

Université Lumière Lyon 2

M2-SISE

Statistique et Informatique pour la Science des données
Support 2 : ANOVA-ANCOVA

Rafik Abdesselam

Courriel : rafik.abdesselam@univ-lyon2.fr

Web : <http://perso.univ-lyon2.fr/~rabdesse/Fr/>

Support pédagogique : <http://perso.univ-lyon2.fr/~rabdesse/Documents/>

Janvier 2023

Plan du cours



ANOVA - ANCOVA

- Introduction
- Analyse de la variance
- Analyse de la covariance

Objectif du cours

- Présenter les **concepts** et les **conditions** d'application de l'analyse de la variance.
- L'analyse de variance abrégée sous le terme anglais **ANOVA** (ANalysis Of VAriance) est une technique statistique permettant de comparer **les moyennes** d'un nombre quelconque de populations, contrairement à ce que pourrait laisser penser son nom.
- Son lien avec la régression est présenté, mais d'une façon générale, elle consiste à **comparer plusieurs moyennes** d'échantillons provenant de populations **normales**.
- Une large place est accordée dans ce cours aux exemples et exercices sur données réelles traités avec les logiciels SAS et SPAD.

Objectif du cours

- **Pré-requis** : Quelques notions de Statistique & Probabilités et de statistique inférentielle.
- **Quelques références bibliographiques** :

[1] Dagnelie, P. Analyse statistique à plusieurs variables. Gembloux, Presses agronomiques, 1986, 362p.

[2] Mervyn, G.Marasinghe, William J. Kennedy, SAS for Data Analysis, Intermediate Statistical Methods. Statistics and Computing, Springer, 2008.

[3] Saporta, G. Probabilités Analyse des données et Statistique. Editions technip, 1990.

Introduction

- Les techniques d'analyse de variance **ANOVA** sont des outils entrant dans le cadre du **Modèle Linéaire Général** et où une **variable quantitative est expliquée** par **une ou plusieurs variables qualitatives**.
- L'objectif est alors de **comparer les moyennes** empiriques de la variable quantitative observée pour les différentes catégories d'unités statistiques.
- Ces catégories sont définies par l'observation des variables qualitatives ou **facteurs** prenant différentes modalités ou encore de variables quantitatives découpées en **classes ou niveaux**.
- Une combinaison de niveaux définit une **cellule**, groupe ou **traitement**.

Introduction

- L'analyse de variance ANOVA est une technique statistique de tests et d'estimation qui permet d'**analyser l'effet** d'une voire plusieurs variables **qualitatives** sur une variable **continue**.
- L'ANOVA est très utilisée dans le contexte des plans d'expériences et des traitements des données expérimentales.
- D'un point de vue modélisation, l'ANOVA n'est autre qu'une **régression multiple sur variables explicatives qualitatives** (nominales).
- Les principes essentiels de l'analyse de la variance à un et à deux facteurs de classification avec/sans interaction seront exposés.

Terminologie

- Les variables qualitatives sont appelées "**facteurs ou critères**" et leurs modalités "**niveaux**" du facteur. En présence de plusieurs facteurs, une combinaison de niveaux est un "**traitement**".
- Statistiquement, l'ANOVA est une **généralisation du test de Student** pour comparer plus de 2 moyennes.
- Il arrive fréquemment que les données soient groupées en classes selon certains critères ou facteurs tels que, par exemple, l'âge, l'appartenance sociale, la région géographique, etc.

Exemple introductif

- Prenons comme exemple, le cas d'une étude sur la fréquence d'utilisation des moyens de transports en commun.
- On peut supposer que celle-ci sera différente en fonction de l'âge des personnes interrogées.
- Il est donc naturel de diviser la population en classes d'âges (par exemple : adolescents, adultes, personnes âgées) avant d'effectuer l'échantillonnage.
- Sur la base des observations des différents échantillons constitués, la question sera de savoir s'il existe effectivement une **différence significative** d'utilisation des transports en commun **entre les classes d'utilisateurs** considérées.
- Ceci revient à effectuer un **test de comparaison multiple de moyennes**.

Exemple illustratif

- Les données représentent la fréquence journalière d'utilisation des moyens de transports en commun de trois groupes d'utilisateurs.

	Adolescents	Adultes	Personnes âgées
	3	5	3
	6	7	3
	5	6	2
	6	7	2
	5	5	5
Moyenne	5	6	3

- Le problème consiste à **détecter les différences**, si elles existent, **entre les moyennes des populations** à partir desquelles ces observations ont été obtenues.
- Comparer la différence entre les moyennes des groupes d'utilisateurs, mesurée en terme de **variabilité**, tout en tenant compte de la variabilité existant entre les usagers à l'intérieur de chaque groupe.

Exemple : Cas particulier 1

- Pour bien distinguer entre ces **deux types de variabilité**, considérons les données des deux tableaux fictifs suivants.
- Variation nulle à l'intérieur des groupes

	Adolescents	Adultes	Personnes âgées
	5	6	3
	5	6	3
	5	6	3
	5	6	3
	5	6	3
Moyenne	5	6	3

- Toutes les observations dans chaque échantillon ont ici la même valeur. Il n'y a donc **aucune variation à l'intérieur des groupes** (ou échantillons d'utilisateurs), mais il y a **une variation entre les groupes** d'utilisateurs, puisque les moyennes d'échantillonnage sont différentes.

Exemple : Cas particulier 2

- Variation nulle entre les groupes

	Adolescents	Adultes	Personnes âgées
	8	4	3
	6	5	7
	5	9	7
	4	7	8
	7	5	5
Moyenne	6	6	6

- Par contre, dans ce cas, la moyenne de chaque groupe d'utilisateurs est identique. Il n'y a donc **pas de variation entre les groupes**, mais il y a **une variation à l'intérieur des groupes** puisque toutes les observations dans chaque groupe n'ont pas la même valeur.
- En pratique, les observations obtenues ne sont ni exactement identiques, ni de moyennes égales ; elles sont hétérogènes comme les données de l'exemple introductif.

ANOVA à un facteur contrôlé

- **Les données** : n observations réparties dans p groupes ou échantillons. Chaque échantillon j ($j = 1, p$) contient n_j observations, correspondant à un niveau différent d'un facteur. Les tailles n_j des échantillons pouvant être égales ou différentes, $n = \sum_j^p n_j$.
- Le tableau suivant illustre un exemple de données.

Ech.1	Ech.2	...	Ech.j	...	Ech.p
x_{11}	x_{12}	...	x_{1j}	...	x_{1p}
x_{21}	:	:	:	:	:
x_{i1}	x_{i2}	...	:	:	x_{ip}
:	x_{n_22}	:	:	:	:
:		:	x_{nj}	:	:
x_{n_11}		...			:
		...			x_{n_pp}
\bar{x}_1	\bar{x}_2	...	\bar{x}_j	...	\bar{x}_p

$\bar{x}_j = \frac{\sum_{i=1}^{n_j} x_{ij}}{n_j}$ la moyenne des n_j observations du groupe j .

$\bar{x} = \frac{\sum_{j=1}^p \sum_{i=1}^{n_j} x_{ij}}{n} = \frac{\sum_{j=1}^p n_j \bar{x}_j}{n}$ la moyenne globale.

ANOVA à un facteur

- Le facteur à p niveaux est supposé avoir un effet uniquement sur les moyennes des distributions et non sur leur variance. Il s'agit d'un **test de comparaison de p moyennes**.
- D'une façon générale, il s'agit de tester s'il existe une différence "significative" entre les moyennes m_j ($j = 1, p$) des p populations dans lesquelles ont été prélevés les p échantillons indépendants de taille n_j ($j = 1, p$) de l'étude.
- En d'autres termes, effectuer le **test statistique** suivant :
$$\begin{cases} H_0 : m_1 = m_2 = \dots = m_p \text{ pas de différence significative.} \\ H_1 : m_j \neq m_k \quad j \neq k \text{ différence significative.} \end{cases}$$
- Il suffit donc qu'une moyenne soit différente de toutes les autres pour que l'hypothèse nulle H_0 soit rejetée.
- Il s'agit là d'une généralisation à p populations du test classique (t de Student) de comparaison de moyennes de 2 échantillons.

Conditions d'application - Ecart

- Conditions d'application à vérifier avant toute utilisation d'une analyse de la variance :
 - 1 Les échantillons doivent être indépendants.
 - 2 Les distributions des populations considérées doivent être normales (hypothèse de normalité - test paramétrique).
 - 3 Les populations d'où sont prélevés les échantillons doivent posséder la même variance : $\sigma^2 = \sigma_1^2 = \sigma_2^2 = \dots = \sigma_p^2$ (hypothèse d'homocédasticité).
- Pour procéder à une analyse de la variance, on s'intéresse à trois types d'écart :
 - 1 Chaque observation par rapport à sa moyenne respective : $(x_{ij} - \bar{x}_j)$.
 - 2 Chaque moyenne d'échantillonnage par rapport à la moyenne globale : $(\bar{x}_j - \bar{x})$.
 - 3 Chaque observation par rapport à la moyenne globale : $(x_{ij} - \bar{x})$.

Tableau d'analyse de la variance à 1 facteur

- Les résultats obtenus sont résumés dans un tableau d'analyse de la variance (ou tableau ANOVA).

Sources de variation	Somme des carrés	Degrés de liberté	Carrés moyens	F
Entre Inter / Between	SC_{ent} $\sum_{j=1}^p n_j (\bar{x}_j - \bar{x})^2$	p - 1	$CM_{ent} = s_{ent}^2$ $\frac{SC_{ent}}{p-1}$	$\frac{s_{ent}^2}{s_{int}^2}$
Intérieur Intra / Within	SC_{int} $\sum_{j=1}^p \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$	n - p	$CM_{int} = s_{int}^2$ $\frac{SC_{int}}{n-p}$	
Totale Total	SC_{tot} $\sum_{j=1}^p \sum_{i=1}^n (x_{ij} - \bar{x})^2$	n - 1		

- La variation totale est la somme de 2 variations : $SC_{tot} = SC_{ent} + SC_{int}$.
- Cette propriété montre pourquoi la technique de comparaison de moyennes est appelée analyse de la variance, car ces sommes de carrés sont utilisées pour estimer des variances.
- En effet, le test réalisé consiste à décomposer la variance (constante) de x en deux parties : une variance interclasse ($CM_{ent} = s_{ent}^2$) et une variance intraclasse ou erreur ($CM_{int} = s_{int}^2$) puis à établir le test de Fisher (rapport de 2 variances $F = \frac{s_{ent}^2}{s_{int}^2}$).

ANOVA à 1 facteur contrôlé

- **Approche régression** : On peut associer à l'analyse de la variance à 1 facteur le modèle de régression linéaire multiple suivant :

$$y_{ij} = \mu + \alpha_j + \varepsilon_{ij}, \quad i = 1, \dots, n_j \text{ et } j = 1, \dots, p.$$

où, y_{ij} désigne la i ème observation du j ème échantillon,

μ est la moyenne générale commune,

α_j est l'effet du niveau j ème du facteur,

ε_{ij} est l'erreur relative à l'observation y_{ij} .

- L'objectif est de tester certains paramètres du modèle notamment l'hypothèse nulle suivante : $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_p$
ce qui signifie que les p facteurs ont un effet identique. L'hypothèse alternative : $H_1 : \text{au moins une différence } \alpha_j \neq \alpha_{j'} \text{ } j \neq j'$
c'est-à-dire que les α_j ne sont pas tous identiques.
- Modèle de régression à $(p + 1)$ paramètres à estimer.

Exemple d'application

- Les données représentent le niveau de production de 27 cadres employés, affectés à une tâche d'assemblage et de vérification, selon leur statut dans l'entreprise.

Niveau de Production selon le **Statut - Niveau de responsabilité** de l'employé.

	Junior	Intermédiaire	Supérieur
	45 49 45	51 51 49	49 50 52
	45 48 45	50 46 49	48 48 51
	47 47 46	50 46 51	51
		48 49	
Effectif	9	11	7
Moyenne	46.333	49.091	49.857

Exemple d'application - Programme SAS

```
/* ANOVA - Données productivité */
data product;
/* Déclaration des variables */
input NUM$ NPRO TAPT SEXE$ NRES$;
label EMPL = 'Employé(e)'
NPRO = 'Niveau de production'
TAPT = 'Test d"aptitude'
NRES = 'Niveau de responsabilité';
/* Lecture des données */
cards;
1 45 105 M J
2 49 120 M S
etc.
26 46 111 M J
27 51 124 F S
run;

/* ANOVA à 1 facteur contrôlé NRES - Vérification des hypothèses de base du modèle - Hypothèse H1 : normalité */
proc sort; by NRES; run;
proc univariate normal; by NRES; run; quit;
/*Comparaison de moyennes - Hypothèse H2 : Homoscédasticité */
proc glm; class NRES; model NPRO = NRES;
means NRES / tukey lines hovtest = levene; /*bartlett */
run;quit;

/* ANOVA à 1 facteur contrôlé SEXE - Vérification des hypothèses de base du modèle - Hypothèse H1 : normalité */
proc sort; by SEXE; run;
proc univariate normal; by SEXE; run; quit;
/*Comparaison de moyennes - Hypothèse H2 : Homoscédasticité */
proc glm; class SEXE; model NPRO = SEXE;
means SEXE / tukey lines hovtest = levene; /*bartlett */
run;quit;

/* ANOVA à 2 facteurs contrôlés SEXE NRES avec interaction*/
proc GLM; class NRES SEXE; model NPRO = NRES SEXE NRES*SEXE;
means NRES SEXE / tukey lines; run; quit;
```

Conditions d'application

- Effet du **Statut** de l'employé sur le niveau de Production.

1 Hypothèse de normalité :

The Univariate Procedure : Tests for Normality

	Test		Statistic		p-Value
Junior	Shapiro-Wilk	W	0.851103	<i>Pr < W</i>	0.0766 ✓
Intermédiaire	Shapiro-Wilk	W	0.869325	<i>Pr < W</i>	0.0760 ✓
Supérieur	Shapiro-Wilk	W	0.913363	<i>Pr < W</i>	0.4197 ✓

2 Hypothèse d'homoscédasticité - Egalité des variances :

The GLM Procedure : Bartlett's Test for Homogeneity of Production Variance

Source	DF	Chi-Square	<i>Pr > ChiSq</i>
Statut	2	0.3344	0.8461 ✓

Test de comparaison de moyennes

- 1 Hypothèses statistiques $\begin{cases} H_0 : \alpha_J = \alpha_I = \alpha_S \\ H_1 : \text{les moyennes ne sont pas toutes égales} \end{cases}$
- 2 Seuil de signification : $\alpha = 5\%$
- 3 Conditions d'application du test : Le niveau de production doit être normalement distribué selon chaque niveau de responsabilité (facteur) - Egalité des variances $\sigma_J^2 = \sigma_I^2 = \sigma_S^2$.
- 4 Statistique de test "Unilatéral - risque à droite" :
La fonction discriminante de Fisher sous H_0 :
$$F = \frac{S_{ext}^2}{S_{int}^2} \rightarrow F_{p-1=2; n-p=24} \text{ d.d.l.}$$
- 5 Calcul de la statistique de test sous l'hypothèse nulle $H_0 : \alpha_J = \alpha_I = \alpha_S$
$$f_0 = \frac{29.26}{2.74} = 10.68$$
- 6 Règle de décision : $f_{5\%} = 3.40$ cf. Table de la loi de Fisher $F_{(2;24)}$.
- 7 Décision et conclusion : f_0 appartient à la zone de rejet de H_0 ($f_0 = 10.68 > f_{5\%} = 3.40$), on peut donc conclure, avec risque d'erreur $\alpha = 5\%$, qu'il y a une différence significative. Les moyennes ne sont pas significativement égales. Il y a des disparités entre les niveaux moyens de production selon le niveau de responsabilité de l'employé dans l'entreprise.

Principaux résultats - SAS

The MEANS Procedure : Analysis Variable : Production

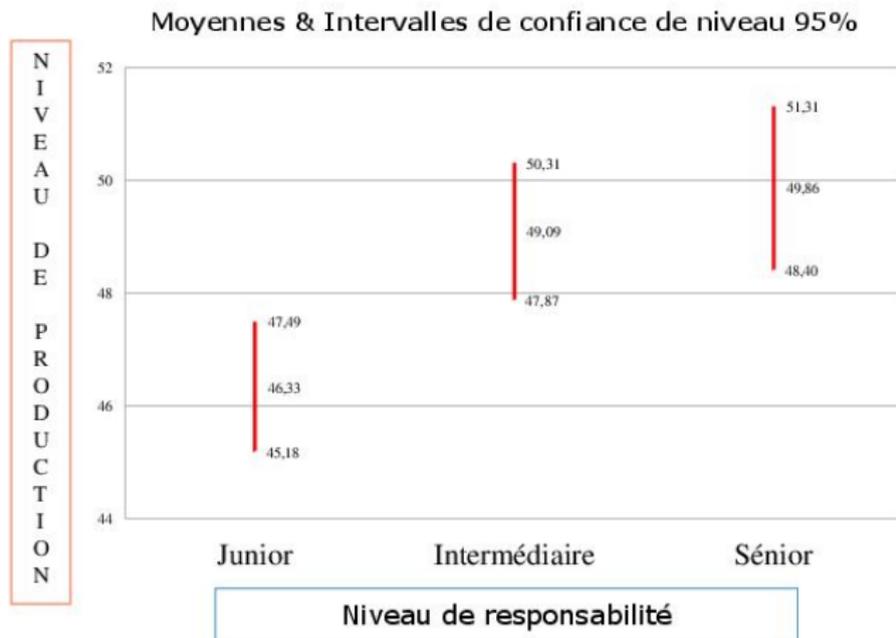
Statut	N	Mean	Std Dev	Minimum	Maximum
Junior	9	46.3333333	1.5000000	45.0000	49.0000
Intermédiaire	11	49.0909091	1.8140863	46.0000	51.0000
Supérieur	7	49.8571429	1.5735916	48.0000	52.0000

The ANOVA Procedure : Effet du Statut sur le niveau de production

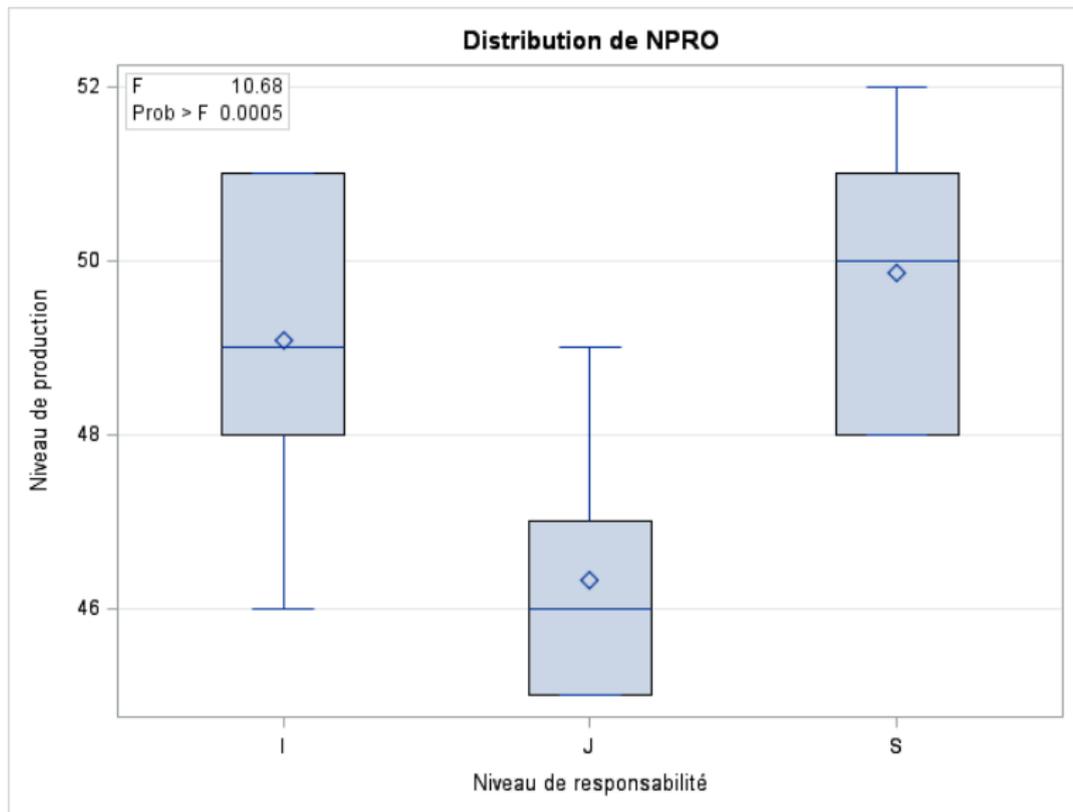
Source	Sum of squares	DF	Mean square	F-value	Pr > F
Model	58.5300625	2	29.2650313	10.68	0.0005 ✓
Error	65.7662338	24	2.7402597		
Total	124.2962963	26			

Principaux résultats

- Intervalles de confiance



Résultats SAS - BoxPlot



Principaux résultats - SAS

- Localiser les disparités

Means with the same letter are not significantly different

Tukey Grouping	Mean	N	Statut
A	49.8571	7	Supérieur
A	49.0909	11	Intermédiaire
B	46.3333	9	Junior

Comparisons significant at the 0.05 level are indicated by ***

Statut Comparison	Difference			
	Between Means	95%	Confidence Limits	
S - I	0.7662	-0.8856	2.4181	
S - J	3.5238	1.8020	5.2456	***
I - S	-0.7662	-2.4181	0.8856	
I - J	2.7576	1.2220	4.2932	***
J - S	-3.5238	-5.2456	-1.8020	***
J - I	-2.7576	-4.2932	-1.2220	***

ANOVA à deux facteurs

- Supposons maintenant le cas où les n observations de l'échantillon sont classées selon 2 facteurs **A** et **B** respectivement à p et q modalités-niveaux.
- Les n observations peuvent être réparties dans un tableau à p lignes (facteur **A**) et q colonnes (facteur **B**).

Effectif Moyenne	B ₁	...	B _k	...	B _q	Ensemble
A ₁	n ₁₁	...	n _{1k}	...	n _{1q}	n _{1.}
	\bar{x}_{11}	...	\bar{x}_{1k}	...	\bar{x}_{1q}	$\bar{x}_{1.}$
:	:	...	:	...	:	:
A _j	n _{j1}	...	n _{jk}	...	n _{jq}	n _{j.}
	\bar{x}_{j1}	...	\bar{x}_{jk}	...	\bar{x}_{jq}	$\bar{x}_{j.}$
:	:	...	:	...	:	:
A _p	n _{p1}	...	n _{pk}	...	n _{pq}	n _{p.}
	\bar{x}_{p1}	...	\bar{x}_{pk}	...	\bar{x}_{pq}	$\bar{x}_{p.}$
Ensemble	n _{.1}	...	n _{.k}	...	n _{.q}	n = n _{..}
	$\bar{x}_{.1}$...	$\bar{x}_{.k}$...	$\bar{x}_{.q}$	$\bar{x} = \bar{x}_{..}$

ANOVA à deux facteurs

- Les trois questions principales que l'on se pose lors d'une analyse de variance à 2 facteurs :
 - 1 Y a-t-il un effet du facteur A : les moyennes mesurées sur les p populations définies par le facteur A sont-elles différentes ?
 - 2 Y a-t-il un effet du facteur B : les moyennes mesurées sur les q populations définies par le facteur B sont-elles différentes ?
 - 3 Y a-t-il un effet conjugué des deux facteurs A et B : une interaction entre les moyennes du facteur A et celles du facteur B ?

Ces 3 effets sont analysés à l'aide de la statistique de test de Fisher-Snedecor : rapport de 2 variances.

ANOVA à 2 facteurs contrôlés

- **Approche régression** : on peut associer à l'analyse de la variance à 2 facteurs le modèle de régression suivant :

$$y_{ijk} = \mu + \alpha_j + \beta_k + \gamma_{jk} + \varepsilon_{ijk}$$

où, ε_{ijk} i.i.d. $\rightarrow N(0; \sigma^2)$ $i = 1, \dots, m$; $j = 1, \dots, p$ et $k = 1, \dots, q$.

- Les trois questions principales évoquées précédemment se traduisent de la façon suivante :

① $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_p$

② $H_0 : \beta_1 = \beta_2 = \dots = \beta_q$

③ $H_0 : \gamma_{11} = \gamma_{12} = \dots = \gamma_{jk} = \dots = \gamma_{pq}$

- Modèle de régression à $(1 + p + q + pq)$ paramètres à estimer.

Tableau d'analyse de la variance à 2 facteurs

Sources de variation	Somme carrés	Degrés de liberté	Carrés moyens	F
Facteur A	SC_{entA}	$p - 1$	$CM_{entA} = \frac{SC_{entA}}{p-1}$	$\frac{CM_{entA}}{CM_{int}}$
Facteur B	SC_{entB}	$q - 1$	$CM_{entB} = \frac{SC_{entB}}{q-1}$	$\frac{CM_{entB}}{CM_{int}}$
Interaction AB	SC_{entAB}	$(p - 1)(q - 1)$	$CM_{entAB} = \frac{SC_{entAB}}{(p-1)(q-1)}$	$\frac{CM_{entAB}}{CM_{int}}$
Intérieur	SC_{int}	$n - pq$	$CM_{int} = \frac{SC_{int}}{n-p}$	
Total	SC_{tot}	$n - 1$		

- $SC_{entA} = \sum_{j=1}^p n_j (\bar{x}_j - \bar{x})^2$ Les variations de la moyenne marginale \bar{x}_j du facteur A autour de sa moyenne α_j .
- $SC_{entB} = \sum_{k=1}^q n_{.k} (\bar{x}_{.k} - \bar{x})^2$ Les variations de la moyenne marginale \bar{x}_j du facteur B autour de sa moyenne β_k .
- $SC_{entAB} = \sum_{j=1}^p \sum_{k=1}^q n_{jk} (\bar{x}_{jk} - \bar{x}_j - \bar{x}_{.k} + \bar{x})^2$ Les fluctuations de \bar{x}_{jk} de A et B autour de la moyenne μ , abstraction faite des variations des moyennes marginales \bar{x}_j et \bar{x}_k ; elle mesure l'influence de l'interaction des facteurs A et B sur la moyenne.
- $SC_{int} = \sum_{i=1}^n \sum_{j=1}^p \sum_{k=1}^q (x_{ijk} - \bar{x}_{jk})^2$ Somme des carrés résiduels - Les fluctuations aléatoires des x_{ijk} autour de \bar{x}_{jk} dans le traitement (A_j, B_k) .
- $SC_{tot} = \sum_{i=1}^n \sum_{j=1}^p \sum_{k=1}^q (x_{ijk} - \bar{x})^2 = SC_{int} + SC_{entA} + SC_{entB} + SC_{entAB}$

Exemple d'application

- Les données représentent le niveau de production de 27 cadres employés, affectés à une tâche d'assemblage et de vérification, selon leur statut - niveau de responsabilité et leur Sexe - genre dans l'entreprise.
- SPAD : Répartition de la production selon le statut et le Sexe de l'employé.

Effectif Moyenne Ecart-type	Féminin	Masculin	Ensemble
Intermédiaire	9 49.778 1.030	2 46.000 0.000	11 49.091 1.730
Junior	2 48.500 0.500	7 45.714 0.881	9 46.333 1.414
Supérieur	5 50.400 1.356	2 48.500 0.500	7 49.857 1.457
Ensemble	16 49.813 1.236	11 46.273 1.286	27 48.370 2.146

Conditions d'application

- Effet du **Sexe** de l'employé sur le niveau de Production.

1 Hypothèse de normalité :

The Univariate Procedure : Tests for Normality

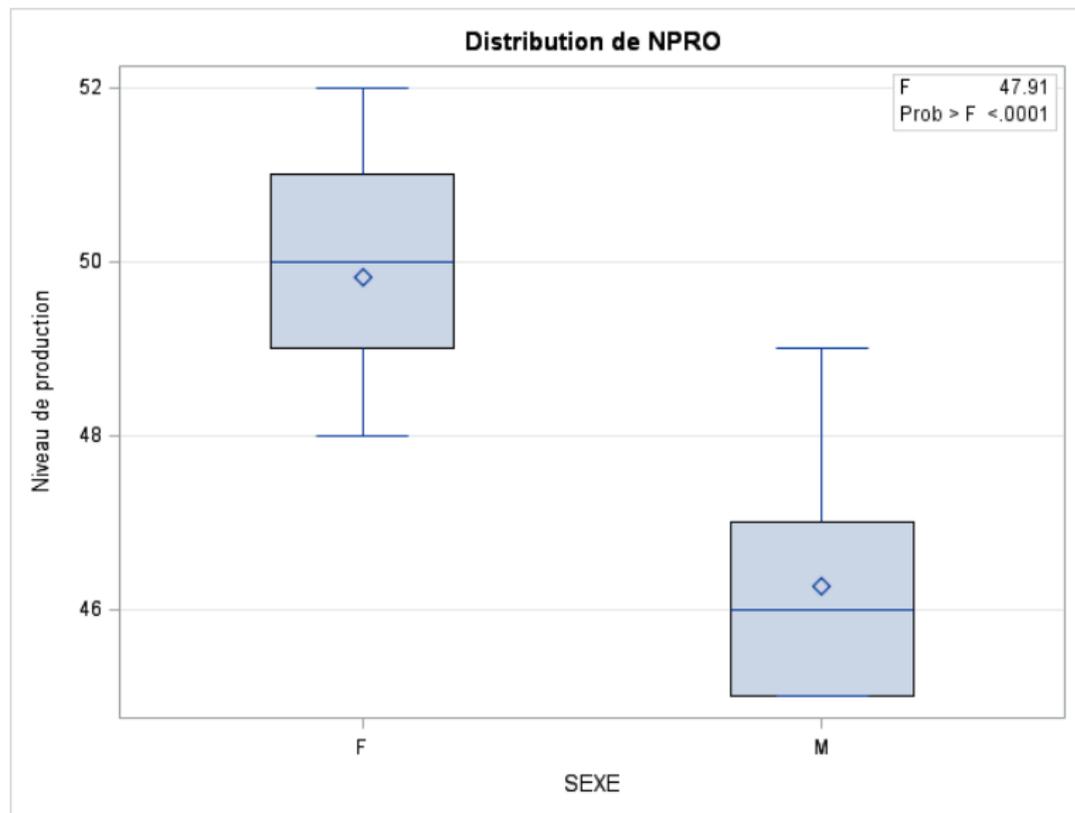
	Test		Statistic		p-Value
Féminin	Shapiro-Wilk	W	0.869325	$Pr < W$	0.0974 ✓
Masculin	Shapiro-Wilk	W	0.851103	$Pr < W$	0.0736 ✓

2 Hypothèse d'homoscédasticité - Egalité des variances :

The ANOVA Procedure : Levene's Test for Homogeneity of Production Variance

Source	DF	Sum of Squares	Mean Square	F-Value	$Pr > F$
Sexe	1	0.1027	0.1027	0.03	0.8546 ✓
Error	25	74.8838	2.9954		

Résultats SAS - BoxPlot



Principaux résultats - SAS

- Effet du **Statut** et du **Sexe** de l'employé sur la production.

1 Significativité du modèle dans son ensemble.

Source	Sum of squares	DF	Mean square	F-value	Pr > F
Model	99.1121693	5	19.8224339	16.53	<.0001 ✓
Error	25.1841270	21	1.1992441		
Total	124.2962963	26			

2 Significativité des facteurs et de leur interaction.

Source	Type III SS	DF	Mean square	F-value	Pr > F
SEXE	36.65908991	1	36.65908991	30.57	<.0001 ✓
STATUT	16.83835896	2	8.41917948	7.02	0.0046 ✓
SEXE*STATUT	2.70344328	2	1.35172164	1.13	0.3428

Trois test de comparaison de moyennes

- Test 1 : Effet du **Statut - Niveau de responsabilité** de l'employé sur la production.

1 Hypothèses statistiques $\begin{cases} H_0 : \alpha_J = \alpha_I = \alpha_S \\ H_1 : \text{les moyennes ne sont pas toutes égales} \end{cases}$

2 Seuil de signification : $\alpha = 5\%$

3 Conditions d'application du test : **Le niveau de production doit être normalement distribué selon chaque niveau de responsabilité (facteur) - Egalité des variances $\sigma_J^2 = \sigma_I^2 = \sigma_S^2$.**

4 Statistique de test "Unilatéral - risque à droite" :
La fonction discriminante de Fisher sous H_0 :

$$F = \frac{S_{entA}^2}{S_{int}^2} \rightarrow F_{p-1=2; n-pq=21} \text{ d.d.l.}$$

5 Calcul de la statistique de test sous l'hypothèse nulle $H_0 : \alpha_J = \alpha_I = \alpha_S$

$$f_0 = \frac{8.419}{1.199} = 7.02$$

6 Règle de décision : $f_{5\%} = 3.47$ cf. Table de la loi de Fisher $F_{(2;21)}$.

7 Décision et conclusion : f_0 appartient à la zone de rejet de H_0 ($f_0 = 7.02 > f_{5\%} = 3.47$), on peut donc conclure, avec risque d'erreur $\alpha = 5\%$, qu'il y a une différence significative. Les moyennes ne sont pas significativement égales. Il y a des disparités entre les niveaux moyens de production selon le niveau de responsabilité de l'employé dans l'entreprise.

Trois test de comparaison de moyennes

- Test 2 : Effet du **Sexe** de l'employé sur la production.

1 Hypothèses statistiques $\begin{cases} H_0 : \beta_H = \beta_F \\ H_1 : \text{les moyennes sont différentes} \end{cases}$

2 Seuil de signification : $\alpha = 5\%$

3 Conditions d'application du test : **Le niveau de production doit être normalement distribué selon chaque niveau du facteur (Homme / Femme) - Egalité des variances $\sigma_H^2 = \sigma_F^2$.**

4 Statistique de test "Unilatéral - risque à droite" :

La fonction discriminante de Fisher sous H_0 :

$$F = \frac{S_{ext}^2}{S_{int}^2} \rightarrow F_{q-1=1; n-pq=21} \text{ d.d.l.}$$

5 Calcul de la statistique de test sous l'hypothèse nulle $H_0 : \beta_H = \beta_F$

$$f_0 = \frac{36.659}{1.199} = 30.57$$

6 Règle de décision : $f_{5\%} = 4.32$ cf. Table de la loi de Fisher $F_{(1;21)}$.

7 Décision et conclusion : f_0 appartient à la zone de rejet de H_0 ($f_0 = 30.57 > f_{5\%} = 4.32$), on peut donc conclure, avec risque d'erreur $\alpha = 5\%$, qu'il y a une différence significative entre les niveaux moyens de production des hommes et des femmes.

Trois test de comparaison de moyennes

- Test 3 : Interaction - Effet du **Statut - Niveau de responsabilité** et du **sexe** de l'employé sur la production.

1 Hypothèses statistiques $\left\{ \begin{array}{l} H_0 : \gamma_{JH} = \gamma_{IH} = \gamma_{SH} = \dots = \gamma_{SF} \\ H_1 : \text{les moyennes ne sont pas toutes égales} \end{array} \right.$

2 Seuil de signification : $\alpha = 5\%$

3 Conditions d'application du test : Le niveau de production doit être normalement distribué selon chaque niveau de responsabilité (facteur A) et chaque niveau du sexe (facteur B) - Egalité des variances selon les niveaux de chaque facteur.

4 Statistique de test "Unilatéral - risque à droite" :

La fonction discriminante de Fisher sous H_0 :

$$F = \frac{S_{enAB}^2}{S_{int}^2} \rightarrow F_{(p-1)(q-1)=2; n-pq=21} \text{ d.d.l.}$$

5 Calcul de la statistique de test sous l'hypothèse nulle $H_0 : \gamma_{JH} = \gamma_{IH} = \dots = \gamma_{SF}$

$$f_0 = \frac{1.352}{1.199} = 1.13$$

6 Règle de décision : $f_{5\%} = 3.47$ cf. Table de la loi de Fisher $F_{(2;21)}$.

7 Décision et conclusion : f_0 appartient à la zone de non rejet de H_0 ($f_0 = 1.13 < f_{5\%} = 3.47$), on peut donc conclure, avec risque d'erreur $\alpha = 5\%$, qu'il n'y a pas une différence significative ; pas d'effet d'interaction ou pas d'effet conjugué du (statut et sexe) sur les niveaux moyens de production.

Principaux résultats - SPAD

- Effet du **Statut** et du **Sexe** de l'employé sur la production.
- **Significativité des niveaux des facteurs.**

Iden Libellé CRITERE(S)	Coeff.	E.type	T	Proba.	V.test
FEMI - Féminin	1.4106	0.255	5.529	0.000 ✓	4.30
MASC - Masculin	-1.4106	0.255	5.529	0.000 ✓	-4.30
INTE - Intermédiaire	-0.2598	0.355	0.731	0.473	-0.72
JUNI - Junior	-1.0415	0.360	2.896	0.009 ✓	-2.63
SUPE - Supérieur	1.3013	0.367	3.541	0.002 ✓	3.10

Principaux résultats - SPAD

- Effet du **Statut** et du **Sexe** de l'employé sur la production.
- **Significativité de l'interaction des niveaux des facteurs.**

Iden Libellé Interaction(s)	Coeff.	E.Type	T	Proba.	V.test
FEMI - Féminin					
INTE - Intermédiaire	0.4783	0.355	1.347	0.192	1.30
FEMI - Féminin					
JUNI - Junior	-0.0177	0.360	0.049	0.961	-0.05
FEMI - Féminin					
SUPE - Supérieur	-0.4606	0.363	1.268	0.219	-1.23
MASC - Masculin					
INTE - Intermédiaire	-0.4783	0.355	1.347	0.192	-1.30
MASC - Masculin					
JUNI - Junior	0.0177	0.360	0.049	0.961	0.05
MASC - Masculin					
SUPE - Supérieur	0.4606	0.363	1.268	0.219	1.23
Constante	48.1487	0.255	188.722	0.000	12.42

Instructions sous SAS & SPAD

- Le tableau suivant résume la syntaxe des procédures et **instructions SAS** nécessaires pour obtenir les principaux résultats de l'ANOVA - ANCOVA.

Hypothèse	Test	Instruction SAS	Procédure
Normalité	histogramme	nom-var / normal	univariate
Normalité	Shapiro-Wilk	normal	univariate
Normalité	qq-plot	nom-var / normal	univariate
Normalité	pp-plot	nom-var / normal	capability
Homoscédasticité	Bartlett ou Levene	means / hovtest =	anova ou glm
ANOVA 1 facteur			anova ou glm
ANOVA p facteurs ($p \geq 2$)			glm
ANCOVA			glm
Estimation des paramètres		model / solution	glm
Comparaisons		means / t cldiff	anova ou glm
Disparités		means / tukey lines	anova ou glm

- Instructions d'un projet SPAD** pour l'ANOVA - ANCOVA.

Groupe de méthodes	Méthode
Scoring et modélisation	Vareg (Régression, analyse de variance-covariance)