qu'**au seuil global α** deux distributions F_i et $F_{i'}$, parmi les $\mathbf{k(k-1)}$ **comparaisons** que nous allons faire, sont significativement différentes si :

$$|\overline{R_{i,\bullet}} - \overline{R_{i',\bullet}}| \geq u \left(1 - \frac{\alpha}{k(k-1)} \right) \sqrt{\frac{N(N+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_{i'}} \right)},$$

où $u \left(1 - \frac{\alpha}{k(k-1)} \right)$ est le $100 \left(1 - \frac{\alpha}{k(k-1)} \right)$ quantile de la loi normale centrée réduite.

Il s'agit d'une application des inégalités de Bonferroni². Cette procédure est plus puissante que la précédente.

Nous décidons qu'**au seuil global α** deux distributions F_i et $F_{i'}$, parmi les $\mathbf{k(k-1)}$ **comparaisons** que nous allons faire, sont significativement différentes si :

$$|\overline{R_{i,\bullet}} - \overline{R_{i',\bullet}}| \geq q(k, +\infty, 1-\alpha) \sqrt{\frac{N(N+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_{i'}} \right)},$$

où $q(k, +\infty, 1-\alpha)$ est le $100(1-\alpha)$ quantile de la loi de l'étendue studentisée pour k moyennes et $+\infty$ degrés de liberté. Il s'agit d'une procédure analogue à celle de Tukey-Kramer² dans le cas paramétrique et valide asymptotiquement. Elle est généralement plus puissante que les deux approches précédentes.

3.2. Les échantillons sont indépendants : Test de Jonckheere-Terpstra

La statistique J_N de Jonckheere-Terpstra permet de raffiner l'approche de la statistique KW_N de Kruskal-Wallis : supposons que les k modalités du facteur pour lequel nous avons réalisé au total N expériences soient naturellement ordonnées. C'est par exemple le cas

² Consulter le chapitre ?? pour plus de détails sur les procédures de comparaisons multiples.

dans la situation suivante : vous souhaitez trouver la dose optimale d'engrais à utiliser pour améliorer un rendement. Vous allez donc réaliser des expériences avec des doses de plus en plus importantes d'engrais et les modalités de votre facteur explicatif seront donc naturellement ordonnées par la quantité croissante d'engrais utilisé. Les entiers n_i , pour $1 \leq i \leq k$, désignent les effectifs associés à chacune des modalités du facteur explicatif.

La statistique J_N de Jonckheere-Terpstra permet de tester les hypothèses suivantes.

Hypothèses :

$$\mathcal{H}_0 : \alpha_1 = \dots = \alpha_k = 0$$

contre

$$\mathcal{H}_1 : \alpha_1 \leq \dots \leq \alpha_k = 0 \text{ et il existe au moins un } i_0 \text{ tel que } \alpha_{i_0} < \alpha_{i_0+1}.$$

3.2.1. Cas où il n'y a pas d'ex æquo.

La statistique J_N est construite à l'aide de toutes les variables de Mann-Whitney $U_{i,j}$, associées à l'échantillon i et l'échantillon j , lorsque $1 \leq i < j \leq k$:

$$J_N = \sum_{1 \leq i < j \leq k} U_{i,j}.$$

Sous l'hypothèse nulle \mathcal{H}_0 : « $\alpha_1 = \dots = \alpha_k = 0$ » :

– L'espérance et la variance de la statistique J_N sont :

$$\begin{aligned} \mathbb{E}[J_N] &= \frac{N^2 - \sum_{i=1}^k n_i^2}{4}, \\ \text{Var}[J_N] &= \frac{1}{72} \left(N^2(3 + 2N) - \sum_{i=1}^k n_i^2(3 + 2n_i) \right). \end{aligned}$$

– Les valeurs critiques de la statistique J_N sont tabulées pour de faibles valeurs de k et des n_i .
– Lorsque $n_i \geq 5$, pour tout $1 \leq i \leq k$, nous avons l'approximation normale avec correction de continuité :

$$\frac{J_N - \mathbb{E}[J_N]}{\sqrt{\text{Var}[J_N]}} \approx \mathcal{N}(0, 1).$$

Nous cherchons l'entier ϕ_α tel que $\mathbb{P}[J_N \geq \phi_\alpha] \approx \alpha$ puis nous rejetons l'hypothèse nulle \mathcal{H}_0 au seuil α si la valeur prise par la statistique J_N est supérieure ou égale à ϕ_α .

3.2.2. Cas où il y a des ex æquo.

Nous pouvons utiliser une méthode de répartition des ex æquo ou des tests de Mann-Whitney basés sur des rangs moyens, l'inconvénient de la seconde méthode étant que nous ne pouvons utiliser les mêmes tables qu'en absence d'ex æquo.

3.3. Les échantillons ne sont pas indépendants : Test de Friedman

Nous nous plaçons ici dans le cas où les échantillons utilisés pour tester l'influence d'un facteur **ne sont pas indépendants**.

Individu	Facteur A		
	1	...	n
1	$x_{1,1}$...	$x_{n,1}$
⋮	⋮	⋮	⋮
k	$x_{1,k}$...	$x_{n,k}$

Nous construisons alors le tableau des rangs :

Individu	Facteur A		Total
	1	...	n
1	$r_{1,1}$...	$r_{n,1}$
⋮	⋮	⋮	⋮
k	$r_{1,k}$...	$r_{n,k}$
Total	$r_{1,\bullet}$...	$r_{n,\bullet}$
			$k\bar{n}(n+1)/2$

Si nous sommes en présence de répétitions $x_{i,j,k}$ nous remplaçons $x_{i,j}$ par la moyenne $\bar{x}_{i,j}$ des valeurs pour chaque cas où il y a des répétitions.

Nous cherchons à tester les hypothèses suivantes.

Hypothèses :

$$\mathcal{H}_0 : \text{Les niveaux du facteur ont tous la même influence}$$

contre

$$\mathcal{H}_1 : \text{Les niveaux du facteur n'ont pas tous la même influence.}$$

3.3.1. Cas où il n'y a pas d'ex æquo.

La statistique de Friedman $F_{k,n}$ est définie par :

$$F_{k,n} = \frac{12k}{n(n+1)} \sum_{i=1}^n \left(\frac{R_{i,\bullet}}{k} - \frac{n+1}{2} \right)^2 = \frac{12}{k\bar{n}(n+1)} \sum_{i=1}^n R_{i,\bullet}^2 - 3k(n+1).$$

Nous admettons que sous l'hypothèse nulle \mathcal{H}_0 : « Les niveaux du facteur ont tous la même influence » les distributions pour chaque individu ne diffèrent que par un paramètre de position, ce que nous pouvons vérifier par un test non paramétrique de Levene par exemple.
– Pour de petites valeurs de k nous utilisons une table spécifique. Il se peut que nous vous fournissions une table du coefficient de concordance $W_{k,n}$ de Kendall car la statistique de Friedman $F_{k,n} = k(n-1)W_{k,n}$.

- Pour des valeurs de k assez grandes nous utilisons l'approximation asymptotique suivante :

$$F_{k,n} \approx \chi_{n-1}^2.$$

Nous rejetons l'hypothèse nulle \mathcal{H}_0 si la valeur prise par $F_{k,n}$ est trop grande.

3.3.2. Cas où il y a des ex æquo.

- *Méthode de répartition des ex æquo*

Nous départageons les ex æquo à l'aide d'une table de nombres aléatoires. À chacune des valeurs égales nous associons un entier au hasard puis nous affectons, par ordre croissant de ces entiers, un rang différent à chaque observation. Ainsi chacun des rangs des observations est différent et nous pouvons directement appliquer les résultats du paragraphe précédent.

- *Méthode des rangs moyens*

Dans chaque classement présentant des ex æquo nous attribuons à chacun de ceux-ci le rang moyen du groupe d'ex æquo auquel il appartient et qui n'est pas nécessairement un entier. Lorsque le classement numéro m a h_m groupes d'ex æquo, nous lui attribuons la somme $T_m = \sum_{l=1}^{h_m} (t_{l,m} - t_{l,m})$ où $t_{l,m}$ désigne le nombre d'éléments du l -ème de ces h_m groupes. S'il n'y a pas d'ex æquo nous avons évidemment $T_m = 0$ puisque la répartition des n entiers du classement en classes de nombres égaux donne $h_m = n$ et $t_{l,m} = 1$ pour tout l . Alors la statistique de Friedman corrigée est définie par :

$$\begin{aligned} F_{k,n}^* &= \frac{12k(n-1)}{(n^3-n) - \frac{1}{k} \sum_{m=1}^n T_m} \sum_{l=1}^n \left(\frac{R_{l,\bullet}}{k} - \frac{n+1}{2} \right)^2 \\ &= \frac{1}{1 - \frac{1}{(n^3-n)k} \sum_{m=1}^n T_m} \frac{12k}{n(n+1)} \sum_{l=1}^n \left(\frac{R_{l,\bullet}}{k} - \frac{n+1}{2} \right)^2. \end{aligned}$$

Nous en déduisons que :

$$F_{k,n}^* = \frac{F_{k,n}}{1 - \frac{1}{(n^3-n)k} \sum_{m=1}^n T_m}.$$

4. Les tests non paramétriques sur nk échantillons : 2 facteurs

4.1. Les échantillons sont indépendants : Test de Friedman

Nous nous plaçons ici dans le cas où les échantillons utilisés pour tester l'influence d'un facteur **sont indépendants**.

Facteur B	Facteur A
1	1 ... n
\vdots	\vdots
k	$x_{1,1}$... $x_{n,1}$
	\vdots
	\vdots
	$x_{1,k}$... $x_{n,k}$

Nous construisons alors le tableau des rangs :

Facteur B	Facteur A	Total
1	1 ... n	
\vdots	\vdots	
k	$r_{1,1}$... $r_{n,1}$	$n(n+1)/2$
	\vdots	
	\vdots	
	$r_{1,k}$... $r_{n,k}$	$n(n+1)/2$
Total	$r_{1,\bullet}$... $r_{n,\bullet}$	$kn(n+1)/2$

Si nous sommes en présence de répétitions $x_{i,j,k}$ nous remplaçons $x_{i,j}$ par la moyenne $\bar{x}_{i,j}$ des valeurs pour chaque cas où il y a des répétitions.

Nous admettons **a priori** que l'influence des couples de niveaux (A_i, B_j) des facteurs A et B , pour $1 \leq i \leq n, 1 \leq j \leq k$, se traduit par une décomposition de la forme :

$$X_{i,j} = \mu + \alpha_i + \beta_j + \epsilon_{i,j}, \quad 1 \leq i \leq n, \quad 1 \leq j \leq k$$

avec $\sum_{i=1}^n \alpha_i = 0$ et $\sum_{j=1}^k \beta_j = 0$. Les $N = \sum_{i=1}^n k = n \times k$ variables aléatoires $\epsilon_{i,j}$ étant indépendantes et ayant une même distribution inconnue et de moyenne nulle.

Nous cherchons à tester les hypothèses suivantes.

Hypothèses :

\mathcal{H}_0 : Les niveaux du facteur A ont tous la même influence

contre

\mathcal{H}_1 : Les niveaux du facteur A n'ont pas tous la même influence.

\mathcal{H}_0 : Les niveaux du facteur B ont tous la même influence

contre

\mathcal{H}_1 : Les niveaux du facteur B n'ont pas tous la même influence.

4.1.1. Cas où il n'y a pas d'ex æquo.

La variable $F_{k,n}$ de Friedman est utilisée pour tester l'hypothèse

$$\mathcal{H}_0 : \alpha_1 = \dots = \alpha_n = 0$$

contre

$$\mathcal{H}_1 : \text{Il existe au moins un } i_0 \text{ tel que } \alpha_{i_0} \neq 0.$$

ou de manière équivalente

$$\mathcal{H}_0 : \text{Les niveaux du facteur } A \text{ ont tous la même influence}$$

contre

$$\mathcal{H}_1 : \text{Les niveaux du facteur } A \text{ n'ont pas tous la même influence.}$$

Nous commençons par calculer le rang $R_{i,j}$ de $X_{i,j}$ parmi les n valeurs de la colonne i ,

puis la somme des rangs associée à chaque colonne : $R_{i,\bullet} = \sum_{j=1}^k R_{i,j}$ et enfin la moyenne

$$\text{des rangs de chaque colonne : } \overline{R_{i,\bullet}} = \frac{R_{i,\bullet}}{k}.$$

La statistique de Friedman $F_{k,n}$ est définie par :

$$F_{k,n} = \frac{12k}{n(n+1)} \sum_{i=1}^n \left(\frac{R_{i,\bullet}}{k} - \frac{n+1}{2} \right)^2 = \frac{12}{kn(n+1)} \sum_{i=1}^n R_{i,\bullet}^2 - 3k(n+1).$$

– Pour de petites valeurs de k nous utilisons une table spécifique. Il se peut que nous vous fournissions une table du coefficient de concordance $W_{k,n}$ de Kendall car $F_{k,n} = k(n-1)W_{k,n}$.

– Pour des valeurs de k assez grandes nous utilisons l'approximation asymptotique suivante :

$$F_{k,n} \approx \chi_{n-1}^2.$$

Nous rejetons l'hypothèse nulle \mathcal{H}_0 si la valeur prise par la statistique de Friedman $F_{k,n}$ est trop grande.

Si nous voulions également tester l'influence du facteur B nous aurions analysé les tableaux ci-dessous avec la même méthode.

Facteur A	Facteur B
	1 ... n
1	$x_{1,1}$... $x_{n,1}$
⋮	⋮ ⋮ ⋮
k	$x_{1,k}$... $x_{n,k}$

Facteur A	Facteur B	Total
	1 ... n	
1	$r_{1,1}$... $r_{n,1}$	$n(n+1)/2$
⋮	⋮ ⋮ ⋮	$n(n+1)/2$
k	$r_{1,k}$... $r_{n,k}$	$n(n+1)/2$
Total	$r_{1,\bullet}$... $r_{n,\bullet}$	$kn(n+1)/2$

Comme nous avons échangé le rôle du facteur A et du facteur B nous testons maintenant :

$$\mathcal{H}_0 : \text{Les niveaux du facteur } B \text{ ont tous la même influence}$$

contre

$$\mathcal{H}_1 : \text{Les niveaux du facteur } B \text{ n'ont pas tous la même influence.}$$

Nous ne pouvons pas tester l'existence d'une interaction par cette méthode puisque le modèle utilisé ne comporte pas de terme d'interaction. Il existe d'autres tests pour étudier l'existence d'une interaction.

4.1.2. Cas où il y a des ex æquo.

- *Méthode de répartition des ex æquo*

Nous départageons les ex æquo à l'aide d'une table de nombres aléatoires. À chacune des valeurs égales nous associons un entier au hasard puis nous affectons, par ordre croissant de ces entiers, un rang différent à chaque observation. Ainsi chacun des rangs des observations est différent et nous pouvons directement appliquer les résultats du paragraphe précédent.

- *Méthode des rangs moyens*

Dans chaque classement présentant des ex æquo nous attribuons à chacun de ceux-ci le rang moyen du groupe d'ex æquo auquel il appartient et qui n'est pas nécessairement un entier. Lorsque le classement numéro m a h_m groupes d'ex æquo, nous lui attribuons la somme $T_m = \sum_{l=1}^{h_m} (t_{l,m}^3 - t_{l,m})$ où $t_{l,m}$ désigne le nombre d'éléments du l -ème de ces h_m groupes. S'il n'y a pas d'ex æquo nous avons évidemment $T_m = 0$ puisque la répartition des n entiers du classement en classes de nombres égaux donne $h_m = n$ et $t_{l,m} = 1$ pour tout l . Alors la statistique de Friedman corrigée est définie par :

$$\begin{aligned} F_{k,n}^* &= \frac{12k(n-1)}{(n^3-n) - \frac{1}{k} \sum_{m=1}^n T_m} \sum_{l=1}^n \left(\frac{R_{l,\bullet}}{k} - \frac{n+1}{2} \right)^2 \\ &= \frac{1}{1 - \frac{1}{(n^3-n)k} \sum_{m=1}^n T_m} \frac{12k}{n(n+1)} \sum_{l=1}^n \left(\frac{R_{l,\bullet}}{k} - \frac{n+1}{2} \right)^2. \end{aligned}$$

Nous en déduisons que :

$$F_{k,n}^* = \frac{F_{k,n}}{1 - \frac{1}{(n^3-n)k} \sum_{m=1}^n T_m}.$$

Feuille de Travaux Pratiques n° 1

Les données des deux premiers exercices sont inspirées du livre de G. Pupion et P.-C. Pupion, éditions Economica, 1998.

Exercice 1. Économie

Nous mesurons un indice économique sur onze entreprises. Nous sommes amenés à nous poser la question suivante : « Pouvons-nous considérer que la médiane associée à cet indice est nulle ? »

Entreprise	Indice
1	1
2	4
3	10
4	20
5	0,5
6	-3
7	-7
8	5
9	4
10	3
11	1

.....

Exercice 2. Étude d'activité

Nous disposons de la variation du chiffre d'affaires de 20 entreprises dans un même secteur d'activité. Le chiffre d'affaires dans ce secteur d'activité est-il resté stable ?

Entreprise	1	2	3	4	5	6	7	8	9	10
x_i	-25	-2156	4525	2697	-379	404	-1123	-1733	-2658	-477
Entreprise	11	12	13	14	15	16	17	18	19	20
x_i	-3568	-12071	165	269	-4306	-983	-582	-1897	-1412	662

.....

Exercice 3. Âge des arbres

Nous souhaitons évaluer une nouvelle méthode permettant de déterminer l'âge d'un arbre sans avoir à l'abattre. Pour ce faire, nous sacrifions 11 arbres pour lesquels nous avons réalisé les deux types mesures : estimation de l'âge de l'arbre à l'aide de la méthode dont nous souhaitons tester l'efficacité puis calcul de l'âge exact de l'arbre après abattage. Nous avons reporté les données dans le tableau ci-dessous :

Arbre	1	2	3	4	5	6	7	8	9	10	11
Âge estimé avant abattage	29	28	42	32	22	32	28	21	30	23	39
Âge réel après abattage	25	24	38	27	19	28	24	22	26	19	34

Pouvons-nous nous fier aux résultats de la nouvelle méthode proposée pour estimer l'âge d'un arbre ?

Exercice 4. Hauteurs des arbres

Nous souhaitons comparer la hauteur des arbres de deux types de hêtres. Pouvons-nous dire, à l'aide des mesures de taille exprimées en m et que nous avons reportées dans le tableau ci-dessous, qu'il y a une différence entre les tailles moyennes des arbres des deux hêtres ?

Type 1	Type 2	Type 1	Type 2
23,4	22,5	24,4	22,9
24,6	23,7	24,9	24,0
25,0	24,4	26,2	24,5
26,3	25,3	26,8	26,0
26,8	26,2	26,9	26,4
27,0	26,7	27,6	26,9
27,7	27,4		28,5

Nous disposons désormais de mesures de taille, exprimées en m , provenant d'une troisième hêtre.

Type 3	18,9	21,1	21,2	22,1	22,5	23,6	24,5	24,6	26,2	26,7
--------	------	------	------	------	------	------	------	------	------	------

Y a-t-il des différences entre les tailles moyennes des arbres provenant des trois différentes hétraies ?

.....

Exercice 5. Rendements fouragers

Nous nous intéressons à l'ensemble des prairies d'une région donnée et nous souhaitons identifier l'importance, absolue ou relative, de la variabilité de la production fourragère, d'une part, d'une prairie à l'autre, et d'autre part, d'un endroit à l'autre, à l'intérieur des différentes prairies. Dans ce but, nous avons tout d'abord choisi au hasard trois prairies, dans l'ensemble du territoire, puis au sein de chacune de ces trois prairies, cinq petites parcelles, de deux mètres carrés. Dans l'optique d'un échantillonnage à deux degrés, les trois prairies constituent trois unités du premier degré, et les quinze petites parcelles quinze unités du deuxième degré.

Dans chacune des petites parcelles, nous avons mesuré les rendements en matière sèche à une date donnée. Les valeurs observées, exprimées en tonne par hectare, figurent dans le tableau ci-dessous.

	Prairie 1	Prairie 2	Prairie 3
Parcelle 1	2,06	1,59	1,92
Parcelle 2	2,99	2,63	1,85
Parcelle 3	1,98	1,98	2,14
Parcelle 4	2,95	2,25	1,33
Parcelle 5	2,70	2,09	1,83

Les rendements sont-ils homogènes ?

.....

Exercice 6. Impact de promotions

Un dirigeant de magasin à succursales multiples envisage de faire trois types de promotions nommées P_1 , P_2 et P_3 qui ont un coût sensiblement égal. Afin de déterminer celle qui sera finalement retenue, il fait tester les trois possibilités de promotion par un total de 16 magasins : 5 pour P_1 , 5 pour P_2 et 6 pour P_3 . Le relevé de δ , le taux d'accroissement du chiffre d'affaires, exprimé en %, de chacun de ces magasins a été reporté dans le tableau ci-dessous.

Promotion	1	2	3	4	5	6
Promotion 1	2,1	3,5	4,0	3,1	2,3	
Promotion 2	1,8	3,6	4,3	2,7	5,1	
Promotion 3	2,2	2,5	3,1	3,8	6,0	3,5

En utilisant la statistique de Jonckheere-Terpstra, déterminer si les promotions ont la même influence sur δ le taux d'accroissement du chiffre d'affaires.

.....

Exercice 7. Comparaison de résultats

Nous disposons de trente échantillons dont nous souhaitons déterminer la teneur en un composé chimique donné. Chacun d'entre eux est analysé avec trois méthodes différentes d'analyse chimique. Les résultats obtenus ont été reproduits dans le tableau ci-dessous.

Échantillon	Méthode			Méthode			
	1	2	3	Échantillon	1	2	3
1	133	129	138	16	153	150	152
2	131	132	138	17	125	123	122
3	119	121	121	18	124	120	124
4	124	124	121	19	127	125	124
5	123	124	124	20	136	132	130
6	122	122	123	21	131	130	133
7	127	131	135	22	136	136	133
8	116	116	115	23	123	120	123
9	116	118	122	24	123	117	116
10	104	101	101	25	122	118	121
11	119	117	115	26	101	104	107
12	126	120	121	27	96	97	98
13	96	93	93	28	108	106	108
14	100	97	99	29	124	122	119
15	103	99	102	30	137	136	134

Observons-nous une différence entre les résultats des différentes méthodes d'analyse chimique ?

.....

Table des matières

1	Les tests non paramétriques sur un échantillon	1
1.1	Test du signe	1
1.2	Test des rangs signés de Wilcoxon	3
1.2.1	Cas où il n'y a pas d'ex æquo.	3
1.2.2	Cas où il y a des ex æquo.	4
2	Les tests non paramétriques sur deux échantillons	5
2.1	Les échantillons sont indépendants : Test de Mann-Whitney	5
2.1.1	Cas où il n'y a pas d'ex æquo.	6
2.1.2	Cas où il y a des ex æquo.	7
2.2	Les échantillons sont indépendants : Test de la médiane de Mood	7
2.3	Les échantillons sont dépendants : Test de Wilcoxon	9
2.3.1	Cas où il n'y a pas d'ex æquo.	9
2.3.2	Cas où il y a des ex æquo.	9
3	Les tests non paramétriques sur k échantillons	10
3.1	Les échantillons sont indépendants : Test de Kruskal-Wallis	10
3.1.1	Cas où il n'y a pas d'ex æquo.	10
3.1.2	Cas où il y a des ex æquo.	11
3.2	Les échantillons sont indépendants : Test de Jonckheere-Terpstra	12
3.2.1	Cas où il n'y a pas d'ex æquo.	13
3.2.2	Cas où il y a des ex æquo.	13
3.3	Les échantillons ne sont pas indépendants : Test de Friedman	14
3.3.1	Cas où il n'y a pas d'ex æquo.	14
3.3.2	Cas où il y a des ex æquo.	15
4	Les tests non paramétriques sur nk échantillons	16
4.1	Les échantillons sont indépendants : Test de Friedman	16
4.1.1	Cas où il n'y a pas d'ex æquo.	17
4.1.2	Cas où il y a des ex æquo.	18

Juge	Ordre	Produit	Note	26	1	delisse	st manet	scoup	andros	26	2	andros	26	2	carrefour	poti	4	26	Juge
26	1	delisse	5	26	2	andros	st manet	4	27	1	poti	7	27	2	carrefour	st manet	5	27	Juge
27	3	andros	6	27	4	scoup	3	27	5	5	delisse	7	27	6	carrefour	carrefour	7	27	Juge
28	1	scoup	3	28	2	carrefour	3	28	3	3	st manet	2	28	4	delisse	carrefour	7	28	Juge
28	5	poti	0	28	6	andros	1	29	5	29	st manet	8	29	2	scoup	6	6	29	Juge
29	3	poti	2	29	4	carrefour	6	29	5	29	andros	4	29	6	delisse	7	7	29	Juge
30	1	carrefour	7	30	2	delisse	5	30	3	30	scoup	3	30	4	andros	5	5	30	Juge
30	5	st manet	4	30	6	poti	2	31	1	31	andros	4	31	2	delisse	9	9	31	Juge
31	3	poti	1	31	4	carrefour	4	31	5	31	scoup	0	31	6	st manet	8	8	31	Juge
36	1	st manet	3	36	2	scoup	1	36	3	36	carrefour	2	36	4	poti	0	2	36	Juge
36	5	delisse	1	36	6	andros	6	36	1	39	scoup	7	39	2	poti	2	2	39	Juge
39	3	st manet	3	39	4	carrefour	4	39	5	39	andros	3	39	6	delisse	6	6	39	Juge
41	1	delisse	8	41	2	andros	2	41	3	41	carrefour	5	41	4	st manet	4	4	41	Juge
41	5	poti	10	41	6	scoup	2	42	1	42	carrefour	5	42	2	delisse	5	5	42	Juge
42	3	poti	7	42	4	andros	4	42	5	42	scoup	3	42	6	st manet	3	3	42	Juge
44	1	andros	6	44	2	scoup	6	44	3	44	delisse	3	44	4	st manet	7	7	44	Juge
44	5	poti	2	44	6	carrefour	2	46	1	46	carrefour	4	46	2	poti	6	6	46	Juge
46	3	st manet	0	46	4	delisse	4	46	5	46	scoup	0	46	6	andros	4	4	46	Juge
47	1	poti	8	47	2	delisse	3	47	3	47	carrefour	1	47	4	andros	1	1	47	Juge
47	5	st manet	4	47	6	scoup	6	47	4	47	delisse	2	47	5	andros	4	4	47	Juge