



**HAL**  
open science

# Étude d'un coefficient d'indépendance d'une variable Y par rapport à une variable aléatoire X

J. Guy, J. Bass

► **To cite this version:**

J. Guy, J. Bass. Étude d'un coefficient d'indépendance d'une variable Y par rapport à une variable aléatoire X. Annales de l'ISUP, 1976, XXI (1-2), pp.43-63. hal-04082327

**HAL Id: hal-04082327**

<https://hal.sorbonne-universite.fr/hal-04082327v1>

Submitted on 26 Apr 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

ETUDE D'UN COEFFICIENT D'INDEPENDANCE  
D'UNE VARIABLE ALEATOIRE Y  
PAR RAPPORT A UNE VARIABLE ALEATOIRE X.

Jean GUY\*

et Jean BASS\*\*

Résumé :

Soit  $(X, Y)$  un couple de variables aléatoires. Soit  $\omega_{ij} = \Pr(X = x_i \text{ et } Y = y_j)$ . On s'intéresse à la manière dont la loi  $\omega_{ij}$  se rapproche d'une loi d'indépendance. Pour cela, on pose :

$$\omega_{ij} = \mu p_i q_j + (1 - \mu) u_{ij},$$

où  $p_i = \sum_j \omega_{ij} = \Pr(X = x_i)$  est la loi marginale relative à X, puis on cherche la plus grande valeur possible  $\mu_0$  du coefficient  $\mu$  ( $0 \leq \mu_0 \leq 1$ ).

$$\text{On montre que } \mu_0 = \sum_j \inf_i \frac{\omega_{ij}}{p_i}$$

Grâce à la définition préalable d'une distance entre matrices stochastiques il apparaît que le rapport  $d = \frac{1 - \mu_0}{\mu_0}$  représente la distance de la matrice A

des probabilités conditionnelles :

$$a_{ij} = \frac{\omega_{ij}}{p_i} = \Pr(Y = y_j / X = x_i)$$

à l'ensemble des matrices stochastiques d'indépendance.

On donne enfin quelques exemples d'application à des problèmes concrets:

\* U.E.R. Etudes médicales et biologiques de l'Université René Descartes.  
45, rue des Saint-Pères. 75006 - Paris.

\*\* U.E.R. 47 de Mathématiques de l'Université Pierre et Marie Curie  
Tour 46, 4, place Jussieu, 75280 Paris Cedex 05



44.

I. Enoncé du problème

Soient  $X$  et  $Y$  deux variables aléatoires prenant respectivement les valeurs  $x_1, x_2, \dots, x_m$  et  $y_1, y_2, \dots, y_n$ . Appelons :

$\omega_{ij}$  la probabilité de  $X = x_i, Y = y_j$

$p_i$  la probabilité de  $X = x_i$

$q_j$  la probabilité de  $Y = y_j$

(les probabilités  $p_i$  et  $q_j$  étant supposées non nulles).

Si  $\omega_{ij} = p_i q_j$ ,  $X$  et  $Y$  sont indépendantes.

Pour  $m \geq n$ , si  $\omega_{ij} = p_i \delta_{f(i), j}$  ( $f(i)$  étant une application de l'ensemble des indices  $i$  sur l'ensemble des indices  $j$ ),  $Y$  devient une fonction de  $X$ .

Dans les autres cas, il y a une dépendance stochastique de  $Y$  par rapport à  $X$  et il convient de chercher un critère simple pour juger la manière dont la dépendance en question se rapproche de l'indépendance.

Le coefficient de corrélation linéaire  $\rho$  entre  $X$  et  $Y$  donne des indications à ce sujet. S'il y a indépendance, il est nul. Mais s'il est nul, on ne peut rien dire. Il peut même arriver que  $\rho = 0$  alors que  $Y$  est fonction de  $X$ , ainsi que cela a lieu si  $Y = X^2$ ,  $X$  étant distribuée suivant une loi normale centrée.

Seul le cas où  $\rho = \pm 1$  est précis :  $Y$  est alors fonction linéaire de  $X$ ; réciproquement,  $Y = aX + b$  impose  $|\rho| = 1$ .

Il est donc naturel de chercher d'autres procédés simples pour tester le degré de dépendance de  $Y$  par rapport à  $X$ . L'objet de ce travail est d'en suggérer un.

Posons :

$$(1) \quad \omega_{ij} = \mu p_i q_j' + (1 - \mu) u_{ij}.$$

$p_i$  est la probabilité pour que  $X = x_i$ .  $\mu$  est un nombre compris entre 0 et 1.  $u_{ij}$  a, relativement aux indices, les caractères d'une probabilité :  $u_{ij} \geq 0$ ,

$\sum_{ij} u_{ij} = 1$ . On suppose enfin que  $q_j$  est aussi une probabilité. On vérifie que



toutes ces conditions sont compatibles. On a alors :

$$(2) \quad p_i = \sum_j u_{ij}$$

$$(3) \quad q_j = \mu q'_j + (1 - \mu) \sum_i u_{ij}$$

En général  $q'_j \neq q_j$ .

Si  $\mu = 1$ , on a  $\omega_{ij} = p_i q'_j = p_i q_j$ . Il y a indépendance. Si  $\mu = 0$ ,  $\omega_{ij} = u_{ij}$ .

Si l'on se donne  $\omega_{ij}$ ,  $\mu$  n'est pas déterminé. En particulier  $\mu$  peut prendre la valeur 0. Mais en général  $\mu$  ne peut pas prendre la valeur 1. Il y a une valeur  $\mu_0$  maximale compatible avec la donnée de  $\omega_{ij}$ . Cette valeur représente jusqu'à un certain point la manière dont la loi  $\omega_{ij}$  donnée se rapproche de l'indépendance. C'est elle qui va servir à tester le degré de dépendance de Y par rapport à X.

Mais, avant de faire le calcul de  $\mu_0$ , nous pouvons interpréter la relation

(1). Introduisons  $m + 2$  variables aléatoires indépendantes

Z prenant les valeurs 1 et 0 avec les probabilités  $\mu$  et  $1 - \mu$ ,

A prenant n valeurs avec les probabilités  $q'_j$ .

$B_i$  prenant n valeurs avec les probabilités  $u_{ij}$  ( $i = 1, 2, \dots, m$ ).

On choisit la valeur de X. Pour trouver la valeur de Y, une fois X choisie, on commence par tirer au sort la valeur de Z. Si  $Z = 1$ , on fait une expérience pour déterminer A et l'on prend pour valeur de Y la valeur trouvée pour A, quelle que soit la valeur initiale de X. Si  $Z = 0$ , on fait une expérience pour déterminer la valeur de la variable  $B_i$  dont l'indice correspond à la valeur initiale de X, et l'on prend pour valeur de Y la valeur trouvée pour  $B_i$ .

Pratiquement, Z résulte d'une épreuve de Bernoulli (pile ou face généralisée). A est le résultat du tirage dans une urne. Aux  $B_i$  correspondent m urnes. On essaie Z. Si  $Z = 1$ , la valeur de Y résulte du tirage dans la première urne. Si  $Z = 0$ , on obtient la valeur de Y en tirant une boule de celle des urnes qui correspond à la valeur initialement choisie pour X.

2. Calcul de la valeur maximale  $\mu_0$ . Compatibilité

L'équation (1) entraîne :

46.

$$q_j^* \leq \frac{\omega_{ij}}{p_i} \quad (\text{quel que soit } i) \text{ et par conséquent :}$$

$$(4) \quad \mu q_j^* \leq \inf_i \frac{\omega_{ij}}{p_i}$$

Il en résulte que :

$$(5) \quad \mu \leq \sum_j \left( \inf_i \frac{\omega_{ij}}{p_i} \right).$$

et la valeur maximale de  $\mu$  est :

$$(6) \quad \mu_0 = \sum_j \left( \inf_i \frac{\omega_{ij}}{p_i} \right)$$

Les valeurs associées des  $q_j^*$  sont alors données par :

$$(7) \quad \mu_0 q_j^* = \inf_i \frac{\omega_{ij}}{p_i}$$

et celles des  $u_{ij}$  par (1). Elles sont déterminées, sauf bien entendu si  $\mu_0 = 1$ .

Montrons que, pour toute valeur de  $\mu$  telle que  $0 < \mu < \mu_0$ , le système (1) permet de calculer des  $q_j^*$  (non complètement déterminés si  $\mu < \mu_0$ ) et des  $u_{ij}$ .

On a :

$$(8) \quad \mu q_j^* = \frac{\omega_{ij}}{p_i} - (1-\mu) \frac{u_{ij}}{p_i} \leq \frac{\omega_{ij}}{p_i} \leq \inf_i \frac{\omega_{ij}}{p_i}$$

Donnons-nous des  $q_j^* \geq 0$  satisfaisant à l'inégalité :

$$(9) \quad \mu q_j^* \leq \inf_i \frac{\omega_{ij}}{p_i}$$



Nous pouvons en outre demander que  $\sum_j q_j' = 1$ , puisque :

$$\mu \leq \mu_0 = \sum_j \left[ \inf_i \frac{\omega_{ij}}{p_i} \right]$$

Les  $q_j'$  étant ainsi choisis, l'égalité (1) donne les  $u_{ij}$ . Ils sont bien tels que  $\sum_{ij} u_{ij} = 1$ . En outre, ils sont bien positifs, car

$$\omega_{ij} - \mu p_i q_j' \geq 0.$$

Le problème a donc en général une infinité de solutions. Il en a une seule si  $\mu = \mu_0$ . Dans ce cas, on peut préciser le comportement des  $u_{ij}$ . Soit  $i_0$  l'indice  $i$  pour lequel  $\frac{\omega_{ij}}{p_i}$  est minimal. On voit que :

$$(10) \quad u_{i_0 j} = 0.$$

### 3. Relation entre $\mu$ et $\rho$

On remarque tout d'abord que  $\mu$ , et en particulier  $\mu_0$ , ne dépend que des probabilités des variables X, Y, et non des valeurs prises par ces variables. On ne change donc pas  $\mu$ , ou  $\mu_0$ , si l'on modifie les valeurs X, Y en conservant les probabilités. Au contraire, le coefficient de corrélation linéaire  $\rho$  dépend des valeurs prises par X et Y. Mais, pour comparer  $\rho$  et  $\mu_0$ , on peut choisir ces valeurs comme on le veut. Nous supposons seulement que :

$$EX = EY = 0.$$

A partir de la relation (1), on obtient :

$$\sigma_1^2 = EX^2 = \sum_{ij} \omega_{ij} x_i^2 = \mu \sigma_1^2 + (1 - \mu) \sum_{ij} u_{ij} x_i^2$$

soit encore :

$$(11) \quad \sigma_1^2 = \sum_{ij} u_{ij} x_i^2$$

et l'on a pour  $\sigma_2^2$  :

48.

$$(12) \quad \sigma_2^2 = \sum_{ij} \omega_{ij} y_j^2 = \mu \sum_{ij} q'_j y_j^2 + (1 - \mu) \sum_{ij} u_{ij} y_j^2$$

La définition même de  $\rho$  donne :

$$(13) \quad \rho \sigma_1 \sigma_2 = \sum_{ij} \omega_{ij} x_i y_j = 0 + (1 - \mu) \sum_{ij} u_{ij} x_i y_j.$$

Donc :

$$(14) \quad \rho = \frac{1 - \mu}{\sigma_1 \sigma_2} \sum_{ij} u_{ij} x_i y_j$$

Posons

$$r = \frac{\sum_{ij} u_{ij} x_i y_j}{\sqrt{\sum_{ij} u_{ij} x_i^2} \cdot \sqrt{\sum_{ij} u_{ij} y_j^2}}$$

$r$  n'est pas en général un coefficient de corrélation car, pour la loi  $u_{ij}$ ,

$Y$  n'est pas de valeur moyenne nulle. Mais on a bien entendu :

$$|r| \leq 1.$$

On obtient :

$$(15) \quad \rho = \frac{(1 - \mu) r}{\sigma_2} \sqrt{\sum_{ij} u_{ij} y_j^2}$$

et d'après (12)

$$(16) \quad \sqrt{1 - \mu} \sqrt{\sum_{ij} u_{ij} y_j^2} \leq \sigma_2.$$

Donc :

$$(17) \quad \rho \leq \sqrt{1 - \mu}$$

Cette dernière inégalité, valable pour tout  $\mu$ , l'est en particulier pour  $\mu_0$ .



4. Exemple du couple de variables aléatoires à deux valeurs

Supposons  $m = n = 2$ . On va a priori choisir les deux valeurs possibles de X et de Y pour que  $EX = 0$ ,  $EY = 0$ . Alors elles sont définies à un facteur près, qui n'a pas d'influence sur les calculs et les résultats.

Nous désignerons par  $p$ ,  $1 - p$  les probabilités des valeurs  $x_1$  et  $x_2$  de X, par  $q$ ,  $1 - q$  celles des valeurs  $y_1$  et  $y_2$  de Y, par  $\lambda$  la probabilité du couple  $(x_1, y_1)$ . A un facteur près, on a nécessairement :

$$x_1 = 1 - p \qquad x_2 = -p$$

$$y_1 = 1 - q \qquad y_2 = -q$$

ce qui permet de construire le tableau usuel suivant des probabilités des couples de valeurs et des probabilités marginales :

X Y	1 - p	- p	
1 - q	λ	q - λ	q
- q	p - λ	1 + λ - p - q	1 - q
	p	1 - p	1

On a :

$$(18) \quad p + q - 1 \leq \lambda \leq \inf(p, q).$$

Il y a indépendance si  $\lambda = pq$ .

Il n'est pas en général possible de choisir  $\lambda$ , une fois  $p$  et  $q$  donnés, pour que Y soit fonction de X. Il faut pour cela que  $p = q$ . Alors, si  $\lambda = p = q$ , le tableau ci-dessus devient :

$$\begin{pmatrix} p & 0 \\ 0 & 1 - p \end{pmatrix}$$



50.

On trouve facilement :

$$\begin{aligned} EX^2 &= p(1-p) & EY^2 &= q(1-q) \\ EXY &= \lambda - pq. \end{aligned}$$

Le coefficient de corrélation entre X et Y est :

$$(19) \quad \rho = \frac{\lambda - pq}{\sqrt{p(1-p)} \cdot \sqrt{q(1-q)}}$$

Pour calculer  $\mu$ , on écrit :

	$\lambda = \mu pr + (1 - \mu) u_{11}$	
	$p - \lambda = \mu p(1 - r) + (1 - \mu) u_{12}$	
	$q - \lambda = \mu(1 - p)r + (1 - \mu) u_{21}$	
	$1 + \lambda - p - q = \mu(1 - p)(1 - r) + (1 - \mu) u_{22}$	
On a :		

$$\mu r = \inf\left(\frac{\lambda}{p}, \frac{q - \lambda}{1 - p}\right), \quad \mu(1 - r) = \inf\left(\frac{p - \lambda}{p}, \frac{1 + \lambda - p - q}{1 - p}\right)$$

La valeur de  $r$  dépend du signe de  $\lambda - pq$ . La valeur maximale  $\mu_0$  de  $\mu$  est dans tous les cas égale à  $\mu_0 = 1 - \frac{|\lambda - pq|}{p(1 - p)}$ , ce qui s'écrit aussi :

$$(20) \quad 1 - \mu_0 = \frac{|\lambda - pq|}{p(1 - p)}$$

Il s'agit du  $\mu_0$  correspondant aux probabilités conditionnelles de Y une fois X choisi. Si l'on permute X et Y, on trouve un  $\mu'_0$  tel que :

$$(21) \quad 1 - \mu'_0 = \frac{|\lambda - pq|}{q(1 - q)}$$

On constate que :

$$(22) \quad (1 - \mu_0)(1 - \mu'_0) = \rho^2$$

On vérifie que  $\mu_0$  et  $\mu'_0$  sont compris entre 0 et 1 et que :

$$(23) \quad \rho^2 \leq \inf(1 - \mu_0, 1 - \mu'_0)$$

5. **Interprétation du coefficient  $\mu_0$  à l'aide de la distance de deux matrices stochastiques.**

Nous poserons  $a_{ij} = \frac{\omega_{ij}}{p_i}$ .  $a_{ij}$  est la probabilité conditionnelle de  $Y = y_j$ , lorsque  $X = x_i$ . La matrice  $A = (a_{ij})$  est une matrice stochastique :

$$a_{ij} \geq 0, \quad \sum_j a_{ij} = 1.$$

Appelons  $E_{mn}$  l'ensemble des matrices stochastiques A à m lignes et n colonnes.

Nous allons maintenant définir une distance  $d(A, B)$  entre deux matrices  $A, B \in E_{mn}$  par :

$$(24) \quad d(A, B) = \frac{1}{m} \sum_{ij} \left| \frac{a_{ij} - b_{ij}}{\mu_{A0} \mu_{B0}} \right|,$$

$\mu_{A0}$  et  $\mu_{B0}$  étant les deux constantes entièrement définies à partir des matrices A et B, suivant :

$$(25) \quad \mu_{A0} = \sum_j (\inf_i a_{ij}), \quad \mu_{B0} = \sum_j (\inf_i b_{ij}).$$

Vérifions tout d'abord que  $d(A, B)$  est bien une distance.

a) La relation de définition (24) demeure inchangée par permutation de A et de B, d'où  $d(A, B) = d(B, A)$ . On vérifie de plus immédiatement que  $d(A, A) = 0$ .

b) Une valeur nulle de  $d(A, B)$  entraîne l'identité des matrices A et B. Remarquons tout d'abord que  $d(A, B) = 0 \Rightarrow \mu_{A0} \neq 0$  et  $\mu_{B0} \neq 0$  (la distance  $d(A, B)$  devient infiniment grande si une et une seulement des deux constantes  $\mu_{A0}$  et  $\mu_{B0}$  est nulle ; elle n'est plus définie si  $\mu_{A0} = \mu_{B0} = 0$ ). Par suite  $d(A, B) = 0$  entraîne :



52.

$$(26) \quad \frac{a_{ij}}{\mu_{A_0}} = \frac{b_{ij}}{\mu_{B_0}} \quad (\forall i, j) \quad \text{c'est-à-dire} \quad a_{ij} = \frac{\mu_{A_0}}{\mu_{B_0}} b_{ij}.$$

Par sommation sur j de la dernière égalité, il vient :

$$(27) \quad 1 = \frac{\mu_{A_0}}{\mu_{B_0}}, \text{ d'où il résulte que } a_{ij} = b_{ij} \quad (\forall i, j).$$

c) L'inégalité triangulaire

$$(28) \quad d(A, C) \leq d(A, B) + d(B, C)$$

est également valable. Ecrivons :

$$d(A, B) = d'(A', B'), \quad d(A, C) = d'(A', C'), \quad d(B, C) = d'(B', C'),$$

où A', B' et C' représentent les matrices (non stochastiques cette fois) définies par :

$$A' = \frac{1}{\mu_{A_0}} A, \quad B' = \frac{1}{\mu_{B_0}} B, \quad C' = \frac{1}{\mu_{C_0}} C,$$

On a :

$$(29) \quad d'(A', B') = \frac{1}{m} \sum_{ij} |a'_{ij} - b'_{ij}|$$

et des formules analogues. Or (29) définit l'une des distances usuelles entre deux matrices de dimensions (m, n), distance pour laquelle on a bien :

$$d'(A', C') \leq d'(A', B') + d'(B', C').$$

Lorsque la matrice B devient une matrice stochastique d'indépendance (les m lignes sont alors identiques), on peut poser  $b_{ij} = b_j (\forall i)$ . Donc :

$$(30) \quad \mu_{B_0} = \sum_j \inf_i (b_j) = \sum_j b_j = 1.$$

Prenons en particulier pour matrice B la matrice d'indépendance Q' dont les éléments  $q'_j$  découlent de ceux de la matrice A par la résolution des équations

$$(31) \quad a_{ij} = \mu_{A_0} q'_j + (1 - \mu_{A_0}) v_{ij}$$

( $v_{ij}$  étant automatiquement les éléments d'une nouvelle matrice stochastique V).  
On obtient :

$$\begin{aligned}
 (32) \quad d(A, Q') &= \frac{1}{m} \sum_{ij} \left| \frac{a_{ij}}{\mu_{Ao}} - q'_j \right| \\
 &= \frac{1}{m} \sum_{ij} \left| q'_j + \frac{1 - \mu_{Ao}}{\mu_{Ao}} v_{ij} - q'_j \right| \\
 &= \frac{1}{m} \frac{1 - \mu_{Ao}}{\mu_{Ao}} \sum_{ij} v_{ij} = \frac{1}{m} \frac{1 - \mu_{Ao}}{\mu_{Ao}} \sum_i 1 \\
 &= \frac{1 - \mu_{Ao}}{\mu_{Ao}}
 \end{aligned}$$

$\frac{1 - \mu_{Ao}}{\mu_{Ao}}$  est donc la distance de la matrice A à la matrice d'indépendance qui

lui est associée par la formule (31).

Nous voulons montrer que, si  $b_j$  représente les éléments d'une matrice stochastique d'indépendance quelconque, on a :

$$(33) \quad \frac{1}{m} \sum_{ij} \left| \frac{a_{ij}}{\mu_{Ao}} - b_j \right| \geq \frac{1 - \mu_{Ao}}{\mu_{Ao}}$$

Or, par l'hypothèse,

$$a_{ij} = \mu_{Ao} q'_j + (1 - \mu_{Ao}) v_{ij}$$

Nous poserons  $b_j = q'_j + r_j$ .

Le premier membre de (33) s'écrit donc :

$$(34) \quad \frac{1}{m} \sum_{ij} \left| \frac{1 - \mu_{Ao}}{\mu_{Ao}} v_{ij} - r_j \right|$$



54.

$1 - \mu_{Ao}$  est différent de 0. Nous savons que certains des  $v_{ij}$  sont nuls. Mais ils ne le sont pas tous. Nous allons d'abord supposer que la matrice  $B = (b_j)$  (à  $n$  lignes identiques) est suffisamment voisine de la matrice  $Q' = q_j'$  pour que, lorsque  $v_{ij} \neq 0$ , on ait :

$$(35) \quad \frac{1 - \mu_{Ao}}{\mu_{Ao}} v_{ij} - r_j > 0.$$

Dans (34), il y a donc des termes qui se réduisent à  $|r_j|$  et d'autres qui sont représentés par (35), au facteur  $\frac{1}{m}$  près.  $i$  étant choisi, faisons dans (34) la sommation par rapport à  $j$ . Il y a deux sortes de termes :

$$(36) \quad \sum_{j'} \left( \frac{1 - \mu_{Ao}}{\mu_{Ao}} v_{ij'} - r_{j'} \right) \quad (j' \text{ tel que } v_{ij'} \neq 0)$$

$$(37) \quad \sum_{j''} |r_{j''}| \quad (j'' \text{ tel que } v_{ij''} = 0)$$

Or (36) est la somme de :

$$(38) \quad \sum_j \frac{1 - \mu_{Ao}}{\mu_{Ao}} \quad (\text{pour toutes les valeurs de } j)$$

et de :

$$(39) \quad - \sum_{j'} r_{j'}$$

Au total, en ajoutant (36) et (37), on obtient :

$$(40) \quad \sum_j \frac{1 - \mu_{Ao}}{\mu_{Ao}} v_{ij} + \sum_{j''} |r_{j''}| - \sum_{j'} r_{j'}$$

Mais  $\sum_j r_j = 0$ . Donc  $\sum_{j'} r_{j'} = - \sum_{j''} r_{j''}$  et l'on obtient :

$$(41) \quad \sum_j \frac{1 - \mu_{Ao}}{\mu_{Ao}} v_{ij} + \sum_{j''} |r_{j''}| + \sum_{j''} r_{j''}$$

soit :

$$(42) \quad \frac{1 - \mu_{Ao}}{\mu_{Ao}} + H,$$

où H est une quantité  $\geq 0$ . On somme enfin par rapport à i, on divise par m, et l'on voit que (34) est la somme de  $\frac{1 - \mu_{Ao}}{\mu_{Ao}}$  et d'une quantité  $\geq 0$ . Cela prouve l'inégalité (33).

La démonstration  
prouve que  $\frac{1 - \mu_{Ao}}{\mu_{Ao}}$  est

un minimum relatif pour la distance entre la matrice stochastique A et l'ensemble E des matrices d'indépendance. Mais E est un ensemble convexe. Si donc la distance entre le point A et l'ensemble E admet un minimum relatif au point Q', c'est aussi le minimum absolu. En effet, s'il existait un point B de E tel que :

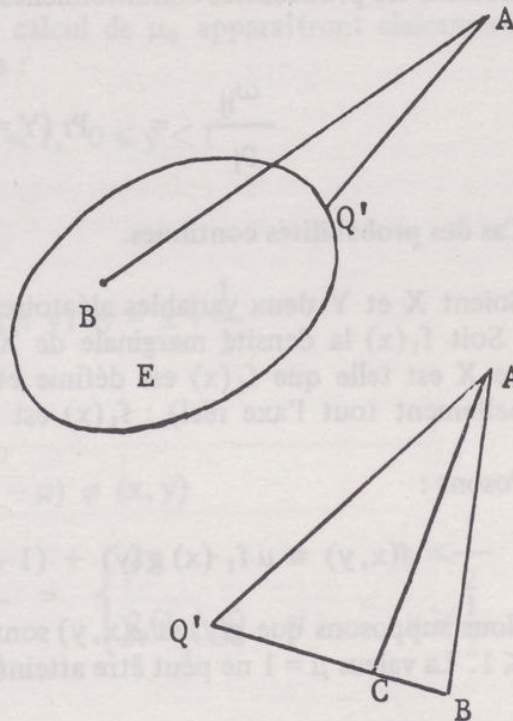


Fig. 1

$$d(A, B) < d(A, Q'),$$

on aurait :

$$d(A, C) < d(A, Q')$$

pour tout C du segment BQ'. Mais  $C \in E$  (cf. fig. 1). Il y aurait donc, aussi près qu'on le veut de Q', un point C tel que  $d(A, C) < d(A, Q')$  et d(A, Q') ne serait pas un minimum relatif.



56.

On peut conclure de (33) que :

$$d(A, Q') = \frac{1 - \mu_{Ao}}{\mu_{Ao}}$$

est la distance de la matrice stochastique A à l'ensemble des matrices stochastiques d'indépendance. La matrice Q', complètement définie par les équations (31), est la matrice d'indépendance la plus proche de la matrice stochastique A qui représente les probabilités conditionnelles :

$$\frac{\omega_{ij}}{p_i} = \Pr(Y = y_j / X = x_i).$$

#### 6. Cas des probabilités continues.

Soient X et Y deux variables aléatoires ayant une densité de probabilité  $f(x, y)$ . Soit  $f_1(x)$  la densité marginale de X. Nous supposons que la variable aléatoire X est telle que  $f_1(x)$  est définie et ne s'annule pas sur un intervalle (éventuellement tout l'axe réel) ;  $f_1(x)$  est nulle en dehors de cet intervalle.

Posons :

$$(43) \quad f(x, y) = \mu f_1(x) g(y) + (1 - \mu) \varphi(x, y)$$

Nous supposons que  $g(y)$  et  $\varphi(x, y)$  sont des densités de probabilité et que  $0 \leq \mu \leq 1$ . La valeur  $\mu = 1$  ne peut être atteinte que si :

$$(44) \quad f(x, y) = f_1(x) g(y)$$

Il y a alors indépendance et  $g(y)$  coïncide avec la densité de probabilité marginale  $f_2(y)$  de Y.

Dans le cas général, on a :

$$(45) \quad \mu f_1(x) g(y) \leq f(x, y)$$

On en déduit :

$$(46) \quad \mu g(y) \leq \frac{f(x, y)}{f_1(x)}$$

(densité de probabilité conditionnelle de Y pour  $X = x$ ). Par suite :

$$(47) \quad \mu g(y) \leq \inf_x \frac{f(x, y)}{f_1(x)}$$

La valeur maximale  $\mu_0$  de  $\mu$  est telle que :

$$(48) \quad \mu_0 g(y) = \inf_x \frac{f(x, y)}{f_1(x)}$$

On voit que :

$$(49) \quad \mu_0 = \int_{-\infty}^{+\infty} \inf_x \frac{f(x, y)}{f_1(x)} dy.$$

**Exemple :** Les modalités de calcul de  $\mu_0$  apparaîtront clairement sur le très simple exemple suivant : Posons :

$$f(x, y) = x + y \text{ pour } 0 < x \leq 1, 0 \leq y < 1$$

$$f(x, y) = 0 \text{ ailleurs.}$$

$$\text{On a } f_1(x) = \int_0^1 (x + y) dy = x + \frac{1}{2}$$

On pose :

$$x + y = \mu(x + \frac{1}{2}) g(y) + (1 - \mu) \varphi(x, y)$$

$$g(y) \mu_0 = \inf_{0 \leq x \leq 1} \frac{x + y}{x + \frac{1}{2}} = \begin{cases} 2y & \text{si } y \leq \frac{1}{2} \\ \frac{2(1+y)}{3} & \text{si } y \geq \frac{1}{2} \end{cases}$$

Donc :

$$\mu_0 = \int_0^{\frac{1}{2}} 2y dy + \int_{\frac{1}{2}}^1 \frac{2(1+y)}{3} dy = \frac{5}{6}$$

## 7. Exemples d'applications.

Dès qu'une variable aléatoire Y est stochastiquement dépendante d'une autre variable aléatoire X, le coefficient d'indépendance  $\mu_0$  défini par (6) et la loi d'indépendance  $q_j^i$  la plus proche de la loi conditionnelle  $\frac{\omega_{ij}}{P_i}$

des caractéristiques intéressantes pour ce couple de variables. Les quelques exemples simples traités ci-après sont destinés à mieux préciser les informations



58.

liées à la connaissance de  $\mu_0$  et de la loi  $q_j$ .

a) Etude de la corrélation existant entre les sexes des jumeaux.

A partir des données statistiques concernant les grossesses gémeillaires, il est possible d'estimer les diverses probabilités pour que les deux jumeaux soient, dans l'ordre de leur naissance, garçon-garçon ( $\omega_{GG}$ ), garçon-fille ( $\omega_{GF}$ ), fille-garçon ( $\omega_{FG}$ ) et fille-fille ( $\omega_{FF}$ ). On peut considérer comme valables les valeurs numériques :

$$\begin{aligned} \omega_{GG} &= 0,329, & \omega_{GF} &= \omega_{FG} = 0,179. \\ \omega_{FF} &= 0,313. \end{aligned}$$

qui conduisent au tableau suivant représentant la matrice stochastique A (probabilités conditionnelle pour que le second jumeau soit un garçon ou une fille lorsque le sexe du premier jumeau est connu).

	2e jumeau	G	F
1er jumeau			
G		0,648	0,352
F		0,364	0,636

Par suite  $\mu_0 = 0,364 + 0,352 = 0,716$  et la distance à la loi d'indépendance la plus proche est

$$d = \frac{1}{\mu_0} - 1 = 0,397.$$

Remarquons ici que l'égalité  $\omega_{GF} = \omega_{FG}$  entraîne  $\rho = 1 - \mu_0 = 0,284$

La loi  $q_j$  apparaît ensuite nettement si nous décomposons la matrice A conformément à (31), soit :

$$(50) \quad A = \begin{pmatrix} 0,648 & 0,352 \\ 0,364 & 0,636 \end{pmatrix} = \mu_0 \begin{pmatrix} 0,508 & 0,492 \\ 0,508 & 0,492 \end{pmatrix} + (1 - \mu_0) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



Le complément de  $\mu_0$  à l'unité ( $1 - \mu_0 = 0,284$ ) avait déjà été interprété par E. Borel [1]. Ce complément précise le nombre de fois où le sexe du second enfant est complètement déterminé par le sexe du premier (jumeaux « vrais »): l'identité des sexes, due alors à l'identité des génotypes, découle analytiquement

de la présence de la matrice unité dans le terme  $(1 - \mu_0) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

de la décomposition (50).

b) Valeurs de  $\mu_0$  dans les phénomènes d'hérédité liée au sexe.

Dans le cas des espèces animales analogues à l'espèce humaine du point de vue du déterminisme sexuel (chromosomes sexuels XX chez la femelle et XY chez le mâle), considérons un gène porté par le chromosome X et susceptible d'exister sous deux états allèles G et g : G (dominant chez les femelles) conduit à un phénotype normal cependant que le génotype g pour le mâle ou gg pour la femelle conduit à un phénotype taré. Si  $\eta$  est la proportion de chromosomes X porteurs de g pour l'ensemble de la population (les croisements étant effectués au hasard) et si  $\alpha$  et  $\beta$  sont les proportions de mâles et de femelles, on obtient sans difficultés le tableau suivant (matrice stochastique A) pour les probabilités conditionnelles qu'un sujet soit mâle ou femelle lorsque nous savons que son phénotype est normal (N) ou taré (T).

Phénotype \ Sexe	$\delta$	$\text{♀}$
	N	$\frac{\alpha(1-\eta)}{1-\alpha\eta-\beta\eta^2}$
T	$\frac{\alpha}{\alpha+\beta\eta}$	$\frac{\beta\eta}{\alpha+\beta\eta}$

Pour  $\alpha = \beta = \frac{1}{2}$ , nous trouvons :

$$(51) \quad \mu_0 = \frac{1}{2+\eta} + \frac{\eta}{1+\eta} = \frac{1+3\eta+\eta^2}{(2+\eta)(1+\eta)}$$

Les courbes représentant  $\mu_0$  et  $d = \frac{1}{\mu_0} - 1$  en fonction de  $\eta$  sont portées sur la fig. 2. On constate notamment que  $0,5 < \mu_0 < 0,833$ , les limites indiquées étant respectivement atteintes de manière asymptotique pour  $\eta$  tendant vers zéro ou vers l'unité. Deux statistiques effectuées séparément sur des sujets



normaux et sur des sujets tarés doivent donc permettre d'atteindre  $\mu_0$  : si la valeur trouvée pour  $\mu_0$  diffère significativement de l'unité, on pourra conclure à l'existence d'une hérédité liée au sexe et un tel phénomène sera d'autant plus facile à mettre en évidence par l'intermédiaire de  $\mu_0$  que la tare sera moins répandue, puisque  $\eta$  très faible conduit à la valeur minimale de  $\mu_0$ .

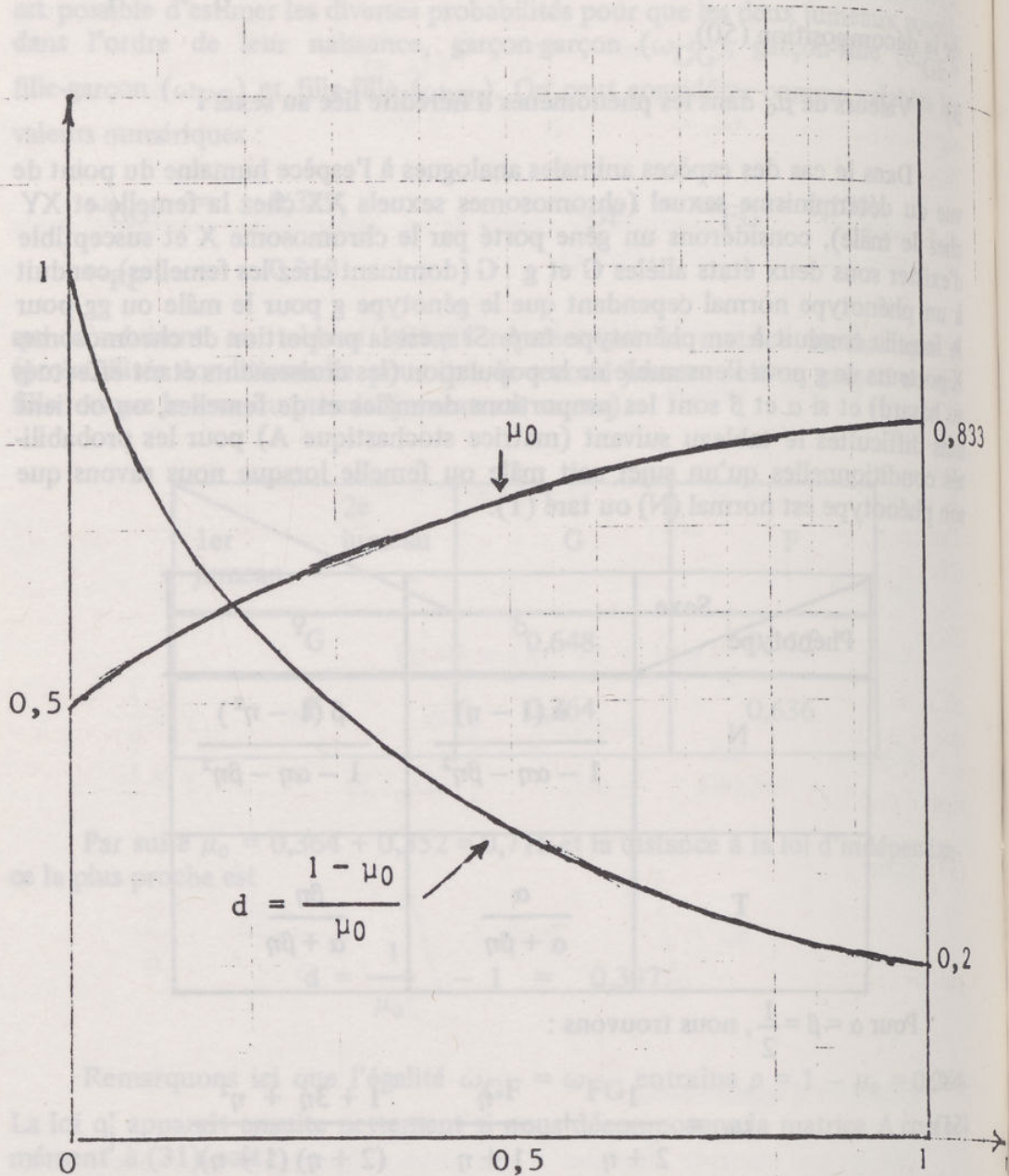


Fig. 2



c) Possibilité de tester les phénomènes de pléiotropie par un coefficient d'indépendance.

Toujours en biologie, il y a pléiotropie lorsqu'un gène comportant plusieurs états allèles participe simultanément à la réalisation de plusieurs caractères phénotypiques. Il conviendra de former ici la matrice stochastique A décrivant les probabilités conditionnelles des phénotypes possibles pour le deuxième caractère lorsque les phénotypes relatifs au premier caractère sont connus. Comme il ne peut y avoir indépendance totale entre les deux séries de phénotypes lorsqu'une liaison existe par suite de la pléiotropie, une différence significative de  $\mu_0$  (coefficient d'indépendance du deuxième caractère par rapport au premier) vis-à-vis de l'unité signalera un tel phénomène.

d) Emploi de  $\mu_0$  dans certains problèmes de mécanique quantique.

Notre dernier exemple concernera les distributions électroniques dans les atomes et les molécules. Pour l'état électronique fondamental de l'atome d'hélium, une approximation usuelle est celle du produit simple de deux orbitales atomiques identiques afin de représenter les deux électrons 1 s. C'est là un modèle d'indépendance totale conduisant à  $\mu_0 = 1$  lorsque nous étudions les densités de probabilité de présence dans tout l'espace du 2<sup>e</sup> électron en fonction de la position du premier électron, supposée préalablement fixée. Toutefois, d'autres types de fonctions propres approchées sont possibles, comme celle d'Eckart et Hylleraas [2], soit :

$$(52) \quad \psi(1, 2) = N [\exp(-\alpha r_1 - \beta r_2) + \exp(-\beta r_1 - \alpha r_2)],$$

où  $r_1$  et  $r_2$  représentent les distances des électrons 1 et 2 au noyau cependant que N est le facteur usuel de normalisation. L'ajustement des paramètres donne  $\alpha = 2,15$  et  $\beta = 1,19$  en unités atomiques de Hartree. Pour le calcul de  $\mu_0$ , il convient d'évaluer les densités de probabilités conditionnelles :

$$(53) \quad f(2/1) = \frac{\psi^2(1, 2)}{\int \psi^2(1, 2) d\tau_2}$$

puis de calculer :

$$(54) \quad \mu_0 = \int \inf_{(1)} [f(2/1)] d\tau_2$$

On trouve de cette manière  $\mu_0 = 0,737$  pour la fonction propre approchée (41).

Dans cette application à des variables aléatoires de nature continue (X et Y caractérisent les positions des électrons numérotés 1 et 2), nous vérifions



que les électrons du modèle proposé par Eckart et Hylleraas ne sont plus totalement indépendants bien que  $\mu_0$  reste relativement élevé. Notons d'ailleurs que la fonction propre rigoureuse devrait conduire à  $\mu_0 = 0$  car  $\inf [f(2/1)]$  (1)

reste constamment nul pour une bonne fonction propre, la répulsion coulombienne interdisant à deux électrons d'être simultanément présents dans un même élément de volume.

On trouve de cette manière  $\mu_0 \approx 0,737$  pour la fonction propre approchée

$$\mu_0 = \int_0^1 \ln |f(2/1)| d\tau$$

$$f(2/1) = \frac{\psi^2(1,2)}{\int \psi^2(1,2) d\tau} = 1$$

Dans cette application à des variables séparées de nature continue (X et Y caractérisent les positions des électrons numérotés 1 et 2), nous vérifions

Fig. 2

REFERENCES

- [1] BOREL (E.) *Eléments de la théorie des probabilités*, p. 155 Hermann édit. (1924)
- [2] cf. PAULING (L.) *Introduction to quantum mechanics* p. 224. Mc Graw Hill édit. (1935).

Reçu en Juin 1973.



Etude d'un coefficient d'indépendance d'une variable aléatoire Y par rapport  
à une variable aléatoire X

1973

### REFERENCES

1. H. DEBIEVE, *Revue de Statistique*, 1972, 1, 1-10.  
2. H. DEBIEVE, *Revue de Statistique*, 1973, 2, 1-10.  
3. H. DEBIEVE, *Revue de Statistique*, 1974, 3, 1-10.  
4. H. DEBIEVE, *Revue de Statistique*, 1975, 4, 1-10.  
5. H. DEBIEVE, *Revue de Statistique*, 1976, 5, 1-10.  
6. H. DEBIEVE, *Revue de Statistique*, 1977, 6, 1-10.  
7. H. DEBIEVE, *Revue de Statistique*, 1978, 7, 1-10.  
8. H. DEBIEVE, *Revue de Statistique*, 1979, 8, 1-10.  
9. H. DEBIEVE, *Revue de Statistique*, 1980, 9, 1-10.  
10. H. DEBIEVE, *Revue de Statistique*, 1981, 10, 1-10.

Reçu en Juin 1973.