

La vérification de l'hypothèse sur la constance des probabilités.

Par

Stanisław Kołodziejczyk.

(Laboratoire Biométrie de l'Institut Nencki, Soc. Scient. de Varsovie).

1. Dans un mémoire récent Mrs. J. Neyman et E. S. Pearson proposent une méthode générale¹⁾ de vérification des hypothèses statistiques.

Ayant donné un fait observé F et une hypothèse statistique H , on calcule ce qu'on appelle la vraisemblance de cette hypothèse λ_H . Si λ_H est petite, par ex. si $\lambda_H \leq \lambda_0$, on convient de conclure que l'hypothèse est probablement fausse. Si au contraire $\lambda_0 < \lambda_H$, on admet qu'on n'a pas de raison suffisante pour une conclusion pareille. Le nombre λ_0 doit être choisi de manière que le danger de rejection d'une hypothèse vraie ne soit pas trop grand. Ce danger est mesuré par P_{λ_0} — la valeur de la probabilité pour qu'on ait $\lambda_H \leq \lambda_0$, déterminée par l'hypothèse considérée H , ou la borne supérieure de cette probabilité, si celle-ci n'est pas déterminée par H . Il est évident d'ailleurs, que si l'on rejette l'hypothèse H lorsque $\lambda_H \leq \lambda_0$, la probabilité π pour qu'on rejette une hypothèse vraie ne surpasse pas P_{λ_0} . Or, si λ_0 est choisi de telle façon que $P = \varepsilon$, ε étant un nombre arbitraire entre zéro et un, et qu'on accepte la règle de rejeter l'hypothèse H lorsque $\lambda_H \leq \lambda_0$ et de l'accepter dans les autres cas, on peut être sûr qu'une hypothèse vraie sera rejetée avec une fréquence moyenne, qui ne surpasse pas ε .

¹⁾ J. Neyman and E. S. Pearson: *On the Use and Interpretation of Certain Test Criteria for Purposes of Statistical Inference*. *Biometrika* Vol. XX—A. P. 175—240 et 264—294.

Le point essentiel dans la méthode de Mrs. J. Neyman et E. S. Pearson consiste dans le principe que le degré de notre confiance en une hypothèse statistique H peut être mesuré par la valeur de la vraisemblance λ_H . Evidemment ce principe peut être admis ou rejeté, selon qu'il semble intuitif ou non. De ce point de vue il est intéressant que toutes les méthodes de vérification des hypothèses, qui sont entrées en usage général et qui ont été examinées¹⁾, sont des conséquences du principe mentionné.

Le but de la note présente est de démontrer, que la méthode bien connue de Lexis-Bortkiewicz de vérification de l'hypothèse sur la constance des probabilités dans plusieurs séries des épreuves indépendantes est — elle aussi — une conséquence du principe des Mrs. J. Neyman et E. S. Pearson²⁾.

2. Rappelons la terminologie. Soit F un fait, déterminé par les coordonnées

$$(1) \quad x_1, x_2, \dots, x_s$$

qui peuvent varier dans certaines limites, et h — une hypothèse concernant F . Si h détermine la probabilité (sensu stricto, ou la probabilité élémentaire) de F , nous dirons, que h est une hypothèse statistique simple. Toute hypothèse H , qui n'est pas simple est dite composée. Il est évident que l'hypothèse composée H peut être transformée en une hypothèse simple. Il suffit pour cela d'adjoindre à H arbitrairement quelques suppositions supplémentaires pour que l'ensemble de ces suppositions détermine la probabilité de F . Si h est une hypothèse simple, qui peut être obtenue de H par ce procédé, nous dirons, que h appartient à H .

Soit Ω l'ensemble des hypothèses simples qu'on considère dans un cas donné comme admissibles. Soit h une de ces hypo-

¹⁾ J. Neyman and E. S. Pearson: *On the Use and Interpretation...* Loc. cit.

J. Neyman and E. S. Pearson: *On the Problem of Two Samples*. Bull. de l'Académie Polonaise des Sciences et des Lettres. A, 1930, p: 73—96.

J. Neyman: *Contribution to the Theory of Certain Test Criteria*. Bull. de l'Institut International de Statistique, 1929. P. 44—88.

J. Neyman: *Sur la limite de la vraisemblance de l'hypothèse*. C. R., t. 188, p. 1360.

J. Neyman: *Sur une méthode de vérification des hypothèses*. Ibidem p. 1467.

²⁾ Ce problème m'a été posé par Mr. J. Neyman, à qui je dois aussi quelques indications sur la méthode.

thèses et P_{hF} — la probabilité (sensu stricto ou la probabilité élémentaire) de F , déterminée par h . Nous allons supposer que l'ensemble des valeurs de P_{hF} , qui correspondent à un même fait F est borné, quel que soit F . Soit P_F la borne supérieure de P_{hF} par rapport à l'ensemble Ω et pour un fait fixé F .

Si l'on a observé F , on appelle la vraisemblance de l'hypothèse simple h le rapport

$$(2) \quad \lambda_h = \frac{P_{hF}}{P_F}.$$

Si H est une hypothèse composée, pour définir la vraisemblance de H lorsqu'on a observé F , on considère le sousensemble Ω_H des hypothèses simples qui appartiennent à H . Soit H_{HF} la borne supérieure des nombres P_{hF} par rapport à l'ensemble Ω_H et le fait F . La vraisemblance λ_H de H est alors

$$(3) \quad \lambda_H = \frac{P_{HF}}{P_F}.$$

3. Le problème, que nous allons considérer, peut être précisé comme suit.

Nous considérons les épreuves indépendantes qui peuvent donner lieu à un des deux événements: E ou sa négation \bar{E} . Ces épreuves sont effectuées en s séries et soit n_i le nombre des épreuves, qui appartiennent à l' i -ème série ($i = 1, 2, \dots, s$). Supposons qu'on sait, que la probabilité de E dans les épreuves appartenant à une même série est constante. Soit p_i la valeur de cette probabilité qui correspond à l' i -ème série des épreuves.

Le fait observé F consiste dans les nombres k_i des épreuves de l' i -ème série, qui ont donné lieu à l'événement E , pour $i = 1, 2, \dots, s$.

L'hypothèse H , qu'on doit vérifier, affirme que

$$(4) \quad p_1 = p_2 = \dots = p_s$$

c'est-à-dire que la probabilité de l'événement E était la même dans toutes les séries des épreuves. Désignons par p la valeur de cette probabilité. Il est à remarquer, que p n'est pas déterminée par l'hypothèse H , qui par conséquent est une hypothèse composée.

Nous allons vérifier cette hypothèse par rapport à l'ensemble Ω des hypothèses admissibles, qui renferme toute hypothèse h précisant les valeurs des probabilités p_i quelconques $0 \leq p_i \leq 1$ pour $i = 1, 2, \dots, s$.

4. La probabilité du fait observé F déterminée par une hypothèse h est

$$(5) \quad P_{hF} = \prod_{i=1}^s C_{n_i}^{k_i} p^{k_i} (1 - p_i)^{n_i - k_i}$$

La valeur maximum de P_{hF} est égale à

$$(6) \quad P_F = \prod_{i=1}^s C_{n_i}^{k_i} q_i^{k_i} (1 - q_i)^{n_i}$$

où

$$(7) \quad q_i = \frac{k_i}{n_i}.$$

La probabilité du fait F , déterminée par une hypothèse simple qui appartient à H est égale à

$$(8) \quad P'_{HF} = p^{k_0} (1 - p)^{n_0 - k_0} \prod_{i=1}^s C_{n_i}^{k_i}$$

où

$$(9) \quad k_0 = \sum_{i=1}^s k_i; \quad n_0 = \sum_{i=1}^s n_i.$$

Posons encore

$$(10) \quad q_0 = \frac{k_0}{n_0}.$$

Alors le maximum de P'_{HF} sera

$$(11) \quad P_{HF} = q_0^{k_0} (1 - q_0)^{n_0 - k_0} \prod_{i=1}^s C_{n_i}^{k_i}$$

et la vraisemblance λ_H de l'hypothèse H

$$(12) \quad \lambda_H = \frac{q_0^{k_0} (1 - q_0)^{n_0 - k_0}}{\prod_{i=1}^s q_i^{k_i} (1 - q_i)^{n_i - k_i}}.$$

La formule (12) présente la solution de la première partie du problème. Pour la compléter il faudrait considérer la probabilité

$$(13) \quad P\{\lambda_H \leq \lambda_0\}$$

pour qu'on ait $\lambda_H \leq \lambda_0$, λ_0 étant un nombre positif quelconque. En cas où cette probabilité serait déterminée par l'hypothèse H , c'est-à-dire si elle était indépendante de la valeur commune p des probabilités de l'événement E dans toutes les s séries des épreuves, qui n'est pas déterminée par H , la solution serait complétée par le calcul de (13). Si au contraire $P\{\lambda_H \leq \lambda_0\}$ n'était pas déterminée par H , il serait nécessaire de calculer la borne supérieure de $P\{\lambda_H \leq \lambda_0\}$.

Evidemment

$$(14) \quad P\{\lambda_H \leq \lambda_0\} = \sum \prod_{i=1}^s C_{n_i}^{k_i} p^{k_i} (1-p)^{n_i-k_i},$$

où la somme Σ s'étend sur tous les systèmes des valeurs des nombres k_i pour lesquels

$$(15) \quad 0 \leq k_i \leq n_i \quad (i = 1, 2, \dots, s)$$

$$(16) \quad \lambda_H = \prod_{i=1}^s \frac{q_0^{k_i} (1-q_0)^{n_i-k_i}}{q_i^{k_i} (1-q_i)^{n_i-k_i}} \leq \lambda_0.$$

Le calcul de la somme (14) semble inabordable et nous nous bornerons au calcul de sa limite, lorsque (a) le nombre n_0 des épreuves effectuées croît indéfiniment de telle manière que

$$(17) \quad n_i \geq \nu n_0$$

ν étant un nombre positif fixe. Nous supposons de plus que (b)

$$(18) \quad 0 < p < 1,$$

ce qui est d'ailleurs évidemment une limitation sans importance.

Désignons par $I(A)$ l'intégrale

$$(19) \quad I(A) = \frac{\prod_{i=1}^s \sqrt{n_i}}{(2\pi p(1-p))^{\frac{s}{2}}} \int \dots \int_A e^{-\frac{\sum_{i=1}^s n_i (q_i - p)^2}{2p(1-p)}} dq_1 dq_2 \dots dq_s$$

étendue sur un domaine A dans l'espace à s dimensions, et où q_1, q_2, \dots, q_s désignent les variables continues. Considérons λ_H comme une fonction de ces variables et désignons par W_0 le domaine défini par l'inégalité

$$(20) \quad \lambda_H > \lambda_0.$$

Si les conditions (a) et (b) sont remplies, et si le nombre n_0 est assez grand, on peut appliquer le théorème de Laplace et écrire

$$(21) \quad P\{\lambda_H > \lambda_0\} = I(W_0) + \eta.$$

$|\eta|$ étant si petit que l'on veut. Fixons un nombre positif arbitraire ε et choisissons N assez grand pour qu'on ait $|\eta| < \varepsilon$, lorsque $n_0 > N$.

Vu les propriétés connues de la fonction sous le signe de l'intégrale (19) on peut trouver un tel nombre positif χ_0 que

$$(22) \quad I(W_1) = \left(\frac{1}{\sqrt{2\pi}}\right)^s \int \dots \int_{W_1'} e^{-\frac{\sum_{i=1}^s y_i^2}{2}} dy_1 dy_2 \dots dy_s > 1 - \varepsilon,$$

où W_1 désigne le domaine défini par l'inégalité

$$(23) \quad \chi^2 = \frac{\sum_{i=1}^s (q_i - p)^2}{p(1-p)} \leq \chi_0^2$$

et W_1' est la transformation de W_1 par les formules

$$(24) \quad q_i = p + y_i \sqrt{\frac{p(1-p)}{n_i}}, \quad i = 1, 2, \dots, s.$$

Remarquons que le domaine W_1' , et par conséquent l'intégrale $I(W_1)$, sont indépendants de n_i , ($i = 1, 2, \dots, s$).

Considérons W_2 — la partie commune des domaines W_0 et W_1 . Il est aisé de voir que

$$(25) \quad P\{\lambda_H > \lambda_0\} = I(W_2) + \eta_2$$

où $|\eta_2| < 2\varepsilon$, lorsque

$$(26) \quad n_0 > N.$$

Les raisonnements qui suivent ont pour but de trouver deux domaines W_5 et W_6 tels que

$$(27) \quad W_5 \subset W_2 \subset W_6$$

et que les intégrales $I(W_5)$ et $I(W_6)$ tendent vers une même limite $1 - P_{\lambda_0}$ lorsque $n_0 \rightarrow \infty$. Il est clair qu'alors le même nombre sera aussi la limite de la probabilité $P\{\lambda_H > \lambda_0\}$.

Dans tout point du domaine W_2 on a

$$(28) \quad |q_i - p| \leq \frac{\chi_0 \sqrt{p(1-p)}}{\sqrt{n_i}} \leq \frac{\chi_0}{\sqrt{n_0}} \sqrt{\frac{p(1-p)}{v}}$$

$$(29) \quad |q_0 - p| \leq \frac{\chi_0}{\sqrt{n_0}} s \sqrt{vp(1-p)}$$

$$(30) \quad |q_i - q_0| \leq \frac{\chi_0}{\sqrt{n_0}} \frac{(sv+1)\sqrt{p(1-p)}}{v}.$$

Posons

$$(31) \quad q_i = q_0 + x_i \sqrt{\frac{q_0(1-q_0)}{n_i}}$$

Si le nombre N est assez grand il résulte de (26) et (29) qu'il existe un nombre positif $\alpha < \frac{1}{2}$ tel que

$$(32) \quad \alpha < q_0 < 1 - \alpha$$

done que

$$(33) \quad \frac{\alpha}{1-\alpha} < \frac{q_0}{1-q_0} < \frac{1-\alpha}{\alpha}$$

et

$$(34) \quad \frac{\alpha}{1-\alpha} < \frac{1-q_0}{q_0} < \frac{1-\alpha}{\alpha}.$$

Supposons que N est assez grand et fixons α satisfaisant (32). Il résulte alors de (30), (33) et (34) que les variables x_i ($i = 1, 2, \dots, s$) sont bornées dans leur ensemble et p. ex. que $|x_i| < M$ pour $i = 1, 2, \dots, s$. Un calcul simple donne alors

$$\begin{aligned} -\log \lambda_H &= \sum_{i=1}^s n_i \left(q_0 + x_i \sqrt{\frac{q_0(1-q_0)}{n_i}} \right) \left[x_i \sqrt{\frac{1-q_0}{n_i q_0}} - \right. \\ &\quad \left. - \frac{1}{2} x_i^3 \frac{1-q_0}{n_i q_0} + \frac{1}{\left(1 + \Theta_i \sqrt{\frac{1-q_0}{n_i q_0}} \right)^3} \frac{x_i^3}{3} \left(\frac{1-q_0}{n_i q_0} \right)^{3/2} \right] - \\ &\quad - \sum_{i=1}^s n_i \left(1 - q_0 - x_i \sqrt{\frac{q_0(1-q_0)}{n_i}} \right) \left[x_i \sqrt{\frac{q_0}{n_i(1-q_0)}} + \right. \end{aligned}$$

$$\begin{aligned}
 & + \frac{1}{2} x_i^2 \frac{q_0}{n_i(1-q_0)} + \left[\frac{1}{\left(1 - \Theta'_i \sqrt{\frac{q_0}{n_i(1-q_0)}}\right)^3} \frac{x_i^3 \left(\frac{q_0}{n_i(1-q_0)}\right)^{3/2}}{3} \right] = \\
 (35) \quad & = \frac{1}{2} \sum_{i=1}^s x_i^2 + \eta_s, \quad 0 \leq \Theta'_i, \Theta''_i \leq 1
 \end{aligned}$$

où $|\eta_s| < \frac{\varepsilon}{2}$ dans tout point du domaine W_2 pourvu que N soit assez grand et l'inégalité (26) — satisfaite.

Posons

$$(36) \quad Q_0^2 = -2 \lg \lambda_0$$

et considérons deux domaines W_3 et W_4 définis par des inégalités

$$(37) \quad Q^2 = \sum_{i=1}^s x_i^2 = \sum_{i=1}^s \frac{(q_i - q_0)^2 n_i}{q_0(1 - q_0)} \leq Q_0^2 - \varepsilon$$

$$(38) \quad \chi^2 = \sum_{i=1}^s \frac{(q_i - p)^2 n_i}{p(1 - p)} \leq \chi_0^2$$

et

$$(39) \quad Q^2 \leq Q_0^2 + \varepsilon$$

$$(40) \quad \chi^2 \leq \chi_0^2$$

respectivement. Il est clair que

$$(41) \quad W_3 \subset W_2 \subset W_4$$

et par conséquent

$$(42) \quad I(W_3) \leq I(W_2) \leq I(W_4).$$

Observons que les inégalités (37) et (39) sont équivalentes à

$$(43) \quad \sum_{i=1}^s \frac{(q_i - q_0)^2 n_i}{p(1 - p)} \leq (Q_0^2 - \varepsilon) \left(1 + \frac{q_0 - p}{p}\right) \left(1 - \frac{q_0 - p}{1 - p}\right)$$

et

$$(44) \quad \sum_{i=1}^s \frac{(q_i - q_0)^2 n_i}{p(1 - p)} \leq (Q_0^2 + \varepsilon) \left(1 + \frac{q_0 - p}{p}\right) \left(1 - \frac{q_0 - p}{1 - p}\right)$$

respectivement. A cause de (29) et au moyen de l'augmentation éventuelle du nombre N on peut satisfaire les inégalités :

$$(45) \quad Q_0^2 > (Q_0^2 - \varepsilon) \left(1 + \frac{q_0 - p}{p}\right) \left(1 - \frac{q_0 - p}{1 - p}\right) > Q_0^2 - 2\varepsilon$$

et

$$(46) \quad Q_0^2 < (Q_0 + \varepsilon) \left(1 + \frac{q_0 - p}{p}\right) \left(1 - \frac{q_0 - p}{1 - p}\right) < Q_0^2 + 2\varepsilon.$$

Désignons par W_5 et W_6 les domaines définis par les inégalités

$$(47) \quad \sum_{i=1}^s \frac{(q_i - q_0)^2 n_i}{p(1-p)} \leq Q_0^2 - 2\varepsilon$$

$$(48) \quad \chi^2 \leq \chi_0^2$$

et

$$(49) \quad \sum_{i=1}^s \frac{(q_i - q_0) n_i}{p(1-p)} \leq Q_0^2 - 2\varepsilon$$

$$(50) \quad \chi^2 \leq \chi_0^2$$

respectivement. Il est clair que si N est assez grand, on a pour tout $n_0 > N$

$$(51) \quad W_5 \subset W_3 \subset W_2 \subset W_4 \subset W_6$$

et par conséquent

$$(52) \quad I(W_5) \leq I(W_2) \leq I(W_6).$$

Pour calculer les limites des intégrales $I(W_5)$ et $I(W_6)$, considérons les domaines, soit W_7 et W_8 , qui correspondent aux inégalités (47) et (49) respectivement. A cause des propriétés du nombre χ_0 nous aurons

$$(53) \quad I(W_7) = I(W_5) + \eta_4$$

$$(54) \quad I(W_8) = I(W_6) + \eta_5$$

où $0 \leq \eta_4 < \varepsilon$ et $0 \leq \eta_5 < \varepsilon$.

Les intégrales $I(W_7)$ et $I(W_8)$ sont des cas particuliers de l'intégrale $I(W_\sigma)$ prise dans le domaine W_σ , où l'on a

$$(55) \quad \sum_{i=1}^s \frac{(q_i - q_0)^2 n_i}{p(1-p)} \leq \sigma^2.$$

Le calcul de l'intégrale $I(W_\sigma)$ ne présente aucune difficulté. La transformation (24) donne immédiatement

$$(56) \quad I(W_\sigma) = \left(\frac{1}{\sqrt{2\pi}} \right)^s \int \dots \int e^{-\frac{1}{2} \sum_{i=1}^s y_i^2} dy_1 dy_2 \dots dy_s$$

où l'intégration s'étend à un domaine, où l'on a

$$(57) \quad \sum_{i=1}^s y_i^2 - \frac{\sqrt{n_i}}{n_0} \left(\sum_{i=1}^s y_i \sqrt{n_i} \right)^2 \leq \sigma^2$$

Or la valeur de (56) est connue, savoir

$$(58) \quad I(W_\sigma) = C \int_0^\sigma t^{s-2} e^{-\frac{1}{2} t^2} dt$$

où

$$(59) \quad \frac{1}{C} = \int_0^\infty t^{s-2} e^{-\frac{1}{2} t^2} dt.$$

En combinant les résultats (25), (48), (51), (52) et (58), on peut écrire

$$(60) \quad C \int_0^{\sqrt{Q_{m+2s}}} t^{s-2} e^{-\frac{1}{2} t^2} dt + \eta_2 - \eta_4 \leq P\{\lambda_H > \lambda_0\} \leq C \int_0^{\sqrt{Q_{s-2s}}} t^{s-2} e^{-\frac{1}{2} t^2} dt,$$

d'où on conclut que

$$(61) \quad \lim_{n_0 \rightarrow \infty} P\{\lambda_H > \lambda_0\} = C \int_0^{Q_0} t^{s-2} e^{-\frac{1}{2} t^2} dt$$

donc que

$$(62) \quad \lim_{n_0 \rightarrow \infty} P\{\lambda_H \leq \lambda_0\} = C \int_0^{Q_0} t^{s-2} e^{-\frac{1}{2} t^2} dt = P_{\lambda_0},$$

ce qu'il fallait démontrer. Il est à remarquer, que la limite de la

probabilité $P\{\lambda_H \leq \lambda_0\}$ est parfaitement déterminée par l'hypothèse H .

Si le nombre des épreuves effectuées n_0 est très grand, l'égalité (62) permet de mesurer la probabilité $P\{\lambda_H \leq \lambda_0\}$ par la valeur de sa limite P_{λ_0} .

La technique de la vérification de l'hypothèse H sur la constance de la probabilité pendant s séries des épreuves indépendantes consiste dans le calcul du nombre

$$(63) \quad Q_H^2 = -2 \lg \lambda_H.$$

Si ce nombre est grand, on conclut que la vraisemblance de l'hypothèse H est petite et on hésite d'accepter H . On considère ensuite les tables ¹⁾ de l'intégrale (62) et on y trouve la valeur de

$$(64) \quad P_H = C \int_{Q_H}^{\infty} t^{s-2} e^{-\frac{1}{2}t^2} dt.$$

Si P_H est jugé petit, on n'hésite plus de rejeter H , vu que la probabilité pour qu'on rejette une hypothèse vraie, est dans ce cas plus petite que P_H .

Pour le calcul de Q_H on peut s'adresser à la formule

$$(65) \quad \lambda_H = \frac{q_0^{k_0} (1 - q_0)^{n_0 - k_0}}{\prod_{i=1}^s q_i^{k_i} (1 - q_i)^{n_i - k_i}}$$

ou à la formule approchée

$$(66) \quad Q_H^2 = \frac{\sum_{i=1}^s (q_i - q_0)^2 n_i}{q_0 (1 - q_0)}.$$

On voit sans peine que cette dernière formule est en connexion simple avec le coefficient de dispersion D introduit pour la vérification de l'hypothèse H par Lexis et étudié par M. L. v. Bortkiewicz, savoir

$$(67) \quad Q_H^2 = (s - 1) D^2.$$

On sait que si D est sensiblement plus grand que l'unité, d'après la méthode de Lexis il faudrait rejeter H . Le verdict serait le même si l'on partait du principe général des Mrs. J. Ney-

¹⁾ Karl Pearson: *Tables for Statisticians and Biometricians*. Cambridge, 1924.

man et E. S. Pearson. Ainsi dans le cas où l'on sait que la probabilité de l'événement E était constante pendant chaque série des épreuves, la méthode classique de la vérification de l'hypothèse H est une conséquence du principe mentionné.

Il est à remarquer que la circonstance que la loi de probabilité élémentaire de l'expression

$$(68) \quad z = \sqrt{\frac{\sum_{i=1}^s (q_i - q_0)^2 n_i}{q_0(1 - q_0)}}$$

peut être représentée approximativement par la fonction

$$(69) \quad Cz^{s-2} e^{-\frac{1}{2}z^2}$$

est connue ¹⁾. Au contraire le résultat (62) semble être nouveau.

¹⁾ Voir p. ex. R. A. Fisher: *Statistical Methods for Research Workers*. Oliver and Boyd. 1925.