

Les alternatives existantes

Welch (1951) a proposé une autre approche, que nous ne présenterons pas ici par manque de temps. Le lecteur intéressé par ce sujet pourra en première lecture ouvrir le livre de Howell (sixième édition) à la page 327, puis aller lire l'article original de Welch. Il est à noter que la procédure de Welch se trouve dans la plupart des logiciels statistiques.

Enfin, à titre d'information, Wilcox (1987), dans son ouvrage « **New statistical procedures for the social sciences** » a un avis tranché sur les conséquences de l'hétérogénéité des variances. Il conseille d'utiliser la procédure de Welch, et en particulier lorsque les échantillons sont de tailles inégales.

Sommaire

- 1 Violation des conditions d'application
- 2 Transformation des variables
- 3 Grandeur de l'effet expérimental
- 4 Puissance
 - Puissance a posteriori
 - Détermination du nombre de répétitions
- 5 Facteur à effets aléatoires

Une dernière remarque

Lorsque l'une des deux conditions (la condition de normalité des variables erreurs ou la condition d'homogénéité des variables erreurs) n'est pas vérifiée au moyen d'un test statistique, il faut s'assurer que cela n'est pas dû à une valeur extrême ou aberrante. Par exemple, pour savoir si une des valeurs recueillies n'est pas représentative, nous pouvons par exemple utiliser **les tests de Grubbs ou de Dixon**. Pour ce sujet, le lecteur pourra consulter un cours recommandé par le rédacteur.

Transformations normalisantes

Il n'est pas conseillé dans un premier temps, d'utiliser les transformations normalisantes, mais plutôt d'avoir une réflexion profonde sur la nature des données à analyser et sur le modèle statistique à utiliser. En dernier recours, nous pourrions les envisager, comme nous l'avons conseillé dans le paragraphe précédent. Il en existe un certain nombre. Voici les principales :

$$\begin{aligned} Y'_i &= \log(Y_i + c) && \text{la transformation logarithmique,} \\ &= Y_i^{\lambda} && \text{la transformation puissance,} \\ &= \Phi^{-1}(Y_i) && \text{la transformation réciproque,} \\ &= \arcsin(\sqrt{Y_i}) && \text{la transformation arc sinus} \\ &= \frac{Y_i^{\lambda}-1}{\lambda}, \text{ avec } \lambda > 0 && \text{de la racine carrée,} \\ & && \text{la transformation Box-Cox} \end{aligned}$$

Sommaire

- 1 Violation des conditions d'application
- 2 Transformation des variables
- 3 Grandeur de l'effet expérimental
- 4 Puissance
 - Puissance a posteriori
 - Détermination du nombre de répétitions
- 5 Facteur à effets aléatoires

Sommaire

Remarque

Il est conseillé, au sujet de ces transformations, de lire les pages 327 à 334 du livre de Howell (sixième édition) pour savoir comment les appliquer et dans quel cas il faut utiliser celle-ci plutôt que celle-la.

Et si rien ne marche

Si nous n'avons toujours pas les conditions requises après ces transformations, il faut alors utiliser le test non paramétrique de Kruskal-Wallis. Ce test sera présenté ici ultérieurement.

- 1 Violation des conditions d'application
- 2 Transformation des variables
- 3 **Grandeur de l'effet expérimental**
- 4 Puissance
 - Puissance a posteriori
 - Détermination du nombre de répétitions
- 5 Facteur à effets aléatoires

Contexte

À l'heure actuelle, il existe au moins six mesures de la grandeur de l'effet expérimental. Elles sont toutes différentes et prétendent toutes être moins biaisées que les autres mesures. Ici, dans ce cours nous présenterons uniquement une des mesures les plus courantes : le eta carré.

Mesure de la taille de l'effet η^2

Quand le test d'égalité des moyennes est rejeté (l'hypothèse nulle \mathcal{H}_0 est rejetée), nous pouvons souhaiter donner une mesure de la taille de la différence entre moyennes.

Nous définissons η^2 comme le « pourcentage » de la variabilité des données Y_{ij} expliquée par la différence entre les groupes :

$$\eta^2 = 1 - \frac{SC_R}{SC_{Tot}} = \frac{SC_F}{SC_{Tot}}$$

Exemple : D'après le livre de Georges Parreins, Techniques statistiques, moyens rationnels de choix et de décision

Nous voulons tester cinq types de carburateurs. Pour chaque type, nous disposons de six pièces que nous montons successivement en parallèle sur des voitures que nous supposons avoir des caractéristiques parfaitement identiques. Le tableau qui va suivre indique pour chaque essai la valeur d'un paramètre lié à la consommation.

Suite de l'exemple

	Car. 1	Car. 2	Car. 3	Car. 4	Car. 5
Essai 1	21	23	18	20	22
Essai 2	24	23	19	21	20
Essai 3	25	32	28	25	24
Essai 4	20	23	19	15	21
Essai 5	34	32	24	29	27
Essai 6	17	15	14	9	26

Suite de l'exemple

Afin de réaliser le test de Fisher de l'analyse de la variance à un facteur fixe, nous allons nous assurer que les conditions du modèle linéaire sont bien vérifiées. Pour cela, nous allons réaliser les différents tests grâce au logiciel R. Les lignes de commande à taper sont les suivantes :

```
> car<-rep(1:5,c(6,6,6,6,6))
> car<-factor(car)
> conso<-c(21,24,25,20,34,17,23,23,32,23,32,
15,18,19,28,19,24,14,20,21,25,15,29,9,22,20,
24,21,27,26)
> modele1<-aov(conso~car)
> shapiro.test(residuals(modele1))
```

Suite de l'exemple

Shapiro-Wilk normality test
data: residuals(modele1)
W = 0.9726, p-value = 0.6119

La p -valeur (0,6119) du test de Shapiro-Wilk étant strictement supérieure à $\alpha = 5\%$, le test n'est pas significatif. Vous conservez donc l'hypothèse nulle H_0 . Le risque d'erreur associé à cette décision est un risque de deuxième espèce β . Vous ne pouvez pas l'évaluer dans le cas d'un test de Shapiro-Wilk.

Suite de l'exemple

Il nous reste à vérifier la condition d'homogénéité des variances. Pour cela, nous allons utiliser le test de Bartlett, qui est un test paramétrique. La ligne de commande à taper est la suivante :

```
> bartlett.residuals(modele1)~car)
Bartlett test of homogeneity of variances
data: residuals(modele1) by car
Bartlett's K-squared = 3.9326, df = 4,
p-value = 0.4152
```

Suite de l'exemple

La p -valeur (0,4152) du test de Bartlett étant strictement supérieure à $\alpha = 5\%$, le test n'est pas significatif. Vous conservez donc l'hypothèse nulle \mathcal{H}_0 . Le risque d'erreur associé à cette décision est un risque de deuxième espèce β . Vous ne pouvez pas l'évaluer dans le cas d'un test de Shapiro-Wilk. Nous pouvons donc désormais construire le tableau de l'analyse de la variance et réaliser le test de Fisher.

Suite de l'exemple

Nous avons donc le tableau de l'analyse de la variance suivant :

Variation	SC	ddl	CM	F_{obs}
Due au facteur	108,3333	4	27,08332	0,8502509
Résiduelle	796,3333	25	31,85333	
Totale	904,6666	29		

Comme $F_{obs} < F_c(2, 75871)$, le test de Fisher n'est pas significatif. Nous décidons de ne pas rejeter l'hypothèse nulle (\mathcal{H}_0) et par conséquent de l'accepter. Le risque d'erreur associé à cette décision est un risque de deuxième espèce β qu'il faudra évaluer (voir le paragraphe suivant).

Suite de l'exemple

Nous calculons η^2 à titre d'exemple. En principe, nous le calculerons lorsque le test est significatif :

$$\eta^2 = \frac{108,3333}{904,6666} \simeq 0,1197.$$

Remarque

Nous retrouvons cette valeur en face de Multiple R-squared, résultat obtenu en tapant les deux lignes de commande suivantes :

- > model2<-lm(conso car)
- > summary(model2)

Sommaire

Rappel

Dans l'analyse de la régression linéaire simple, nous utilisons le coefficient de détermination R^2 pour mesurer le pourcentage de la variance de la variable Y expliquée par le modèle. Rappelons ici sa définition :

$$R^2 = 1 - \frac{SC_R}{SC_{Tot}} = \frac{SC_{Regression}}{SC_{Tot}}$$

Cette égalité ressemble beaucoup à celle qui définit le eta carré. Nous pouvons donc faire un parallèle entre ces deux mesures.

- 1 Violation des conditions d'application
- 2 Transformation des variables
- 3 Grandeur de l'effet expérimental
- 4 Puissance
 - Puissance a posteriori
 - Détermination du nombre de répétitions
- 5 Facteur à effets aléatoires



Cas de l'analyse de la variance à un facteur fixe

Nous nous intéressons à la puissance $1 - \beta$, où β est le risque de commettre une erreur de deuxième espèce, du test F d'analyse de la variance pour le test de l'hypothèse nulle

$$(\mathcal{H}_0) : \alpha_1 = \alpha_2 = \dots = \alpha_J = 0$$

contre l'hypothèse alternative

$$(\mathcal{H}_1) : \text{Il existe } i_0 \in \{1, 2, \dots, J\} \text{ tel que } \alpha_{i_0} \neq 0.$$

Calcul de la puissance

Cette puissance $1 - \beta$ est donnée par la formule suivante :

$$1 - \beta = \mathbb{P} \left[F'(I - 1, I(J - 1); \lambda) > F(I - 1, I(J - 1); 1 - \alpha) \right],$$

où $F(I - 1, I(J - 1); 1 - \alpha)$ est le $100(1 - \alpha)$ quantile de la loi de Fisher à $I - 1$ et $I(J - 1)$ degrés de liberté et $F'(I - 1, I(J - 1); \lambda)$ est une variable aléatoire qui suit une loi de Fisher non-centrale à $I - 1$ et $I(J - 1)$ degrés de liberté et de paramètre de non-centralité λ .



Suite et fin de l'exemple

Ensuite, il est facile de calculer

$$\beta = \mathbb{P} \left[F(4, 25; 1, 700427) < F(4, 25; 0, 95) \right] \text{ et d'en déduire } 1 - \beta.$$

Pour cela, il suffit de taper la ligne de commande suivante :

$$\text{pF}(\text{qf}(0.95, 4, 25), 4, 25, \text{nCP}=1.700427)$$

et nous obtenons $\beta = 0, 8678891$.

Par conséquent, la puissance est égale à 0, 1321109, qui est une puissance très faible.

Calcul du paramètre ϕ dans le cas déséquilibré

Si le nombre de répétitions n_i effectué pour chaque modalité i du facteur α n'est pas constant, c'est-à-dire si le plan expérimental n'est pas équilibré, le paramètre de non-centralité λ devient :

$$\lambda = \frac{1}{2\sigma^2} \sum_{i=1}^I n_i \alpha_i^2.$$

Le paramètre de non-centralité normalisé ϕ est alors :

$$\phi = \frac{1}{\sigma} \sqrt{\frac{1}{I} \sum_{i=1}^I n_i \alpha_i^2}.$$

Remarques

- Il faut avoir à l'esprit que nous sommes dans l'impossibilité de calculer exactement le paramètre ϕ ou le paramètre λ . (Il y a cette relation que nous venons d'exposer qui lie les deux paramètres.)
Au mieux, nous serons capable de donner une estimation de ϕ car nous ne pourrions jamais connaître la variance σ^2 de la population.
- Il est d'usage de travailler sur ϕ car les abaques que nous allons utiliser pour calculer les puissances se servent du paramètre ϕ et non du paramètre λ .

Sommaire

Calcul de la puissance - Suite et fin

Nous utilisons la formule ci-dessus pour déterminer les valeurs de J pour lesquelles la puissance $1 - \beta$ est supérieure à une valeur $1 - \beta_0$ fixée à l'avance, généralement 0,8 soit 80%.
Remarquons que là encore il est nécessaire de connaître σ^2 ou au moins d'avoir une idée précise de la valeur de ce paramètre ce qui n'est malheureusement généralement pas le cas. Dans cette situation, nous considérons plutôt le paramètre de sensibilité Δ/σ à la place de Δ .

- 1 Violation des conditions d'application
- 2 Transformation des variables
- 3 Grandeur de l'effet expérimental
- 4 Puissance
 - Puissance a posteriori
 - Détermination du nombre de répétitions
- 5 Facteur à effets aléatoires

Rappel

Dans l'analyse de la variance à un facteur à effets fixes avec l modalités, nous observons pour chaque modalité du facteur n_i réalisations indépendantes d'une variable aléatoire Y . Nous savons que le modèle utilisé dans cette analyse s'écrit :

$$Y_{ij} = \mu + \alpha_j + \mathcal{E}_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, l,$$

où \mathcal{E}_{ij} sont indépendantes et $\mathcal{L}(\mathcal{E}_{ij}) = \mathcal{N}(0; \sigma^2)$. Cette variable représente l'erreur commise lors des observations, μ désigne l'effet « global » ou moyenne générale de la variable aléatoire Y

et les effets α_j satisfont la contrainte : $\sum_{j=1}^l \alpha_j = 0$.

Un nouveau modèle

Mais ce modèle ne correspond pas toujours à la réalité. Dans certains cas, en particulier quand les modalités sont choisies au hasard, le fait de supposer que les effets sont fixes n'est pas adapté. Nous sommes amenés à considérer que chaque contribution α_j est une réalisation, indépendante des autres réalisations, d'une variable aléatoire A_j de loi $\mathcal{N}(0; \sigma_A^2)$, elle même indépendante de \mathcal{E} . Dans ces conditions, le modèle s'écrit :

$$Y_{ij} = \mu + A_j + \mathcal{E}_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, l.$$

Mise en place du test de l'effet du facteur aléatoire

Nous nous proposons de tester l'hypothèse nulle

$$(\mathcal{H}_0) : \sigma_A^2 = 0$$

contre l'hypothèse alternative

$$(\mathcal{H}_1) : \sigma_A^2 \neq 0.$$

Remarque

Ce test ne compare plus les moyennes mais teste au moyen de la variance du facteur A, si il y a un effet de ce facteur aléatoire.

Notations et propriétés

Si \bar{Y}_j et \bar{Y} désignent respectivement la moyenne des Y_{ij} où $j = 1, \dots, n_i$ et la moyenne de toutes les variables Y_{ij} , un calcul simple nous montre que les lois des trois variables sont :

$$\mathcal{L}(Y_{ij}) = \mathcal{N}(\mu_i; \sigma^2 + \sigma_A^2), \quad \mathcal{L}(\bar{Y}_j) = \mathcal{N}\left(\mu_i; \frac{\sigma^2}{n_i} + \sigma_A^2\right),$$

$$\mathcal{L}(\bar{Y}) = \mathcal{N}\left(\mu; \frac{\sigma^2}{n} + \frac{\sigma_A^2}{n^2} \sum_{i=1}^l n_i^2\right).$$

Tableau de l'ANOVA

Variation	SC	ddl	CM	F_{obs}	c
Due au facteur A	$\sum (\bar{Y}_j - \bar{Y})^2$	$l - 1$	s_A^2	$\frac{s_A^2}{s_R^2}$	c
Résiduelle	$\sum (Y_{ij} - \bar{Y}_j)^2$	$n - l$	s_R^2		
Totale	$\sum (Y_{ij} - \bar{Y})^2$	$n - 1$			

Remarque

Nous retrouvons strictement les mêmes formules que celles du cas de l'analyse de la variance à un facteur à effets fixes.

Propriété

Si les trois conditions sont satisfaites et si l'hypothèse nulle (\mathcal{H}_0) est vraie alors

$$F_{obs} = \frac{s_A^2}{s_R^2}$$

est une réalisation d'une variable aléatoire F qui suit une loi de Fisher à $l - 1$ degrés de liberté au numérateur et $n - l$ degrés de liberté au dénominateur. Cette loi est notée $\mathcal{F}_{l-1, n-l}$.

