

## EXAMEN D'ANALYSE STATISTIQUE

## DEUXIEME SESSION

Tous documents interdits - Calculatrice non programmable autorisée

## Questions de cours

15

1) Voici plusieurs variables. Indiquez, pour chacune d'entre elles, la notation utilisée en cours et s'il s'agit d'une variable aléatoire ou pas. / 1

- Ecart-type empirique
- Fréquence empirique

2) Expliquez la logique d'un test statistique de conformité à une valeur de type / 4

$$\begin{cases} H_0: \theta = 0 \\ H_1: \theta > 0 \end{cases}$$

Plus précisément à quoi correspondent les notions d'erreurs de type I, II, la puissance du test et expliquez comment la région critique est calculée.

## Exercice 1 : Construction d'un estimateur par la méthode du maximum de vraisemblance

16

On dispose d'un échantillon simple au hasard  $(X_1, X_2, \dots, X_n)$  issu du Processus Générateur de Données suivant :  $X_i \rightarrow N(\theta; 1)$  pour  $i=1, \dots, n$ .

1) Indiquez la fonction de vraisemblance correspondant à ce cas. / 1

2) Après avoir résolu les conditions nécessaires d'optimisation, donnez l'expression de l'estimateur de l'espérance  $\theta$  par la méthode du maximum de vraisemblance. / 2

3) Résolvez les conditions suffisantes. / 1

4) Est-ce un estimateur sans biais ? / 1

5) L'estimateur est-il efficace ? (démontrez). / 1

## Indications :

La fonction de densité de la loi normale d'espérance  $\mu$  et d'écart-type  $\sigma$  s'écrit :

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$$

La quantité d'information de Fisher se calcule comme  $I_n(\theta) = -E\left[\frac{d^2 \ln L}{d\theta^2}\right]$ , la borne de Rao-Cramer correspondant à l'inverse de cette quantité.

## Exercice 2 : Choix de marques

112

Afin d'optimiser son positionnement stratégique à New York, une grande marque de vêtements (Ralf) réalise une étude comparative. Elle réalise un tirage aléatoire simple d'une population de référence constituée de plusieurs milliers de clients New Yorkais potentiels et les interroge sur leurs préférences en termes de marques. Les effectifs enquêtés se répartissent dans le tableau de contingence ci-dessous.

Préférence	Marque Ralf	Marque Lacorte	Autres Marques
Catégorie Socio – Prof (CSP)			
Etudiants	352	300	243
Autres Catégories	280	310	123

Une campagne publicitaire de la marque Ralf cible spécifiquement les étudiants. Ainsi, la marque Ralf souhaiterait savoir si elle parvient à attirer davantage la préférence des jeunes que des personnes issues d'autres catégories socio-professionnelles. Peut-on considérer que la marque Ralf a un positionnement privilégié auprès des étudiants ? En d'autres termes est-ce que la proportion de personnes préférant Ralf au sein de la population étudiante est différente de la proportion de personnes préférant Ralf au sein des populations ayant une autre CSP ?

Vous répondrez à cette question en construisant un intervalle de confiance représentant la différence entre la proportion d'étudiants préférant Ralf et la proportion d'autres catégories préférant Ralf.

Pour rédiger votre réponse respectez les étapes et notations suivantes :

Etape 0 : Déterminez les 2 populations mères de référence ici. Indiquez la taille des échantillons respectifs *pour ce problème spécifique*. Indiquez la proportion d'étudiants préférant la marque Ralf (cette fréquence empirique sera notée  $F_{nA}$ ) et la proportion de personnes d'autres catégories préférant la marque Ralf (cette fréquence empirique sera notée  $F_{nB}$ )

Etape 1 : Modélisez. Définissez les variables aléatoires  $X^A$  et  $X^B$  indiquez leur loi de probabilité.

Etape 2 : En supposant les échantillons indépendants, donnez un estimateur adéquat de la différence entre la proportion d'étudiants préférant Ralf et la proportion d'autres catégories préférant Ralf et indiquez sa loi de probabilité (preuve *complète*)

Etape 3 : Après avoir centré et réduit, construisez l'intervalle de confiance de niveau 99%.

Etape 4 : Faire une phrase d'interprétation.

2) Avec un niveau de confiance de 90%, le résultat serait-il le même ?

3) Lorsque le niveau de confiance de l'intervalle baisse, que constatez-vous et pourquoi ?

4) Discutez de la validité de l'hypothèse d'indépendance dans ce cas-là.



# Correction Session 2 (exo Fi)

## Questions de Cours

- ① - Écart type empirique :  $S_n$ , aléatoire  
- Fréquence empirique :  $F_n$ , aléatoire
- ② VOIR COURS

## Exercice 1

①  $L(x_1, \dots, x_n; \mu) = \prod_{i=1}^n f(x_i, \mu)$  si l'échantillon est i.i.d.

où  $f(x_i; \mu) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}\right)$

②  $\hat{\mu}$  est solution de  $\text{Max}_{\mu} L(x_1, \dots, x_n; \mu) \Leftrightarrow \text{Max}_{\mu} \ln L(x_1, \dots, x_n; \mu)$

où  $\ln L(\cdot)$  est la log-vraisemblance notée  $l(x_1, \dots, x_n; \mu)$

$$l(x_1, \dots, x_n; \mu) = \sum_{i=1}^n \ln f(x_i; \mu) = \sum_{i=1}^n \ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}$$

$$l(x_1, \dots, x_n; \mu) = -n \ln \sigma - \frac{n}{2} \ln 2\pi - \frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}$$

(CN)  $\frac{dl(x_1, \dots, x_n; \hat{\mu})}{d\mu} = 0 \Leftrightarrow -\frac{1}{2} \sum_{i=1}^n 2(x_i - \mu) \times (-1) = 0$

$\Leftrightarrow \sum (x_i - \hat{\mu}) = 0 \Leftrightarrow \sum x_i - n\hat{\mu} = 0 \Leftrightarrow$

$$\hat{\mu} = \frac{\sum x_i}{n}$$

③ (CS)  $\frac{d^2 l(x_1, \dots, x_n; \hat{\mu})}{d\mu^2} < 0$   $\frac{d^2 l(\cdot)}{d\mu^2} = -n < 0$  o.t.s

## Exercice 1 (suite)

④  $\hat{\mu}$  est sans biais si  $E(\hat{\mu}) = \mu$

$$E(\hat{\mu}) = E\left(\frac{1}{m} \sum_{i=1}^m X_i\right) = \frac{1}{m} \sum_{i=1}^m E(X_i) \quad \text{car } E(\cdot) \text{ est un opérateur linéaire.}$$

$$E(\hat{\mu}) = \frac{1}{m} \sum_{i=1}^m \mu \quad \text{car } E(X_i) = \mu$$

$$E(\hat{\mu}) = \frac{1}{m} m \times \mu = \mu \quad \underline{\hat{\mu} \text{ est sans biais.}}$$

⑤  $\hat{\mu}$  est efficace si sa variance atteint la borne de Rao-Cramer

$$\bullet V(\hat{\mu}) = V\left(\frac{1}{m} \sum_{i=1}^m X_i\right) = \frac{1}{m^2} \sum_{i=1}^m V(X_i) \quad \text{car les } X_i \text{ sont indépendants par hypothèse.}$$

$$\text{or } V(X_i) = 1$$

$$\text{D'où } \underline{V(\hat{\mu})} = \frac{1}{m^2} \sum 1 = \frac{m \times 1}{m^2} = \boxed{\frac{1}{m}}$$

$$\bullet \text{Borne de Rao-Cramer} = \frac{1}{I_m(\mu)}$$

$$I_m(\mu) = -E\left(\frac{d^2 \ln \mathcal{L}}{d\mu^2}\right) = -E\left(\frac{d^2 \mathcal{L}}{d\mu^2}\right) = -E(-m) = m$$

Borne =  $\frac{1}{m} = V(\hat{\mu})$  / Donc  $\hat{\mu}$  est efficace sa variance est la plus petite possible puisqu'elle atteint la borne de Rao-Cramer.

## Exercice 2

① Étape 0:

On a 2 populations mixtes

- Pop. A : ensemble des étudiants à New York
- Pop. B : ensemble des personnes d'autres catégories à New York

• On dispose d'échantillons iid tirés de ces 2 populations mixtes

$$\text{de tailles } \underline{m_A} = 352 + 300 + 243 = \underline{895}$$

$$\underline{m_B} = 280 + 310 + 123 = \underline{713}$$

• on dispose de statistiques d'échantillon : proportion d'étudiants préférant la marque Ralf :  $\underline{F_{m_A}} = \frac{352}{895} \approx 0,3933 \approx \underline{39,33\%}$

et proportion d'autres personnes préférant la marque Ralf ~~39,27%~~

$$\underline{F_{m_B}} = \frac{280}{713} \approx 0,3927 \approx \underline{39,27\%}$$

Étape 1:

Sont les variables aléatoires  $X_i^A$  et  $X_i^B$  représentant respectivement le fait de préférer la marque Ralf pour un étudiant (respectivement une personne d'une autre catégorie).

$X_i^A = 1$  si l'individu  $i$  préfère Ralf  
0 sinon

idem pour  $X_i^B$ .

Ces variables aléatoires binaires suivent une loi de Bernoulli

## Exercice 2 (suite)

$X_i^A \sim \mathcal{B}(p_A)$     où  $p_A$  et  $p_B$  représentent la proba qu'un  
 $X_i^B \sim \mathcal{B}(p_B)$     individu de la pop A (et B pour  $p_B$ )  
prépare la marque Ralf.  
ces probas sont inconnus.

on sait que  $E(X_i^A) = p_A$      $V(X_i^A) = p_A(1-p_A)$   
 $E(X_i^B) = p_B$      $V(X_i^B) = p_B(1-p_B)$ .  
(Caractéristiques de la loi de proba - théorie).

étape 2:

les fréquences empiriques sont des estimateurs de la proportion  
dans la population mère (ou de la proba)

ainsi  $F_m^A$  est défini la moyenne empirique  $\bar{X}_m^A = \frac{1}{m_A} \sum_{i=1}^m X_i^A = F_m^A$

estimateur de  $p_A$  noté  $\hat{p}_A = \bar{X}_{m_A}^A = F_m^A$

$$\begin{aligned} \text{en effet } E(F_m^A) &= E(\bar{X}_m^A) = E\left(\frac{1}{m_A} \sum_{i=1}^{m_A} X_i^A\right) = \frac{1}{m_A} \sum_{i=1}^{m_A} \underbrace{E(X_i^A)}_{p_A} \\ &= \frac{1 \times m_A}{m_A} \times p_A = p_A \end{aligned}$$

$F_m^A$  est un estimateur sans biais de  $p_A$ .

Par le théorème central limite nous connaissons la loi  
de proba asymptotique (quand la taille de l'échantillon  
tend vers plus l'infini) de la moyenne empirique  $\bar{X}_m^A$   
et donc de la fréquence empirique  $F_m^A$ .

## Exercice 2 (suite)

Étape 2 (suite):

TCL:  $\bar{X}_{m_A}^A \underset{\infty}{\sim} \mathcal{N}\left(\mu_A, \frac{\sigma_A}{\sqrt{m_A}}\right)$  ici  $\bar{X}_{m_A}^A = \bar{F}_m^A$   
 $\mu_A = p_A$

Donc dans le cas de la fréquence

$$\sigma_A = \sqrt{V(X_i^A)} = \sqrt{p_A(1-p_A)}$$

asymptotique on a:  $\boxed{\bar{F}_m^A \underset{\infty}{\sim} \mathcal{N}\left(p_A; \sqrt{\frac{p_A(1-p_A)}{m_A}}\right)}$

de la même façon pour  $\boxed{\bar{F}_m^B \underset{\infty}{\sim} \mathcal{N}\left(p_B; \sqrt{\frac{p_B(1-p_B)}{m_B}}\right)}$

- La différence de proportion dans la population mère s'écrit  $p_A - p_B$  et s'estime par  $\bar{F}_m^A - \bar{F}_m^B$

en effet  $E(\bar{F}_m^A - \bar{F}_m^B) = p_A - p_B$

$$\text{car } E(\bar{F}_m^A - \bar{F}_m^B) = \underbrace{E(\bar{F}_m^A)}_{p_A} - \underbrace{E(\bar{F}_m^B)}_{p_B}$$

$\bar{F}_m^A - \bar{F}_m^B$  est un estimateur sans biais de  $p_A - p_B$ .

- La loi de proba asymptotique de  $\bar{F}_m^A - \bar{F}_m^B$  est la loi normale car c'est une combi. lin. de lois normales.

J'ai donc  $\boxed{\bar{F}_m^A - \bar{F}_m^B \underset{\infty}{\sim} \mathcal{N}(p_A - p_B, ?)}$

$V(\bar{F}_m^A - \bar{F}_m^B) = V(\bar{F}_m^A) + V(\bar{F}_m^B)$  si les échantillons A et B sont indépendants ce que l'on suppose.



## Exercice 2 (suite)

Étape 2 suite

$$\text{on } V(F_m^A) = \frac{p_A(1-p_A)}{m_A} \quad V(F_m^B) = \frac{p_B(1-p_B)}{m_B}$$

$$\text{donc } V(F_m^A - F_m^B) = \frac{p_A(1-p_A)}{m_A} + \frac{p_B(1-p_B)}{m_B}$$

$$\text{D'où } F_m^A - F_m^B \underset{\infty}{\sim} N(p_A - p_B ; \left| \frac{p_A(1-p_A)}{m_A} + \frac{p_B(1-p_B)}{m_B} \right|)$$

Étape 3

on cherche un intervalle  $[A, B]$  tel que  $P(A < p_A - p_B < B) = 0,99$   
on centre et on réduit la loi de l'estimation de  $p_A - p_B$  :

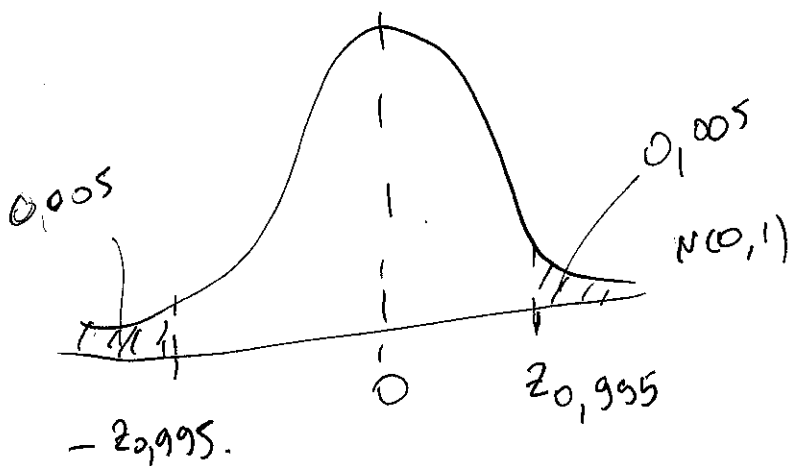
$$Z = \frac{F_m^A - F_m^B - (p_A - p_B)}{\sqrt{\frac{p_A(1-p_A)}{m_A} + \frac{p_B(1-p_B)}{m_B}}} \underset{\infty}{\sim} N(0, 1)$$

$$\text{Donc } P(-z_{0,995} < Z < z_{0,995}) = 0,99$$

Dans la table on voit que

$$z_{0,995} = \Phi^{-1}(0,995) = \underline{2,575}$$

$$P(-2,575 < Z < 2,575) = 0,99$$



## Exercice 2 (suite)

Étape 3 (suite)

les écarts types des fréquences corrigées sont estimés en remplaçant  $p_A(1-p_A)$  par  $F_{m^A}(1-F_{m^A})$ .

$$\text{D'où } Z = \frac{F_{m^A} - F_{m^B} - p_A - p_B}{\sqrt{\frac{F_{m^A}(1-F_{m^A})}{m_A} + \frac{F_{m^B}(1-F_{m^B})}{m_B}}}$$

on déduit l'intervalle :  $P(-2,575 < Z < 2,575) = 0,99$

$$\Leftrightarrow P\left(-2,575 < \frac{F_{m^A} - F_{m^B} - p_A - p_B}{\sqrt{\frac{F_{m^A}(1-F_{m^A})}{m_A} + \frac{F_{m^B}(1-F_{m^B})}{m_B}}} < 2,575\right) = 0,99$$

$$\sqrt{\frac{F_{m^A}(1-F_{m^A})}{m_A} + \frac{F_{m^B}(1-F_{m^B})}{m_B}}$$

$$\Leftrightarrow P\left[F_{m^A} - F_{m^B} - 2,575 \times \sqrt{\frac{F_{m^A}(1-F_{m^A})}{m_A} + \frac{F_{m^B}(1-F_{m^B})}{m_B}} < p_A - p_B\right]$$

$$\left[ \underbrace{F_{m^A} - F_{m^B}}_{\approx 0,0059} + 2,575 \times \sqrt{\frac{F_{m^A}(1-F_{m^A})}{m_A} + \frac{F_{m^B}(1-F_{m^B})}{m_B}} \right] = 0,99$$

$$= \sqrt{0,0026 + 0,0033} = 0,02451$$

$$\Leftrightarrow \mathbb{P}(p_A - p_B) = [-0,0625 ; 0,0637]$$

## Exercice 2 (suite)

Étape 4

Il y a 99% de chances pour que la différence entre la proportion d'étudiants préférant Ralf et la proportion d'autres catégories préférant Ralf soit comprise entre  $-6,25\%$  et  $+6,37\%$ . Cet intervalle comprenant le 0 on ne peut pas dire que Ralf parvienne spécifiquement à attirer la préférence des étudiants.

② Avec un niveau de confiance de 0,90 on utilise

$$z_{0,95} = 1,645$$

$$IC(p_A - p_B)_{90\%} = [0,00059 - 1,645 \times 0,02651 ; 0,00059 + 1,645 \times 0,02651]$$

$$= [-0,0337 ; 0,0409]$$

③ Lorsque le niveau de confiance baisse l'intervalle s'éclaircit car on est moins restrictif.

④ Nous avons supposé l'indépendance entre les échantillons A et B : il n'y avait donc pas de lien entre la proba qu'un étudiant préfère Ralf et la proba que quelqu'un d'une autre catégorie préfère Ralf. Cette hypothèse, plausible, ne semble pas plus forte que de supposer l'indépendance au sein de chaque échantillon.