

Analyse de variance multivariée — MANOVA

Résumé

Le modèle linéaire gaussien standard a pour objectif de modéliser la dépendance (supposée linéaire) d'une variable aléatoire réelle (Y) par rapport soit à d'autres variables quantitatives contrôlées (cas de la régression), soit à des facteurs (cas de l'analyse de variance), soit à un mélange des deux (cas de l'analyse de covariance). Ce modèle s'étend sans difficulté majeure au cas où la réponse n'est pas unidimensionnelle, mais multidimensionnelle : la variable aléatoire réelle Y est alors remplacée par un vecteur aléatoire.

Retour au [plan du cours](#)

1 Introduction

Dans ce cours, nous avons déjà repris et développé les méthodes d'analyse de la variance ; par contre, nous ne sommes revenus ni sur les méthodes de régression ni sur les méthodes d'analyse de la covariance. De la même manière, dans le cadre du modèle linéaire gaussien multivarié, nous détaillerons uniquement la généralisation de l'analyse de variance. Pour la régression et l'analyse de covariance, le passage au cas multidimensionnel se fait de la même façon. En particulier, on trouve les mêmes tests que ceux présentés ici.

Concernant la bibliographie, l'ouvrage de référence pour ce chapitre est celui de Seber (1984). On peut également indiquer l'ouvrage de Anderson (2003) et celui de Rencher (1995).

Le chapitre 5 est donc consacré aux plans factoriels avec réponse multidimensionnelle, autrement dit à l'analyse de variance multivariée (ou multidimensionnelle), encore appelée MANOVA (acronyme pour *Multivariate ANalysis Of VAriance*). Seuls seront traités dans ce chapitre les cas de un facteur et de deux facteurs croisés. Dans ce contexte, nous n'aborderons pas la généralisation des intervalles de confiance et nous nous consacrerons seulement à l'estimation ponctuelle des paramètres et aux tests d'hypothèses.

Pour ce qui est de l'estimation ponctuelle des paramètres, les principes et les résultats sont de même nature que ceux vus dans le cas unidimensionnel. Toutefois, la méthode du maximum de vraisemblance se trouve compliquée par le fait que chaque observation est maintenant la réalisation d'une loi gaussienne multidimensionnelle, ce qui alourdit l'écriture de la vraisemblance et nécessite des dérivations matricielles. Par ailleurs, l'expression des paramètres est maintenant matricielle et non plus vectorielle. Ainsi, si nous notons D la dimension du vecteur aléatoire réponse Y ($D \geq 2$), on retrouve l'expression habituelle des estimateurs des paramètres

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y},$$

dans laquelle \mathbf{Y} est maintenant une matrice $n \times D$, de sorte que $\hat{\beta}$ est une matrice $p \times D$ (n est le nombre total d'observations et p est le nombre de colonnes de $X : J$ dans le cas d'un seul facteur, JK dans le cas de deux facteurs croisés, etc.).

La loi normale centrée unidimensionnelle prise jusqu'à présent comme modèle pour les erreurs avait pour variance σ^2 . Cette variance est ici remplacée par une matrice de variances-covariances Σ , $D \times D$, pour la loi normale centrée, multidimensionnelle d'ordre D , des erreurs. Si on pose $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta}$ (matrice $n \times D$ des valeurs prédites) et $\hat{\mathbf{U}} = \mathbf{Y} - \hat{\mathbf{Y}}$ (matrice $n \times D$ des résidus), la matrice Σ est estimée par $\frac{1}{n-p}\hat{\mathbf{U}}'\hat{\mathbf{U}}$, où $\hat{\mathbf{U}}'\hat{\mathbf{U}}$ est distribuée selon une loi de Wishart, généralisation multidimensionnelle de la loi de khi-deux.

Pour ce qui est des tests, le test de Fisher, permettant de tester différentes hypothèses nulles en ANOVA unidimensionnelle, est maintenant remplacé par plusieurs tests (quatre dans SAS) dont les statistiques sont calculées à partir des valeurs propres des deux matrices remplaçant numérateur et dénominateur de la statistique de Fisher. Les tests fournis par SAS sont les tests de Wilks, de Lawley-Hotelling, de Pillai et de Roy. Dans les cas simples, ils sont tous les quatre équivalents. Dans les autres cas, les trois premiers sont voisins et très rarement contradictoires. Par contre, le quatrième est moins précis et est déconseillé. S'il faut en privilégier un, nous recommandons plus particulièrement le test de Wilks.

C'est encore la procédure GLM de SAS qui est utilisée pour mettre en œuvre la MANOVA.

Dans tout ce chapitre, l'objectif est de modéliser un vecteur aléatoire Y de \mathbb{R}^D ($D \in \mathbb{N}$, $D \geq 2$) au moyen d'une loi gaussienne sur \mathbb{R}^D .

2 Écriture du modèle à un seul facteur

2.1 Les données

- On considère ici un unique facteur, encore noté F , possédant J niveaux ($J \geq 2$), indicés par j ($j = 1, \dots, J$).
- Pour chaque niveau j de F , on réalise n_j observations du vecteur aléatoire Y de \mathbb{R}^D ($n_j \geq 1$); on pose $n = \sum_{j=1}^J n_j$.
- On note Y_{ij} le vecteur aléatoire associé à la i -ième observation réalisée au niveau j de F : $Y_{ij} \in \mathbb{R}^D$.

L'objectif de la MANOVA est d'étudier l'influence des niveaux du facteur F sur les valeurs du vecteur réponse Y . Cette influence va être étudiée globalement, dans \mathbb{R}^D , d'où la nécessité d'avoir recours à des techniques multidimensionnelles, différentes de celles vues en ANOVA.

Remarque. — Parallèlement à la MANOVA, il est habituel de faire une ANOVA pour chacune des D composantes du vecteur Y (le logiciel SAS le fait automatiquement). C'est un complément intéressant pour la MANOVA, mais cela ne la remplace pas. En particulier, les tests à regarder pour le choix d'un modèle adapté à un jeu de données sont les tests multidimensionnels.

2.2 Le modèle

2.2.1 Écriture initiale

Pour chaque expérience (i, j) (i -ième observation réalisée au niveau j de F), on écrit le vecteur aléatoire réponse Y_{ij} de \mathbb{R}^D sous la forme :

$$Y_{ij} = \beta_j + U_{ij}.$$

Attention, les trois éléments de cette écriture doivent être vus comme des **vecteurs-lignes** de \mathbb{R}^D , comme précisé ci-dessous.

- Le vecteur $\beta_j = (\beta_j^1 \cdots \beta_j^d \cdots \beta_j^D)$ est un paramètre à estimer ; il modélise la valeur de la réponse Y au niveau j de F .
- Le terme $U_{ij} = (U_{ij}^1 \cdots U_{ij}^d \cdots U_{ij}^D)$ est le vecteur aléatoire des erreurs. On sup-

pose que les U_{ij} sont i.i.d., de loi $\mathcal{N}_D(0_D, \Sigma)$, où Σ est une matrice symétrique et strictement définie-positives ; on doit également estimer Σ . On notera que Σ ne dépend pas de j , autrement dit on est toujours dans le cadre d'un modèle homoscédastique.

- Les vecteurs aléatoires Y_{ij} sont donc indépendants, de loi $\mathcal{N}_D(\beta_j^t, \Sigma)$.

Finalement, il y a $J \times D$ paramètres de moyenne β_j^d à estimer, ainsi que $\frac{D(D+1)}{2}$ paramètres de variance $(\Sigma)_d^{d'}$ ($1 \leq d \leq D$; $1 \leq d' \leq D$). Comme on dispose de nD observations, on doit veiller à ce que la taille n de l'échantillon utilisé vérifie : $n \geq J + \frac{D+1}{2}$.

2.2.2 Écriture matricielle

L'ensemble des nD observations réalisées peut se mettre sous la forme matricielle suivante :

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{U}.$$

- Dans l'écriture ci-dessus, \mathbf{X} et β sont des matrices réelles (non aléatoires) de dimensions respectives $n \times J$ et $J \times D$.
- Comme dans le cas unidimensionnel, les colonnes de la matrice d'incidence \mathbf{X} sont les indicatrices Z^j des niveaux du facteur F , de sorte que \mathbf{X} ne comporte que des 0 et des 1.
- Les termes \mathbf{Y} et \mathbf{U} sont des matrices aléatoires de dimension $n \times D$. Elles sont gaussiennes et vérifient :

$$\mathbb{E}(\mathbf{U}) = \mathbf{0}_{n \times D} ; \quad \mathbb{E}(\mathbf{Y}) = \mathbf{X}\beta ; \quad \text{Var}(\mathbf{U}) = \text{Var}(\mathbf{Y}) = \mathbf{I}_n \otimes \Sigma.$$

Dans cette dernière écriture, \mathbf{I}_n désigne la matrice identité d'ordre n et \otimes le produit matriciel direct, ou produit de Kronecker. En fait, on a

$$\mathbf{I}_n \otimes \Sigma = \begin{pmatrix} \Sigma & \mathbf{0}_{D \times D} & \cdots & \mathbf{0}_{D \times D} \\ \mathbf{0}_{D \times D} & \Sigma & \cdots & \mathbf{0}_{D \times D} \\ \cdots & \cdots & \cdots & \cdots \\ \mathbf{0}_{D \times D} & \mathbf{0}_{D \times D} & \cdots & \Sigma \end{pmatrix},$$

où chacun des n^2 termes de cette matrice est lui-même une matrice (un bloc matriciel), de dimension $D \times D$. La matrice $\mathbf{I}_n \otimes \Sigma$ est donc carrée d'ordre nD .

2.2.3 Paramétrage centré

Comme dans le cas unidimensionnel, ce paramétrage consiste à décomposer chaque vecteur-ligne β_j sous la forme :

$$\beta_j = \mu + \alpha_j, \text{ avec } \mu = \frac{1}{J} \sum_{j=1}^J \beta_j \text{ et } \alpha_j = \beta_j - \mu.$$

Le paramètre μ est l'effet (moyen) général et le paramètre α_j est l'effet principal (ou différentiel) du niveau j de F . Ces deux paramètres sont des vecteurs de \mathbb{R}^D et on notera que l'on a encore $\sum_{j=1}^J \alpha_j = 0_D$.

2.2.4 Paramétrage SAS

Pour ce paramétrage, on pose $m = \beta_J$ et $a_j = \beta_j - \beta_J$ (de sorte que, encore une fois, $a_J = 0_D$). Les paramètres m et a_j sont également des vecteurs de \mathbb{R}^D .

3 Estimation des paramètres du modèle à un facteur

3.1 Vraisemblance et log-vraisemblance

La vraisemblance de l'échantillon des y_{ij} s'écrit

$$\begin{aligned} L(y_{ij}, \beta, \Sigma) &= \prod_{j=1}^J \prod_{i=1}^{n_j} \frac{1}{(2\pi)^{D/2} (\det \Sigma)^{1/2}} \exp\left[-\frac{1}{2} (y_{ij} - \beta_j) \Sigma^{-1} (y_{ij} - \beta_j)'\right] \\ &= C_1 (\det \Sigma)^{-n/2} \exp\left[-\frac{1}{2} \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - X_j \beta) \Sigma^{-1} (y_{ij} - X_j \beta)'\right], \end{aligned}$$

où C_1 est une constante et X_j un vecteur-ligne à J éléments, comportant un 1 en j -ième colonne et des 0 partout ailleurs (en fait, X_j est n'importe laquelle des lignes de la matrice \mathbf{X} correspondant aux observations du niveau j de F).

La log-vraisemblance s'écrit

$$\begin{aligned} l(y_{ij}, \beta, \Sigma) &= \log[L(y_{ij}, \beta, \Sigma)] \\ &= C_2 - \frac{n}{2} \log(\det \Sigma) - \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|_{\mathbf{I}_n, \Sigma^{-1}}^2, \end{aligned}$$

où \log désigne le logarithme népérien et où $C_2 = \log(C_1)$. Par ailleurs, on rappelle que si \mathbf{N} est une matrice $n \times n$, symétrique et strictement définie-positive, si \mathbf{P} est une matrice $p \times p$, également symétrique et strictement définie-positive, alors, pour toute matrice \mathbf{A} de dimension $n \times p$, on peut définir sa norme carrée par

$$\|\mathbf{A}\|_{\mathbf{N}, \mathbf{P}}^2 = \text{tr}(\mathbf{A}\mathbf{P}\mathbf{A}'\mathbf{N}),$$

où tr désigne la trace de la matrice correspondante (cette norme est appelée norme de Hilbert-Schmidt).

3.2 Estimation maximum de vraisemblance

Pour estimer les matrices β et Σ , on doit ici faire des dérivations matricielles. On admettra les résultats ci-dessous (qui sont également les résultats de l'estimation moindres carrés, en l'absence de l'hypothèse de normalité). Pour des précisions, on pourra se reporter à Seber (1984). Citons également le site "The Matrix Cookbook", très intéressant, à l'adresse suivante :

<http://matrixcookbook.com/>

3.2.1 Estimation des paramètres d'intérêt

On obtient $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\bar{y}_{\bullet 1} \dots \bar{y}_{\bullet j} \dots \bar{y}_{\bullet J})'$ (matrice $J \times D$), où $\bar{y}_{\bullet j} = \frac{1}{n_j} \sum_{i=1}^{n_j} y_{ij}$ est le vecteur-ligne de \mathbb{R}^D , moyenne des observations de Y au niveau j de F . On notera $\hat{\mathbf{B}}$ l'estimateur correspondant défini par $\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$.

Remarque. — Signalons que $\hat{\beta}$ peut s'obtenir colonne par colonne, au moyen des résultats, pour chaque colonne, d'une ANOVA unidimensionnelle à un facteur : $\hat{\beta}^d = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}^d$ est la solution de l'ANOVA de Y^d sur F . Ainsi, pour obtenir $\hat{\beta}$, on pourra utiliser les estimations unidimensionnelles de $\hat{\beta}^d$ ($d = 1, \dots, D$) fournies par SAS au début des sorties de la MANOVA.

3.2.2 Valeurs prédites

Pour un niveau j donné, et pour toute observation i faite à ce niveau, la valeur prédite correspondante est : $\hat{y}_{ij} = \hat{\beta}_j = \bar{y}_{\bullet j}$ (vecteur de \mathbb{R}^D). On notera $\hat{\mathbf{Y}}$ la matrice aléatoire $n \times D$ de l'ensemble des valeurs prédites.

3.2.3 Résidus

Il vient : $\hat{u}_{ij} = y_{ij} - \hat{y}_{ij} = y_{ij} - \bar{y}_{\bullet j}$ (vecteur de \mathbb{R}^D). On notera $\hat{\mathbf{U}}$ la matrice aléatoire $n \times D$ des résidus ainsi définis.

3.2.4 Estimation de la matrice des variances-covariances

Elle est obtenue, comme dans le cas unidimensionnel, à partir de la matrice des résidus : $\hat{\Sigma} = \frac{1}{n - J} \hat{\mathbf{U}}' \hat{\mathbf{U}}$ (matrice $D \times D$).

3.3 Propriétés des estimateurs maximum de vraisemblance

- Les matrices $\hat{\mathbf{B}}$, de dimension $J \times D$, $\hat{\mathbf{Y}}$, de dimension $n \times D$, et $\hat{\mathbf{U}}$, de dimension $n \times D$, sont des matrices aléatoires gaussiennes, d'espérances respectives β , $\mathbf{X}\beta$ et $\mathbf{0}_{n \times D}$ (leurs matrices de variances-covariances ne seront pas explicitées).
- La matrice aléatoire $\hat{\mathbf{U}}$ est indépendante des matrices aléatoires $\hat{\mathbf{B}}$ et $\hat{\mathbf{Y}}$.
- Enfin, $(n - J)\hat{\Sigma} = \hat{\mathbf{U}}' \hat{\mathbf{U}}$ est une matrice aléatoire distribuée selon une loi de Wishart de dimension D , à $n - J$ degrés de liberté et de matrice associée Σ ; cette loi de probabilité est notée $W_D(n - J, \Sigma)$.

3.4 Indications sur la loi de Wishart

Il s'agit, en quelques sortes, d'une généralisation multidimensionnelle de la loi de khi-deux. Dans \mathbb{R}^D ($D \geq 2$), considérons m ($m \geq D$) vecteurs (colonnes) aléatoires notés T_i ($i = 1, \dots, m$), supposés i.i.d. selon une loi normale centrée, de matrice de variances-covariances Σ ($D \times D$, symétrique et strictement définie-positif). Alors, la matrice aléatoire $\mathbf{W} = \sum_{i=1}^m T_i T_i'$ (de dimension $D \times D$) définit une loi de Wishart de dimension D , à m d.d.l., de matrice associée Σ . Elle est notée $W_D(m, \Sigma)$.

Remarque. — On notera que Σ n'est pas la matrice des variances-covariances de \mathbf{W} . En effet, la matrice aléatoire \mathbf{W} est constituée de $D \times D = D^2$ éléments aléatoires et admet donc une matrice de variances-covariances de dimension $D^2 \times D^2$ qui ne sera pas explicitée (mais qui ne peut être Σ).

3.4.1 Quelques propriétés de la loi de Wishart

- Telle qu'elle est définie ci-dessus, la loi de Wishart $W_D(m, \Sigma)$ apparaît comme la loi de m fois la matrice des variances-covariances empiriques d'une loi normale centrée de \mathbb{R}^D , de matrice de variances-covariances Σ .
- $\mathbb{E}(\mathbf{W}) = m\Sigma$ (immédiat).
- Supposons : $\mathbf{W}_1 \sim W_D(m_1, \Sigma)$; $\mathbf{W}_2 \sim W_D(m_2, \Sigma)$; \mathbf{W}_1 et \mathbf{W}_2 indépendantes; alors :
 $\mathbf{W}_1 + \mathbf{W}_2 \sim W_D(m_1 + m_2, \Sigma)$. (Cette propriété est admise.)

Pour plus de détails sur la loi de Wishart, on pourra encore se reporter à Seber (1984).

4 Tests dans le modèle à un facteur

La seule hypothèse considérée ici est la significativité du facteur F , autrement dit la significativité du modèle lui-même (puisque F est le seul facteur pris en compte, pour l'instant, dans le modèle). L'hypothèse nulle s'écrit sous l'une des formes suivantes :

$$\{H_0 : F \text{ n'a pas d'effet sur } Y\} \iff \{H_0 : \beta_1 = \dots = \beta_J\} \iff \{H_0 : \alpha_1 = \dots = \alpha_J = 0\} \iff \{H_0 : a_1 = \dots = a_J = 0\}$$
 (il y a chaque fois $J - 1$ contraintes vectorielles dans \mathbb{R}^D).

La mise en œuvre d'un test permettant de tester H_0 contre son contraire H_1 , avec un niveau α fixé, nécessite de généraliser le test de Fisher qui ne peut plus être utilisé ici. Cette généralisation peut se faire de différentes manières et conduit à différents tests. Avant de les introduire, nous devons définir les deux matrices à partir desquelles ils sont construits.

4.1 Les matrices H et E

4.1.1 Retour sur le cas où Y est unidimensionnelle

Revenons sur le cas de l'ANOVA à un seul facteur F avec une variable réponse Y unidimensionnelle. Pour tester la significativité du modèle, donc du facteur, on a utilisé le test de Fisher, qui fait intervenir les quantités ci-dessous.

- $SSF = \sum_{j=1}^J n_j (\bar{y}_{\bullet j} - \bar{y}_{\bullet\bullet})^2$: c'est la somme des carrés expliquée par le facteur F , ou somme des carrés inter-groupes, ou *between sum of squares*. Sous l'hypothèse nulle H_0 , cette somme est nécessairement

“petite” ; nous la noterons H , car elle est liée à l’hypothèse testée : $H = SSF$. Son d.d.l. est $J - 1$.

- $SSE = \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\bullet j})^2$: c’est la somme des carrés résiduelle (non expliquée par F , donc par le modèle), ou somme des carrés intra-groupes, ou *pooled within sum of squares*. Elle représente la somme des carrés liée à l’erreur du modèle, pour cette raison notée E : $E = SSE$. Son d.d.l. est $n - J$.

La statistique du test de Fisher peut s’écrire sous la forme :

$$F = \frac{SSF}{SSE} \frac{n - J}{J - 1} = \frac{H}{E} \frac{n - J}{J - 1}.$$

(Ne pas confondre les deux notations F , tantôt pour le facteur, tantôt pour la statistique de Fisher.)

Enfin, rappelons la relation $SST = \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\bullet\bullet})^2 = SSF + SSE = H + E$, où SST est la somme des carrés totale, de d.d.l. $n - 1$.

4.1.2 Généralisation au cas où Y est multidimensionnelle

Dans le cas de la MANOVA en dimension D et à un seul facteur, on généralise les quantités H et E comme indiqué ci-dessous.

- La somme H est remplacée par la matrice \mathbf{H} , $D \times D$, définie par :

$$\mathbf{H} = \sum_{j=1}^J n_j (\bar{y}_{\bullet j} - \bar{y}_{\bullet\bullet})' (\bar{y}_{\bullet j} - \bar{y}_{\bullet\bullet}).$$

Le d.d.l. associé, qui vaut toujours $J - 1$, sera noté ν_H .

- La somme E est remplacée par la matrice \mathbf{E} , $D \times D$, définie par :

$$\mathbf{E} = \sum_{j=1}^J \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\bullet j})' (y_{ij} - \bar{y}_{\bullet j}).$$

Le d.d.l. associé, qui vaut toujours $n - J$, sera noté ν_E . On notera que $\mathbf{E} = \hat{\mathbf{U}}' \hat{\mathbf{U}}$.

- La somme des carrés totale, SST , est remplacée par la somme des deux matrices définies ci-dessus : $\mathbf{H} + \mathbf{E}$; son d.d.l. est $n - 1$. Nous n’utiliserons pas de notation particulière pour cette matrice.

En ANOVA à un seul facteur, la statistique de Fisher est proportionnelle au rapport $\frac{H}{E}$. En MANOVA à un seul facteur, pour tester la significativité de ce facteur, les tests utilisés s’appuient sur l’un des produits matriciels $\mathbf{H}\mathbf{E}^{-1}$ ou $\mathbf{H}(\mathbf{H} + \mathbf{E})^{-1}$.

Remarque. — Nous avons déjà donné l’expression de l’estimateur de la matrice des variances-covariances : $\hat{\Sigma} = \frac{1}{n - J} \hat{\mathbf{U}}' \hat{\mathbf{U}}$. On voit qu’on peut la réécrire sous la forme : $\hat{\Sigma} = \frac{1}{n - J} \mathbf{E} = \frac{\mathbf{E}}{\nu_E}$.

Remarque. — La matrice $\hat{\Sigma}$ contient des variances et des covariances empiriques résiduelles, c’est-à-dire conditionnées par le facteur F . Si, à partir des éléments de $\hat{\Sigma}$, on calcule des coefficients de corrélations linéaires empiriques entre composantes de Y , il s’agira de corrélations résiduelles, plus couramment appelées *corrélations partielles* (conditionnelles à F). On trouve ces corrélations partielles en sortie du logiciel SAS (en pratique, elles ont peu d’intérêt).

4.2 Le test de Wilks

4.2.1 Principe

Il s’agit du test le plus courant dans le contexte de la MANOVA qui est, en fait, une adaptation du test du rapport des vraisemblances.

Notons θ le vecteur de tous les paramètres du modèle, de dimension $(J \times D) + \frac{D(D + 1)}{2}$, $\hat{\theta}$ son estimation maximum de vraisemblance et $\hat{\theta}_0$ son estimation maximum de vraisemblance sous la contrainte définie par H_0 (autrement dit sous la contrainte linéaire : $\beta_1 = \dots = \beta_J$). La statistique du test du rapport des vraisemblances est $\frac{L(y, \hat{\theta}_0)}{L(y, \hat{\theta})}$, dont on peut vérifier qu’elle

vaut : $[\frac{\det(\mathbf{E})}{\det(\mathbf{H} + \mathbf{E})}]^{n/2}$ (voir Seber, 1984). Le test de Wilks consiste à considérer la puissance $2/n$ de cette quantité, autrement dit sa statistique est définie par

$$\Lambda = \frac{\det(\mathbf{E})}{\det(\mathbf{H} + \mathbf{E})} = \prod_{k=1}^s \frac{1}{1 + \lambda_k},$$

où les λ_k sont les valeurs propres de la matrice \mathbf{HE}^{-1} (ou $\mathbf{E}^{-1}\mathbf{H}$) et où $s = \inf(D, J - 1)$ est le nombre de valeurs propres non nulles de cette matrice.

La mise en œuvre de ce test va dépendre du cas de figure.

4.2.2 Cas où on se ramène à un test de Fisher exact

Cela se produit dans trois cas particuliers.

- Cas d'un facteur à 2 niveaux : $J = 2 \iff \nu_H = J - 1 = 1$ (D quelconque). On peut alors montrer :

$$\frac{1 - \Lambda}{\Lambda} \frac{\nu_E - D + 1}{D} = \frac{1 - \Lambda}{\Lambda} \frac{n - (D + 1)}{D} \sim F_{D; n - (D + 1)}.$$

Les tables de la distribution de Fisher permettent donc de faire un test exact.

- Cas d'un facteur à 3 niveaux : $J = 3 \iff \nu_H = 2$ (D quelconque). Il vient dans ce cas :

$$\frac{1 - \sqrt{\Lambda}}{\sqrt{\Lambda}} \frac{\nu_E - D + 1}{D} = \frac{1 - \sqrt{\Lambda}}{\sqrt{\Lambda}} \frac{n - (D + 2)}{D} \sim F_{2D; 2(n - (D + 2))}.$$

Même chose, on peut faire un test de Fisher exact.

- Cas où Y est à 2 dimensions : $D = 2$ (J quelconque). Il vient maintenant :

$$\frac{1 - \sqrt{\Lambda}}{\sqrt{\Lambda}} \frac{\nu_E - 1}{\nu_H} = \frac{1 - \sqrt{\Lambda}}{\sqrt{\Lambda}} \frac{n - (J + 1)}{J - 1} \sim F_{2\nu_H; 2(\nu_E - 1)} (F_{2(J - 1); 2(n - (J + 1))}).$$

Toujours la même chose.

4.2.3 Cas où on dispose de tables

Des tables du test de Wilks on été établies et permettent de faire encore un test exact dans de nombreux autres cas (on les trouve dans les ouvrages de statistique multidimensionnelle). Pour les niveaux 10%, 5%, 2, 5%, 1% et 0, 5%, on dispose de tables pour D variant de 3 à 10, pour ν_H variant de 3 à 13 (et souvent plus) et pour ν_E variant de 1 à 20, ainsi que pour les valeurs 30, 40, 60 et 120. On trouvera ces tables, par exemple, dans Seber (1984).

4.2.4 Approximation de Fisher

Dans les cas où on ne dispose pas de tables permettant de faire un test exact, on pourra faire un test de Fisher approché (d'autant meilleur que n est grand) en utilisant le résultat suivant :

$$\phi = \frac{1 - \Lambda^{1/t}}{\Lambda^{1/t}} \frac{ft - g}{D\nu_H} \sim F_{D\nu_H; ft - g} \text{ (approximativement).}$$

Dans l'expression ci-dessus, on a les définitions suivantes :

$$\begin{aligned} f &= \nu_H + \nu_E - \frac{\nu_H + D + 1}{2} = (n - 1) - \frac{J + D}{2}; \\ g &= \frac{D\nu_H}{2} - 1; \\ t &= \left[\frac{D^2\nu_H^2 - 4}{D^2 + \nu_H^2 - 5} \right]^{1/2}. \end{aligned}$$

Remarque. — Dans chacun des trois cas particuliers $J = 2$, $J = 3$ et $D = 2$, on pourra vérifier que l'expression ci-dessus redonne celle fournie plus haut. Dans ces trois cas, la distribution de Fisher n'est donc pas une approximation de la loi de la statistique de test, mais sa distribution exacte.

Remarque. — Concernant le test de Wilks (ainsi, d'ailleurs, que les suivants), le logiciel SAS fournit, en sortie de la procédure GLM, la valeur de Λ , la valeur $\phi = \frac{1 - \Lambda^{1/t}}{\Lambda^{1/t}} \frac{ft - g}{D\nu_H}$, les d.d.l. $D\nu_H$ et $ft - g$, et une p -value représentant la probabilité qu'une loi de Fisher à $D\nu_H$ et $ft - g$ d.d.l. dépasse ϕ ; le test réalisé à partir de cette p -value est donc, selon le cas, soit un test exact, soit l'approximation de Fisher indiquée ci-dessus.

4.3 Autres tests

Dans la littérature statistique, on trouve d'autres tests permettant d'éprouver la même hypothèse nulle. Ces tests sont automatiquement fournis par le logiciel SAS, en même temps que le test de Wilks. Nous donnons leur principe ci-dessous.

4.3.1 Le test de la trace de Lawley-Hotelling

La statistique de ce test est :

$$T^2 = \nu_E \text{trace}(\mathbf{HE}^{-1}) = (n - J) \sum_{k=1}^s \lambda_k.$$

Pour un niveau de test $\alpha = 5\%$, pour des valeurs de D variant de 2 à 6, pour $\nu_H = J - 1$ variant de D à 6, puis prenant les valeurs 8, 10, 12, 15, 20, 25, 40 et 60, enfin pour $\nu_E = n - J$ variant de D à 8, puis prenant les valeurs 10, 20, 30 \dots 100 et 200, on dispose de tables pour la statistique $\frac{T^2}{\nu_E} = \sum_{k=1}^s \lambda_k$, permettant de faire un test exact.

Dans les autres cas, on utilise l'approximation suivante

$$\frac{T^2}{c \nu_E} = \frac{1}{c} \text{trace}(\mathbf{HE}^{-1}) \sim F_{a; b} \text{ (approximativement)}$$

avec :

$$\begin{aligned} a &= D \nu_H ; b = 4 + \frac{a + 2}{B - 1} ; \\ B &= \frac{(\nu_E + \nu_H - D - 1)(\nu_E - 1)}{(\nu_E - D - 3)(\nu_E - D)} ; \\ c &= \frac{a(b - 2)}{b(\nu_E - D - 1)}. \end{aligned}$$

Remarque. — Compte-tenu de l'expression de la statistique de ce test, celui-ci est la généralisation multidimensionnelle la plus naturelle du test de Fisher.

4.3.2 Le test de la trace de Pillai

La statistique de ce test est

$$V = \text{trace} [\mathbf{H}(\mathbf{H} + \mathbf{E})^{-1}] = \sum_{k=1}^s \mu_k = \sum_{k=1}^s \frac{\lambda_k}{1 + \lambda_k},$$

où $s = \inf(D, J - 1)$, où les μ_k sont les valeurs propres de la matrice $\mathbf{H}(\mathbf{H} + \mathbf{E})^{-1}$ et les λ_k celles de la matrice \mathbf{HE}^{-1} .

Si l'on pose $k_1 = \frac{1}{2}(|D - \nu_H| - 1)$ et $k_2 = \frac{1}{2}(\nu_E - D - 1)$, des tables permettent de réaliser un test exact de Pillai pour $\alpha = 5\%$, s variant de 2 à 6, k_1 et k_2 variant de 0 à 10 ou bien prenant les valeurs 15, 20 ou 25.

Dans les autres cas, on utilisera l'approximation suivante :

$$\begin{aligned} \frac{V}{s - V} \frac{2k_2 + s + 1}{2k_1 + s + 1} &= \frac{V}{s - V} \frac{n + \inf(D, J - 1) - (D + J)}{\sup(D, J - 1)} \\ &\sim F_{s(2k_1 + s + 1); s(2k_2 + s + 1)} \end{aligned}$$

(approximativement).

4.3.3 Le test de la plus grande racine de Roy

La statistique de ce dernier test est λ_{\max} , la plus grande des valeurs propres de \mathbf{HE}^{-1} . On trouve diverses approximations qui permettent de mettre en œuvre ce test. Celle utilisée par SAS est la suivante :

$$S = \frac{\lambda_{\max}(\nu_H + \nu_E - r)}{r} \sim F_{r; \nu_H + \nu_E - r} \text{ (approximativement)},$$

où $r = \max(D, \nu_H)$.

On notera que, dans ce cas, la loi de Fisher est un minorant pour la loi de S , ce qui signifie que la p -value calculée par SAS pour ce test est toujours plus petite que celle des autres tests. Pour cette raison, **nous déconseillons ce test.**

4.4 Cas particulier : $J = 2$

Il n'y a qu'une seule valeur propre non nulle dans ce cas particulier, et les différentes approximations par une loi de Fisher données ci-dessus sont toutes identiques (le résultat est simple, bien qu'un peu fastidieux, à vérifier). De plus, elles correspondent toutes les quatre au test exact de Fisher donné en 5.3.2. En fait, dans ce cas, la statistique utilisée vaut

$$\lambda_1 \frac{n - D - 1}{D}$$

et est distribuée selon une loi de Fisher à D et $n - D - 1$ degrés de liberté.

5 Illustration

Les données

Il s'agit d'un exemple fictif d'analyse de variance multidimensionnelle, de dimension 3, à un seul facteur. Le facteur est à deux niveaux, notés 1 et 2, et figure en première colonne. Les variables réponses figurent dans les trois colonnes suivantes et prennent des valeurs entières comprises entre 8 et 24. Il y a 8 individus observés, donc 8 lignes dans le fichier des données reproduit ci-dessous.

```
1 10 12 14
1 11 13 15
1 8 9 8
1 9 10 8
2 15 17 16
2 19 18 17
2 21 20 19
2 23 22 24
```

5.0.1 Le programme SAS

Le programme SAS ci-dessous permet de faire la MANOVA de ces données de façon standard.

```
* ----- ;
* options facultatives pour la mise en page des sorties ;
* ----- ;
options pagesize=64 linesize=76 nodate;
title;
footnote 'MANOVA - donnees fictives - 1 facteur';
* ----- ;
*          lecture des donnees          ;
*          (le fichier "fic1.don" contient les donnees ;
*          et se trouve dans le repertoire de travail) ;
* ----- ;
data fic1;
infile 'fic1.don';
input f $ y1 y2 y3;
run;
* ----- ;
```

```
*          procedure GLM pour la MANOVA          ;
* ----- ;
proc glm data=fic1;
class f;
model y1-y3 = f / ss3 solution;
manova H = f;
run;
quit;
```

5.0.2 Les sorties de la procédure GLM

```
PAGE 1                                The GLM Procedure
-----

Class Level Information

Class          Levels      Values
f                2        1 2

Number of observations      8

PAGE 2                                The GLM Procedure
-----
Dependent Variable: y1

Source          DF          Sum of Squares      Mean Square      F Value      Pr > F
Model          1          200.0000000      200.0000000      30.00      0.0015
Error          6          40.0000000           6.6666667
Corrected Total 7          240.0000000

R-Square      Coeff Var      Root MSE      y1 Mean
0.8333333      17.80682      2.581989      14.50000

Source          DF      Type III SS      Mean Square      F Value      Pr > F
f                1      200.0000000      200.0000000      30.00      0.0015

Parameter          Estimate          Standard Error      t Value      Pr > |t|
Intercept          19.50000000 B      1.29099445      15.10      <.0001
f 1                -10.00000000 B      1.82574186      -5.48      0.0015
f 2                 0.00000000 B      .              .              .
```

PAGE 3

 The GLM Procedure

 Dependent Variable: y2

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Parameter	Estimate	Standard Error	t Value	Pr > t
Model	1	136.1250000	136.1250000	33.00	0.0012	Intercept	19.00000000 B	1.83428006	10.36	<.0001
Error	6	24.7500000	4.1250000			f 1	-7.75000000 B	2.59406374	-2.99	0.0244
						f 2	0.00000000 B	.	.	.
Corrected Total	7	160.8750000								

PAGE 5

R-Square Coeff Var Root MSE y2 Mean
 0.846154 13.42816 2.031010 15.12500

Characteristic Roots and Vectors of: E Inverse * H, where
 H = Type III SSCP Matrix for f
 E = Error SSCP Matrix

Source	DF	Type III SS	Mean Square	F Value	Pr > F	Characteristic Root	Percent	Characteristic Vector y1	V'EV=1 y2	y3
f	1	136.1250000	136.1250000	33.00	0.0012	26.4943529	100.00	-0.23580920	1.15536589	-0.45600123
						0.0000000	0.00	-0.36273387	0.42959731	0.01073044
						0.0000000	0.00	0.18534079	-0.44542771	0.23501557

Parameter	Estimate	Standard Error	t Value	Pr > t
Intercept	19.25000000 B	1.01550480	18.96	<.0001
f 1	-8.25000000 B	1.43614066	-5.74	0.0012
f 2	0.00000000 B	.	.	.

MANOVA Test Criteria and Exact F Statistics
 for the Hypothesis of No Overall f Effect
 H = Type III SSCP Matrix for f
 E = Error SSCP Matrix

S=1 M=0.5 N=1

PAGE 4

 The GLM Procedure

 Dependent Variable: y3

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F	Statistic	Value	F Value	Num DF	Den DF	Pr > F
Model	1	120.1250000	120.1250000	8.93	0.0244	Wilks' Lambda	0.03637111	35.33	3	4	0.0025
Error	6	80.7500000	13.4583333			Pillai's Trace	0.96362889	35.33	3	4	0.0025
Corrected Total	7	200.8750000				Hotelling-Lawley Trace	26.49435290	35.33	3	4	0.0025
						Roy's Greatest Root	26.49435290	35.33	3	4	0.0025

6 Modèle à deux facteurs croisés

6.1 Données, modèle et paramétrages

On considère maintenant deux facteurs explicatifs notés F_1 et F_2 , à J et K niveaux respectivement. Les niveaux de F_1 sont indicés par j ($j = 1, \dots, J$) et ceux de F_2 par k ($k = 1, \dots, K$). Pour chaque couple (j, k) obtenu par croisement des deux facteurs, on réalise n_{jk} observations ($n_{jk} \geq 1$) d'un vecteur aléatoire réponse à valeurs dans \mathbb{R}^D ($D \geq 2$), ces vecteurs étant notés Y_{ijk}

($i = 1, \dots, n_{jk}$). On pose $n = \sum_{j=1}^J \sum_{k=1}^K n_{jk}$.

Le modèle se met sous la forme

$$Y_{ijk} = \beta_{jk} + U_{ijk},$$

où chaque β_{jk} est un paramètre de \mathbb{R}^D et où $U_{ijk} \sim \mathcal{N}_D(0, \Sigma)$, les U_{ijk} étant indépendants (et donc i.i.d.). On a ainsi : $Y_{ijk} \sim \mathcal{N}_D(\beta_{jk}, \Sigma)$.

On peut réécrire le modèle sous la forme matricielle

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{U},$$

où \mathbf{Y} est une matrice aléatoire $n \times D$, \mathbf{X} est la matrice d'incidence, de dimension $n \times JK$, dont les colonnes sont les indicatrices des cellules (j, k), β est la matrice $JK \times D$ des paramètres et \mathbf{U} la matrice aléatoire $n \times D$ des erreurs.

Encore une fois, les paramètres β_{jk} peuvent être décomposés selon le paramétrage centré (en écrivant $\beta_{jk} = \mu + \alpha_j^1 + \alpha_k^2 + \gamma_{jk}$, avec les contraintes usuelles de centrage) ou encore selon le paramétrage SAS (en écrivant maintenant $\beta_{jk} = m + a_j^1 + a_k^2 + c_{jk}$, avec tous les paramètres dont au moins un des indices j ou k est maximum égaux à 0), tous les paramètres intervenant étant des vecteurs-lignes de \mathbb{R}^D .

6.2 Tests et estimations

Tout ce qui a été exposé dans le cas d'un seul facteur se généralise ici sans autre difficulté que celle due aux écritures. En particulier, on retrouve les mêmes tests.

- Pour les tests de nullité des effets d'interactions, on notera que l'on a maintenant $\nu_H = (J - 1)(K - 1)$, $\nu_E = n - JK$ et que le nombre de valeurs propres non nulles à prendre en compte est $s = \inf(D, \nu_H)$.
- Par ailleurs, on notera que lorsque le test de significativité de chaque facteur F_1 et F_2 est fait dans le cadre du modèle complet, ν_E demeure égal à $n - JK$ et les simplifications indiquées dans la remarque 46 ne s'appliquent plus.
- Dans un modèle additif à deux facteurs croisés, ν_E vaut $n - (J + K - 1)$.
- Lorsqu'un facteur ne possède que $J = 2$ niveaux, les 4 tests multidimensionnels sont encore tous identiques, mais l'expression de la statistique de test est plus compliquée que celle indiquée en 5.3.4, car le nombre de niveaux du second facteur intervient.

- Enfin, concernant les estimations de $\hat{\beta}$, on les obtient colonne par colonne, comme indiqué dans la remarque 42.

6.3 Généralisation

On peut encore envisager, avec une variable réponse Y multidimensionnelle d'ordre D , des modèles à trois facteurs croisés ou plus. Il n'y a aucune difficulté théorique, mais seulement des difficultés formelles (complexité des écritures) et pratiques : nombre très important de paramètres à estimer, donc d'observations à réaliser. Nous ne détaillons pas davantage ces extensions.

6.4 Illustration

6.4.1 Les données

Les données, toujours fictives, sont de même nature que les précédentes. La variable réponse est à 3 dimensions et figure dans les 3 dernières colonnes du fichier. Il y a maintenant 2 facteurs, le premier à 2 niveaux (notés 1 et 2), le second à 4 niveaux (notés 1, 2, 3 et 4). Les facteurs figurent dans les 2 premières colonnes du fichier. Pour chaque cellule (il y en a 8), on a réalisé 4 observations, de sorte que l'on dispose de 32 observations. Les données sont reproduites ci-dessous.

1	1	8	7	10
1	1	9	13	11
1	1	8	9	8
1	1	9	10	8
1	2	10	12	14
1	2	11	13	15
1	2	11	10	12
1	2	13	12	12
1	3	13	16	19
1	3	15	17	20
1	3	12	16	16
1	3	14	18	18
1	4	12	18	19
1	4	18	19	23
1	4	13	11	19
1	4	16	21	20
2	1	15	17	16

```

2 1 17 18 17
2 1 15 16 18
2 1 18 17 18
2 2 21 20 24
2 2 23 22 24
2 2 20 23 22
2 2 22 26 23
2 3 25 25 30
2 3 28 25 29
2 3 23 28 25
2 3 25 30 27
2 4 28 27 32
2 4 29 26 29
2 4 26 29 27
2 4 27 34 29
    
```

6.4.2 Le programme SAS

L'option `nouni` du programme ci-dessous permet d'éviter les traitements unidimensionnels.

```

options pagesize=64 linesize=76 nodate;
title;
footnote 'MANOVA - donnees fictives - 2 facteurs';
* ----- ;
*           lecture des donnees           ;
*   (le fichier "fic2.don" contient les donnees ;
*   et se trouve dans le repertoire de travail) ;
* ----- ;
data fic2;
infile 'fic2.don';
input f1 $ f2 $ y1 y2 y3;
run;
* ----- ;
*           procedure GLM pour la MANOVA           ;
* ----- ;
proc glm data=fic2;
class f1 f2;
model y1-y3 = f1 | f2 / nouni;
manova H = f1 | f2 / printh printe;
    
```

```

run;
quit;
    
```

6.4.3 Les sorties de la procédure GLM

```

PAGE 1
-----
                                The GLM Procedure

                                Class Level Information

                                Class          Levels   Values

                                f1              2       1 2

                                f2              4       1 2 3 4

                                Number of observations      32

PAGE 2
-----
                                The GLM Procedure
                                Multivariate Analysis of Variance
    
```

```

                                E = Error SSCP Matrix

                                y1              y2              y3

                                y1              63              13              35
                                y2              13              159.75          -2.25
                                y3              35              -2.25          66
    
```

```

                                Partial Correlation Coefficients from the Error SSCP Matrix / Prob > |r|

                                DF = 24              y1              y2              y3

                                y1              1.000000          0.129584          0.542782
                                                0.5370              0.0051

                                y2              0.129584          1.000000          -0.021912
                                                0.5370              0.9172

                                y3              0.542782          -0.021912          1.000000
                                                0.0051              0.9172
    
```

```

PAGE 3
-----
                                The GLM Procedure
                                Multivariate Analysis of Variance

                                tests de f1
                                -----
    
```

```

                                H = Type III SSCP Matrix for f1

                                y1              y2              y3

                                y1              903.125          855.3125          775.625
                                y2              855.3125          810.03125          734.5625
    
```

y3 775.625 734.5625 666.125

H = Type III SSCP Matrix for f2
E = Error SSCP Matrix

Characteristic Roots and Vectors of: E Inverse * H, where
H = Type III SSCP Matrix for f1
E = Error SSCP Matrix

S=3 M=-0.5 N=10

Characteristic Root	Percent	Characteristic Vector		V'EV=1
		y1	y2	
19.8676273	100.00	0.07148178	0.03487179	0.05101439
0.0000000	0.00	-0.11562870	-0.00337882	0.13836212
0.0000000	0.00	-0.06841360	0.07223796	0.00000000

Statistic	Value	F Value	Num DF	Den DF	Pr > F
Wilks' Lambda	0.06751086	12.09	9	53.693	<.0001
Pillai's Trace	0.97150442	3.83	9	72	0.0005
Hotelling-Lawley Trace	13.23494405	31.40	9	31.536	<.0001
Roy's Greatest Root	13.19120030	105.53	3	24	<.0001

NOTE: F Statistic for Roy's Greatest Root is an upper bound.

MANOVA Test Criteria and Exact F Statistics
for the Hypothesis of No Overall f1 Effect

tests des interactions

H = Type III SSCP Matrix for f1
E = Error SSCP Matrix

H = Type III SSCP Matrix for f1*f2

Statistic	Value	F Value	Num DF	Den DF	Pr > F
Wilks' Lambda	0.04792112	145.70	3	22	<.0001
Pillai's Trace	0.95207888	145.70	3	22	<.0001
Hotelling-Lawley Trace	19.86762728	145.70	3	22	<.0001
Roy's Greatest Root	19.86762728	145.70	3	22	<.0001

	y1	y2	y3
y1	28.375	23.0625	6.125
y2	23.0625	23.34375	7.6875
y3	6.125	7.6875	4.375

tests de f2

Characteristic Roots and Vectors of: E Inverse * H, where
H = Type III SSCP Matrix for f1*f2
E = Error SSCP Matrix

H = Type III SSCP Matrix for f2

Characteristic Root	Percent	Characteristic Vector		V'EV=1
		y1	y2	
0.57694177	85.26	0.13350497	0.02373892	-0.04939697
0.09184462	13.57	-0.06443491	0.05012982	0.12343212
0.00788528	1.17	0.03441843	-0.05804521	0.06380435

	y1	y2	y3
y1	352.375	408.5625	474.125
y2	408.5625	479.59375	553.4375
y3	474.125	553.4375	640.375

MANOVA Test Criteria and F Approximations for
the Hypothesis of No Overall f1*f2 Effect

H = Type III SSCP Matrix for f2
E = Error SSCP Matrix

H = Type III SSCP Matrix for f1*f2
E = Error SSCP Matrix

Characteristic Root	Percent	Characteristic Vector		V'EV=1
		y1	y2	
13.1912003	99.67	0.02098208	0.03734450	0.09571189
0.0429976	0.32	-0.14092901	0.04193749	0.06806198
0.0007461	0.01	-0.05346802	-0.05737958	0.08918153

Statistic	Value	F Value	Num DF	Den DF	Pr > F
Wilks' Lambda	0.57625194	1.52	9	53.693	0.1660
Pillai's Trace	0.45780353	1.44	9	72	0.1872
Hotelling-Lawley Trace	0.67667167	1.61	9	31.536	0.1564
Roy's Greatest Root	0.57694177	4.62	3	24	0.0110

MANOVA Test Criteria and F Approximations
for the Hypothesis of No Overall f2 Effect

NOTE: F Statistic for Roy's Greatest Root is an upper bound.

6.4.4 Application

À titre d'application, on pourra, à partir des valeurs propres données ci-dessus dans le cadre du test de significativité des interactions, retrouver les valeurs de Λ (Wilks' Lambda), de F (F Value pour le test de Wilks), ainsi que les degrés de liberté (on remarquera les d.d.l. décimaux).