

Calculs sur les pourcentages

0 Loi normale

La loi normale $N(m; \hat{\sigma})$, c'est la loi définie par les fonctions de densité :

$$f(x) = \frac{1}{\hat{\sigma} \sqrt{2\pi}} * e^{-\frac{(x-m)^2}{2\hat{\sigma}^2}} \text{ où } m \text{ et } \hat{\sigma} \text{ i } p^{**}$$

Si X suit une variable aléatoire qui suit la loi normale $N(m; \hat{\sigma})$ ou $X \sim N(m; \hat{\sigma})$ alors : $m = E(X)$ et $\hat{\sigma}^2 = V(X)$.

Lorsque l'on lit la table de $\Phi(t)$ de $N(0; 1)$, on a :

$$P(X < t) = F(t) = \Phi(t) \text{ et } \Phi(a) = 1 - \Phi(-a)$$

Lorsque l'on lit le quantile $\tilde{\sigma}_{\tilde{N}}$:

Pour $\tilde{N} = 0,05$, $\tilde{\sigma}_{0,05} = 1,96$ et pour $\tilde{N} = 0,01$, $\tilde{\sigma}_{0,01} = 2,58$

Calcul : $P(X < \tilde{\sigma}_{\tilde{N}}) = 1 - \tilde{N}/2$

Soit P , une population constituée uniquement de deux catégories A et B. Un échantillon E de P est supposé obtenu par tirage au sort simple, aléatoire et non exhaustif (= avec remise).

I Intervalle de confiance

A. Intervalle de confiance et risque

p_0 est le pourcentage (= la fréquence) des individus de la catégorie A dans E .

p est le pourcentage de la population initiale P .

Définition : L'intervalle de confiance au risque \tilde{N} pour p est un intervalle où p a la probabilité $(1 - \tilde{N})$ de se trouver.

Remarques :

– \tilde{N} est la probabilité ou le risque de se tromper en affirmant que p est dans l'intervalle de confiance au risque \tilde{N} pour p .

– $(1 - \tilde{N})$ est dit le coefficient de sécurité ou est aussi dit le seuil de signification.

à On pose $I_{\tilde{N}} = p_0 \pm \tilde{\sigma}_{\tilde{N}} * \sqrt{(p_0 q_0 / n)}$

où $q_0 = 1 - p_0$ et $\tilde{\sigma}_{\tilde{N}}$ est la valeur donnée par la table de la loi normale : $X \sim N(0;1)$ et $P(|X| > \tilde{\sigma}_{\tilde{N}}) = \tilde{N}$.

Définition : \bar{E} est dit un grand échantillon si : $np_0, nq_0, n\tilde{N}_i, n(1 - \tilde{N}_i) > 5 ; i = 1, 2$.

On supposera par la suite que \bar{E} est un grand échantillon si $n > 30$.

Remarques : Plus le risque est faible, plus l'intervalle est large.

à Si d est la précision et ne doit pas dépasser d_0 :

$$d = |p - p_0| = \tilde{O}_{\tilde{N}} * t(p_0q_0 / n)$$

è $d^2 = \tilde{O}_{\tilde{N}}^2 * p_0q_0 / n < d_0^2$ Ainsi $n > \tilde{O}_{\tilde{N}}^2 * p_0q_0 / d_0^2$

B. Test de Conformité

On veut comparer une valeur théorique p à une valeur p_0 observée sur un échantillon de taille n .

On teste l'hypothèse nulle : $H_0 = "p \text{ est la proportion des individus de la population de la catégorie A."}$

à On pose : $\tilde{O}_0 = |p_0 - p| / t(pq / n)$ avec $n > 30$

Test :

q Si $\tilde{O}_{\tilde{N}} \leq \tilde{O}_0$, on refuse H_0 au seuil de signification \tilde{N}

q Si $\tilde{O}_{\tilde{N}} > \tilde{O}_0$, on ne peut refuser H_0 au risque \tilde{N}

C. Test d'homogénéité

On considère ici deux populations d'individus de catégories A et B et on extrait un échantillon de chacune. On observe une proportion p_0 et p'_0 d'individus de la catégorie A sur les deux échantillons. On veut comparer p_0 et p'_0 . Pour cela on formule l'hypothèse $H_0 = "Les deux populations ont la même proportion p d'individus de la catégorie A."$

à On pose $\tilde{O}_0 = |p_0 - p'_0| / \hat{a}_d$ et $\hat{a}_d = t(p_0q_0 / n_0 + p'_0q'_0 / n'_0)$

On effectue ensuite le même test que précédemment.

II Test du Khi-2

Soient $x_1 \dots x_Y$ avec Y des variables aléatoires indépendantes telles que :
 $x_i \sim N(0 ; 1)$ et $X_Y = \sum x_i^2$.

Définition : X_Y suit une loi dite du χ^2 à Y degrés de liberté.

On peut montrer que la densité du X_Y est donnée par :

$$f(x) = 1/2^{Y/2} \Gamma(Y/2) * x^{Y/2-1} e^{-x/2} \text{ où } \Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$$

Propriété : Si $x_1 \dots x_n$ sont n variables aléatoires avec $x_i \sim N(0 ; 1)$ et si $x_1 \dots x_n$ sont liés par k relations, alors : $\sum x_i^2 \sim \chi^2_{n-k}$.

Si $X \sim \chi^2_n$, $Y \sim \chi^2_m$ et X et Y sont indépendantes, alors $(X + Y) \sim \chi^2_{n+m}$.

A. Valeur de Khi-2

Si $\tilde{N} \in [0; 1[$, il existe un unique nombre positif noté $\chi^2_{\tilde{N}, \tilde{Y}}$ tel que :

$$P(\chi^2_{\tilde{N}, \tilde{Y}} \leq \chi^2_{\tilde{Y}}) = \tilde{N}$$

à $\chi^2_{\tilde{N}, \tilde{Y}}$ sera donné par la table pour un \tilde{N} donné et un coefficient de liberté \tilde{Y} défini :
 $\tilde{Y} = (\text{nb de lignes} - 1) * (\text{nb de colonnes} - 1)$

Contexte : On considère une distribution observée sur un échantillon E de taille n_T réparti en k classes. O , dispose :

- d'une série d'effectifs observés : $O_1, O_2 \dots$
- d'une série d'effectifs théoriques : $C_1 = n_T * p_1 \dots C_k = n_T * p_k$

Objectif : Comparer la distribution donnée par E et la distribution théorique (à un seuil de risque \tilde{N}).

Validité : On suppose que E est non exhaustif et aléatoire, le nombre de classes est peu élevé ($k < 20$) et les effectifs $C_i > 5$.

B. Test de conformité

On pose : $\chi^2_0 = \sum (O_i - C_i)^2 / C_i$

H_0 : "Dans la population observée, la distribution des fréquences observées est celle décrite par le modèle théorique."

Propriété : Si $\chi^2_0 \leq \chi^2_{\tilde{N}; k-1}$, on refuse l'hypothèse au seuil de risque \tilde{N} , sinon on ne peut refuser l'hypothèse.

Remarque : Dans le premier cas, la probabilité de se tromper en refusant l'hypothèse n'étant que du risque \tilde{N} , on prend ce risque. Dans le deuxième cas, le test ne fournit aucune exigence contre l'hypothèse.

à Si $k = 2$ et p_0 (respectivement p) le pourcentage observé (respectivement théorique) de l'une des deux classes, le test sur les pourcentage avec \hat{O}_0 est équivalent au test χ^2 .

à Lorsque le degré de liberté est 1, on effectue la correction de Yates en remplaçant :

$$\chi^2_0 = \sum (|O_i - C_i| - \frac{1}{2})^2 / C_i$$

C. Test d'échantillonnage

Le modèle théorique est connu et on veut tester l'hypothèse nulle suivante : H_0
 "L'échantillon E extrait de la population est représentatif."

à On la test comme dans le B

Remarque : On note x_{moy} et $\hat{\sigma}_e$ la moyenne et l'écart-type donnés par \bar{E} si l'espérance mathématique m et l'écart-type s de la distribution théorique ne sont pas connus. On les estime par : $m = x_{moy}$ et $s = \sqrt{(n/n-1) * \hat{\sigma}_e}$.

Lorsque $\bar{Y} > 30$: $(2 X \bar{Y})^{1/2} \times N((2\bar{Y} - 1)^{1/2}; 1)$

D. Test d'homogénéité

On observe un caractère sur une population D . On constitue m échantillons E_1, \dots, E_m de tailles n_i .

On considère un partage en k classes pour chaque $i \in \{1; \dots; m\}$. On note O_{ij} , l'effectif des observations de E_i dans la classe j .

On fait ensuite un tableau de contingence :

Échantillon	Classe j	...	Classe k	Total
E_1	O_{1j}			n_1
...				
E_m			O_{mk}	n_m
Total	t_j		t_k	n_T

$$\chi^2_0 = \sum (O_{ij} - C_{ij})^2 / C_{ij}$$

On teste $H_0 =$ "Les échantillons E_1, \dots, E_m sont extraits d'une même population."
à Pour tester H_0 , on commence par estimer les probabilités de chaque classe en utilisant H_0 elle-même.

On pose $E = \sum E_i$ et $|E| = n_T$.

- q La probabilité pour la classe j est $p_j = t_j / n_T$
- q L'effectif théorique pour la classe j est $C_{ij} = n_i t_j / n_T$
- q $\chi^2_i = \sum (O_{ij} - C_{ij})^2 / C_{ij}$
- q $\chi^2_0 = \sum \chi^2_i$

Propriétés : Si $\chi^2_0 \leq \chi^2_{n;\gamma}$, on refuse H_0 au seuil \bar{N} .
Si $\chi^2_0 > \chi^2_{n;\gamma}$, on ne peut refuser H_0 au seuil \bar{N} .

III Test de Student

Soit $X \sim N(0; 1)$ et $Y \sim \chi^2_{\gamma}$, deux variables aléatoires indépendantes.

A. Définition

La variable aléatoire $t_{\gamma} = X / \sqrt{Y / \gamma}$, elle suit une loi de probabilité dite loi de Student de degré de liberté γ .

Cette loi est définie par la fonction de densité suivante :

$f(x) = A[1 + x^2 / \hat{Y}]^{-(\hat{Y}+1)/2}$, A est une constante positive

B. Propriétés

$E(t_{\hat{Y}}) = 0$ et $\hat{\sigma}(t_{\hat{Y}}) = \sqrt{[\hat{Y} / (\hat{Y}-2)]}$ ainsi que $\hat{Y} \geq 3$.

Remarque : si \hat{Y} tend vers l'infini, $t_{\hat{Y}} \sim N(0; 1)$.

C. Table

Pour $\hat{N} \in]0; 1]$, il existe un nombre positif unique $t_{\hat{N}; \hat{Y}}$ tel que :

$$P(|t_{\hat{Y}}| \leq t_{\hat{N}; \hat{Y}}) = \hat{N}.$$

à Connaissant \hat{N} et \hat{Y} , la table donne $t_{\hat{N}; \hat{Y}}$.

Ex : $t_{0,05;6} = 2,447$ et $t_{0,05;8} = 2,306$

VI Calcul sur les moyennes

On prélève un échantillon non exhaustif de taille n. On note $x_1 \dots x_n$ les observations relatives du caractère étudié.

à On note : $\bar{x}_{moy} = (\sum x_i) / n$ et $\hat{\sigma}_e^2 = (\sum x_i^2) / n - \bar{x}_{moy}^2$.

A. Intervalle de confiance

Définition : L'intervalle de confiance au risque \hat{N} pour la moyenne m est un intervalle où m a la probabilité $1 - \hat{N}$ de se trouver.

à On estime : $m = \bar{x}_{moy}$ et $s = \hat{\sigma}_e \cdot \sqrt{[n / (n - 1)]}$.

è On suppose que $X \sim N(0; 1)$ et $n > 10$.

Propriété : On a l'intervalle au risque \hat{N} pour m :

$$I_{\hat{N}} = \bar{x}_{moy} \pm t_{\hat{N}; n-1} \cdot s / \sqrt{n}$$

Remarques : _ Lorsque $n > 30$, on peut remplacer $t_{\hat{N}; n-1}$ par $\hat{O}_{\hat{N}}$.

_ Le risque de se tromper en affirmant que $m \in I_{\hat{N}}$ est \hat{N} .

_ $I_{\hat{N}}$ est aussi dit l'indice de confiance au seuil de sécurité $1 - \hat{N}$.

B. Comparaison d'une moyenne théorique à une moyenne observée

On s'intéresse à un caractère dans une population et on suppose a priori que la moyenne de la variable aléatoire X correspondant à ce caractère est m et on veut tester cette supposition, c'est-à-dire l'hypothèse H_0 : "m est la moyenne du caractère dans la population".

On extrait un échantillon non exhaustif de taille n. On suppose que $n > 10$ ou que X suit une loi normale.

à On note \bar{x}_{moy} la moyenne observée sur l'échantillon. On pose :

$$t_0 = | \bar{x}_{\text{moy}} - m | / (s / \sqrt{n})$$

Test : _ Si $t_0 \geq t_{\tilde{N};n-1}$, on refuse H_0 au seuil de risque \tilde{N}
 _ Si $t_0 < t_{\tilde{N};n-1}$, on ne peut refuser H_0 au seuil de risque \tilde{N}

C. Composition de deux moyennes

I Test paramétrique :

On suppose que l'on a deux échantillons prélevés sur deux populations distribuées, chacune suivant une loi normale de moyenne m_i et de variance $\hat{\sigma}_i^2$.

On suppose que $\hat{\sigma}_1 = \hat{\sigma}_2 = s$.

à On veut tester l'hypothèse H_0 : "m1 et m2 sont égales".

Remarque : Si H_0 est rejetée, on conclut que les moyennes sont différentes au seuil prédéterminé sinon, on conclut qu'aucune différence significative n'est montrée entre les deux moyennes.

On estime $m_i = \bar{x}_{i\text{moy}}$ où $\bar{x}_{i\text{moy}}$ est la moyenne de l'échantillon 1 puis 2 et $s^2 = s_1^2(n_1 - 1) + s_2^2(n_2 - 1) / (n_1 + n_2 - 2)$.

$$t_0 = \frac{|\bar{x}_1 - \bar{x}_2|}{s \times \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Et on pose :

Propriété : soit $\tilde{N} \in]0 ; 1[$:

_ $t_0 \geq t_{\tilde{N};n_1+n_2-2}$ à Il existe une différence entre les deux moyennes au seuil \tilde{N} .
 _ $t_0 < t_{\tilde{N};n_1+n_2-2}$ à Aucune différence n'a pu être observée.

Remarques :

$$t_0 = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2 + s_2^2}{n}}}$$

à Si $n_1 = n_2 = n$, le calcul est simplifié :

$$t_0 = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2 + s_2^2}{n_1 + n_2}}}$$

à Si n_1 et $n_2 > 30$, alors t_0 peut être remplacé par :

I Test de Wilcoxon :

Il s'agit de tester l'homogénéité de deux échantillons E et E' de tailles n et n' extraits de deux populations P et P' de façon indépendante et non exhaustive. P présente une moyenne m et P' présente une moyenne m' .

Le test de Wilcoxon ne suppose pas que les populations sont distribuées suivant une loi binomiale. Il est dit non paramétrique.

On suppose par la suite que $n \geq n'$. On veut tester l'hypothèse H_0 : "Il n'y a pas de différence entre les moyennes m et m' ".

Le rang :

On range les observations sur $E \cup E'$ par ordre croissant de 1 à $n+n'$. On attribue à chaque valeur x_i (observée) le nombre \hat{a}_i défini par :

$\hat{a}_i(x_i) =$ la moyenne des rangs de valeurs identiques à x_i .

Table de Wilcoxon :

Pour $n \geq n' \geq 0$; $1 \leq h \leq n+n'$ et n, n' , deux entiers avec $n \geq n'$, la table de Wilcoxon fournit un intervalle $]\hat{e}_1 ; \hat{e}_2[$.

Test : _ Si $h \in]\hat{e}_1 ; \hat{e}_2[$, il existe une différence significative entre m et m' .

_ Si $h \notin]\hat{e}_1 ; \hat{e}_2[$, aucune différence significative n'est observée.