

# REVUE DE STATISTIQUE APPLIQUÉE

E. MORICE

## **De quelques difficultés, pour les praticiens d'utiliser la documentation statistique**

*Revue de statistique appliquée*, tome 15, n° 1 (1967), p. 71-81

[http://www.numdam.org/item?id=RSA\\_1967\\_\\_15\\_1\\_71\\_0](http://www.numdam.org/item?id=RSA_1967__15_1_71_0)

© Société française de statistique, 1967, tous droits réservés.

L'accès aux archives de la revue « *Revue de statistique appliquée* » (<http://www.sfds.asso.fr/publicat/rsa.htm>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme  
Numérisation de documents anciens mathématiques  
<http://www.numdam.org/>

## DE QUELQUES DIFFICULTÉS, POUR LES PRATICIENS D'UTILISER LA DOCUMENTATION STATISTIQUE

E. MORICE

Il y a quelque vingt années deux professeurs avaient été chargés, à l'Institut de Statistique de l'Université de Paris, l'un d'un enseignement des méthodes statistiques élémentaires (c'était l'auteur de ces lignes), l'autre d'un enseignement de statistique mathématique. Ayant par ailleurs une occupation principale, ils se rencontraient rarement, chacun ne connaissant le cours de l'autre que par le contenu du programme qui avait été préalablement fixé, les cours rédigés n'étant pas encore publiés.

Vint le moment où un tirage ronéoté des deux cours fut disponible : je constatai alors que nous utilisions quelquefois une terminologie et des symboles différents : j'en parlai à mon collègue et lui demandai s'il n'estimait pas utile d'homogénéiser notre terminologie et nos notations.

Ce n'était certes pas chose difficile à réaliser, en principe, car nous aurions pu aisément nous mettre d'accord, mais il y avait maintenant l'obstacle créé par les cours écrits, utilisés par les étudiants suivant les cours oraux : la modification de l'un aurait dû logiquement entraîner la correction de l'autre.

L'affaire en resta là, chacun se contentant d'indiquer l'emploi possible d'autres termes ou d'autres symboles que ceux qu'il utilisait normalement. D'ailleurs, disait mon collègue, un cours est destiné à être suivi régulièrement : dès l'instant que termes et symboles ont été préalablement définis et que l'on continue à les employer, conformément à leur définition, peu importe le choix qui a été fait. De plus, les étudiants qui poursuivront leurs études en statistique, auront à étudier d'autres auteurs français et étrangers, utilisant peut-être d'autres notations : il est bon qu'ils sachent dès maintenant qu'ils auront nécessairement à faire un effort d'adaptation et que faire appel sur tel ou tel point particulier à l'ouvrage de M. X ne peut pas se réduire à appliquer sans comprendre la formule N° Y de la page N° Z.

Visant des étudiants, et compte tenu de la quasi impossibilité de normaliser, sur le plan international, et même sur le plan national, tous les symboles nécessaires à l'enseignement et à la recherche statistique ce point de vue peut se défendre.

Mais le développement de l'emploi des méthodes statistiques dans tous les domaines techniques, l'emploi, par des techniciens qui ne sont pas toujours des spécialistes de l'étude statistique, de techniques codifiées de manière systématique ou d'abaques donnant par simple lecture la réponse à une question précise, nécessitent que toutes précautions soient prises pour éviter des erreurs d'interprétation qui, sans conséquence pratique dans bien des cas, risquent d'être graves dans certains cas particuliers (interprétation des résultats de petits échantillons).

A cet égard les points suivants semblent mériter d'être pris particulièrement en considération ; notion de variance, fractiles (ou quantiles),

distribution cumulative (ou fonction de répartition) et intervalles de confiance (dans le cas de variables discrètes).

## 1 - NOTION DE VARIANCE

Pour le probabiliste, il n'y a pas de difficulté, la variance d'une variable aléatoire  $X$  ou moment centré d'ordre 2 est définie comme l'espérance mathématique (ou valeur moyenne) des carrés des écarts des valeurs de cette variable à son espérance mathématique (ou valeur moyenne)

$$V(X) = E [X - E(X)]^2 = \begin{cases} \sum p_i (x_i - m) \\ \text{ou} \\ \int (x - m)^2 f(x) dx \end{cases}$$

la sommation étant étendue à tout le domaine de variation de  $X$

$$m = E(X)$$

$$p_i \sim \Pr(X = x_i) \quad , \quad \text{variable discrète}$$

$$f(x) dx = \Pr[x < X < x + dx] \quad , \quad \text{variable continue.}$$

Il s'agit d'un concept théorique défini à partir d'une loi de probabilité et pour lequel on utilise généralement l'un des symboles  $\sigma^2$  ou  $\mu_2$ .

Le statisticien doit relier ce concept formel à la réalité généralement représentée par un nombre fini  $n$  d'observations constituant un échantillon d'une population généralement finie, d'effectif  $N$  à l'instant où l'on s'y intéresse, mais que l'on considère le plus souvent comme représentative d'une population  $P$  non finie engendrée par le même mécanisme, opérant dans les mêmes conditions, ce mécanisme étant considéré comme générateur d'une certaine loi de probabilité.

Ainsi par exemple un calculateur électronique, programmé pour produire des chiffres au hasard (loi de probabilité  $p_i = 1/10$ ,  $i = 0, 1, \dots, 9$ ), produit une population non finie de chiffres - s'il continue à fonctionner - dont une table constitue une population finie d'effectif  $N$ , dans laquelle on pourra prélever des échantillons d'effectif  $n$ .

Pratiquement on peut en général considérer cet échantillon comme représentant - pour l'information cherchée, la population  $N$  à laquelle on s'intéresse (population mère), celle-ci étant confondue avec la population  $P$  d'où au moins trois concepts liés à la variabilité dans cette population.

a) la variance dans la population ou plus exactement la variance relative à la loi de probabilité à laquelle on estime pouvoir rattacher cette population. Cette variance est généralement notée  $\sigma^2$ .

b) la variance effectivement constatée dans l'échantillon et définie comme carré moyen des écarts à la moyenne  $\bar{x}$  soit :

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (1)$$

c) l'estimation sans biais de la variance  $\sigma^2$  à partir de l'échantillon, soit :

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad , \quad (2)$$

( Dans le cas d'un tirage sans remise, ce serait :

$$\frac{N-1}{N} \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (2')$$

C'est surtout à propos des expressions (1) et (2) que se pose le problème des notations. Si les notations  $s_n^2$ ,  $\sigma'^2$ ,  $s'^2$  ou  $\sigma^2$  ou même  $\sigma^2$  sont utilisées par certains auteurs, pour l'expression (1), on constate que la notation  $s^2$  est employée suivant les auteurs, tantôt pour l'expression (1) (Indjoudjian, Kenney and Keeping, Lloyd and Lipow, Deming, Risser et Traynard, ...) mais plus souvent pour l'expression (2) (Bazowski, Calot, Chartier, Cochran, Dixon and Masey, Ehrenfeld and Littauer, Fourgeaud, Fraser, Hansen Hurwitz and Madow, Morice, Owen, E. Pearson, Tippet, Wilks ...).

La notation  $s$  intervenant dans de nombreuses formules de tests couramment utilisés pour l'estimation et la comparaison de moyennes (test  $t$  de Student, test de Behrens-Fisher, test de Welch, ...), et pour l'estimation et la comparaison de variances (test  $\chi^2$ , test  $F$  de Fisher, test de Neyman et Pearson, test de Cochran ...), il est regrettable que le praticien, non toujours bien familiarisé avec la littérature statistique théorique, soit obligé, avant d'utiliser telle ou telle formule de rechercher, dans l'ouvrage auquel il s'adresse, la signification de symbole  $s$ .

Un effort de normalisation dans ce domaine semble indispensable : en attendant qu'il se réalise - s'il se réalise - il semble que l'on puisse demander aux auteurs d'adjoindre à leurs ouvrages un glossaire des symboles utilisés avec une définition et une formulation précises.

Le Secrétariat français (AFNOR) du Comité Technique de l'Organisation Internationale de normalisation (Comité ISO TC/69 : Procédés statistiques d'interprétation des séries d'observations) a, en particulier, lors de la dernière réunion internationale (Paris, Juin 1966) proposé de normaliser le symbole

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x - \bar{x})^2$$

avec la définition "Estimateur de la variance  $\sigma^2$ ". Cette proposition a été adoptée et figurera ultérieurement dans les recommandations I. S. O.

Remarque - Dans le problème qui vient d'être envisagé, on a négligé de tenir compte de l'élément intermédiaire : population réelle finie d'effectif  $N$ , entre l'échantillon d'effectif  $n$  et la population conceptuelle  $P$ .

Certains auteurs, en particulier des spécialistes de la théorie des sondages, tirages sans remise, dans une population finie d'effectif  $N$ , considérée pour elle-même, utilisent les notations suivantes :

$$\begin{aligned} \sigma^2 &= \frac{1}{N} \sum_{j=1}^N (X^j - \bar{X})^2 \\ S^2 &= \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2 \\ s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned}$$

$X_j$  ( $j = 1, 2, \dots, N$ ) étant un élément quelconque de la population, de moyenne  $\bar{X}$  et  $x_i$  ( $i = 1, 2, \dots, n$ ) un élément quelconque de l'échantillon de moyenne  $\bar{x}$  (voir par exemple Hansen, Hurwitz et Madow "Sample survey Methods and Theory).

Ces formules conduisent alors à

$$E[s^2] = S^2$$

Alors que dans les tirages avec remise on avait

$$E(s^2) = \sigma^2$$

Le rapport  $\frac{N}{N-1}$  étant en général très voisin de l'unité, cette distinction entre  $S^2$  et  $\sigma^2$  est, le plus souvent sans intérêt pratique.

## 2 - FRACTILES (ou QUANTILES)

Pour une variable aléatoire continue  $X$ , le fractile  $x_p$  est défini tantôt par :

$$P = \Pr [X < x_p] \quad (3)$$

tantôt par :

$$P = \Pr [X > x_p] \quad (4)$$

(les signes  $<$  et  $>$  pouvant évidemment, dans ce cas, être remplacés par  $\leq$  ou  $\geq$ , sans modifier la signification.

La définition (3) paraît être la plus fréquemment utilisée par les auteurs d'ouvrages de statistique théorique ou appliquée (Cramer, Kendall, Hald, Brownlee, Dixon and Masey, Johnson and Leone, Calot, Littauer and Ehrenfeld, Centre de Formation aux Applications Industrielles de la Statistique - Paris).

Par contre des auteurs américains (Mood) et plus particulièrement les spécialistes de la fiabilité (Epstein, Bazowski, Lloyd and Lipow) ont, dans des ouvrages récents sur la fiabilité, utilisé des définitions du type (4), qui se présente comme la mise en formule de la terminologie de R. A. Fisher : " $x_p$  is the upper P percentage point of the X' distribution".

La traduction littérale de cette terminologie est parfois utilisée par des auteurs français.

Ainsi on trouvera l'expression  $\chi^2 = 5,99$  est le "point 5 % de la distribution de  $\chi^2$  pour 2 degrés de liberté, c'est-à-dire :

$$0,05 = \Pr(\chi^2 > 5,99) \quad \text{pour } \nu = 2 \text{ d. d. l.}$$

Les deux points de vue ci-dessus peuvent évidemment être envisagés, mais il semble plus logique d'adopter la définition (1), qui associe des variations de même sens pour  $P$  et  $x_p$  et qui a, de plus, le mérite d'être d'accord avec la terminologie utilisée pour les quartiles : il n'est jamais venu à l'idée de qui que ce soit d'appeler premier quartile (un quart) ; celui qui avait une probabilité  $P = 1/4$  d'être dépassé.

Dans le cas d'une variable *discrète*, il convient de plus, pour rester en accord avec la définition ci-après de la fonction cumulative, d'utiliser la définition (3) modifiée sous la forme

$$P = \Pr[X \leq x_p] \quad (5)$$

qui reste valable aussi bien pour une variable continue que pour une variable discrète.

Dans ce cas la définition (3) n'a de sens précis, d'un point de vue concret, que si  $x_p$  est une des valeurs possibles de la variable discrète; cette définition avec le signe  $\leq$  est d'accord avec la présentation habituelle des tables (lois binomiale, Poisson) qui donnent la valeur de  $P = \Pr[X \leq x_p]$  en fonction des valeurs possibles  $X = x_p$ .

En ce qui concerne les tables statistiques, on retrouve l'inconvénient signalé ci-dessus : certaines tables donnent, aussi bien pour une variable continue  $X$ , la valeur  $x_p$  telle que l'on ait  $P = \Pr(X \leq x_p)$ , par exemple :

HALD - "Statistical Tables and Formulas", Wiley (1952), pour toutes les variables qui y figurent (  $u$  ,  $t$  ,  $\chi^2$  ,  $F$  et  $\frac{W}{\sigma}$  )

ROMIG - 50-100 Binomial Tables, Wiley (1952), pour la variable binomiale  $\mathcal{B}(n, p)$  pour  $50 < n < 100$ .

MOLINA - Poisson's exponential limits. Van Nostrand (1947) (pour la loi de Poisson).

CENTRE DE FORMATION AUX APPLICATIONS INDUSTRIELLES DE LA STATISTIQUE - Pour les tables des lois : Normale, Binomiale et de Poisson.

- D'autres tables utilisent la correspondance entre  $P$  et  $x_p$ , définie par  $P = \Pr(X \geq x_p)$  par exemple.

FISHER and YATES - Statistical tables for biological, agricultural and medical research, Hafner (1953) pour les variables  $\chi^2$  et  $F$ .

A. M. S. 8 - Tables of the binomial probability distribution. National Bureau of Standards (1947).

CENTRE DE FORMATION POUR LES VARIABLES  $\chi^2$  ET  $F$ .

Enfin pour des lois symétriques, par exemple la loi de Student, certaines tables associent à une valeur  $x$  de la variable, la probabilité que  $X$  soit extérieure à l'intervalle  $(-x, +x)$ .

Evidemment, le choix a été fait en fonction de ce que les auteurs considéraient comme devant être l'usage le plus fréquent de la table : il n'en résulte pas de sérieuses difficultés, mais un simple effort d'attention. On peut se demander cependant si une présentation un peu plus homogène ne simplifierait pas le travail des utilisateurs plus ou moins "conditionnés" par l'emploi d'une certaine table et qui, pour un problème particulier, sont amenés à utiliser une table de présentation différente.

### 3 - DISTRIBUTION CUMULATIVE - FONCTION CUMULATIVE

Là encore, on trouve le plus généralement la définition

$$F(x) = \Pr(X \leq x) \quad (6)$$

aussi bien pour une variable discrète soit :

$$F(x_k) = \sum_{i=1}^k p_i = p_1 + p_2 + \dots + p_k ,$$

que pour une variable continue, mais divers auteurs d'ouvrages théoriques de Calcul des probabilités (Chapelon, Calot, Fourgeaud, Girault) utilisent la définition générale

$$F(x) = \Pr(X < x) \quad (7)$$

soit, pour une variable discrète :

$$F(x_k) = \sum_{i=1}^{k-1} p_i = p_1 + p_2 + \dots + p_{k-1}$$

Ces deux définitions (6) et (7), identiques dans le cas d'une variable continue, peuvent devenir très différentes dans le cas d'une variable discrète, en particulier pour les petites valeurs de celle-ci.

Ainsi, pour une loi de Poisson de moyenne égale à 1, on a :

$$F(1) = 0,7358 \quad \text{définition (6)}$$

$$F(1) = 0,3679 \quad \text{définition (7)}$$

Il paraît préférable d'utiliser la définition (6) qui est celle qui correspond à la présentation habituelle des tables de lois d'une variable discrète qui, généralement, donnent la probabilité cumulée *jusqu'à et y<sub>1</sub> compris* la valeur *x*.

Pour ces raisons, le Secrétariat français du Comité technique ISO/TC. 69 de l'Organisation internationale de normalisation, a proposé lors de sa dernière réunion à Paris, Juin 1966, les définitions suivantes :

- Fractile  $x_p$ , défini dans tous les cas par :

$$P = \Pr[X \leq x_p]$$

- Fonction de répartition (ou distribution cumulative) définie dans tous les cas par :

$$F(x) = \Pr[X \leq x]$$

Ces deux définitions sont d'ailleurs résumées en une seule par la relation

$$F(x_p) = P = \Pr[X \leq x_p]$$

Bien que ces définitions aient été accueillies favorablement par la grande majorité des pays présents à la réunion ou ayant fait part de leurs observations sur l'ensemble du projet de normalisation présenté par le secrétariat français, le Comité ISO/TC. 69 a estimé qu'il n'y avait pas lieu, pour l'instant, de les inclure dans une norme internationale.

(Peut-on ajouter, en manifestant quelque regret, d'une telle lenteur, que ce Comité international, depuis sa première réunion (Genève 1951) n'a normalisé qu'une demi-douzaine de symboles qui, à part celui relatif à l'estimateur de la variance  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ , étaient ceux qui ne donnaient lieu dans les documents existants, à aucune équivoque pour leur interprétation).

#### 4 - INTERVALLES DE CONFIANCE ET TESTS D'HYPOTHESE

Des problèmes analogues se posent, concernant la définition de l'intervalle de confiance ou de l'intervalle d'acceptation d'une hypothèse relativement au paramètre d'une loi de probabilité d'une variable discrète (loi binomiale, loi de Poisson...).

Considérons par exemple l'intervalle de confiance du paramètre  $p$  d'une loi binomiale, défini à partir d'un échantillon de  $n$  éléments ayant donné lieu à l'observation de  $c$  d'entre eux possédant la propriété correspondante à la proportion  $p$  dans la population.

Si  $1 - \alpha = 1 - \alpha_1 - \alpha_2$  est le coefficient de confiance associé à un intervalle  $p'$ ,  $p''$  ( $p' < p''$ ) les limites de l'intervalle de confiance sont définies par la plupart des auteurs par les conditions

$$I \quad \left\{ \begin{array}{l} \alpha_1 = \sum_{x=0}^c C_n^x p'' (1-p'')^{n-x} = \Pr[x \leq c \mid p = p''] \\ \alpha_2 = \sum_{x=c}^n C_n^x p' (1-p')^{n-x} = \Pr[x > c \mid p = p'] \end{array} \right.$$

Ces relations sont d'ailleurs simplement la formulation de la question : quelles sont les valeurs extrêmes  $p'$  et  $p''$  que l'on peut encore considérer comme compatibles avec le résultat de l'expérience, aux risques  $\alpha_1$  et  $\alpha_2$  respectivement associés aux limites  $p''$  et  $p'$ .

Cependant j'ai rencontré récemment d'une part chez un auteur américain, le système de conditions :

$$II \quad \left\{ \begin{array}{l} \alpha_1 = \sum_{x=0}^{c-1} (n, p'') \\ 1 - \alpha_2 = \sum_{x=0}^{c-1} (n, p') \end{array} \right.$$

avec des abaques construits à partir de ces définitions

et d'autre part, dans une grande entreprise française des abaques construits avec le système de conditions

$$III \quad \left\{ \begin{array}{l} \alpha_1 = \sum_{x=0}^c (n, p'') \\ 1 - \alpha_2 = \sum_{x=0}^c (n, p') \end{array} \right.$$

systèmes qui diffèrent par l'une ou l'autre des deux conditions du système (I) qui peut s'écrire aussi

$$I \quad \left\{ \begin{array}{l} \alpha_1 = \sum_{x=0}^c (n, p'') \\ 1 - \alpha = \sum_{x=0}^{c-1} (n, p') \end{array} \right.$$

Evidemment, ces divergences n'ont aucune importance, si dans le texte ou dans la note explicative de l'abaque, on a bien précisé quelle était la définition de l'intervalle de confiance qui avait été adoptée.

Malheureusement les abaques publiés dans un ouvrage, dans un article ou utilisés dans un bureau ne portent que bien rarement de telles indications et risquent d'être utilisés avec une signification autre que celle qui leur convient.

Si  $n$  et  $c$  ne sont pas petits, l'erreur qui résulte d'une interprétation erronée est pratiquement sans importance.

Pour  $n = 200$ ,  $c = 10$  par exemple, l'intervalle de confiance à 0,95 ( $\alpha_1 = \alpha_2 = 0,025$ ), correspondant aux 3 systèmes précédents est suivant la définition utilisée :

I	$p' = 0,024$	$p'' = 0,089$
II	$p' = 0,024$	$p'' = 0,083$
III	$p' = 0,027$	$p'' = 0,089$

Mais il n'en va pas de même dans le cas de petits échantillons utilisés de plus en plus fréquemment pour des raisons d'économie, dans des essais destructifs, par exemple essais de durée de vie dans les études de fiabilité.

Ainsi pour  $n = 20$ ,  $c = 5$ , les trois définitions précédentes donneraient :

I	$p' = 0,086$	$p'' = 0,490$
II	$p' = 0,086$	$p'' = 0,436$
III	$p' = 0,119$	$p'' = 0,490$

On peut remarquer que l'emploi par certains techniciens de la définition (III) n'est peut être pas sans rapport avec les définitions habituellement adoptées pour un plan d'échantillonnage simple ( $\alpha, p_1, \beta, p_2$ ) dans lesquels l'effectif  $n$  et le critère d'acceptation  $c$  sont définis théoriquement par

$$IV \quad \left\{ \begin{array}{l} \sum_0^c C_n^x p_1^x (1 - p_1)^{n-x} = 1 - \alpha \\ \sum_0^c C_n^x p_2^x (1 - p_2)^{n-x} = \beta \end{array} \right.$$

A cet égard, il convient aussi de remarquer que  $n$  et  $c$  devant être entiers, ces conditions ne peuvent être satisfaites que d'une manière approximative, avec deux valeurs différant d'une unité pour  $c$ , les valeurs de  $n$  qui leur correspondent étant choisies de manière que les risques réels  $\alpha'$  et  $\beta'$  qui en résultent soient

$$\begin{array}{lll} \alpha' < \alpha & \text{et} & \beta' > \beta & \text{si } c \text{ est pris par excès} \\ \alpha' > \alpha & \text{et} & \beta' < \beta & \text{si } c \text{ est pris par défaut} \end{array}$$

Il en résulte que certains tableaux de plans tout préparés à l'avance pour la commodité des utilisateurs en correspondent que d'assez loin aux indications qu'ils contiennent lorsque  $n$  et  $c$  sont petits.

Ainsi, par exemple, le document H. 108 "Quality Control and Reliability Handbook" (Assistant Secretary of Defense U. S. A.), donne une série de plans pour des essais de durée de vie (essais de durée préfixée  $T$  sans remplacement des défectueux). Ces plans sont définis théoriquement par les conditions :

$$\sum_{x=0}^c C_n^x (1 - e^{-T/\theta_0})^x (e^{-T/\theta_0})^{n-x} = 1 - \alpha$$

$$\sum_{x=0}^c C_n^x (1 - e^{-T/\theta_1})^x (e^{-T/\theta_1})^{n-x} = \beta ,$$

l'acceptation ayant lieu, si le nombre des défectueux est inférieur à  $r = c + 1$ , avec des probabilités respectives  $1 - \alpha$  ou  $\beta$  si la durée de vie moyenne dans le lot est égale à  $\theta_0$  ou  $\theta_1$ .

Le tableau donne divers plans en fonction de valeurs particulières de  $\alpha$ ,  $\beta$ ,  $\frac{T}{\theta_0}$  et  $\frac{\theta_1}{\theta_0}$  ( $\theta_1 < \theta_0$ ).

Si pour  $\alpha = \beta = 0,01$ ,  $T/\theta_0 = \frac{1}{20}$ ,  $\theta_1/\theta_0 = 2/3$ , le plan proposé  $n = 2275$ ,  $r = 136$ , donne bien lieu à des risques très voisins des valeurs théoriques de  $\alpha$  et  $\beta$ , on y trouve aussi des plans avec de très petites valeurs de  $n$  et  $r$ .

Par exemple pour  $\alpha = \beta = 0,25$ ,  $T/\theta_0 = 1/3$ ,  $\theta_1/\theta_0 = 1/10$ , le tableau propose le plan  $n = 1$ ,  $r = 1(c = 0)$ .

Les risques réels, correspondant à ce plan sont  $\alpha' \approx 0,28$  et  $\beta' \approx 0,04$  au lieu de  $0,25$  et  $0,25$ . Or nulle part le document H. 108 n'attire l'attention de l'utilisateur sur les variations des risques associés aux informations en  $n$  et  $r$  données par ce tableau : peut-être a-t-on pensé que le bon sens y suffirait :

Des remarques analogues pourraient être faites à propos d'autres tableaux publiés dans H. 108.

## 5 - TESTS D'HYPOTHESES

Relativement à une distribution binomiale, par exemple, le problème se pose de la manière suivante : étant donné un échantillon de  $n$  observations pouvant donner lieu à  $x$  éléments de la catégorie ( $p$ ), quel est l'intervalle ( $k'$ ,  $k''$ ) tel que si on accepte l'hypothèse  $p = p_0$  lorsque  $k' \leq x \leq k''$ , le risque de rejeter une hypothèse vraie soit  $\alpha = \alpha' + \alpha''$  (intervalle bilatéral d'acceptation).

Posé sous la forme

$$\left\{ \begin{array}{l} \alpha' = \Pr[x < k' \mid p = p_0] \\ \alpha'' = \Pr[x > k'' \mid p = p_0] \end{array} \right.$$

le problème est insoluble,  $k'$  et  $k''$  devant être entiers.

Pour chaque limite on trouvera en réalité deux entiers consécutifs  $k'$  et  $k' + 1$ , d'une part,  $k''$  et  $k'' - 1$  d'autre part tels que

$$\Pr[x < k'] = \alpha'_1 < \alpha'$$

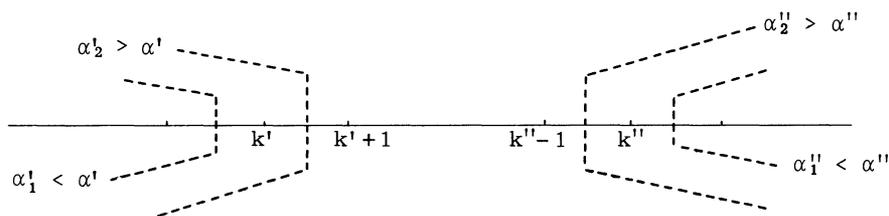
$$\Pr[x < k' + 1] = \alpha'_2 > \alpha'$$

et

$$\Pr[x > k''] = \alpha''_1 < \alpha''$$

$$\Pr[x > k'' - 1] = \alpha''_2 > \alpha''$$

l'égalité des probabilités à  $\alpha'$  et  $\alpha''$  étant pratiquement impossible



Un problème d'approximation se pose, c'est-à-dire un problème de choix entre les 4 combinaisons

$$(k' , k'') - (k' , k'' - 1) - (k' + 1 , k'') - (k' + 1 , k'' - 1)$$

définies par les inégalités ci-dessus.

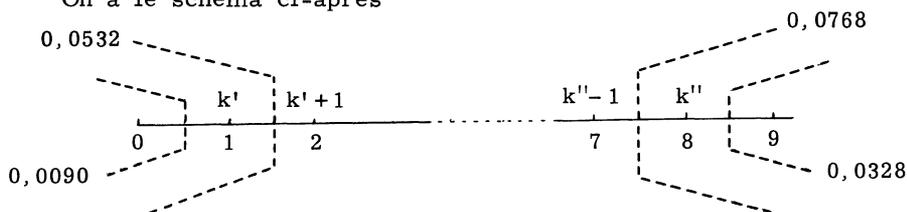
Certains auteurs indiquent comme règle systématique le choix de l'intervalle  $(k' , k'')$  pour lequel on a  $\alpha'_1 < \alpha'$ ,  $\alpha''_1 < \alpha''$ , donc un risque total  $\alpha'_1 + \alpha''_1 < \alpha$ .

Cette solution est non seulement celle pour laquelle le risque total de rejeter l'hypothèse  $p = p_0$  - si elle est vraie - est au plus égal à  $\alpha$ , mais celle pour laquelle les risques unilatéraux sont au plus égaux aux valeurs théoriques fixées  $\alpha'$  et  $\alpha''$ . Cette exigence supplémentaire correspond-elle vraiment aux soucis de l'expérimentateur ? Si, par exemple,  $\alpha'_1$  est beaucoup plus petit que  $\alpha'$ , alors que  $\alpha'_2 > \alpha'$  est presque égal à  $\alpha'$ .

Un tel problème ne se posant pratiquement avec des différences appréciables entre les diverses solutions, que lorsque  $n$  est petit, il semble donc que chaque cas particulier doive être résolu en fonction des informations données par les tables dans lesquelles on trouvera tous les éléments du schéma ci-dessus.

**Exemple** : Test de l'hypothèse  $p_0 = 0,08$  à partir d'un échantillon de  $n = 50$ , au risque  $\alpha = 0,05 + 0,05 = 0,10$ .

On a le schéma ci-après



qui donne les résultats suivants :

<u>Acceptation</u>	<u>risque total</u>
$1 \leq x \leq 8$	$0,0090 + 0,0328 = 0,0418$
$1 \leq x \leq 7$	$0,0090 + 0,0768 = 0,0858$
$2 \leq x \leq 8$	$0,0532 + 0,0328 = 0,0860$
$2 \leq x \leq 7$	$0,0532 + 0,0768 = 0,1300$

Des 4 solutions ci-dessus, trois donnent un risque total inférieur à la valeur fixée  $\alpha = 0,10$  : la première n'est pas nécessairement celle qui correspond au mieux aux risques de la décision prise, risques qui comprennent aussi bien celui d'accepter une hypothèse fausse que de refuser une hypothèse exacte (voir ci-après les courbes puissance du test correspondant aux diverses solutions).

Remarque - Les mêmes problèmes (intervalles de confiance et tests d'hypothèse) se posent évidemment à propos de la loi de Poisson.

