

Mesures d'association.¹

1. Le cas paramétrique sous hypothèse de multinormalité

1.1. Coefficient de corrélation simple

Le coefficient de corrélation simple $\rho(X, Y)$ mesure le degré d'association linéaire entre deux variables aléatoires X et Y .

$$\rho(X, Y) = \text{Cor}(X, Y) = \frac{\text{Cov}[X, Y]}{(\text{Var}[X] \times \text{Var}[Y])^{\frac{1}{2}}}.$$

On a toujours $-1 \leq \rho(X, Y) \leq 1$.

Lorsque $\rho(X, Y) = 0$, on dit que les deux variables aléatoires X et Y sont non corrélées. Attention, ne pas être corrélé ne signifie pas généralement être indépendants. Mais l'indépendance implique toujours l'absence de corrélation.

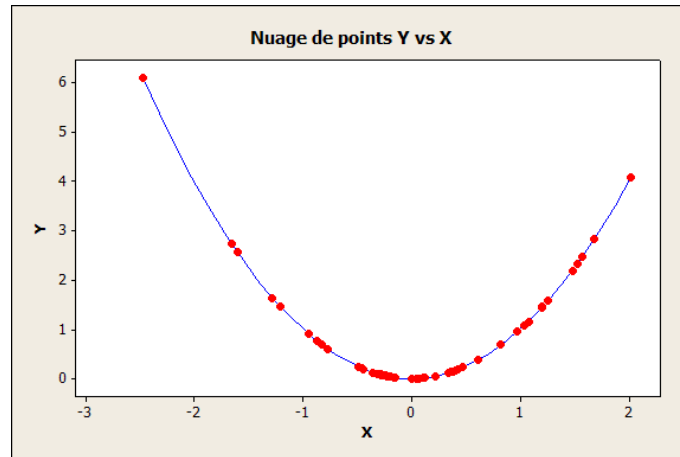
Exemple 1.1.

Prenons $X \sim \mathcal{N}(0, 1)$ et $Y = X^2$, on sait alors que $Y \sim \chi_1^2$ une loi du χ^2 à un degré de liberté. Ainsi les propriétés classiques d'une loi du χ^2 impliquent que $\mathbb{E}[Y] = 1$ et $\text{Var}[Y] = 2$. Remarquer que nous ne sommes donc pas dans le cas qui sera développé plus tard où le couple (X, Y) suit une loi normale bidimensionnelle. On a alors :

$$\begin{aligned} \rho(X, Y) &= \text{Cor}(X, Y) = \frac{\text{Cov}[X, Y]}{(\text{Var}[X] \times \text{Var}[Y])^{\frac{1}{2}}} = \frac{\mathbb{E}[(X - 0)(X^2 - 1)]}{(1 \times 2)^{\frac{1}{2}}} \\ &= \frac{1}{\sqrt{2}} \mathbb{E}[X^3 - X] = \frac{1}{\sqrt{2}} (\mathbb{E}[X^3] - \mathbb{E}[X]) = 0 - 0 = 0. \end{aligned}$$

Car X suit une loi normale centrée donc symétrique par rapport à 0. Pourtant il est clair que les deux variables aléatoires X et $Y = X^2$ ne sont pas indépendantes !

¹Les références [11], [8] et [12] ayant servi à l'élaboration de ce document sont mentionnées dans la bibliographie.



En bleu la relation entre les deux variables aléatoires X et $Y = X^2$, en rouge 50 réalisations du couple (X, X^2) . On calculera au paragraphe 1.2.2 la valeur de la statistique de corrélation à partir de l'échantillon formé des points en rouge.

.....

Il faut donc toujours garder en mémoire que le coefficient de corrélation ne mesure que le **degré d'association linéaire** entre deux variables aléatoires. Nous verrons dans les sections suivantes quels outils utiliser lorsque l'on suspecte que l'association n'est pas nécessairement linéaire.

Lorsque $|\rho(X, Y)| = 1$ alors, avec une probabilité de 1, $Y = aX + b$ pour des constantes réelles a , du même signe que $\rho(X, Y)$, et b .

Le coefficient de corrélation est invariant par transformation linéaire. Si $\text{Cor}(X, Y) = \rho(X, Y)$ et que l'on pose :

$$\begin{aligned} X' &= aX + b, \\ Y' &= cY + d, \end{aligned}$$

alors $\text{Cor}(X', Y') = \rho(X, Y)$. Ceci a une conséquence extrêmement importante : on peut aussi bien travailler avec les variables brutes qu'avec les variables centrées réduites pour l'étude de la corrélation.

1.2. Le cas bidimensionnel

1.2.1. Loi normale bidimensionnelle

Supposons que l'on dispose d'un vecteur aléatoire (X, Y) gaussien. Nous verrons dans la suite, au moment de tester cette hypothèse, qu'un vecteur aléatoire (X, Y) n'est pas

nécessairement gaussien si X et Y suivent des lois normales². L'hypothèse que l'on fait porte sur le couple et est de ce fait plus délicate à tester que simplement tester la normalité de l'échantillon \mathbf{x} issu de X et celle de l'échantillon \mathbf{y} issu de Y . Commençons par rappeler ce que l'on entend par la loi d'un vecteur aléatoire gaussien à deux dimensions.

Définition 1.1. Loi normale bidimensionnelle

La densité de la loi normale à deux dimensions est :

$$f(x, y) = K \exp \left[-\frac{1}{2} P(x, y) \right]$$

où :

$$K = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho(X, Y)^2}},$$

$$P(x, y) = \frac{1}{1-\rho(X, Y)^2} \left[\frac{(x-m_X)^2}{\sigma_X^2} - 2\rho(X, Y) \frac{(x-m_X)(y-m_Y)}{\sigma_X\sigma_Y} + \frac{(y-m_Y)^2}{\sigma_Y^2} \right].$$

Cette loi dépend donc des cinq paramètres : m_X , m_Y , σ_X , σ_Y et $\rho(X, Y)$.

Les quatre premiers sont les moyennes et les écarts types des deux lois marginales, la loi de X et celle de Y , qui sont des lois normales, $X \sim \mathcal{N}(m_X, \sigma_X^2)$ et $Y \sim \mathcal{N}(m_Y, \sigma_Y^2)$. Les lois liées, c'est-à-dire la probabilité d'obtenir une valeur $Y = y$ sachant que l'on a obtenu $X = x$ ou d'obtenir une valeur $X = x$ sachant que l'on a obtenu une valeur $Y = y$ sont également des lois normales dont les paramètres respectifs ($m_{Y|X=x}$, $\sigma_{Y|X=x}$) et ($m_{X|Y=y}$, $\sigma_{X|Y=y}$) s'expriment ainsi :

$$m_{Y|X=x} = m_Y + \rho(X, Y) \frac{\sigma_Y}{\sigma_X} (x - m_X),$$

$$\sigma_{Y|X=x}^2 = \sigma_Y^2 (1 - \rho(X, Y)^2).$$

$$m_{X|Y=y} = m_X + \rho(X, Y) \frac{\sigma_X}{\sigma_Y} (y - m_Y),$$

$$\sigma_{X|Y=y}^2 = \sigma_X^2 (1 - \rho(X, Y)^2).$$

On remarque que :

- La variance de chacune des lois liées est constante : elle ne dépend pas de la valeur de X ou de Y prise en compte pour la calculer, c'est-à-dire du nombre x , qui intervient dans $\sigma_{Y|X=x}$, respectivement y , qui intervient dans $\sigma_{X|Y=y}$, pour lequel elle est calculée.
- La moyenne de chacune des lois liées varie linéairement avec la variable liée. On peut donc parler de deux droites de régression :

$$D_{YX} : y - m_Y = \rho(X, Y) \frac{\sigma_Y}{\sigma_X} (x - m_X),$$

$$D_{XY} : x - m_X = \rho(X, Y) \frac{\sigma_X}{\sigma_Y} (y - m_Y).$$

²Voir l'ouvrage de J.-Y. Ouvrard [10], exercice 13.1 page 276, pour un contre-exemple.

- $\rho(X, Y)$ est le coefficient de corrélation linéaire :

$$\begin{aligned}\rho(X, Y) &= \frac{\text{Cov}[X, Y]}{(\text{Var}[X] \times \text{Var}[Y])^{\frac{1}{2}}} \\ &= \frac{\mathbb{E}[(X - m_X) \times (Y - m_Y)]}{[\mathbb{E}[(X - m_X)^2] \times \mathbb{E}[(Y - m_Y)^2]]^{\frac{1}{2}}}.\end{aligned}$$

$\rho(X, Y) = 0$ est une condition **nécessaire et suffisante** pour que les deux variables aléatoires X et Y soient indépendantes car ici on a une **hypothèse** de normalité bidimensionnelle **supplémentaire**, voir le paragraphe 1.1 pour la situation générale.

Si $\rho(X, Y) = \pm 1$, les deux droites sont confondues et il y a une relation fonctionnelle entre X et Y . Plus précisément, $(X - m_X)$ et $(Y - m_Y)$ sont colinéaires, c'est-à-dire qu'il existe λ et μ deux nombres réels tels que $\lambda(X - m_X) + \mu(Y - m_Y) = 0$. On reconnaît dans la formule précédente l'équation d'une droite. Ainsi si $\rho(X, Y) = \pm 1$, les deux variables X et Y sont liées par une relation fonctionnelle linéaire.

La valeur absolue de $\rho(X, Y)$ est toujours inférieure ou égale à 1. Plus $\rho(X, Y)$ est proche de cette valeur, plus la liaison entre X et Y est serrée. Si l'on considère un échantillon $((x_1, y_1), \dots, (x_n, y_n))$ de réalisations du couple (X, Y) , la dispersion des couples de points (x_i, y_i) autour des droites est d'autant plus faible que les deux droites sont plus proches l'une de l'autre.

- $\rho(X, Y)^2$ est égal au coefficient de détermination R^2 que l'on trouve dans le contexte de la régression linéaire simple.

1.2.2. Estimation des paramètres

On considère un n -échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ indépendant et identiquement distribué suivant la loi de (X, Y) , qui est une loi normale bidimensionnelle.

Les estimateurs des moyennes m_X et m_Y et des écarts types σ_X et σ_Y que nous allons utiliser sont :

$$\begin{aligned}\widehat{m}_X &= \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \\ \widehat{m}_Y &= \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, \\ \widehat{\sigma}_X^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, \\ \widehat{\sigma}_Y^2 &= \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2, \\ \widehat{\sigma}_X &= \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2},\end{aligned}$$

$$\widehat{\sigma}_Y = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2}.$$

Par convention et par soucis d'alléger les notations, on note \bar{X} à la place de \widehat{m}_X . Les quatre premiers estimateurs ci-dessus sont des estimateurs **sans biais** et convergents :

$$\begin{aligned} \mathbb{E}[\widehat{m}_X] &= m_X & \text{Var}[\widehat{m}_X] &= \frac{\sigma_X^2}{n}, \\ \mathbb{E}[\widehat{m}_Y] &= m_Y & \text{Var}[\widehat{m}_Y] &= \frac{\sigma_Y^2}{n}, \\ \mathbb{E}[\widehat{\sigma}_X^2] &= \sigma_X^2 & \text{Var}[\widehat{\sigma}_X^2] &= \frac{2\sigma_X^4}{n-1}, \\ \mathbb{E}[\widehat{\sigma}_Y^2] &= \sigma_Y^2 & \text{Var}[\widehat{\sigma}_Y^2] &= \frac{2\sigma_Y^4}{n-1}. \end{aligned}$$

Les deux derniers sont **biaisés** et convergents :

$$\begin{aligned} \mathbb{E}[\widehat{\sigma}_X] &= \sqrt{\frac{2}{n-1}} \frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right)} \sigma_X = \left(1 - \frac{3}{4n} + o\left(\frac{1}{n}\right)\right) \sigma_X, \\ \text{Var}[\widehat{\sigma}_X] &= \frac{1}{n-1} \left[n-1 - \frac{2\Gamma\left(\frac{n}{2}\right)^2}{\Gamma\left(\frac{n-1}{2}\right)^2} \right] \sigma_X^2, \\ \mathbb{E}[\widehat{\sigma}_Y] &= \sqrt{\frac{2}{n-1}} \frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{n-1}{2}\right)} \sigma_Y = \left(1 - \frac{3}{4n} + o\left(\frac{1}{n}\right)\right) \sigma_Y, \\ \text{Var}[\widehat{\sigma}_Y] &= \frac{1}{n-1} \left[n-1 - \frac{2\Gamma\left(\frac{n}{2}\right)^2}{\Gamma\left(\frac{n-1}{2}\right)^2} \right] \sigma_Y^2 \end{aligned}$$

où $o(1/n)$ est une fonction qui tend vers 0 plus vite que $1/n$ et $\Gamma(x)$ est une fonction mathématique définie à l'aide d'une intégrale !

Il serait alors possible de corriger le biais³ de $\widehat{\sigma}_X$. La situation est beaucoup plus compliquée que pour les estimateurs ci-dessus. On comprend ainsi pourquoi on préfère estimer

³Il suffit de poser $\widehat{\sigma}_X^{\text{corrigé}} = \sqrt{\frac{n-1}{2}} \frac{\Gamma\left(\frac{n-1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \widehat{\sigma}_X$. Néanmoins l'obtention d'intervalle de confiance

la variance d'une loi normale à la place de son écart type. « Heureusement », les deux paramètres qui ont été choisis pour définir une loi normale sont sa moyenne et sa variance. Au passage on remarque que la racine carrée d'un estimateur sans biais n'est pas nécessairement sans biais.

Le coefficient de corrélation linéaire $\rho(X, Y)$ est un rapport. Pour estimer ce rapport nous allons estimer le numérateur $\text{Cov}[X, Y]$ et le dénominateur $(\text{Var}[X] \times \text{Var}[Y])^{\frac{1}{2}}$. En ce qui concerne le dénominateur, nous verrons qu'il suffit d'utiliser les estimateurs $\widehat{\sigma}_X$ et $\widehat{\sigma}_Y$ ci-dessus. Un estimateur du numérateur peut être défini par la quantité :

$$\begin{aligned}\widehat{\text{Cov}}[X, Y] &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\ &= \frac{1}{n-1} \sum_{i=1}^n (X_i \times Y_i) - \bar{X}\bar{Y}.\end{aligned}$$

C'est un estimateur sans biais et convergent de $\text{Cov}[X, Y]$:

$$\mathbb{E}[\widehat{\text{Cov}}[X, Y]] = \text{Cov}[X, Y].$$

Notez la similitude avec l'estimateur corrigé de la variance. Ici aussi le dénominateur est $n-1$.

L'estimateur du coefficient de corrélation $\rho(X, Y)$, que nous allons utiliser, est $\widehat{\rho}(X, Y)$ défini par :

$$\begin{aligned}\widehat{\rho}(X, Y) &= \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}) \times (Y_i - \bar{Y})}{\left(\widehat{\sigma}_X^2 \times \widehat{\sigma}_Y^2\right)^{\frac{1}{2}}} \\ &= \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}) \times (Y_i - \bar{Y})}{\widehat{\sigma}_X \times \widehat{\sigma}_Y} \\ &= \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}) \times (Y_i - \bar{Y})}{\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \times \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2\right)^{\frac{1}{2}}}\end{aligned}$$

reste délicate puisque la loi de cet estimateur n'est pas classique. P. Chapouille, dans son livre [4], propose la formule approximative :

$$\widehat{\sigma}_{X\text{ corrigé}} \approx \widehat{\sigma}_X \sqrt{\frac{2n-2}{2n-3}} \approx \widehat{\sigma}_X \left(1 + \frac{1}{4n-6}\right).$$

$$\begin{aligned}
&= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}) \times (Y_i - \bar{Y}) \\
&= \frac{1}{n-1} \left(\sum_{i=1}^n (X_i - \bar{X})^2 \times \sum_{i=1}^n (Y_i - \bar{Y})^2 \right)^{\frac{1}{2}} \\
&= \frac{\sum_{i=1}^n (X_i - \bar{X}) \times (Y_i - \bar{Y})}{\left(\sum_{i=1}^n (X_i - \bar{X})^2 \times \sum_{i=1}^n (Y_i - \bar{Y})^2 \right)^{\frac{1}{2}}} \\
&= \frac{\sum_{i=1}^n (X_i \times Y_i) - \bar{X} \times \bar{Y}}{\left(\sum_{i=1}^n X_i^2 - \bar{X}^2 \right)^{\frac{1}{2}} \left(\sum_{i=1}^n Y_i^2 - \bar{Y}^2 \right)^{\frac{1}{2}}}.
\end{aligned}$$

L'étude des propriétés de cet estimateur est plus complexe puisqu'il s'agit d'un rapport de variables aléatoires qui ne sont pas indépendantes⁴. Il s'agit d'un estimateur **biaisé** et convergent de $\rho(X, Y)$:

$$\begin{aligned}
\mathbb{E} \left[\widehat{\rho(X, Y)} \right] &= \rho(X, Y) - \frac{\rho(X, Y)(1 - \rho(X, Y)^2)}{2n} \\
\text{Var} \left[\widehat{\rho(X, Y)} \right] &= \frac{(1 - \rho(X, Y)^2)^2}{n} + o\left(\frac{1}{n}\right),
\end{aligned}$$

où $o(1/n)$ est une fonction qui tend vers 0 lorsque n tend vers $+\infty$ plus vite que $1/n$. La distribution asymptotique de $\widehat{\rho(X, Y)}$ est :

$$\sqrt{n}(\widehat{\rho(X, Y)} - \rho(X, Y)) \approx \mathcal{N}\left(0, (1 - \rho(X, Y)^2)^2\right).$$

Ce résultat n'est pas assez précis pour permettre de résoudre les problèmes qui nous intéressent généralement. On l'utilisera uniquement si l'on ne peut pas faire autrement et si $n \geq 30$.

Ainsi lorsque l'on dispose d'un échantillon $(x_1, y_1), \dots, (x_n, y_n)$ de réalisations du couple (X, Y) on obtient les estimations des paramètres de la loi normale bidimensionnelle suivie par (X, Y) en utilisant les formules ci-dessous qui sont les réalisations des estimateurs ci-dessus sur notre échantillon. On note \mathbf{x} l'échantillon (x_1, \dots, x_n) , \mathbf{y} l'échantillon (y_1, \dots, y_n) et (\mathbf{x}, \mathbf{y}) l'échantillon $(x_1, y_1), \dots, (x_n, y_n)$.

⁴Voir la note 6

Les estimations des moyennes m_X et m_Y et des écarts types σ_X et σ_Y sont :

$$\begin{aligned}\widehat{m}_X(\mathbf{x}) &= \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \\ \widehat{m}_Y(\mathbf{y}) &= \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \\ \widehat{\sigma}_X^2(\mathbf{x}) &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \\ \widehat{\sigma}_Y^2(\mathbf{y}) &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2, \\ \widehat{\sigma}_X(\mathbf{x}) &= \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}, \\ \widehat{\sigma}_Y(\mathbf{y}) &= \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}.\end{aligned}$$

On note \bar{x} bien qu'il s'agisse de la moyenne de l'échantillon \mathbf{x} et que de ce fait on pourrait également écrire $\bar{\mathbf{x}}$. Il s'agit là aussi d'une convention mais si l'on voulait être cohérent avec les notations utilisées pour les autres estimateurs on devrait écrire $\widehat{m}_X(\mathbf{x})$.

L'estimation de la covariance $\text{Cov}[X, Y]$ est donnée par :

$$\widehat{\text{Cov}}[X, Y](\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y}).$$

L'estimation du coefficient de corrélation $\rho(X, Y)$ est donnée par :

$$\begin{aligned}\widehat{\rho}(X, Y)(\mathbf{x}, \mathbf{y}) &= \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{\left(\widehat{\sigma}_X^2(\mathbf{x}) \times \widehat{\sigma}_Y^2(\mathbf{y})\right)^{\frac{1}{2}}} \\ &= \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{\widehat{\sigma}_X(\mathbf{x}) \times \widehat{\sigma}_Y(\mathbf{y})} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x}) \times (y_i - \bar{y})}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \times \sum_{i=1}^n (y_i - \bar{y})^2\right)^{\frac{1}{2}}}.\end{aligned}$$

$$= \frac{\sum_{i=1}^n (x_i \times y_i) - \bar{x} \times \bar{y}}{\left(\sum_{i=1}^n x_i^2 - \bar{x}^2 \right)^{\frac{1}{2}} \times \left(\sum_{i=1}^n y_i^2 - \bar{y}^2 \right)^{\frac{1}{2}}}$$

Exemple 1.2.

On évalue la formule ci-dessus pour l'échantillon représenté en rouge au paragraphe 1.1 :

Pearson correlation of Y and X = 0,006

.....

1.2.3. Procédure de test

Sous l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendantes », la variable :

$$\sqrt{n-2} \frac{\widehat{\rho(X, Y)}}{\sqrt{1 - \widehat{\rho(X, Y)}^2}}$$

suit une loi de Student T_{n-2} à $n-2$ degrés de liberté.

Ceci permet de réaliser le test suivant :

$$\boxed{\mathcal{H}_0 : \rho(X, Y) = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho(X, Y) \neq 0.}$$

Pour prendre une décision associée à un test de niveau $(1 - \alpha)$ %, il suffit donc de trouver la valeur critique d'une loi de Student à $n-2$ degrés de liberté pour $\alpha/2$.

Il est également possible, en réalisant un test unilatéral, de tester les hypothèses suivantes :

$$\boxed{\mathcal{H}_0 : \rho(X, Y) = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho(X, Y) > 0.}$$

$$\boxed{\mathcal{H}_0 : \rho(X, Y) = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho(X, Y) < 0.}$$

Notons que, si $\rho_0 \neq 0$, la propriété ci-dessus ne permet pas de tester des hypothèses du type :

$$\boxed{\mathcal{H}_0 : \rho(X, Y) = \rho_0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho(X, Y) \neq \rho_0.}$$

Pour procéder à ces types de test, on utilise la transformation de Fisher⁵. En effet, Fisher, le premier, a démontré que si l'on introduit les deux variables :

$$\begin{aligned} Z(X, Y) &= \frac{1}{2} \ln \left(\frac{1 + \rho(X, Y)}{1 - \rho(X, Y)} \right) \\ &= \tanh^{-1}(\rho(X, Y)), \\ \widehat{Z}(X, Y) &= \frac{1}{2} \ln \left(\frac{1 + \widehat{\rho}(X, Y)}{1 - \widehat{\rho}(X, Y)} \right) \\ &= \tanh^{-1}(\widehat{\rho}(X, Y)), \end{aligned}$$

alors la différence $\widehat{Z}(X, Y) - Z(X, Y)$ suit une loi voisine de la loi normale de moyenne $\rho(X, Y)/(2(n-1))$ et de variance proche de $1/(n-3)$. On approximera donc dans le cas général la loi de cette différence par une loi $\mathcal{N}(\rho(X, Y)/(2(n-1)), 1/(n-3))$ et sous une hypothèse d'indépendance entre X et Y , qui implique donc $\rho(X, Y) = 0$, par une loi $\mathcal{N}(0, 1/(n-3))$.

Pour passer des résultats obtenus pour $Z(X, Y)$ à des résultats concernant $\rho(X, Y)$, on utilise la transformation inverse de celle de Fisher :

$$\begin{aligned} \rho(X, Y) &= \frac{\exp(2Z(X, Y)) - 1}{\exp(2Z(X, Y)) + 1} \\ \rho(X, Y) &= \tanh(Z(X, Y)). \end{aligned}$$

La plupart des logiciels de statistique permettent de calculer les estimations par intervalles de $\rho(X, Y)$. Certains proposent même des intervalles de confiance exacts.

1.2.4. Remarques sur les liaisons

- L'existence de deux droites de régression dans les études de corrélation est embarrassante car elle amène à la question suivante : laquelle des deux droites représente-t-elle le mieux la réalité ? On ne peut trancher sans utiliser la notion de causalité : si l'on pense que X est la cause de Y , on retiendra la droite D_{YX} qui est la régression de la variable dépendante par rapport à la variable cause.

⁵Il existe d'autres transformations pour obtenir des intervalles de confiance pour $\rho(X, Y)$, par exemple l'approximation de Ruben qui semble être plus précise mais encore plus complexe que celle de Fisher. On peut par exemple montrer que, quelque soit la valeur du coefficient de corrélation $\rho(X, Y)$, la variable aléatoire $(R - \rho(X, Y))\sqrt{(n-2)/[(1-R^2)(1-\rho(X, Y)^2)]}$ suit approximativement une loi \mathcal{T}_{n-2} de Student à $n-2$ degrés de liberté. On retrouve le résultat du paragraphe 1.2.3 établi lorsque $\rho(X, Y) = 0$. Pour plus de détails voir [2] et [6].

Dans les cas où l'on ne peut établir une relation de cause à effet, il n'y a pas lieu de fixer son attention sur l'une des droites de régression plutôt que sur l'autre. Le coefficient de corrélation linéaire $\rho(X, Y)$ est alors le paramètre le plus intéressant ; quelque soit sa valeur, positive ou négative, il ne préjuge en rien d'un quelconque rapport de cause à effet entre les variables étudiées.

Dans tous les cas, l'interprétation de toute étude de liaison mérite beaucoup de réflexions, et seules les raisons physiques ou biologiques et non statistiques pourront permettre de porter des jugements de cause à effet.

- Lorsque l'on trouve une courbe de régression qui serre de près le phénomène étudié, il faut se garder de conclure que cette formule traduit la loi exacte qui gouverne le phénomène.
- Enfin signalons que souvent deux variables X et Y sont fortement corrélées avec une troisième Z , le temps par exemple, et on peut conclure à une corrélation significative entre X et Y , alors qu'à priori il n'y a aucune relation entre ces deux grandeurs si ce n'est leur liaison avec Z . On sent alors la nécessité d'introduire une mesure de liaison qui permettra de connaître l'association de X et Y en éliminant l'influence de Z : il s'agit de la notion de corrélation partielle qui sera étudiée aux paragraphes 1.5, 4 et 6.

1.3. Le cas général ($n \geq 2$)

1.3.1. Loi multinormale

On considère un vecteur gaussien de dimension $p \geq 1$: $\mathbf{X} = (X_1, \dots, X_p)$. Un tel vecteur est entièrement déterminé par la connaissance de sa moyenne $\boldsymbol{\mu}(\mathbf{X})$ et de sa matrice de variance-covariance $\boldsymbol{\Sigma}(\mathbf{X})$. En termes plus clairs, il suffit de connaître :

$$\begin{aligned} \boldsymbol{\mu}(\mathbf{X}) &= (\mu_1, \dots, \mu_p) = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_p]), \\ \boldsymbol{\Sigma}(\mathbf{X}) &= ((\sigma_{i,j} = \text{Cov}[X_i, X_j]))_{1 \leq i, j \leq p} = \begin{pmatrix} \text{Var}[X_1] & \dots & \text{Cov}[X_1, X_p] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_p, X_1] & \dots & \text{Var}[X_p] \end{pmatrix}. \end{aligned}$$

Il s'agit d'une extension de ce que nous venons de voir pour le cas bidimensionnel.

On introduit alors la matrice des corrélations $\boldsymbol{\rho}(\mathbf{X})$ entre les composantes X_i , $1 \leq i \leq p$ de \mathbf{X} :

$$\begin{aligned} \boldsymbol{\rho}(\mathbf{X}) &= ((\rho_{i,j} = \text{Cor}[X_i, X_j]))_{1 \leq i, j \leq p} \\ &= \begin{pmatrix} 1 & \dots & \frac{\text{Cov}[X_1, X_p]}{(\text{Var}[X_1] \text{Var}[X_p])^{\frac{1}{2}}} \\ \vdots & \ddots & \vdots \\ \frac{\text{Cov}[X_p, X_1]}{(\text{Var}[X_p] \text{Var}[X_1])^{\frac{1}{2}}} & \dots & 1 \end{pmatrix}. \end{aligned}$$

1.3.2. Estimation

Afin de construire un estimateur simple de la matrice $\boldsymbol{\rho}(\mathbf{X})$ on doit disposer d'un n -échantillon $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ indépendant et identiquement distribué du vecteur \mathbf{X} .

Ce n -échantillon est donc composé de n éléments, qui sont des vecteurs de taille p , que l'on note $(X_{1,1}, \dots, X_{p,1}), \dots, (X_{1,n}, \dots, X_{p,n})$. $X_{i,j}$ est la variable aléatoire associée à la valeur prise par la i -ème composante X_i lors de la réalisation de la j -ème expérience. Par exemple $X_{1,3}$ serait la variable aléatoire associée à la valeur prise par X_1 lors de la troisième expérience et $X_{5,2}$ serait la variable aléatoire associée à la valeur prise par X_5 lors de la seconde expérience.

Un estimateur $\hat{\boldsymbol{\mu}}$ de la moyenne $\boldsymbol{\mu}$ est :

$$\begin{aligned}\hat{\boldsymbol{\mu}} &= \frac{\mathbf{X}_1 + \dots + \mathbf{X}_n}{n} \\ &= \frac{\sum_{i=1}^n \mathbf{X}_i}{n}.\end{aligned}$$

On note, comme d'habitude, $\hat{\boldsymbol{\mu}}$ par $\bar{\boldsymbol{\mu}}$. Il s'agit d'un estimateur sans biais et convergent, comme au paragraphe 1.2.2 :

$$\begin{aligned}\mathbb{E}[\hat{\boldsymbol{\mu}}] &= \mathbb{E}\left[\frac{\mathbf{X}_1 + \dots + \mathbf{X}_n}{n}\right] \\ &= \frac{\mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i\right]}{n} \\ &= \frac{\sum_{i=1}^n \mathbb{E}[\mathbf{X}_i]}{n} \\ &= \frac{\sum_{i=1}^n \boldsymbol{\mu}}{n} \\ &= \boldsymbol{\mu}.\end{aligned}$$

$$\text{Var}[\hat{\boldsymbol{\mu}}] = \frac{\boldsymbol{\Sigma}}{n}.$$

Un estimateur $\widehat{\boldsymbol{\rho}(\mathbf{X})}$ de la matrice des corrélations $\boldsymbol{\rho}(\mathbf{X})$ est alors donné par :

$$\widehat{\boldsymbol{\rho}(\mathbf{X})} = \left(\left(\widehat{\rho}_{i,j} = \widehat{\rho}_{i,j} = \text{Cor}[\widehat{X}_i, \widehat{X}_j] \right) \right)_{1 \leq i, j \leq p}$$

$$= \begin{pmatrix} & & & \frac{1}{n-1} \sum_{i=1}^n (X_{1,i} - \bar{X}_1) \times (X_{p,i} - \bar{X}_p) \\ & 1 & \dots & \frac{\hat{\sigma}_1 \times \hat{\sigma}_p}{\hat{\sigma}_1 \times \hat{\sigma}_p} \\ & \vdots & \ddots & \vdots \\ \frac{1}{n-1} \sum_{i=1}^n (X_{p,i} - \bar{X}_p) \times (X_{1,i} - \bar{X}_1) & & & \\ \frac{\hat{\sigma}_p \times \hat{\sigma}_1}{\hat{\sigma}_p \times \hat{\sigma}_1} & \dots & & 1 \end{pmatrix}$$

On obtient alors la forme explicite suivante :

$$\begin{pmatrix} & & & \frac{\sum_{i=1}^n (X_{1,i} - \bar{X}_1) \times (X_{p,i} - \bar{X}_p)}{\left(\sum_{i=1}^n (X_{1,i} - \bar{X}_1)^2 \times \sum_{i=1}^n (X_{i,p} - \bar{X}_p)^2 \right)^{\frac{1}{2}}} \\ & 1 & \dots & \\ & \vdots & \ddots & \vdots \\ \frac{\sum_{i=1}^n (X_{p,i} - \bar{X}_p) \times (X_{1,i} - \bar{X}_1)}{\left(\sum_{i=1}^n (X_{p,i} - \bar{X}_p)^2 \times \sum_{i=1}^n (X_{1,i} - \bar{X}_1)^2 \right)^{\frac{1}{2}}} & \dots & & 1 \end{pmatrix}.$$

La distribution asymptotique de ρ_{ij} est :

$$\sqrt{n}(\widehat{\rho}_{ij} - \rho_{ij}) \approx \mathcal{N}\left(0, (1 - \rho_{ij}^2)^2\right).$$

Ce résultat n'est pas assez précis pour permettre de résoudre les problèmes qui nous intéressent généralement. On l'utilisera uniquement si l'on ne peut pas faire autrement et si $n \geq 30$.

On a donc une estimation $\widehat{\boldsymbol{\rho}}(\mathbf{X})(\mathbf{x})$ de $\boldsymbol{\rho}(\mathbf{X})$ à l'aide d'un échantillon \mathbf{x} de n réalisations de \mathbf{X} . En notant \mathbf{x}_1 l'échantillon formé par les n réalisations, $(x_{1,1}, \dots, x_{1,n})$, de X_1, \dots , \mathbf{x}_p l'échantillon formé par les n réalisations, $(x_{p,1}, \dots, x_{p,n})$, de X_p on a :

$$\widehat{\boldsymbol{\rho}}(\mathbf{X})(\mathbf{x}) = \left(\left(\widehat{\rho}(\mathbf{X})(\mathbf{x})_{i,j} = \widehat{\rho}_{i,j}(\mathbf{x}_i, \mathbf{x}_j) = \text{Cor}[\widehat{X}_i, \widehat{X}_j](\mathbf{x}_i, \mathbf{x}_j) \right) \right)_{1 \leq i,j \leq p}$$

$$= \begin{pmatrix} 1 & \dots & \frac{1}{n-1} \sum_{i=1}^n (x_{1,i} - \bar{x}_1) \times (x_{p,i} - \bar{x}_p) \\ \vdots & \ddots & \vdots \\ \frac{1}{n-1} \sum_{i=1}^n (x_{p,i} - \bar{x}_p) \times (x_{1,i} - \bar{x}_1) & \dots & 1 \end{pmatrix} \cdot \begin{pmatrix} \hat{\sigma}_1(\mathbf{x}_1) \times \hat{\sigma}_p(\mathbf{x}_p) \\ \vdots \\ \hat{\sigma}_p(\mathbf{x}_p) \times \hat{\sigma}_1(\mathbf{x}_1) \end{pmatrix}.$$

L'estimation $\hat{\rho}(\mathbf{X})(\mathbf{x})$ de $\rho(\mathbf{X})$ est alors :

$$\begin{pmatrix} 1 & \dots & \frac{\sum_{i=1}^n (x_{1,i} - \bar{x}_1) \times (x_{p,i} - \bar{x}_p)}{\left(\sum_{i=1}^n (x_{1,i} - \bar{x}_1)^2 \times \sum_{i=1}^n (x_{i,p} - \bar{x}_p)^2 \right)^{\frac{1}{2}}} \\ \vdots & \ddots & \vdots \\ \frac{\sum_{i=1}^n (x_{p,i} - \bar{x}_p) \times (x_{1,i} - \bar{x}_1)}{\left(\sum_{i=1}^n (x_{p,i} - \bar{x}_p)^2 \times \sum_{i=1}^n (x_{i,1} - \bar{x}_1)^2 \right)^{\frac{1}{2}}} & \dots & 1 \end{pmatrix}.$$

1.3.3. Test de l'hypothèse $\rho_{ij} = 0$

On peut maintenant s'intéresser au test d'indépendance des composantes d'un vecteur gaussien.

$$\boxed{\mathcal{H}_0 : \rho_{ij} = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho_{ij} \neq 0.}$$

Sous l'hypothèse nulle \mathcal{H}_0 , alors :

$$t_{ij} = \sqrt{n-2} \frac{\hat{\rho}_{ij}(\mathbf{x}_i, \mathbf{x}_j)}{\sqrt{1 - \hat{\rho}_{ij}(\mathbf{x}_i, \mathbf{x}_j)^2}}$$

est la réalisation d'une variable aléatoire T qui suit une loi de Student à $n - 2$ degrés de liberté, i.e. $T \sim \mathcal{T}_{n-2}$.

On rejette l'hypothèse nulle \mathcal{H}_0 au niveau $(1 - \alpha) \%$ lorsque $|t_{ij}| > \mathcal{T}_{n-2, \alpha/2}$.

1.3.4. Test de l'hypothèse $\rho_{ij} = \rho_0$, $\rho_0 \neq 0$

Si $\rho_0 \neq 0$, la propriété du paragraphe 1.3.3 ci-dessus ne permet pas de tester des hypothèses du type :

$$\boxed{\mathcal{H}_0 : \rho_{ij} = \rho_0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho_{ij} \neq \rho_0.}$$

.

Il existe alors deux possibilités :

- Si l'on souhaite faire un test de ce type à la main, on utilise la transformation de Fisher introduite ci-dessus, au paragraphe 1.2.3, pour obtenir un intervalle de confiance pour ρ_{ij} de niveau approximatif $(1 - \alpha) \%$:

$$\left[\tanh \left(\hat{z} - \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right), \tanh \left(\hat{z} + \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right) \right]$$

où $\hat{z} = \tanh(\widehat{\rho}_{ij}(\mathbf{x}_i, \mathbf{x}_j))$ et $z_{\alpha/2}$ est le quantile de la loi $\mathcal{N}(0, 1)$.

- Certains logiciels, comme SPSS, R, SAS ou StatXact, proposent des versions *exactes* de l'intervalle de confiance dont une approximation est donnée ci-dessus. En effet la distribution exacte de ρ_{ij} n'est pas une distribution de probabilité classique⁶. Si historiquement, les approximations présentées ci-dessus étaient les seules possibilités auxquelles l'expérimentateur pouvaient recourir pour obtenir des intervalles de confiance pour ρ_{ij} , il est devenu possible avec l'augmentation des capacités de calcul dont vous disposez désormais de déterminer les intervalles de confiance exacts pour ρ_{ij} . Ceci présente l'avantage de ne pas avoir à utiliser des résultats asymptotiques qui peuvent s'avérer incorrects si la taille de l'échantillon est petite.

Exemple 1.3.

Considérons deux échantillons \mathbf{x}_1 et \mathbf{x}_2 .

Procédons au test de l'hypothèse \mathbf{x}_1 est issu d'une variable aléatoire X_1 qui suit une loi normale et au test de l'hypothèse \mathbf{x}_2 est issu d'une variable aléatoire X_2 qui suit une loi normale. Compte tenu des effectifs on procède au test de **Shapiro-Wilk** ou plutôt de sa variante disponible dans Minitab : le test de **Ryan-Joiner**.

\mathbf{x}_1	\mathbf{x}_2
0,19577	-0,81685
0,92204	0,58257
-0,04690	1,24359
0,45863	0,31088
-0,58454	0,95124
1,51722	0,02028
-0,97862	-0,91823
-0,04557	-0,18670
-0,03707	-0,96033
0,84505	1,11031

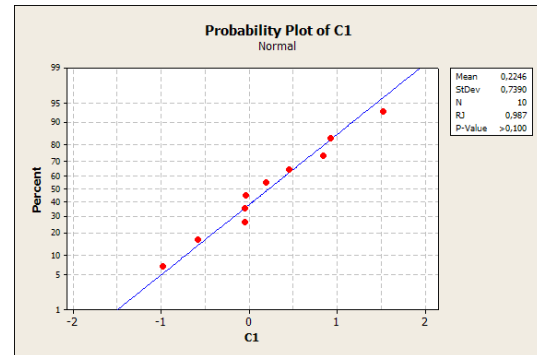
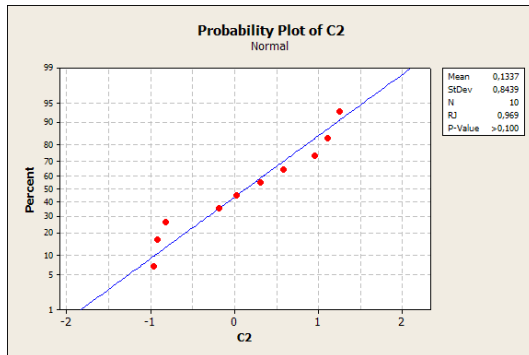
⁶La densité du coefficient de corrélation $\rho(X, Y)$ d'un couple gaussien (X, Y) au point $-1 \leq r \leq 1$ est $\frac{1}{\pi}(n-2)(1-r^2)^{(n-4)/2}(1-\rho^2)^{(n-1)/2} \int_0^{+\infty} \frac{d\beta}{(\cosh \beta - \rho r)^{n-1}}$, voir [14] pour plus de détails.

On commence par tester la normalité du couple (X_1, X_2) .

$\mathcal{H}_0 : (X_1, X_2)$ suit une loi normale bidimensionnelle

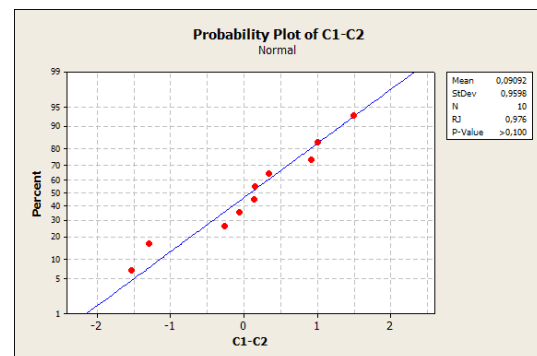
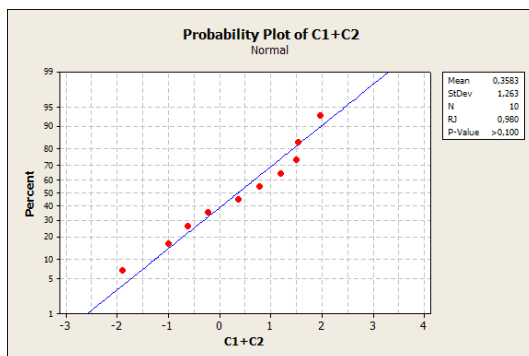
contre

$\mathcal{H}_1 : (X_1, X_2)$ ne suit pas une loi bidimensionnelle.



Les p -valeurs sont toutes les deux supérieures à 0,1 : on ne peut donc pas rejeter l'hypothèse de normalité pour X_1 et X_2 .

Notons que l'on a uniquement vérifié la normalité de X_1 et celle de X_2 . Ceci n'implique pas celle du couple (X_1, X_2) . Il s'agit néanmoins d'une **condition nécessaire**. Nous verrons au paragraphe 9 comment vérifier l'hypothèse de normalité du couple (X_1, X_2) . Une des premières choses à faire est de tester la normalité de $X_1 + X_2$ et de $X_1 - X_2$, voir [5].



Les p -valeurs sont supérieures à 0,100, on ne peut donc pas rejeter l'hypothèse de normalité de $X_1 + X_2$ et de $X_1 - X_2$.

Attention, nous testons ici la multinormalité de (X_1, X_2) et nous avons procédé à quatre tests de Shapiro-Wilk. Nous avons, comme d'habitude, fixé un seuil de $\alpha_{global} = 5\%$ pour le test global de multinormalité. Pour garantir ce risque global de 5%, il faut fixer un risque différent de 5% pour chaque test de Shapiro-Wilk, vous avez déjà rencontré ce problème dans le contexte des comparaisons multiples en analyse de la variance. En effet si le risque pour chacun des tests de Shapiro-Wilk est fixé à $\alpha_{ind} = 5\%$, alors le

risque global est au maximum de $\alpha_{global} = 1 - (1 - 0,05)^4 = 0,185 > 0,05$. Ainsi on doit fixer un seuil de $\alpha_{ind} = 1 - \sqrt[4]{1 - \alpha_{global}} = 1 - \sqrt[4]{0,95} = 0,013$ pour chacun des tests de Shapiro-Wilk, ce qui revient à être moins exigeant pour chacun des tests individuels puisque $0,013 < 0,05$ et donc à accepter la normalité dans plus de cas.

Dans cet exemple cette correction ne change rien à la conclusion puisqu'aucune des p -valeurs n'étaient comprises entre 0,013 et 0,05.

Ainsi on ne peut pas rejeter l'hypothèse nulle \mathcal{H}_0 de normalité du couple (X_1, X_2) .

Testons maintenant l'hypothèse d'indépendance de X_1 et X_2 . Puisque le couple (X_1, X_2) suit une loi normale, l'indépendance de X_1 et X_2 est équivalente à l'absence de corrélation linéaire entre X_1 et X_2 . Il suffit donc de tester :

$$\boxed{\mathcal{H}_0 : \rho(X_1, X_2) = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho(X_1, X_2) \neq 0.}$$

Correlations: X1; X2

Pearson correlation of X1 and X2 = 0,270

P-Value = 0,450

Avec SPSS 13.0 on obtient :

Correlations

		X1	X2
X1	Pearson Correlation	1	,270
	Sig. (2-tailed)		,450
	N	10	10
X2	Pearson Correlation	,270	1
	Sig. (2-tailed)	,450	
	N	10	10

En utilisant le module Tests Exactes de SPSS 13.0 :

	Value	Approx. Sig.	Exact Sig.
Pearson's R	0,270	0,450	0,446
N of Valid Cases	10		

Notons que dans notre cas la significativité exacte est proche de la significativité approchée. Néanmoins, dans de nombreux cas il n'en va pas de même.

Ainsi avec une p -valeur de 0,446, le test n'est pas significatif au seuil $\alpha = 5\%$. On ne peut rejeter l'hypothèse nulle \mathcal{H}_0 d'absence de corrélation entre X_1 et X_2 . On en déduit que l'on ne peut rejeter l'hypothèse d'indépendance de X_1 et X_2 .

.....

1.4. Corrélation multiple

1.4.1. Définition

Le coefficient de corrélation multiple R est la corrélation maximale possible entre une variable réelle X_1 et toutes les combinaisons linéaires de composantes d'un vecteur aléatoire \mathbf{X}_2 .

Rappelons rapidement ce que l'on entend par **combinaison linéaire**.

Si \mathbf{U} et \mathbf{V} sont deux vecteurs de \mathbb{R}^k , $k \geq 1$, une combinaison linéaire de \mathbf{U} et \mathbf{V} est un vecteur $\text{CL}_{(\mathbf{U}, \mathbf{V})}(\alpha, \beta)$ défini comme la somme $\alpha\mathbf{U} + \beta\mathbf{V}$ avec $(\alpha, \beta) \in \mathbb{R}^2$.

Cette notion s'étend au cas de n vecteurs : si $\mathbf{U}_1, \dots, \mathbf{U}_n$ sont n vecteurs de \mathbb{R}^k , $k \geq 1$, une combinaison linéaire des $(\mathbf{U}_i)_{1 \leq i \leq n}$ est un vecteur $\text{CL}_{(\mathbf{U}_1, \dots, \mathbf{U}_n)}(\alpha_1, \dots, \alpha_n)$ défini comme la somme $\sum_{i=1}^n \alpha_i \mathbf{U}_i$ avec $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$.

Mise en équation : on dispose de $X_1 \in \mathbb{R}$ et $\mathbf{X}_2 \in \mathbb{R}^{p-1}$ tels que $(X_1, \mathbf{X}_2) \in \mathbb{R}^p$ ait une distribution dans \mathbb{R}^p de moyenne $\boldsymbol{\mu}$ et de matrice de variance-covariance $\boldsymbol{\Sigma}$.

La moyenne $\boldsymbol{\mu}$ du couple (X_1, \mathbf{X}_2) est reliée à la moyenne μ_1 de X_1 et $\boldsymbol{\mu}_2$ de \mathbf{X}_2 par :

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}.$$

On adopte une écriture particulière pour $\boldsymbol{\Sigma}$ pour faire ressortir la variance de X_1 , notée $\sigma_{11}(X_1)$, la matrice de variance-covariance de \mathbf{X}_2 , notée $\boldsymbol{\Sigma}_{22}(\mathbf{X}_2)$, et les covariances des composantes de \mathbf{X}_2 avec X_1 , qui sont les composantes du vecteur $\boldsymbol{\Sigma}_{21}(X_1, \mathbf{X}_2)$.

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11}(X_1) & \boldsymbol{\sigma}_{21}(X_1, \mathbf{X}_2)^T \\ \boldsymbol{\Sigma}_{21}(X_1, \mathbf{X}_2) & \boldsymbol{\Sigma}_{22}(\mathbf{X}_2) \end{pmatrix}$$

où $\boldsymbol{\sigma}_{21}(X_1, \mathbf{X}_2)^T$ est la transposée de $\boldsymbol{\sigma}_{21}(X_1, \mathbf{X}_2)$. Remarquons que $\text{Var}[X_1] = \sigma_{11}(X_1) = \text{Cov}[X_1, X_1] = \sigma_{X_1}^2$ avec les notations du paragraphe 1.2.2.

On montre par un peu de calcul matriciel que la corrélation maximale R vérifie :

$$R^2(X_1, \mathbf{X}_2) = \frac{\boldsymbol{\sigma}_{21}(X_1, \mathbf{X}_2)^T \boldsymbol{\Sigma}_{22}(\mathbf{X}_2)^{-1} \boldsymbol{\sigma}_{21}(X_1, \mathbf{X}_2)}{\sigma_{11}(X_1)}.$$

On peut vérifier qu'il s'agit bien d'un nombre réel.

Le cas de la corrélation simple est légèrement différent de celui-ci puisqu'ici on ne peut trouver facilement que la valeur absolue de $R(X_1, \mathbf{X}_2)$. A nouveau on doit faire le lien entre le coefficient de détermination en régression linéaire multiple et $R^2(X_1, \mathbf{X}_2)$: ils ont la même valeur.

1.4.2. Estimation

On se donne, comme précédemment, un n -échantillon indépendant et identiquement distribué $((X_{1,1}, \mathbf{X}_{2,1}), \dots, (X_{1,n}, \mathbf{X}_{2,n}))$. En utilisant les estimateurs définis au paragraphe 1.3.2, un estimateur du coefficient $R^2(X_1, \mathbf{X}_2)$ est donné par la formule suivante :

$$\begin{aligned} R^2(\widehat{X_1}, \widehat{\mathbf{X}_2}) &= \frac{\widehat{\boldsymbol{\sigma}}_{21}(\widehat{X_1}, \widehat{\mathbf{X}_2})^T \widehat{\boldsymbol{\Sigma}}_{22}(\widehat{\mathbf{X}_2})^{-1} \widehat{\boldsymbol{\sigma}}_{21}(\widehat{X_1}, \widehat{\mathbf{X}_2})}{\widehat{\sigma}_{11}(X_1)} \\ &= \frac{\widehat{\boldsymbol{\sigma}}_{21}(\widehat{X_1}, \widehat{\mathbf{X}_2})^T \widehat{\boldsymbol{\Sigma}}_{22}(\widehat{\mathbf{X}_2})^{-1} \widehat{\boldsymbol{\sigma}}_{21}(\widehat{X_1}, \widehat{\mathbf{X}_2})}{\widehat{\sigma}_{11}(X_1)}. \end{aligned}$$

On note désormais \mathbf{x}_1 un échantillon de n réalisations de X_1 et \mathbf{x}_2 un échantillon de n réalisations de \mathbf{X}_2 .

Une estimation de $R^2(X_1, \mathbf{X}_2)$ est alors :

$$R^2(\widehat{X_1}, \widehat{\mathbf{X}_2})(\mathbf{x}_1, \mathbf{x}_2) = \frac{\left(\widehat{\boldsymbol{\sigma}}_{21}(\widehat{X_1}, \widehat{\mathbf{X}_2})(\mathbf{x}_1, \mathbf{x}_2) \right)^T \left(\widehat{\boldsymbol{\Sigma}}_{22}(\widehat{\mathbf{X}_2})(\mathbf{x}_2) \right)^{-1} \widehat{\boldsymbol{\sigma}}_{21}(\widehat{X_1}, \widehat{\mathbf{X}_2})(\mathbf{x}_1, \mathbf{x}_2)}{\widehat{\sigma}_{11}(X_1)(\mathbf{x}_1)}.$$

1.4.3. Asymptotique

Sans hypothèse supplémentaire on ne pourrait rien dire⁷ sur la distribution de $R^2(\widehat{X_1}, \widehat{\mathbf{X}_2})$ et on ne pourrait donc pas effectuer de test.

Supposons désormais que nos $X_1 \in \mathbb{R}$ et $\mathbf{X}_2 \in \mathbb{R}^{p-1}$ soient tels que $(X_1, \mathbf{X}_2) \in \mathbb{R}^p$ ait une distribution **multinormale** dans \mathbb{R}^p de moyenne $\boldsymbol{\mu}$ et de matrice de variance-covariance $\boldsymbol{\Sigma}$. En conservant les notations du paragraphe 1.4.1 on peut écrire :

$$(X_1, \mathbf{X}_2) \sim \mathcal{N} \left(\begin{pmatrix} \mu_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \boldsymbol{\sigma}_{21}^T \\ \boldsymbol{\sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \right).$$

La distribution asymptotique de $R(X_1, \mathbf{X}_2)$ est alors identique à celle du coefficient de corrélation simple :

$$\sqrt{n}(R(\widehat{X_1}, \widehat{\mathbf{X}_2}) - R(X_1, \mathbf{X}_2)) \approx \mathcal{N} \left(0, (1 - R(X_1, \mathbf{X}_2)^2)^2 \right).$$

Ce résultat n'est pas assez précis pour permettre de résoudre les problèmes qui nous intéressent généralement. On l'utilisera uniquement si l'on ne peut pas faire autrement et si $n \geq 30$.

⁷Le théorème central limite permet d'avoir des informations sur la loi limite de $R(\widehat{X_1}, \widehat{\mathbf{X}_2})$ mais ces renseignements ne sont valables que pour des effectifs très importants, $n \rightarrow +\infty$ et de ce fait ne présentent que peu d'intérêt dans la pratique.

1.4.4. Test de l'hypothèse $R(X_1, \mathbf{X}_2) = 0$

Toujours sous l'hypothèse de multinormalité du vecteur (X_1, \mathbf{X}_2) comme précisé au paragraphe 1.4.3, on peut connaître la loi exacte d'une fonction de l'estimateur $\widehat{R}^2(X_1, \mathbf{X}_2)$ sous l'hypothèse nulle « \mathcal{H}_0 : Le coefficient de corrélation multiple $R(X_1, \mathbf{X}_2)$ de X_1 et \mathbf{X}_2 est nul » :

$$\frac{n-p-1}{p} \frac{R^2(\widehat{X}_1, \widehat{\mathbf{X}}_2)}{1 - R^2(\widehat{X}_1, \widehat{\mathbf{X}}_2)} \sim \mathcal{F}_{p, n-p-1}$$

où $\mathcal{F}_{p, n-p-1}$ est une loi de Fisher à p et $n-p-1$ degrés de liberté⁸.

On se sert de cette propriété pour tester les hypothèses :

\mathcal{H}_0 : Le coefficient de corrélation multiple $R(X_1, \mathbf{X}_2)$ de X_1 et \mathbf{X}_2 est nul

contre

\mathcal{H}_1 : Le coefficient de corrélation multiple $R(X_1, \mathbf{X}_2)$ de X_1 et \mathbf{X}_2 est non nul.

Ces hypothèses sont équivalentes aux hypothèses :

\mathcal{H}_0 : X_1 et \mathbf{X}_2 sont indépendants

contre

\mathcal{H}_1 : X_1 et \mathbf{X}_2 sont liés.

Pour obtenir des intervalles de confiance pour $R(X_1, \mathbf{X}_2)$ on peut aussi utiliser la transformation en « argument tangente hyperbolique », i.e. en \tanh^{-1} , voir le paragraphe 1.2.3.

1.4.5. Test de l'hypothèse $R(X_1, \mathbf{X}_2) = R_0, R_0 \neq 0$

On se place encore sous l'hypothèse de multinormalité du vecteur (X_1, \mathbf{X}_2) comme précisé au paragraphe 1.4.3. Notons que la propriété du paragraphe 1.4.4 ci-dessus ne permet pas de tester, si $R_0 \neq 0$, des hypothèses du type :

⁸Si $p = 1$, on retrouve exactement le cas de la corrélation simple. En effet dans cette situation, on étudie la corrélation entre deux variables X_1 et X_2 à valeurs réelles. La formule ci-dessus pour $R^2(X_1, X_2)$ est alors identique à celle de $\rho(X_1, X_2)^2$. De plus, on sait que $\frac{\sqrt{n-2}\rho(\widehat{X}_1, \widehat{X}_2)}{\sqrt{1-\rho(\widehat{X}_1, \widehat{X}_2)^2}}$ suit une loi de Student à $n-2$ degrés de liberté. On montre qu'alors son carré $\left(\sqrt{n-2} \frac{\rho(\widehat{X}_1, \widehat{X}_2)}{\sqrt{1-\rho(\widehat{X}_1, \widehat{X}_2)^2}}\right)^2 = (n-2) \frac{\rho(\widehat{X}_1, \widehat{X}_2)^2}{1-\rho(\widehat{X}_1, \widehat{X}_2)^2} = (n-2) \frac{R^2(\widehat{X}_1, \widehat{X}_2)}{1-R^2(\widehat{X}_1, \widehat{X}_2)}$ suit une loi de Fisher à 1 et $n-2$ degrés de liberté. Ceci est bien conforme au résultat de ce paragraphe puisque si $p = 1$ on a vu que $\frac{n-p-1}{p} \frac{R^2(\widehat{X}_1, \widehat{\mathbf{X}}_2)}{1-R^2(\widehat{X}_1, \widehat{\mathbf{X}}_2)} = (n-2) \frac{R^2(\widehat{X}_1, \widehat{X}_2)}{1-R^2(\widehat{X}_1, \widehat{X}_2)}$ suit une loi de Fisher à p et $n-p-1$, c'est-à-dire 1 et $n-2$, degrés de liberté. Réciproquement, si $(n-2) \frac{R^2(\widehat{X}_1, \widehat{X}_2)}{1-R^2(\widehat{X}_1, \widehat{X}_2)}$ suit une loi de Fisher à 1 et $n-2$, on sait qu'alors $\sqrt{(n-2) \frac{R^2(\widehat{X}_1, \widehat{X}_2)}{1-R^2(\widehat{X}_1, \widehat{X}_2)}} = \frac{\sqrt{n-2}\rho(\widehat{X}_1, \widehat{X}_2)}{\sqrt{1-\rho(\widehat{X}_1, \widehat{X}_2)^2}}$ suit une loi de Student à $n-2$ degrés de liberté.

$$\boxed{\mathcal{H}_0 : R(X_1, \mathbf{X}_2) = R_0}$$

contre

$$\boxed{\mathcal{H}_1 : R(X_1, \mathbf{X}_2) \neq R_0.}$$

La démarche est alors la même que pour le coefficient de corrélation simple, voir paragraphe 1.3.4. Deux approches sont possibles :

- Utiliser la transformation de Fisher introduite ci-dessus, au paragraphe 1.2.3, pour obtenir un intervalle de confiance pour $R(X_1, \mathbf{X}_2)$ de niveau approximatif $(1 - \alpha) \%$:

$$\left[\tanh \left(\hat{z} - \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right), \tanh \left(\hat{z} + \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right) \right]$$

où $\hat{z} = \tanh \left(R(\widehat{X}_1, \widehat{\mathbf{X}}_2)(\mathbf{x}_1, \mathbf{x}_2) \right)$ et $z_{\alpha/2}$ est le quantile de la loi $\mathcal{N}(0, 1)$.

- Utiliser des logiciels, comme SPSS, R, SAS ou StatXact, qui proposent des versions *exactes* de l'intervalle de confiance dont une approximation est donnée ci-dessus. Cette approche est particulièrement intéressante si l'effectif commun n aux échantillons \mathbf{x}_1 et \mathbf{x}_2 est petit.

1.5. Corrélation partielle

Pour introduire formellement le concept de corrélation partielle, un petit détour mathématique est nécessaire. Les applications seront généralement beaucoup plus simples.

1.5.1. Définition

Considérons deux vecteurs $\mathbf{X}_1 \in \mathbb{R}^q$ et $\mathbf{X}_2 \in \mathbb{R}^{m-q}$ qui ont une distribution conjointe multinormale :

$$(\mathbf{X}_1, \mathbf{X}_2) \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \right).$$

La distribution conditionnelle de \mathbf{X}_1 sachant \mathbf{X}_2 , on étend les résultats du paragraphe 1.1 qui concernaient le cas de deux variables réelles au cas de deux vecteurs, est :

$$\mathbf{X}_1 | \mathbf{X}_2 \sim \mathcal{N} \left(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{X}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11|2} \right),$$

où $\boldsymbol{\Sigma}_{11|2} = \boldsymbol{\Sigma}_{1,1} - \boldsymbol{\Sigma}_{1,2} \boldsymbol{\Sigma}_{2,2}^{-1} \boldsymbol{\Sigma}_{2,1}$.

En écrivant $\boldsymbol{\Sigma}_{11|2} = ((\sigma_{i,j|q+1,\dots,m}))_{1 \leq i,j \leq p}$, il s'agit d'une matrice carrée de taille p , on

définit le coefficient de corrélation partielle entre les composantes i et j de \mathbf{X}_1 sachant \mathbf{X}_2 comme :

$$\rho_{i,j|q+1,\dots,m} = \frac{\sigma_{i,j|q+1,\dots,m}}{\sigma_{i,i|q+1,\dots,m}^{\frac{1}{2}} \sigma_{j,j|q+1,\dots,m}^{\frac{1}{2}}}.$$

Son interprétation pratique est la suivante :

Le coefficient de corrélation partielle entre les composantes i et j du vecteur \mathbf{X}_1 sachant \mathbf{X}_2 représente la corrélation entre les composantes i et j après avoir éliminé l'effet des variables de \mathbf{X}_2 sur les composantes i et j du vecteur \mathbf{X}_1 .

1.5.2. Estimation

On se donne, comme précédemment, un n -échantillon indépendant et identiquement distribué $((\mathbf{X}_{1,1}, \mathbf{X}_{2,1}), \dots, (\mathbf{X}_{1,n}, \mathbf{X}_{2,n}))$. En utilisant les estimateurs définis au paragraphe 1.3.2, un estimateur du coefficient $\rho_{i,j|q+1,\dots,m}$ est donné par la formule suivante :

$$\widehat{\rho_{i,j|q+1,\dots,m}} = \frac{\widehat{\sigma_{i,j|q+1,\dots,m}}}{\widehat{\sigma_{i,i|q+1,\dots,m}}^{\frac{1}{2}} \widehat{\sigma_{j,j|q+1,\dots,m}}^{\frac{1}{2}}},$$

où $((\widehat{\sigma_{i,j|q+1,\dots,m}}))_{1 \leq i,j \leq p} = \widehat{\Sigma_{11|2}} = \widehat{\Sigma_{1,1}} - \widehat{\Sigma_{1,2}} \widehat{\Sigma_{2,2}}^{-1} \widehat{\Sigma_{2,1}}$.

On note désormais \mathbf{x}_1 un échantillon de n réalisations de \mathbf{X}_1 et \mathbf{x}_2 un échantillon de n réalisations de \mathbf{X}_2 .

Une estimation de $\rho_{i,j|q+1,\dots,m}$ est alors :

$$\widehat{\rho_{i,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2) = \frac{\widehat{\sigma_{i,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2)}{\widehat{\sigma_{i,i|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2)^{\frac{1}{2}} \widehat{\sigma_{j,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2)^{\frac{1}{2}}},$$

où

$$\begin{aligned} ((\widehat{\sigma_{i,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2)))_{1 \leq i,j \leq p} &= \widehat{\Sigma_{11|2}}(\mathbf{x}_1, \mathbf{x}_2) \\ &= \widehat{\Sigma_{1,1}}(\mathbf{x}_1, \mathbf{x}_2) - \widehat{\Sigma_{1,2}}(\mathbf{x}_1, \mathbf{x}_2) \left(\widehat{\Sigma_{2,2}}(\mathbf{x}_1, \mathbf{x}_2) \right)^{-1} \widehat{\Sigma_{2,1}}(\mathbf{x}_1, \mathbf{x}_2). \end{aligned}$$

1.5.3. Asymptotique

La distribution asymptotique d'un coefficient de corrélation partielle est la même que pour une corrélation simple, c'est-à-dire :

$$\sqrt{n} (\widehat{\rho_{i,j|q+1,\dots,m}} - \rho_{i,j|q+1,\dots,m}) \approx \mathcal{N}(0, (1 - \rho_{i,j|q+1,\dots,m}^2)^2).$$

Ce résultat n'est applicable que si $n \geq 30$. Ne s'en servir que si l'on ne peut pas faire autrement.

1.5.4. Test de l'hypothèse $\rho_{i,j|q+1,\dots,m} = 0$

On peut maintenant s'intéresser au test de nullité du coefficient de corrélation partielle des composantes du vecteur gaussien \mathbf{X}_1 connaissant le vecteur \mathbf{X}_2 .

$$\boxed{\mathcal{H}_0 : \rho_{i,j|q+1,\dots,m} = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho_{i,j|q+1,\dots,m} \neq 0.}$$

Ces hypothèses sont équivalentes à :

$$\boxed{\mathcal{H}_0 : \text{Les composantes } X_{1,i} \text{ et } X_{1,j} \text{ de } \mathbf{X}_1 \text{ sont indépendantes sans l'effet de } \mathbf{X}_2}$$

contre

$$\boxed{\mathcal{H}_1 : \text{Les composantes } X_{1,i} \text{ et } X_{1,j} \text{ de } \mathbf{X}_1 \text{ sont liées sans l'effet de } \mathbf{X}_2.}$$

Sous l'hypothèse nulle \mathcal{H}_0 , alors :

$$t_{i,j|q+1,\dots,m}(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{n-m+q-2} \frac{\widehat{\rho_{i,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2)}{\sqrt{1 - \widehat{\rho_{i,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2)^2}}$$

est la réalisation d'une variable aléatoire T qui suit une loi de Student à $n-m+q-2$ degrés de liberté, i.e. $T \sim \mathcal{T}_{n-m+q-2}$.

On rejette l'hypothèse nulle \mathcal{H}_0 au niveau $(1-\alpha)\%$ lorsque $|t_{i,j|q+1,\dots,m}(\mathbf{x}_1, \mathbf{x}_2)| > t_{n-m+q-2, \alpha/2}$.

On peut aussi appliquer la transformation z dans ce contexte, voir le paragraphe 1.2.3.

1.5.5. Test de l'hypothèse $\rho_{i,j|q+1,\dots,m} = \rho_0$, $\rho_0 \neq 0$

Si $\rho_0 \neq 0$, la propriété du paragraphe 1.5.4 ci-dessus ne permet pas de tester des hypothèses du type :

$$\boxed{\mathcal{H}_0 : \rho_{i,j|q+1,\dots,m} = \rho_0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho_{i,j|q+1,\dots,m} \neq \rho_0.}$$

Il existe alors deux possibilités :

- Si l'on souhaite faire un test de ce type à la main, on utilise la transformation de Fisher introduite ci-dessus, au paragraphe 1.2.3, pour obtenir un intervalle de confiance pour $\rho_{i,j|q+1,\dots,m}$ de niveau approximatif $(1-\alpha)\%$:

$$\left[\tanh \left(\widehat{z} - \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right), \tanh \left(\widehat{z} + \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right) \right]$$

où $\widehat{z} = \tanh(\widehat{\rho_{i,j|q+1,\dots,m}}(\mathbf{x}_1, \mathbf{x}_2))$ et $z_{\alpha/2}$ est le quantile de la loi $\mathcal{N}(0, 1)$.

- Certains logiciels, comme SPSS, R, SAS ou StatXact, proposent des versions *exactes* de l'intervalle de confiance.

1.5.6. Cas de trois variables réelles

Intéressons-nous au cas le plus simple. Soient X_1, X_2 et X_3 trois variables aléatoires réelles telles que la loi jointe de (X_1, X_2, X_3) soit multinormale de paramètres $\boldsymbol{\mu}$ et $\boldsymbol{\Sigma}$. Ainsi on a ici $m = 3$ et $q = 2$. On montre alors la relation suivante :

$$\rho_{12|3} = \frac{\rho_{12} - \rho_{13}\rho_{23}}{\sqrt{1 - \rho_{13}^2}\sqrt{1 - \rho_{23}^2}}$$

où ρ_{12} est le coefficient de corrélation simple entre les variables X_1 et X_2 , ρ_{13} est le coefficient de corrélation simple entre les variables X_1 et X_3 et ρ_{23} est le coefficient de corrélation simple entre les variables X_2 et X_3 , voir le paragraphe 1.3.1.

Dans la plupart des situations expérimentales que vous rencontrerez cette formule sera suffisante. On remarque également que la définition est symétrique en X_1 et X_2 , c'est-à-dire :

$$\rho_{12|3} = \frac{\rho_{12} - \rho_{13}\rho_{23}}{\sqrt{1 - \rho_{13}^2}\sqrt{1 - \rho_{23}^2}} = \frac{\rho_{21} - \rho_{23}\rho_{13}}{\sqrt{1 - \rho_{23}^2}\sqrt{1 - \rho_{13}^2}} = \rho_{21|3}.$$

La corrélation partielle de X_1 et X_2 sachant X_3 est heureusement la même que celle de X_2 et X_1 sachant X_3 .

On considère un échantillon \mathbf{x} indépendant et identiquement distribué suivant la loi de X_1, X_2, X_3 . On note \mathbf{x}_1 l'échantillon des réalisations de X_1 , \mathbf{x}_2 l'échantillon des réalisations de X_2 et \mathbf{x}_3 l'échantillon des réalisations de X_3 .

On tire de \mathbf{x} une estimation de $\rho_{12|3}$ en calculant des estimations de $\rho_{12}, \rho_{13}, \rho_{23}$ comme expliqué au paragraphe 1.3.2.

On peut maintenant s'intéresser au test de nullité du coefficient de corrélation partielle de X_1 et X_2 connaissant la variable X_3 .

$$\boxed{\mathcal{H}_0 : \rho_{12|3} = 0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho_{12|3} \neq 0.}$$

Ces hypothèses sont équivalente à :

$$\boxed{\mathcal{H}_0 : \text{Les variables } X_1 \text{ et } X_2 \text{ sont indépendantes sans l'effet de } X_3}$$

contre

$$\boxed{\mathcal{H}_1 : \text{Les variables } X_1 \text{ et } X_2 \text{ sont liées sans l'effet de } X_3.}$$

Sous l'hypothèse nulle \mathcal{H}_0 , alors :

$$t_{12|3}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) = \sqrt{n-3} \frac{\widehat{\rho}_{12|3}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)}{\sqrt{1 - \widehat{\rho}_{12|3}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)^2}}$$

est la réalisation d'une variable aléatoire T qui suit une loi de Student à $n-3$ degrés de liberté, i.e. $T \sim \mathcal{T}_{n-3}$.

On rejette l'hypothèse nulle \mathcal{H}_0 au niveau $(1-\alpha)\%$ lorsque $|t_{12|3}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)| > t_{n-3, \alpha/2}$. On peut aussi appliquer la transformation z dans ce contexte, voir le paragraphe 1.2.3.

Si $\rho_0 \neq 0$, la propriété du paragraphe ci-dessus ne permet pas de tester des hypothèses du type :

$$\boxed{\mathcal{H}_0 : \rho_{12|3} = \rho_0}$$

contre

$$\boxed{\mathcal{H}_1 : \rho_{12|3} \neq \rho_0.}$$

Il existe alors deux possibilités :

- Si l'on souhaite faire un test de ce type à la main, on utilise la transformation de Fisher introduite ci-dessus, au paragraphe 1.2.3, pour obtenir un intervalle de confiance pour $\rho_{12|3}$ de niveau approximatif $(1-\alpha)\%$:

$$\left[\tanh \left(\widehat{z} - \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right), \tanh \left(\widehat{z} + \frac{z_{\alpha/2}}{(n-3)^{\frac{1}{2}}} \right) \right]$$

où $\widehat{z} = \tanh(\widehat{\rho}_{12|3}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3))$ et $z_{\alpha/2}$ est le quantile de la loi $\mathcal{N}(0, 1)$.

- Certains logiciels, comme SPSS, R, SAS ou StatXact, proposent des versions *exactes* de l'intervalle de confiance.

Exemple 1.4.

Détaillons le traitement de la corrélation partielle dans la situation suivante où l'on dispose de trois variables continues réelles.

- Le BMI X_1 .
- Le poids X_2 .
- La taille X_3 .

On rappelle que le BMI d'un individu est un indice qui se calcule à partir de la taille et du poids par la formule suivante :

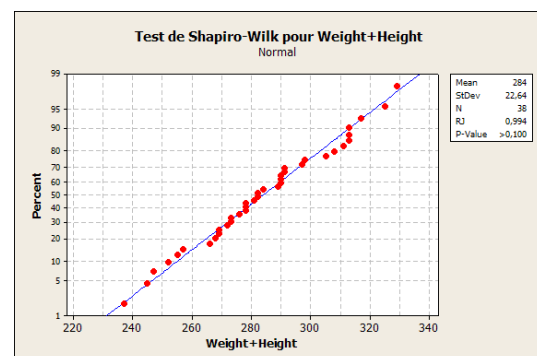
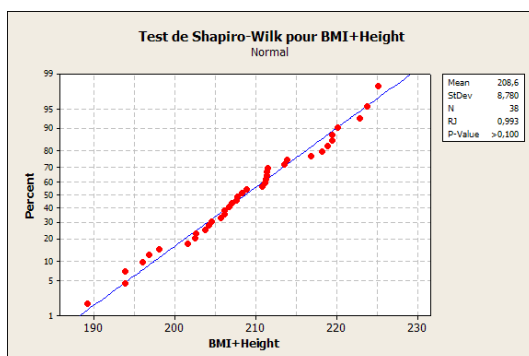
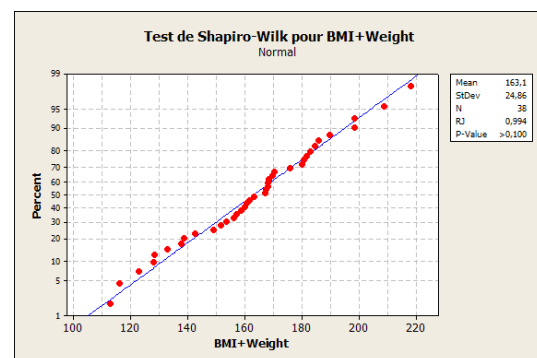
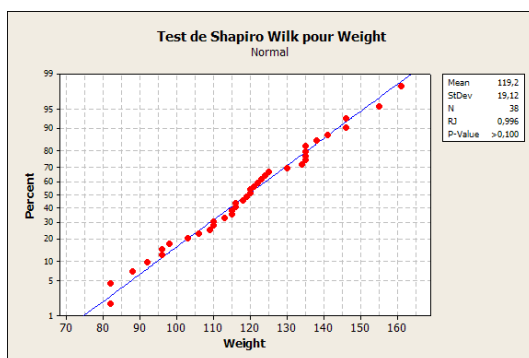
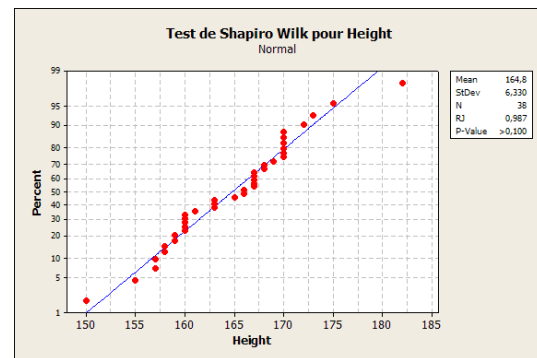
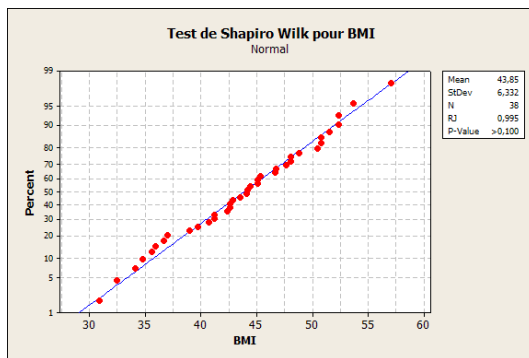
$$\text{BMI} = \frac{\text{Weight}}{\text{Height}^2}$$

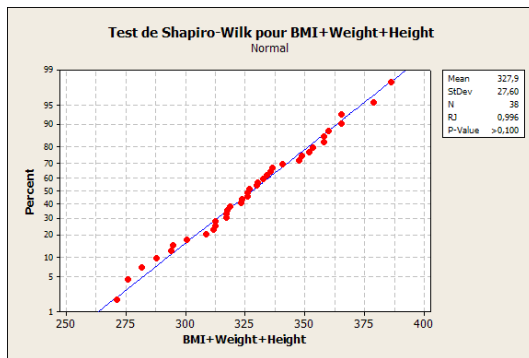
où la masse s'exprime en *kg* et la taille en *m*.

L'échantillon étudié a un effectif $n = 38$ et est exclusivement constitué d'individus de sexe féminin.

Comme nous l'avons fait remarquer, l'hypothèse nécessaire à l'approche paramétrique de la corrélation simple, multiple ou partielle qui vous a été présentée plus haut sont que la loi du vecteur (X_1, X_2, X_3) est une loi multinormale sur \mathbb{R}^3 .

Ceci a pour conséquence qu'il faut que X_1 , X_2 et X_3 aient une distribution normale mais comme nous l'avons déjà souligné ce n'est pas suffisant. On doit par exemple aussi tester la normalité de $X_1 + X_2$, $X_2 + X_3$, $X_1 + X_3$ et $X_1 + X_2 + X_3$. Se référer au paragraphe 9 pour un exposé détaillé des tests de multinormalité ou au livre de R. Christensen [5].





Les p -valeurs sont toutes supérieures à 0,100, on ne peut donc pas rejeter l'hypothèse de normalité de X_1 , X_2 , X_3 , $X_1 + X_2$, $X_1 + X_3$, $X_2 + X_3$ et de $X_1 + X_2 + X_3$.

Attention, nous testions ici la multinormalité de (X_1, X_2, X_3) et nous avons procédé à sept tests de Shapiro-Wilk. Nous avons, comme d'habitude, fixé un seuil de $\alpha_{global} = 5\%$ pour le test global de multinormalité. Pour garantir ce risque global de 5%, il faut fixer un risque différent de 5% pour chaque test de Shapiro-Wilk, nous avons déjà évoqué ce problème à l'exemple 1.3. En effet si le risque pour chacun des tests de Shapiro-Wilk est fixé à $\alpha_{ind} = 5\%$, alors le risque global est au maximum de $\alpha_{global} = 1 - (1 - 0,05)^7 = 0,302 \gg 0,05$. Ainsi on doit fixer un seuil de $\alpha_{ind} = 1 - \sqrt[7]{1 - \alpha_{global}} = 1 - \sqrt[7]{0,95} = 0,007$ pour chacun des sept tests de Shapiro-Wilk, ce qui revient à être moins exigeant pour chacun des tests individuels puisque $0,007 < 0,05$ et donc à accepter la normalité dans plus de cas.

Dans cet exemple cette correction ne change rien à la conclusion puisqu'aucune des p -valeurs n'étaient comprises entre 0,007 et 0,05.

Suite à cette analyse sommaire rien ne permet de rejeter l'hypothèse de multinormalité de (X_1, X_2, X_3) .

On peut donc utiliser une approche paramétrique. Commençons par étudier les coefficients de corrélation simple, dits de Pearson.

Correlations: BMI; Weight; Height

	BMI	Weight
Weight	0,878	0,000

Height	-0,038	0,441
	0,820	0,006

Cell Contents: Pearson correlation
P-Value

On constate ainsi que :

- Au niveau $\alpha = 5 \%$, on rejette l'hypothèse nulle \mathcal{H}_0 d'indépendance des variables BMI et Weight, car $0,000 < 0,05$.
- Au niveau $\alpha = 5 \%$, on ne peut pas rejeter l'hypothèse nulle \mathcal{H}_0 d'indépendance des variables BMI et Height, car $0,820 > 0,05$.
- Au niveau $\alpha = 5 \%$, on rejette l'hypothèse nulle \mathcal{H}_0 d'indépendance des variables Height et Weight, car $0,006 < 0,05$.

Ces résultats sont en accord avec la formule permettant de calculer le BMI à l'aide de la taille et de la masse.

- La liaison est linéaire et positive pour la masse, ce que l'on retrouve bien ici avec la valeur de $\rho_{BW} = 0,878$ et une p -valeur de $0,000$.
- La liaison est négative pour la taille mais, vu sa faiblesse, $\rho_{BH} = -0,038$, on ne peut pas rejeter l'hypothèse d'indépendance entre ces deux variables, p -valeur de $0,820$. La formule ci-dessus implique une **relation non linéaire** entre le BMI et la taille c'est pourquoi le coefficient de corrélation linéaire ne peut la mettre en évidence.
- La liaison positive entre la taille et la masse, $\rho_{WH} = 0,441$ et p -valeur de $0,006$, est en accord avec notre connaissance a priori de l'existence d'une relation linéaire positive entre la taille et la masse d'un individu.

Passons maintenant au calcul de la corrélation multiple R du BMI sur (Height, Weight), de la corrélation multiple R du Height sur (BMI, Weight) et de la corrélation multiple R du Weight sur (Height, BMI).

Pour effectuer ce calcul on utilise le fait que R^2 est égal au coefficient de détermination obtenu lorsque l'on fait la régression linéaire multiple de BMI sur (Height, Weight). On obtient par exemple avec Minitab :

Regression Analysis: BMI versus Weight; Height

The regression equation is

BMI = 87,1 + 0,368 Weight - 0,529 Height

Predictor	Coef	SE Coef	T	P
Constant	87,090	1,899	45,85	0,000
Weight	0,367955	0,004146	88,75	0,000
Height	-0,52855	0,01252	-42,21	0,000

S = 0,432656 R-Sq = 99,6% R-Sq(adj) = 99,5%

Pour calculer les autres coefficients de corrélation multiple, on procède de même en changeant les rôles de BMI, Height et Weight.

Regression Analysis: Weight versus BMI; Height

The regression equation is

$$\text{Weight} = -236 + 2,71 \text{ BMI} + 1,44 \text{ Height}$$

Predictor	Coef	SE Coef	T	P
Constant	-236,083	5,253	-44,95	0,000
BMI	2,70570	0,03049	88,75	0,000
Height	1,43599	0,03050	47,09	0,000

$$S = 1,17324 \quad R\text{-Sq} = 99,6\% \quad R\text{-Sq}(\text{adj}) = 99,6\%$$

Regression Analysis: Height versus BMI; Weight

The regression equation is

$$\text{Height} = 164 - 1,86 \text{ BMI} + 0,686 \text{ Weight}$$

Predictor	Coef	SE Coef	T	P
Constant	164,436	0,933	176,17	0,000
BMI	-1,85552	0,04396	-42,21	0,000
Weight	0,68556	0,01456	47,09	0,000

$$S = 0,810650 \quad R\text{-Sq} = 98,4\% \quad R\text{-Sq}(\text{adj}) = 98,4\%$$

On a ainsi trouvé successivement :

$$\rho_{B/s(H,W)} = 99,6 \%$$

$$\rho_{W/s(B,H)} = 99,6 \%$$

$$\rho_{H/s(B,W)} = 98,4 \%$$

Calculons désormais les trois corrélations partielles $\rho_{(\text{Height,Weight})|\text{BMI}}$, $\rho_{(\text{Weight,BMI})|\text{Height}}$ et $\rho_{(\text{Height,BMI})|\text{Weight}}$.

On applique la formule ci-dessus :

$$\rho_{HW|B} = \frac{\rho_{HW} - \rho_{HB}\rho_{WB}}{\sqrt{1 - \rho_{HB}^2}\sqrt{1 - \rho_{WB}^2}} \approx 0,992$$

$$\rho_{WB|H} = \frac{\rho_{WB} - \rho_{WH}\rho_{BH}}{\sqrt{1 - \rho_{WH}^2}\sqrt{1 - \rho_{BH}^2}} \approx 0,998$$

$$\rho_{HB|W} = \frac{\rho_{HB} - \rho_{HW}\rho_{BW}}{\sqrt{1 - \rho_{HW}^2}\sqrt{1 - \rho_{BW}^2}} \approx -0,990.$$

Ces résultats sont à comparer avec les valeurs de corrélation simple. En particulier la corrélation entre le poids et la taille en éliminant l'effet du poids semble significative. Pour aller plus loin et décider de la significativité de ces corrélations, on peut se référer à une table ou utiliser la loi du coefficient de corrélation partielle énoncée aux paragraphes 1.5.3 et 1.5.6.

Certains logiciels, comme SPSS, nous renseignent directement sur les p -valeurs associées aux corrélations partielles :

Partial Correlations

Control Variables			Weight	Height
BMI	Weight	Correlation	1,000	,992
		Significance (2-tailed)	.	,000
		df	0	35
	Height	Correlation	,992	1,000
		Significance (2-tailed)	,000	.
		df	35	0

Partial Correlations

Control Variables			Weight	BMI
Height	Weight	Correlation	1,000	,998
		Significance (2-tailed)	.	,000
		df	0	35
	BMI	Correlation	,998	1,000
		Significance (2-tailed)	,000	.
		df	35	0

Partial Correlations

Control Variables			BMI	Height
Weight	BMI	Correlation	1,000	-,990
		Significance (2-tailed)	.	,000
		df	0	35
	Height	Correlation	-,990	1,000
		Significance (2-tailed)	,000	.
		df	35	0

On constate ainsi que les trois tests sont significatifs au seuil $\alpha = 5 \%$.

On ne peut conserver l'hypothèse d'indépendance de deux variables sachant la troisième et ce dans les trois cas qui se présentent ici : (Height, Weight)|BMI, (Weight, BMI)|Height et (Height, BMI)|Weight.

Avez-vous une critique à formuler quant au modèle que nous avons utilisé ici ? Par exemple l'hypothèse de multinormalité est-elle vraisemblable ? Bien que les tests utilisés ne l'infirmement pas, notre connaissance de la relation entre le BMI et le couple (Height, Weight) ne devait-elle pas nous faire exclure ce modèle ?

.....

Exemple 1.5.

On réalise dans un collège deux tests d'évaluation communs à tous les élèves quel que soit leur âge mais exclusivement de sexe masculin. L'un porte sur leur compétence en mathématiques et l'autre en sport, le temps mis à parcourir 100 m en course à pied. Les performances ont été notées sur 100 dans les deux tests puis ramenées à des notes sur 20. On dispose donc de deux variables aléatoires :

- l'une appelée Math et associée à la note en mathématiques qu'a reçu un élève,
- l'autre appelée Sport et associée à la note en sport qu'a reçu un élève.

On les abrégera parfois en M et S.

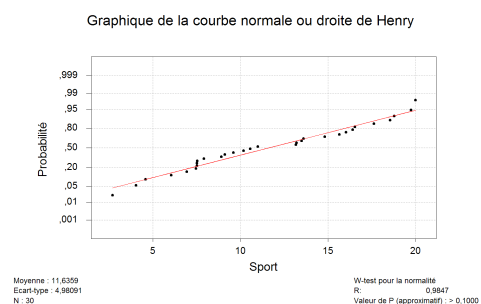
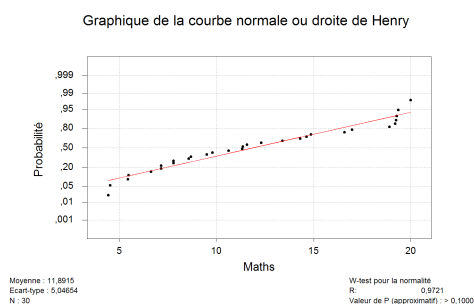
Les données ont été reportées dans le tableau suivant :

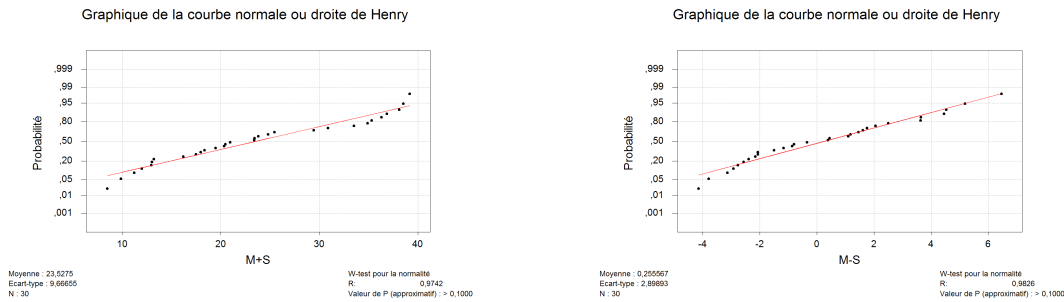
Élève	Math	Sport	Élève	Math	Sport
1	7,15	2,69	16	12,29	13,16
2	7,79	10,56	17	11,56	7,92
3	8,57	8,91	18	14,31	9,11
4	7,14	6,04	19	11,33	13,49
5	5,47	7,53	20	9,49	10,99
6	4,43	4,04	21	16,99	16,54
7	4,52	7,45	22	14,86	16,03
8	6,63	4,58	23	18,91	16,41
9	5,43	7,50	24	19,35	14,81
10	7,78	10,18	25	19,24	17,63
11	9,79	13,59	26	14,62	18,78
12	13,40	6,92	27	20,00	18,54
13	8,68	7,51	28	19,19	20,00
14	11,35	9,59	29	19,28	15,64
15	10,62	13,19	30	16,59	19,73

L'échantillon étudié a un effectif $n = 30$ et est exclusivement constitué d'individus de sexe masculin.

Comme nous l'avons fait remarquer, l'hypothèse nécessaire à l'approche paramétrique de la corrélation simple, multiple ou partielle qui vous a été présentée plus haut est que la loi du vecteur (M, S) est une loi multinormale sur \mathbb{R}^2 .

Ceci a pour conséquence qu'il faut que M, S aient une distribution normale mais comme nous l'avons déjà souligné ce n'est pas suffisant. On doit par exemple aussi tester la normalité de $M + S$ et $M - S$. Se référer au paragraphe 9 pour un exposé détaillé des tests de multinormalité ou au livre de R. Christensen [5].





Les p -valeurs sont toutes supérieures à 0,100, on ne peut donc pas rejeter l'hypothèse de normalité de M , S , $M + S$ et $M - S$.

Attention, nous testions ici la multinormalité de (M, S) et nous avons procédé à quatre tests de Shapiro-Wilk. Attention, nous testions ici la multinormalité de (X, Y) et nous avons procédé à quatre tests de Shapiro-Wilk. Nous avons, comme d'habitude, fixé un seuil de $\alpha_{global} = 5\%$ pour le test global de multinormalité. Pour garantir ce risque global de 5%, il faut fixer un risque différent de 5% pour chaque test de Shapiro-Wilk. En effet si le risque pour chacun des tests de Shapiro-Wilk est fixé à $\alpha_{ind} = 5\%$, alors le risque global est au maximum de $\alpha_{global} = 1 - (1 - 0,05)^4 = 0,185 > 0,05$. Ainsi on doit fixer un seuil de $\alpha_{ind} = 1 - \sqrt[4]{1 - \alpha_{global}} = 1 - \sqrt[4]{0,95} = 0,013$ pour chacun des tests de Shapiro-Wilk, ce qui revient à être moins exigeant pour chacun des tests individuels puisque $0,013 < 0,05$ et donc à accepter la normalité dans plus de cas.

Dans cet exemple cette correction ne change rien à la conclusion puisqu'aucune des p -valeurs n'étaient comprises entre 0,013 et 0,05. Suite à cette analyse sommaire rien ne permet de rejeter l'hypothèse de multinormalité de (M, S) .

On peut donc utiliser une approche paramétrique et utiliser le coefficient de corrélation de M et S .

Corrélation de Pearson de Maths et Sport = 0,833

Valeur de $p = 0,000$

On constate ainsi qu'au seuil $\alpha = 5\%$, on rejette l'hypothèse nulle \mathcal{H}_0 d'indépendance des variables Math et Sport, car $0,000 < 0,05$. On décide donc qu'il y a une corrélation significative entre les résultats en mathématiques et en sport. La valeur de l'estimation, 0,833, du coefficient de corrélation étant positive, l'association est positive, c'est-à-dire, plus on est fort en sport plus on est fort en mathématiques. Qu'en pensez-vous? **De quel autre facteur faudrait-il tenir compte?**

Il est évident que les résultats des élèves à ces tests dépend de leur âge puisque le même questionnaire est posé à tous les élèves quelque soit leur niveau, donc aussi bien en classe de 6^{ème} qu'en classe de 3^{ème}.

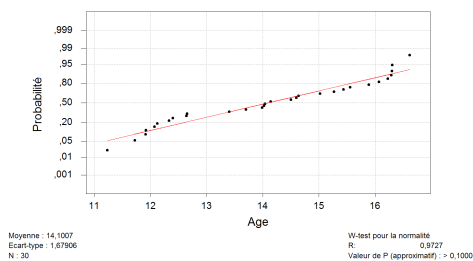
Fort heureusement il a été possible de retrouver l'âge, et même une information encore plus précise, la date de naissance, de chacun des élèves qui a participé à l'évaluation. On a alors décidé de coder l'âge comme une variable aléatoire quantitative continue, appelée Age et parfois abrégée en A. Cette décision correspond à l'idée, relativement raisonnable, suivante : le développement d'un enfant ne serait pas exactement le même si celui-ci est né en début ou en fin d'année civile.

Élève	Math	Sport	Age	Élève	Math	Sport	Age
1	7,15	2,69	11,72	16	12,29	13,16	14,02
2	7,79	10,56	12,33	17	11,56	7,92	15,02
3	8,57	8,91	11,91	18	14,31	9,11	13,40
4	7,14	6,04	11,92	19	11,33	13,49	13,70
5	5,47	7,53	12,64	20	9,49	10,99	14,14
6	4,43	4,04	12,07	21	16,99	16,54	15,56
7	4,52	7,45	12,12	22	14,86	16,03	16,31
8	6,63	4,58	11,23	23	18,91	16,41	16,30
9	5,43	7,50	12,65	24	19,35	14,81	16,23
10	7,78	10,18	12,40	25	19,24	17,63	15,89
11	9,79	13,59	13,99	26	14,62	18,78	15,44
12	13,40	6,92	14,50	27	20,00	18,54	16,29
13	8,68	7,51	14,04	28	19,19	20,00	16,62
14	11,35	9,59	14,64	29	19,28	15,64	16,07
15	10,62	13,19	14,60	30	16,59	19,73	15,27

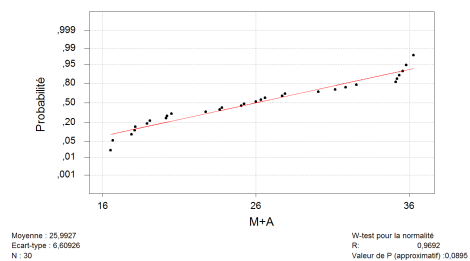
Puisque nous introduisons une nouvelle variable aléatoire Age, pour étudier les corrélations du vecteur (Math, Sport, Age) nous devons tester la normalité du vecteur (Math, Sport, Age).

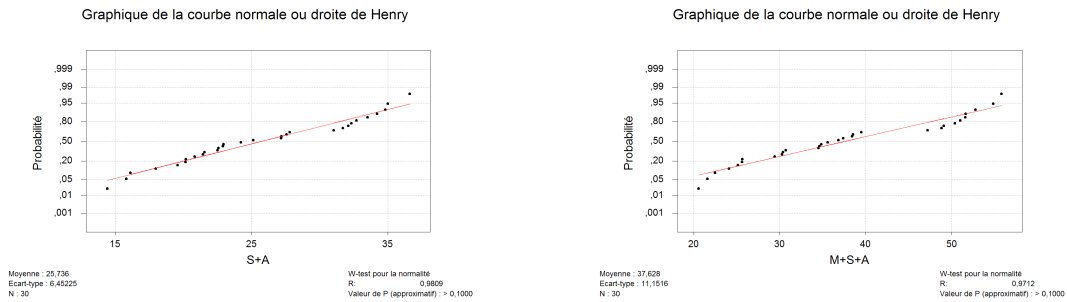
On procède comme précédemment. En s'appuyant sur les résultats des tests de normalité ci-dessus, on voit qu'il ne reste plus qu'à s'intéresser à la normalité de Math + Age, Sport + Age et de Math + Sport + Age.

Graphique de la courbe normale ou droite de Henry



Graphique de la courbe normale ou droite de Henry





Aucun de ces tests n'est significatif au seuil 5 % et donc **a fortiori** ils ne seront pas significatifs au seuil inférieur qu'il faudrait fixer pour garantir un risque global de 5 %. En ajoutant ces résultats aux précédents, rien ne permet de rejeter l'hypothèse de normalité du vecteur (Math, Sport, Age).

On calcule donc les corrélations simples entre les composantes de ce vecteur :

Corrélations : Maths; Sport; Age

	Maths	Sport
Sport	0,833	0,000
Age	0,909	0,837
	0,000	0,000

Contenu de la cellule : corrélation de Pearson
Valeur de p

On note une corrélation positive et significative, au seuil $\alpha = 5\%$ pour chacune des associations possibles. Ainsi, comme on pouvait le prévoir, plus l'on est âgé, mieux l'on réussit en mathématiques et en sport.

Exerçons-nous encore au calcul de coefficients de corrélation multiple. Ainsi déterminons la corrélation multiple R du Math sur (Sport, Age), la corrélation multiple R du Sport sur (Math, Age) et la corrélation multiple R du Age sur (Math, Sport).

Pour effectuer ce calcul on utilise le fait que R^2 est égal au coefficient de détermination obtenu lorsque l'on fait la régression linéaire multiple de Math sur (Sport, Age). On obtient par exemple avec Minitab :

Analyse de régression : Maths en fonction de Sport; Age

L'équation de régression est

$$\text{Maths} = -20,9 + 0,244 \text{ Sport} + 2,13 \text{ Age}$$

Régresseur	Coef	Er-T coef	T	P
Constante	-20,919	4,635	-4,51	0,000
Sport	0,2437	0,1412	1,73	0,096
Age	2,1258	0,4189	5,07	0,000

$$S = 2,071 \quad R\text{-carré} = 84,3\% \quad R\text{-carré (ajust)} = 83,2\%$$

Pour calculer les autres coefficients de corrélation multiple, on procède de même en changeant les rôles de Math, Sport et Age.

Analyse de régression : Sport en fonction de Maths; Age

L'équation de régression est

$$\text{Sport} = -12,5 + 0,408 \text{ Maths} + 1,37 \text{ Age}$$

Régresseur	Coef	Er-T coef	T	P
Constante	-12,533	7,565	-1,66	0,109
Maths	0,4077	0,2362	1,73	0,096
Age	1,3701	0,7099	1,93	0,064

$$S = 2,679 \quad R\text{-carré} = 73,1\% \quad R\text{-carré (ajust)} = 71,1\%$$

Analyse de régression : Age en fonction de Maths; Sport

L'équation de régression est

$$\text{Age} = 10,3 + 0,230 \text{ Maths} + 0,0885 \text{ Sport}$$

Régresseur	Coef	Er-T coef	T	P
Constante	10,3401	0,3338	30,97	0,000
Maths	0,22965	0,04525	5,07	0,000
Sport	0,08848	0,04585	1,93	0,064

$$S = 0,6807 \quad R\text{-carré} = 84,7\% \quad R\text{-carré (ajust)} = 83,6\%$$

On a ainsi trouvé successivement :

$$\rho_{M/s(S,A)} = 84,3 \%$$

$$\rho_{S/s(M,A)} = 73,1 \%$$

$$\rho_{A/s(M,S)} = 84,7 \%$$

Calculons désormais les trois corrélations partielles $\rho_{(\text{Maths,Sport})|\text{Age}}$, $\rho_{(\text{Maths,Age})|\text{Sport}}$ et $\rho_{(\text{Sport,Age})|\text{Maths}}$.

On applique la formule ci-dessus :

$$\rho_{MS|A} = \frac{\rho_{MS} - \rho_{MA}\rho_{SA}}{\sqrt{1 - \rho_{MA}^2}\sqrt{1 - \rho_{SA}^2}} \approx 0,315$$

$$\rho_{MA|S} = \frac{\rho_{MA} - \rho_{MS}\rho_{SA}}{\sqrt{1 - \rho_{MS}^2}\sqrt{1 - \rho_{SA}^2}} \approx 0,699$$

$$\rho_{SA|M} = \frac{\rho_{SA} - \rho_{MS}\rho_{MA}}{\sqrt{1 - \rho_{MS}^2}\sqrt{1 - \rho_{MA}^2}} \approx 0,348.$$

Ces résultats sont à comparer avec les valeurs de corrélation simple. En particulier la corrélation entre le poids et la taille en éliminant l'effet du poids semble significative. Pour aller plus loin et décider de la significativité de ces corrélations on peut se référer à une table ou utiliser la loi du coefficient de corrélation partielle énoncée aux paragraphes 1.5.3 et 1.5.6.

Certains logiciels, comme SPSS, nous renseignent directement sur les p -valeurs associées aux corrélations partielles :

Correlations

Control Variables				
Age	Maths	Correlation	1,000	,315
		Significance (2-tailed)	.	,096
		df	0	27
	Sport	Correlation	,315	1,000
		Significance (2-tailed)	,096	.
		df	27	0

Correlations

Sport				
Age	Correlation	1,000	,699	
	Significance (2-tailed)	.	,000	
	df	0	27	

	Maths	Correlation	,699	1,000
		Significance (2-tailed)	,000	.
		df	27	0

Correlations

Control Variables			Age	Sport
Maths	Age	Correlation	1,000	,348
		Significance (2-tailed)	.	,064
		df	0	27
	Sport	Correlation	,348	1,000
		Significance (2-tailed)	,064	.
		df	27	0

On constate ainsi que les trois tests ne sont pas significatifs au seuil $\alpha = 5\%$.

On conserve l'hypothèse d'indépendance de deux variables sachant la troisième et ce dans les trois cas qui se présentent ici : (Maths, Sport)|Age, (Maths, Age)|Sport et (Sport, Age)|Maths.

On constate que les résultats à l'épreuve de sport et à celle de mathématiques sont indépendants si l'on tient compte de l'âge des élèves, ce qui est bien plus conforme à ce que l'on pouvait penser avant de réaliser cette expérience.

.....

2. Mesures non paramétriques

Les techniques développées dans ce paragraphe permettent de tester l'indépendance ou de mesurer le degré de dépendance entre deux variables aléatoires appariées continues quelconques. Certains de ces tests permettent de faire apparaître des relations en tendance entre les variables et non pas des relations de type fonctionnel comme c'est le cas avec les techniques présentées précédemment. En effet celles-ci se placent systématiquement dans un cadre d'étude d'une liaison de type linéaire entre les différentes variables mises en jeu.

Ainsi le non paramétrique apparaît à deux niveaux : en premier lieu dans le fait que l'on n'aura pas à vérifier d'hypothèse sur la loi des variables étudiées à partir du moment où celles-ci sont quantitatives, continues ou discrètes, ou qualitatives dans certains cas, et en deuxième lieu dans le fait que l'on ne fait d'hypothèse sur le type de relation pouvant exister entre les variables étudiées.

Comme à l'accoutumée, on utilisera, lorsque cela est possible des tests paramétriques conjointement aux tests non paramétriques.

3. La statistique $\rho_{S,n}$ de Spearman

Le coefficient de corrélation de Spearman est basé sur l'étude de la corrélation des rangs.

3.1. Cadre d'application

La mesure de la dépendance au sens de Spearman entre deux variables aléatoires continues X et Y sera notée $\rho_S(X, Y)$. On note F la fonction de répartition de X et G celle de Y . $\rho_S(X, Y)$ pourra être estimé par la statistique de Spearman $\rho_{S,n}(X, Y)$ définie sur un n -échantillon $((X_1, Y_1), \dots, (X_n, Y_n))$ indépendant et identiquement distribué suivant la loi de (X, Y) .

La statistique $\rho_{S,n}(X, Y)$ permet de réaliser plusieurs tests :

$\mathcal{H}_0 : X$ et Y sont indépendantes

contre

$\mathcal{H}_1 : X$ et Y sont liées.

$\mathcal{H}_0 : X$ et Y sont indépendantes

contre

$\mathcal{H}_1 :$ Les valeurs prises par X et Y ont tendance à être concordantes.

$\mathcal{H}_0 : X$ et Y sont indépendantes

contre

$\mathcal{H}_1 :$ Les valeurs prises par X et Y ont tendance à être discordantes.

On se reportera au paragraphe 5.2 pour une définition des termes concordants et discordants.

Ce coefficient a le même champ d'application que la statistique de Kendall $\tau_n(X, Y)$ qu'on lui préférera généralement.

3.2. Le coefficient de corrélation $\rho_S(X, Y)$ de Spearman

Le coefficient de corrélation de Spearman est un nombre associé à (X, Y) qui sert à mesurer le degré de dépendance qui lie X et Y . Il peut être défini comme étant le coefficient de corrélation simple du couple aléatoire $(F(X), G(Y))$:

$$\rho_S(X, Y) = \rho(F(X), G(Y)).$$

Il possède les propriétés suivantes :

- $-1 \leq \rho_S(X, Y) \leq 1$,
- « X et Y indépendantes » implique $\rho_S(X, Y) = 0$,
- $\rho_S(X, Y) = 1$ (resp. $\rho_S(X, Y) = -1$) si et seulement si il existe une fonction ϕ croissante (resp. décroissante) de \mathbb{R} dans \mathbb{R} telle que $Y = \phi(X)$,
- Si ϕ et ψ désignent deux fonctions croissantes de \mathbb{R} dans \mathbb{R} alors $\rho_S(\phi(X), \psi(Y)) = \rho_S(X, Y)$.

3.3. Estimation de $\rho_S(X, Y)$

À chaque couple (x_i, y_i) de l'échantillon on associe le couple d'entiers (r_i, s_i) où r_i est le rang de x_i dans x_1, \dots, x_n et s_i est le rang de y_i dans y_1, \dots, y_n . On appelle R_i la valeur aléatoire associée au rang d'un X_i , et S_i celle associée au rang d'un Y_i , puis on pose $R = (R_1, \dots, R_n)$ et $S = (S_1, \dots, S_n)$. On calcule alors simplement le coefficient de corrélation des rangs :

$$\begin{aligned} \rho_{S,n}(X, Y) &= \widehat{\rho(R, S)} = \frac{\sum_{i=1}^n (R_i - \bar{R}) \times (S_i - \bar{S})}{\widehat{\sigma}_R \widehat{\sigma}_S} \\ &= \frac{12}{n^3 - n} \sum_{i=1}^n \left[\left(R_i - \frac{n+1}{2} \right) \left(S_i - \frac{n+1}{2} \right) \right] \\ &= 1 - \frac{6}{n^3 - n} \sum_{i=1}^n (R_i - S_i)^2. \end{aligned}$$

La réalisation $\rho_{S,n}(X, Y)(\mathbf{x}, \mathbf{y})$ de la statistique $\rho_{S,n}(X, Y)$ sur l'échantillon $((x_1, y_1), \dots, (x_n, y_n))$, noté (\mathbf{x}, \mathbf{y}) , possède les propriétés suivantes :

- $-1 \leq \rho_{S,n}(X, Y)(\mathbf{x}, \mathbf{y}) \leq 1$,
- $\rho_{S,n}(X, Y)(\mathbf{x}, \mathbf{y}) = 1$ si et seulement si $\forall i = 1 \dots n, r_i = s_i$,
- $\rho_{S,n}(X, Y)(\mathbf{x}, \mathbf{y}) = -1$ si et seulement si $\forall i = 1 \dots n, r_i = n + 1 - s_i$.

On introduit souvent une variable auxiliaire $D_{S,n}^2(X, Y) = \sum_{i=1}^n (R_i - S_i)^2$. On a alors la relation suivante :

$$\rho_{S,n}(X, Y) = 1 - \frac{6D_{S,n}^2(X, Y)}{n^3 - n}.$$

Ainsi dans certaines tables vous trouverez les valeurs de $D_{S,n}$.

La statistique $\rho_{s,n}(X, Y)$ possède les propriétés suivantes :

- $\rho_{s,n}(X, Y)$ est un estimateur convergent avec biais de $\rho_s(X, Y)$.

$$\mathbb{E}[\rho_{s,n}(X, Y)] = \rho_s(X, Y) + \frac{3(\tau(X, Y) - \rho_s(X, Y))}{n+1}.$$

Pour la définition de $\tau(X, Y)$, voir la section suivante. On note donc désormais $\rho_{s,n}(X, Y)$ par $\widehat{\rho_s(X, Y)}$.

- Sous l'hypothèse nulle \mathcal{H}_0 « X et Y sont indépendantes » on a :

- ★ $\mathbb{E}[\widehat{\rho_s(X, Y)}] = 0,$

- ★ $\text{Var}[\widehat{\rho_s(X, Y)}] = \frac{1}{n-1}.$

- ★ Pour $4 \leq n \leq 19$, la distribution de $\widehat{\rho_s(X, Y)}$ se déduit de celle de $D_{s,n}^2(X, Y)$ qui est tabulée.

- ★ Pour $20 \leq n \leq 30$, on utilise la variable $R_n(X, Y) = \sqrt{n-2} \frac{\widehat{\rho_s(X, Y)}}{\sqrt{1 - \widehat{\rho_s(X, Y)}^2}}$ qui suit approximativement une loi de Student à $n-2$ degrés de liberté.

- ★ Pour $n > 30$, on a l'approximation normale $\sqrt{n-1} \widehat{\rho_s(X, Y)} \approx \mathcal{N}(0, 1)$.

On rappelle qu'un estimateur $\widehat{\rho_s(X, Y)}$ de $\rho_s(X, Y)$ est une variable aléatoire $\widehat{\rho_s(X, Y)}$ qui prend pour valeur une estimation $\widehat{\rho_s(X, Y)}(\mathbf{x}, \mathbf{y})$ sur chaque échantillon (\mathbf{x}, \mathbf{y}) de taille n . Ainsi un estimateur est une variable aléatoire alors qu'une estimation est un nombre réel construit pour être proche de la vraie valeur du paramètre inconnu $\rho_s(X, Y)$.

3.4. Procédure de test

La statistique $\widehat{\rho_s(X, Y)}$ permet de tester l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendants » contre plusieurs alternatives, cf plus haut.

$$\boxed{\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}}$$

contre

$$\boxed{\mathcal{H}_1 : X \text{ et } Y \text{ sont liées.}}$$

Si l'on cherche un niveau de signification de α , on cherche r_α tel que :

$$\mathbb{P}[\widehat{\rho_s(X, Y)} \geq r_\alpha] \leq \frac{\alpha}{2}.$$

On rejette l'hypothèse nulle \mathcal{H}_0 si, pour notre échantillon (\mathbf{x}, \mathbf{y}) , $\widehat{\rho_s(X, Y)}(\mathbf{x}, \mathbf{y}) \notin]-r_\alpha, r_\alpha[.$

$\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}$

contre

 $\mathcal{H}_1 : \rho_s(X, Y) > 0$, i.e. les valeurs prises par X et Y ont tendance à être concordantes.

Si l'on cherche un niveau de signification de α , on cherche r_α tel que :

$$\mathbb{P} \left[\widehat{\rho_s(X, Y)} \geq r_\alpha \right] \leq \alpha.$$

On rejette l'hypothèse nulle \mathcal{H}_0 si, pour notre échantillon (\mathbf{x}, \mathbf{y}) , $\widehat{\rho_s(X, Y)}(\mathbf{x}, \mathbf{y}) \geq r_\alpha$.

 $\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}$

contre

 $\mathcal{H}_1 : \rho_s(X, Y) < 0$, i.e. les valeurs prises par X et Y ont tendance à être discordantes.

Si l'on cherche un niveau de signification de α , on cherche r_α tel que :

$$\mathbb{P} \left[\widehat{\rho_s(X, Y)} \leq r_\alpha \right] \leq \alpha.$$

On rejette l'hypothèse nulle \mathcal{H}_0 si, pour notre échantillon (\mathbf{x}, \mathbf{y}) , $\widehat{\rho_s(X, Y)}(\mathbf{x}, \mathbf{y}) \leq r_\alpha$.

Exemple 3.1.

On reprend les données de l'exemple 1.5. On remarque qu'il n'y a pas d'ex æquo. On peut donc calculer les estimations des coefficients de corrélation de Spearman ρ_s (Maths, Sport), ρ_s (Maths, Age) et ρ_s (Sport, Age).

En utilisant SPSS 13.0 on obtient :

Correlations

			Maths	Sport	Age
Spearman's rho	Maths	Correlation Coefficient	1,000	,819(**)	,890(**)
		Sig. (2-tailed)	.	,000	,000
		N	30	30	30
	Sport	Correlation Coefficient	,819(**)	1,000	,818(**)
Sig. (2-tailed)		,000	.	,000	
N		30	30	30	

		N	30	30	30
	Age	Correlation Coefficient	,890(**)	,818(**)	1,000
		Sig. (2-tailed)	,000	,000	.
		N	30	30	30

** Correlation is significant at the 0.01 level (2-tailed).

Tous les tests sont significatifs au seuil $\alpha = 5\%$. Il y a donc une association significative entre toutes les variables prises deux à deux. Voir le paragraphe 5.1 pour un calcul *exact* des p -valeurs avec SPSS 13.0 et le module Tests Exactes.

.....

3.5. Cas des ex æquo

Deux méthodes sont adaptées au cas où la variable n'est pas continue ou au cas où l'on a observé des ex æquo :

- on modifie la valeur de $\widehat{\rho_S(X, Y)}$,
- on départage les ex æquo à l'aide d'une table de nombres aléatoires.

3.6. Statistique corrigée $\rho_{S,n}^*$

On se donne un échantillon (\mathbf{x}, \mathbf{y}) qui sont les réalisations d'un n -échantillon indépendant et identiquement distribué suivant la loi de (X, Y) .

Les n valeurs observées x_1, \dots, x_n sont regroupées en h classes d'ex æquo C_1, \dots, C_h . Certaines de ces classes peuvent ne comporter qu'un seul élément, si cet élément n'a pas d'ex æquo. On regroupe de même les valeurs y_i en k classes d'ex æquo C'_i . Au couple de rangs réels (r_i, s_i) associé à (x_i, y_i) on substitue le couple de rangs fictifs (r_i^*, s_i^*) où r_i^* est le rang moyen du groupe d'ex æquo auquel appartient x_i et s_i^* est le rang moyen du groupe d'ex æquo auquel appartient y_i . On note $d_i = \text{Card}(C_i)$ et $d'_i = \text{Card}(C'_i)$.

On calcule alors $\delta(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n (d_i^3 - d_i)$ et $\delta'(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n ((d'_i)^3 - d'_i)$.

Dans le cas d'ex æquo la valeur prise par la statistique $\rho_{S,n}^*$ est :

$$\begin{aligned} \rho_{S,n}^*(X, Y)(\mathbf{x}, \mathbf{y}) &= \frac{12}{\sqrt{(n^3 - n - \delta)(n^3 - n - \delta')}} \sum_{i=1}^n \left[\left(r_i^* - \frac{n+1}{2} \right) \left(s_i^* - \frac{n+1}{2} \right) \right] \\ &= \frac{12 \sum_{i=1}^n (r_i \times s_i) - 3n(n+1)^2}{\sqrt{(n^3 - n - \delta)(n^3 - n - \delta')}}. \end{aligned}$$

$\rho_{s,n}^*(X, Y)$ est donc la variable aléatoire associée à $\rho_{s,n}^*(X, Y)(\mathbf{x}, \mathbf{y})$.

Lorsque $n > 20$ et $(\delta(\mathbf{x}, \mathbf{y}) + \delta'(\mathbf{x}, \mathbf{y}))/n^3 < 0, 1$ on utilise l'approximation normale :

$$\sqrt{n-1}\rho_{s,n}^*(X, Y)^* \approx \mathcal{N}(0, 1).$$

Dans les autres situations, il n'y a pas de table numérique.

Exemple 3.2.

On dispose de 7 couples de valeurs entières prises par le couple de variables aléatoires discrètes (X, Y) . On a classé ces valeurs suivant les valeurs croissantes de X . Les couples (x_i, y_i) , (r_i, s_i) et (r_i^*, s_i^*) ont été reportés dans le tableau ci-dessous :

x_i	0	1	1	1	2	2	5
r_i	1	2	2	2	5	5	7
r_i^*	1	3	3	3	5,5	5,5	7
y_i	3	2	0	3	4	5	3
s_i	3	2	1	3	6	7	3
s_i^*	4	2	1	4	6	7	4

On obtient :

$$\sum_{i=1}^7 (r_i^* \times s_i^*) = 124,5$$

$$\delta = (3^3 - 3) + (2^3 - 2) = 30$$

$$\delta' = (3^3 - 3) = 24$$

$$\rho_s^* = 0,485.$$

Comme $n = 7 < 20$, on ne peut utiliser l'approximation du paragraphe 3.6. On ne dispose pas non plus de table. On ne peut rien conclure quant à la significativité d'un test de corrélation partielle.

.....

3.7. Départition des ex æquo

On départage les ex æquo à l'aide d'une table de nombres aléatoires de la manière suivante : si on a r répétitions, on lit r nombres de suite dans la table de nombres aléatoires. La valeur des nombres permettent alors de départager les ex æquo, et donc de se ramener au cas où il n'y a pas de répétition.

Exemple 3.3.

Si par exemple on est face au cas où $x_1 = x_4 = x_7 = x_9 = 1$. On lit dans une table de nombres aléatoires 13407, 50270, 84980 et 22116. On décide alors $x_1 < x_9 < x_4 < x_7$ et l'on note $x_1 = 1^+$, $x_4 = 1^{+++}$, $x_7 = 1^{++++}$ et $x_9 = 1^{++}$. On aura alors, par exemple, $(x_7 - x_4) > 0$ et $(x_1 - x_9) < 0$.

.....

Exemple 3.4.

Appliquer ce procédé à l'exemple 3.2 ci-dessus.

.....

4. Corrélation partielle de Spearman

L'observation d'une corrélation entre deux variables X et Y peut être due à leur association avec une troisième Z .

4.1. Coefficient de corrélation partielle de Spearman $\rho_S(X, Y|Z)$

Le coefficient de corrélation de Spearman est lié à l'analyse de la corrélation de Pearson des rangs d'un n -échantillon.

Si l'on a trois variables aléatoires X , Y et Z réelles et continues, on définit le coefficient de corrélation partielle de Spearman comme étant le coefficient de corrélation partielle de Pearson des rangs :

$$\frac{\rho_S(X, Y) - \rho_S(X, Z) \times \rho_S(Y, Z)}{\sqrt{(1 - \rho_S(X, Z)^2)(1 - \rho_S(Y, Z)^2)}}.$$

4.2. Estimation de $\rho_S(X, Y|Z)$

On considère trois variables aléatoires X , Y et Z ainsi qu'un n -échantillon indépendant identiquement distribué suivant la loi de (X, Y, Z) , noté $((X_1, Y_1, Z_1), \dots, (X_n, Y_n, Z_n))$. On suppose que (X, Y, Z) suit une loi continue.

On définit la statistique de corrélation partielle de Spearman $\rho_{S,n}(X, Y|Z)$, qui s'exprime à l'aide des statistiques de corrélation de Spearman des n couples (X_i, Y_i) , des n couples (X_i, Z_i) et des n couples (Y_i, Z_i) , par l'expression suivante :

$$\rho_{S,n}(X, Y|Z) = \frac{\rho_{S,n}(X, Y) - \rho_{S,n}(X, Z) \times \rho_{S,n}(Y, Z)}{\sqrt{(1 - \rho_{S,n}(X, Z)^2)(1 - \rho_{S,n}(Y, Z)^2)}}.$$

$\rho_{S,n}(X, Y|Z)$ a les propriétés suivantes :

- $-1 \leq \rho_{S,n}(X, Y|Z) \leq 1$,
- $\rho_{S,n}(X, Y|Z)$ est un estimateur convergent du coefficient de corrélation partielle de rang de Spearman $\rho_S(X, Y|Z)$. On notera donc désormais $\rho_{S,n}(X, Y|Z)$ par $\widehat{\rho_S(X, Y|Z)}$.

Sous l'hypothèse nulle $\mathcal{H}_0 \ll \rho_{S,n}(X, Y|Z) = 0 \gg$, la distribution de $\widehat{\rho_S(X, Y|Z)}$ est tabulée.

4.3. Méthode d'utilisation

Identique au cas non partiel.

Exemple 4.1.

On souhaite étudier la corrélation de Spearman des variables de l'exemple 1.5, Maths et Sport, en éliminant l'influence de la variable Age. On calcule donc une estimation de la corrélation partielle de Spearman $\rho_S(\text{Maths}, \text{Sport}|\text{Age})$:

$$\widehat{\rho_S(\text{Maths}, \text{Sport}|\text{Age})}(\mathbf{x}) = \frac{0,819 - 0,890 * 0,818}{\sqrt{(1 - 0,890^2)(1 - 0,818^2)}} \approx 0,347.$$

On doit alors chercher dans une table adéquate la valeur critique pour $\alpha = 5 \%$. Si vous parvenez à en trouver une, comparer ce résultat à celui obtenu dans le cadre paramétrique à l'exemple 1.5.

.....

5. La statistique τ_n de Kendall

5.1. Cadre d'application

La mesure de la dépendance au sens de Kendall entre deux variables aléatoires continues X et Y sera notée $\tau(X, Y)$.

Ce nombre pourra être estimé par la statistique de Kendall τ_n définie sur un échantillon $((x_1, y_1), \dots, (x_n, y_n))$ indépendant et identiquement distribué.

La statistique τ_n permet de réaliser plusieurs tests :

$$\boxed{\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}}$$

contre

$$\boxed{\mathcal{H}_1 : X \text{ et } Y \text{ sont liées.}}$$

$$\boxed{\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}}$$

contre

$$\boxed{\mathcal{H}_1 : \text{Les valeurs prises par } X \text{ et } Y \text{ ont tendance à être concordantes.}}$$

$\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}$

contre

 $\mathcal{H}_1 : \text{Les valeurs prises par } X \text{ et } Y \text{ ont tendance à être discordantes.}$

5.2. Le coefficient de corrélation $\tau(X, Y)$ de Kendall

Considérons deux paires (x_i, y_i) et (x_j, y_j) issues de l'échantillon. Elles sont dites :

- concordantes si $(x_i - x_j)(y_i - y_j) > 0$, c'est-à-dire si l'on a simultanément $(x_i > x_j)$ et $(y_i > y_j)$ ou $(x_i < x_j)$ et $(y_i < y_j)$.
- discordantes si $(x_i - x_j)(y_i - y_j) < 0$, c'est-à-dire si l'on a simultanément $(x_i > x_j)$ et $(y_i < y_j)$ ou $(x_i < x_j)$ et $(y_i > y_j)$.

Considérons deux couples de variables aléatoires (X_1, Y_1) et (X_2, Y_2) de même loi que celle du couple étudié (X, Y) .

Une concordance parfaite est telle que $X_2 > X_1 \iff Y_2 > Y_1$, c'est-à-dire :

$$\begin{aligned} \mathbb{P}[(X_2 > X_1) \text{ et } (Y_2 > Y_1)] + \mathbb{P}[(X_2 < X_1) \text{ et } (Y_2 < Y_1)] &= 1 \\ &\Downarrow \\ \mathbb{P}[(X_2 - X_1)(Y_2 - Y_1) > 0] &= 1. \end{aligned}$$

Une discordance parfaite est telle que $X_2 > X_1 \iff Y_2 < Y_1$, c'est-à-dire :

$$\begin{aligned} \mathbb{P}[(X_2 > X_1) \text{ et } (Y_2 < Y_1)] + \mathbb{P}[(X_2 < X_1) \text{ et } (Y_2 > Y_1)] &= 1 \\ &\Downarrow \\ \mathbb{P}[(X_2 - X_1)(Y_2 - Y_1) < 0] &= 1. \end{aligned}$$

On introduit donc $\tau^+(X, Y)$ et $\tau^-(X, Y)$ qui mesurent respectivement la concordance et la discordance du couple (X, Y) :

$$\begin{aligned} \tau^+(X, Y) &= \mathbb{P}[(X_2 - X_1)(Y_2 - Y_1) > 0] \\ \tau^-(X, Y) &= \mathbb{P}[(X_2 - X_1)(Y_2 - Y_1) < 0] \end{aligned}$$

Le coefficient $\tau(X, Y)$ de Kendall est défini par :

$$\tau(X, Y) = \tau^+(X, Y) - \tau^-(X, Y).$$

Il mesure le degré de concordance si $\tau(X, Y) > 0$ ou au contraire le degré de discordance si $\tau(X, Y) < 0$.

Il possède les propriétés suivantes :

- $-1 \leq \tau(X, Y) \leq 1$,
- X et Y indépendantes implique $\tau(X, Y) = 0$,
- $\tau(X, Y) = 1$ (resp. $\tau(X, Y) = -1$) si et seulement si il existe une fonction ϕ croissante (resp. décroissante) de \mathbb{R} dans \mathbb{R} telle que $Y = \phi(X)$,

- Si ϕ et ψ désignent deux fonctions croissantes de \mathbb{R} dans \mathbb{R} alors $\tau(\phi(X), \psi(Y)) = \tau(X, Y)$,
- Si (X, Y) suit une loi normale bivariée alors $\tau(X, Y)$ et $\rho(X, Y)$ sont liés par la relation : $\rho(X, Y) = \sin(\pi/2 \times \tau(X, Y))$.

5.3. Estimation de $\tau(X, Y)$

L'estimation se fait de manière « naturelle » : on commence par compter le nombre de paires de couples concordants c et le nombre de paires de couples discordants d dans l'échantillon $(x_1, y_1), \dots, (x_n, y_n)$. Il apparaît ici une difficulté supplémentaire par rapport à la théorie. En effet la loi de (X, Y) est continue dont la probabilité que $X_1 = X_2$ ou que $Y_1 = Y_2$ est nulle. Mais dans l'échantillon il se peut néanmoins que vous observiez plusieurs fois la même valeur. On note alors e le nombre de ces paires de couples. Ce sont celles pour lesquelles on a $(x_j - x_i)(y_j - y_i) = 0$. Dans un premier temps on supposera qu'il n'y a pas d'ex æquo.

c = le nombre de paires de couples $(x_i, y_i), (y_i, y_j)$ tels que $(x_j - x_i)(y_j - y_i) > 0$
avec $1 \leq i < j \leq n$.

d = le nombre de paires de couples $(x_i, y_i), (y_i, y_j)$ tels que $(x_j - x_i)(y_j - y_i) < 0$
avec $1 \leq i < j \leq n$.

e = le nombre de paires de couples $(x_i, y_i), (y_i, y_j)$ tels que $(x_j - x_i)(y_j - y_i) = 0$
avec $1 \leq i < j \leq n$.

Rappelons que pour le moment on suppose que $e = 0$. Toute paire de couples est forcément du type c ou du type d :

$$c + d = \frac{n(n-1)}{2}$$

qui est le nombre total de paires de couples qu'il est possible de faire.

On note C_n, D_n et E_n les variables aléatoires associées à c, d et e .

On définit alors :

$$\begin{aligned} \tau_n(X, Y) &= \frac{2C_n - \frac{n(n-1)}{2}}{\frac{n(n-1)}{2}} \\ &= \frac{4C_n}{n(n-1)} - 1. \end{aligned}$$

Sous l'hypothèse nulle \mathcal{H}_0 « X et Y sont indépendantes » la distribution de τ_n a les propriétés suivantes :

- $-1 \leq \tau_n \leq 1$, une valeur proche de 1 suggérant une forte concordance entre les valeurs prises par X et Y , une valeur proche de -1 suggérant une forte discordance entre les valeurs prises par X et Y .

- $\tau_n(X, Y)$ est un estimateur sans biais et convergent de $\tau(X, Y)$. On note donc désormais :

$$\tau_n(X, Y) = \widehat{\tau(X, Y)}$$

- Sous l'hypothèse nulle \mathcal{H}_0 « X et Y sont indépendantes », la distribution de $\widehat{\tau(X, Y)}$ a les caractéristiques suivantes :

$$\star \mathbb{E} \left[\widehat{\tau(X, Y)} \right] = 0,$$

$$\star \text{Var} \left[\widehat{\tau(X, Y)} \right] = \frac{2(2n + 5)}{9n(n - 1)},$$

$$\star \frac{\widehat{\tau(X, Y)}}{\sqrt{\text{Var}[\widehat{\tau(X, Y)}]}} \approx \mathcal{N}(0, 1).$$

Pour réaliser des tests avec des effectifs inférieurs à 20, on se reportera donc à des tables spécifiques ou l'on utilisera un logiciel disposant de statistiques exactes.

5.4. Procédure de test

La statistique $\widehat{\tau(X, Y)}$ permet de tester l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendantes » contre plusieurs alternatives, cf plus haut.

\mathcal{H}_0 : X et Y sont indépendantes

contre

\mathcal{H}_1 : X et Y sont liées.

Si l'on cherche un niveau de signification de α , on cherche t_α tel que :

$$\mathbb{P} \left[\widehat{\tau(X, Y)} \geq t_\alpha \right] \leq \frac{\alpha}{2}.$$

On rejette l'hypothèse nulle \mathcal{H}_0 si $\widehat{\tau(X, Y)} \notin]-t_\alpha, t_\alpha[$.

\mathcal{H}_0 : X et Y sont indépendantes

contre

\mathcal{H}_1 : $\tau(X, Y) > 0$, i.e. les valeurs prises par X et Y ont tendance à être concordantes.

Si l'on cherche un niveau de signification de α , on cherche t_α tel que :

$$\mathbb{P} \left[\widehat{\tau(X, Y)} \geq t_\alpha \right] \leq \alpha.$$

On rejette l'hypothèse nulle \mathcal{H}_0 si $\widehat{\tau(X, Y)} \geq t_\alpha$.

$\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}$

contre

 $\mathcal{H}_1 : \tau(X, Y) < 0$, i.e. les valeurs prises par X et Y ont tendance à être discordantes.

Si l'on cherche un niveau de signification de α , on cherche t_α tel que :

$$\mathbb{P} \left[\widehat{\tau(X, Y)} \leq t_\alpha \right] \leq \alpha.$$

On rejette l'hypothèse nulle \mathcal{H}_0 si $\widehat{\tau(X, Y)} \leq t_\alpha$.

Exemple 5.1.

On reprend les données de l'exemple 1.5. On remarque qu'il n'y a pas d'ex æquo. On peut donc calculer les estimations des coefficients de corrélation de Kendall τ (Maths, Sport), τ (Maths, Age) et τ (Sport, Age).

En utilisant SPSS 13.0 on obtient :

Correlations

			Maths	Sport	Age
Kendall's	Maths	Correlation Coefficient	1,000	,605(**)	,701(**)
tau_b		Sig. (2-tailed)	.	,000	,000
		N	30	30	30
	Sport	Correlation Coefficient	,605(**)	1,000	,600(**)
		Sig. (2-tailed)	,000	.	,000
		N	30	30	30
	Age	Correlation Coefficient	,701(**)	,600(**)	1,000
		Sig. (2-tailed)	,000	,000	.
		N	30	30	30

** Correlation is significant at the 0.01 level (2-tailed).

Tous les tests sont significatifs au seuil $\alpha = 5 \%$. Il y a donc une association significative entre toutes les variables prises deux à deux.

On calcule alors les p -valeurs exactes associées aux estimations des corrélations de Spearman et de Kendall. Or l'obtention de ces valeurs serait trop longue pour nos moyens de calcul actuels et l'on doit donc utiliser une méthode de Monte-Carlo pour trouver un intervalle de confiance à 99 % sur ces p -valeurs. Si le seuil α du test n'appartient pas à l'intervalle obtenu on conclut comme à l'accoutumée en fonction de la place de l'intervalle par rapport à α . Si α appartient à l'intervalle obtenu, il faut relancer le calcul en augmentant la taille de l'échantillon utilisé pour calculer l'intervalle de confiance sur la p -valeur par la méthode de Monte Carlo.

Symmetric Measures

	Monte Carlo Sig.		
	Sig.	99% Confidence Interval	
		Lower Bound	Upper Bound
Kendall's tau-b	,000(c)	,000	,000
Spearman Correlation	,000(c)	,000	,000
Pearson's R	,000(c)	,000	,000

c Based on 10000 sampled tables with starting seed 2000000.

Symmetric Measures

	Monte Carlo Sig.		
	Sig.	99% Confidence Interval	
		Lower Bound	Upper Bound
Kendall's tau-b	,000(c)	,000	,000
Spearman Correlation	,000(c)	,000	,000
Pearson's R	,000(c)	,000	,000

c Based on 10000 sampled tables with starting seed 1314643744.

Symmetric Measures

	Monte Carlo Sig.		
	Sig.	99% Confidence Interval	
		Lower Bound	Upper Bound

	Sig.	99% Confidence Interval	
		Lower Bound	Upper Bound
Kendall's tau-b	,000(c)	,000	,000
Spearman Correlation	,000(c)	,000	,000
Pearson's R	,000(c)	,000	,000

c Based on 10000 sampled tables with starting seed 1502173562.

Ici, pour $\alpha = 5\%$, tous les intervalles sont entièrement inclus dans $[0, \alpha]$ et donc tous les tests sont significatifs au seuil $\alpha = 5\%$.

.....

5.5. Cas des ex æquo

Deux méthodes sont adaptées au cas où la variable n'est pas continue ou au cas où l'on a observé des ex æquo :

- on modifie la valeur de $\tau(\widehat{X}, \widehat{Y})$,
- on départage les ex æquo à l'aide d'une table de nombres aléatoires.

5.6. Statistique corrigée τ_n^*

Les n valeurs observées x_1, \dots, x_n sont regroupées en h classes d'ex æquo C_1, \dots, C_h . Certaines de ces classes peuvent ne comporter qu'un seul élément, si cet élément n'a pas d'ex æquo. On regroupe de même les valeurs y_j en k classes d'ex æquo C'_j . On note $d_i = \text{Card}(C_i)$ et $d'_j = \text{Card}(C'_j)$. Enfin $d = \sum_{i=1}^h d_i(d_i - 1)$ et $d' = \sum_{j=1}^k d'_j(d'_j - 1)$.

On calcule alors :

$$s^* = 1 \times \text{Nombre de concordants} + (-1) \times \text{Nombre de discordants} \\ + 0 \times \text{Nombre de cas d'égalité.}$$

Puis on pose :

$$\tau^* = \frac{2s^*}{\sqrt{(n(n-1) - d)(n(n-1) - d')}}.$$

τ_n^* est alors la variable aléatoire associée à τ^* et S_n^* celle associée à s^* .

Lorsque $n > 10$, $d/n^2 < 0,1$ et $d'/n^2 < 0,1$ on utilise l'approximation normale :

$$\frac{\sqrt{18}S_n^*}{\sqrt{n(n-1)(2n+5) - \delta - \delta'}} \approx \mathcal{N}(0, 1)$$

où $\delta = \sum_{i=1}^h d_i(d_i - 1)(2d_i + 5)$ et $\delta' = \sum_{j=1}^k d'_j(d'_j - 1)(2d'_j + 5)$.

Exemple 5.2.

Appliquer ce procédé à l'exemple 3.2.

.....

5.7. Départition des ex æquo

On départage les ex æquo à l'aide d'une table de nombres aléatoires de la manière suivante : si on a r répétitions, on lit r nombres de suite dans la table de nombres aléatoires. La valeur des nombres permettent alors de départager les ex æquo, et donc de se ramener au cas où il n'y a pas de répétition.

Exemple 5.3.

Si par exemple on est face au cas où $x_1 = x_4 = x_7 = x_9 = 1$. On lit dans une table de nombres aléatoires 13407, 50270, 84980 et 22116. On décide alors $x_1 < x_9 < x_4 < x_7$ et l'on note $x_1 = 1^+$, $x_4 = 1^{+++}$, $x_7 = 1^{++++}$ et $x_9 = 1^{++}$. On aura alors, par exemple, $(x_7 - x_4) > 0$ et $(x_1 - x_9) < 0$.

.....

Exemple 5.4.

Appliquer ce procédé à l'exemple 3.2.

.....

6. Corrélration partielle de Kendall

L'observation d'une corrélation entre deux variables X et Y peut être due à leur association avec une troisième Z .

6.1. Coefficient de corrélation partiel de Kendall $\tau(X, Y|Z)$

On commence par chercher les probabilités de concordance et discordance pour deux triplets (X_1, Y_1, Z_1) , (X_2, Y_2, Z_2) indépendants et identiquement distribués suivant la loi de (X, Y, Z) . On suppose que la loi de (X, Y, Z) est continue.

- La concordance de X et Z sans tenir compte de Y est :

$$p_{1,\bullet} = \mathbb{P}[(X_2 - X_1)(Z_2 - Z_1) > 0].$$

- La discordance de X et Z sans tenir compte de Y est :

$$p_{2,\bullet} = \mathbb{P}[(X_2 - X_1)(Z_2 - Z_1) < 0].$$

- La concordance de Y et Z sans tenir compte de X est :

$$p_{\bullet,1} = \mathbb{P}[(Y_2 - Y_1)(Z_2 - Z_1) > 0].$$

- La discordance de Y et Z sans tenir compte de X est :

$$p_{\bullet,2} = \mathbb{P}[(Y_2 - Y_1)(Z_2 - Z_1) < 0].$$

- La concordance simultanée de X et Y avec Z est :

$$p_{1,1} = \mathbb{P}[(X_2 - X_1)(Z_2 - Z_1) > 0 \text{ et } (Y_2 - Y_1)(Z_2 - Z_1) > 0].$$

- La probabilité que, simultanément, (X_1, Z_1) et (X_2, Z_2) soient concordants et que (Y_1, Z_1) et (Y_2, Z_2) soient discordants est :

$$p_{2,1} = \mathbb{P}[(X_2 - X_1)(Z_2 - Z_1) < 0 \text{ et } (Y_2 - Y_1)(Z_2 - Z_1) > 0].$$

- La probabilité que, simultanément, (X_1, Z_1) et (X_2, Z_2) soient discordants et que (Y_1, Z_1) et (Y_2, Z_2) soient concordants est :

$$p_{1,2} = \mathbb{P}[(X_2 - X_1)(Z_2 - Z_1) > 0 \text{ et } (Y_2 - Y_1)(Z_2 - Z_1) < 0].$$

- La discordance simultanée de X et Y avec Z est :

$$p_{2,2} = \mathbb{P}[(X_2 - X_1)(Z_2 - Z_1) < 0 \text{ et } (Y_2 - Y_1)(Z_2 - Z_1) < 0].$$

Le coefficient de corrélation partielle de Kendall est alors :

$$\tau(X, Y|Z) = \frac{p_{1,1}p_{2,2} - p_{1,2}p_{2,1}}{\sqrt{(p_{1,\bullet}p_{2,\bullet}p_{\bullet,1}p_{\bullet,2})}}.$$

Ses propriétés sont les suivantes :

- $-1 \leq \tau(X, Y|Z) \leq 1$,
- on peut exprimer le coefficient de corrélation partielle de Kendall du couple (X, Y) connaissant Z , $\tau(X, Y|Z)$, à l'aide des coefficients de corrélation de Kendall $\tau(X, Y)$, $\tau(X, Z)$ et $\tau(Y, Z)$ des couples (X, Y) , (X, Z) et (Y, Z) :

$$\tau(X, Y|Z) = \frac{\tau(X, Y) - \tau(X, Z) \times \tau(Y, Z)}{\sqrt{(1 - \tau(X, Z)^2)(1 - \tau(Y, Z)^2)}}$$

- « (X, Y) et Z indépendants » donne $\tau(X, Y|Z) = \tau(X, Y)$,
- « X, Y et Z indépendantes » donne $\tau(X, Y|Z) = 0$.

6.2. Notations

On considère trois variables aléatoires X , Y et Z ainsi qu'un n -échantillon indépendant identiquement distribué de (X, Y, Z) , noté $((X_1, Y_1, Z_1), \dots, (X_n, Y_n, Z_n))$. On suppose que (X, Y, Z) suit une loi continue.

On va à nouveau utiliser la notion de paire de couples concordants ou discordants mais en la modifiant pour tenir compte de la présence de la troisième variable aléatoire Z . On dispose de $n(n-1)/2$ paires de triplets $(X_i, Y_i, Z_i), (X_j, Y_j, Z_j)$ avec $1 \leq i < j \leq n$.

On considère donc pour commencer :

- $C_{x,z}$ l'ensemble des paires de triplets pour lesquels X et Z sont concordants.
 - $D_{x,z}$ l'ensembles des paires de triplets pour lesquels X et Z sont discordants.
- On note $N_{1\bullet}$ le nombre d'éléments de $C_{x,z}$ et $N_{2\bullet}$ le nombre d'éléments $D_{x,z}$.

Puis on s'intéresse à :

- $C_{y,z}$ l'ensemble des paires de triplets pour lesquels Y et Z sont concordants.
 - $D_{y,z}$ l'ensemble des paires de triplets pour lesquels Y et Z sont discordants.
- On note $N_{\bullet 1}$ le nombre d'éléments de $C_{y,z}$ et $N_{\bullet 2}$ le nombre d'éléments $D_{y,z}$.

Pour prendre en compte simultanément les variables X et Y on partage les $n(n-1)/2$ paires de triplets en quatre groupes disjoints.

- $C_{1,1}$ correspond donc aux paires de triplets pour lesquelles à la fois X et Z sont concordants et Y et Z sont concordants :

$$C_{1,1} = C_{x,z} \cap C_{y,z}.$$

On note $N_{1,1}$ le nombre d'éléments de $C_{1,1}$.

- $C_{1,2}$ correspond donc aux paires de triplets pour lesquelles à la fois X et Z sont concordants et Y et Z sont discordants :

$$C_{1,2} = C_{x,z} \cap D_{y,z}.$$

On note $N_{2,1}$ le nombre d'éléments de $C_{2,1}$.

- $C_{2,1}$ correspond donc aux paires de triplets pour lesquelles à la fois X et Z sont discordants et Y et Z sont concordants :

$$C_{2,1} = D_{x,z} \cap C_{y,z}.$$

On note $N_{1,2}$ le nombre d'éléments de $C_{1,2}$.

- $C_{2,2}$ correspond donc aux paires de triplets pour lesquelles à la fois X et Z sont discordants et Y et Z sont discordants :

$$C_{2,2} = D_{x,z} \cap D_{y,z}.$$

On note $N_{2,2}$ le nombre d'éléments de $C_{2,2}$.

On résume la situation dans le tableau suivant :

	$C_{x,z}$	$D_{x,z}$	Total
$C_{y,z}$	$N_{1,1}$	$N_{1,1}$	$N_{\bullet,1}$
$D_{y,z}$	$N_{1,1}$	$N_{1,1}$	$N_{\bullet,2}$
Total	$N_{1,\bullet}$	$N_{2,\bullet}$	$n(n-1)/2$

6.3. Estimateur de $\tau(X, Y|Z)$

On définit la statistique de corrélation partielle de Kendall $\tau_{XY|Z,n}$ par l'expression suivante :

$$\tau_{XY|Z,n} = \frac{N_{1,1}N_{2,2} - N_{1,2}N_{2,1}}{\sqrt{(N_{1,\bullet}N_{2,\bullet}N_{\bullet,1}N_{\bullet,2})}}$$

On peut exprimer cette statistique à l'aide des statistiques de corrélation de Kendall des n couples (X_i, Y_i) , des n couples (X_i, Z_i) et des n couples (Y_i, Z_i) :

$$\tau_{XY|Z,n} = \frac{\tau_n(X, Y) - \tau_n(X, Z) \times \tau_n(Y, Z)}{\sqrt{(1 - \tau_n(X, Z)^2)(1 - \tau_n(Y, Z)^2)}}$$

$\tau_{XY|Z,n}$ a les propriétés suivantes :

- $-1 \leq \tau_{XY|Z,n} \leq 1$,
- $\tau_{XY|Z,n}$ est un estimateur convergent du coefficient de corrélation partiel de rang de Kendall $\tau(X, Y|Z)$.

On notera donc désormais $\tau_{XY|Z,n}$ par $\widehat{\tau(X, Y|Z)}$.

Sous l'hypothèse nulle « $\mathcal{H}_0 : \tau(X, Y|Z) = 0$ », la distribution de $\widehat{\tau(X, Y|Z)}$ a les caractéristiques suivantes :

- $\mathbb{E} \left[\widehat{\tau(X, Y|Z)} \right] = 0$, $\text{Var} \left[\widehat{\tau(X, Y|Z)} \right] = \frac{2(2n+5)}{9n(n-1)}$,
- elle est symétrique $\mathbb{P} \left[\widehat{\tau(X, Y|Z)} \leq t \right] = \mathbb{P} \left[\widehat{\tau(X, Y|Z)} \geq t \right]$, pour tout $t \in \mathbb{R}$,
- Pour $n < 30$, on se reporte à une table. Sinon on utilise l'approximation suivante :

$$3\sqrt{\frac{(n^2 - n)}{(4n + 10)}} \widehat{\tau(X, Y|Z)} \approx \mathcal{N}(0, 1).$$

6.4. Méthode d'utilisation

Identique au cas non partiel.

Exemple 6.1.

On souhaite étudier la corrélation de Kendall des variables Maths et Sport en éliminant l'influence de la variable Age. On calcule donc une estimation de la corrélation partielle de Kendall $\tau(\widehat{\text{Maths, Sport}}|\text{Age})(\mathbf{x})$:

$$\tau(\widehat{\text{Maths, Sport}}|\text{Age})(\mathbf{x}) = \frac{0,605 - 0,701 * 0,600}{\sqrt{(1 - 0,701^2)(1 - 0,600^2)}} \approx 0,323.$$

On trouve dans la table que la valeur critique pour $\alpha = 5\%$ est de 0,253 qui est inférieure à notre estimation et de ce fait le test est significatif. On rejette donc l'hypothèse nulle \mathcal{H}_0 d'indépendance Maths et Sport lorsque l'on élimine l'influence de la variable Age et on décide qu'il y a une association significative de Maths et Sport lorsque l'on élimine l'influence de la variable Age. Comparer ce résultat à celui obtenu dans le cadre paramétrique à l'exemple 1.5. Lequel des deux faut-il préférer ?

.....

7. Test du χ^2 d'indépendance

7.1. Les classes

Ce test peut être réalisé quelque soit la loi du couple (X, Y) , discrète, continue ou mixte. Toutefois il est **peu puissant**.

Le domaine \mathcal{D}_X des valeurs de X est divisé en h classes adjacentes $C_{i,\bullet}$. Le domaine \mathcal{D}_Y des valeurs de Y est divisé en k classes adjacentes $C_{\bullet,j}$.

Ainsi les valeurs de l'échantillon se répartissent dans les $h \times k$ classes $C_{i,j} = C_{i,\bullet} \times C_{\bullet,j}$. Le nombre de couples (x_l, y_l) tels que $(x_l, y_l) \in C_{i,j}$ est noté $n_{i,j}$. Le nombre de valeurs x_i appartenant à la classe $C_{i,\bullet}$ est noté $n_{i,\bullet}$ et le nombre de valeurs x_j appartenant à la classe $C_{\bullet,j}$ est noté $n_{\bullet,j}$. On a $n_{i,\bullet} = \sum_{j=1}^k n_{i,j}$ et $n_{\bullet,j} = \sum_{i=1}^h n_{i,j}$.

$X \setminus Y$	1	...	j	...	k	Total
1	$n_{1,1}$...	$n_{1,j}$...	$n_{1,k}$	$n_{1,\bullet}$
\vdots	\vdots		\vdots		\vdots	\vdots
i	$n_{i,1}$...	$n_{i,j}$...	$n_{i,k}$	$n_{i,\bullet}$
\vdots	\vdots		\vdots		\vdots	\vdots
h	$n_{h,1}$...	$n_{h,j}$...	$n_{h,k}$	$n_{h,\bullet}$
Total	$n_{\bullet,1}$...	$n_{\bullet,j}$...	$n_{\bullet,k}$	n

7.2. Statistique du test

Pour tester l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendantes », on introduit la statistique χ_n^2 suivante :

$$\chi_n^2 = \sum_{i=1}^h \sum_{j=1}^k \frac{\left(N_{i,j} - \frac{N_{i,\bullet} \times N_{\bullet,j}}{n} \right)^2}{\frac{N_{i,\bullet} \times N_{\bullet,j}}{n}}$$

où $N_{i,j}$ est le nombre aléatoire de couple (X_i, Y_j) qui dans un échantillon indépendant identiquement distribué de taille n appartient à la classe $C_{i,j}$, $N_{i,\bullet} = \sum_{j=1}^k N_{i,j}$ et $N_{\bullet,j} = \sum_{i=1}^h N_{i,j}$.

Sous l'hypothèse nulle \mathcal{H}_0 la loi de χ_n^2 converge vers la loi du χ^2 à $(h-1)(k-1)$ degrés de liberté lorsque que n tend vers $+\infty$.

Il est raisonnable de faire cette approximation lorsque $n \geq 50$ et que chaque effectif théorique $n_{i,j}^* = n \frac{n_{i,\bullet}}{n} \frac{n_{\bullet,j}}{n}$ est supérieur à 5. Dans la situation où les effectifs théoriques ne seraient pas tous supérieurs à 5, il est nécessaire de procéder au regroupement de classes adjacentes, C_i et C_{i+1} ou C'_j et C'_{j+1} pour $1 \leq i \leq h-1$ et $1 \leq j \leq k-1$.

7.3. Méthode d'utilisation

La valeur critique χ_α^2 du test au seuil α % est le plus petit nombre χ_α^2 tel que $\mathbb{P}[\chi_n^2 \geq \chi_\alpha^2] = \alpha$ où χ_n^2 suit une loi $\chi_{(h-1)(k-1)}^2$.

Si la valeur $\chi_{\text{obs},n}^2$, qui est la réalisation de la statistique du test sur l'échantillon :

$$\chi_{\text{obs},n}^2 = \sum_{i=1}^h \sum_{j=1}^k \frac{\left(n_{i,j} - \frac{n_{i,\bullet} \times n_{\bullet,j}}{n} \right)^2}{\frac{n_{i,\bullet} \times n_{\bullet,j}}{n}}$$

est supérieure à χ_α^2 on rejette l'hypothèse nulle \mathcal{H}_0 d'indépendance.

Ce test pose plusieurs difficultés :

- Comment faire les regroupements en classe lorsque cela s'avère nécessaire, par exemple si les variables étudiées sont continues ?
- Les conditions d'application sont très contraignantes. Un effectif total n **supérieur à 50** et des fréquences d'apparition toutes supérieures à 5. Dans le livre de J. Bouyer, [3], ainsi que dans celui de G. Pupion et P.-C. Pupion, [12], il est indiqué que le test est encore utilisable si les effectifs théoriques sont tous supérieurs à 3. J. Bouyer, [3], évoque même la possibilité de se contenter du fait qu'il y ait moins de 20 % des cellules

pour lesquelles les effectifs théoriques soient inférieurs à 5. Néanmoins tous les auteurs s'entendent pour dire que si dans une telle situation vous obteniez des valeurs proches de la significativité, il est impératif de compléter l'étude par l'utilisation de certains des autres tests présentés ici.

- La nécessité de fondre plusieurs modalités en une seule pour que les conditions d'applications, mentionnées au paragraphe ci-dessus, soient remplies modifie les variables sur lesquelles porte le test.
- Ce test ne tient pas compte de l'éventuelle présence d'un ordre sur les lignes ou les colonnes du tableau de contingence. Si l'on peut ordonner les modalités de l'un des deux facteurs, on préférera utiliser un test de Kruskal-Wallis et si l'on peut ordonner les modalités des deux facteurs on utilisera un test de Jonckheere-Terpstra. On pourra alors, si l'on rejette l'hypothèse nulle \mathcal{H}_0 : « Le facteur X n'a pas d'effet sur la réponse Y », étudier les raisons à l'origine de la non-indépendance à l'aide d'un test post-hoc. Le cas d'un tableau à deux lignes et k colonnes ou à h lignes et deux colonnes peut également être étudié à l'aide d'un test de Mann-Whitney.

7.4. Correction de Yates

Lorsque l'on étudie l'indépendance de deux variables et que certaines des fréquences attendues sous l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendantes » sont inférieures à 5, on peut corriger la statistique du test pour prendre en compte cette situation. Attention il faut néanmoins que toutes les fréquences attendues soient supérieures à 3. La correction de Yates est une correction de continuité qui consiste à utiliser la statistique de test modifiée de la manière suivante :

$$\chi_n^2 = \sum_{i=1}^h \sum_{j=1}^k \frac{\left(\left| N_{i,j} - \frac{N_{i,\bullet} \times N_{\bullet,j}}{n} \right| - \frac{1}{2} \right)^2}{\frac{N_{i,\bullet} \times N_{\bullet,j}}{n}}$$

où $N_{i,j}$ est le nombre aléatoire de couple (X_i, Y_j) qui dans un échantillon indépendant identiquement distribué de taille n appartient à la classe $C_{i,j}$, $N_{i,\bullet} = \sum_{j=1}^k N_{i,j}$ et $N_{\bullet,j} = \sum_{i=1}^h N_{i,j}$.

Attention, on n'utilisera la correction que si l'une des fréquences attendues est strictement inférieure à 5. Si au contraire elles sont toutes supérieures ou égales à 5, on montre que la correction de Yates ne modifie que peu la valeur de la réalisation de la statistique du test. Comme le signale J Bouyer, [3], ce point de vue, c'est-à-dire le fait de réserver cette correction à de petits échantillons, n'est pas partagé par tous les auteurs. Si l'un des effectifs théoriques est inférieur à 3, on n'a pas d'autre solution que d'appliquer le test exact de Fisher décrit au paragraphe 8. Plus généralement, P. Dagnélie, [6], conseille l'utilisation systématique du test exact de Fisher lorsque l'effectif total n est inférieur ou égal à 40.

7.5. Étude des résidus

Lorsque l'hypothèse d'indépendance est vérifiée, les termes dont les carrés sont les contributions à la valeur χ_n^2 :

$$\frac{N_{i,j} - \frac{N_{i,\bullet} \times N_{\bullet,j}}{n}}{\sqrt{\frac{N_{i,\bullet} \times N_{\bullet,j}}{n}}}$$

sont approximativement des variables aléatoires qui suivent des lois normales centrées réduites. On a montré qu'une meilleure approximation de la loi normale centrée réduite est obtenue si l'on considère les valeurs ci-dessus et que l'on les divise par des estimations des écarts types correspondants :

$$\sqrt{\left(1 - \frac{N_{i,\bullet}}{n}\right) \left(1 - \frac{N_{\bullet,j}}{n}\right)}.$$

On obtient finalement les **écarts réduits** :

$$\frac{N_{i,j} - \frac{N_{i,\bullet} \times N_{\bullet,j}}{n}}{\sqrt{\left(1 - \frac{N_{i,\bullet}}{n}\right) \left(1 - \frac{N_{\bullet,j}}{n}\right) \frac{N_{i,\bullet} \times N_{\bullet,j}}{n}}}$$

On peut alors étudier les écarts réduits comme on le ferait pour des résidus obtenus après une régression par exemple : étude de la normalité, via une droite de Henry par exemple, identification d'éventuelles valeurs aberrantes ou non représentatives, voir le cours sur ce sujet.

Remarques :

- Il existe une version « exacte » de ce test.
- Pour des modélisations plus complexes et, par exemple, l'étude de corrélations partielles, on utilise le **modèle log-linéaire**, voir le cours associé.

8. Test exact de Fisher

Pour présenter le principe de ce test, on commence par étudier le cas de deux variables X et Y , ayant deux modalités, puis on traite le cas général. Pour plus de détails sur les tests exacts sur les tableaux de contingence, on se reportera à l'article de A. Agresti [1].

8.1. Deux variables à deux modalités

Soit X et Y deux variables quantitatives discrètes ou qualitatives. On suppose dans ce paragraphe que X et Y ne peuvent prendre que 2 valeurs différentes. On se donne un

n -échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ indépendant et identiquement distribué suivant la loi de (X, Y) ainsi qu'un échantillon (\mathbf{x}, \mathbf{y}) formé d'une réalisation de chaque (X_i, Y_i) , $1 \leq i \leq n$. Enfin on note \mathbf{x} l'échantillon des réalisations de X et \mathbf{y} l'échantillon des réalisations de Y . On considère le tableau des effectifs $n_{i,j}$ suivant chacune des 2 modalités de X et des 2 modalités de Y apparaissant dans l'échantillon (\mathbf{x}, \mathbf{y}) :

$X \backslash Y$	1	2	Total
1	$n_{1,1}$	$n_{1,2}$	$n_{1,\bullet}$
2	$n_{2,1}$	$n_{2,2}$	$n_{2,\bullet}$
Total	$n_{\bullet,1}$	$n_{\bullet,2}$	n

Dans ce tableau les marges $(n_{1,\bullet}, n_{2,\bullet}, n_{\bullet,1}, n_{\bullet,2})$ n'apportent pas d'information sur l'éventuelle dépendance de X et de Y . En effet elles n'indiquent que la répartition des effectifs entre les deux modalités de X , indépendamment de la valeur de Y , et la répartition des effectifs entre les deux modalités de Y , indépendamment de la valeur de X . Ce sont les valeurs présent par $n_{1,1}, n_{1,2}, n_{2,1}$ et $n_{2,2}$ qui servent pour étudier la dépendance de X et de Y .

L'idée du test exact de Fisher est de considérer l'ensemble Γ des tableaux ayant les mêmes marges $(n_{1,\bullet}, n_{2,\bullet}, n_{\bullet,1}, n_{\bullet,2})$ que le tableau des effectifs observés ci-dessus :

$$\Gamma_{\substack{(n_{1,\bullet}, n_{2,\bullet}) \\ (n_{\bullet,1}, n_{\bullet,2})}} = \left\{ \begin{array}{|c|c|} \hline a & b \\ \hline c & d \\ \hline \end{array} \right\}, \text{ avec } a + b = n_{1,\bullet}, c + d = n_{2,\bullet}, a + c = n_{\bullet,1} \text{ et } b + d = n_{\bullet,2} \left. \right\}.$$

Les marges $(n_{1,\bullet}, n_{2,\bullet}, n_{\bullet,1}, n_{\bullet,2})$ étant fixées, la connaissance de a détermine par différence les valeurs de b, c et d :

$$\Gamma_{\substack{(n_{1,\bullet}, n_{2,\bullet}) \\ (n_{\bullet,1}, n_{\bullet,2})}} = \left\{ \gamma(a) = \begin{array}{|c|c|} \hline a & n_{1,\bullet} - a \\ \hline n_{\bullet,1} - a & n_{\bullet,2} - n_{\bullet,1} + a \\ \hline \end{array}, \text{ avec } 0 \leq a \leq \inf(n_{1,\bullet}, n_{\bullet,1}) \right\}.$$

On montre que, sous l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendantes », la probabilité d'obtenir un tableau de type $\gamma(a)$ suit une loi hypergéométrique, voir [7] et [6] :

$$\begin{aligned} \mathbb{P}[n_{1,1} = a] &= \frac{C_{n_{1,\bullet}}^a \times C_{n_{2,\bullet}}^b}{C_{n_{\bullet,2} + n_{\bullet,1}}^{n_{1,\bullet}}} \\ &= \frac{C_{a+c}^a \times C_{b+d}^b}{C_{a+b+c+d}^{a+b}} \\ &= \frac{(a+c)! \times (b+d)! \times (a+b)! \times (c+d)!}{a! \times b! \times c! \times d! \times (a+b+c+d)!} \\ &= \frac{n_{\bullet,1}! \times n_{\bullet,2}! \times n_{1,\bullet}! \times n_{2,\bullet}!}{a! \times (n_{\bullet,1} - a)! \times (n_{1,\bullet} - a)! \times (n_{\bullet,2} - n_{\bullet,1} + a)! \times n_{\bullet,2}!} \end{aligned}$$

où $z!$ désigne la factorielle du nombre entier positif z qui vaut $z! = z \times (z - 1) \times \dots \times 1$ avec la convention $0! = 1$.

Procédure de test :

On teste l'hypothèse suivante :

$$\boxed{\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}}$$

contre

$$\boxed{\mathcal{H}_1 : X \text{ et } Y \text{ sont liées.}}$$

Le principe du test consiste à évaluer la probabilité de rencontrer une distribution aussi anormale ou plus anormale que celle que l'on a observée.

Supposons par exemple que $\frac{a}{a+c} \leq \frac{b}{b+d}$. On additionne alors les probabilités d'obtenir des tableaux ayant des valeurs de $n_{1,1}$ comprises en 0 et a .

Supposons par contre que $\frac{a}{a+c} \geq \frac{b}{b+d}$. On additionne alors les probabilités d'obtenir des tableaux ayant des valeurs de $n_{1,1}$ comprises en b et $\inf(a+b, a+c)$.

Le test bilatéral de niveau α conduit alors au rejet de l'hypothèse nulle \mathcal{H}_0 quand la probabilité totale calculée par la méthode ci-dessus est inférieure ou égale à $\alpha/2$.

Il s'agit en fait d'un niveau de signification maximal α car la distribution hypergéométrique sur laquelle repose le test est discrète.

- La dénomination « exact » du test vient du fait que l'on ne fait appel à aucune approximation pour calculer la p -valeur du test.
- Le niveau de signification du test est d'au plus α , il se peut qu'il soit beaucoup plus faible, ce qui fait de ce test un test conservatif.
- Ce test peut également être utilisé pour comparer deux pourcentages.

Exemple 8.1. On cherche à savoir s'il existe une association entre la couleur des yeux, clairs ou foncés, et la couleur des cheveux, blonds ou bruns, chez les êtres humains. On considère un groupe de 20 individus de sexe féminin ou masculin.

Cheveux \ Yeux	Clairs	Foncés	Total
Blonds	5	1	6
Bruns	0	14	14
Total	5	15	20

On souhaite tester l'hypothèse :

\mathcal{H}_0 : La couleur des yeux est indépendante de celle des cheveux

contre

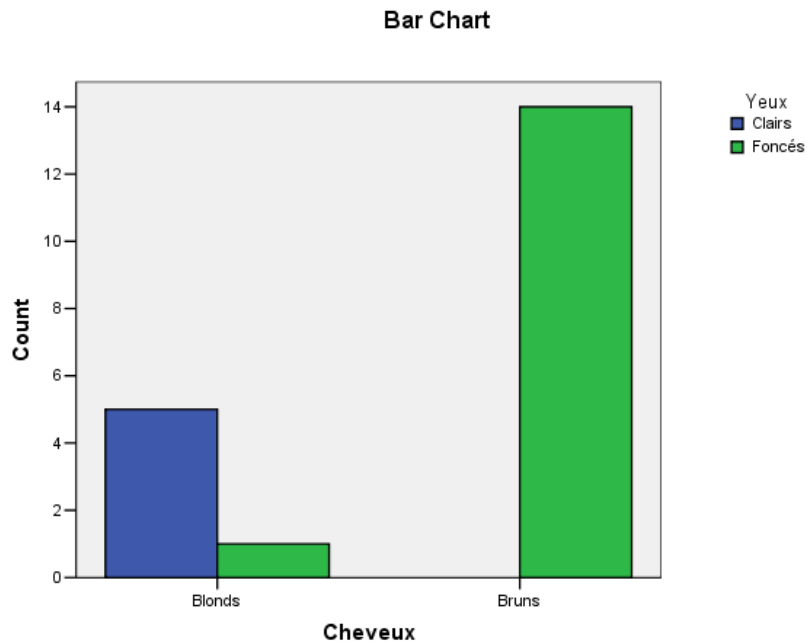
\mathcal{H}_1 : La couleur des yeux n'est pas indépendante de celle des cheveux.

On est dans la situation où $\frac{a}{a+c} \geq \frac{b}{b+d}$ car $\frac{5}{5} = 1 \geq \frac{1}{15}$. On cherche ainsi la p -valeur associée au test en additionnant la probabilité d'obtenir des tableaux pour lesquels $n_{1,1}$ est compris entre 5 et $\inf(6, 5) = 5$. La p -valeur est donc égale à la probabilité d'obtenir un tableau pour lequel $n_{1,1} = 5$ connaissant les marges (6, 14, 5, 15) :

$$\begin{aligned}
 \mathbb{P}[n_{1,1} = 5] &= \frac{(5+0)! \times (1+14)! \times (5+1)! \times (0+14)!}{5! \times 1! \times 0! \times 14! \times 20!} \\
 &= \frac{5! \times 15! \times 6! \times 14!}{5! \times 1! \times 0! \times 14! \times 20!} \\
 &= \frac{5 \times 4 \times 3 \times 2 \times 1 \times 6}{20 \times 19 \times 18 \times 17 \times 16} \\
 &= \frac{1}{2 \times 19 \times 17 \times 4} \\
 &= \frac{1}{2584}
 \end{aligned}$$

Donc la p -valeur associée au test vaut : $\frac{1}{2584}$. Or pour un niveau de signification de $\alpha = 5\%$, c'est-à-dire, $\alpha = \frac{1}{20}$ on a obtenu une p -valeur inférieure ou égale au seuil du test : le test est significatif. On rejette l'hypothèse nulle \mathcal{H}_0 et on décide l'hypothèse alternative \mathcal{H}_1 : « La couleur des yeux n'est pas indépendante de celle des cheveux ». Il existe une association significative au seuil $\alpha = 5\%$ entre la couleur des yeux et celle des cheveux.

Voici ce que l'on obtient lorsque l'on utilise SPSS 13.0 :



Cheveux * Yeux Crosstabulation

		Yeux		Total
		Clairs	Foncés	
Cheveux	Blonds	5	1	6
	Bruns	0	14	14
Total		5	15	20

Chi-Square Tests

	Value	Exact Sig. (2-sided)	Exact Sig. (1-sided)
Fisher's Exact Test		,000	,000
N of Valid Cases	20		

Ces résultats sont en accord avec les calculs que l'on a réalisés à la main ci-dessus : le test

est significatif au seuil $\alpha = 5\%$.

Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	15,556(b)	1	,000
Continuity Correction(a)	11,429	1	,001
N of Valid Cases	20		

a Computed only for a 2x2 table

b 3 cells (75,0%) have expected count less than 5.

The minimum expected count is 1,50.

On constate que SPSS propose un test du χ^2 , avec ou sans la correction de Yates mentionnée au paragraphe 7.4. SPSS signale que 75 % des cellules ont des effectifs attendus de moins de 5 ce qui rend impossible l'utilisation du test du χ^2 dans notre situation : on ne pouvait qu'utiliser le test exact de Fisher.

.....

8.2. Deux variables ayant un nombre fini quelconque de modalités

On reprend l'idée du paragraphe 8.1 ci-dessus et on l'étend au cas où les deux variables étudiées ont un nombre fini quelconque, mais supérieur à deux, de modalités. Cette extension a été réalisée en premier par G. H. Freeman et J.H. Halton en 1951, voir [9], c'est pourquoi ce test est aussi parfois appelé test de Freeman-Halton ou de Fisher-Freeman-Halton.

Soit X et Y deux variables quantitatives discrètes ou qualitatives. On se donne un n -échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ indépendant et identiquement distribué suivant la loi de (X, Y) ainsi qu'un échantillon (\mathbf{x}, \mathbf{y}) formé d'une réalisation de chaque (X_i, Y_i) , $1 \leq i \leq n$. Enfin on note \mathbf{x} l'échantillon des réalisations de X et \mathbf{y} l'échantillon des réalisations de Y . On considère le tableau des effectifs $n_{i,j}$ suivant chacune des h modalités de X et des k modalités de Y apparaissant dans l'échantillon (\mathbf{x}, \mathbf{y}) :

$X \backslash Y$	1	...	j	...	k	Total
1	$n_{1,1}$...	$n_{1,j}$...	$n_{1,k}$	$n_{1,\bullet}$
\vdots	\vdots		\vdots		\vdots	\vdots
i	$n_{i,1}$...	$n_{i,j}$...	$n_{i,k}$	$n_{i,\bullet}$
\vdots	\vdots		\vdots		\vdots	\vdots
h	$n_{h,1}$...	$n_{h,j}$...	$n_{h,k}$	$n_{h,\bullet}$
Total	$n_{\bullet,1}$...	$n_{\bullet,j}$...	$n_{\bullet,k}$	n

Dans ce tableau les marges $(n_{1,\bullet}, \dots, n_{h,\bullet}, n_{\bullet,1}, \dots, n_{\bullet,k})$ n'apportent pas d'information sur l'éventuelle dépendance de X et de Y . En effet elles n'indiquent que la répartition des effectifs entre les h modalités de X , indépendamment de la valeur de Y , et la répartition des effectifs entre les k modalités de Y , indépendamment de la valeur de X . Ce sont les valeurs prises par $n_{1,1}, n_{1,2}, \dots, n_{1,k}, n_{2,1}, \dots, n_{2,k}, \dots, n_{i,j}, \dots, n_{h,k-1}$ et $n_{h,k}$ qui servent pour étudier la dépendance de X et de Y .

L'idée du test exact de Fisher est de considérer l'ensemble Γ des tableaux ayant les mêmes marges $(n_{1,\bullet}, \dots, n_{h,\bullet}, n_{\bullet,1}, \dots, n_{\bullet,k})$ que le tableau des effectifs observés ci-dessus :

$$\Gamma_{\substack{(n_{1,\bullet}, \dots, n_{h,\bullet}) \\ (n_{\bullet,1}, \dots, n_{\bullet,k})}} = \left\{ \begin{array}{|c|c|c|c|c|} \hline n_{1,1} & \cdots & n_{1,j} & \cdots & n_{1,k} \\ \hline \vdots & & \vdots & & \vdots \\ \hline n_{i,1} & \cdots & n_{i,j} & \cdots & n_{i,k} \\ \hline \vdots & & \vdots & & \vdots \\ \hline n_{h,1} & \cdots & n_{h,j} & \cdots & n_{h,k} \\ \hline \end{array} \right\}, \text{ avec } \left\{ \begin{array}{l} n_{1,1} + \dots + n_{1,k} = n_{1,\bullet} \\ \vdots = \vdots \\ n_{1,1} + \dots + n_{h,1} = n_{h,\bullet} \\ n_{h,1} + \dots + n_{h,k} = n_{\bullet,1} \\ \vdots = \vdots \\ n_{1,k} + \dots + n_{h,k} = n_{\bullet,k} \end{array} \right\}.$$

Les marges $(n_{1,\bullet}, \dots, n_{h,\bullet}, n_{\bullet,1}, \dots, n_{\bullet,k})$ étant fixées, la connaissance de $hk - h - k + 1 = (h-1)(k-1)$ valeurs détermine celle de toutes les autres.

On montre que, sous l'hypothèse nulle \mathcal{H}_0 : « X et Y sont indépendantes », la probabilité d'obtenir un tableau γ appartenant à Γ suit une loi hypergéométrique généralisée, voir [7] :

$$\begin{aligned} \mathbb{P}[\gamma] &= \frac{(n_{1,1} + \dots + n_{1,k})! \times \dots \times (n_{h,1} + \dots + n_{h,k})! \times (n_{1,1} + \dots + n_{h,1})! \times (n_{1,k} + \dots + n_{h,k})!}{n_{1,1}! \times \dots \times n_{h,k}! (n_{1,1} + \dots + n_{h,k})!} \\ &= \frac{\prod_{i=1}^h n_{i,\bullet}! \prod_{j=1}^k n_{\bullet,j}!}{\left(\prod_{i=1}^h \left(\prod_{j=1}^k n_{i,j}! \right) \right) \left(\sum_{i=1}^h \sum_{j=1}^k n_{i,j} \right)!}, \end{aligned}$$

où $z!$ désigne la factorielle du nombre entier positif z qui vaut $z! = z \times (z-1) \times \dots \times 1$ avec la convention $0! = 1$.

Procédure de test :

On teste l'hypothèse suivante :

$$\boxed{\mathcal{H}_0 : X \text{ et } Y \text{ sont indépendantes}}$$

contre

$$\boxed{\mathcal{H}_1 : X \text{ et } Y \text{ sont liées.}}$$

Le principe du test consiste à évaluer la probabilité de rencontrer une distribution aussi anormale ou plus anormale dans un tableau γ de Γ que celle que l'on a observée γ_{obs} :

$$\begin{aligned} p_2 &= \sum_{D(\gamma) \geq D(\gamma_{obs})} \mathbb{P}[\gamma] \\ &= \mathbb{P}[D(\gamma) \geq D(\gamma_{obs})]. \end{aligned}$$

où la fonction $\gamma \rightarrow D(\gamma)$ est définie pour tout $\gamma \in \Gamma$ par :

$$D(\gamma) = -2 \log \left(\frac{(2\pi)^{\frac{(h-1)(k-1)}{2}}}{\frac{(hk-1)^{\frac{h-1}{2}}}{N^{\frac{h-1}{2}}} \prod_{i=1}^h n_{i,\bullet}^{\frac{k-1}{2}} \prod_{j=1}^k n_{\bullet,j}^{\frac{h-1}{2}}} \mathbb{P}[\gamma] \right).$$

On définit alors la région critique du test comme la partie Γ^* de l'ensemble de référence Γ :

$$\Gamma^* = \{y \in \Gamma \text{ tels que } D(y) \geq D(x)\}.$$

Le test bilatéral de niveau α conduit alors au rejet de l'hypothèse nulle \mathcal{H}_0 quand la probabilité totale p_2 calculée par la méthode ci-dessus est inférieure ou égale à α^9 .

Il s'agit en fait d'un niveau de signification maximal α car la distribution hypergéométrique généralisée sur laquelle repose le test est discrète.

- La dénomination « exact » du test vient du fait que l'on ne fait appel à aucune approximation pour calculer la p -valeur du test.
- Le niveau de signification du test est d'au plus α , il se peut qu'il soit beaucoup plus faible, ce qui fait de ce test un test conservatif.
- Ce test peut également être utilisé pour comparer k pourcentages sur l populations.

⁹On remarque que la p -valeur calculée au paragraphe 8.1 pour le cas d'un test exact de Fisher sur un tableau de taille 2×2 n'est pas obtenue de la même manière qu'ici. Ceci est dû au fait que dans le cas du paragraphe 8.1, on pouvait déduire toutes les valeurs du tableaux en connaissant uniquement ses marges et une de valeurs des effectif $n_{1,1}$, $n_{1,2}$, $n_{2,1}$ ou $n_{2,2}$ et ainsi ranger « naturellement » les tableaux dans l'ordre croissant ce qui permettait alors de définir facilement ce que l'on entendait par une distribution aussi ou plus anormale que celle du tableau observée.

9. Tests de multinormalité

9.1. Utilisation de tests multiples unidimensionnels

On montre, voir le livre de J.-Y. Ouvrard [10], qu'un **vecteur aléatoire** \mathbf{X} à n composantes (X_1, \dots, X_n) est **gaussien**, c'est-à-dire qu'il suit une loi normale de dimension n , parfois également appelée loi multinormale sur \mathbb{R}^n , **si et seulement si toutes** les combinaisons linéaires de ses composantes, voir le paragraphe 1.4.1 pour une définition de cette notion mathématique, sont des variables aléatoires gaussiennes sur \mathbb{R} , c'est-à-dire qu'elles suivent des lois normales unidimensionnelles.

Le résultat ci-dessus permet donc apparemment de ramener au cas de la normalité d'une **variable aléatoire réelle** Y l'étude de la normalité d'un **vecteur aléatoire** \mathbf{X} . Malheureusement, si l'on voulait se servir dans la pratique de ce résultat, il faudrait tester une infinité de combinaisons linéaires de composantes de \mathbf{X} ! C'est bien sûr impossible non seulement pour des raisons de temps de calcul mais aussi pour des problèmes de risque de première espèce α . En effet une telle démarche, si l'on pouvait l'accomplir se solderait par un risque de première espèce global α_{global} de 100 % quelque soit le risque individuel α_{ind} non nul fixé au départ¹⁰.

On a utilisé dans les exemples 1.3, 1.4 et 1.5 la propriété exposée ci-dessus pour vérifier sommairement la validité de l'hypothèse de multinormalité d'un vecteur aléatoire \mathbf{X} . Ne pouvant réaliser des tests pour toutes les combinaisons linéaires des composantes (X_1, \dots, X_n) des vecteurs que l'on a étudiés, on s'est contenté de concentrer notre étude sur certaines combinaisons linéaires :

- Les composantes du vecteur \mathbf{X} : X_1, \dots, X_n .
- Les sommes des composantes du vecteur \mathbf{X} : $X_1 + X_2, X_1 + X_3, \dots, X_1 + X_n, X_2 + X_3, \dots, X_i + X_j, \dots, X_{n-1} + X_n$, où $1 \leq i < j \leq n$.

Si le vecteur \mathbf{X} n'a que deux composantes (X_1, X_2) on peut rajouter la combinaison linéaire suivante :

- La différence des composantes du vecteur \mathbf{X} : $X_1 - X_2$.

Ainsi si l'on considère un vecteur à deux composantes (X_1, X_2) on réalise 4 tests de normalité individuels, celui de X_1 , celui de X_2 , celui de $X_1 + X_2$ et celui de $X_1 - X_2$. On adapte¹¹ donc la valeur du risque de première espèce individuel α_{ind} pour obtenir un

¹⁰En effet on doit faire une **infinité** de tests et $\alpha_{global} = \lim_{n \rightarrow \infty} (1 - (1 - \alpha_{ind})^n) = 1$.

¹¹On utilise ici une procédure de Sidak [7], il ne s'agit que d'une possibilité parmi d'autres, on aurait aussi bien pu se servir de l'inégalité de Bonferroni, voir [7] et [6] et la recommandation de H.C. Thode dans [13] à ce sujet. On utilise ces procédures car le calcul du niveau exact du test est très complexe.

risque de première espèce global de α_{global} de la manière suivante :

$$\begin{aligned}\alpha_{global} &= 1 - (1 - \alpha_{ind})^4, \\ \alpha_{ind} &= 1 - \sqrt[4]{(1 - \alpha_{global})}.\end{aligned}$$

Par exemple si on souhaite obtenir un $\alpha_{global} = 5 \%$ on trouve $\alpha_{ind} = 1,27 \%$.

Plus généralement, supposons que n est supérieur ou égal 3 et considérons un vecteur à n composantes (X_1, \dots, X_n) . On réalise alors $2^n - 1$ tests de normalité individuels, ceux de X_1, \dots, X_n , ceux de $X_1 + X_2, \dots, X_{n-1} + X_n$, ceux des sommes de trois composantes, \dots , ceux des sommes de j composantes, \dots et celui de la somme de toutes les composantes $X_1 + \dots + X_n$. On adapte, voir la note de bas de page numéro 11, donc la valeur du risque de première espèce individuel α_{ind} pour obtenir un risque de première espèce global de α_{global} de la manière suivante :

$$\begin{aligned}\alpha_{global} &= 1 - (1 - \alpha_{ind})^{\frac{k(k+1)}{2}}, \\ \alpha_{ind} &= 1 - \sqrt{\frac{k(k+1)}{2}(1 - \alpha_{global})}.\end{aligned}$$

Par exemple si on souhaite obtenir un $\alpha_{global} = 5 \%$ on trouve :

k	$\frac{k(k+1)}{2}$	$\alpha_{ind} \%$
3	7	0,6391
4	15	0,3201
5	31	0,1602
10	1023	0,0050
20	1048576	$0,4891 \times 10^{-5}$
50	1125899906842624	$0,4556 \times 10^{-14}$
100	1267650600228229401496703205376	$0,4046 \times 10^{-29}$

$\alpha_{ind} = 1,27 \%$.

On voit que cette procédure n'est adaptée que pour des cas pratiques de petite dimension $n \leq 5$. Dans ce cas elle permet d'évaluer rapidement la possible multinormalité de la distribution sans avoir à recourir à des outils spécifiques qui sont plus complexes et donc plus long, voir le paragraphe 9.2 suivant.

Exemple 9.1.

Pour des applications de cette procédure voir les exemples 1.3, 1.4 et 1.5.

.....

9.2. Une test basée sur la modification du test de Shapiro-Wilk unidimensionnel

Malkovitch (1971), Malkovitch and Afifi (1973) ont présenté un test de multinormalité pour un vecteur aléatoire \mathbf{X} de dimension n basé sur une modification du test de Shapiro-Wilk. Puisque le vecteur \mathbf{X} est normal si et seulement si toutes les combinaisons linéaires des composantes (X_1, \dots, X_n) de \mathbf{X} sont normales le test de multinormalité est accepté si

$$\inf_{\substack{\text{Sur les combinaisons linéaires } Y \\ \text{des composantes de } \mathbf{X}}} W_Y \geq \text{Valeur critique du test,}$$

où W_Y est la valeur de la réalisation de la statistique de Shapiro-Wilk calculée pour Y .

On détermine alors une combinaison linéaire susceptible de réaliser le minimum apparaissant dans le membre de droite de l'équation ci-dessus.

Les détails étant complexes, on n'exposera pas ici la procédure du test. Pour plus de détails voir l'ouvrage de H.C. Thode [13].

Malheureusement, ce test n'est pas disponible sous Minitab, ni sous la plupart des logiciels de statistique.

9.3. Asymétrie et aplatissement multivariés de Mardia

Mardia (1970) et (1974) a proposé un test de multinormalité basé sur l'asymétrie et l'aplatissement. Ce test est par exemple disponible sous S-Plus, R et la plupart des logiciels de statistique. On se donne donc un vecteur aléatoire \mathbf{X} de dimension p .

On considère un n -échantillon $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ indépendant et identiquement distribué suivant la loi de \mathbf{X} . On se donne également \mathbf{x} une réalisation du n -échantillon $(\mathbf{X}_1, \dots, \mathbf{X}_n)$.

9.4. Centrer-réduire.

Dans cette section on considère un vecteur aléatoire \mathbf{X} qui admet une moyenne $\boldsymbol{\mu}$ et une matrice de variance-covariance $\boldsymbol{\Sigma}$.

On commence par calculer les valeurs centrées-réduites¹² :

$$\mathbf{Z}_i = \boldsymbol{\Sigma}^{-\frac{1}{2}} (\mathbf{X}_i - \boldsymbol{\mu}).$$

¹²« Scaled residuals » en anglais. Cette transformation est la version multivariée du centrer-réduire que vous connaissez en dimension 1 : on soustrait la moyenne $\boldsymbol{\mu}$ puis on divise par la racine carrée de la variance notée $\boldsymbol{\Sigma}^{-\frac{1}{2}}$.

Si l'on suppose que la loi de \mathbf{X} est multinormale, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, les vecteurs aléatoires \mathbf{Z}_i suivent alors une loi normale multidimensionnelle « centrée-réduite » $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$.

On calcule alors les rayons au carré¹³ de la manière suivante :

$$\begin{aligned} R_i^2 &= \mathbf{Z}_i^T \mathbf{Z}_i \\ &= (\mathbf{X}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\mu}), \end{aligned}$$

$(\mathbf{X}_i - \boldsymbol{\mu})^T$ est la transposée de $(\mathbf{X}_i - \boldsymbol{\mu})$.

Si l'on suppose que la loi de \mathbf{X} est multinormale, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, la variable aléatoire R_i^2 suit alors une loi du χ^2 à m degrés de liberté.

Or les paramètres de la loi, $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, sont généralement inconnus. On doit donc les estimer pour pouvoir calculer les valeurs centrées réduites ainsi que les rayons au carré. On procède donc ainsi :

- i- On estime la moyenne $\boldsymbol{\mu}$ de \mathbf{X} à l'aide de l'estimateur défini au paragraphe 1.3.2 ; on obtient la valeur $\widehat{\boldsymbol{\mu}}(\mathbf{x})$. Puis on estime la matrice de variance-covariance $\boldsymbol{\Sigma}$ à l'aide de l'estimateur défini au paragraphe 1.3.2 ; on obtient la valeur $\widehat{\boldsymbol{\Sigma}}(\mathbf{x})$.
- ii- On estime alors les valeurs centrées réduites et les rayons au carré :

$$\begin{aligned} \widehat{\mathbf{Z}}_i(\mathbf{x}) &= \widehat{\boldsymbol{\Sigma}}(\mathbf{x})^{-\frac{1}{2}} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}(\mathbf{x})). \\ \widehat{R}_i^2(\mathbf{x}) &= \widehat{\mathbf{Z}}_i(\mathbf{x})^T \widehat{\mathbf{Z}}_i(\mathbf{x}) \\ &= (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}(\mathbf{x}))^T \widehat{\boldsymbol{\Sigma}}(\mathbf{x})^{-1} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}(\mathbf{x})). \end{aligned}$$

De nombreuses procédures d'études de la multinormalité se basent sur les propriétés des rayons au carré et des rayons au carré généralisés définis par :

$$\begin{aligned} R_{i,j} &= \mathbf{Z}_i^T \mathbf{Z}_j \\ &= (\mathbf{X}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X}_j - \boldsymbol{\mu}), \end{aligned}$$

On obtient des estimations de ces variables ainsi :

$$\begin{aligned} \widehat{R}_{ij}(\mathbf{x}) &= \widehat{\mathbf{Z}}_i(\mathbf{x})^T \widehat{\mathbf{Z}}_j(\mathbf{x}) \\ &= (\mathbf{x}_i - \widehat{\boldsymbol{\mu}}(\mathbf{x}))^T \widehat{\boldsymbol{\Sigma}}(\mathbf{x})^{-1} (\mathbf{x}_j - \widehat{\boldsymbol{\mu}}(\mathbf{x})). \end{aligned}$$

On remarque que $R_{i,i} = R_i^2$ et que $\widehat{R}_{i,i}(\mathbf{x}) = \widehat{R}_i^2(\mathbf{x})$.

¹³« Squared radii » en anglais.

9.5. Asymétrie et aplatissement

La mesure d'asymétrie multivariée d'un vecteur aléatoire \mathbf{X} est :

$$b_{1,p} = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n R_{ij}^3.$$

On l'estime en utilisant les estimateurs de $R_{i,j}$ du paragraphe 9.4 :

$$\widehat{b}_{1,p}(\mathbf{x}) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \widehat{R}_{ij}(\mathbf{x})^3.$$

La mesure d'aplatissement multivariée d'un vecteur aléatoire \mathbf{X} est :

$$b_{2,p} = \frac{1}{n} \sum_{i=1}^n R_{ii}^2.$$

On l'estime en utilisant les estimateurs de $R_{i,j}$ du paragraphe 9.4 :

$$\widehat{b}_{2,p}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \widehat{R}_{ii}(\mathbf{x})^2.$$

9.6. Étude sous l'hypothèse de multinormalité de \mathbf{X} .

Si l'on se place sous l'hypothèse nulle \mathcal{H}_0 : « La loi de \mathbf{X} est multinormale », $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, on a les résultats suivants :

-i- L'espérance de $\widehat{b}_{1,p}$ est :

$$\mathbb{E} \left[\widehat{b}_{1,p} \right] = \frac{p(p+2)((n+1)(p+1) - 6)}{(n+1)(n+3)}.$$

-ii- On a un résultat sur la distribution asymptotique de $\widehat{b}_{1,p}$:

$$\frac{n\widehat{b}_{1,p}}{6} \text{ suit asymptotiquement une loi du } \chi^2_{\frac{p(p+1)(p+2)}{6}}.$$

-iii- L'espérance de $\widehat{b}_{2,p}$ est :

$$\mathbb{E} \left[\widehat{b}_{2,p} \right] = \frac{p(p+2)(n-1)}{(n+1)}.$$

-iv- La variance de $\widehat{b}_{2,p}$ est :

$$\text{Var} \left[\widehat{b}_{2,p} \right] = \frac{8p(m+2)(n-3)(n-m-1)(n-m+1)}{(n+1)^2(n+3)(n+5)}.$$

-v- On a un résultat sur la distribution asymptotique de $\widehat{b}_{2,p}$:

$$\widehat{b}_{2,p} \text{ suit asymptotiquement une loi } \mathcal{N}\left(p(p+2), \frac{8p(p+2)}{n}\right).$$

Certaines des valeurs critiques sont tabulées, sinon l'on utilise l'approximation suivante pour $\widehat{b}_{1,p}$. On pose :

$$A = \frac{(p+1)(n+1)(n+3)\widehat{b}_{1,p}}{((n+1)(p+1)-6)6}.$$

A suit alors une loi approximativement une loi du χ^2 à $\frac{p(p+1)(p+2)}{6}$ degrés de liberté. Ce test est disponible sous SAS et R.

Il est recommandé par H.C. Thode dans [13] pour sa bonne puissance.

9.7. Un test basé sur la fonction caractéristique empirique

Il s'agit d'un test proposé par Henze et Zirkler (1990).

On mentionne ici ce test pour ces bonnes propriétés de puissance. En particulier celles de $T_{0,5}$.

La statistique du test est alors :

$$D_{n,\beta} = \frac{1}{n^2} \sum_{j=1}^n \sum_{k=1}^n \exp\left(-\frac{\beta^2}{2} \|\widehat{\mathbf{Z}}_i - \widehat{\mathbf{Z}}_j\|^2\right) + (1 + 2\beta^2)^{-\frac{m}{2}} \\ - \frac{2(1 + \beta^2)^{-\frac{m}{2}}}{n} \sum_{j=1}^n \exp\left(-\frac{\beta^2}{2(1 + \beta^2)} r_j^2\right).$$

Ce test est disponible sous SAS.

Références

- [1] A. Agresti. A survey of exact inference for contingency tables. *Statistical Science*, 7(1) :131–152, 1992.
- [2] A. Boomsma. Comparing approximations of confidence intervals for the product-moment correlation coefficient. *Statistica Neerlandica*, 31 :179–186, 1977.
- [3] J. Bouyer. *Méthodes Statistiques*. Editions INSERM, 1996.
- [4] P. Chapouille. *Planification et analyse des expériences*. Masson, Paris, 1973.

- [5] R. Christensen. *Linear Models for Multivariate, Time Series, and Spatial Data*. Springer Texts in Statistics. Springer-Verlag, 1991.
- [6] P. Dagnélie. *Statistique Théorique et Appliquée*, volume 2. De Boeck & Larcier, Bruxelles, 1998.
- [7] P. Dagnélie. *Statistique Théorique et Appliquée*, volume 1. De Boeck & Larcier, Bruxelles, 1998.
- [8] B. Falissard. *Comprendre et utiliser les statistiques dans les sciences de la vie*. Abrégés. Masson, Paris, 3^{ème} édition, 2005.
- [9] G. H. Freeman and J.H. Halton. Note on an exact treatment of contingency, goodness of fit and other problems of significance. *Biometrika*, 38 :141–149, 1952.
- [10] J.-Y. Ouvrard. *Probabilités*, volume 2. Cassini, Paris, 2000.
- [11] G. Parreins. *Techniques Statistiques : moyens rationnels de choix et de décision*. Dunod technique, Paris, 1974.
- [12] G. Pupion and P.-C. Pupion. *Tests non paramétriques*. Statistique mathématique et probabilité. Economica, Paris, 1998.
- [13] H.C. Thode. *Testing for Normality*. Number 164 in Statistics : textbooks and monographs. Marcel Dekker, New-York, 2002.
- [14] E. Weisstein. Correlation coefficient-bivariate normal distribution from mathworld—a wolfram web resource. Lien internet : <http://mathworld.wolfram.com/Correlation-CoefficientBivariateNormalDistribution.html>.

Table des matières

1	Le cas paramétrique sous hypothèse de multinormalité	1
1.1	Coefficient de corrélation simple	1
1.2	Le cas bidimensionnel	2
1.2.1	Loi normale bidimensionnelle	2
1.2.2	Estimation des paramètres	4
1.2.3	Procédure de test	9
1.2.4	Remarques sur les liaisons	10
1.3	Le cas général ($n \geq 2$)	11
1.3.1	Loi multinormale	11
1.3.2	Estimation	12
1.3.3	Test de l'hypothèse $\rho_{ij} = 0$	14
1.3.4	Test de l'hypothèse $\rho_{ij} = \rho_0, \rho_0 \neq 0$	15
1.4	Corrélation multiple	18
1.4.1	Définition	18
1.4.2	Estimation	19
1.4.3	Asymptotique	19

1.4.4	Test de l'hypothèse $R(X_1, \mathbf{X}_2) = 0$	20
1.4.5	Test de l'hypothèse $R(X_1, \mathbf{X}_2) = R_0, R_0 \neq 0$	20
1.5	Corrélation partielle	21
1.5.1	Définition	21
1.5.2	Estimation	22
1.5.3	Asymptotique	22
1.5.4	Test de l'hypothèse $\rho_{i,j q+1,\dots,m} = 0$	23
1.5.5	Test de l'hypothèse $\rho_{i,j q+1,\dots,m} = \rho_0, \rho_0 \neq 0$	23
1.5.6	Cas de trois variables réelles	24
2	Mesures non paramétriques	38
3	La statistique $\rho_{S,n}$ de Spearman	39
3.1	Cadre d'application	39
3.2	Le coefficient de corrélation $\rho_S(X, Y)$ de Spearman	40
3.3	Estimation de $\rho_S(X, Y)$	40
3.4	Procédure de test	41
3.5	Cas des ex æquo	43
3.6	Statistique corrigée $\rho_{S,n}^*$	43
3.7	Départition des ex æquo	44
4	Corrélation partielle de Spearman	45
4.1	Coefficient de corrélation partielle de Spearman $\rho_S(X, Y Z)$	45
4.2	Estimation de $\rho_S(X, Y Z)$	45
4.3	Méthode d'utilisation	46
5	La statistique τ_n de Kendall	46
5.1	Cadre d'application	46
5.2	Le coefficient de corrélation $\tau(X, Y)$ de Kendall	47
5.3	Estimation de $\tau(X, Y)$	48
5.4	Procédure de test	49
5.5	Cas des ex æquo	52
5.6	Statistique corrigée τ_n^*	52
5.7	Départition des ex æquo	53
6	Corrélation partielle de Kendall	53
6.1	Coefficient de corrélation partiel de Kendall $\tau(X, Y Z)$	53
6.2	Notations	55
6.3	Estimateur de $\tau(X, Y Z)$	56
6.4	Méthode d'utilisation	56

7	Test du χ^2 d'indépendance	57
7.1	Les classes	57
7.2	Statistique du test	58
7.3	Méthode d'utilisation	58
7.4	Correction de Yates	59
7.5	Étude des résidus	60
8	Test exact de Fisher	60
8.1	Deux variables à deux modalités	60
8.2	Deux variables ayant un nombre fini quelconque de modalités	65
9	Tests de multinormalité	68
9.1	Utilisation de tests multiples unidimensionnels	68
9.2	Une test basée sur la modification du test de Shapiro-Wilk unidimensionnel	70
9.3	Asymétrie et aplatissement multivariés de Mardia	70
9.4	Centrer-réduire.	70
9.5	Asymétrie et aplatissement	72
9.6	Étude sous l'hypothèse de multinormalité de \mathbf{X}	72
9.7	Un test basé sur la fonction caractéristique empirique	73