

M. GIRAULT

**Liaisons, corrélation, régression**

*Mathématiques et sciences humaines*, tome 5 (1964), p. 11-22

[http://www.numdam.org/item?id=MSH\\_1964\\_\\_5\\_\\_11\\_0](http://www.numdam.org/item?id=MSH_1964__5__11_0)

© Centre d'analyse et de mathématiques sociales de l'EHESS, 1964, tous droits réservés.

L'accès aux archives de la revue « Mathématiques et sciences humaines » (<http://msh.revues.org/>) 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/>

M. GIRAULT

LIAISONS, CORRELATION, REGRESSION

## INTRODUCTION

1.1.- La physique classique a utilisé la notion de liaison fonctionnelle entre deux ensembles de grandeurs:

Exemple: Soit une barre métallique portée à différentes températures:  $t_1 t_2 \dots t_n$ : elle prend les longueurs  $l_1 l_2 \dots l_n$ .

Chaque fois que la barre est portée à la température  $t_2$  par exemple, elle prend la longueur  $l_2$ . Ainsi à toute température (d'un certain intervalle de températures) est associée une longueur et une seule. Telle est la notion de fonction ou de liaison fonctionnelle  $t \rightarrow l$ .

Ce type de modèle, bien adapté aux phénomènes étudiés en physique classique ne convient généralement pas dans les sciences biologiques, économiques, humaines.

1.2.- L'idée de liaison en probabilité a été introduite pour la première fois par Francis Galton dans ses travaux sur l'hérédité dans les espèces végétales puis animales. Les principaux résultats furent publiés en 1889 sous le titre "Natural inheritance".

L'auteur s'intéresse tout spécialement aux mesures  $X$  et  $Y$  d'une même partie du corps chez deux êtres dont l'un est fils de l'autre. Appelons "taille" la mesure retenue:

$$\left. \begin{array}{l} X_i \\ Y_j \end{array} \right\} \text{ est un couple de tailles observées } \left\{ \begin{array}{ll} X_i & \text{celle de l'individu (i)} \\ Y_j & \text{" " (j)} \end{array} \right.$$

où (j) est fils de (i).

Galton montre que la loi de fréquence des couples de nombres  $(X, Y)$  est bien représentée par une loi de distribution à deux dimensions généralisant celle de Laplace-Gauss et étudiée précédemment par Bravais(1) (Mémoire de 1846).

Le modèle utilisé par Galton est très fréquemment employé. Ce modèle présente des particularités qu'il faut bien comprendre et, pour cela, il importe de distinguer plusieurs notions fondamentales concernant la liaison en probabilité, notions qui dans le modèle particulier de Galton se trouvent soit confondues soit liées entre elles d'une manière particulière. En d'autres termes, il faut bien distinguer les propriétés générales de toute distribution des propriétés particulières de la distribution de Laplace-Gauss à plusieurs dimensions.

---

(1) Bravais ne semble pas s'être intéressé à la liaison en probabilité.

Remarque préliminaire: Les notions que nous allons étudier se rattachent à des distributions sur un espace à plusieurs dimensions (tout spécialement deux): ce sont des distributions sur des couples  $(X, Y)$  de nombres.

Ces distributions peuvent être indifféremment des distributions statistiques ou des distributions de probabilité; ces dernières elles-mêmes pouvant être discrètes ou continues.

a) Distribution statistique: échantillon de taille  $(n)$ : ensemble de  $n$  couples  $\{x_1, y_1\}; x_2, y_2), \dots, (x_n, y_n)$  chacun affecté de la "masse"  $\frac{1}{n}$

b) Distribution de probabilité

b.1 - Discrète: ensemble de couples  $(x_i, y_i)$

$$\begin{aligned} & \text{ou} \\ & x_i \in \{x_1, x_2, \dots, x_n, \dots\} \\ & y_j \in \{y_1, y_2, \dots, y_m, \dots\} \end{aligned}$$

le couple  $(x_i, y_j)$  étant affecté de la probabilité  $P_{ij}$

b.2 - Continue: ensemble de couples  $(x, y)$ , où  $x$  et  $y$  appartiennent à des ensembles continus (intervalles où même  $\mathbf{R}$ ).

Chaque pavé  $\left\{ \begin{array}{l} (x, x + \Delta x) \\ (y, y + \Delta y) \end{array} \right\}$  est affecté d'une probabilité.

En particulier tout pavé élémentaire  $\left\{ \begin{array}{l} (x, x + dx) \\ (y, y + dy) \end{array} \right\}$  est affecté de la probabilité  $f(x, y) dx dy$ .

Dans chaque cas, nous choisirons pour fixer les idées un type de distribution tantôt discret, tantôt continu; étant bien entendu que les notions étudiées se définissent pour les deux types de distributions.

#### 4.- LIAISON EN PROBABILITE (ou liaison stochastique)

4.1.- Considérons donc une épreuve à laquelle est associé le couple de nombres  $(x, y)$  distribués, par exemple, suivant une loi discrète.

$$\left\{ \begin{array}{l} x \in \{x_1, x_2, \dots, x_n\} = A \\ y \in \{y_1, y_2, \dots, y_m\} = B \end{array} \right. \quad (x, y) \in A \times B$$

Posons

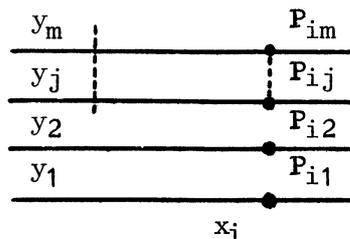
$$P_{ij} = P(X = x_i, Y = y_j) \quad i = 1, \dots, n; j = 1, \dots, m$$

Rappelons les définitions des lois marginales et des droites conditionnelles:

Distribution marginale de  $X$ : soit  $a_i = P(X = x_i) \quad i = 1, \dots, n$  On a, (axiome d'additivité):

$$a_i = \sum_j P_{ij}$$

que l'on note encore  $P_i$ .



La donnée des  $x_i$  et celle des  $p_i$  définit ( $p_i$ ) la distribution marginale de X.

Distribution marginale de Y : on a de façon analogue

$$b_j = p_{.j} = \sum_i p_{ij}$$

Les  $y_j$  et les  $p_{.j}$  définissent la distribution marginale de Y.

Remarque: Si la distribution est continue, de densité superficielle  $f(x, y)$ . Les densités marginales sont :

$$a(x) = \int f(x, y) dy \quad \text{pour X}$$

$$\text{et } b(y) = \int f(x, y) dx \quad \text{pour Y}$$

ou, si l'on préfère, les lois élémentaires sont respectivement

$$\text{pour X} \quad a(x) dx = \int_y f(x, y) dx dy$$

$$\text{pour Y} \quad b(y) dy = \int_x f(x, y) dx dy$$

## 2.2.- Distribution conditionnelle de Y si $X = x_i$

On effectue l'épreuve donnant le couple (X, Y) et l'on sait que  $X = x_i$ . La probabilité d'avoir  $Y = y_j$  est alors

$$b_{(j)}^{(i)} = \frac{p_{ij}}{p_i} \quad \text{notée} \quad \text{Prob.} \left\{ Y = y_j / X = x_i \right\}$$

On aurait de façon analogue:

$$\text{Prob.} \left\{ X = x_i / Y = y_j \right\} = a_i^{(j)} = \frac{p_{ij}}{p_{.j}}$$

Et, dans le cas d'une distribution continue:

loi élémentaire de Y si  $X = x$  :

$$b_{(y)}^{(x)} dy = \frac{f(x, y) dx dy}{a(x) dx} = \frac{f(x, y)}{a(x)} dy$$

## 2.3.- Dépendance en probabilité

On peut considérer non seulement les lois conditionnelles de Y si X prend une valeur particulière.

$$\left[ \begin{array}{l} X = x_i \quad \text{dans le cas discret} \\ x < X < x + dx \quad \text{dans le cas continu} \end{array} \right]$$

mais les lois conditionnelles de Y lorsque X appartient à un sous-ensemble de valeurs possibles: un intervalle  $(x_0, x_1)$  par exemple.

On dit que Y dépend de X (en probabilité) si les lois conditionnelles de Y dépendent des conditions imposées à X. Il y a liaison en probabilité si

$$b_{(j)}^{(i)} \text{ dépend de } a_i$$

## 2.4.- Indépendance en probabilité

Au contraire on dit de Y qu'il est indépendant de X (en probabilité) si toute loi conditionnelle de Y est indépendante des conditions imposées à X; c'est-à-dire si toutes les lois conditionnelles sont identiques.

On aura en particulier:

(distribution discrète):  $b_{(j)}^{(1)} = b_{(j)}^{(2)} = \dots = b_{(j)}^{(n)}$  et cela quel que soit j.

soit :

$$\frac{p_{1j}}{p_{1.}} = \frac{p_{2j}}{p_{2.}} = \frac{p_{3j}}{p_{3.}} = \dots = \frac{p_{nj}}{p_{n.}} = \frac{\sum_k p_{kj}}{\sum_k p_{k.}} = p_{.j}$$

d'où :

$$\forall i, \forall j : \boxed{p_{ij} = p_{i.} \cdot p_{.j}} \quad (1)$$

De la même manière, dans le cas continu, on démontre qu'il faut avoir, quels que soient x et y

$$\boxed{f(x,y) = a(x) \cdot b(y)} \quad (2)$$

Des relations (1) et (2) on déduit plusieurs conséquences. On a simplement supposé, pour obtenir ces relations, que les lois conditionnelles de Y pour des conditions élémentaires de X  $[X = x \text{ ou } x < X < x + dx]$  étaient indépendantes de x :

On déduit de (1) ou de (2)

- Que toute loi conditionnelle de Y pour des conditions exprimées sur X est identique à la loi marginale de Y,
- Que toute loi conditionnelle de X pour des conditions exprimées sur Y est identique à la loi marginale de X.

On exprime cette situation en disant que les v.a. X et Y sont mutuellement indépendantes en probabilité.

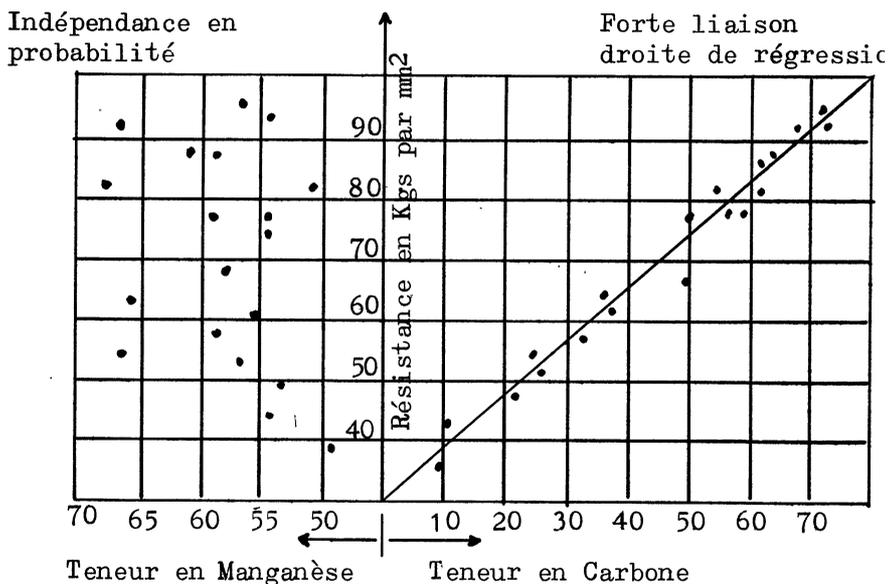
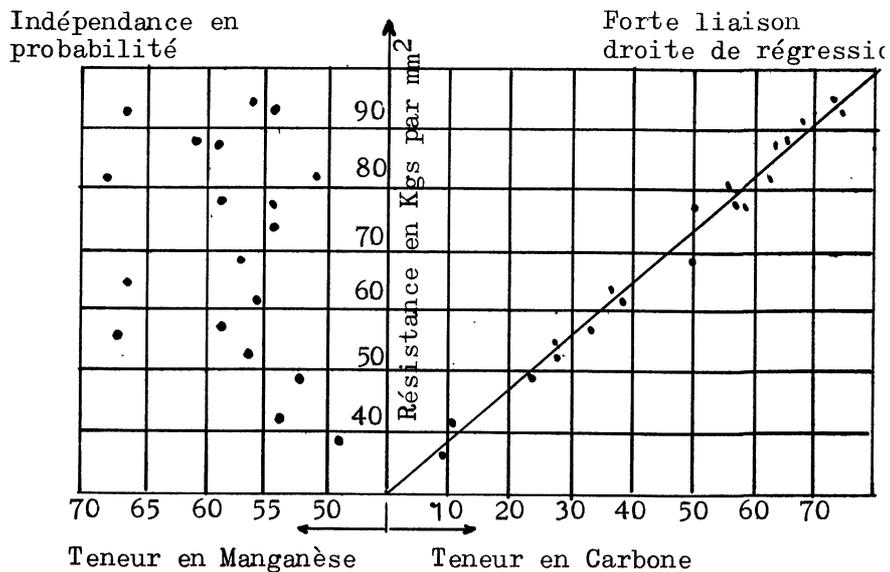
Dans ce cas, la distribution est complètement déterminée par la donnée des distributions marginales.

Une connaissance quelconque sur X, ne modifie en rien l'information que l'on possède sur Y.

## 3.- MOMENTS DU 1<sup>er</sup> ORDRE - Ligne de régression

Pour varier les notations, nous supposons ici que la distribution est continue, de densité superficielle  $f(x,y)$ .

Illustration:



La loi élémentaire de Y

$$\text{si } X = x \quad \text{est } b_{(y)}^{(x)} \quad dy = \frac{f(x,y)}{a(x)} \quad dy$$

$$\text{où } a(x) = \int_y f(x,y) \quad dy$$

La moyenne de (Y si X = x) est (quand elle existe)

$$m^{(x)} = y(x) = \int_y y \quad b_{(y)}^{(x)} \quad dy$$

$\bar{y}(x)$  est une fonction de x; sa courbe figurative est dite ligne de régression de Y en x.

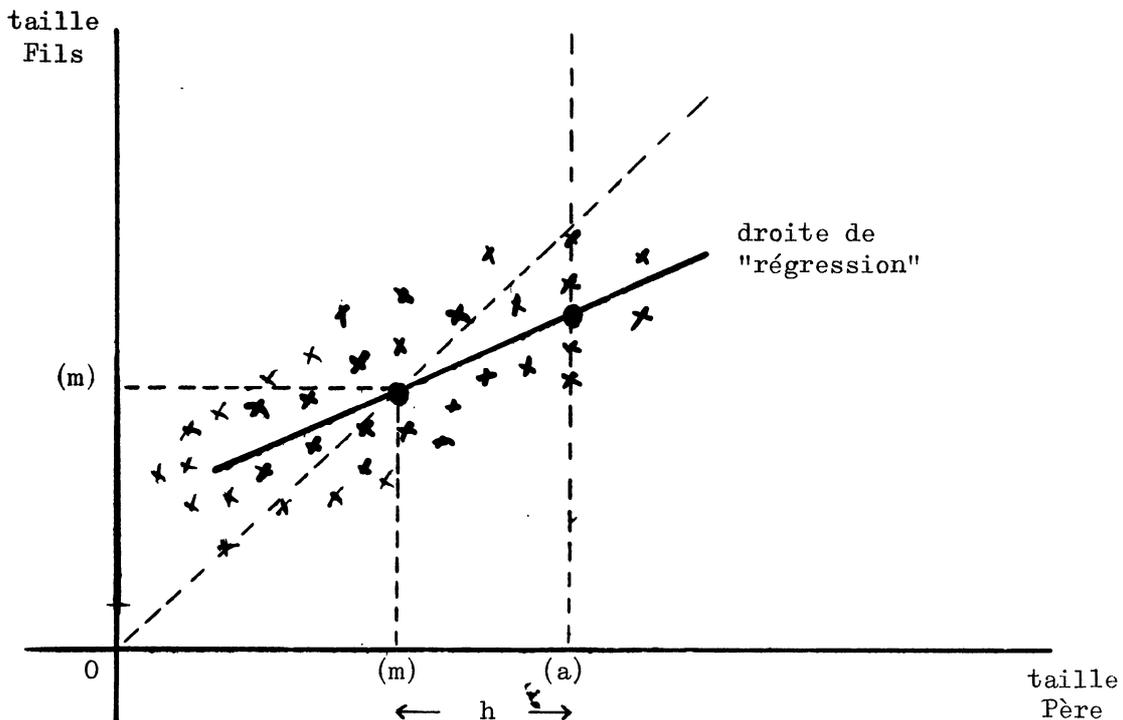
On définit d'une manière analogue la ligne de régression de X en y; c'est le lieu des moyennes conditionnelles de X pour  $Y = y$ .

Si les v.a. X et Y sont indépendantes, les lignes de régression sont des droites parallèles aux axes.

Une ligne de régression (de Y en x par exemple) est un indice de la liaison en probabilité de Y avec le "facteur" X. L'origine du terme "régression" est la suivante.

### Origine du terme "régression"

Etude de l'hérédité des tailles (F. GALTON puis K. PEARSON).



Soit Y la taille d'un individu, X celle de son père. Les couples (X,Y) sont décrits par une loi de probabilité où la ligne de régression de Y en x est une droite de pente r positive inférieure à 1.

Soit m la taille moyenne de la population étudiée. Les "grands" individus (de taille  $a = m + h$  par exemple) ont des fils dont les tailles sont en moyenne:  $a + rh$ . Ces fils sont des individus (en moyenne) grands mais ce caractère "grand" est moins accusé qu'il ne l'est chez leurs pères: on dit que ce caractère est en régression d'une génération à la suivante.

#### 4.- MOMENTS DU SECOND ORDRE

4.1.- Pour décrire la liaison des grandeurs X et Y, les indices donnés par les moments du premier ordre sont insuffisants. On s'intéresse donc aux moments du second ordre (en supposant là encore qu'ils existent).

$$\left. \begin{aligned} \text{Soit } \mu_{20} &= \sigma_x^2 \text{ la variance de X} \\ \mu_{02} &= \sigma_y^2 \text{ la variance de Y} \\ \mu_{11} &= E \left[ (X-\bar{X})(Y-\bar{Y}) \right] \end{aligned} \right\} \begin{array}{l} \text{moments du 2ème ordre} \\ \text{du couple X, Y.} \end{array}$$

Soit enfin  $\left[ \sigma_{(y)}^{(x)} \right]^2$  la variance conditionnelle de Y / si X = x.

On appelle fonction scédastique (la fonction  $z(x) = \sigma_{(y)}^{(x)}$ ). Si  $\sigma_{(y)}^{(x)} = c^{te}$  on dit qu'il y a homoscédasticité de Y par rapport à x.

#### 4.2.- Indices de corrélation

Plusieurs indices ont été proposés, mais en fait le plus utilisé est le coefficient de corrélation

$$\rho(x, y) = \frac{\mu_{11}}{\sigma_x \sigma_y}$$

La signification du coefficient de corrélation s'interprète commodément sur l'ellipse des variances que nous définirons plus loin.

Un autre indice intéressant est le rapport de corrélation que nous allons maintenant définir.

#### 4.3.- Analyse de variance - Rapport de corrélation

Sur le couple aléatoire (X,Y) on considère:

la moyenne générale de Y, soit  $E(Y) = \bar{y} = \iint yf(x,y) dx dy$

les moyennes conditionnelles de Y si X = x

$$\text{soit } E \left[ \frac{Y}{X = x} \right] = \bar{y}_{(x)} = \int yf(x,y) dy.$$

Soit (x,y) un couple de valeurs obtenues.

L'écart de y peut se décomposer de la manière suivante:

$$y - \bar{y} = (y - \bar{y}_x) + (\bar{y}_x - \bar{y}) \quad (1)$$

La variance de Y est  $\sigma^2 = E \left[ (y - \bar{y})^2 \right] = \iint (y - \bar{y})^2 f(x,y) dx dy$

En remplaçant  $(y - \bar{y})$  par sa valeur (1) on obtient:

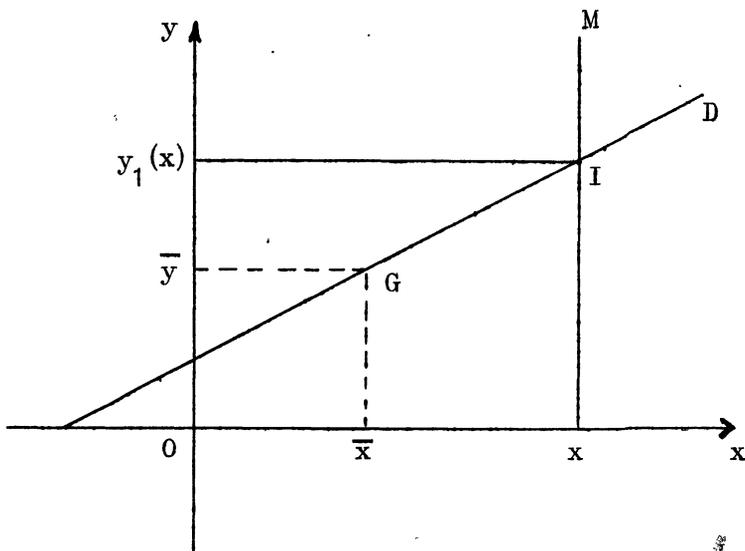


Son équation est  $y_1(x)$  tel que  $\iint [y - y_1(x)]^2 f(x,y) dx dy$  soit minimum.

La recherche de ce minimum montre que D passe par G, centre de la distribution. D a pour pente

$$\rho \frac{\sigma_y}{\sigma_x}$$

Remarque: Certains auteurs appellent la droite D: "droite de régression" de Y en x.



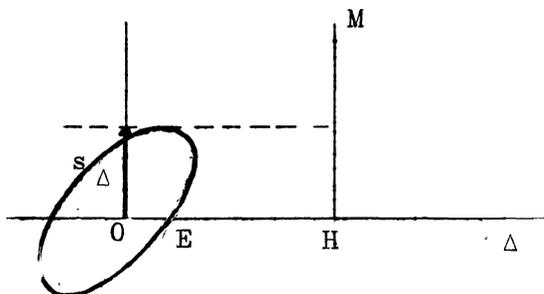
#### 4.5.- Ellipse des variances

Toutes les notions précédentes: 4-1, 4-2, 4-3 et 4-4 se rattachant aux moments du 2ème ordre. Les moments des formes linéaires en X, Y (expressions de la forme  $aX + bY$ ) s'expriment facilement à partir de l'ellipse des variances.

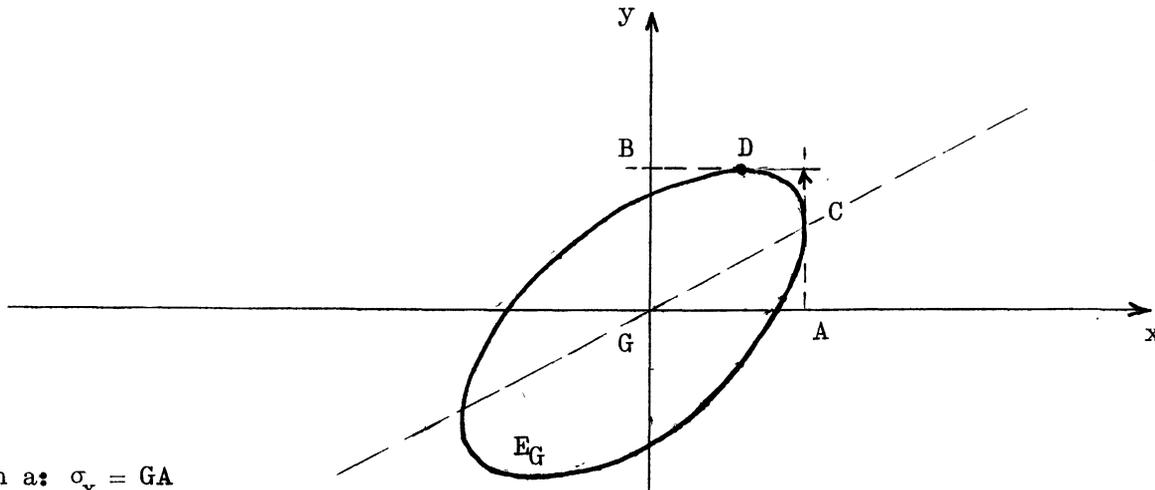
Soit une distribution sur le plan (ox,oy) admettant des moments du second ordre:  $E(X^2)$  et  $E(Y^2)$  (donc aussi  $E(X.Y)$ ).

Il existe une ellipse E de centre O qui a la propriété suivante:

Soit  $\Delta$  une droite issue de l'origine. Le moment de la distribution par rapport à  $\Delta$  (espérance mathématique  $\overline{MH^2}$ ) est égal à  $s_\Delta^2$  ou  $s_\Delta$  est la distance de O à la tangente à E parallèle à  $\Delta$ .



En particulier, si l'on rapporte la distribution à son centre G l'ellipse associée:  $E_G$  est dite "ellipse des variances"



On a:  $\sigma_x = GA$   
 $\sigma_y = GB$

La droite GC, de pente  $\rho \frac{\sigma_y}{\sigma_x}$  est la droite des moindres carrés de Y en x. La droite GD, de pente  $\frac{1}{\rho} \frac{\sigma_y}{\sigma_x}$  est la droite des moindres carrés de X en y.

$\rho = 1 \iff$  droites GC et GD confondues:  $E_G$  est une ellipse aplatie: X et Y sont liés par une relation fonctionnelle linéaire.

$\rho = 0 \iff$  ellipse des variances a pour axes ox et oy.

Cette circonstance peut se produire non seulement lorsque X et Y sont liés en probabilité mais également lorsqu'il existe entre eux une liaison fonctionnelle non monotone.

Donc: Lorsqu'on étudie des distributions très générales:  $\rho$  est un mauvais indice de corrélation; le rapport de corrélation R est plus significatif.

### 5.- CAS DES DISTRIBUTIONS A REGRESSION LINEAIRE

Soit  $\mathcal{D}_1$  la famille des distributions sur (X,Y) telles que la courbe de régression de Y en x (voir 3) soit une droite. On dit alors que la régression (de Y en X) est linéaire.

On démontre alors que:

La droite des moindres carrés s'identifie à la droite de régression.

- La variance conditionnelle de Y si x est alors  $\sigma_2^2 = \overline{GF^2}$

- Enfin le rapport de régression R s'identifie au coefficient de corrélation:

$$\rho = \frac{\mu_{11}}{\sigma_x \sigma_y} = 1 - \frac{\sigma_2^2}{\sigma^2}$$

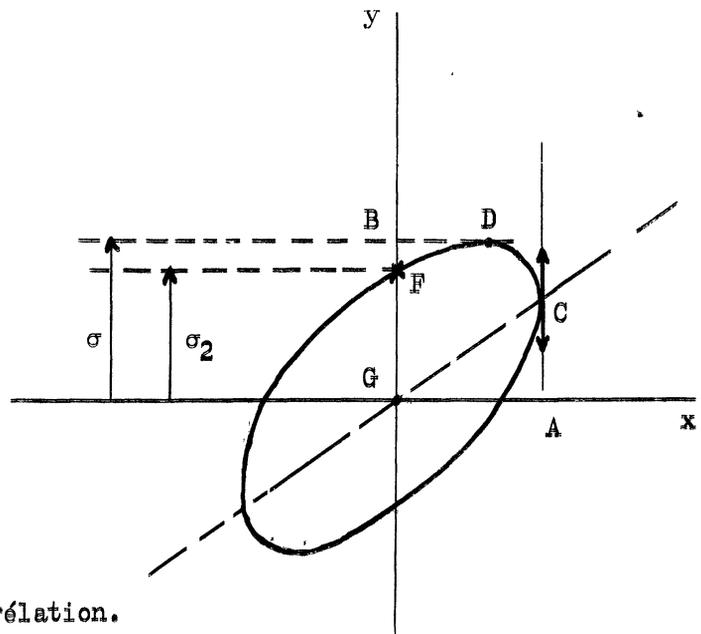
où  $\sigma^2 = \sigma_y^2$  est la variance totale de Y.

$\sigma_2^2$  est la variance conditionnelle de Y.

Conclusion: Au sein de la famille des distributions  $\mathcal{D}_1$ , le coefficient de corrélation

$$\rho = \frac{\mu_{11}}{\sigma_x \sigma_y}$$

est un indice significatif de la corrélation.



### 6.- UN MODELE PARTICULIER DE LIAISON: LA DISTRIBUTION DE LAPLACE GAUSS A DEUX DIMENSIONS

La distribution dite "de Laplace-Gauss à deux dimensions", étudiée pour la première fois par BRAVAIS (1846), fut utilisée par F. GALTON; K. PEARSON et après eux par de très nombreux auteurs. Elle présente une structure relativement simple et est susceptible de nombreuses applications. Il importe toutefois, avant de l'adopter comme modèle, de faire preuve d'esprit critique.

Expression de la loi de Laplace-Gauss à deux dimensions: loi élémentaire:

$$(1) \quad \frac{1}{2n} \frac{1}{\sigma_x \sigma_y \sqrt{1-\rho^2}} e^{-\frac{1}{2} \frac{1}{1-\rho^2} \left[ \frac{x^2}{\sigma_x^2} - 2 \frac{xy}{\sigma_x \sigma_y} + \frac{y^2}{\sigma_y^2} \right]} dx dy$$

Les courbes réunissant les points d'égale densité sont des ellipses homothétiques et concentriques. Le centre commun à toutes ces ellipses est le centre de la distribution, de coordonnées  $m_x = E(X)$  et  $m_y = E(Y)$ . On a choisi ce point pour origine des axes pour écrire la relation (1).

L'une de ces ellipses est l'ellipse des variances (E) dont l'équation est :

$$(2) \quad \boxed{\frac{x^2}{\sigma_x^2} - 2\rho \frac{xy}{\sigma_x \sigma_y} + \frac{y^2}{\sigma_y^2} = 1 - \rho^2}$$

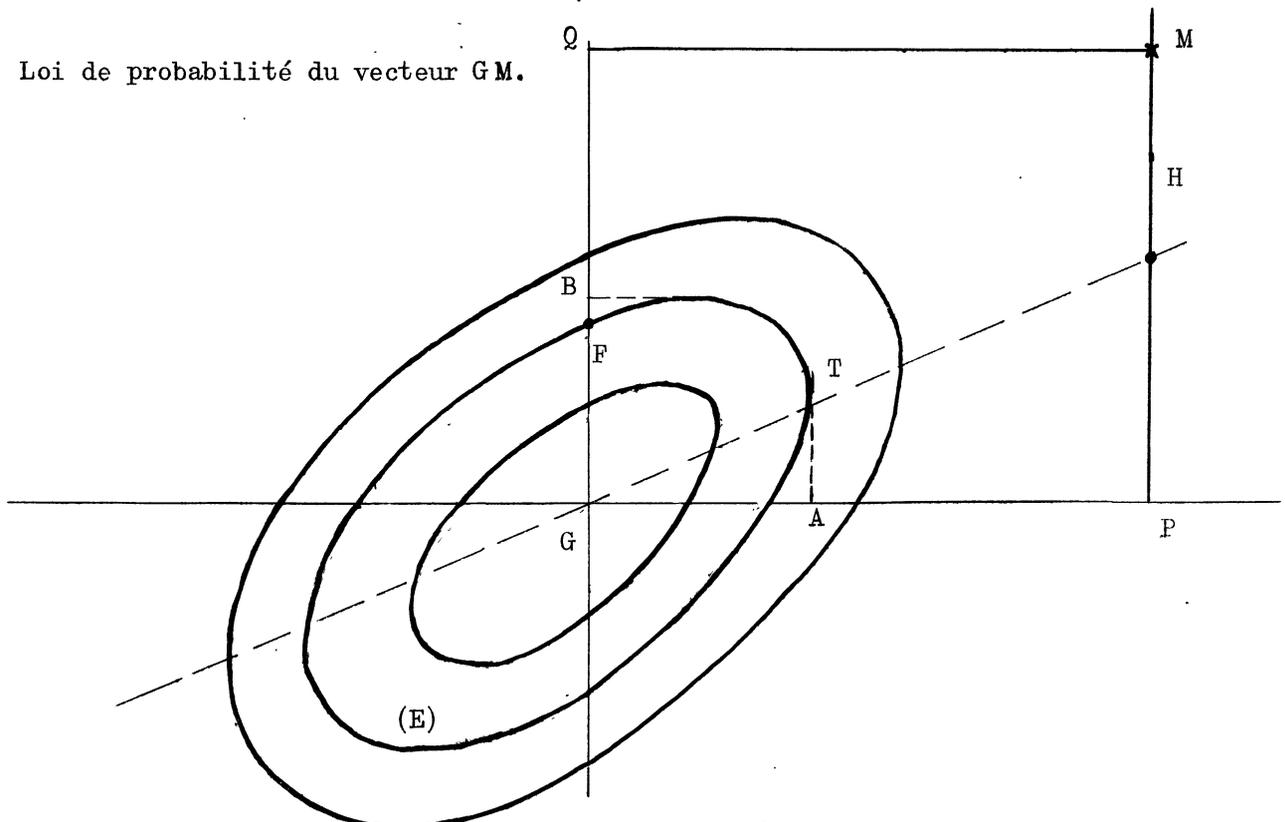
Dans ce modèle; les courbes de régression sont des droites; on dit que les régressions sont linéaires et par conséquent toutes les particularités (5) sont réalisées.

Ces propriétés (5) se rattachent aux moments du 1er et du 2e ordre; et ne font pas état des lois de probabilité des variables étudiées.

La famille des distributions de Laplace-Gauss étudiées maintenant ne contient que 5 paramètres:  $m_x$ ;  $m_y$ ;  $\sigma_x$ ;  $\sigma_y$  et  $\rho$  par exemple. Or la donnée de l'ellipse E des variances définit ces 5 paramètres. Donc dans ce cas, non seulement les moments du 2ème ordre des formes linéaires en (X,Y) sont déterminées, mais les lois de probabilité de ces formes ainsi que les lois des variables liées.

### Distribution de Laplace Gauss à 2 dimensions

Loi de probabilité du vecteur GM.



Les trois ellipses sont des courbes joignant les points de même densité (courbes de niveau de la surface figurative des densités). E est l'ellipse des variances.

Le vecteur aléatoire  $\vec{GP}$  obéit à la loi de L.G. (centre G, écart type GA)

Le vecteur aléatoire  $\vec{GQ}$  obéit à la loi de L.G. (centre G, écart type GB)

Si P est donné :  $\vec{PM}$  obéit à la loi de L.G. (centre H, écart type GF).

La droite GT est à la fois : la ligne de régression de Y en x et la droite des moindres carrés de Y en x.

La v.a. X obéit à la loi de Laplace-Gauss de moyenne  $m_x$  et d'écart type  $\sigma_x$ . De même pour Y.

Si  $X = x$  la v.a. Y obéit à la loi de Laplace-Gauss de moyenne :

$$\bar{y}(x) = m_y + \rho \frac{\sigma_y}{\sigma_x} (x - m_x)$$

son écart type est :

$$\sigma_{y/x} = \sigma_y \sqrt{1 - \rho^2}$$

Cet écart-type est indépendant de x :

il y a homos célasticité

Il y a même plus : la loi de l'écart de Y (rapporté à la moyenne conditionnelle) c'est-à-dire la loi de

$$Y_1 = Y - y_1(x) \quad \text{où} \quad y_1(x) = E[Y \text{ si } X = x]$$

est indépendante de x.

On dit alors (définition due à S. BERNSTEIN) que la corrélation entre X et Y est dure.

-----

#### BIBLIOGRAPHIE

- R. FERON "Information - Régression - Corrélation" - thèse - Publication ISUP - Vol. V fasc. 34 - 1956 -
- G. DARMOIS - Sur certaines formes de liaisons en probabilité - Colloque CNRS - Le Calcul des Probabilités et ses applications - Lyon - 1948 - Publ. CNRS -