

## Chap 4 / Tests non paramétriques de comparaison d'échantillons.

### I | Test de Smirnov (à 2 échantillons)

Deux échantillons proviennent-ils de la même population ?

Exemples :

- 1/ Double correction : Deux correcteurs d'un examen, corrigent-ils de la même manière ?
- 2/ Efficacité de pièces : Les pièces produites par deux machines différentes, sont-elles de même consistance ?
- 3) Deux champs d'oliviers produisent-ils la même quantité d'olives.

Le test consiste à vérifier si deux échantillons sont issus de la même population, c. à d. si ils suivent la même loi de fonction de répartition  $F$  inconnue.

Soient deux échantillons  $X_1, \dots, X_{n_1}$  et  $Y_1, \dots, Y_{n_2}$  issus de deux v. a.  $X$  et  $Y$  respectivement de fonctions de répartition  $F_1$  et  $F_2$  inconnues.

Notons par  $F_{n_1}^*$  et  $F_{n_2}^*$  les fonctions de répartition empiriques des deux échantillons respectivement.

Test :  $H_0: "F_1 = F_2"$  vs  $H_1: "F_1 \neq F_2"$

ou  $H_0: "\forall x F_1(x) = F_2(x)"$  vs  $H_1: "\exists x F_1(x) \neq F_2(x)"$

La statistique de ~~Smirnov~~ Smirnov définie par Smirnov est donnée par :

$$T = \sup_x p \left| F_{n_1}^*(x) - F_{n_2}^*(x) \right| \sqrt{\frac{n_1 \cdot n_2}{n_1 + n_2}}$$

Pour  $n_1$  et  $n_2$  assez grands,  $T$  suit asymptotiquement une loi tabulée. (voir la table de Smirnov à 2 échantillons)

La région critique du test est donnée par :

$$RC = [T > k_\alpha]$$

où  $k_\alpha$  est fixé.

## II / Test de Wilcoxon.

Les deux échantillons  $X_1, \dots, X_{n_1}$  et  $Y_1, \dots, Y_{n_2}$  proviennent-ils de la même population (ont-ils la même loi)

1/ Considérons la statistique

$$U = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \mathbb{1}_{(y_j < x_i)}$$

(Somme des  $y_j$  avant chaque  $x_i$ )

$U$  varie de 0 à  $n_1 n_2$ .

test :  $H_0: F_1 = F_2$  vs  $H_1: F_1 \neq F_2$

Sous  $H_0$ , ( $F_1 = F_2$ ),  $U$  est uniformément répartie sur  $\{0, \dots, n_1 n_2\}$

On a  $E(U) = \frac{n_1 n_2}{2}$ . En effet,

$$E(U) = \sum_{k=0}^{n_1 n_2} k \frac{1}{n_1 n_2 + 1} = \frac{1}{n_1 n_2 + 1} \sum_{k=0}^{n_1 n_2} k = \frac{1}{n_1 n_2 + 1} \left( \frac{n_1 n_2 (n_1 n_2 + 1)}{2} \right)$$

(suite arithmétique de raison 1)

La région critique du test est donnée par :

$$RC = \left[ \left| U - \frac{n_1 n_2}{2} \right| > k_\alpha \right] \text{ à un seuil donné.}$$

2/ Autre formulation du test de Wilcoxon

Considérons la somme des rangs des  $X_i$  dans l'échantillon global notée  $W_X$ , on a

$$\min W_X = \frac{n_1(n_1+1)}{2} \text{ et } \max W_X = \frac{n_1(n_1+1)}{2} + n_1 n_2.$$

Exp:

$X$  : 0 0,5 1 1,5 2

$Y$  : 0,1 0,6 1,1 3

L'échantillon global ordonné : 0 0,1 0,5 0,6 1 1,1 1,5 2 3

$$W_X = 1 + 3 + 5 + 7 + 8 = 24$$

$$W_Y = 2 + 4 + 6 + 9 = 21$$

Rq: Il est facile de voir que  $W_x = U \rightarrow \frac{n_1(n_2+1)}{2}$

Sous  $H_0$ ,  $W_x$  est uniformément répartie sur  $\left\{ \frac{n_1(n_2+1)}{2}, \dots, \frac{n_2(n_1+1)}{2} \right\}$

$$E(W_x) = \frac{n_1 n_2}{2} + \frac{n_1(n_2+1)}{2}, \quad \text{Var}(W_x) = \frac{n_1 n_2 (n_1 + n_2 + 1)}{12}$$

La région critique du test de Wilcoxon est donnée par :

$$RC = [ |W_x| > k_\alpha ] \quad \alpha \text{ seuil donné}$$

La loi  $W_x$  est tabulée (voir la table de Wilcoxon)

Rq:

Si  $n_1$  et  $n_2$  sont supérieurs à 18, alors la région critique :

$$RC = \left[ |W_x - \frac{n_1(n_1+n_2+1)}{2}| > U_{\alpha/2} \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}} \right]$$

ou  $U_{\alpha/2} = \Phi_{N(0,1)}^{-1}(\alpha/2)$  (approximation par la loi normale  $N(0,1)$ )

### III / Comparaison de plusieurs échantillons issus d'une v.a. qualitative

Exemple

Considérons les 4 sections d'étudiants de la 1<sup>ère</sup> année MI

La variable qualitative est la mention obtenue au BAC

les résultats sont présentés dans le tableau suivant.

	S.M	AB	B	TB
Section A	$n_{11}$	$n_{12}$	$n_{13}$	$n_{14}$
Section B	$n_{21}$	$n_{22}$	$n_{23}$	$n_{24}$
Section C	$n_{31}$	$n_{32}$	$n_{33}$	$n_{34}$
Section D	$n_{41}$	$n_{42}$	$n_{43}$	$n_{44}$

$i = 1, \dots, 4$      $j = 1, \dots, 4$

$n_{ij}$  est le n<sup>bre</sup> d'étudiants de la section  $i$  ayant la mention  $j$

Les 4 échantillons proviennent-ils de la même population ?

Soient  $k$  échantillons  $éch_1, \dots, éch_k$ , et  $X$  une v.a. observée sur les  $k$  échantillons.

On note par  $Mod_1, \dots, Mod_r$  les  $r$  modalités de  $X$

# Le tableau des résultats

Mod \ Ech	Mod 1	Mod 2	...	Mod r	Somme des lignes
ech 1	$n_{11}$	$n_{12}$	...	$n_{1r}$	$n_{1\cdot}$
ech 2					
...					
ech k	$n_{k1}$	$n_{k2}$	...	$n_{kr}$	$n_{k\cdot}$
Somme ligne	$n_{\cdot 1}$	$n_{\cdot 2}$	...	$n_{\cdot r}$	$n$

## Notations

$n_{i\cdot} = \sum_{j=1}^r n_{ij}$  est la taille du  $i^{\text{ème}}$  échi

$n_{\cdot j} = \sum_{i=1}^k n_{ij}$  est la taille effectif de la  $j^{\text{ème}}$  Mod

$n = \sum_{i=1}^k n_{i\cdot} = \sum_{j=1}^r n_{\cdot j} = \sum_{i=1}^k \sum_{j=1}^r n_{ij}$  = taille globale de k échantillons

Test:  $H_0$ : "les k échantillons ont la  $m$  loi."

vs  $H_1$ : "les k échantillons sont significativement différents."

1<sup>ère</sup> Cas: On suppose que l'on connaît pour  $H_0$ , les proportions  $p_1, \dots, p_r$  de posséder les modalités Mod 1, ..., Mod r.

La statistique de décision est donnée par:

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^r \left( \frac{n_{ij} - n_{i\cdot} p_j}{n_{i\cdot} p_j} \right)^2$$

$\chi^2$  suit la distribution de  $\chi^2_{kr-k}$  sous  $H_0$

La région critique  $RC = [\chi^2 > k\alpha]$   $\alpha$  seuil fixé

2<sup>ème</sup> Cas: Les proportions  $p_1, \dots, p_r$  sont inconnues.

Souvent, en pratique  $p_1, \dots, p_r$  sont inconnues, alors on estime  $p_j$  pour  $j=1, \dots, r$  par:  $\hat{p}_j = \frac{n_{\cdot j}}{n}$

on obtient ainsi  $(r-1)$  estimateurs indépendants  $\sum_{j=1}^r p_j = 1$

La statistique de décision devient :

$$d_0^2 = \sum_{i=1}^k \sum_{j=1}^r \left( \frac{n_{i \cdot} \hat{p}_j - n_{ij}}{n_{i \cdot} \hat{p}_j} \right)^2$$
$$= \sum_{i=1}^k \sum_{j=1}^r \left( \frac{n_{i \cdot} n_{\cdot j} - n_{ij}}{n_{i \cdot} n_{\cdot j}} \right)^2$$

Et  $d_0^2 \sim \chi_{k(r-k)-(r-1)}^2 = \chi_{(r-1)(k-1)}^2$

Rq. Si la variable  $X$  est quantitative, les modalités sont remplacées par des classes.