

## Chapitre 03 : Analyse de la variance d'un facteur — ANOVA

### I. Définition et exemple:

1) Définition et principe de l'ANOVA

L'ANOVA d'un facteur permet de tester l'effet d'un facteur qualitatif ( $x$ ) ayant  $h$  modalités ou niveaux ( $h \geq 3$ ) sur une variable quantitative  $y$

⇒ objectif : Comparer les moyennes de la variable  $y$  pour chaque modalité du facteur

Le principe de l'ANOVA repose sur la dispersion des données autour de la moyenne (variance)

Cette dispersion peut avoir deux origines :

L'effet du facteur étudié :

Dans ce cas, une partie de la dispersion est imputable au niveau du facteur étudié. Cette partie de dispersion est appelée variabilité factuelle ou variabilité inter-classe

Variabilité résiduelle ou intra-classe :

Il s'agit de la variabilité qui reste lorsque la variabilité factuelle est soustraite à la variabilité totale, elle correspond à la part qui n'est pas expliquée  
ainsi le principe de l'ANOVA est de déterminer à l'aide d'un test statistique que si les moyennes de la variable quantitative  $y$  sont toutes égales ou non.

2) Exemple:

On considère une étude sur la plantation d'arbre dans 3 forêts.  
on se propose de comparer les hauteurs moyennes des arbres.  
pour ce faire on dispose de

forêt 1	forêt 2	forêt 3
23,3	18,9	22,5
24,4	21,1	22,9
24,6	21,1	23,7
24,9	22,1	24,0
25,0	22,5	24,0
26,2	23,5	24,5

soient les forêts les variables qualitatives ayant 3 modalités  
1, 2, 3.

La hauteur des arbres : la variable quantitative  $y$

Problématique : peut-on considérer que la hauteur moyenne  
des arbres diffère selon les forêts ?

II Modèles statistiques et hypothèses:

1) Modèle statistique:

ANOVA est un modèle de régression linéaire qui s'écrit

$$y_{ij} = \mu_i + \varepsilon_{ij} \quad \forall \begin{cases} i = 1, \dots, k \\ j = 1, \dots, n_i \end{cases}$$

Avec :  $\mu_i$  : est le paramètre qui représente l'effet de la  
modalité  $i$  de la variable qualitative  $x$

$\varepsilon_{ij}$  : ce sont les résidus c'est-à-dire les écarts entre  
les observations et les moyennes des groupes  
auxquelles elles sont relatives  $\varepsilon_{ij} = y_{ij} - \bar{y}_i$

## 2) Hypothèses de l'ANOVA:

$H_1$ : les résidus  $\varepsilon_{ij}$  sont indépendants

$H_2$ : les résidus sont homogènes ou encore homoscedastiques

$$\mu_{ij} = \mu_i + \varepsilon_{ij} \quad \begin{cases} i = 1, \dots, 3 \\ j = 1, \dots, 6 \end{cases}$$

## III Notation et calcul de la dispersion totale:

### 1) Notation:

\* Le facteur étudié (ici, les fûts) comporte 3 modalités

\* le nombre d'observation pour chacune des modalités est noté  $n_i$  (ici,  $n_1 = n_2 = n_3 = 6$ )

\* le nombre totale d'observations est noté  $n$  tel que  $n = \sum_{i=1}^3 n_i$ ,  
(ici,  $n = 6 + 6 + 6 = 18$ )

\* les observations sont notées  $y_{ij}$   
(ici  $y_{ij}$  correspond à la hauteur des arbres)

← c'est l'observation relative à la hauteur du  $j^{\text{e}}$  arbre dans les modalités

\*  $i$ : est l'indice des modalités telle que  
 $i = 1, \dots, 3$  (ici  $i = 1, \dots, 3$ )

\*  $j$ : est l'indice de l'observation au sein d'une modalité tel que  $j = 1, \dots, n_i$  (ici  $j = 1, \dots, 6$ )

les moyennes des observations de chaque modalité sont

noté  $\bar{y}_i = \bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$



\* la moyenne générale ou globale des observations est noté  $\bar{y}$

$$\Rightarrow \bar{y} = \frac{1}{n} \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij}) = \frac{1}{n} \sum_{i=1}^h n_i \bar{y}_i$$

2) Calcul de la dispersion totale:

Pour mesurer la dispersion / variabilité totale on utilise la somme des carrés totale SCT qui correspond à la somme des distances au carré entre chaque valeur ~~notée~~ observée et la moyenne générale ou globale. la SCT se calcule

$$SCT = \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2$$

IV Décomposition de la variance ou de la dispersion totale  
La variabilité totale mesurée par la SCT est décomposée en + la variabilité factorielle (due au facteur qualitatif)

notée par la somme des carrés factorielle

$$\Rightarrow SCF = \sum_{i=1}^h \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2 = \sum_{i=1}^h n_i (\bar{y}_i - \bar{y})^2$$

=> variabilité inter classe.

\* la variabilité résiduelle ou intra classe mesurée par la somme des carrés résiduelles noté SCR

$$\Rightarrow SCR = \sum_{i=1}^h \sum_{j=1}^{n_i} \epsilon_{ij}^2$$

$$= \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

l'équation de l'ANOVA s'écrit comme suit:

$$SCT = SCF + SCR \Rightarrow \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2 = \sum_{i=1}^h \sum_{j=1}^{n_i} (\bar{y}_i - \bar{y})^2 + \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

# V Tests d'hypothèses : Teste de Fisher

## 1) Tableau d'ANOVA de la variance

Source de variation	Somme des carrés	d.d.l	carrés moyens	Statistique F de Fisher	
				calculé	théorique
actuel	$SCF = \sum_{i=1}^b n_i (\bar{y}_i - \bar{y})^2$	$b-1$	$CM_F = \frac{SCF}{b-1}$	$F_c = \frac{CM_F}{CM_R}$	$F_\alpha(b-1, n-b)$
résidu	$SCR = \sum_{i=1}^b \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$	$n-b$	$CM_R = \frac{SCR}{n-b}$		
total	$SCT = \sum_{i=1}^b \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2$	$n-1$	variance empirique $= \frac{SCT}{n-1}$		

## 2) Teste de Fisher :

### \* Hypothèses :

$H_0 : \mu_1 = \mu_2 = \dots = \mu_b$  contre  $H_1 : \text{il existe au moins un } \mu_i \text{ différent des autres.}$

### \* Statistique et loi :

Dans  $H_0$  vraie on a 
$$F_c = \frac{SCF/b-1}{SCR/n-b} = \frac{\sum_{i=1}^b n_i (\bar{y}_i - \bar{y})^2 / b-1}{\sum_{i=1}^b \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 / n-b} \sim F_{b-1, n-b}$$

### \* Règle de décision

• si  $F_c \leq F_\alpha = F_\alpha(b-1, n-b)$  alors  $H_0$  est vraie

• si  $F_c > F_\alpha$  alors  $H_1$  est vraie

1  
2  
3  
4

## VI Application à l'exemple :

\* moyenne de chaque modalité :

$$\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij})$$

$$* \bar{y}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} (y_{1j})$$

$$\Rightarrow \boxed{\bar{y}_1 = 24,733}$$

$$* \bar{y}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} (y_{2j})$$

$$\Rightarrow \boxed{\bar{y}_2 = 21,533}$$

$$* \bar{y}_3 = \frac{1}{n_3} \sum_{j=1}^{n_3} (y_{3j})$$

$$\boxed{\bar{y}_3 = 23,6}$$

\* moyenne générale des observations :

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij}) \\ &= \frac{1}{n} \sum_{i=1}^h n_i \times \bar{y}_i \end{aligned}$$

$$\Rightarrow \bar{y} = \frac{1}{18} (n_1 \times \bar{y}_1 + n_2 \times \bar{y}_2 + n_3 \times \bar{y}_3)$$

$$= \frac{1}{18} (6 \times 24,733 + 6 \times 21,533 + 6 \times 23,6)$$

$$\boxed{\bar{y} = 23,289}$$

\* Règle de décision :

$$* \text{ si } F_c \leq F_t \Rightarrow \text{ad ( } \cancel{h-1} \text{ )}$$



# \* Equation d'analyse de la variance

2

$$SCT = SCF + SCA \Leftrightarrow$$

$$\sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2 = \sum_{i=1}^h n_i (\bar{y}_i - \bar{y})^2 + \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

$$\begin{aligned} * SCF &= \sum_{i=1}^h n_i (\bar{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^3 n_i (\bar{y}_i - \bar{y})^2 \end{aligned}$$

$$\Rightarrow SCF = n_1 (\bar{y}_1 - \bar{y})^2 + n_2 (\bar{y}_2 - \bar{y})^2 + n_3 (\bar{y}_3 - \bar{y})^2$$

$$\boxed{SCF = 31,592}$$

$$\begin{aligned} * SCA &= \sum_{i=1}^h \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 \\ &= \sum_{i=1}^3 \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 \end{aligned}$$

$$\Rightarrow \boxed{SCA = 19,707}$$

$$\begin{aligned} * SCT &= SCF + SCA \\ &= 31,592 + 19,707 \\ &= 51,3 \end{aligned}$$

## \* Tableau d'analyse de la variance:

Source de variation	Somme des carrés	d.d.l	Carrés moyens	Statistiques de Fisher	
				calculée	théorique
facteur	SCF = 31,592	2	CM <sub>F</sub> = 15,796	F <sub>c</sub> = 12,021	F <sub>t</sub> = F <sub>5%</sub> (2,15) = 3,68
résidu	SCA = 19,707	15	CM <sub>A</sub> = 1,314		
totale	SCT = 51,3	17	3,018		

\* Test de Fisher :

\* Hypothèses

$H_0 : \mu_1 = \mu_2 = \mu_3$  contre  $H_1$  : il existe au moins un  $\mu_i$  différent des autres

\* Sous  $H_0$  vraie, on a :

$$F_0 = \frac{SCF / k - 1}{SCM / n - k} \rightsquigarrow F_{\alpha}(k-1, n-k) = F_{5\%}(2, 15)$$

\* Règle de décision

Si  $F_0 \leq F_t$  alors  $H_0$  est vraie sinon  $H_1$  est vraie

Conclusion :

$$\text{on a } F_0 = 12,021 > F_t = 3,68$$

Donc  $H_1$  est vraie  $\Rightarrow$  oui, la hauteur moyenne des arbres diffère selon la forêt.

Question 2 ex chapitre 2 :

Test d'hypothèse jointes

$$* \frac{\hat{\beta}_1 - 2\hat{\beta}_2}{\sqrt{\hat{V}(\hat{\beta}_1 - 2\hat{\beta}_2)}} \rightsquigarrow st(T-k) \approx st(20)$$

$$\Rightarrow t_{\frac{\alpha}{2}}^{T-k} = t_{0,025}^{20} = 2,086$$

$$\hat{V}(\hat{\beta}_1 - 2\hat{\beta}_2) = \hat{V}\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 4\hat{V}\left(\begin{pmatrix} 1 \\ 2 \end{pmatrix}\right) - 4\text{Cov}\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \end{pmatrix}\right)$$

$$\hat{\Sigma}_{\hat{\beta}} = \hat{V}_{\text{var}}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

$$\text{matrice } \hat{\Sigma}_{\hat{\beta}} \text{ et } \sigma^2 = \frac{SCM}{T-k}$$



$$R^2 = 1 - \frac{SCR}{SCT}$$

$$\Leftrightarrow SCR = (1 - R^2) \times SCT$$

$$= (1 - 0,937) \times 317,46$$

$$= 20$$

$$\Rightarrow \hat{R}^2 = \frac{20}{20} = 1$$

$$\Rightarrow \hat{\beta} = (X'X)^{-1}$$

$$\Rightarrow \hat{\beta} = \begin{pmatrix} V(\hat{\beta}_1) & \text{cov}(\hat{\beta}_1, \hat{\beta}_2) & \text{cov}(\hat{\beta}_1, \hat{\beta}_0) \\ V(\hat{\beta}_2) & \text{cov}(\hat{\beta}_2, \hat{\beta}_1) & V(\hat{\beta}_0) \end{pmatrix}$$

$$\hat{y}_t = 0,51x_{1t} + 0,35x_{2t} + 27,3$$

$$\hat{V}(\hat{\beta}_1) = 0,008$$

$$\hat{V}(\hat{\beta}_2) = 0,0025$$

$$\text{cov}(\hat{\beta}_1, \hat{\beta}_2) = -0,08$$

$$\text{Si } t_c = \left| \frac{\hat{\beta}_1 - 2\hat{\beta}_2}{\sqrt{\hat{V}(\hat{\beta}_1 - 2\hat{\beta}_2)}} \right| < t_{\frac{\alpha}{2}}^{T-K} \text{ alors } H_1 \text{ est vraie}$$

$$\text{Si non, } H_1 \text{ est vraie} \Rightarrow t_{\frac{\alpha}{2}}^{T-K} = t_{0,025}^{20} = 2,086$$

$$\hat{V}(\hat{\beta}_1 - 2\hat{\beta}_2) = 0,33$$

$$t_c = \left| \frac{0,51 + 2 \times 0,35}{\sqrt{0,33}} \right| = 2,106 > 2,086$$

$$\Rightarrow \text{on accepte } H_1: \beta_1 \neq 2\beta_2$$