

## III / Comparaisons de proportions <sup>chp 3</sup> Test non paramétrique

### 1/ Comparaison d'une proportions à une valeur donnée

soit un  $n$ -éch d'une population présentant un caractère A.

On note par  $X$  le nombre d'individus ayant le caractère A.

$$X = \sum_{i=1}^n \mathbb{1}_A(X_i)$$

si  $X_i$  présente A  $\mathbb{1}_A(X_i) = 1$  sinon  $\mathbb{1}_A(X_i) = 0$ .

On note par  $p$  la proportion d'individus présentant le caractère A.

$\hat{p} = \frac{1}{n} X$  un estimateur sans biais de  $p$ .

La question qu'on se pose : La proportion est elle égale à une  $p_0$  ou pas ?  
Test bilatéral :  $H_0 : p = p_0$  vs  $H_1 : p \neq p_0$

On rejette  $H_0$  si  $U_{\text{obs}} = \frac{|X - np_0|}{\sqrt{np_0(1-p_0)}} > U_{1-\alpha/2}$ .

on  $X \sim B(n, p)$  (somme de bernoulli ind)

$$E(X) = np \quad \text{et} \quad \text{var}(X) = np(1-p)$$

pour  $n$  assez grand  $X \sim N(np, \sqrt{np(1-p)})$  (En general,  $n$  et  $n-p$ )

Donc sous  $H_0 : U_{\text{obs}} \sim N(0,1)$

$$U_{1-\alpha/2} = F^{-1}_{N(0,1)}(1-\alpha/2)$$

[Rq: Dans certains cas, on utilise la stat de décision suivante]

$$U_{\text{obs}} = 2\sqrt{n} \left| \arcsin \sqrt{\frac{x}{n}} - \arcsin \sqrt{p_0} \right| > U_{1-\alpha/2} \quad U_{\text{obs}} \sim N(0,1)$$

[arc sin est calculé pour les angles exprimés en radians (pas en degré)]

• Test unilatérale  $H_0 : p = p_0$  vs  $H_1 : p > p_0$ .

La règle de décision  $U_{\text{obs}} = \frac{X - np_0}{\sqrt{np_0(1-p_0)}} > U_{1-\alpha}$

[pour le second test  $U_{\text{obs}} = 2\sqrt{n} \left| \arcsin \sqrt{\frac{x}{n}} - \arcsin \sqrt{p_0} \right| > U_{1-\alpha}$

## II) Comparaison de deux proportions :

Deux échantillons de tailles  $n_1$  et  $n_2$  respect

$p_1$  La proportion d'inds de l'ech 1 ayant le caractère A

$p_2$  ————— 2. ————— A

On compare les deux proportions  $p_1$  et  $p_2$ .

On peut présenter les résultats de la table de contingence suivant

	Ech 1	Ech 2	Totaux
A	a	b	a+b
$\bar{A}$	$n_1 - a$	$n_2 - b$	$n_1 + n_2 - a - b$
Totaux	$n_1$	$n_2$	$n_1 + n_2$

$$H_0 : p_1 = p_2$$

Méthode asymptotique : les tailles d'éch suffisamment grandes.

$H_0: p_1 = p_2 = p_0$  contre  $H_1: p_1 \neq p_2$ .

La statistique de décision :  $U = \frac{1}{n_1} \sum_{i=1}^{n_1} 1_A(x_i) - \frac{1}{n_2} \sum_{i=1}^{n_2} 1_A(y_i)$

$RC = [ |U| > k_\alpha ]$

On rejette le fait que  $p_1 = p_2$  au risque  $\alpha$  pour  $k_\alpha$  donnée tp  $P(RC) = \alpha$

$U \sim N\left(0, \sqrt{p_0(1-p_0)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}\right)$   $P\left(\frac{U}{\sqrt{p_0(1-p_0)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} > \frac{k_\alpha}{\sqrt{p_0(1-p_0)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}\right) = \alpha/2$

Pour un test unilatéral  $H_0: p_1 = p_2 = p_0$  contre  $H_1: p_1 > p_2$  (ou  $p_1 < p_2$ )

au seuil  $\alpha$ , on rejette  $H_0$  si  $P(U > k_\alpha) = \alpha$ .

$P\left(\frac{U}{\sqrt{p_0(1-p_0)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} > \frac{k_\alpha}{\sqrt{p_0(1-p_0)\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}\right) = \alpha$

Exp

Sur 96 pièces provenant d'un fournisseur A, il y a 12 qui sont défectueuses et sur 55 pièces provenant d'un autre fournisseur B, il y a 15 défectueuses. Peut-on dire, que la proportion de pièces défectueuses de A est la même que celle de B, au seuil  $\alpha = 5\%$ .

$H_0: p_1 = p_2 = p$  contre  $H_1: p_1 \neq p_2$ .

$\hat{p}_1 = \frac{1}{n} \sum_{i=1}^{96} 1_{\text{def}}(x_i) = \frac{12}{96} = \frac{1}{8}$   $\hat{p}_2 = \frac{15}{55} = \frac{3}{11}$

$U = \hat{p}_1 - \hat{p}_2 \sim N\left(0, \sqrt{p(1-p)\left(\frac{1}{96} + \frac{1}{55}\right)}\right)$

$p$  inconnue on l'estime par  $\hat{p} = \frac{12+15}{96+55} = \frac{27}{151} = 0,18$

$P\left(\frac{|U|}{\sqrt{p(1-p)\left(\frac{1}{96} + \frac{1}{55}\right)}} > k_\alpha\right) = 0,05$

$k_\alpha = F_{N(0,1)}^{-1}\left(1 - \frac{0,05}{2}\right) = 1,96$

$\frac{U}{\sqrt{0,18(1-0,18)\left(\frac{1}{96} + \frac{1}{55}\right)}} \text{ obs} = \frac{\frac{1}{8} - \frac{3}{11}}{\sqrt{\dots}} = 2,28$

rejet  $H_0$ .

270/15  
150/0/1  
42

0,05/2  
10 top

## Test de Cramer von Mises

Soit  $\{x_1, \dots, x_n\}$  n-éch. de  $X$  de f.r.  $F$

$F^*(x)$  f.r. connue.

$H_0: F(x) = F^*(x) \quad \forall x$  contre  $H_1: \exists x, F(x) \neq F^*(x)$

On suppose que tout les paramètres de  $F^*$  sont connus.  
Indicateur d'écart de test est

$$I = \int_{-\infty}^{+\infty} |F^*(x) - F(x)|^2 dF(x)$$

La distribution de cet indicateur a été tabulée. On démonte que

$$I = \frac{1}{2n} + \sum_{i=1}^n \left[ \frac{x_i - 1}{2n} - F^*(x_i) \right]^2$$

où les valeurs de l'éch. ordonnées  $\uparrow$ .

Au seuil  $\alpha$ , on rejette  $H_0$  si la valeur de  $I$  dépasse  $k_\alpha$ .

Rq: Le test de Cramer von Mises a les m<sup>^</sup>es applications que le test de Kolmogorov. La différence entre ces deux tests réside dans le fait que pour le test de Kolmogorov seul l'écart max entre la distribution empirique et la distribution d'ajustement est en considération alors que le test de Cramer V.M., l'indicateur d'écart prend mieux en compte <sup>l'ensemble</sup> les données en ce sens que la somme des écarts intervient. Le test de Kolmogorov, est donc beaucoup plus sensible à l'existence de pics aberrants d'un échantillon que le test de Cramer V.M. On pense généralement que ce dernier est plus puissant, mais cela n'a pas été démontré théoriquement.

## 2/ Test de $\chi^2$

Il permet de comparer la densité de la loi à l'histogramme construit à partir des obs.

Le pb avec l'histogramme est le choix ts arbitraire des classes.

on suppose néanmoins que  $k$  classes sont choisies.

Le principe de test de  $\chi^2$  est de comparer le pourcentage d'obs observé ds la classe numéro  $i$ , que nous noterons par  $\hat{p}_i$  au pourcentage d'obs contenues ds cette classe, que nous noterons par  $p_i$

Soient les  $k$  classes  $C_1 = [a_1, a_2[ ; \dots ; C_k = [a_{k-1}, a_k[$

$$p_1 = P(C_1) = F_0(a_2) - F_0(a_1); \dots ; p_k = P(C_k) = F_0(a_{k+1}) - F_0(a_k)$$

Test :  $H_0: F(x) = F_0(x) \forall x$  contre  $H_1: \exists x F(x) \neq F_0(x)$

$$\text{La stat de décision } \chi_{obs}^2 = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}$$

La distribution de  $\chi_{obs}^2$  est  $\chi_{k-1}^2$  ( $\sum N_i = n$ )

$N_i$  = effectif de la classe  $C_i$

$$N_i \sim N(np_i; np_i(1-p_i)) \quad \left( \frac{N_i - np_i}{\sqrt{np_i(1-p_i)}} \right) \sim N(0, 1)$$

On admet aussi que si  $p_i$  est petit  $np_i(1-p_i) \approx np_i$  par suite

$$\sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i} \sim \chi_{k-1}^2 \quad (n \text{ assez grand})$$

le test de  $\chi^2$  est applicable pour  $np_i > 10$  généralement si non on regroupe les classes.

Exp

Soit un dé à 6 faces dont on veut vérifier la régularité.

$X$  le nbre apparu après lancement.

Tester  $X \sim \mathcal{U}_{\{1, \dots, 6\}}$  contre  $F_X(x) \neq \frac{x}{6}$  pour  $x$

c'  $H_0: p_i = 1/6 \forall i$  contre  $H_1: \exists i \text{ tq } p_i \neq 1/6$ .

Pour cela on lance le dé 100 fois on obtient

$i$	1	2	3	4	5	6
$n_i$	17	11	22	15	11	31

Ces données sont-elles compatibles avec la distr sous  $H_0$ ?

$C_i = \text{Face } i^\circ$

$$p_{\text{ref}} \sum_{i=1}^6 \frac{(n_i - n p_i)^2}{n p_i} = 24,56$$

Au seuil  $\alpha = 0,05$ , on trouve  $\chi_{5}^2$   $k = 11,1$

$\chi_{\text{obs}}^2 > \chi_{5}^2$  on rejette  $H_0$ .

Re: si sous  $H_0$  la forme de la distribution est connue à  $r$  paramètres inconnus. On estime les paramètres.  $\chi_{\text{obs}}^2$  la statistique de décision suit ds ce cas  $\chi_{k-1-r}^2$ .

### III / Tests d'indépendance de deux caractères.

#### 1) Test de $\chi^2$ et test de Mann.

Soit une population sur la quelle on observe de caractères  $X$  et  $Y$ . On s'intéresse et propose de tester l'indépendance des deux caractères. On suppose que les individus de la population peuvent être regroupés en  $k$  classes selon  $X$  et  $s$  classes selon  $Y$ .

Re: si  $X$  ou  $Y$  sont qualitatifs, les classes sont les modalités. On note par  $N_i$  l'effectif de classe  $C_i$  sur  $X$   $f_i$  fréquence.  $N_j$  sur  $Y$   $f_j$  fréquence.

Sous  $H_0$ :  $X \perp Y \Leftrightarrow f_{ij} = f_i \times f_j$ .  
La statistique de décision  $d^2 = \sum_{i,j} \frac{(f_{ij} - f_i \times f_j)^2}{n \cdot f_i \times f_j} \sim \chi_{(k-1)(s-1)}^2$

On rejette  $H_0$  au seuil  $\alpha$ , si  $d > d_{\alpha}$ .

Exp: On interroge 100 personnes de catégories professionnelles  $X_1$  et  $X_2$ , et lisant le journal  $Y_1$  ou  $Y_2$ . On obtient

$Y \backslash X$	$X_1$	$X_2$	$n_{i\cdot}$
$Y_1$	27	9	36
$Y_2$	12	52	64
$n_{\cdot j}$	39	61	100

A-t-on indépendance entre  $X$  et  $Y$  au seuil 5%.

Indep de X et Y ( $\Leftrightarrow \rho = 0$ )

$$r = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad \text{comparer } r \text{ à } 0$$

Proposition

Si X et Y sont indep<sup>te</sup> et si R est le coefficient de corrélation empirique on a

$$\frac{R}{\sqrt{1-R^2}} \sqrt{n-2} \sim T_{n-2}$$

Test  $X \perp Y$  RC =  $[|R| > k]$ .

Rq: on donne

$$Z = \frac{1}{2} \log\left(\frac{1+r}{1-r}\right) \sim N\left(\frac{1}{2} \log\left(\frac{1+\rho}{1-\rho}\right); \frac{1}{n-3}\right)$$

Rq

si X et Y ne sont pas gaussiennes, le test précédent reste valable mais il test la corrélation linéaire non pas l'indépendance.

$$\frac{N_{ij}}{n} \sim N\left(p_{ij}, \sqrt{\frac{p_{ij}}{n}}\right) = N\left(p_{ij}, \sqrt{\frac{p_{ij}}{n}}\right)$$

$$\left(\frac{\frac{N_{ij}}{n} - p_{ij}}{\sqrt{\frac{p_{ij}}{n}}}\right)^2 \sim \chi_1^2$$