

• TEST ANOVA :

Soient A un facteur à p modalités A_1, \dots, A_p et X un caractère quantitatif, avec $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$, on suppose que $\sigma_1^2 = \dots = \sigma_p^2 = \sigma^2$.
On observe la valeur de X_i pour chaque n_i individus d'un échantillon indépendants.

\mathcal{H}_0 : " $\forall i ; \mu_i = \mu$ " VS \mathcal{H}_1 : " $\exists i ; \mu_i \neq \mu$ "

	sce (SS)	ddl (DF)	cm (MS)	f_{obs} (F)
Total	$\sum_{i=1}^p \sum_{j=1}^{n_i} x_{i,j}^2 - n\bar{x}^2 = sce_F + sce_R$	$n - 1$		
Factoriel	$\sum_{i=1}^p n_i \bar{x}_i^2 - n\bar{x}^2$	$p - 1$	$\frac{sce_F}{ddl_F}$	$\frac{cm_F}{cm_R}$
Résiduel	$\sum_{i=1}^p \sum_{j=1}^{n_i} x_{i,j}^2 - \sum_{i=1}^p n_i \bar{x}_i^2$	$n - p$	$\frac{sce_R}{ddl_R}$	

On calcule f_α tq $\mathbb{P}(F \geq f_\alpha) = \alpha$, où $F \sim \mathcal{F}(p-1, n-p)$. On a $\mathcal{R}_\alpha = [f_\alpha; \infty[$.

Ou si $p - \text{value} = \mathbb{P}(F \geq f_{obs}) \leq \alpha$, alors on rejette \mathcal{H}_0 .

• TEST ANOVA (égalité des variances) : $s_i^2 = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (x_{i,j} - \bar{x}_i)^2 = \frac{1}{n_i-1} (\sum_{j=1}^{n_i} x_{i,j}^2 - n_i \bar{x}_i^2)$ et $s^2 = \frac{1}{n-p} \sum_{i=1}^p (n_i - 1) s_i^2$

\mathcal{H}_0 : " $\forall i ; \sigma_i^2 = \sigma^2$ " VS \mathcal{H}_1 : " $\exists i ; \sigma_i^2 \neq \sigma^2$ "

Tests	Statistique de test	Quantile	Zone rejet \mathcal{R}_α
Bartlett (si $\min(n_1, \dots, n_p) \geq 4$)	$K = \frac{(n-p) \ln(s^2) - \sum_{i=1}^p (n_i - 1) \ln(s_i^2)}{1 + \frac{1}{3(p-1)} (\sum_{i=1}^p \frac{1}{n_i - 1} - \frac{1}{n-p})}$	k_α tq $\mathbb{P}(K \geq k_\alpha) = \alpha$, où $K \sim \chi^2(p-1)$	$[k_\alpha ; +\infty[$
Cochran (si $n_1 = \dots = n_p$)	$C_{obs} = \frac{\max(s_i^2)}{\sum_{i=1}^p s_i^2}$	$c(m, p)$ (voir table)	$[c(m, p) ; +\infty[$

• TEST ANOVA (comparaison 2 moyennes μ_k et μ_ℓ) : $s_R = \sqrt{cm_R}$

\mathcal{H}_0 : " $\mu_k = \mu_\ell$ " VS \mathcal{H}_1 : " $\mu_k \neq \mu_\ell$ "

Test	Statistique de test	Quantile	Zone rejet \mathcal{R}_α
Bonferroni	$T = \frac{\bar{x}_k - \bar{x}_\ell}{s_R \sqrt{\frac{1}{n_k} + \frac{1}{n_\ell}}}$	t_α^{**} tq $\mathbb{P}(T \geq t_\alpha^{**}) = \frac{2\alpha}{p(p-1)}$, où $T \sim \mathcal{Stu}(n-p)$	$]-\infty ; -t_\alpha^{**}] \cup [t_\alpha^{**} ; +\infty[$

LOI GRANDS NOMBRES : Si on répète N fois une expérience avec une proba p d'apparition d'un événement A , la fréquence de cet événement tend vers p quand $N \rightarrow \infty$.

Si $p\text{-val} \geq \alpha \Rightarrow$ non-rejet \mathcal{H}_0

Si $p\text{-val} < \alpha \Rightarrow$ rejet \mathcal{H}_0

Cas unilatéral gauche : $p - \text{val} = \mathbb{P}(X \leq x_{obs})$

Cas unilatéral droit : $p - \text{val} = 1 - \mathbb{P}(X \leq x_{obs})$

Cas bilatéral : $p - \text{val} = 2(1 - \mathbb{P}(X \leq x_{obs}))$

	$\mathcal{U}(1/n)$	$\mathcal{Ber}(p)$	$\mathcal{B}(n, p)$	$\mathcal{P}(\lambda)$	$\mathcal{U}([a, b])$	$\chi^2(n)$	$\mathcal{Stu}(n)$
\mathbb{E}	$(n+1)/2$	p	np	λ	$(b+a)/2$	n	$\begin{cases} 0, & \text{si } n > 1 \\ \text{FI sinon} \end{cases}$
\mathbb{V}	$(n^2-1)/12$	$p(1-p)$	$np(1-p)$	λ	$(b-a)^2/12$	$2n$	$\begin{cases} \frac{n}{n-2}, & \text{si } n > 2 \\ +\infty, & \text{si } 1 < n \leq 2 \\ \text{FI sinon} \end{cases}$

Soit $X \sim \mathcal{B}(n, p)$. Si $n \geq 30$ et $np \leq 10$, alors $X \sim \mathcal{P}(\lambda)$, avec $\lambda = np$.

Soit $X \sim \mathcal{B}(n, p)$. Si $np \geq 10$ et $n(1-p) \geq 10$, alors $X \sim \mathcal{N}(np, \sqrt{np(1-p)})$.

Soit $X \sim \mathcal{P}(\lambda)$. Si $\lambda \geq 10$, alors $X \sim \mathcal{N}(\lambda, \sqrt{\lambda})$.

Variable qualitative \rightarrow nominale = plusieurs modalités non mesurables s'excluant mutuellement / ordinale = degrés d'un état

Variable quantitative \rightarrow résultat d'une mesure ou d'un comptage

$Y = \hat{a}X + \hat{b} \Leftrightarrow \hat{a} = \frac{\text{cov}(x,y)}{s_{\text{ech}}^2(x)} = r \frac{s_{\text{ech}}(y)}{s_{\text{ech}}(x)}$ et $\hat{b} = \bar{y} - \hat{a}\bar{x} / X = \hat{a}Y + \hat{b} \Leftrightarrow \hat{a} = \frac{\text{cov}(x,y)}{s_{\text{ech}}^2(y)} = r \frac{s_{\text{ech}}(x)}{s_{\text{ech}}(y)}$ et $\hat{b} = \bar{y} - \hat{a}\bar{x} \Rightarrow$ droites MCO

Etude des résidus

On appelle **valeur ajustée** de la $i^{\text{ème}}$ observation de la variable Y l'approximation

$$\hat{y}_i = \hat{a}x_i + \hat{b}$$

On appelle **résidu** e_i , l'erreur observée que l'on commet en approchant y_i par \hat{y}_i : $e_i = y_i - \hat{y}_i$

angle + ouvert entre 2 MCOS \Rightarrow liaison moins fortes

Si nouvelle valeur x^* de X , on prédit $\hat{y}^* = \hat{a}x^* + \hat{b}$

Démarche Préviation avec la droite des MCO

- calcul de la droite des MCO
- validation du modèle \Rightarrow étude des résidus et détection des valeurs aberrantes et influentes
- qualité de l'ajustement \Rightarrow décomposition de la variance, coefficient de détermination et test significativité globale
- qualité de prédiction (PRESS) et prédiction

Coefficient de détermination

Part de la variance de y expliquée par la relation $\hat{y} = \hat{a}x + \hat{b}$

$$R^2 = \frac{\text{Var}(\hat{y})}{\text{Var}(y)}$$

Dans le cas d'un ajustement linéaire, on peut montrer que $R^2 = r^2(x, y)$ (où r est le coefficient de corrélation linéaire)

- $R^2 \in [0, 1]$
- Plus R est proche de 1, plus le modèle explique correctement la variabilité de Y .

Validité du modèle

- vérifier la normalité des résidus
- vérifier que les résidus ne contiennent pas d'information structurée
- vérifier que les résidus ne sont pas auto-corrélés entre eux

Croisement de deux variables quantitatives

Représentation graphique (nuage de points)

Coefficient de corrélation

- Calcul de l'indicateur statistique
- Test de nullité du coefficient de corrélation

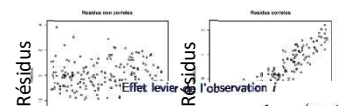
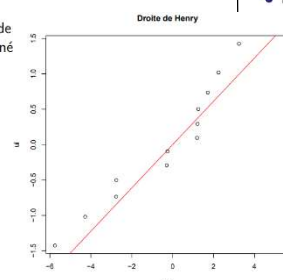
Régression linéaire

- Estimation des coefficients
- Validité du modèle (Etude des résidus et des observations influentes)
- Qualité d'ajustement (R^2 , significativité globale)
- Prédiction

Vérification de l'homoscédasticité des résidus

Les résidus sont **homoscédastiques** si leur répartition est homogène et ne dépend pas des valeurs de la variable explicative (et donc pas non plus des valeurs prédites).

On vérifie que les résidus n'ont pas de structure particulière en traçant un graphe des résidus :



$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

Lever grand \Rightarrow observation atypique.
Remarque : Même si l'hypothèse d'homoscédasticité est vérifiée, les résidus n'ont pas la même variance.
 $E(e_i) = 0$ et $\text{Var}(e_i) = \sigma^2(1 - h_i)$

Equation d'analyse de la variance

$$y_i - \bar{y} = (\bar{y}_i - \bar{y}) + (y_i - \bar{y}_i)$$

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{Somme des carrés totale SCT}} = \underbrace{\sum_{i=1}^n (\bar{y}_i - \bar{y})^2}_{\text{Somme des carrés expliquée SCE}} + \underbrace{\sum_{i=1}^n (y_i - \bar{y}_i)^2}_{\text{Somme des carrés résiduelle SCR}}$$

$$\underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}_{\text{Variance totale}} = \underbrace{\frac{1}{n} \sum_{i=1}^n (\bar{y}_i - \bar{y})^2}_{\text{Variance expliquée}} + \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_i)^2}_{\text{Variance résiduelle}}$$

Droite de Henry :

$P(\epsilon \leq e_i)$	0.077	0.154	0.231	0.308	0.385	0.462	0.538
e_i	-5.75	-4.27	-2.76	-2.75	-0.29	-0.25	1.19
u_i	-1.43	-1.02	-0.74	-0.50	-0.29	-0.10	0.10

Ex : $P(Z \leq u_i) = P(\epsilon \leq e_i) = \frac{7}{13} = 0.538 \Rightarrow u_i = 0.0954$

Propriété : Un estimateur sans biais de la variance de l'erreur du modèle est

$$s_e^2 = \frac{1}{n-2} \sum e_i^2$$

Vérification de la normalité des résidus

- histogramme \Rightarrow la distribution doit être unimodale et symétrique autour de 0.
- tests (Kolmogorov-Smirnov, Shapiro Wilks, ...)
- droite de Henry \Rightarrow compare les quantiles théoriques de normale et la distribution cumulée estimée sur les donnés

Validité du modèle

- vérifier la normalité des résidus
- vérifier que les résidus ne contiennent pas d'information structurée
- vérifier que les résidus ne sont pas auto-corrélés entre eux

Tableau de distribution des fréquences

Modalités	Effectif	Fréquence
x_1	n_1	$f_1 = n_1/n$
\vdots	\vdots	\vdots
x_p	n_p	$f_p = n_p/n$
Total	n	1

On définit l'angle α_i de la modalité x_i par $\alpha_i = 360 \times f_i$.

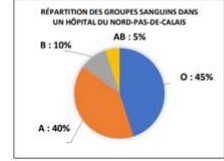


Tableau de contingence

	T	Somme
p	$n_{11} \cdots n_{1j} \cdots n_{1t}$	n_{1*}
	\vdots	\vdots
	$n_{i1} \cdots n_{ij} \cdots n_{it}$	\vdots
	\vdots	\vdots
	$n_{p1} \cdots n_{pj} \cdots n_{pt}$	n_{p*}
Somme	$n_{*1} \cdots n_{*j} \cdots n_{*t}$	n

Variables qualitative nominale
Biais : $B(\widehat{\theta}_n) = \mathbb{E}(\widehat{\theta}_n - \theta)$

Tableau de distribution des effectifs et fréquences cumulés

Valeur des modalités	Effectif n_i	Fréquence f_i	Effectif cumulé N_i	Fréquence cumulée F_i	Fréquence cumulée G_i
x_1	n_1	f_1	n_1	f_1	1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_i	n_i	f_i	$n_1 + \cdots + n_i$	$f_1 + \cdots + f_i$	$f_p + \cdots + f_{p-i}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_p	n_p	f_p	n	1	f_p

Tableau de distribution des effectifs et fréquences cumulés

Classes	Effectif n_i	Fréquence f_i	Densité	Effectif cumulé N_i	Fréquence cumulée F_i	Fréquence cumulée G_i
$[a_1, b_1]$	n_1	f_1	$d_1 = f_1/(b_1 - a_1)$	n_1	f_1	1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$[a_i, b_i]$	n_i	f_i	$d_i = f_i/(b_i - a_i)$	$n_1 + \cdots + n_i$	$f_1 + \cdots + f_i$	$f_p + \cdots + f_{p-i}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$[a_p, b_p]$	n_p	f_p	$d_p = f_p/(b_p - a_p)$	n	1	f_p

Variables qualitative ordinale et quantitative discrète

Variables quantitatives continue

Intervalle confiance moyenne si σ inconnue : $\left[\bar{x} - t_{1-\alpha/2}^{n-1} \frac{\widehat{\sigma}}{\sqrt{n}}; \bar{x} + t_{1-\alpha/2}^{n-1} \frac{\widehat{\sigma}}{\sqrt{n}} \right]$ / si σ connue : $\left[\bar{x} - u_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}; \bar{x} + u_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$

Conditions : Si $n \leq 30$: il faut que $X \sim \mathcal{N}$ / Si $n > 30$: pas de condition sur la loi de X .

Intervalle confiance proportion : on pose $\hat{\pi} = p$, si $n\hat{\pi} \geq 10$ et $n(1 - \hat{\pi}) \geq 10$: $\left[\hat{\pi} - u_{1-\alpha/2} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}; \hat{\pi} + u_{1-\alpha/2} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}} \right]$

2 variables qualitatives \Rightarrow test χ^2 | 2 variables quantitatives \Rightarrow test coeff corrélation | 1 quali & 1 quanti \Rightarrow t-test ou ANOVA

Test comparaison proportion à π_0 connue (1 échantillon) : $U = \sqrt{n} \frac{p - \pi_0}{\sqrt{\pi_0(1-\pi_0)}} \sim \mathcal{N}(0,1)$ / $n_1, n_2 \geq 30$ et $n\pi_0 \geq 5$ et $n(1 - \pi_0) \geq 5$.

Test comparaison proportions (2 échantillons) : $U = \frac{p_1 - p_2}{\sqrt{\frac{p(1-p)}{n_1} + \frac{p(1-p)}{n_2}}} \sim \mathcal{N}(0,1)$, où $p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$ / $n_1, n_2 \geq 30$ et $\min(n_i p; n_i(1 - p)) \geq 5$

Test comparaison moyenne à μ_0 connue (1 échantillon) : $T = \sqrt{n} \frac{\bar{x} - \mu_0}{\widehat{\sigma}} \sim \text{Stu}(n - 1)$ / $n \geq 30$ ou $X \sim \mathcal{N}(\mu, \sigma)$

Test Fisher (comparaison 2 variances observées s_1^2, s_2^2) : $F = \frac{\widehat{\sigma}_1^2}{\widehat{\sigma}_2^2} \sim \mathcal{F}(n_1 - 1; n_2 - 1)$ / $n_1, n_2 \geq 30$ ou $X_{1,2} \sim \mathcal{N}(\mu_{1,2}, \sigma_{1,2})$

t-test comparaison moyennes (2 échantillons indépendants) : $T = \frac{\bar{x}_1 - \bar{x}_2}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim \text{Stu}(n_1 + n_2 - 2)$, où $S = \sqrt{\frac{(n_1 - 1)\widehat{\sigma}_1^2 + (n_2 - 1)\widehat{\sigma}_2^2}{n_1 + n_2 - 2}}$ si $\sigma_1^2 = \sigma_2^2$ / $n_1, n_2 \geq 30$ ou $X_{1,2} \sim \mathcal{N}(\mu_{1,2}, \sigma_{1,2})$
 $T = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\widehat{\sigma}_1^2}{n_1} + \frac{\widehat{\sigma}_2^2}{n_2}}} \sim \text{Stu}(v)$, si $\sigma_1^2 \neq \sigma_2^2$ (Si NON conditions \Rightarrow test Mann & Whitney)

t-test séries appariées :

Sujets	1	2	...	n
X_A	X_1^A	X_2^A	...	X_n^A
X_B	X_1^B	X_2^B	...	X_n^B
$D = X_A - X_B$	$D_1 = X_1^A - X_1^B$	$D_2 = X_2^A - X_2^B$...	$D_n = X_n^A - X_n^B$

$T = \sqrt{n} \frac{\bar{d}}{\widehat{\sigma}_D} \sim \text{Stu}(n - 1)$ / $n \geq 30$ ou $D \sim \mathcal{N}(\mu_D, \sigma_D)$
(Si NON conditions \Rightarrow Test Wilcoxon)

Test χ^2 ajustement (1 quali / 1 échantillon)

Test χ^2 indépendance (2 qualis / 1 échantillon)

Test χ^2 homogénéité (1 quali / $l > 2$ échantillons)

Modalités	1	2	...	c
O_i	O_{1i}	O_{2i}	...	O_{ci}
π_i	π_1	π_2	...	π_c
T_i	$n \times \pi_1$	$n \times \pi_2$...	$n \times \pi_c$

$K = \sum_{i=1}^c \frac{O_i^2}{T_i} - n \sim \chi^2(c - 1)$ / $\forall i; T_i \geq 5$

(Si NON conditions \Rightarrow regrouper des T_i si possible)

		Y				Somme
		1	2	...	c	
	X	1	2	...	c	n_1
	1	O_{11}	O_{12}	...	O_{1c}	n_1
	2	O_{21}	O_{22}	...	O_{2c}	n_2
	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
	l	O_{l1}	O_{l2}	...	O_{lc}	n_l
Somme		m_1	m_2	...	m_c	n

		X				Somme
		1	2	...	c	
	Echantillon	1	2	...	c	n_1
	1	O_{11}	O_{12}	...	O_{1c}	n_1
	2	O_{21}	O_{22}	...	O_{2c}	n_2
	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
	l	O_{l1}	O_{l2}	...	O_{lc}	n_l
Somme		m_1	m_2	...	m_c	n

Test de Mann & Whitney :

- Etape 0 : \mathcal{H}_0 : " $M_{ech1} = M_{ech2}$ " VS \mathcal{H}_1 : " $M_{ech1} \neq M_{ech2}$ "
- Etape 1 : Ranger par ordre croissant l'ensemble des 2 échantillons
- Etape 1bis : Déterminer les rangs. Si doublon, faire rang moyen
- Etape 2 : Calcul de $T_1 = \sum \text{rang}_{ech1}$ et $T_2 = \sum \text{rang}_{ech2}$
- Etape 3 : Calcul de $U_{12} = T_2 - \frac{n_2(n_2+1)}{2}$ et $U_{21} = T_1 - \frac{n_1(n_1+1)}{2}$
- Etape 4 : Statistique de test $U = \min(U_{12}; U_{21})$
- Etape 5 : Recherche de la zone de rejet sur la table

Si $n_1, n_2 \geq 20$, on pose $Z = \frac{U - \mu}{\sigma} \sim \mathcal{N}(0,1)$ avec $\mu = \frac{n_1 n_2}{2}$ et $\sigma = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$

Réécriture de la zone de rejet avec Z.

Test de Wilcoxon

➤ Tableau

Sujets	1	2	...	n
X_A	X_1^A	X_2^A	...	X_n^A
X_B	X_1^B	X_2^B	...	X_n^B
$D = X_A - X_B$	$D_1 = X_1^A - X_1^B$	$D_2 = X_2^A - X_2^B$...	$D_n = X_n^A - X_n^B$

- Etape 0 : \mathcal{H}_0 : " $M_D = 0$ " VS \mathcal{H}_1 : " $M_D \neq 0$ " (ou " $M_D > 0$ " ou " $M_D < 0$ ")
- Etape 1 : Eliminer les $D_i = 0$
- Etape 2 : Ranger par ordre croissant l'ensemble des $|D_i|$
- Etape 2bis : Déterminer les rangs. Si doublon, faire rang moyen.
- Etape 3 : Calcul de $P = \sum \text{rang}_{D_i > 0}$ et $M = \sum \text{rang}_{D_i < 0}$
- Etape 4 : Statistique de test : $T = \min(M; P)$
- Etape 5 : Zone de rejet $\mathcal{R}_\alpha = \{T \leq t\}$

Si $n_1, n_2 \geq 20$, on pose $Z = \frac{T - \mu}{\sigma} \sim \mathcal{N}(0,1)$ avec $\mu = \frac{n(n+1)}{4}$ et $\sigma = \sqrt{\frac{n(n+1)(2n+1)}{24}}$

Réécriture de la zone de rejet avec Z.

Test coefficient corrélation : $T = \frac{r(X,Y)\sqrt{n-2}}{\sqrt{1-r(X,Y)^2}} \sim \text{Stu}(n - 2)$ / $X, Y \sim \mathcal{N}(\mu_{X,Y}, \sigma_{X,Y})$ / $\mathcal{R}_\alpha =]-\infty; -t_{1-\alpha/2}^{n-2} \text{ddl}] \cup [t_{1-\alpha/2}^{n-2} \text{ddl}, \infty[$

Test