

EXAMEN DE STATISTIQUES

DEUXIEME SESSION

Tous documents interdits - Calculatrice non programmable autorisée

Questions de cours

14

Formulez des réponses courtes aux 3 questions suivantes.

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.

- Taille de l'échantillon
- Fréquence empirique
- Ecart-type du caractère statistique

11,5

2) Que signifie « le quantile de niveau 0,90 » d'une loi normale centrée réduite (phrase et équation) ? Quelle fonction utiliseriez-vous pour calculer la valeur de la variable notée $z_{1-\alpha}$ dans les formules du cours, en posant $\alpha=10\%$, décrivez le contenu de la formule tel que vous le taperiez dans EXCEL.

11,5

- LOI.NORMAL(x ; espérance ; écart-type ; cumulative)
- LOI.NORMAL.INVERSE (probabilité ; espérance ; écart-type)
- LOI.STUDENT(x ; degrés_liberté ; uni/bilatéral)
- LOI.STUDENT.INVERSE(probabilité ; degrés_liberté)
- LOI.STUDENT.INVERSE.BILATERALE(probabilité ; degrés_liberté)

3) Quelle est la différence entre un intervalle de pari et un intervalle de confiance ?

11

Exercice 1

17

Une étude a montré que la durée de réponse à un test d'anglais contenant 20 questions parmi l'ensemble des lycéens de France est en moyenne de 45 minutes avec un écart-type de 14 minutes. On souhaite calibrer la durée de la prochaine épreuve. En supposant que les durées de réponse sont identiquement et indépendamment distribués selon une loi normale, indiquez (sans démontrer les lois) :

- 1) Quelle est la probabilité qu'un lycéen tiré au hasard réponde en plus de 60 minutes. 12
- 2) Quelle est la probabilité pour que dans un échantillon aléatoire simple de taille $n=30$, la durée moyenne de réponse soit inférieure à 40 minutes 12
- 3) Comment calibrer la durée maximale de l'épreuve pour que 99% des étudiants aient suffisamment de temps pour répondre au test ? 12
- 4) Définir un seuil de durée tel que la durée moyenne de réalisation de l'épreuve soit inférieure à ce seuil avec une probabilité de 99%. 12

Indication :

Pour cette dernière question et uniquement pour celle-ci vous devez définir un intervalle de pari unilatéral.

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_i^A et X_i^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 G)

Questions de Cours

- ① - Taille de l'échantillon : n , non aléatoire
- Fréquence empirique : F_n , aléatoire
- Écart type du caractère statistique : e , non aléatoire

② $z_{0,90}$ ou 90% quantile de niveau 0,90 d'une loi normale centrée réduite est la valeur prise par la variable aléatoire Z telle qu'il y a 90% de chances pour que la variable aléatoire se réalise en dessous de cette valeur.
 $z_{0,90} = \Phi^{-1}(0,90)$ où Φ^{-1} est l'inverse de la loi normale cumulée centrée réduite. Pour le calculer dans Excel il faut taper la formule :

Loi. NORMALE. INVERSE (0,9 ; 0 ; 1)

③ L'intervalle de pari est centré sur la stat. d'échantillon, la moyenne empirique, tandis que l'intervalle de confiance est centré sur un paramètre, en l'occurrence, l'espérance le plus souvent.

Exercice 1

- la population : ensemble des lycéens de France, est connue.
- le caractère statistique : X^s correspond à la durée de réponse en t , il est de moyenne connue $m = 45$ min et d'écart type $e = 14$ min.
- soit la variable aléatoire X^i : durée de réponse d'un lycéen i tiré au hasard. $X^i \sim \mathcal{N}(\mu, \sigma)$ où $\mu = m = 45$
 $\sigma = e = 14$

$$\textcircled{1} P(X_i > 60) = P\left(Z > \frac{60 - \mu}{\sigma}\right) \text{ où } Z = \frac{X_i - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

$$\Leftrightarrow P\left(Z > \frac{60 - 45}{14}\right) = P\left(Z > \frac{15}{14}\right) = P(Z > 1,07) = 1 - P(Z < 1,07)$$

$$= 1 - \Phi(1,07) \quad \text{Dans la table } \Phi(1,07) = 0,8577$$

$$\text{Donc } \underline{P(X_i > 60) = 0,1423}$$

$$\textcircled{2} \text{ On sait que } \bar{X}_m \sim \mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{m}}\right) \text{ où } \mu = 45 \quad \sigma = 14 \text{ et } m = 30$$

$$P(\bar{X}_m < 40) = P\left(\frac{\bar{X}_m - \mu}{\frac{\sigma}{\sqrt{m}}} < \frac{40 - \mu}{\frac{\sigma}{\sqrt{m}}}\right) = P(Z < -1,96) = 1 - \Phi(1,96)$$

$$\text{Table: } \Phi(1,96) = 0,9750 \quad \text{Donc } \underline{P(\bar{X}_m < 40) = 0,0250 = 2,5\%}$$

$$\textcircled{3} \text{ Soit } \eta \text{ la durée telle que } P(X_i < \eta) = 0,99$$

$$\Leftrightarrow P\left(\frac{X_i - \mu}{\sigma} < \frac{\eta - \mu}{\sigma}\right) = P\left(Z < \frac{\eta - \mu}{\sigma}\right) = 0,99$$

$$\frac{\eta - \mu}{\sigma} = \Phi^{-1}(0,99) \Rightarrow \eta = \Phi^{-1}(0,99) \times \sigma + \mu$$

$$\eta = 2,33 \times 14 + 45 = \underline{77,62 \text{ environ}}$$

Exercice 1 (suite)

(4) Soit S la valeur seil telle que $P(\bar{X}_n < S) = 0,99$

$$\Leftrightarrow P\left(\frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} < \frac{S - \mu}{\frac{\sigma}{\sqrt{n}}}\right) = P\left(Z < \frac{S - \mu}{\frac{\sigma}{\sqrt{n}}}\right) = 0,99.$$

$$\frac{S - \mu}{\frac{\sigma}{\sqrt{n}}} = \Phi^{-1}(0,99) \quad \Leftrightarrow \quad S = \Phi^{-1}(0,99) \times \frac{\sigma}{\sqrt{n}} + \mu$$

$$S = 2,33 \times \frac{14}{\sqrt{30}} + 45 = \underline{\underline{50,9 \text{ euros}}}$$

Exercice 2

① Étape 0:

On a 2 populations mixtes

$\left\{ \begin{array}{l} \text{Pop. A : ensemble des étudiants à New York} \\ \text{Pop. B : ensemble des personnes d'autres catégories à New York} \end{array} \right.$

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

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

représentant la marque Ralf : $\underline{F_{m_A}} = \frac{352}{895} \approx 0,3933 \approx \underline{39,33\%}$

et proportion d'autres personnes représentant la marque Ralf ~~713~~

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

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).

$\left\{ \begin{array}{l} X_i^A = 1 \text{ si l'individu } i \text{ préfère Ralf} \\ 0 \text{ sinon} \end{array} \right.$

idem pour X_i^B

Ces variables aléatoires binaires suivent une loi de Bernoulli

Exercice 2 (suite)

$X_i^A \text{ vs } \mathcal{B}(p_A)$ où p_A et p_B représentent la proba qu'un
 $X_i^B \text{ vs } \mathcal{B}(p_B)$ individu de la pop A (et B pour p_B)
pôse 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 $\bar{X}_m^A = \frac{1}{m_A} \sum_{i=1}^m Y_i^A = F_m^A$

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

en effet $E(F_m^A) = E(\bar{X}_m^A) = E\left(\frac{1}{m_A} \sum_{i=1}^{m_A} Y_i^A\right) = \frac{1}{m_A} \sum_{i=1}^{m_A} \underbrace{E(Y_i^A)}_{p_A}$
 $= \frac{1}{m_A} \times m_A \times p_A = p_A$

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 l'infini) de la moyenne empirique \bar{X}_m
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

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

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

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

- La différence de proportion dans la population n'est pas $p_A - p_B$ et n'est estimée que $\bar{F}_m^A - \bar{F}_m^B$

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)

Etape 2 suite

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

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

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

Etape 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 $p_A - p_B$:

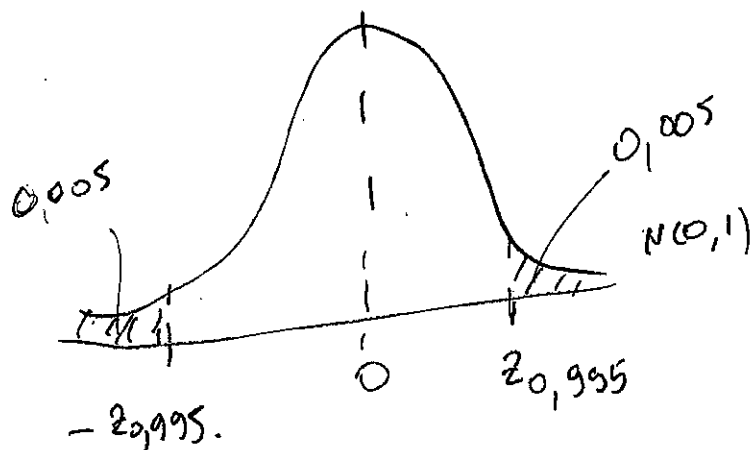
$$Z = \frac{F_m^A - F_m^B - (p_A - p_B)}{\sqrt{\frac{p_A(1-p_A)}{m_A} + \frac{p_B(1-p_B)}{m_B}}} \underset{\infty}{\sim} N(0, 1)$$

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

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 empiriques 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 : $P(-2,575 < z < 2,575) = 0,99$

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

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

$$< \underbrace{F_{m_A} - F_{m_B}}_{\approx 0,00059} + 2,575 \times \left[\frac{F_{m_A}(1-F_{m_A})}{m_A} + \frac{F_{m_B}(1-F_{m_B})}{m_B} \right] = 0,99$$

$$= \sqrt{0,00026 + 0,00033} = 0,02451$$

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

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 parviendrait spécifiquement à attirer la préférence des étudiants.

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

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

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

$$= [-0,0397 ; 0,0409]$$

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 aurait 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, fautive, ne semble pas plus forte que de supposer l'indépendance au sein de chaque échantillon.