

Zeitschrift: Archives des sciences et compte rendu des séances de la Société
Band: 48 (1995)
Heft: 2: Archives des Sciences

Artikel: Propriétés de quelques test d'ajustement pour données groupées : un chapitre choisi de la méthodologie statistique
Autor: Driancourt, Catherine / Streit, Franz
DOI: <https://doi.org/10.5169/seals-740255>

Nutzungsbedingungen

Die ETH-Bibliothek ist die Anbieterin der digitalisierten Zeitschriften. Sie besitzt keine Urheberrechte an den Zeitschriften und ist nicht verantwortlich für deren Inhalte. Die Rechte liegen in der Regel bei den Herausgebern beziehungsweise den externen Rechteinhabern. [Siehe Rechtliche Hinweise.](#)

Conditions d'utilisation

L'ETH Library est le fournisseur des revues numérisées. Elle ne détient aucun droit d'auteur sur les revues et n'est pas responsable de leur contenu. En règle générale, les droits sont détenus par les éditeurs ou les détenteurs de droits externes. [Voir Informations légales.](#)

Terms of use

The ETH Library is the provider of the digitised journals. It does not own any copyrights to the journals and is not responsible for their content. The rights usually lie with the publishers or the external rights holders. [See Legal notice.](#)

Download PDF: 01.11.2024

ETH-Bibliothek Zürich, E-Periodica, <https://www.e-periodica.ch>

Archs Sci. Genève	Vol. 48	Fasc. 2	pp. 185-196	Septembre 1995
-------------------	---------	---------	-------------	----------------

Communication présentée à la séance du 2 février 1995

**PROPRIÉTÉS DE QUELQUES TESTS D'AJUSTEMENT
POUR DONNÉES GROUPEES.
UN CHAPITRE CHOISI DE LA MÉTHODOLOGIE STATISTIQUE**

PAR

Catherine DRIANCOURT* & Franz STREIT**

ABSTRACT

Properties of some goodness-of-fit tests for grouped data. - After a short review of the aims of statistics and of some notions concerning statistical tests and in particular goodness-of-fit tests we propose a general class of such goodness-of-fit tests and show an optimality property of these tests. We determine the mean, the variance and the asymptotic distributions of their test statistics. We study special cases in calculating the actual level of significance and the power function and by indicating the alternative hypotheses against which the tests are unbiased.

§ 1. Types d'expériences statistiques et leurs buts

La statistique s'occupe des phénomènes aléatoires, et plus particulièrement du cas où on analyse de tels phénomènes aléatoires, en observant les valeurs que prennent des variables aléatoires (v.a.) liées à ce phénomène aléatoire. Observer une telle v.a. équivaut à effectuer une expérience à issue incertaine, qui est liée à ce phénomène aléatoire et dont le résultat est un nombre réel. Si l'on choisit habilement les v.a. ainsi observées, on arrive à décrire beaucoup de phénomènes aléatoires, même des phénomènes aléatoires assez complexes et des phénomènes aléatoires à composantes qualitatives (cela, en utilisant des codages appropriés).

Les analyses statistiques se distinguent par leur niveau de complexité:

- 1) Il se peut que l'on veuille uniquement donner une description succincte des observations que l'on a prises du phénomène aléatoire en question. Dans ce cas, on fait de l'analyse des données d'une situation expérimentale concrète particulière ou de la statistique descriptive. On résume les observations en regroupant les données dans des tableaux, en visualisant l'ensemble de ces données graphiquement et en caractérisant les aspects typiques et la structure de cet ensemble par des indices ou

* Ecole Supérieure de Commerce, 14, avenue Trembley, case postale 118, CH-1211 Genève 28.

** Section de Mathématiques, 2-4, rue du Lièvre, case postale 240, CH-1211 Genève 24.

coefficients. On arrive ainsi à un abstrait précis des résultats des expériences faites dans le cas concret analysé.

- 2) Pour les besoins de la recherche scientifique actuelle, des analyses statistiques du type statistique descriptive ne suffisent d'habitude pas, parce qu'elles expriment uniquement des résultats de cas particuliers. Sans considérations additionnelles qui permettent de modéliser le phénomène aléatoire mathématiquement, et de cette façon d'expliquer le mécanisme aléatoire, on n'arrive pas à des résultats applicables à tous les phénomènes aléatoires semblables à celui que l'on a examiné. C'est là où les modèles stochastiques de la théorie des probabilités interviennent. Ces modèles résument les tendances systématiques que l'on pense reconnaître dans les données observées et les combinent d'une façon précise avec des éléments purement stochastiques, qui représentent les fluctuations aléatoires (les déviations non-prédictibles des données individuelles de ces tendances). Chacun de ces modèles représente une interprétation possible des données, un modèle idéalisé de la réalité. En se demandant lesquels des modèles stochastiques sont aptes à expliquer un phénomène aléatoire, on va a priori prendre en considération plusieurs modèles susceptibles de décrire le mécanisme aléatoire. Il faut alors effectuer un choix parmi ces possibilités. En utilisant les méthodes de la statistique inférentielle ou inductive, on arrive à faire ce choix d'une façon aussi optimale que possible (du point de vue minimisation des effets néfastes d'un choix incorrect).

§ 2. Tests d'ajustement

Ci-après, on ne va pas considérer des phénomènes aléatoires du type échantillon aléatoire, c'est-à-dire on considère n v.a. X_1, \dots, X_n et on suppose qu'elles sont engendrées par répétition d'une expérience aléatoire du même type, ainsi que les chances de prendre les différentes valeurs sont les mêmes pour chaque X_j et que les X_j ne s'influencent pas mutuellement. Il suffit donc d'indiquer le modèle probabiliste commun pour tous les X_j . Cela se fait par spécification de la fonction de répartition F_X , la même pour tous les X_j [$j = 1, \dots, n$]. F_X est la fonction qui indique pour chaque valeur réelle x la probabilité de l'événement $X_j \leq x$, donc $F_X(x) = P(X \leq x)$. A l'aide de la fonction de répartition, on peut calculer la probabilité de n'importe quel événement aléatoire (donc n'importe quelle condition probabilisable) se rapportant à X_j . On a que

(*) $F_X(-\infty) = 0$, $F_X(+\infty) = 1$ et que F_X est non-décroissante et continue à droite.

Il reste donc un grand choix de modèles stochastiques a priori possibles. Chaque fonction F_X satisfaisant aux conditions (*) en représente une possibilité. Comme les X_j ne s'influencent pas mutuellement, on a

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = F_X(x_1) \cdot \dots \cdot F_X(x_n).$$

Le but de l'inférence statistique dans ce contexte est de connaître la vraie fonction de répartition F_X .

Parfois la situation expérimentale suggère un choix privilégié F_0^* de F_X , qui est justifié par une théorie potentiellement correcte. La question à résoudre prend alors la forme suivante, caractéristique pour les tests d'ajustement : Décider laquelle des hypothèses

$$H_0^{**} : F_X \equiv F_0^* \quad (\text{hypothèse nulle})$$

$$H_1^{**} : F_X \neq F_0^* \quad (\text{Hypothèse alternative})$$

est correcte. On peut montrer que tester H_0^{**} versus H_1^{**} revient pour F_0^* continue à la même chose que de tester

$$H_0^* : F_X(x) = F_0(x) := x \quad [0 \leq x \leq 1] \quad \text{versus}$$

$$H_1^* : F_X \neq F_0$$

pour les variables transformées $F_0^*(X_1), \dots, F_0^*(X_n)$. (Witting et al., 1970). En ce qui suit, nous allons donc supposer sans perte de généralité que la fonction de répartition de X_1, \dots, X_n à approuver ou rejeter est F_0 , comme définie sous H_0^* (version standard du test d'ajustement). Il est alors raisonnable de limiter le choix de l'alternative à des fonctions de répartition satisfaisant $F_X(0) = 0$ et $F_X(1) = 1$.

La situation de test décrite jusqu'à présent se réfère au cas où l'on dispose de données observées exactes à plusieurs décimales. Mais souvent les données que l'on obtient sont déjà groupées, c'est-à-dire est seulement indiqué pour chaque observation l'intervalle dans lequel elle se trouve, et non pas sa valeur précise. On analyse dans la suite la subdivision de $(0, 1]$ dans k sous-intervalles de longueurs égales, en posant $I_i = (k^{-1}(i-1), k^{-1}i]$ pour le i -ème sous-intervalle.

On écrira

$$p_i = P(X_j \in I_i) = F_X(k^{-1}i) - F_X(k^{-1}(i-1))$$

$$[i = 1, \dots, k, j = 1, \dots, n]$$

pour la probabilité que la j -ème observation tombe dans le i -ème intervalle et

$N_i = N_i(n, k)$ pour le nombre des données observées à valeur dans I_i .

Pour cette manière de rapporter les résultats de l'observation, le problème du test d'ajustement se modifie un peu, dû au fait qu'il est impossible de contrôler les valeurs

de F_X pour tout argument x . Seules les hypothèses concernant les probabilités p_i attribuées par F_X aux intervalles I_i peuvent être contrôlées à cause du regroupement des données. Le problème du test d'ajustement pour données groupées dans sa version standard se réfère donc aux hypothèses

$$H_0 : p_i = k^{-1} \quad [i = 1, \dots, k]$$

donc $\underline{p} = (p_1, \dots, p_k)' = k^{-1}\underline{1}$ où $\underline{1} = (1, \dots, 1)'$ e \mathbb{R}^k et $'$ désigne la transposition d'un vecteur et

$$H_1 : \text{il existe un } i_0 \text{ tel que } p_{i_0} \neq k^{-1} [i_0 \in \{1, \dots, k\}].$$

Un test est une règle de décision qui nous indique pour chaque réalisation de l'échantillon aléatoire, donc pour chaque n -tuple de nombres réels admissibles comme résultat s'il faut accepter H_0 ou H_1 . D'habitude on spécifie un test en spécifiant une statistique (donc une fonction t de X_1, \dots, X_n qui prend des valeurs réelles et pour laquelle on est capable de déterminer la fonction de répartition $F_t(X_1, \dots, X_n)$ à partir de F_X) et on spécifie une valeur (critique) c . Dans la version standard du test, on va comparer la valeur obtenue pour t , à savoir $t(x_1, \dots, x_n)$, à c et rejeter H_0 en faveur de H_1 si $t(x_1, \dots, x_n)$ dépasse c . Naturellement la qualité du test dépend du choix de t et c . Cette qualité peut être mesurée par deux types de probabilités : La probabilité de commettre une erreur de première espèce en acceptant H_1 quand effectivement H_0 est correcte et la probabilité de commettre une erreur de seconde espèce en acceptant H_0 quand effectivement H_1 est correcte. La dernière probabilité mentionnée dépend du vecteur $\underline{p} = (p_1, \dots, p_k)'$. On a $\alpha_{eff} = P(t(X_1, \dots, X_n) > c : k^{-1}\underline{1})$ (probabilité d'accepter H_1 quand H_0 est correcte; α_{eff} est appelée niveau de signification effectif du test) et plus généralement on écrit $P_{ui}((p_1, \dots, p_k)') = P(t(X_1, \dots, X_n) > c : (p_1, \dots, p_k)')$ pour la probabilité d'accepter H_1 si $\underline{p} = (p_1, \dots, p_k)'$ (P_{ui} est appelée fonction de puissance). Un test (uniformément) optimal va avoir pour α prescrit et proche de 0 – on appelle α le niveau de signification (nominal) – l'inégalité $\alpha_{eff} \leq \alpha$ satisfaite et sous cette restriction les valeurs de sa fonction de puissance vont être aussi élevées que possible pour $\underline{p} \neq k^{-1}\underline{1}$. Dans notre problème modifié du test d'ajustement où $\underline{X} = (X_1, \dots, X_n)'$ est réduit à $\underline{N} = (N_1, \dots, N_k)'$ dû au regroupement des observations, le résultat fondamental de la théorie des tests (Lehmann, 1986, Chapitre 3) nous indique qu'un procédé optimal pour tester H_0 versus H_1 (\underline{p}^*) : $\underline{p} = \underline{p}^*$, où \underline{p}^* est fixé, existe; comme statistique du test de ce procédé, on peut choisir $\tilde{t}(\underline{N})$ aux valeurs

$$\begin{aligned} t(\underline{x}) = \tilde{t}(\underline{n}) &= P(N_1 = n_1, \dots, N_k = n_k : \underline{p}^*) [P(N_1 = n_1, \dots, N_k = n_k : \underline{p} = k^{-1}\underline{1})]^{-1} \\ &= k^n p_1^{*n_1} \cdot \dots \cdot p_k^{*n_k} \end{aligned}$$

$\tilde{t}(\underline{N})$ suit une distribution discrète et pour une borne α prescrit $c = c(\alpha)$ doit être choisi de telle façon que $P(t(X_1, \dots, X_n) \geq c : k^{-1}\underline{1}) \leq \alpha$ et que cette inégalité soit aussi proche que possible d'une égalité.

Comme $\tilde{t}(N)$ et $\ln(k^{-n}\tilde{t}(N))$ sont des statistiques équivalentes, on peut – en ajustant c – exprimer le test de manière équivalente à l'aide de la statistique

$$N_1 \ln(p_1^*) + \dots + N_k \ln(p_k^*).$$

§ 3. Quelques propriétés d'une classe de tests

On est ainsi amené à considérer des tests basés sur l'utilisation des statistiques

$$V_{\tilde{c}}(n, k) = \tilde{V}(N_1, \dots, N_k) = \sum_{i=1}^k c_i N_i(n, k),$$

où c_1, \dots, c_k sont des constantes non-positives réelles. Nous allons déduire quelques propriétés de $V_{\tilde{c}}(n, k)$.

A.) Calcul de l'espérance mathématique et de la variance de $V_{\tilde{c}}(n, k)$.

A cette fin, il est important de constater que $V_{\tilde{c}}(n, k)$ admet une représentation comme somme de contributions indépendantes provenant des variables aléatoires X_j individuelles. On a

$$V_{\tilde{c}}(n, k) = \sum_{j=1}^n V_{X_j}(1, k) \quad \text{où}$$

$$V_{X_j}(1, k) = \sum_{i=1}^k c_i 1_{\{I_i\}}(X_j) \quad \text{où } 1_{I_i}(u) = 1 \text{ pour } u \in I_i \text{ et } 0 \text{ pour } u \notin I_i.$$

Pour tout $p \in \mathbb{R}^k$ tel que $p_i \geq 0$ et $p_1 + \dots + p_k = 1$, on trouve

$$E[V_{X_j}(1, k) : p] = \sum_{i=1}^k c_i p_i$$

et

$$\text{Cov}[V_{X_{j_1}}(1, k), V_{X_{j_2}}(1, k) : p] = 0 \quad \text{si } j_1 \neq j_2$$

(cela est dû à l'indépendance stochastique de X_{j_1} et X_{j_2} qui entraîne celle des variables $V_{X_{j_1}}(1, k)$ et $V_{X_{j_2}}(1, k)$) et d'autre part

$$\text{Var}[V_{X_j}(1, k) : p] = \sum_{i=1}^k c_i^2 \text{Var}[1_{I_i}(X_j) : p]$$

$$\begin{aligned}
& + \sum_{i_1=1}^k \sum_{i_2=1}^k c_{i_1} c_{i_2} \text{Cov}[1_{I_{i_1}}(X_j), 1_{I_{i_2}}(X_j) : \underline{p}] \\
& \quad i_1 \neq i_2 \\
& = \sum_{i=1}^k c_i^2 p_i (1 - p_i) - \sum_{i_1=1}^k \sum_{i_2=1}^k c_{i_1} c_{i_2} p_{i_1} p_{i_2} \\
& \quad \quad \quad i_1 \neq i_2 \\
& = \sum_{i=1}^k c_i^2 p_i - \sum_{i_1=1}^k \sum_{i_2=1}^k c_{i_1} c_{i_2} p_{i_1} p_{i_2} \\
& = (\sum_{i=1}^k c_i^2 p_i - (\sum_{i=1}^k c_i p_i)^2).
\end{aligned}$$

Notons que l'on obtient ce résultat également à partir de la formule développée pour la variance, à savoir en prenant comme point de départ la relation

$\text{Var}[V_{X_j}(1, k) : \underline{p}] = E[V_{X_j}(1, k)^2 : \underline{p}] - (E[V_{X_j}(1, k) : \underline{p}])^2$. A l'aide de ces considérations, on obtient donc

$$E[V_{\underline{C}}(n, k) : \underline{p}] = \sum_{i=1}^n E[V_{X_i}(1, k) : \underline{p}] = n \sum_{i=1}^k c_i p_i.$$

$$\text{Var}[V_{\underline{C}}(n, k) : \underline{p}] = \sum_{i=1}^n \text{Var}[V_{X_i}(1, k) : \underline{p}] = n (\sum_{i=1}^k c_i^2 p_i - (\sum_{i=1}^k c_i p_i)^2).$$

B.) Calcul de la distribution asymptotique de la $V_{\underline{C}}(n, k)$.

Nous avons constaté sous A.) que $V_{\underline{C}}(n, k)$ admet la représentation comme somme de v.a. indépendantes et identiquement distribuées à variance positive, pourvu qu'au moins deux composantes de \underline{p} soient positives et que les coefficients c_i correspondants soient différents et différents de 0. Sous cette condition, on peut donc appliquer le théorème limite central sous sa forme la plus classique (Chung, 1974) et on trouve pour le cas où la taille n de l'échantillon aléatoire tend vers ∞ que

$$V_{\underline{C}}(n, k) \xrightarrow{L} N(n \sum_{i=1}^k c_i p_i, n [\sum_{i=1}^k c_i^2 p_i - (\sum_{i=1}^k c_i p_i)^2])$$

c'est-à-dire $V_{\underline{C}}(n, k)$ tend vers une distribution normale quand le nombre de v.a. observées augmente vers l'infini. La relation indiquée signifie que

$$\lim_{n \rightarrow \infty} F_{V_{\underline{C}}(n, k)}(x) = F_W(x) \text{ où } x \in \mathbb{R} \quad \text{et } W$$

est une v.a. Gaussienne à distribution normale N de paramètres comme indiqués. Cela

entraîne que pour un échantillon aléatoire de grande taille (n au moins 40 est recommandé), $V_{\tilde{c}}(n, k)$ se comporte approximativement comme une v.a. normale, et par conséquent une table de la distribution normale centrée réduite peut être utilisée pour déterminer d'une manière approximative les probabilités des événements aléatoires concernant $V_{\tilde{c}}(n, k)$ sans commettre des erreurs considérables.

§ 4. Cas spéciaux

I.) $c_i = -i \ln(d_k)$ [$i = 1, \dots, k$].

On a $V_{\tilde{c}}(n, k) = -\ln(d_k) \sum_{i=1}^k i N_i(n, k) = -\ln(d_k) R^*(n, k)$.

a.) Propriété d'optimalité

Le test est optimal pour discerner les hypothèses H_0 et $H_1: \tilde{p} = (d_k^{-1}, \dots, d_k^{-k})'$.

En effet $e^{c_i} = d_k^{-i} = p_i^*$. Il s'agit donc d'un test optimal de H_0 versus une situation spécifique où les p_i^* décroissent pour i croissant d'une manière géométrique.

b.) Détermination de d_k .

De la relation $\sum_{l=1}^k d_k^{-l} = 1$ découle

$$2d_k^{-1} - d_k^{-(k+1)} - 1 \equiv 0.$$

Cette équation admet une solution unique du problème considéré. On trouve

pour $k = 2$ $d_2 = 1,62$

$k = 3$ $d_3 = 1,83$

$k \rightarrow \infty$ $d_{\infty} = 2$.

Selon le choix de k d_k prend différentes valeurs dans l'intervalle $[1, 62 ; 2]$.

c.) Etude de la statistique $R^*(n, k)$ équivalente à $V_{\tilde{c}}(n, k)$ dans le cas I.

Considérons maintenant la statistique

$$R^*(n, k) = \sum_{i=1}^k i N_i(n, k).$$

Elle ne se distingue de $V_c(n, k)$ avec $c_i = -i \ln(d_k)$ que par une constante de proportionnalité et $R^*(n, k)$ est donc une statistique de test équivalente à $V_c(n, k)$. Evidemment il faut ajuster la région critique du test si l'on teste H_0 versus H_1 comme indiqué sous a.), la région critique de ce test a maintenant la forme $R^*(n, k) < c$, donc elle est unilatérale du côté gauche. Notons que la statistique $R^*(n, k)$ est liée à la statistique $R(n, k)$ mentionnée dans (Maag *et al.*, 1973) et sa forme spécialisée $R(n, n) = S^*(n)$ mentionnée dans (Riedwyl, 1967) par la relation

$$R(n, k) = n^{-1} R^*(n, k) - 2^{-1}(k + 1).$$

Les formules pour l'espérance mathématique et la variance $R^*(n, k)$ sont

$$E[R^*(n, k) : p] = n \sum_{i=1}^k i p_i \quad \text{et}$$

$$\text{Var}[R^*(n, k) : p] = n \left(\sum_{i=1}^k i^2 p_i - \left(\sum_{i=1}^k i p_i \right)^2 \right).$$

En particulier on trouve

$$E[R^*(n, k) : H_0] = n 2^{-1}(k + 1) \quad \text{et}$$

$$\text{Var}[R^*(n, k) : H_0] = n(12)^{-1} (k^2 - 1).$$

Ces moments servent à déterminer les distributions asymptotiques normales pour n grand.

- d.) Détermination du niveau de signification effectif et de la fonction de puissance.

Déterminons maintenant le niveau de signification exact et la puissance exacte du test à région critique $R^*(n, k) < c$ pour n finie. En tenant compte du fait que $\tilde{N} = (N_1, \dots, N_k)'$ suit une distribution hypergéométrique

$$\text{Pui}(\tilde{p}) = \sum_{\Delta} n! [n_1! \dots n_k!]^{-1} p_1^{n_1} \dots p_k^{n_k}$$

où $\Delta = \{(n_1, \dots, n_k) \in \{0, 1, 2, \dots\}^k \text{ tels que:}$

$$n_1 + n_2 + \dots + n_k = n \text{ et } n_1 + 2n_2 + \dots + kn_k < c\}.$$

Un calcul simple pour n fixé tenant compte de la formule

$$(p_1 x + p_2 x^2 + \dots + p_k x^k)^n = x^n (p_1 + p_2 x + \dots + p_k x^{k-1})^n$$

$$= \sum_{m=n}^{nk} P(R^*(n, k) = m) x^m = \sum_{m=n}^{nk} \gamma_m(\tilde{p}) x^m$$

et de la récursion

où $\gamma_n(\underline{p}) = p_1^n$

et $\gamma_{n+j}(\underline{p}) = \frac{1}{j p_1} \sum_{v=1}^{\min(j, k-1)} (vn - j + v) p_{v+1} \gamma_{n+j-v}(\underline{p})$

$[j = 1, 2, \dots, (k - 1)n]$

conduit à

$Pui(\underline{p}) = \sum_{m=n}^{[c-0]} \gamma_m(\underline{p})$

où $[c - 0] := l - 1$ si $l - 1 < c \leq l$ et $l \in \{0, 1, 2, \dots\}$.

α_{eff} s'obtient en substituant $p_i = k^{-1}$, donc $\underline{p} = k^{-1} \underline{1}$, dans les formules précédentes. Nous écrivons $\tilde{\gamma}_m$ comme abréviation de $\gamma_m(k^{-1} \underline{1})$.

e.) Détermination des alternatives pour lesquelles le test est sans biais.

Nous disons que le test statistique à région critique $R^*(n, k) < c$ est sans biais pour tester $H_0 : \underline{p} = k^{-1} \underline{1}$ versus $H_1(\underline{p}^*) : \underline{p} = \underline{p}^*$ quand $Pui(\underline{p}^*) \geq \alpha_{eff}$. Cette propriété désirable d'un test nous garantit que le test ne rejette H_0 pas plus fréquemment dans le cas où H_0 est correcte que dans le cas où H_0 est incorrecte. Selon le résultat d.) le test considéré est sans biais si et seulement si

$$\sum_{m=n}^{[c-0]} \gamma_m(\underline{p}^*) \geq \sum_{m=n}^{[c-0]} \tilde{\gamma}_m.$$

Ce critère peut être utilisé pour délimiter la classe de toutes les hypothèses alternatives uniponctuelles par rapport auxquelles le $R^*(n, k)$ -test est sans biais. Nous illustrons la technique à l'aide de l'exemple suivant qui traite le cas $n = 3$ et $k = 4$.

Table des valeurs possibles de $R^*(n, k)$ désignées par m , de $\gamma_m(\underline{p}^*)$ et de $\tilde{\gamma}_m$.

m	$\gamma_m(\underline{p}^*)$	$\tilde{\gamma}_m$
3	p_1^{*3}	1/64
4	$3p_1^{*2} p_2^*$	3/64

5	$3 p_1^* p_2^{*2} + 3 p_1^{*2} p_3^*$	6/64
6	$6 p_1^* p_2^* p_3^* + 3 p_1^{*2} p_4 + p_2^{*3}$	10/64
7	$6 p_1^* p_2^* p_4^* + 3 p_1^* p_3^{*2} + 3 p_2^{*2} p_3^*$	12/64
8	$6 p_1^* p_3^* p_4^* + 3 p_2^* p_3^{*2} + 3 p_2^{*2} p_4^*$	12/64
9	$6 p_2^* p_3^* p_4^* + 3 p_1^* p_4^{*2} + p_3^{*3}$	10/64
10	$3 p_2^* p_4^{*2} + 3 p_3^{*2} p_4^*$	6/64
11	$3 p_3^* p_4^{*2}$	3/64
12	p_4^{*3}	1/64

Supposons que l'on désire que $\alpha = \alpha_{eff} = 5\%$. On constate

$$P(R^*(3, 4) < 4 : \underline{p} = k^{-1} \underline{1}) = 1/64 < 0,05$$

$$P(R^*(3, 4) < 5 : \underline{p} = k^{-1} \underline{1}) = 1/16 > 0,05.$$

On peut cependant réaliser $\alpha_{eff} = 0,05$ exactement en choisissant le test de région critique $r^*(3, 4) < 4$ qui rejette H_0 avec probabilité 11/15 quand $r^*(3, 4) = 4$ (test randomisé).

La puissance de ce test vaut

$$\text{Pui}(\underline{p}^*) = P(R^*(3, 4) < 4; \underline{p} = \underline{p}^*) + (11/15) P(R^*(3, 4) = 4 : \underline{p} = \underline{p}^*) = p_1^{*3} + (11/5) p_1^{*2} p_2^*.$$

Le test est sans biais si et seulement si

$$p_1^{*3} + (11/5) p_1^{*2} p_2^* \geq 0,05.$$

On peut donc faire les remarques suivantes:

- 1.) Quand $p_1^* \in [0; 0, 146]$, l'inégalité indiquée n'est pas satisfaite et donc le test est biaisé.
- 2.) Le test est sans biais quand

$$p_1^* \in (0, 146; 0, 368] \quad \text{et} \quad p_2^* \geq 5(0,05 - p_1^{*3}) p_1^{*-2} / 11$$

avec $p_2^*, p_3^*, p_4^* \in [0, 1]$ tels que $\sum_{i=1}^4 p_i^* = 1$. Il a du biais quand la condition indiquée pour p_1^* est satisfaite, mais celle pour p_2^* ne l'est pas.

3.) Le test est sans biais quand $p_1^* \in (0, 368; 1]$ et $p_2^*, p_3^*, p_4^* \in [0, 1]$ tels que

$$\sum_{i=1}^4 p_i^* = 1.$$

II.) $c_i = \ln(id_k) \quad [i = 1, \dots, k]$.

$$\begin{aligned} \text{On a } V_{\tilde{c}}(n, k) &= \sum_{i=1}^k \ln(i) N_i(n, k) + n \ln(d_k) \\ &= L^*(n, k) + n \ln(d_k). \end{aligned}$$

a.) Propriété d'optimalité.

Le test est optimal pour discerner les hypothèses H_0 et H_1 :

$\tilde{p} = (d_k, 2d_k, \dots, kd_k)$. En effet $e^{c_i} = id_k = p_i^*$.

Il s'agit donc d'un test optimal pour discerner les hypothèses H_0 et $H_1 : p_i^* = id_k$ où les p_i^* croissent avec i de la façon $p_i^* = ip_1^*$.

b.) Détermination de d_k .

De la relation $\sum_{i=1}^k p_i^* = 1$ découle $d_k = 2 [k(k+1)]^{-1}$.

c.) Etude de la statistique $L^*(n, k)$ équivalente à $V_{\tilde{c}}(n, k)$ dans le cas II.

Pour l'espérance mathématique et la variance, on trouve

$$\begin{aligned} E[L^*(n, k) : \tilde{p}] &= n \sum_{i=1}^k \ln(i) p_i \text{ et} \\ \text{Var}[L^*(n, k) : \tilde{p}] &= n \left(\sum_{i=1}^k p_i \ln^2(i) - \left(\sum_{i=1}^k p_i \ln(i) \right)^2 \right). \end{aligned}$$

En particulier, on obtient pour $\tilde{p} = k^{-1} \mathbf{1}$

$$E[L^*(n, k) : H_0] = k^{-1} n \ln(k!) \text{ et}$$

$$\text{Var}[L^*(n, k) : H_0] = n(k^{-1} \sum_{i=1}^k \ln^2(i) - k^{-2} \ln^2(k!)).$$

Pour n grand, les distributions asymptotiques de $L^*(n, k)$ sont Gaussiennes et déterminées par ces moments.

Une étude plus approfondie de quelques propriétés du $R^*(n, k)$ -test et de quelques autres tests d'ajustement se trouve dans la thèse du premier auteur (Driancourt, 1994).

RÉSUMÉ

Après un bref rappel des buts de la statistique et de quelques notions relatives aux tests statistiques, et en particulier aux tests d'ajustement, nous proposons une classe générale de tests d'ajustement et indiquons une propriété d'optimalité de ces tests. Nous étudions l'espérance mathématique, la variance et les distributions asymptotiques des statistiques concernées et traitons des cas particuliers en abordant la question de la détermination du niveau de signification effectif et de la fonction de puissance, et en analysant pour lesquelles des hypothèses alternatives les tests sont sans biais.

RÉFÉRENCES BIBLIOGRAPHIQUES

- CHUNG, K.L. (1974). *A course in probability theory*. Academic Press, New York.
- DRIANCOURT, C. (1994). *Analyse concrète de propriétés relatives à des tests d'ajustement*. Thèse. Université de Genève.
- LEHMANN, E.L. (1986). *Testing statistical hypotheses*. J. Wiley, New York.
- MAAG, U., STREIT, F. & DROULLY, P. (1973). Goodness-of-fit tests for grouped data. *J. Americ. Statist. Assoc.* 68: 462-465.
- RIEDWYL, H. (1967). Goodness of fit. *J. Americ. Statist. Assoc.* 62: 390-398.
- WITTING, H. & NÖLLE, G. (1970). *Angewandte Mathematische Statistik*, Teubner, Stuttgart, p. 106 (théorème 3.5).