



Munich Personal RePEc Archive

A method for avoiding tables for centiles in the case of confidence regions and statistical tests

Ciuiu, Daniel

Technical University of Civil Engineering, Bucharest, Romania,
Romanian Institute for Economic Forecasting

6 October 2004

Online at <https://mpra.ub.uni-muenchen.de/15029/>
MPRA Paper No. 15029, posted 06 May 2009 13:22 UTC

UNE MODALITÉ D'ÉVITER LES TABLES DES CENTILES DANS LE CAS DES RÉGIONS DE CONFIANCE ET DES TESTS STATISTIQUES

DANIEL CIUIU

RÉSUMÉ. Dans cet article on va déterminer des régions de confiance pour les paramètres d'une répartition et des versions de quelques tests sans utiliser les tables qui contiennent des centiles.

Classification AMS 2000 des sujets. 62F25, 62F03.

Mots clefs et phrases. intervalles de confiance, tests statistiques.

1. Introduction

Les intervalles de confiance sont déterminés d'habitude en utilisant des tables statistiques qui contiennent des centiles pour quelques répartitions (par exemple les centiles de la répartition normale réduite $N(0, 1)$, les centiles de la répartition de Student à n degrés de liberté, t_n , les centiles de la répartition de χ^2 à n degrés de liberté ou les centiles de la répartition de Snedecor—Fisher d'ordres m et n).

Les tests statistiques, à voir le test de concordance de χ^2 utilisent aussi ces tables avec des centiles. Mais cette chose est difficile à implanter aux ordinateurs, parce qu'on a besoin d'un fichier avec les centiles [1], [2], [3].

Dans tout l'article on note par $\overline{X^j}$ la moyenne de l'échantillon pour X^j et par θ le vecteur des paramètres qui définissent la répartition de la variable aléatoire X . On considère qu'on connaît pour $j = \overline{1, k}$ des formules pour la moyenne du X^j et pour la variance du X^j dépendant de θ . On note ces valeurs par $m_j(\theta)$ et respectivement par $D_j(\theta)$.

Pour déterminer les régions de confiance on utilise l'inégalité de Chébychev et on va résoudre un système d'inéquations.

Pour les tests, on va calculer, s'il est possible, la fonction de répartition évaluée pour la statistique concernée qui va être comparée avec certains centiles. On utilise le lemme de Hincin qui dit que si la variable aléatoire X a la fonction de répartition $F(\cdot)$, $F(X)$ suit une loi uniforme sur $[0, 1]$.

2. Régions de confiance

À présent on a une variable aléatoire X avec la fonction de répartition $F(x; \theta_1, \dots, \theta_k)$ dépendant de k paramètres. Nous voulons déterminer une région

de confiance pour $\theta = (\theta_j)_{j=\overline{1,k}}$ avec l'erreur ε_1 . On note par $\varepsilon = \frac{\varepsilon_1}{k}$. Si on écrit l'inégalité de Chébychev pour X^j on obtient

$$P\left(\left|\overline{X^j} - m_j(\theta)\right| > \sqrt{\frac{D_j(\theta)}{n\varepsilon}}\right) \leq \varepsilon. \quad (1)$$

En utilisant (1), la probabilité d'exister j de sorte que $1 \leq j \leq k$ et $\left|\overline{X^j} - m_j(\theta)\right| > \sqrt{\frac{D_j(\theta)}{n\varepsilon}}$ est plus petit ou égale à ε_1 .

Il résulte qu'on doit résoudre le système d'inéquations

$$\left|\overline{X^j} - m_j(\theta)\right| \leq \sqrt{\frac{D_j(\theta)}{n\varepsilon}}, \quad j = \overline{1,k}, \quad (2)$$

qui est équivalent à

$$m_j^2(\theta) - 2 \cdot \overline{X^j} \cdot m_j(\theta) - \frac{D_j(\theta)}{n \cdot \varepsilon} + \overline{X^j}^2 \leq 0, \quad j = \overline{1,k}. \quad (2')$$

On considère la variable aléatoire X normale $N(m, \sigma^2)$ et on a $\theta = (m, \sigma^2)^T$, $k = 2$. Par des calculs on obtient $m_1(\theta) = m$, $D_1(\theta) = \sigma^2$, $m_2(\theta) = m^2 + \sigma^2$ et $D_2(\theta) = 4 \cdot m^2 \cdot \sigma^2 + 2 \cdot \sigma^4$.

On note par $\alpha = m^2$ et par $\beta = \sigma^2$. L'inéquation du système (2') pour $j = 1$ devient

$$\alpha - \frac{\beta}{n\varepsilon} + \overline{X}^2 \leq 2m\overline{X}. \quad (3)$$

Si $\overline{X} = 0$, la région (3) est

$$\beta \geq n \cdot \varepsilon \cdot \alpha, \quad (3')$$

qui est dans le système d'axes $\alpha O \beta$ un angle déterminé par $\alpha = 0$, $\beta \geq 0$ et $\beta = n \cdot \varepsilon \cdot \alpha$, $\alpha \geq 0$. Dans le système $m O \beta$ (3') est l'intérieur de la parabole $\beta = n \cdot \varepsilon \cdot m^2$.

Si $\overline{X} \neq 0$, (3) est

$$\alpha^2 + \frac{1}{n^2\varepsilon^2}\beta^2 - \frac{2}{n\varepsilon}\alpha\beta - 2\overline{X}^2\alpha - \frac{2\overline{X}^2}{n\varepsilon}\beta + \overline{X}^4 \leq 0. \quad (3'')$$

La frontière du (3'') est

$$\alpha^2 + \frac{1}{n^2\varepsilon^2}\beta^2 - \frac{2}{n\varepsilon}\alpha\beta - 2\overline{X}^2\alpha - \frac{2\overline{X}^2}{n\varepsilon}\beta + \overline{X}^4 = 0. \quad (4)$$

On calcule les invariants pour (4) et on obtient $I = 1 + \frac{1}{n^2\varepsilon^2}$, $\delta = 0$ et $\Delta = -\frac{4\overline{X}^4}{n^2\varepsilon^2}$. Donc (4) est une parabole.

L'axe de la parabole est

$$\beta = n\varepsilon \cdot \alpha - \frac{n\varepsilon(n^2\varepsilon^2 - 1)}{n^2\varepsilon^2 + 1} \overline{X}^2. \quad (5)$$

Le sommet de la parabole est $V\left(\frac{n^4\varepsilon^4}{(n^2\varepsilon^2+1)^2} \overline{X}^2, \frac{n\varepsilon}{(n^2\varepsilon^2+1)^2} \overline{X}^2\right)$, et les intersections avec les axes sont $A(\overline{X}^2, 0)$ et $B(0, n\varepsilon \cdot \overline{X}^2)$. L'autre axe est

$$\beta = -\frac{1}{n\varepsilon} \alpha + \frac{n\varepsilon}{n^2\varepsilon^2 + 1} \overline{X}^2. \quad (5')$$

Parce que $(0, 0)$ n'est pas dans $(3'')$, la région de confiance déterminée par la première inéquation est l'intérieur de la parabole.

L'autre inéquation est

$$\alpha^2 + \left(1 - \frac{2}{n\varepsilon}\right) \beta^2 + 2\left(1 - \frac{2}{n\varepsilon}\right) \alpha\beta - 2\overline{X}^2\alpha - 2\overline{X}^2\beta + \overline{X}^2 \leq 0. \quad (6)$$

La frontière du (6) est

$$\alpha^2 + \left(1 - \frac{2}{n\varepsilon}\right) \beta^2 + 2\left(1 - \frac{2}{n\varepsilon}\right) \alpha\beta - 2\overline{X}^2\alpha - 2\overline{X}^2\beta + \overline{X}^2 = 0 \quad (7)$$

On calcule les invariants pour (7) et on obtient $I = \frac{2(n\varepsilon-1)}{n\varepsilon}$, $\delta = \frac{2(n\varepsilon-2)}{n^2\varepsilon^2}$ et $\Delta = -\frac{4\overline{X}^2}{n^2\varepsilon^2}$. Si $n\varepsilon > 2$ (un fait statistique raisonnable si on tient compte que n est le volume de l'échantillon) (7) est une ellipse.

Le centre de l'ellipse est $C\left(0, \frac{n\varepsilon}{n\varepsilon-2} \overline{X}^2\right)$ et l'angle de rotation est φ de sorte que $\tan(\varphi) = -\frac{1+\sqrt{(n\varepsilon-2)^2+1}}{n\varepsilon-2}$.

Les axes de l'ellipse sont a et b et on a

$$a^2 = \frac{n\varepsilon(n\varepsilon - 1 + \sqrt{n^2\varepsilon^2 - 4n\varepsilon + 5})}{(n\varepsilon - 2)^2} \overline{X}^2 \quad \text{et} \quad (8)$$

$$b^2 = \frac{n\varepsilon(n\varepsilon - 1 - \sqrt{n^2\varepsilon^2 - 4n\varepsilon + 5})}{(n\varepsilon - 2)^2} \overline{X}^2. \quad (8')$$

Les points d'intersection avec les axes sont $A\left(0, \frac{1-\sqrt{\frac{2}{n\varepsilon}}}{1-\frac{2}{n\varepsilon}} \overline{X}^2\right)$, $B\left(0, \frac{1+\sqrt{\frac{2}{n\varepsilon}}}{1-\frac{2}{n\varepsilon}} \overline{X}^2\right)$ et $C(\overline{X}^2, 0)$. Donc $O\alpha$ est tangente à l'ellipse.

Quelle que soit \overline{X} ($= 0$ ou $\neq 0$) et pour tout $\beta = \sigma^2$ on a une intervalle $[a_\beta, b_\beta]$ avec $a_\beta \geq 0$ de telle sorte que $\alpha = m^2 \in [a_\beta, b_\beta]$. Si $\overline{X} = 0$, on a $m \in [\sqrt{a_\beta}, \sqrt{b_\beta}] \cup [-\sqrt{b_\beta}, -\sqrt{a_\beta}]$. Si $\overline{X} \neq 0$, on a $m \in [\sqrt{a_\beta}, \sqrt{b_\beta}]$ ou $m \in [-\sqrt{b_\beta}, -\sqrt{a_\beta}]$, dépendant si $\overline{X} > 0$, respectivement si $\overline{X} < 0$.

Exemple 2.1. Soit X qui suit une répartition normale $N(m, \sigma^2)$. On considère un échantillon de volume $n = 1000$ et l'erreur maximale $\varepsilon_1 = 0.2$.

Il resulte $\varepsilon = 0.1$ et $n \cdot \varepsilon = 100$. On a dans cet exemple $\overline{X^2} = 0.231$.

La ligne oblique de l'angle est $\beta = 100\alpha$.

L'ellipse est $\alpha^2 + (1 - \frac{1}{50})\beta^2 + 2(1 - \frac{1}{50})\alpha\beta - 2 \cdot 0.231 \cdot \alpha - 2 \cdot 0.231 \cdot \beta + 0.231^2 = 0$.

La région de confiance pour (α, β) se trouve dans la figure suivante.

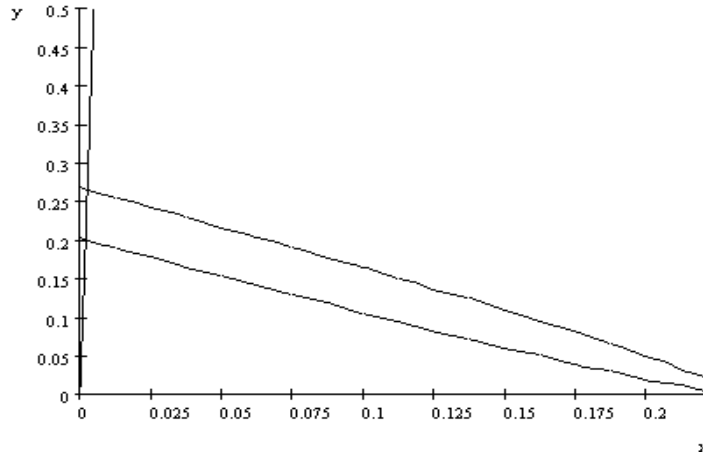


Fig. 1 : La région de confiance pour (α, β) .

La région de confiance pour (m, σ^2) est déterminée par la parabole $\beta = 100m^2$ et $|0.231 - m^2 - \beta| = \frac{\sqrt{4 \cdot m^2 \cdot \beta + 2 \cdot \beta^2}}{10}$.

Elle se trouve dans la figure suivante.

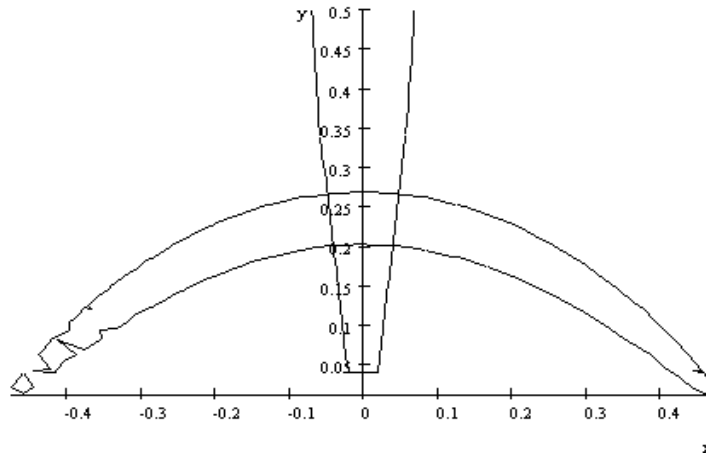


Fig. 2 : La région de confiance pour (m, σ^2)

Exemple 2.2. Soit X qui suit une répartition normale $N(m, \sigma^2)$. On considère un échantillon de volume $n = 1000$ et l'erreur maximale $\varepsilon_1 = 0.2$.

Il resulte $\varepsilon = 0.1$ et $n \cdot \varepsilon = 100$. On a dans cet exemple $\overline{X} = 5.293$ et $\overline{X^2} = 28.1$.

La parabole est

$$\alpha^2 + \frac{1}{100}\beta^2 - \frac{1}{50}\alpha\beta - 2 \cdot 5.293^2 \cdot \alpha - \frac{2 \cdot 5.293^2}{100} \cdot \beta + 5.293^4 = 0.$$

L'ellipse est

$$\alpha^2 + (1 - \frac{1}{50})\beta^2 + 2(1 - \frac{1}{50})\alpha\beta - 2 \cdot 28.1 \cdot \alpha - 2 \cdot 28.1 \cdot \beta + 28.1^2 = 0.$$

La région de confiance pour (α, β) se trouve dans la figure suivante.

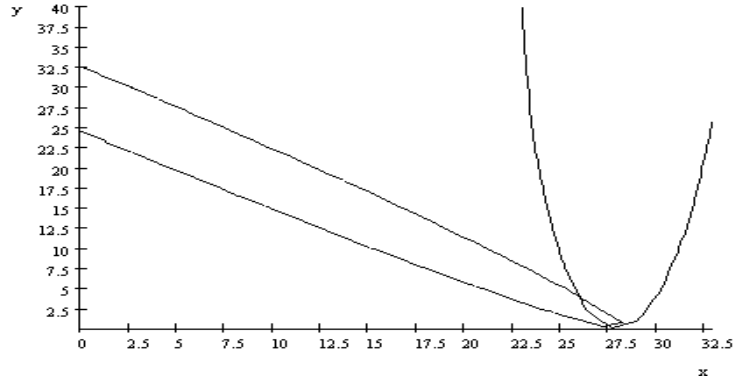


Fig. 3 : La région de confiance pour (α, β) .

La région de confiance pour (m, σ^2) est déterminée par $|m - 5.293| = \frac{\sqrt{m^2 + \beta}}{10}$ et $|\beta + m^2 - 28.1| = \frac{\sqrt{4 \cdot m^2 \cdot \beta + 2 \cdot \beta^2}}{10}$. Elle se trouve dans la figure suivante.

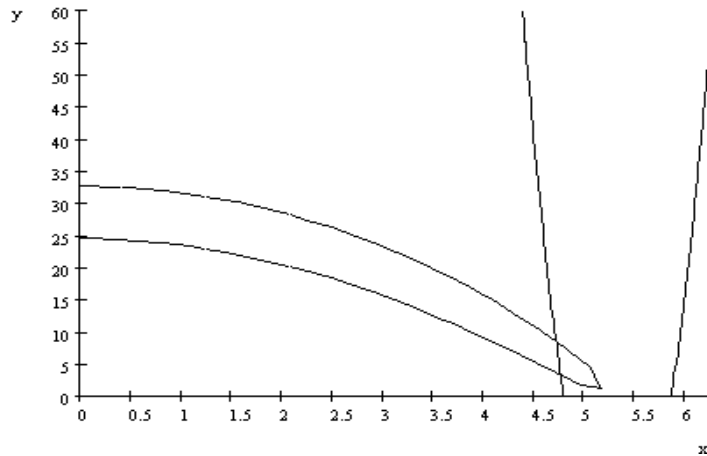


Fig. 4 : La région de confiance pour (m, σ^2)

3. Versions sans centiles pour quelques tests statistiques

Le test T bilatéral vérifie l'hypothèse nulle $H_0 : m = m_0$ contre l'hypothèse alternative $H_1 : m \neq m_0$ avec l'erreur de premier ordre ε . On a un échantillon de volume n de la variable X . On calcule $T = \frac{\bar{X} - m_0}{S} \sqrt{n - 1}$, où \bar{X} est la moyenne de l'échantillon et S^2 est la variance de l'échantillon [1], [3]. On accepte H_0 si $|T| < T_{n-1}(1 - \frac{\varepsilon}{2})$.

Le test T unilatéral gauche vérifie l'hypothèse nulle $H_0 : m = m_0$ contre l'hypothèse alternative $H_1 : m < m_0$ avec l'erreur de premier ordre ε . On accepte H_0 si $T > T_{n-1}(\varepsilon)$.

Le test T unilatéral droit vérifie l'hypothèse nulle $H_0 : m = m_0$ contre l'hypothèse alternative $H_1 : m > m_0$ avec l'erreur de premier ordre ε . On accepte H_0 si $T < T_{n-1}(1 - \varepsilon)$.

Mais la densité de répartition de Student à n degrés de liberté est

$$f(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{1}{2}\right)\sqrt{n}} \cdot \frac{1}{\left(1 + \frac{x^2}{n}\right)^{\frac{n+1}{2}}}. \quad (9)$$

On calcule la fonction de répartition et on obtient

$$F(x) = \frac{2^n \cdot \Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{1}{2}\right)} \cdot \int_0^{\frac{x+\sqrt{x^2+n}}{\sqrt{n}}} \frac{v^{n-1}}{(v^2+1)^n} dv. \quad (10)$$

Dans (10) on remarque la possibilité de calculer $F(x)$ pour chaque x .

Dans le test bilatéral on accepte H_0 si $F(T) \in \left(\frac{\varepsilon}{2}, 1 - \frac{\varepsilon}{2}\right)$. Dans le test unilatéral gauche on accepte H_0 si $F(T) > \varepsilon$. Dans le test unilatéral droit on accepte H_0 si $F(T) < 1 - \varepsilon$.

Le test de χ^2 bilatéral vérifie l'hypothèse nulle $H_0 : \sigma^2 = \sigma_0^2$ contre l'hypothèse alternative $H_1 : \sigma^2 \neq \sigma_0^2$. On calcule $X^2 = \frac{(n-1) \cdot S'^2}{\sigma_0^2}$. On accepte H_0 si $X^2 \in \left(\chi_{n-1}^2\left(\frac{\varepsilon}{2}\right), \chi_{n-1}^2\left(1 - \frac{\varepsilon}{2}\right)\right)$. Le test de χ^2 unilatéral gauche vérifie l'hypothèse nulle $H_0 : \sigma^2 = \sigma_0^2$ contre l'hypothèse alternative $H_1 : \sigma^2 < \sigma_0^2$. On accepte H_0 si $X^2 > \chi_{n-1}^2(\varepsilon)$ au risque d'erreur de premier ordre ε . Le test de χ^2 unilatéral droit vérifie l'hypothèse nulle $H_0 : \sigma^2 = \sigma_0^2$ contre l'hypothèse alternative $H_1 : \sigma^2 > \sigma_0^2$. On calcule $X^2 = \frac{(n-1) \cdot S'^2}{\sigma_0^2}$. On accepte H_0 si $X^2 < \chi_{n-1}^2(1 - \varepsilon)$.

Mais la répartition de χ_{2n}^2 coïncide avec la répartition Erlang $E_{n, \frac{1}{2}}$. La fonction répartition Erlang $E_{n, \lambda}$ d'ordre n et paramètre λ est

$$F_{n, \lambda}(x) = 1 - e^{-\lambda x} \cdot \sum_{j=0}^{n-1} \frac{\lambda^j \cdot x^j}{j!}. \quad (11)$$

Donc si n n'est pas pair ($n-1$ est pair) on peut calculer la fonction de répartition du X^2 . Dans le test de χ^2 bilatéral on accepte H_0 si $F_{\frac{n-1}{2}, \frac{1}{2}}(X^2) \in \left(\frac{\varepsilon}{2}, 1 - \frac{\varepsilon}{2}\right)$. Dans le test de χ^2 unilatéral gauche on accepte H_0 si $F_{\frac{n-1}{2}, \frac{1}{2}}(X^2) > \varepsilon$. Dans le test de χ^2 unilatéral droit on accepte H_0 si $F_{\frac{n-1}{2}, \frac{1}{2}}(X^2) < 1 - \varepsilon$.

Pour le test de Tukey pour l'égalité des moyennes on a k échantillons des volumes n $(X_{ij})_{1 \leq i \leq k, 1 \leq j \leq n}$ concernant k variables aléatoires normales $N(\mu_i, \sigma^2)$. On vérifie l'hypothèse nulle $H_0 : \mu_i = \mu_j$ pour chaque i, j contre l'hypothèse alternative H_1 : il existent i, j de telle sorte que $\mu_i \neq \mu_j$ au risque d'erreur de premier ordre ε . On note par \bar{X}_i la moyenne de l'échantillon i , par \bar{X}_{\min} le minimum des \bar{X}_i et par \bar{X}_{\max} le maximum des \bar{X}_i . On calcule S'^2 un estimateur pour σ^2 qui suit la répartition de χ_r^2 et $q = \frac{\bar{X}_{\max} - \bar{X}_{\min}}{S' \sqrt{\frac{2}{n}}}$. On accepte

H_0 si $q \leq t_r\left(\frac{\varepsilon}{2}\right)$.

Mais comment on peut voir dans les tests de Student pour la moyenne si on ne connaît pas la variance, on peut calculer la fonction de répartition pour q , $F(q)$. On accepte H_0 si $F(q) \leq \frac{\varepsilon}{2}$.

Pour le test de Hartley pour l'égalité des variances on a k échantillons des volumes $n+1$ $(X_{ij})_{1 \leq i \leq k, 1 \leq j \leq n+1}$ sur k variables aléatoires normales $N(\mu_i, \sigma_i^2)$. On vérifie l'hypothèse nulle $H_0 : \sigma_i^2 = \sigma_j^2$ pour chaque i, j contre l'hypothèse alternative H_1 : il existent i, j de sorte que $\sigma_i^2 \neq \sigma_j^2$ au risque d'erreur de premier ordre ε .

On note par $\bar{X}_i = \frac{1}{n+1} \cdot \sum_{j=1}^{n+1} X_{ij}$ et par $S_i^2 = \frac{1}{n} \cdot \sum_{j=1}^{n+1} (X_{ij} - \bar{X}_i)^2$. On note aussi par S_{\min}^2 le minimum des S_i et par S_{\max}^2 le maximum. On calcule $F_{\max} = \frac{S_{\max}^2}{S_{\min}^2}$. On accepte H_0 si $F_{\max} \leq F_{n,n}(1 - \frac{\varepsilon}{2})$, où $F_{n,n}(1 - \frac{\varepsilon}{2})$ est la centile d'ordre $1 - \frac{\varepsilon}{2}$ pour la répartition de Snedecor—Fisher d'ordres (n, n) .

Mais la densité de cette répartition est

$$g_{n,n}(x) = \frac{\Gamma(n)}{\Gamma^2(\frac{n}{2})} \cdot \frac{x^{\frac{n}{2}-1}}{(1+x)^n}. \quad (12)$$

On peut calculer la fonction de répartition et on obtient

$$G_{n,n}(x) = \frac{2 \cdot \Gamma(n)}{\Gamma^2(\frac{n}{2})} \cdot \frac{x^{\frac{n}{2}-1}}{(1+x)^n}. \quad (12')$$

Le test de concordance de χ^2 vérifie l'hypothèse nulle H_0 : la variable aléatoire X a la fonction de répartition $F(x; \theta_1, \dots, \theta_k)$, où $\theta_1, \dots, \theta_k$ sont les paramètres de la répartition, contre l'hypothèse alternative H_1 : la variable aléatoire X n'a pas la fonction de répartition $F(x; \theta_1, \dots, \theta_k)$ avec l'erreur de premier ordre ε . On a un échantillon de volume n de X et r intervalles, $r > k$. On note par n_i le numéro des valeurs de l'échantillon dans l'intervalle I_i et par $n'_i = n \cdot F(I_i; \hat{\theta}_1, \dots, \hat{\theta}_k)$, où $\hat{\theta}_1, \dots, \hat{\theta}_k$ sont des estimations pour $\theta_1, \dots, \theta_k$.

On calcule $X^2 = \sum_{j=1}^r \frac{(n_i - n'_i)^2}{n'_i}$ et on accepte H_0 si $X^2 < \chi_{r-k-1}^2(1 - \varepsilon)$, où $\chi_{r-k-1}^2(1 - \varepsilon)$ est la centile de χ^2 à $r - k - 1$ degrés de liberté.

Si $r - k - 1$ est pair, on peut calculer la fonction de répartition de χ_{r-k-1}^2 en utilisant (11). Donc on accepte H_0 si $F_{\frac{r-k-1}{2}, \frac{1}{2}}(X^2) < 1 - \varepsilon$.

4. Conclusions

Pour les régions de confiance l'erreur ε est maximale (elle peut être plus petite que ε). Dans le cas de k paramètres on doit résoudre un système de k inéquations.

Pour les tests on calcule la fonction de répartition (s'il est possible) pour la statistique et on compare cette valeur avec l'ordre de centile.

Ces considérations peuvent aider le programmeur à faire des programmes pour résoudre des problèmes de statistique. On doit choisir si on utilise l'inégalité de Chébychev ou les tables des centiles (si le système est trop difficile à résoudre, on utilise les centiles).

RÉFÉRENCES

- [1] G. Ciucu, V. Craiu, *Inférence statistique*. Editure Didactique et Pedagogique, Bucarest, 1974. (en roumain)
- [2] V. Craiu, V. Preda, *Des tests pour vérifier la normalité*. Editure de l'Université de Bucarest, 1981. (en roumain)
- [3] A. Popescu, V. Petrehuş, *Probabilitées et statistique*. Editure de l'Université Technique des Constructions, Bucarest, 1997. (en roumain)

(Daniel Ciuiu) DEPARTEMENT DES MATHÉMATIQUES,
UNIVERSITÉ TECHNIQUE DES CONSTRUCTIONS, BUCAREST,
BD. LACUL TEI NO. 124, BUCAREST, ROUMANIE
E-mail address: dciuiu@yahoo.com