

N° 28

23034

COURS DE STATISTIQUES

. D. DACUNHA CASTELLE

Réédition de 1973 - 1974



Fascicule 1

- CHAPITRE 0 - Statistique descriptive.
- CHAPITRE I - Modèles statistiques. Théorie de la décision.
- CHAPITRE II - Théorie élémentaire de l'estimation et sondages.
- CHAPITRE III - Introduction à la théorie des tests et des régions de confiance.
- CHAPITRE IV - La statistique non paramétrique.
- CHAPITRE V - Estimations et tests en variable gaussienne.
- CHAPITRE VI - Le test du  $\chi^2$ .

## CHAPITRE 0

### STATISTIQUE DESCRIPTIVE.

#### A - INTRODUCTION.

La statistique descriptive est une méthode de description des données recueillies à propos de l'étude de certains phénomènes d'ordre économique, sociologique ou expérimental. La description se fait essentiellement dans deux directions.

a) Direction géométrique. Ceci consiste d'abord à classer les données et à les disposer de manière la plus parlante possible (tableaux, images, graphiques etc...). Dans un deuxième temps (cf. fascicule 3), on utilise des techniques plus élaborées, qui consistent par exemple à projeter les données supposées ou rendues numériques (à valeurs dans  $\mathbb{R}^k$ ) sur des plans convenablement choisis, et à interpréter les figures ainsi obtenues. Cette analyse des données moderne (car nécessitant de puissants moyens de calculs) est donc assez voisine de la géométrie descriptive ou de la géométrie côtelée (dessin de machines). La différence essentielle vient de ce que les règles de lecture de figures obtenues sont complexes, et amènent à des déductions de caractère probabiliste.

b) Direction numérique. Elle consiste à résumer l'ensemble des données par la valeur de certaines grandeurs numériques typiques.

Les deux directions ne sont pas évidemment exclusives, et sont toujours utilisées simultanément.

Population. Les données recueillies portent (terme démographique) sur une population, c'est-à-dire sur un certain ensemble d'unités statistiques ou d'individus, éléments de cette population qui doit donc être définie avec précision.

Caractères.

Chacun des individus de la population peut être décrit par un ou plusieurs caractères.

Par exemple pour le personnel d'une entreprise : sexe, âge, qualification, nationalité, nombre de personnes à charge, appartenance à tel syndicat...

Les caractères étudiés peuvent être qualitatifs ou quantitatifs.

Qualitatif si ses diverses modalités (modalité : différentes situations possibles du caractère) ne sont pas mesurables

ex : sexe, couleur de cheveux, qualification ...

Les modalités d'un caractère qualitatif sont les différentes rubriques d'une nomenclature établie de telle façon que l'individu figure dans une et une seule rubrique (rubriques exhaustives et incompatibles). Chaque nomenclature peut être plus ou moins détaillée.

Quantitatif si les diverses modalités sont mesurables ou repérables, c.à.d. si à chacune des modalités est attachée un nombre. Ce nombre est la variable statistique.

Variables statistiques.

Discrètes, si ses valeurs possibles sont isolées (le nombre de valeurs possibles pouvant être fini ou infini).

Continues, si ses valeurs possibles sont en nombre infini et en général à valeur dans un intervalle de  $\mathbb{R}$ .

Souvent, la distinction est difficile à faire entre variable continue et variable discrète.

Par exemple, toute mesure est discrète du fait d'une précision toujours limitée, alors que la nature intrinsèque de la variable (par exemple le diamètre d'une pièce) en fait une variable continue. Réciproquement, on considèrera qu'une variable qui peut prendre un très grand nombre de valeurs possibles, sera une variable continue. Par exemple :

"le salaire d'un ouvrier"

"le bénéfice annuel d'une entreprise"

"le diamètre mesuré au  $1/100$  de mm d'une pièce".

Dans le cas du salaire, la plus faible unité monétaire est le centime. Ainsi, assimiler le salaire à une variable continue, c'est assimiler la plus faible unité monétaire à la précision des mesures.

Pour étudier une variable statistique continue, on définira des classes de valeurs possibles, pouvant avoir une amplitude constante ou variable.

Le nombre de classes à adopter dépend de la précision des mesures et des effectifs de la population étudiée

- . trop de classes  $\Rightarrow$  irrégularités accidentelles provenant du faible nombre d'individus par classe.
- . pas assez de classes  $\Rightarrow$  perte d'information.

On justifiera pourquoi, plus tard, il est intéressant d'obtenir dans les différentes classes des effectifs comparables.

Concluons cette introduction par une remarque importante. La statistique descriptive porte sur une population donnée. Le problème de l'échantillonnage et du sondage (cf. chapitre 1) ne se pose pas. Il n'y a pas de modèle statistique en statistique descriptive, uniquement un ensemble de données, sans structuration a priori.

NOMENCLATURE INSEE  
DES CATEGORIES SOCIO-PROFESSIONNELLES

0. AGRICULTEURS EXPLOITANTS.

0. 0. Agriculteurs exploitants

1. SALARIES AGRICOLES.

1. 0. Salariés agricoles

2. PATRONS DE L'INDUSTRIE ET DU COMMERCE.

2. 1. Industriels  
2. 2. Artisans  
2. 3. Patrons pêcheurs  
2. 6. Gros commerçants  
2. 7. Petits commerçants

3. PROFESSIONS LIBÉRALES ET CADRES SUPÉRIEURS.

3. 0. Professions libérales  
3. 2. Professeurs; professions littéraires et scientifiques  
3. 3. Ingénieurs  
3. 4. Cadres administratifs supérieurs

4. CADRES MOYENS.

4. 1. Instituteurs; professions intellectuelles diverses  
4. 2. Services médicaux et sociaux  
4. 3. Techniciens  
4. 4. Cadres administratifs moyens

5. EMPLOYÉS.

5. 1. Employés de bureau  
5. 3. Employés de commerce

6. OUVRIERS.

6. 0. Contremaîtres  
6. 1. Ouvriers qualifiés  
6. 3. Ouvriers spécialisés  
6. 5. Mineurs  
6. 6. Marins et pêcheurs  
6. 7. Apprentis ouvriers  
6. 8. Manœuvres

7. PERSONNELS DE SERVICE.

7. 0. Gens de maison  
7. 1. Femmes de ménage  
7. 2. Autres personnels de service

8. AUTRES CATÉGORIES.

8. 0. Artistes  
8. 1. Clergé  
8. 2. Armée et Police

9. PERSONNES NON ACTIVES.

9. 1. Étudiants et élèves  
9. 2. Militaires du contingent  
9. 3. Anciens agriculteurs (exploitants et salariés)  
9. 4. Retirés des affaires  
9. 5. Retirés du secteur public  
9. 6. Anciens salariés du secteur privé  
9. 9. Autres personnes non actives.

B - DISTRIBUTIONS STATISTIQUES A UN CARACTERE

1 - TABLEAUX STATISTIQUES.

Soit une population de  $n$  individus, décrite par le caractère  $C$  prenant  $k$  modalités  $C_1, C_2, \dots, C_k$ .

$n_i$  est le nombre d'individus présentant la modalité  $C_i$  ou effectif de la modalité  $C_i$ .

$f_i = \frac{n_i}{n}$  est la fréquence de la modalité  $C_i$ .

Le tableau suivant décrit la population étudiée,

Modalité de $C$	Effectif	ou fréquence (Total connu)
$C_1$	$n_1$	
$C_2$	$n_2$	
$C_i$	$n_i$	
$C_k$	$n_k$	
Total	$n$	

Caractère qualitatif

Même forme que celle-ci-dessus indiquée.

⇒

Exemple : Répartition des Etrangers qui vivent en France (Recensement général de 1962)

(Modalités)	Nombre de pièces à rebuter par lot $\alpha$	Nombre correspondant de lots ou effectif
	1	2
	2	9
	3	14
	4	20
	5	18
	6	15
	7	9
	8	6
	9	4
	10	2
	11	1
Total		100

(Modalités)	Effectif
Nationalité	
Italiens	645 000
Belges	78 000
Allemands	46 000
Hollandais, Luxembourgeois	17 000
Espagnols	431 000
Portugais	50 000
Polonais	177 000
Autres Européens	124 000
Algériens	335 000
Marocains, Tunisiens	54 000
Autres Etrangers	194 000
Total	2 151 000

Caractère quantitatif discret :

Sur 100 lots, chacun de 100 pièces, on a observé le nombre de pièces défectueuses.

La variable statistique est le nombre de lots ayant donné  $\alpha$  pièces défectueuses.

Cas d'une statistique continue.

Dans ce cas, les modalités du caractère sont définies à partir des classes de valeurs possibles (la division en classe est arbitraire).

Exemple : Distribution des ouvriers ayant travaillé toute l'année chez le même employeur selon le revenu annuel (année : 1952)

Tranche de salaire annuel net en milliers de francs	Nombre d'ouvriers
Moins de 100 .....	1 721
de 100 à moins de 125 .....	2 413
--- 125 --- 150 .....	4 342
--- 150 --- 175 .....	8 234
--- 175 --- 200 .....	13 300
--- 200 --- 225 .....	16 053
--- 225 --- 250 .....	16 774
--- 250 --- 300 .....	33 251
--- 300 --- 350 .....	29 211
--- 350 --- 400 .....	22 453
--- 400 --- 500 .....	24 005
--- 500 --- 600 .....	9 477
--- 600 --- 800 .....	4 093
--- 800 --- 1 000 .....	443
--- 1 000 --- 1 500 .....	125
--- 1 500 --- 2 000 .....	12
--- 2 000 --- 5 000 .....	14
Total .....	185 951

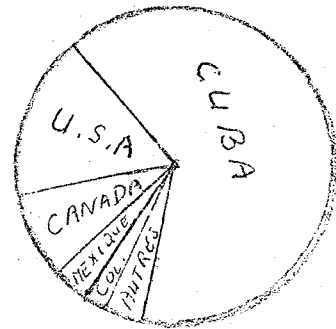
Le problème concret du choix des classes est souvent très important pour les applications numériques. Pour calculer une moyenne il n'est pas toujours souhaitable d'assimiler une classe à son "milieu".

2 - REPRESENTATION GRAPHIQUE

1. Caractères qualitatifs.

Par secteur angulaire : Exemple : Médailles obtenues aux jeux Pan américains de Cali (1971) par pays et par tête d'habitant.

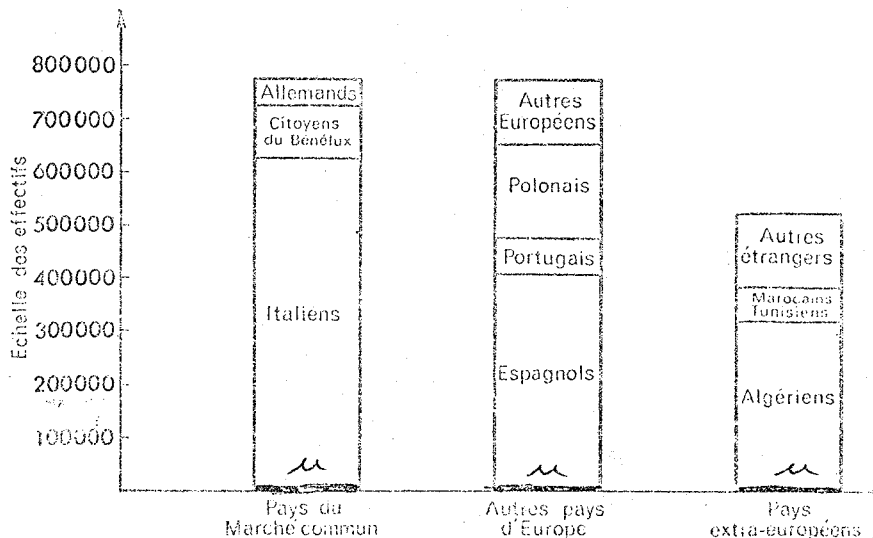
Il y a proportionnalité des aires aux effectifs.



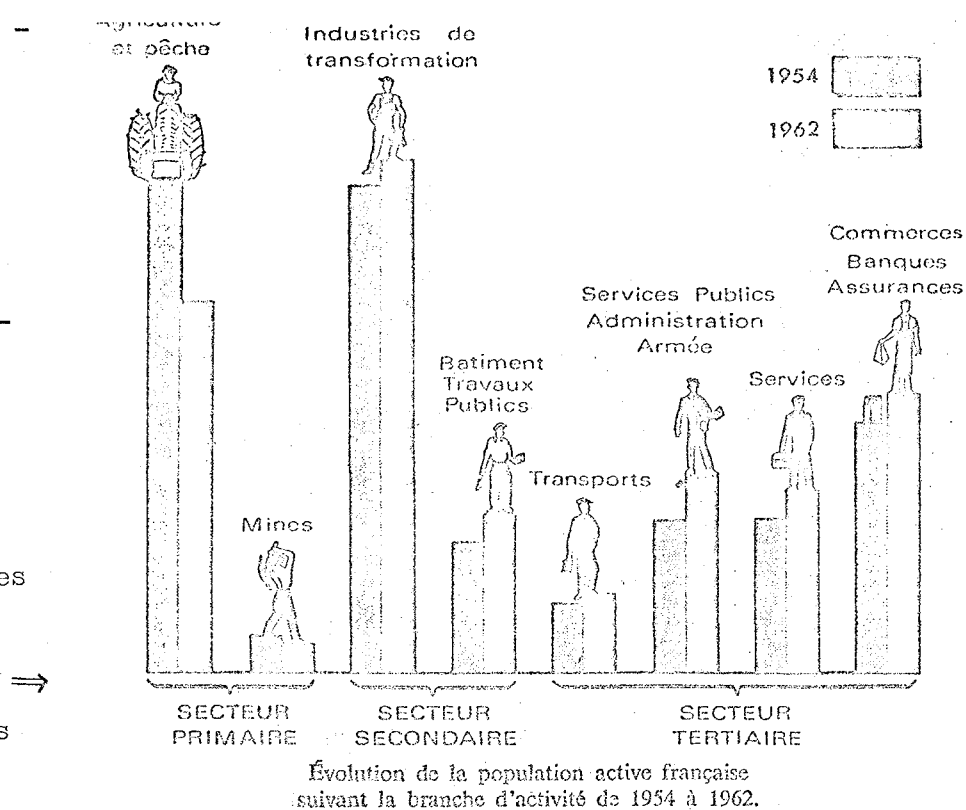
Par tuyaux d'orgue.

Exemple : répartition des étrangers en France (1962)

Les tuyaux d'orgue ont une base constante et une hauteur proportionnelle à l'effectif correspondant. Leur surface est donc proportionnelle à l'effectif.



La représentation par tuyaux d'orgue est particulièrement agréable quand on veut figurer une variation entre deux variables statistiques prenant les mêmes modalités



(Extrait de : Documents pour la classe I P N 1964)

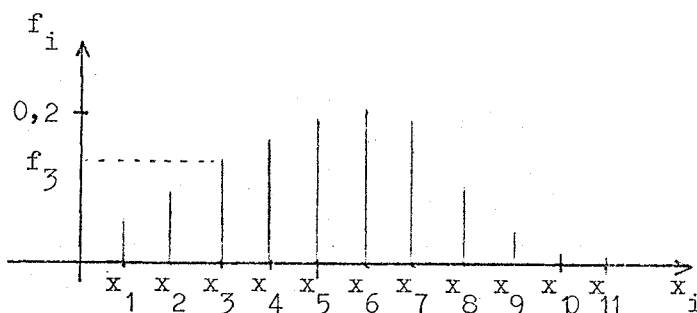
2. Caractères quantitatifs.

→ Variable discrète ou discretisée

Diagramme en bâtons

on figure la fréquence  $f_i$  pour la modalité  $x_i$

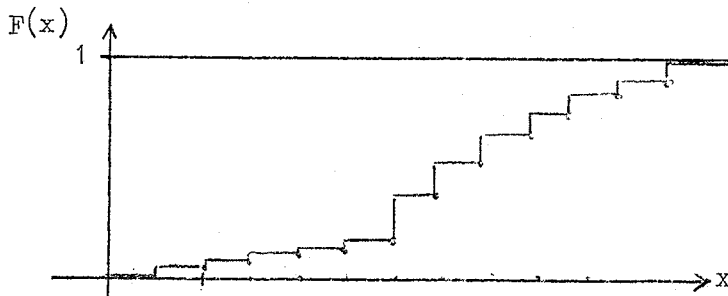
$$(\sum f_i = 1)$$



Courbe cumulative ou fonction de répartition

$$F(x) = \sum_{j=1}^i f_j$$

pour  $x_i < x < x_{i+1}$



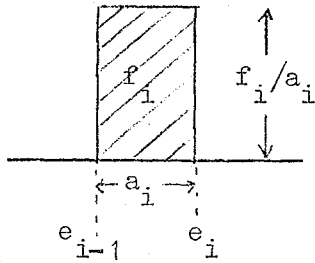
(Discontinuité aux points  $x_i$ , continuité à droite en  $x_i$ ).

Dans le cas continu étudié (premier tableau, page 2), les amplitudes des différentes classes sont des multiples de 25.000 frs.

Les effectifs sont directement comparables si ils correspondent à la même amplitude.

Si les amplitudes sont différentes, les surfaces figurées sur une classe devront représenter la fréquence  $f_i$

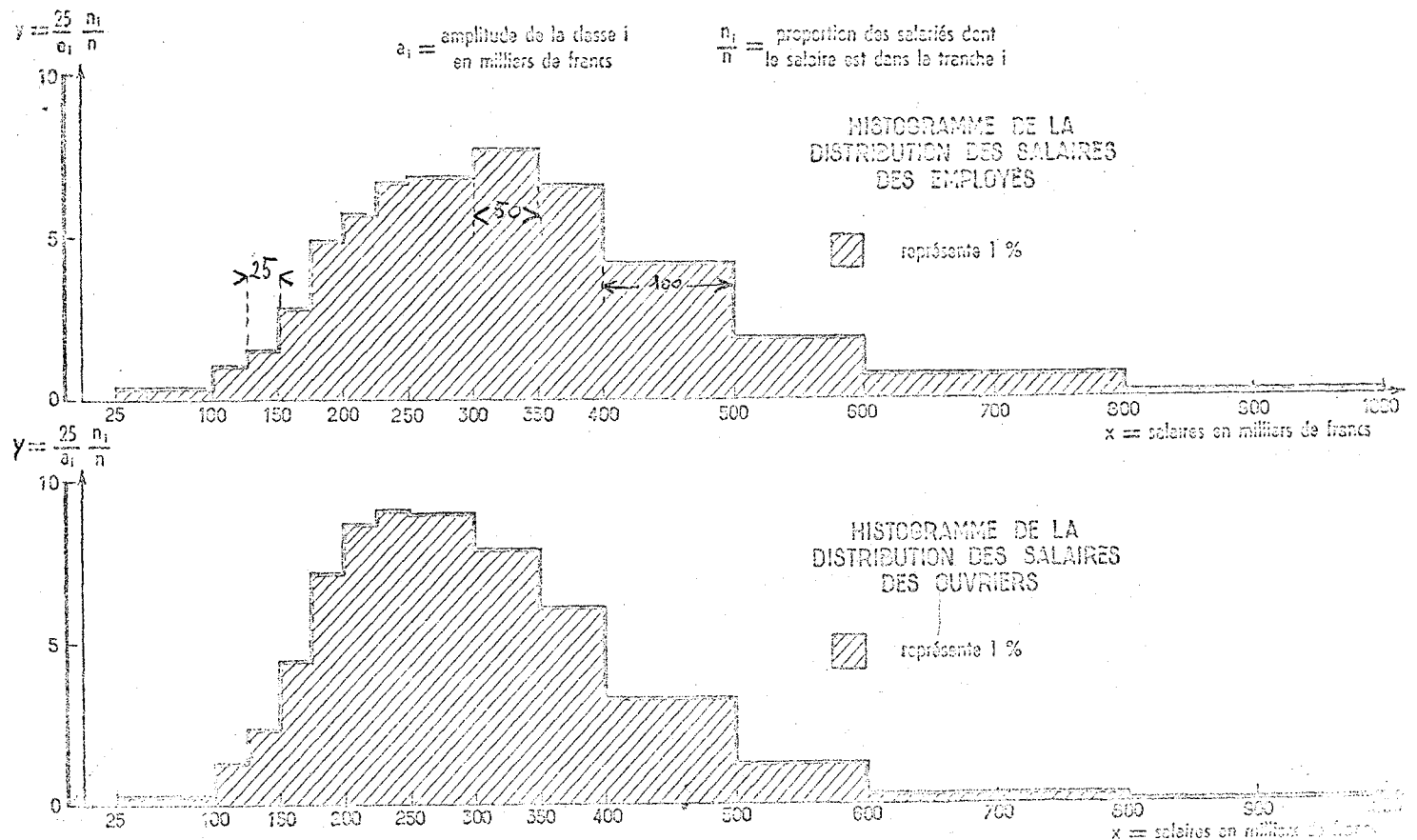




Si les amplitudes de classes diminuent, certaines irrégularités, dues à la faiblesse des effectifs, apparaissent. Cependant, si le nombre d'observations augmente, alors en faisant diminuer l'amplitude des classes, on se "rapprochera" de la courbe de densité théorique.

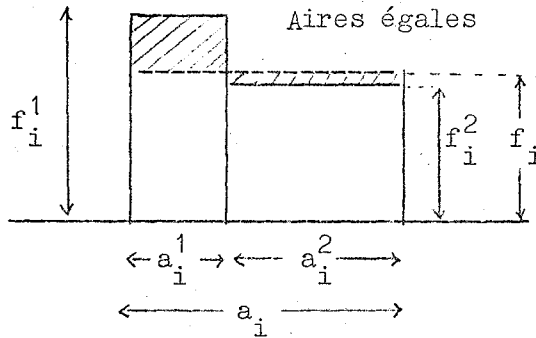
Si l'amplitude des classes augmente, on fait disparaître les irrégularités, en même temps, on approche de la courbe de densité théorique avec moins de précision.

### HISTOGRAMMES COMPARÉS DES DISTRIBUTIONS DES SALAIRES DES EMPLOYÉS ET OUVRIERS



Remarque : Introduction d'une classe supplémentaire ou réunion de deux classes adjacentes.

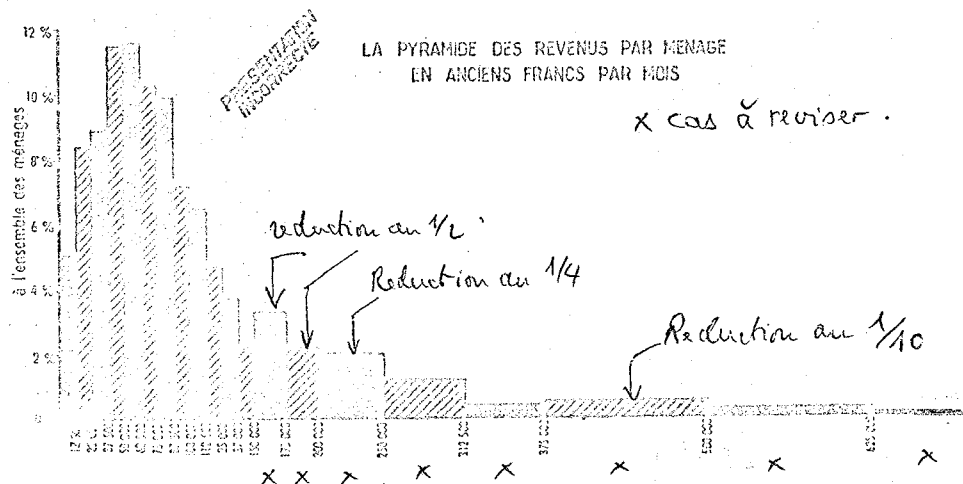
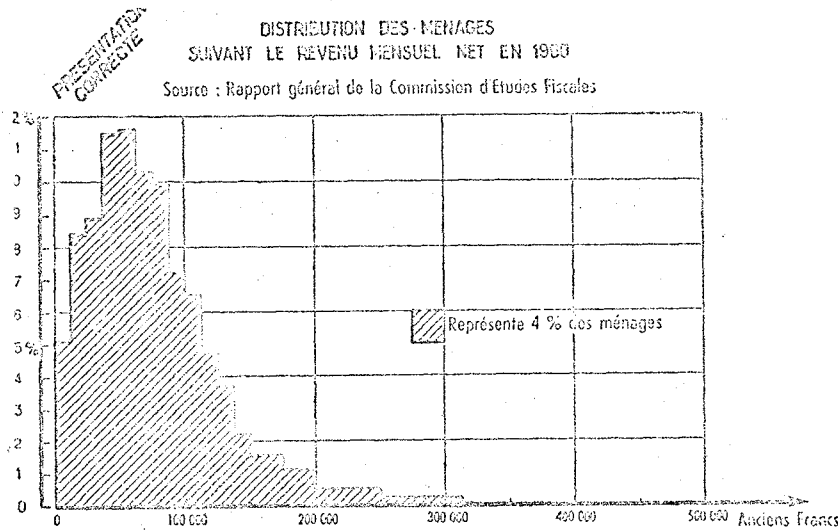
Cette opération est schématisée par le dessin suivant :



On doit avoir :  $a_i^1 f_i^1 + a_i^2 f_i^2 = (a_i^1 + a_i^2) f_i \Rightarrow f_i$ .

Exemple de présentation incorrecte :

Par exemple, pour la classe [175 000, 200 000], la représentation incorrecte figure en ordonnée  $f_i$  ; les surfaces hachurées représentent donc  $f_i \times a_i$ .



3 - CARACTERISTIQUES DE VALEURS CENTRALES.

Elles doivent :

- être définies de façon aussi objectives que possible,
- dépendre de toutes les observations, mais pas de leur ordre (sauf dans certains problèmes spécifiques faisant intervenir le temps, que nous n'étudierons pas)
- être peu sensibles aux fluctuations de l'échantillonnage.

La Médiane :

C'est une valeur de la variable statistique qui partage en deux effectifs égaux les individus de la population.

Si  $F$  est la fonction de répartition de l'échantillon

$$\begin{aligned}
 M \text{ médiane} &\iff F(M) \geq \frac{1}{2} \\
 &1 - F(M) \leq \frac{1}{2}
 \end{aligned}
 \tag{1}$$

Si  $F$  est strictement croissante,  $M$  est unique comme solution de (1)

Si  $F$  est discrète, deux cas

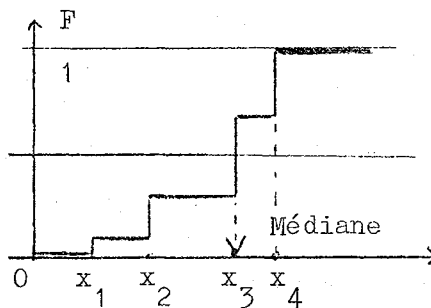
. Aucune valeur  $x_i$  ne donne

$$F(x_i) = \frac{1}{2} . \text{ On retient le } x_i$$

tel que :

$$F(x_i) < \frac{1}{2} < F(x_{i+1})$$

$\implies$

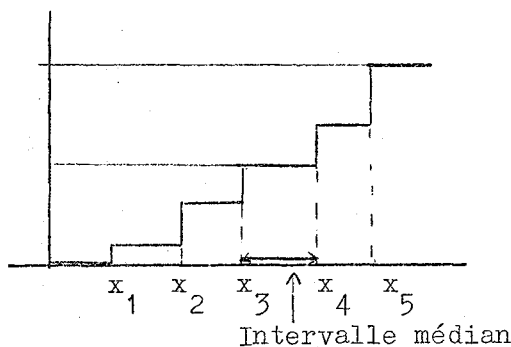


. Si  $\exists x_i$  tel que  $F(x_i) = \frac{1}{2}$  alors

tout point de  $[x_i, x_{i+1}[$  satisfait (1)

$\implies$

On a un intervalle médian.



Règle d'obtention : Classement dans l'ordre croissant des observations  $x_i$ .

Calcul des effectifs cumulés  $nF(x_i)$ .

Si  $nF(x_i) < 50 < nF(x_{i+1}) \rightarrow x_i$  valeur médiane

$nF(x_i) = 50$   $[x_i, x_{i+1}[$  intervalle médian

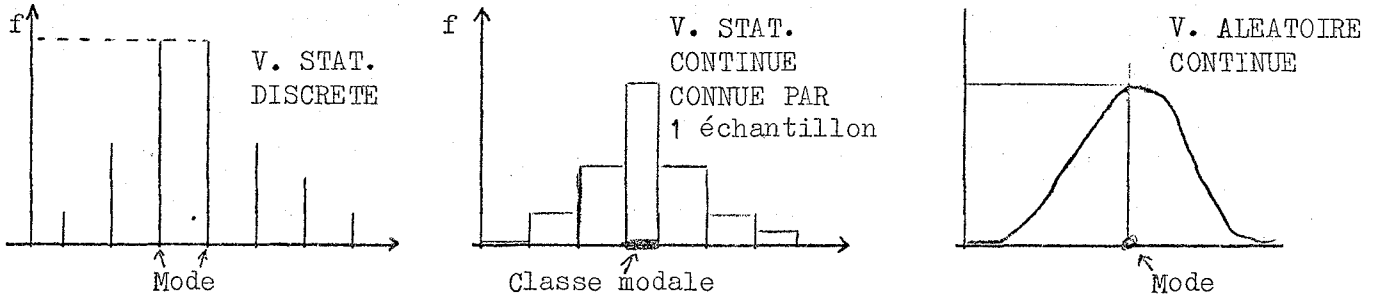
Classement des obs.	$x_0$	$x_1$	...	$x_6$	$x_7$	.....	$x_{25}$
Effectifs cumulés				45	63		

$$\frac{n}{2} = 50$$

Le Mode :

. C'est la valeur la plus "fréquente" de la variable statistique, on dit aussi valeur dominante. C'est pratiquement une valeur "centrale", quand on considère des distributions "régulières".

. Le mode d'une distribution continue à densité correspond au point où la densité prend la plus grande valeur



Moyenne : c'est la caractéristique la plus importante. Elle n'est définie que pour des variables quantitatives à valeurs dans un espace vectoriel.

. C'est :  $\bar{x} = \sum_{i=1}^k f_i x_i$  si  $\{x_1, \dots, x_k\}$  est l'ensemble des valeurs prises par l'échantillon.

4 - CARACTERISTIQUES DE DISPERSION.

Ecart quadratique moyen ou écart-type.

. C'est  $\sigma = \sqrt{\sum_i f_i |x_i - \bar{x}|^2}$   
 $\sigma^2$  est la variance.

5 - CARACTERISTIQUES DE FORME. (A titre d'exemple, elles sont très arbitraires).

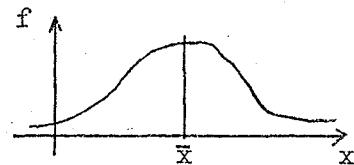
Coefficient d'assymétrie.

Si une distribution est symétrique, ces différents moments centrés d'ordre impair sont nuls.

En effet :

$$\int_R (x-\bar{x})^r f(x) dx = \int_{-\infty}^{\bar{x}} + \int_{\bar{x}}^{\infty}$$

et  $\int_{-\infty}^{\bar{x}} (x-\bar{x})^2 f(x) dx = - \int_{\bar{x}}^{\infty} (x-\bar{x})^2 f(x) dx$



(Symétrie de f par rapport à  $x = \bar{x}$ )

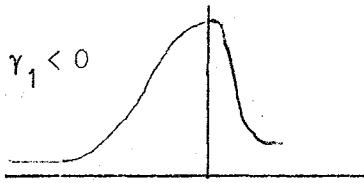
. Fischer a proposé le coefficient :

$$\gamma_1 = \frac{\mu_3}{\mu_2^{3/2}} = \frac{\mu_3}{\sigma^3}$$

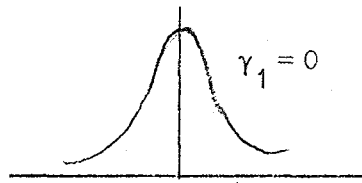
où  $\begin{cases} \mu_3 = \int (x-\bar{x})^3 f(x) dx \\ \mu_2 = \int (x-\bar{x})^2 f(x) dx \end{cases}$

→ sans dimension (c'est pour cela que l'on divise par  $\mu_2^{3/2}$ )

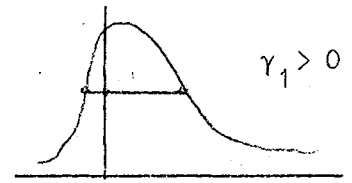
→ nul si la distribution a un axe de symétrie



Distribution dissymétrique étalée à gauche.



symétrique



Distribution dissymétrique étalée à droite.

$$d = \frac{Q_1 + Q_3 - 2M}{M}$$
 est également un coefficient d'assymétrie, sans dimension. (d = 0 si la distribution est unimodale et symétrique).

## 6 - CARACTERISTIQUES DE CONCENTRATION.

Considérons le tableau des salaires (page 2, B).

A chaque casé (classe) associons la masse des salaires ; par exemple :

$$\text{Classe } (125 < \text{salaire} < 150) \Rightarrow 4342 \times \frac{125 \times 150}{2} = 601 \text{ Millions}$$

$$\text{Classe } e_i \Rightarrow S_i = \text{masse des salaires de la classe } e_i$$

$$\text{Classe } e_i \Rightarrow q_i = \frac{\sum_{j=1}^i S_j}{\sum_{j=1}^k S_j}$$

Ainsi les ouvriers dont le salaire est inférieur à 200 représentent 16,5 % de l'effectif des ouvriers et se partagent 8,29 % de la masse totale des salaires.

Notons :  $p_i$  : proportion des ouvriers dont le salaire est  $\leq e_i$

$q_i$  : proportion de la masse salariale gagnée par les ouvriers dont le salaire est  $\leq e_i$

Limites de classe $e_i$ (en milliers de francs)	Effectifs $n_i$	Masses des salaires par classe : $S_i$ (en millions de francs)	$p_i = F(e_i)$ en %	$q_i = \frac{\sum_{j=1}^i S_j}{\sum_{j=1}^k S_j}$ (en %)
100	1 721	114	0,93	0,20
125	2 413	273	2,22	0,66
150	4 342	601	4,56	1,70
175	8 264	1 349	9,00	4,01
200	13 300	2 494	16,15	8,29
225	16 053	3 404	24,79	14,14
250	16 774	3 972	33,81	20,96
300	33 251	9 120	51,69	36,62
350	29 211	9 461	67,40	52,86
400	22 453	8 391	79,47	67,27
500	24 005	10 623	92,38	85,51
600	9 477	5 133	97,48	94,33
800	4 093	2 717	99,68	98,99
1 000	443	386	99,92	99,65
1 500	125	144	99,99	99,90
2 000	12	20,3	99,99	99,93
5 000	14	37,7	100,00	100,00
Total	185 951	58 240	—	—

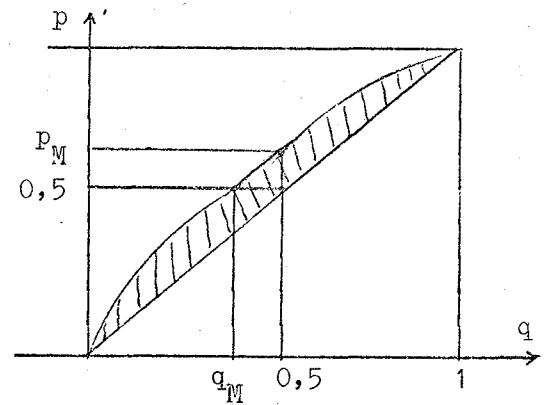
Courbe de concentration.

C'est la courbe

$$p_i = f(q_i)$$

. Comme  $p \leq q$ , elle se situe dans la moitié supérieure du carré  $1 \times 1$  construit sur les axes.

. Elle est croissante



C - DISTRIBUTIONS STATISTIQUES A DEUX CARACTERES

---

Considérons une population de  $n$  individus décrits simultanément par deux caractères A et B, A ayant  $k$  modalités  $\{A_1, A_2, \dots, A_k\}$ , B ayant  $p$  modalités  $\{B_1, B_2, \dots, B_p\}$ .

On cherchera à décrire la distribution statistique associée

soit : - par des tableaux et une représentation graphique

soit : - par une description numérique.

Les caractères peuvent être qualitatifs ou quantitatifs.

I - TABLEAUX STATISTIQUES.

Notons  $n_{ij}$  = nombre d'individus présentant à la fois le caractère  $A_i$  et  $B_j$ .

Du fait que les modalités de A

de B sont incompatibles et exhaustives.

On a : 
$$\sum_{i,j} n_{ij} = n .$$

Modali- tés de B Modali- tés de A	B <sub>1</sub>	B <sub>2</sub>	...	B <sub>j</sub>	...	B <sub>p</sub>	Totaux (Marges) horizontaux
A <sub>1</sub>	$n_{11}$	$n_{12}$	...	$n_{1j}$		$n_{1p}$	$n_{1.}$
A <sub>2</sub>	$n_{21}$	$n_{22}$	...	$n_{2j}$		$n_{2p}$	$n_{2.}$
⋮							
A <sub>i</sub>	$n_{i1}$	$n_{i2}$	...	$n_{ij}$		$n_{ip}$	$n_{i.}$
⋮							
A <sub>k</sub>	$n_{k1}$	$n_{k2}$	...	$n_{kj}$	...	$n_{kp}$	$n_{k.}$
Tableaux verticaux	$n_{.1}$	$n_{.2}$	...	$n_{.j}$		$n_{.p}$	$n$

La fréquence de  $(A_i, B_j)$  est  $\frac{n_{ij}}{n} = f_{ij}$ .

Distributions marginales.

C'est la distribution des deux marges

distribution marginale associée à A :  $n_{1.}, n_{2.}, \dots, n_{k.}$

" " " à B :  $n_{.1}, n_{.2}, \dots, n_{.p}$

C'est la distribution associée à l'un des caractères, une fois l'autre oublié.

Distributions conditionnelles.

Considérons la sous-population réalisant le caractère  $B_j$ . Sur cette sous-population, la distribution du caractère  $A$  est appelée la distribution conditionnelle de  $A$ , à condition que  $B_j$  soit réalisé.

On a :  $f_i^j(A_i/B_j) = \frac{n_{ij}}{n_{.j}}$ .

Il y a  $p$  distributions conditionnelles suivant le caractère  $A$   
 $k$  " " " " "  $B$

Si on note  $(A_i, B_j)$  l'évènement :  $A_i$  et  $B_j$  réalisé on a :

$P(A_i \cap B_j) = P(A_i) P(B_j/A_i) = P(B_j) P(A_i/B_j)$ .

Indépendance.

Les caractères  $A$  et  $B$  sont indépendants si <sup>récapitulé</sup>

$P(A/B) = P(A)$  (et alors  $P(B/A) = P(B)$ ).

Alors, cela est équivalent au fait que les  $p$  distributions conditionnelles de  $A$  sont les mêmes, en particulier identiques à la distribution marginale de  $A$ . En effet :  $f_i^j$  ind de  $j$

$f_i^j = \frac{n_{i1}}{n_{.1}} = \dots = \frac{n_{ij}}{n_{.j}} = \dots = \frac{n_{ip}}{n_{.p}} = \frac{n_{i.}}{n_{..}} = f_{i.}$

Liaison fonctionnelle.

Le caractère  $A$  est lié fonctionnellement au caractère  $B$  si à chaque modalité  $B_j$  de  $B$  correspond une telle modalité du caractère  $A$  (on noterait :  $A = f(B)$ ).

Cette notion n'est pas réciproque comme celle d'indépendance. Elle est particulièrement intéressante dans le cas de caractères quantitatifs.

Si  $f$  est bijective, la notion est réciproque.

$A \setminus B$	$B_1$	$B_2$	$B_3$	$B_4$	$B_5$
$A_1$	4	0	2	0	0
$A_2$	0	1	0	1	0
$A_3$	0	0	0	0	3

. Distributions conditionnelles de  $A$  dégénérées.

Il y a globalement 6 types de distributions à 2 caractères

- Situation plus favorable à une description graphique
- Les deux qualitatifs
  - 1 qualitatif, l'autre quantitatif discret
  - " " " continu
  - 2 quantitatifs discrets
  - 1 quantitatif discret, 1 continu
  - 2 quantitatifs continus
- } Situation plus favorable à la description numérique



## II - REPRESENTATIONS GRAPHIQUES.

Une représentation graphique d'une distribution statistique à 2 caractères doit s'attacher à figurer, d'une part la distribution globale (représenter toute l'information), et si possible les distributions conditionnelles ainsi que marginales, les distributions conditionnelles étant définies par  $f_{ij} = n_{ij}/n_{.j}$ .

1<sup>er</sup> cas : A et B qualitatifs.

On peut représenter la distribution globale et une famille de distributions conditionnelles  $(A/B_j)$  (ou l'inverse).

$n_{ij}$  est "représenté" par un rectangle de base  $n_{.j}$  et dont la hauteur est proportionnelle à la fréquence conditionnelle  $f_{ij}^j$ .

$$n_{ij} = n_{.j} f_{ij}^j .$$

On a donc : → effectifs marginaux caractère B  $n_{.j}$  (base rectangles)  
 → effectifs du tableau (aire des rectangles)  
 → fréquences conditionnelles  $(A/B_j)$  soit  $f_{ij}^j$ .

Voir tableaux et graphiques pages 4 et 5.

2<sup>ème</sup> cas : 1 caractère qualitatif, 1 caractère quantitatif.

A désigne le caractère qualitatif

B " " " quantitatif.

On peut déjà utiliser un mode de représentation comme si les deux caractères étaient qualitatifs.

On peut également représenter les distributions conditionnelles à partir de diagrammes en bâtons (autant de diagrammes que de modalités du caractère qualitatif). Voir page 6.

Distributions conditionnelles en fréquences cumulées (Catégorie / Epoque construction) ⇒

Logements (résidences principales) suivant l'époque de construction (1) et la catégorie socio-professionnelle du chef de ménage (Recensement général de la Population ; mars 1968. Sondage au 1/200).

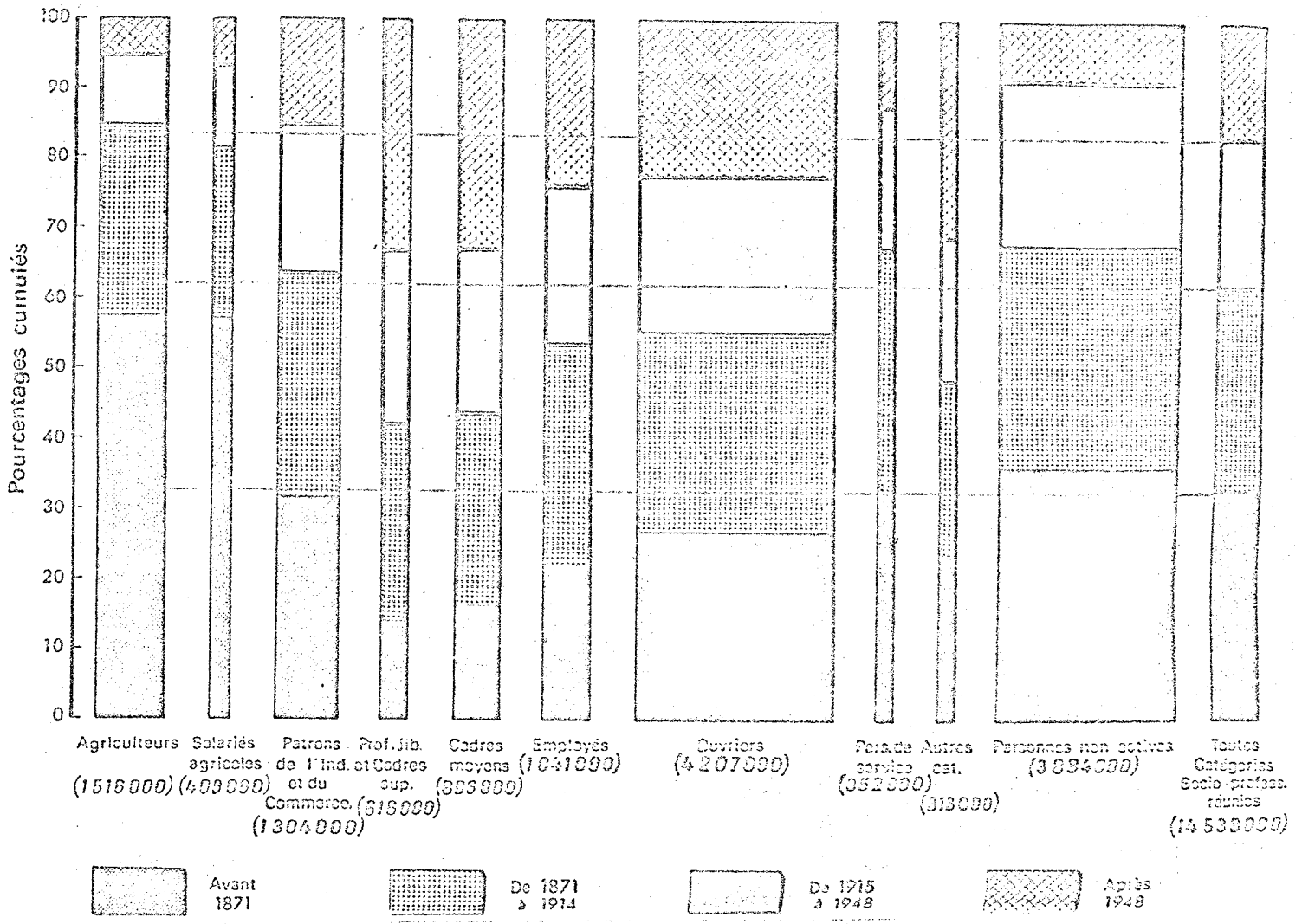
Epoque de construction	Catégorie socio-professionnelle				Total
	Avant 1871	De 1871 à 1914	De 1915 à 1948	Après 1948	
Agriculteurs.....	873 340	410 040	158 380	74 620	1 516 380
Salariés agricoles.....	233 060	100 160	48 600	27 280	409 100
Patrons de l'industrie et du commerce.....	415 380	413 000	280 520	195 380	1 304 280
Professions libérales et Cadres supérieurs ..	87 120	175 660	148 760	204 440	615 980
Cadres moyens.....	144 560	247 800	210 640	293 180	896 240
Employés.....	231 760	322 700	237 800	249 180	1 041 440
Ouvriers.....	1 118 440	1 177 820	954 500	956 040	4 206 800
Personnel de service ..	112 560	124 260	72 400	43 720	351 940
Autres catégories.....	73 240	77 960	65 360	95 960	312 520
Personnes non actives	1 598 840	1 282 120	932 340	322 220	3 883 520
<b>Total</b>	<b>4 686 300</b>	<b>4 280 580</b>	<b>3 109 300</b>	<b>2 462 020</b>	<b>14 518 200</b>

(1) La lecture rigoureuse de cette table est un caractère quantitatif continu. Par suite de la méthode d'échantillonnage du recensement, un écart de 1% est admis pour chaque classe définissant une époque de construction et la fréquence d'un caractère qualitatif.

Distributions conditionnelles (Epoque / Catégorie socio.) ⇒

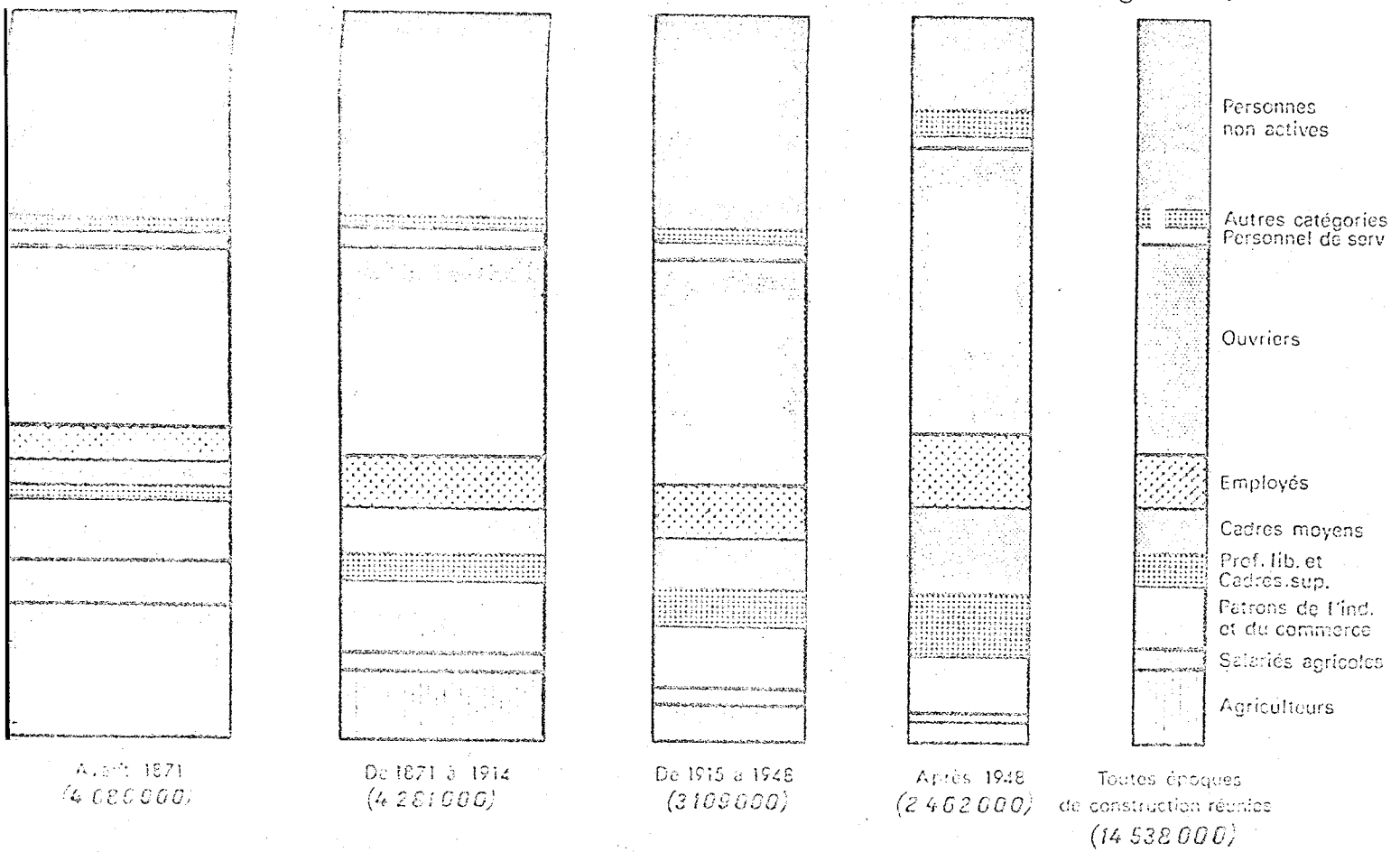
Epoque de construction	Catégorie socio-professionnelle				Toutes époques de construction réunies
	Avant 1871	De 1871 à 1914	De 1915 à 1948	Après 1948	
Agriculteurs.....	57,6	84,7	95,1	100,0	100,0
Salariés agricoles.....	57,0	81,5	93,4	100,0	100,0
Patrons de l'industrie et du commerce.....	31,8	63,5	85,0	100,0	100,0
Professions libérales et cadres supérieurs.....	14,1	42,6	66,8	100,0	100,0
Cadres moyens.....	16,1	43,8	67,3	100,0	100,0
Employés.....	22,3	53,3	76,1	100,0	100,0
Ouvriers.....	26,6	54,6	77,3	100,0	100,0
Personnel de service.....	32,0	67,0	87,6	100,0	100,0
Autres catégories.....	23,4	48,4	69,3	100,0	100,0
Personnes non actives ..	36,0	67,7	91,7	100,0	100,0
<b>Toutes catégories socio-professionnelles réunies</b>	<b>32,2</b>	<b>61,7</b>	<b>83,1</b>	<b>100,0</b>	<b>100,0</b>

Epoque de construction	Catégorie socio-professionnelle				Toutes époques de construction réunies
	Avant 1871	De 1871 à 1914	De 1915 à 1948	Après 1948	
Agriculteurs.....	18,6	9,6	5,1	3,0	10,4
Salariés agricoles.....	23,6	11,9	6,7	4,1	13,2
Patrons de l'industrie et du commerce.....	32,5	21,6	15,7	12,1	22,2
Professions libérales et cadres supérieurs.....	34,4	25,7	20,5	20,4	26,4
Cadres moyens.....	37,5	31,5	27,3	32,3	32,6
Employés.....	42,4	39,0	34,9	42,4	39,8
Ouvriers.....	66,3	66,5	65,6	81,2	68,7
Personnel de service ..	68,7	69,4	67,9	83,0	71,1
Autres catégories.....	70,2	71,2	70,0	86,9	73,3
Personnes non actives ..	100,0	100,0	100,0	100,0	100,0



Distributions conditionnelles (Catégorie/Epoque)

Les dernières colonnes représentent les distributions marginales.

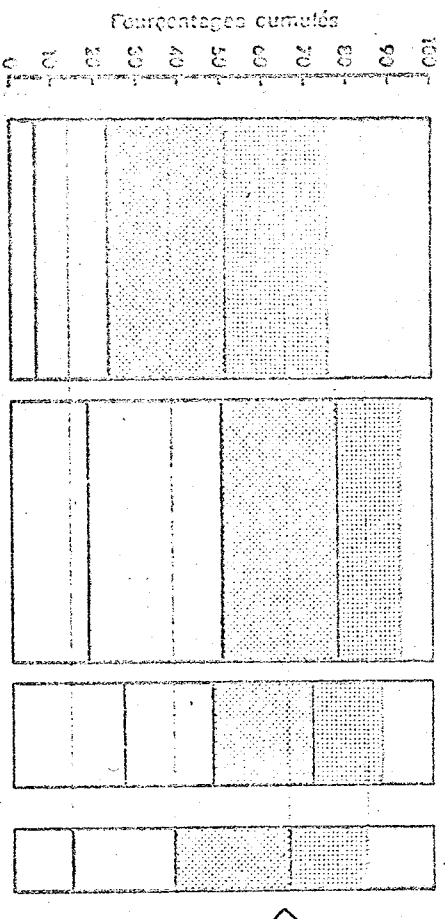


Distribution des logements suivant la catégorie socio-professionnelle de chef de ménage et l'époque de construction.

Exemple: Distribution du logement suivant le statut d'occupation et le nombre de pièces.

Statut d'occupation	Nombre de pièces					Total
	1	2	3	4	5 et plus	
Propriétaires ...	371 600	1 078 420	1 639 960	1 460 700	1 491 800	6 042 480
Locataires (1)	1 099 040	1 923 280	1 673 760	868 140	477 140	6 041 360
Autres (2).....	638 540	509 560	576 720	431 440	278 100	2 454 360
Total .....	2 129 180	3 511 260	3 890 440	2 760 280	2 247 040	14 538 200

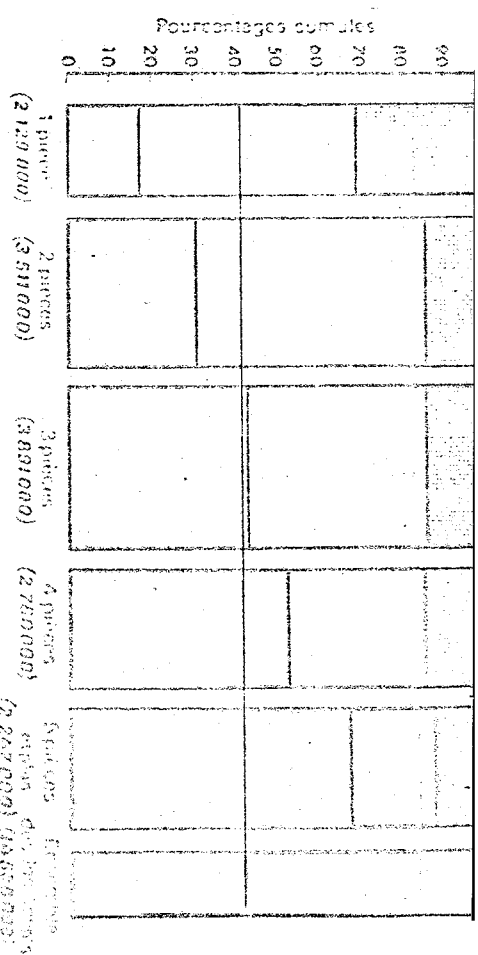
(1) Locataires d'un logement loué vide.  
 (2) Personnes logées par leur employeur (1 280 600), personnes logées à titre gracieux (670 000), locataires ou sous-locataires d'un local meublé (500 000).



Propriétaires (6043000)  
 Locataires (6041000)  
 Autres (2454000)

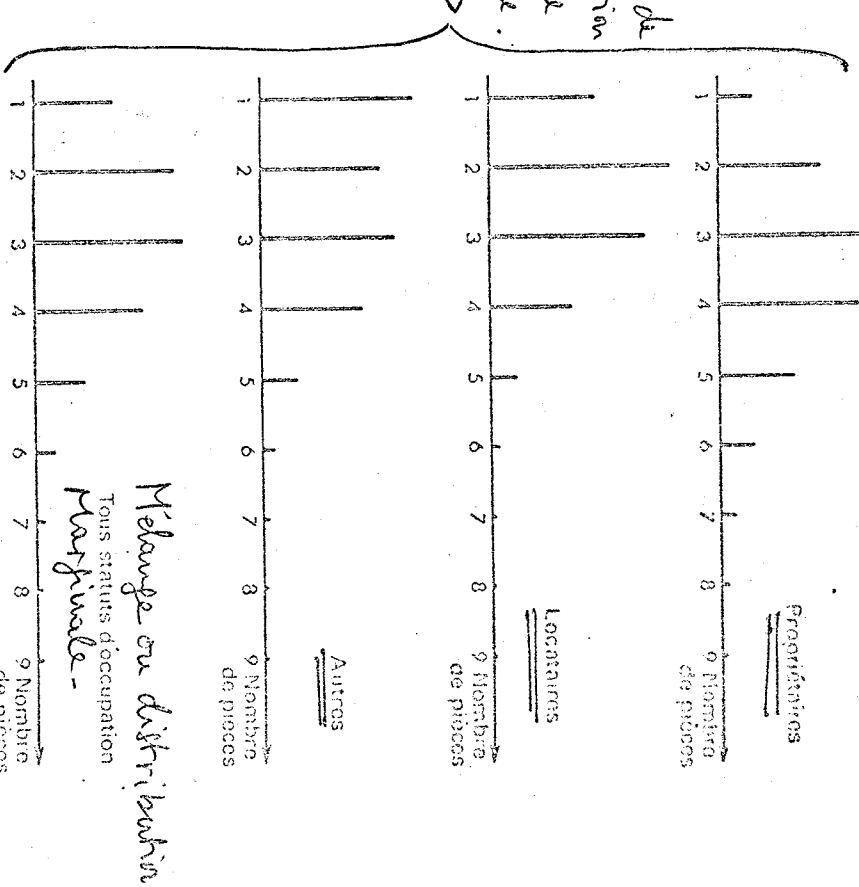
Legend:  
 1 pièce  
 2 pièces  
 3 pièces  
 4 pièces  
 5 pièces et plus

possibilité de représentation par un axe numérique.



Propriétaires  
 Locataires  
 Autres

b) Représentation de la distribution globale et des distributions conditionnelles suivant le statut d'occupation.



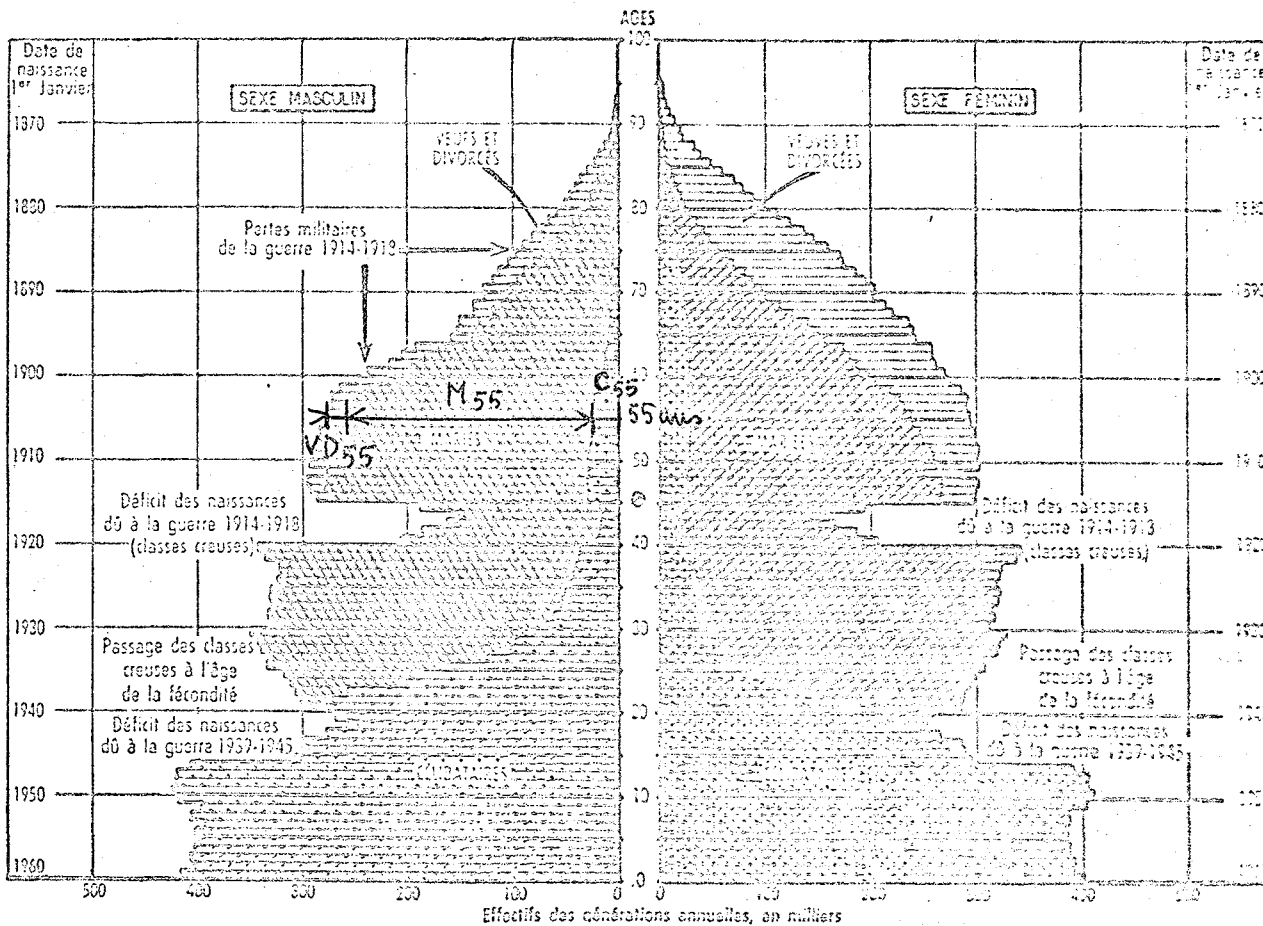
c) Représentation des distributions conditionnelles suivant le nombre de pièces.

Autre exemple. Distribution de la population française masculine suivant l'âge et l'état matrimonial.

- 1) Pyramide, où apparaît, à âge fixé les distributions conditionnelles de l'état matrimonial.
- 2) Distribution conditionnelle suivant l'âge pour les 3 modalités de l'état matrimonial.

PYRAMIDE DES AGES DE LA POPULATION FRANÇAISE PAR ÉTAT MATRIMONIAL - ÉVALUATION AU 1-1-1960

Source : I. N. S. E. E.

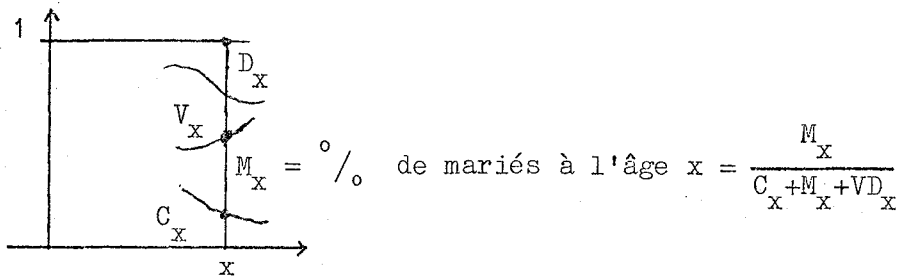


Notons  $C_x$  = effectifs des célibataires à l'âge  $x$

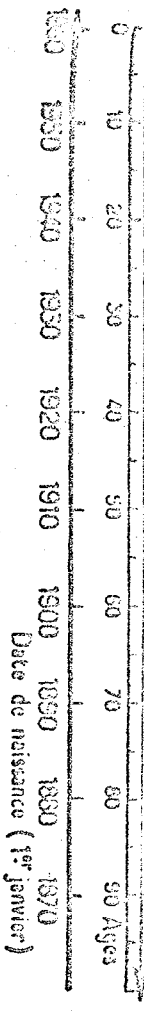
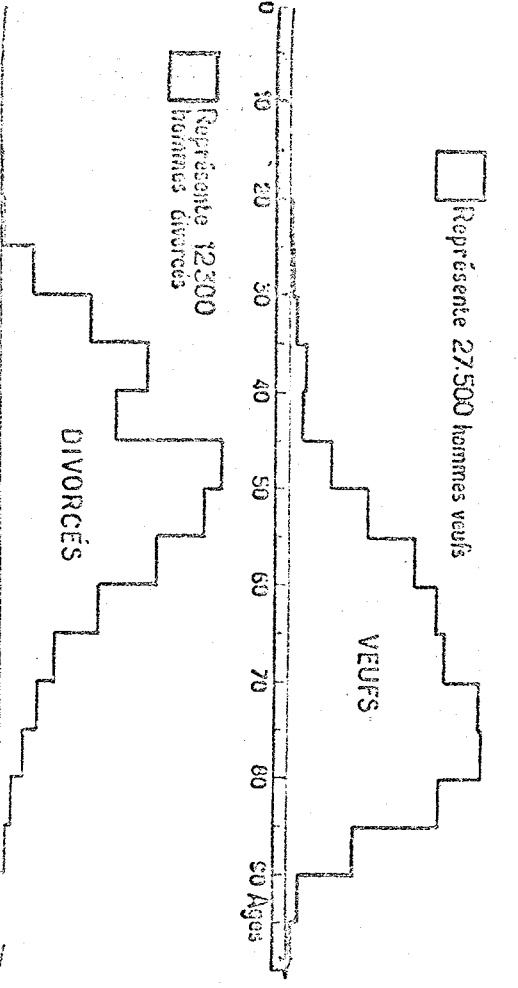
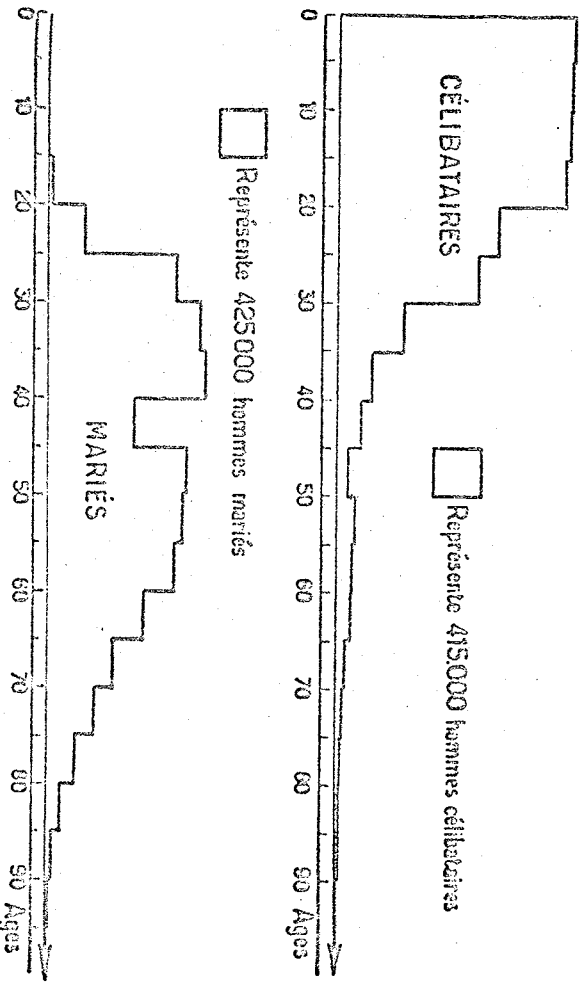
$M_x$  = " " mariés à l'âge  $x$

$VD_x$  = " " veufs ou divorcés "  $x$  .

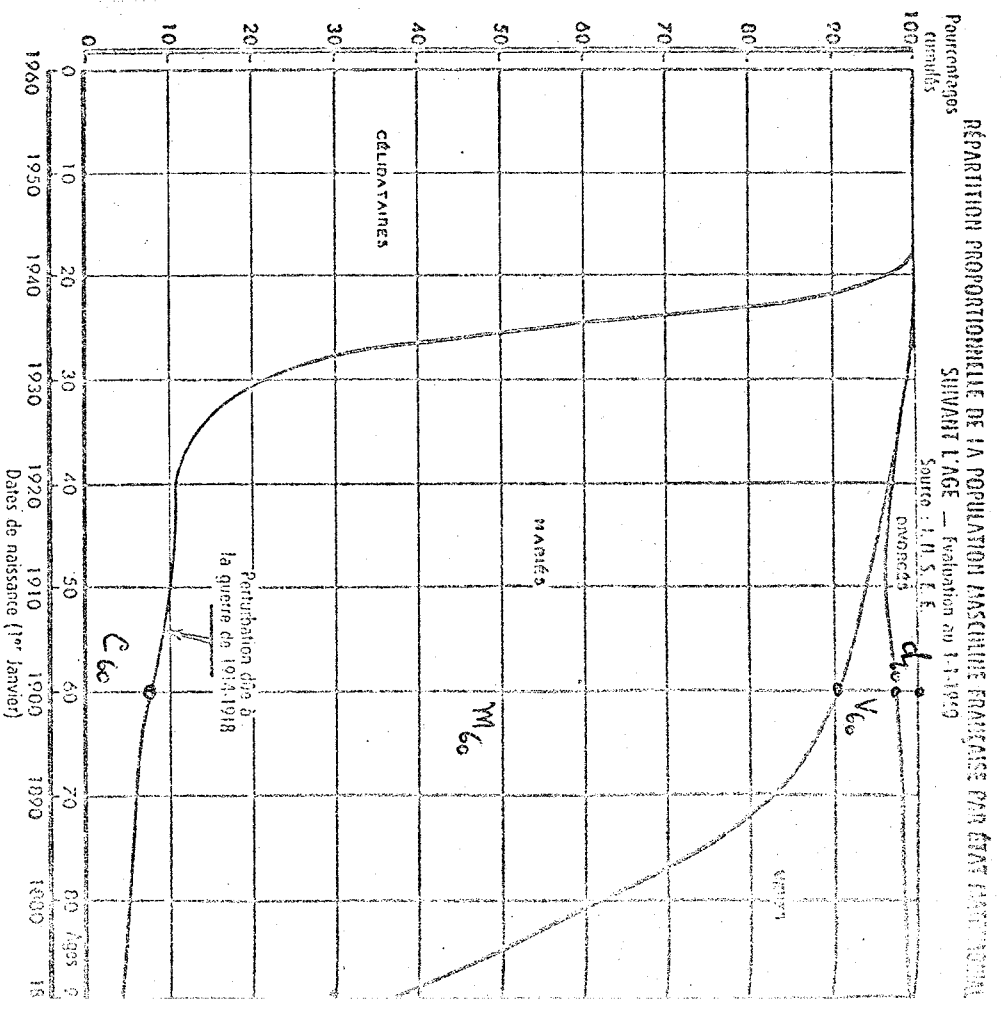
On peut considérer  $x$  comme une variable continue, et ainsi représenter les distributions des 4 modalités de l'état matrimonial pour tout  $x$  (Si on sépare  $V$  et  $D$ )

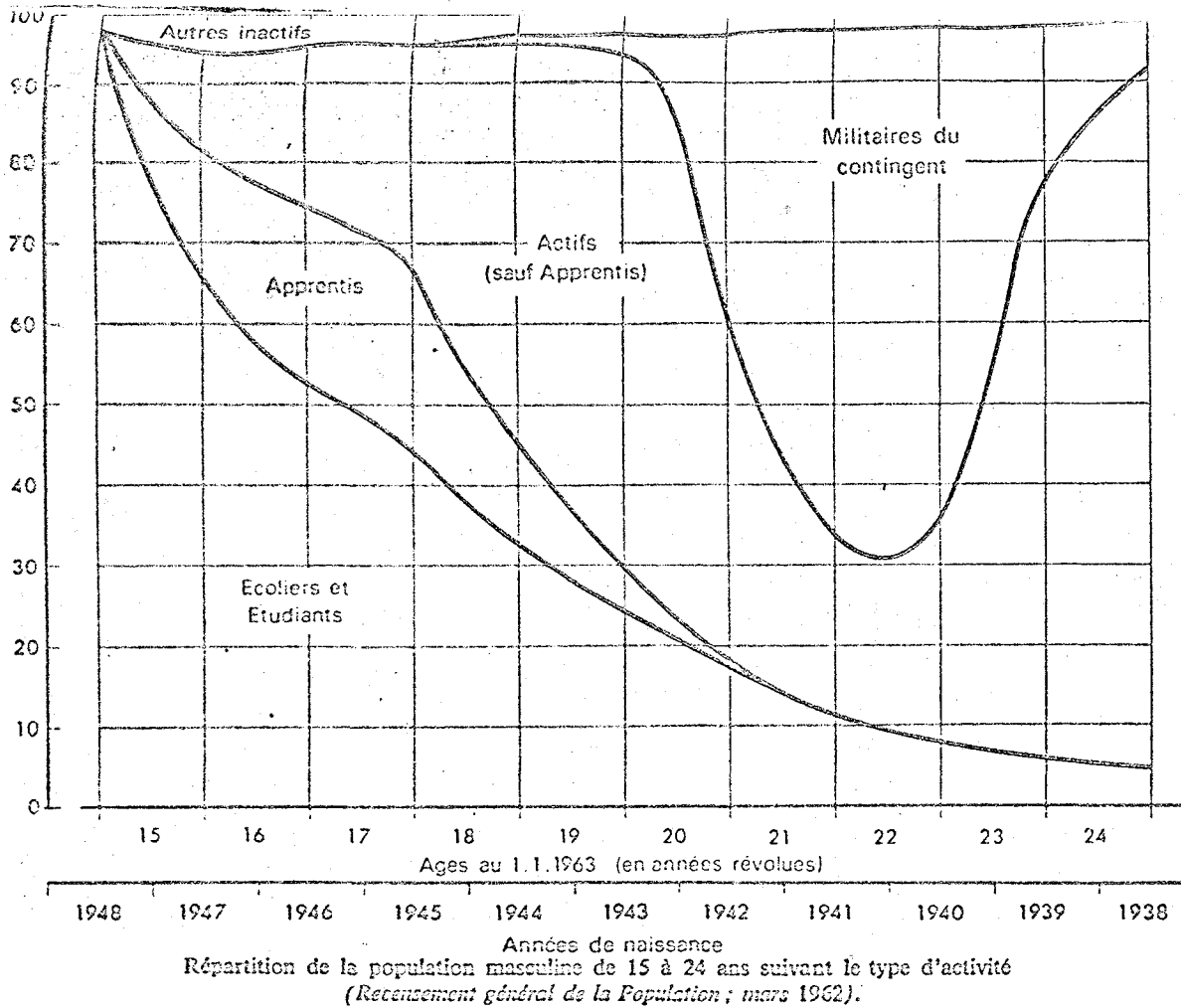


Evaluation du 1<sup>er</sup> janvier 1960  
(Source : I.N.S.E.E.)



1  
8  
1  
2





3 derniers cas : Deux caractères quantitatifs.

On peut utiliser déjà les différents modes de représentation énoncés dans le 1<sup>er</sup> et 2<sup>ème</sup> cas.

Notons  $X$  et  $Y$  les deux variables définissant les caractères. Par exemple, si  $X$  et  $Y$  sont discrètes et  $n_{ij}$  est l'effectif  $(x_i, y_j)$ , on peut représenter cet effectif par un cercle centré en  $(x_i, y_j)$  de surface  $\%$  à  $n_{ij}$ .

Deux exemples sont donnés :

- 1) Répartition des ménages français suivant le nombre de personnes et le nombre de pièces du logement occupé.
- 2) Répartition des ménages suivant l'âge du chef de ménage et le nombre d'enfants de moins de 16 ans.

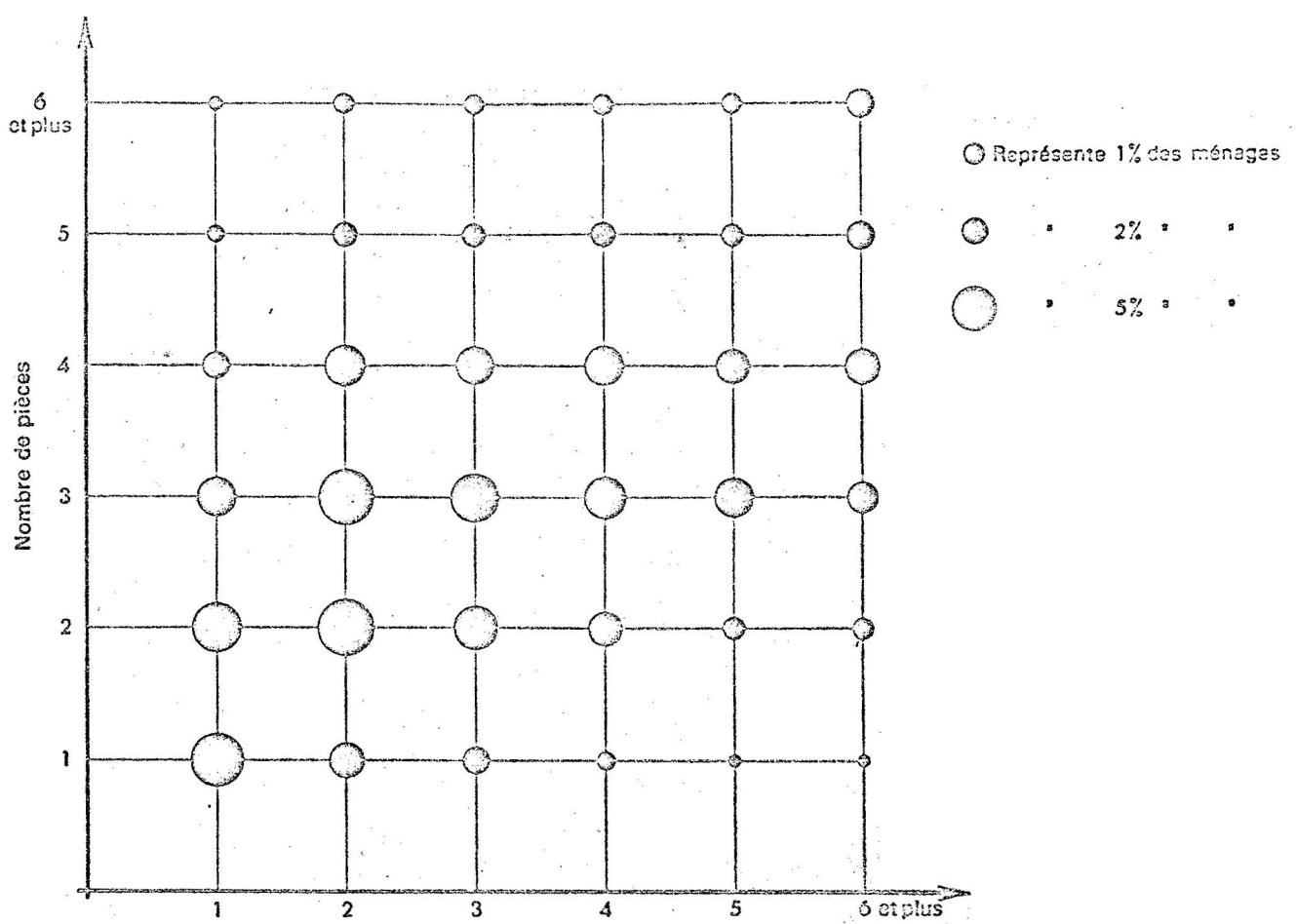
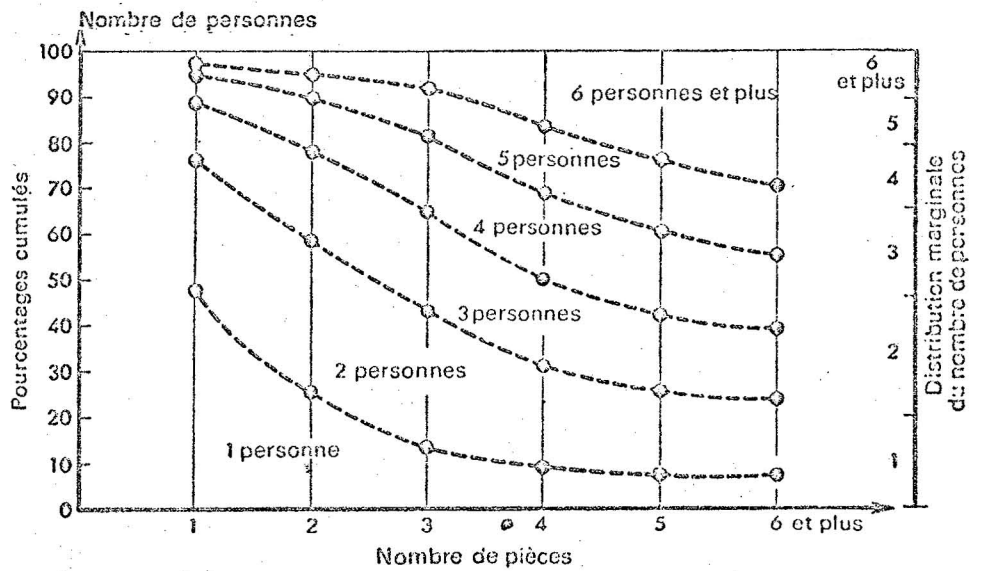
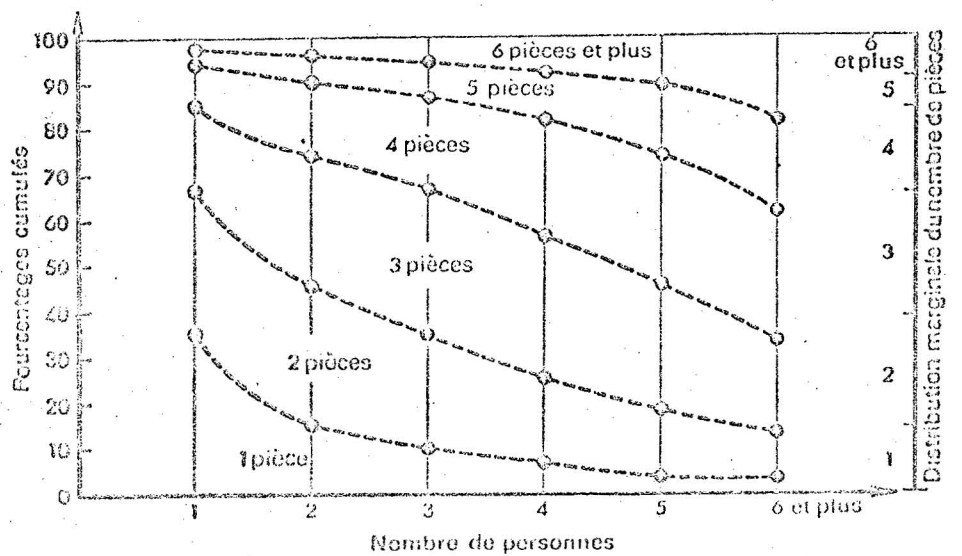


Tableau des données :  
 $n_{ij}$  = nombre de familles

Nbre. de pièces \ Nbre. de personnes	Total					
	1	2	3	4	5	6 et plus
1	1 019 420	880 860	536 320	248 180	95 140	71 680
2	597 120	1 181 480	1 128 640	607 420	230 180	156 920
3	274 920	691 180	856 000	524 460	215 960	151 600
4	134 160	403 620	674 040	534 020	230 860	159 940
5	57 280	192 180	375 840	396 300	201 260	147 140
6 et plus	46 280	161 940	319 620	449 900	301 020	285 340
Total	2 129 180	3 511 260	3 890 440	2 760 280	1 274 420	972 620



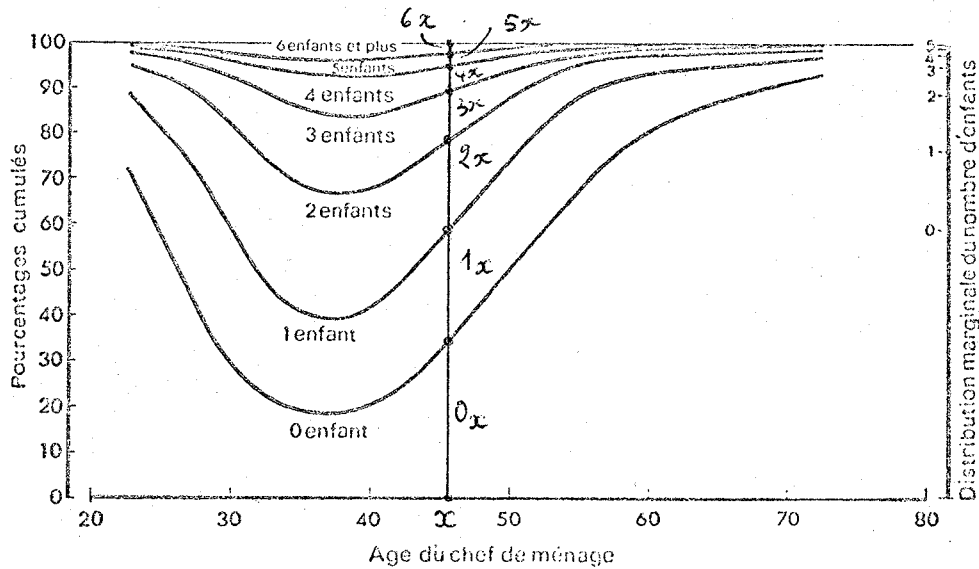
Distributions conditionnelles des logements suivant le nombre de personnes en fonction du nombre de pièces.



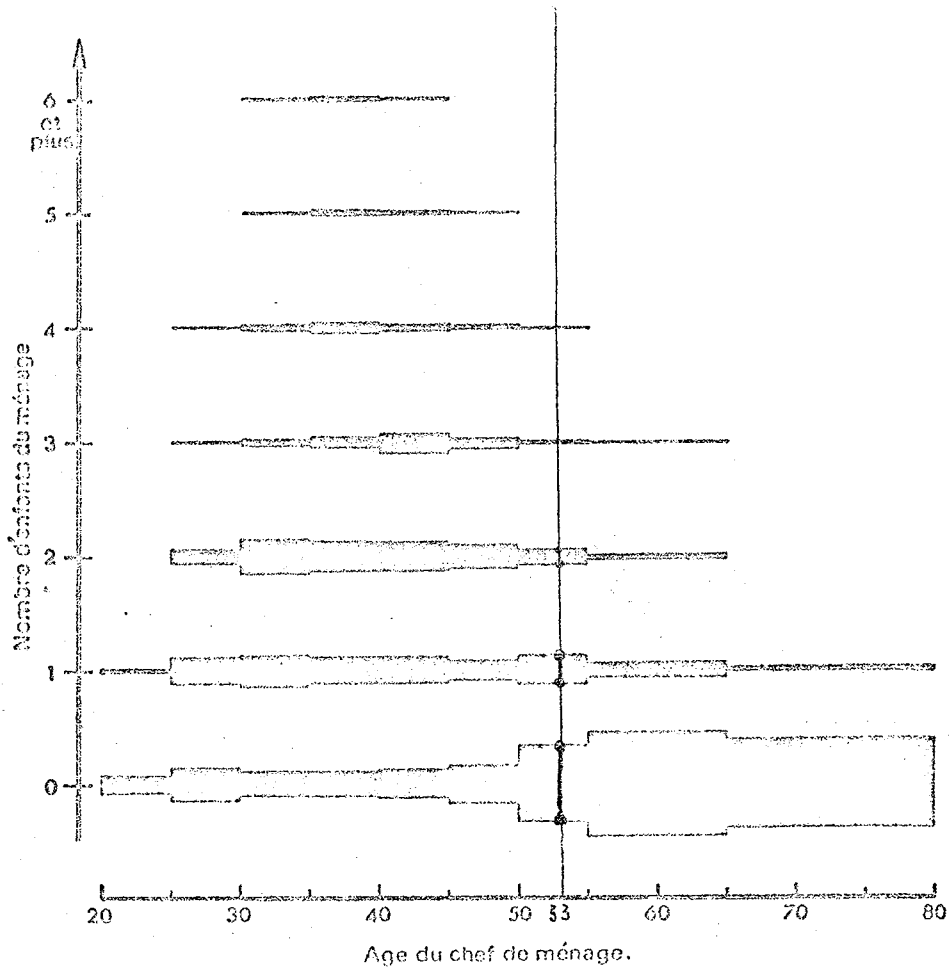


Répartition proportionnelle des ménages suivant l'âge du chef de ménage et le nombre d'enfants de 16 ans et moins (Recensement général de 1962).

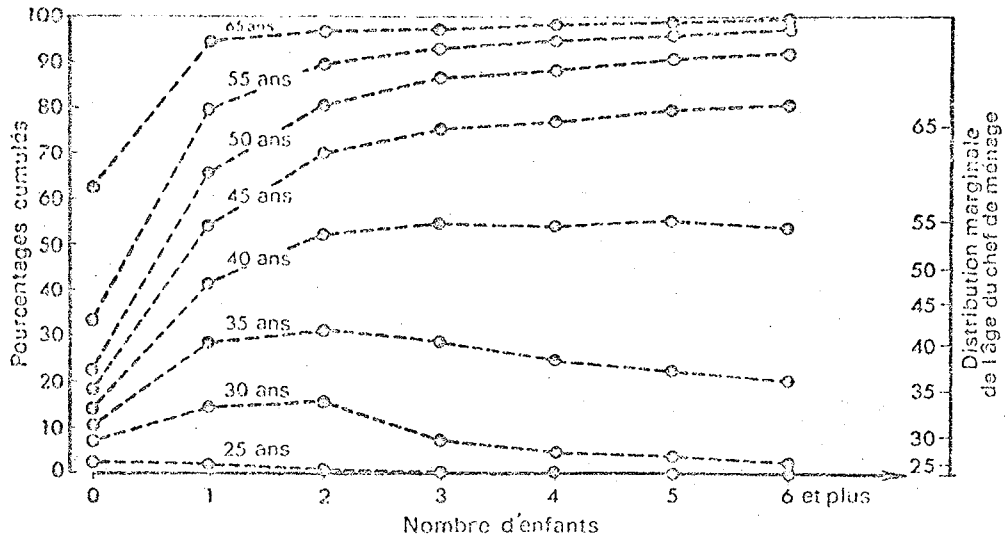
Age du chef de ménage	Nbre d'enfants du ménage							Total
	0	1	2	3	4	5	6 et plus	
25 ans	1 477	376	131	35	11	3	1	2 034
30 ans	2 736	2 142	1 192	430	133	43	24	6 700
35 ans	2 120	2 372	2 496	1 345	597	249	177	9 356
40 ans	1 983	2 130	2 612	1 688	892	424	536	10 065
45 ans	2 315	2 159	2 137	1 309	690	308	263	9 181
50 ans	3 086	1 913	1 313	685	345	147	111	7 600
55 ans	6 275	2 420	1 072	433	193	68	53	10 514
65 ans	17 000	2 489	886	281	107	34	23	20 820
Total	59 131	16 952	12 245	6 364	3 023	1 292	993	100 000



Distributions conditionnelles des ménages suivant le nombre d'enfants en fonction de l'âge du chef de ménage.



Distribution des ménages suivant l'âge du chef de ménage et le nombre d'enfants de 16 ans et moins.



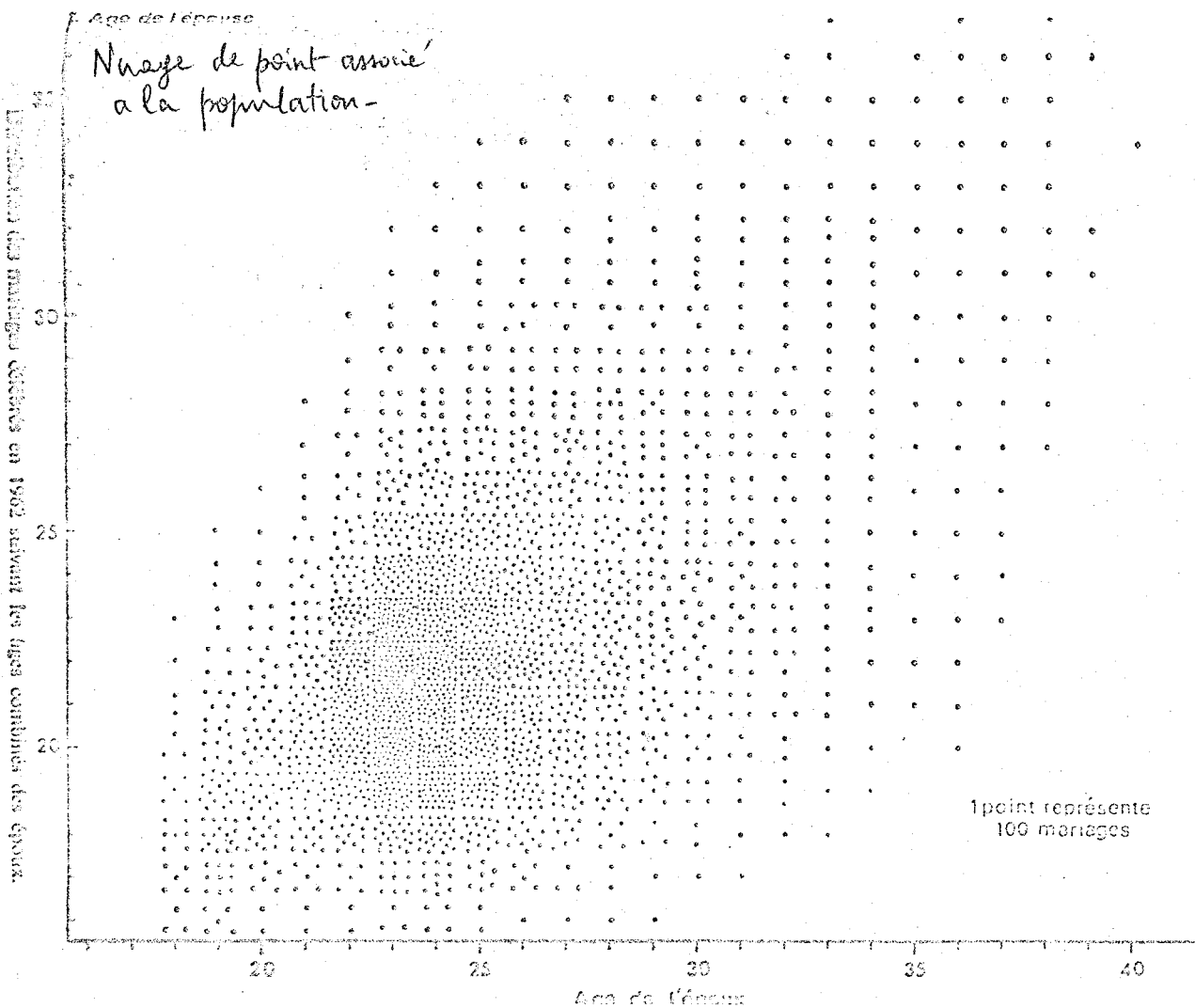
Distributions conditionnelles des ménages suivant l'âge du chef de ménage en fonction du nombre d'enfants.

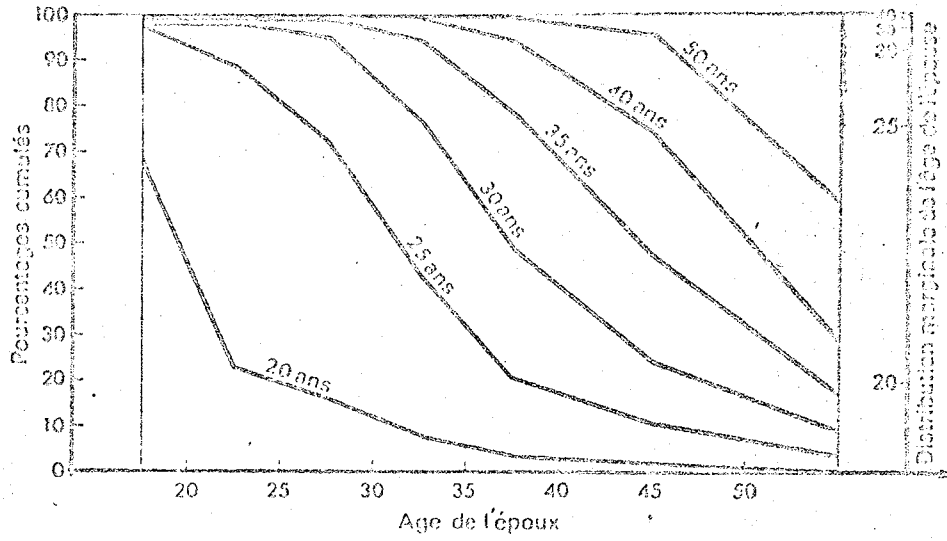
Autre exemple : Distribution des mariages suivant l'âge de l'époux et l'âge de l'épouse.

Exemple 3.

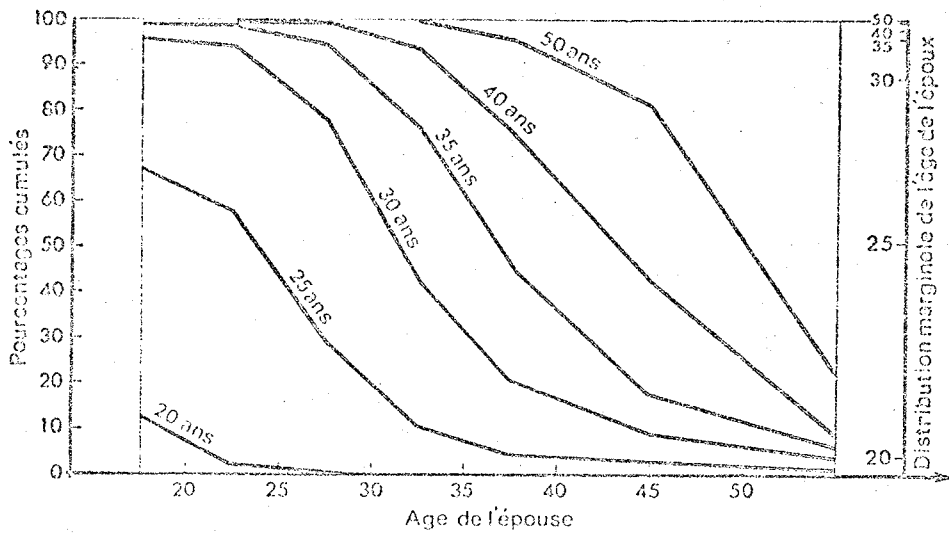
Distribution des mariages célébrés en 1962 suivant les âges combinés des époux (1).

Age de l'épouse \ Age de l'époux	20	25	30	35	40	50	Total
20	6 756	3 051	180	15	3	3	10 017
25	29 416	84 556	13 430	1 205	168	50	128 835
30	15 893	54 978	22 774	3 890	651	113	98 313
35	1 789	8 289	7 809	4 111	1 021	244	23 278
40	253	1 304	1 996	2 078	1 232	362	7 247
50	66	283	447	733	852	697	3 268
	6	46	59	83	145	336	1 147
Total	54 190	152 507	46 695	12 115	4 072	1 807	272 037





Distributions conditionnelles de l'âge de l'épouse en fonction de l'âge de l'époux.

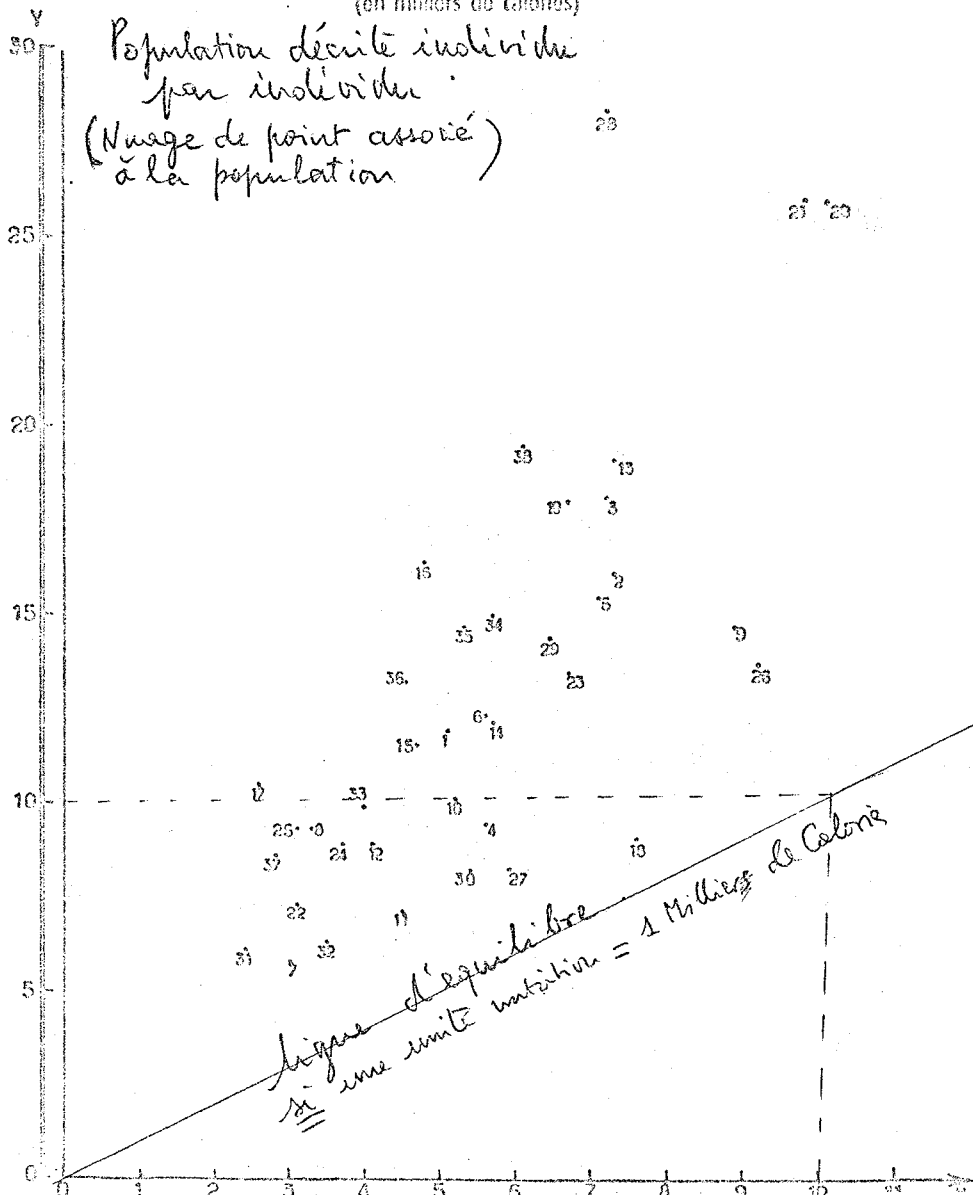


Distributions conditionnelles de l'âge de l'époux en fonction de l'âge de l'épouse.

Autre exemple : Description dans un plan de tous les individus

Ménage n°	x	y	Ménage n°	x	y	Ménage n°	x	y
1	5,1	11,9	14	5,7	12,1	27	5,9	8,2
2	7,3	10,0	15	4,7	11,5	28	7,2	29,2
3	7,2	13,0	16	4,8	10,3	29	6,4	14,5
4	5,5	9,4	17	2,6	10,5	30	5,4	8,2
5	7,1	15,4	18	7,6	9,0	31	2,4	6,1
6	5,6	12,3	19	6,7	17,9	32	3,5	6,3
7	3,0	5,8	20	10,1	25,8	33	4,0	9,9
8	3,3	9,5	21	9,8	25,8	34	5,7	14,9
9	8,9	14,6	22	3,1	7,3	35	5,3	14,6
10	5,2	10,1	23	6,7	13,4	36	4,6	13,2
11	4,5	7,1	24	3,7	8,9	37	2,8	8,6
12	4,1	8,9	25	3,1	9,3	38	6,1	19,4
13	7,3	19,0	26	9,2	13,6			

DISTRIBUTION DE 30 MÉNAGES SUIVANT LE NOMBRE D'UNITÉS DE NUTRITION (x) ET LA CONSOMMATION (y) (en milliers de calories)



- Consommation : nombre de milliers de calories consommées par ménage (y)
  - Unités de nutrition : besoins alimentaires du ménage (calculés théoriquement suivant nombre de personnes du ménage, sexe, âge, travail effectué ...) = (x)
- (Sondage effectué dans une ville du Proche Orient).

Cas particulier : Séries chronologiques.

Lorsque l'un des caractères est le temps, la série statistique s'appelle série chronologique (l'autre caractère est quelconque).

La liaison est fonctionnelle

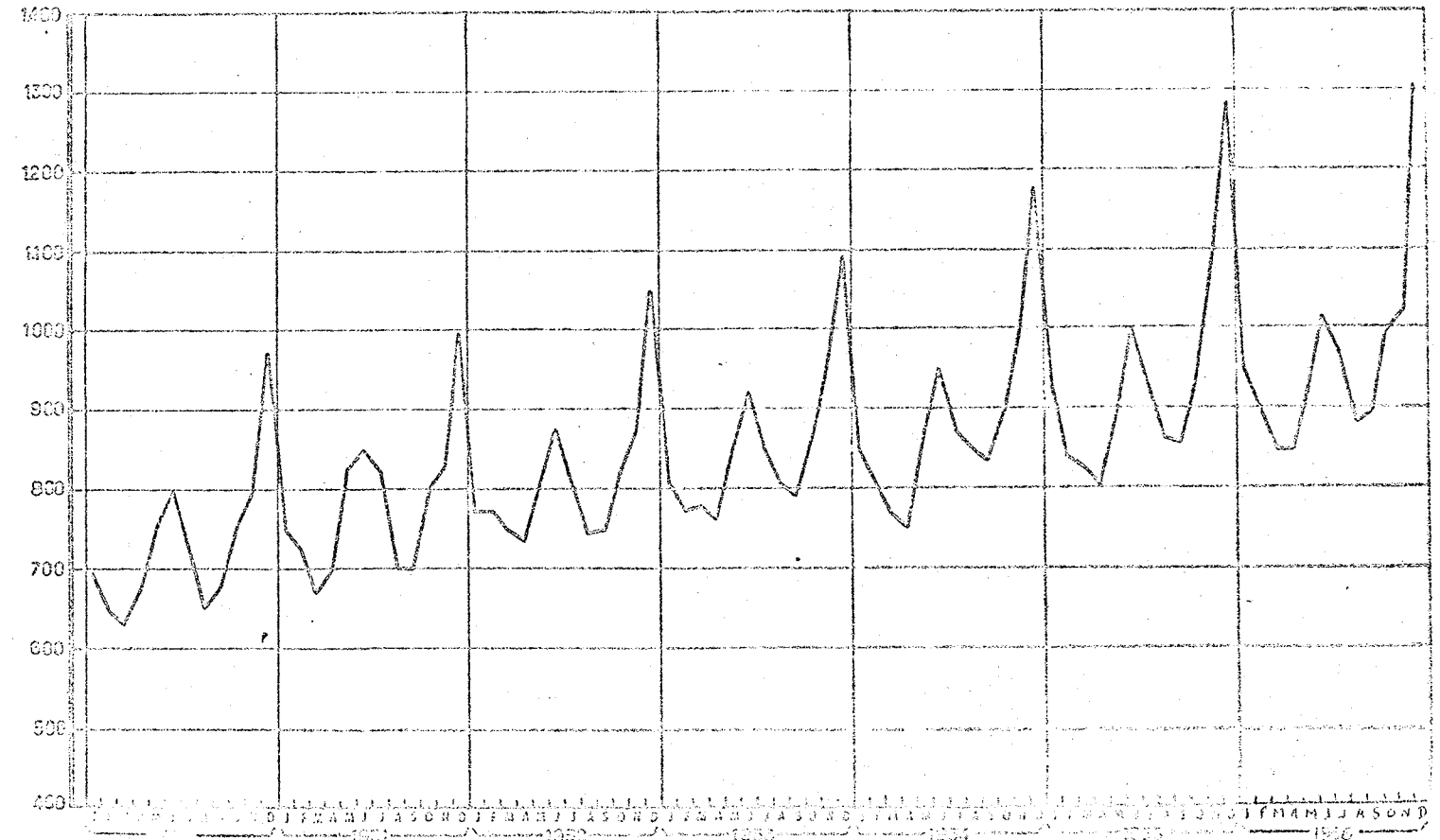
*Série des chiffres d'affaires mensuels d'un rayon d'un grand magasin (exemple) cité par R. Dumas dans L'Entreprise et la Statistique; Dunod éditeur, Paris, 1954).*

Unité : 10 000 francs

Mois \ Année	Année							
	1950	1951	1952	1953	1954	1955	1956	
Janvier .....	700	750	775	815	850	925	945	
Février .....	650	725	775	775	810	810	855	
Mars .....	635	675	750	780	765	825	840	
Avril .....	675	700	735	760	750	800	845	
Mai .....	730	825	810	850	870	850	915	
Juin .....	800	850	870	920	950	1 000	1 015	
Juillet .....	725	825	805	855	875	920	950	
Août .....	650	700	745	810	850	860	875	
Septembre .....	675	700	750	795	835	865	895	
Octobre .....	750	800	825	865	895	920	925	
Novembre .....	800	825	875	960	1 010	1 000	1 120	
Décembre .....	975	1 080	1 050	1 090	1 175	1 285	1 300	

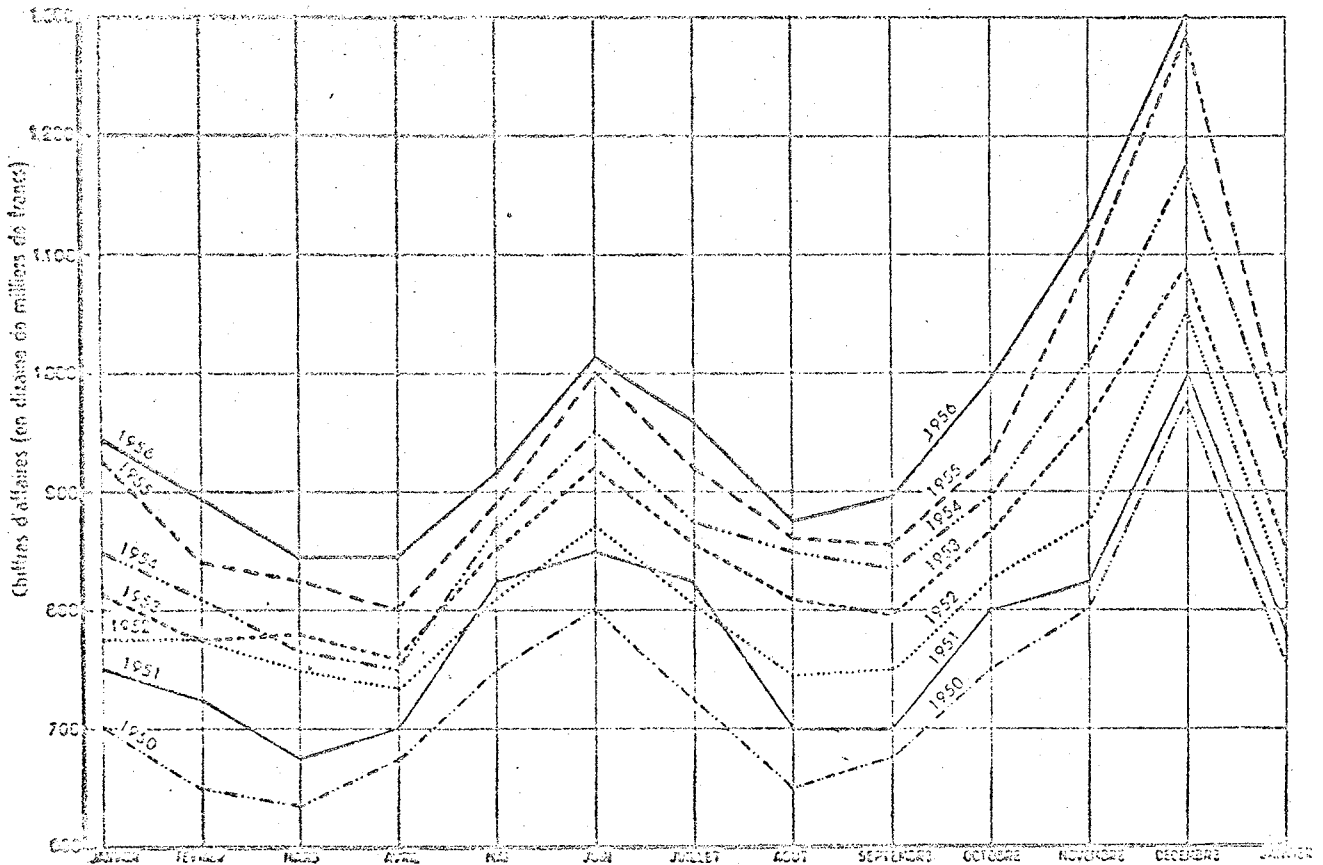
SÉRIE CHRONOLOGIQUE DES CHIFFRES D'AFFAIRES MENSUELS D'UN RAYON D'UN GRAND MAGASIN

Chiffres en 10<sup>4</sup> francs

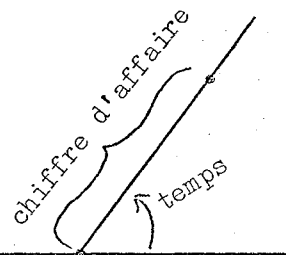


## SÉRIE CHRONOLOGIQUE DES CHIFFRES D'AFFAIRES D'UN RAYON D'UN GRAND MAGASIN

Courbes annuelles superposées

Représentation polaire.

Sachant qu'une année correspond à un angle de  $2\pi$  on peut représenter de façon polaire la série chronologique précédente.



Nota. Il existe bien d'autres représentations graphiques de distributions statistiques à deux caractères (cartogrammes dans la description d'unités géographiques suivant un caractère, représentation graphiques utilisant des papiers fonctionnels, graphique triangulaire...).

## CHAPITRE I

### Modèles statistiques. Théorie de la décision.

#### 1. BUTS DE LA STATISTIQUE.

La démarche essentielle de la statistique consiste à préciser à partir de données expérimentales un modèle probabiliste mal connu.

Il y aura de très nombreuses façons de traduire le verbe préciser : estimer, tester, classer, décider, etc...

Par exemple, si l'on considère deux traitements par engrais d'une même plante, s'ils donnent les tailles suivantes pour 10 pousses

30	37	39	34	36	38	32	32	31
31	38	38	31	31	34	35	32	32

peut-on dire et en quel sens que l'un est meilleur que l'autre si le modèle probabiliste consiste à associer à chaque traitement une loi de probabilité, soit  $F_1$  et  $F_2$ , a-t-on  $F_1$  "meilleur" que  $F_2$  en un certain sens (par exemple moyenne de  $F_1 \geq$  moyenne de  $F_2$  ou autre critère, médiane  $F_1 \geq$  médiane  $F_2$ ).

Le travail du statisticien est donc d'examiner, à l'aide d'outils mathématiques, des données, qui sont les résultats d'une expérience et situées dans un certain ensemble de modèles probabilistes ; soit  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta}$  (on a l'habitude de toujours prendre  $\Omega, \mathcal{A}$  identiques pour tous les modèles car en fait  $(\Omega, \mathcal{A})$  représente la structure des données qui est évidemment indépendante du modèle réel et inconnu [fixé par une valeur  $\theta_0$  de  $\theta, \theta_0 \in \Theta$ ]).

Le choix de l'ensemble  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta}$  des modèles relève en principe, de l'expérimentateur de l'utilisateur



mais en pratique, le statisticien est fréquemment amené à participer au choix de la collection de modèles.

La réalisation de l'expérience elle-même relève aussi du statisticien, en ce qu'il participe à l'élaboration du plan (déroulement) de l'expérience. Ceci est dû au fait suivant, la quantité d'information, la précision et la richesse des déductions que l'on peut faire à partir du traitement de données (résultats d'une expérience) dépend de la structure mathématique de ces données. On peut gagner du temps et du matériel en organisant l'expérience d'une certaine manière, c'est là l'objet de l'étude des plans d'expériences.

De plus, il ne sert à rien d'élaborer des modèles complexes si on ne peut pas travailler dessus concrètement. Donc même au point 1, les techniques statistiques d'extraction de l'information à partir des données sont des plus importantes. Cependant dans certains domaines comme les sciences humaines et l'économie, il est fréquent d'avoir à traiter des données brutes car l'expérience n'est pas modifiable par le statisticien.

La statistique peut être développée à partir de plusieurs points de vue :

a) Celui de la théorie de l'information (très lié à des concepts physiques de type entropie) qui a l'avantage de conduire à une théorie riche, mais l'inconvénient de reposer sur des concepts peu intuitifs pour un étudiant actuel en mathématiques. Nous ne l'aborderons pratiquement pas ici.

b) Celui de la théorie de la décision, lié à la pratique des sciences expérimentales et sociales est celui que nous choisirons.

De toute manière, une fois posés les concepts de base, les différents développements se recouvrent et l'ensemble de la théorie garde une certaine unité.

## 2. MODELE STATISTIQUE. STATISTIQUE.

Définition. Soit  $(\Omega, \mathcal{A})$  un espace mesuré,  $(P_\theta)_{\theta \in \Theta}$  une famille de probabilités sur cet espace, indexée par  $\theta$ . L'ensemble  $(\Omega, \mathcal{A}, P_\theta; \theta \in \Theta)$  est appelé modèle statistique. C'est donc une famille d'espaces de probabilité. Ayant choisi un modèle statistique  $(\Omega, \mathcal{A}, P_\theta; \theta \in \Theta)$  et réalisé une expérience, le statisticien doit indiquer quelles sont les valeurs de  $\theta$  qui sont non contradictoires, en un sens probabiliste à bien préciser, avec les résultats de l'expérience.

L'exemple suivant est à la fois simple et non trivial.

Exemple 1. On joue à pile ou face, et on appelle  $\theta$  la probabilité de tirer pile,  $1-\theta$ , celle de tirer face. On peut faire  $n$  tirages ( $n$ -expériences) indépendants.

Soit  $X_i$  la v.a. qui représente la  $i$ -ème expérience :  $X_i = 1$  si on a obtenu pile au  $i$ -ème coup,  $X_i = 0$  sinon. Ce qu'on connaît, c'est donc une suite de nombres  $+1, 0$ , qui représente les résultats de ces  $n$ -expériences,  $(x_1, x_2, \dots, x_n)$ .

D'après la loi des grands nombres, on sait que :

$$P_\theta: \left( \left| \frac{X_1 + X_2 + \dots + X_n}{n} - \theta \right| > \varepsilon \right) = A(\theta, \theta', \varepsilon) \text{ tend vers } 0 \text{ si } n \rightarrow +\infty$$

si et seulement si  $\theta = \theta'$ .

Considérons le résultat de l'expérience faite  $(X_1(\omega) \dots X_n(\omega)) = (x_1, x_2, \dots, x_n)$ , et  $\theta_0$  une valeur telle que

$$\left| \frac{x_1 + x_2 + \dots + x_n}{n} - \theta_0 \right| > \varepsilon.$$

Si  $\theta$  est égal à  $\theta_0$ , on aurait réalisé l'évènement

$$B(\varepsilon) = \left\{ \left| \frac{X_1 + X_2 + \dots + X_n}{n} - \theta_0 \right| > \varepsilon \right\}$$

avec la probabilité  $P_{\theta_0}(B\varepsilon) = A(\theta_0, \theta_0, \varepsilon) = \eta$ .

Si  $\eta$  est petit, on admet qu'il est plus que douteux que l'on puisse se trouver dans la situation  $\theta = \theta_0$ .

Intuitivement, on rejette la valeur  $\theta_0$ , et aussi les valeurs très voisines de  $\theta_0$ .

Statistique - Image d'un modèle statistique par une statistique.

Définition : On appelle statistique sur le modèle statistique  $\{\Omega, \mathcal{A}, (P_\theta)_{\theta \in \Theta}\}$  toute fonction mesurable  $T$  de  $(\Omega, \mathcal{A})$  dans un espace mesurable  $(\Omega', \mathcal{A}')$ .

Observer une statistique  $T$ , c'est observer un résultat déduit de l'expérience ; ou c'est remplacer le résultat  $x$  de l'expérience par la fonction  $T(x)$ . On peut associer à  $T$ , considéré maintenant comme le résultat de l'expérience, un autre modèle statistique

$$\{\Omega', \mathcal{A}', \{T(P_\theta)\}_{\theta \in \Theta}\}$$

appelé image du 1er modèle par  $T$ . Il est fréquent qu'une statistique  $T$  donne sur le paramètre autant d'information que l'expérience : le modèle image par  $T$  du 1er modèle sera alors le modèle naturel. Nous précisons plus tard cette idée (Chapitre 8 : Exhaustivité). Indiquons ici quelques exemples, qui seront précisés au Chapitre 8.

Exemple : Soit  $X = (X_1, X_2, \dots, X_n)$  un  $n$ -échantillon d'une famille de lois  $(\nu_\theta)_{\theta \in \Theta}$  sur  $\mathbb{R}^d$ .

Statistique d'ordre.

Intuitivement, l'ordre des observations n'a pas d'importance puisque les  $X_i$  sont indépendantes, et on peut par exemple remplacer la suite des observations  $(x_1, \dots, x_n)$  par la suite ordonnée  $(x_{\sigma(1)}, \dots, x_{\sigma(n)})$  qui s'en déduit par la permutation  $\sigma$  telle que  $x_{\sigma(1)} \leq x_{\sigma(2)} \leq \dots \leq x_{\sigma(n)}$ . La statistique  $T$  correspondante, appelée statistique d'ordre, a une loi qui dépend de la valeur de  $\theta$ , mais il est intuitif que la loi de  $X$  conditionnée par  $T$  est indépendante de  $\theta$ : si on sait que  $T$  appartient à un borélien  $A$  de  $(\mathbb{R}^d)^n$  la probabilité pour que  $X$  appartienne à un borélien  $B$  de  $(\mathbb{R}^d)^n$  est égale, intuitivement à la proportion des permutés des points de  $B$  qui sont dans  $A$ .

### 3. COMMENT FORMULER UN PROBLEME DE STATISTIQUE. VOCABULAIRE ET FORMALISME DE LA THEORIE DE LA DECISION.

Il existe un certain nombre de divisions traditionnelles de la statistique. L'une concerne la nature de l'ensemble  $\theta$  : si cet ensemble est "simple", par exemple un ouvert de  $\mathbb{R}^d$ , on dit que l'on a un modèle paramétrique à traiter, sinon on parle de cas non-paramétriques. Nous reviendrons plus loin et plus à fond sur cette distinction.

Mais la division la plus importante consiste en la manière d'énoncer la conclusion : si celle-ci est du type " $\theta = \theta_0$ ", on parle de problème d'estimation. Si par contre  $\theta = \theta_1 + \theta_2$ , et que la conclusion est du type  $\theta \in \theta_1$ , on parle de problème de test. Il y a évidemment bien d'autres "conclusions" possibles ; pour préciser ce que nous entendons par conclusion, nous allons introduire le formalisme de la théorie de la décision.

Les décisions que va prendre le statisticien vont être fonction du résultat de ses expériences. Plus précisément, s'il a observé  $X(\omega)$ , sa décision va être une fonction  $h[X(\omega)] = d(\omega)$ .

Définitions. Soit  $(\Omega, \mathcal{A})$  l'espace mesuré du modèle statistique et  $X$  une v.a. sur  $(\Omega, \mathcal{A})$ .

Soit  $(A, \mathcal{D})$  un espace mesuré, appelé espace des décisions.

Une règle de décision,  $d$ , (relative à  $X$ ) est une v.a. de  $(\Omega, \mathcal{A})$  dans  $(A, \mathcal{D})$  mesurable par rapport à la tribu engendrée par  $X$ .

On va avoir besoin de classer les décisions. Pour cela, on se donne au départ une fonction de perte  $W$  : de  $\Theta \times A \rightarrow \mathbb{R}^+$ , mesurable :

$W(\theta_0, d(\omega))$  est la perte qu'entraîne pour le statisticien (ou son employeur) le fait de prendre la décision  $d(\omega)$  lorsque la vraie valeur du paramètre inconnu est  $\theta_0$ .

Les fonctions de perte classiques sont :

$W(\theta, d) = |\theta - d|^2$  (perte quadratique qui intervient dans les problèmes d'estimation)

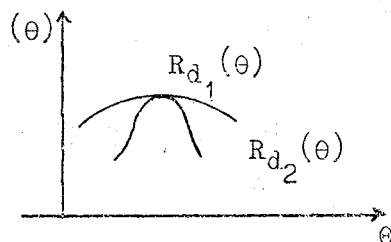
$W(\theta, d) = 1_{\theta_2}(d)$  si  $\theta \in \theta_1$   
 $= 1_{\theta_1}(d)$  si  $\theta \in \theta_2$  (pénalisation de 1 pour toute mauvaise décision).

Comme on ne connaît pas  $\theta$ , on calcule, pour chaque valeur de  $\theta$  la perte de moyenne dite risque

Définition. On appelle fonction de risque pour la décision  $d$ , la fonction  $R_d$  de  $\Theta \rightarrow \mathbb{R}^+$ , définie par :

$$R_d(\theta) = E_{\theta}[W(\theta, d(\omega))] \\ = \int_{\Omega} W(\theta, d(\omega)) dP_{\theta}(\omega).$$

L'ordre naturel sur les fonctions de domaine  $\Theta$  établit une relation de pré-ordre sur les règles de décision :



une décision  $d_1$  est dite meilleure que la décision  $d_2$  si  $R_{d_1}(\theta) \leq R_{d_2}(\theta)$  pour tout  $\theta \in \Theta$ .

Définitions. Une règle de décision est dite admissible, s'il n'en existe pas de strictement meilleure.

S'il existe un plus grand élément pour cet ordre, on parle de règle optimale.

En général, pour des raisons pratiques (calcul) ou liées au problème, on impose aux règles de décision des contraintes limitatives ; on se limite donc par avance à un ensemble  $\mathcal{D}$  de règles.

Exemple : une décision  $d_0$  est dite minimax, si elle minimise le risque maximum :

$$\inf_d \sup_{\theta} R_d(\theta) = \sup_{\theta} R_{d_0}(\theta).$$

On peut ne s'intéresser qu'à ce type de décisions. Nous verrons plus tard d'autres classes de décisions.

En résumé, un problème statistique, vu dans le cadre de la théorie de la décision statistique est la donnée d'un modèle  $(\Omega, \mathcal{A}, P_{\theta})_{\theta \in \Theta}$  ; le choix d'une v.a.  $X$  qui correspond à l'expérience, d'un ensemble  $(A, \mathcal{A})$  de valeurs de décisions possibles, éventuellement la donnée d'une classe  $\mathcal{D}$  de v.a. à valeurs dans  $A$ ,  $\mathcal{B}(X)$ -mesurables (ensemble des règles de décision auxquelles on se limite), d'une fonction de perte  $W$  de  $\Theta \times A \rightarrow \mathbb{R}^+$

Résoudre le problème, c'est trouver les règles de décision qui sont admissibles. Toutefois, le formalisme donné ne permet pas d'aller plus loin : c.à.d. il ne permet pas de choisir entre plusieurs règles admissibles s'il en existe. Le choix dépend alors du problème concret et du type de qualité qu'il vaut mieux imposer à la règle cherchée de ce point de vue (cf. exemple en fin du Chapitre III).

On peut généraliser la notions de règle de décision à celle de stratégie (ou stratégie aléatoire).

Une stratégie est une application (X-mesurable)  $(\Omega, \mathcal{A}) \xrightarrow{S} \pi(A, \mathcal{B})$  où  $\pi$  est l'ensemble des probabilités sur A.

Si  $B \subset A$ , alors on choisit  $d \in B$  avec la probabilité  $S(X(\omega))(B)$ .

Une règle de décision est aussi dite stratégie pure (ou non aléatoire). Elle correspond au cas où  $S(X(\omega)) = \delta_d(X(\omega))$ ,  $\delta_a$  mesure de Dirac en a.

Le risque d'une telle stratégie pour la fonction de perte  $W(\theta, \cdot)$  est

$$R_S(\theta) = \int_{\pi(A)} W(\theta, u) S(\omega, du)$$

si  $S(\omega, du) = S(X(\omega))$ . (Elément de  $\pi(A)$ ).

On note  $\hat{\mathcal{F}}$  l'ensemble des stratégies et  $\mathcal{F}$  l'ensemble des règles de décision. Le lecteur définira sans peine, par analogie avec les décisions des stratégies admissibles, optimales, minimax.

Remarque : Soit  $\phi \in \hat{\mathcal{F}}$  une stratégie admissible de risque constant ; alors  $\phi$  est minimax .

En effet, si  $\phi' \in \hat{\mathcal{F}}$ , il est nécessaire que  $\sup_{\theta \in \Theta} R(\theta, \phi')$  soit supérieur ou égal à la valeur constante

$k = R(\theta, \phi)$  ( $\theta \in \Theta$ ) sinon  $\phi'$  serait meilleure que  $\phi$  et  $\phi$  ne serait pas admissible. Donc :

$$k = R(\theta, \phi) = \sup_{\theta \in \Theta} R(\theta, \phi) = \inf_{\phi' \in \hat{\mathcal{F}}} \sup_{\theta \in \Theta} R(\theta, \phi')$$

Risques bayésiens.  $\theta$  est supposé mesuré.  
 Supposons que l'on munisse  $\theta$  d'une probabilité  $\mu$  ou  
"loi a priori" du paramètre. Cela veut dire que l'on a a priori des  
 idées sur les valeurs possibles de  $\theta$ . Alors le risque moyen de  
 $\phi \in \mathcal{F}$  est :

$$r(\mu, \phi) = \int d\mu(\theta) R(\theta, \phi) .$$

en notant  $R(\theta, \phi)$  pour  $R_{\theta}(\phi)$

Définitions: Soit  $\mu$  une probabilité sur  $\theta$ .  
 Soit  $\phi \in \mathcal{F}$ . On appelle risque bayésien de  $\phi$  par rapport à la  
 loi a priori  $\mu$ , le nombre :

$$r(\mu, \phi) = \int_{\theta} d\mu(\theta) R(\theta, \phi)$$

On dit que  $d$  est une fonction de décision bayésienne, par rapport à  
 $\mu$ , si :

$$r(\mu, d) = \inf_{d' \in \mathcal{F}} r(\mu, d')$$

(On dit que  $\phi$  est une stratégie bayésienne par rapport à  $\mu$ , si :

$$r(\mu, \phi) = \inf_{\phi' \in \mathcal{F}} r(\mu, \phi')$$

au lieu de dire "bayésienne par rapport à  $\mu$ " on dira aussi  
" $\mu$  admissible".

Remarque : Il est souvent difficile de justifier théoriquement la  
 donnée a priori d'une loi  $\mu$ . C'est cependant parfois possible : dans  
 le cadre de la théorie des jeux, si  $\theta$  est l'ensemble des jeux  
 possibles du joueur ennemi,  $\mu$  correspond à la stratégie de ce joueur.  
 L'intérêt des notions bayésiennes est aussi technique ; dans de nombreux  
 cas (cf. chap. 2) on sait calculer des stratégies bayésiennes, alors  
 qu'il est difficile d'obtenir directement des stratégies admissibles.  
 Cela donne souvent des stratégies admissibles.

Soient en effet  $\mu$  une loi a priori et  $\phi$  une stratégie bayésienne par  
 rapport à  $\mu$ . Soit  $\phi'$  une stratégie aussi bonne que  $\phi$

$$\forall \theta, \quad R(\theta, \phi') \leq R(\theta, \phi)$$

D'où  $r(\mu, \phi') \leq r(\mu, \phi)$ .

Cela implique souvent qu'il n'existe aucun  $\theta_0$  tel que  
 $R(\theta_0, \phi') < R(\theta_0, \phi)$ , c'est à dire que  $\phi$  est admissible.



En particulier (vérifiez le) si :

- $\phi$  est l'unique stratégie bayésienne par rapport à  $\mu$   
( $P_\theta$  p.s. pour tout  $\theta$ )
- $\Theta$  est dénombrable et  $\mu$  charge tous les points de  $\Pi$
- $\Theta$  est un ouvert de  $\mathbb{R}^d$ ,  $\mu$  charge tout ouvert inclus dans  $\Theta$  (le support de  $\mu$  est  $\Theta$ ) et la fonction  $\theta \mapsto R(\theta, \phi')$  est continue pour tout  $\phi' \in \hat{\mathcal{F}}$ .

Tout cela est vrai en remplaçant les stratégies (bayésiennes ou admissibles) par les fonctions de décision. (Cf. Chapitre II)

Un exemple de problème décisionnel bayésien.

Soit 2 populations de lois normales sur  $\mathbb{R}^2$ , respectivement  $N_{\theta_1} = N(0, I)$ ,  $N_{\theta_2} = N(m, 2I)$ ,  $I$  identité. On donne comme mesure a priori sur  $(\theta_1, \theta_2)$  la mesure  $(1/2, 1/2)$ . Tracer une droite  $D$  dans le plan  $\mathbb{R}^2$ , telle que la règle de décision qui consiste à classer un point dans la première population s'il est d'un côté de la droite, dans la deuxième sinon, soit optimale parmi ces règles.

Soit  $\lambda X + \mu y = c$  la droite  $D$  avec  $\lambda^2 + \mu^2 = 1$ .

Soit  $X$  un point de coordonnées  $(X_1, X_2)$  la loi de  $\lambda X_1 + \mu X_2$  est si  $\theta_1$  est vrai  $N(0, I)$ , si  $\theta_2$  est vrai  $N(\langle m, a \rangle, 2I)$ , si  $a$  est le vecteur  $(\lambda, \mu)$ .

Les risques sont donc

$$R(\theta_1, D) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^c e^{-x^2/2} dx$$

$$R(\theta_2, D) = \frac{1}{2\sqrt{\pi}} \int_c^{\infty} e^{-\frac{(x - \langle m, a \rangle)^2}{2}} dx$$

et le risque bayésien est  $1/2 |R(\theta_1, D) + R(\theta_2, D)|$

Minimisons d'abord en  $a$ , puis en  $c$ .

$\min_a R(\theta_2, D)$  est atteint pour  $\langle m, a \rangle$  maximum, soit comme  $a$  est normé, pour  $\lambda = \frac{m_1}{m_1 + m_2}$ ,  $\mu = \frac{m_2}{m_1 + m_2}$

$D$  est donc orthogonale à  $m$ ,  $c$  est alors défini

$$e^{-c^2/2} - e^{-\frac{(c - \|m\|)^2}{2}} = 0$$

soit  $c = \frac{\|m\|}{2}$ .

COMPLEMENT DE COURS

Modèle fini de la théorie de la décision

Les notations sont celles du paragraphe précédent. Supposons que  $\Theta$  est un ensemble fini  $\Theta = \{\theta_1, \theta_2, \dots, \theta_k\}$ . A chaque  $\phi \in \hat{\mathcal{Y}}$ , on peut associer un point  $M_\phi \in \mathbb{R}^k$  où  $M_\phi = \{R(\theta_i, \phi)\}_{1 \leq i \leq k}$ .

Soient  $\hat{S} = \{M_\phi ; \phi \in \hat{\mathcal{Y}}\}$  et  $S = \{M_\phi ; \phi \in \mathcal{Y}\}$ .

L'ensemble  $\hat{S}$  est un convexe de  $\mathbb{R}^k$ , car si  $\phi$  et  $\phi'$  sont dans  $\hat{\mathcal{Y}}$  et si  $\alpha \in ]0, 1[$ ,  $\alpha\phi + (1-\alpha)\phi' \in \hat{\mathcal{Y}}$  et  $M_{\alpha\phi + (1-\alpha)\phi'} = \alpha M_\phi + (1-\alpha)M_{\phi'} \in \hat{S}$ .

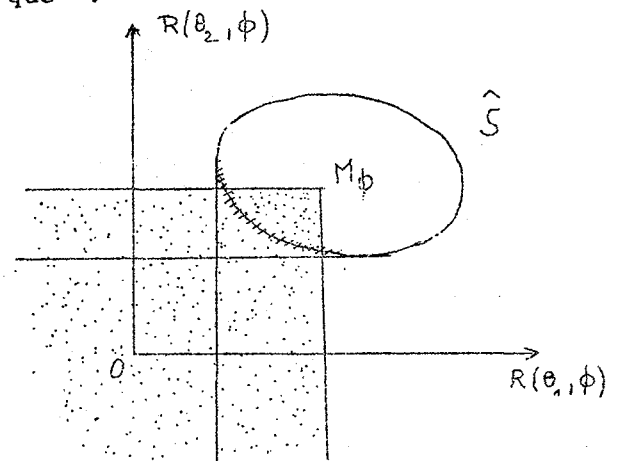
Supposons  $\phi' \leq \phi$ . Alors le point  $M_{\phi'}$  est dans l'ensemble  $Q[M_\phi]$ , si l'on associe à tout point  $M = \{x_1, x_2, \dots, x_k\}$  de  $\mathbb{R}^k$  l'ensemble  $Q[M] = \{y_1, y_2, \dots, y_k\}$ ;  $y_1 \leq x_1, y_2 \leq x_2, \dots, y_k \leq x_k$ .

Donc dire que  $\phi$  est admissible c'est dire que :

$$Q[M_\phi] \cap \hat{S} = M_\phi$$

exemple ( $k = 2$ )

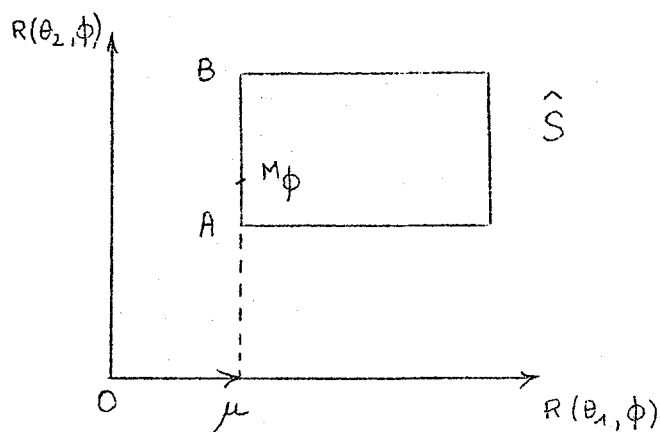
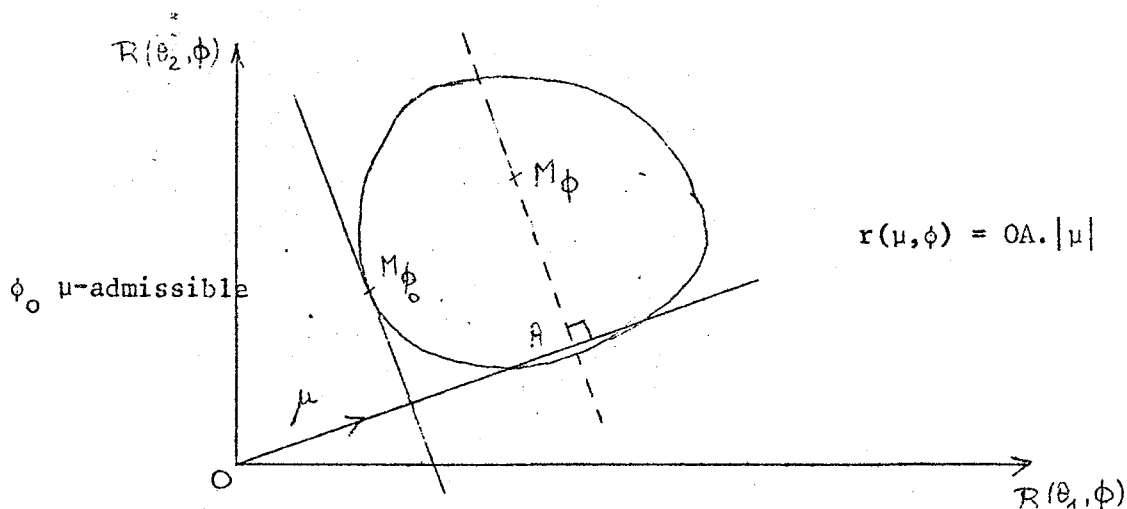
- la zone sombre est  $Q(M_\phi)$
- si  $\hat{S}$  est fermé, la zone hachurée est l'ensemble des  $M_\phi$  pour  $\phi$  admissible.



Soit  $\mu$  une loi a priori sur  $\theta$  ; posons  $\mu_i = \mu(\theta_i)$  ( $1 \leq i \leq k$ ) ;  
 $\mu$  est le vecteur  $(\mu_i)_{1 \leq i \leq k}$ . Le risque bayésien de  $\phi \in \hat{S}$  est :

$$r(\mu, \phi) = \sum_{i=1}^k \mu_i R(\theta_i, \phi) = \langle \mu, M_\phi \rangle$$

Si une stratégie  $\phi$  est  $\mu$ -admissible,  $\hat{S}$  est d'un seul côté de la droite passant par  $M_\phi$  et orthogonale à  $\mu$  ; on sait que si les  $\mu_i$  sont tous non nuls une telle stratégie est admissible ; c'est faux dans le cas contraire comme le montre le 2ème schéma ci-dessous :

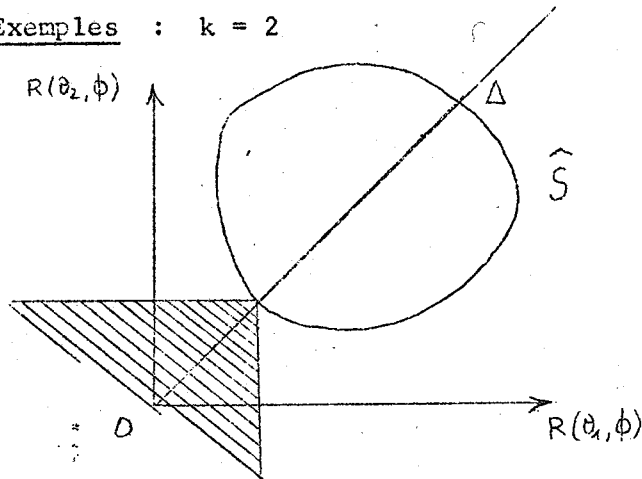


$$\mu = (1, 0)$$

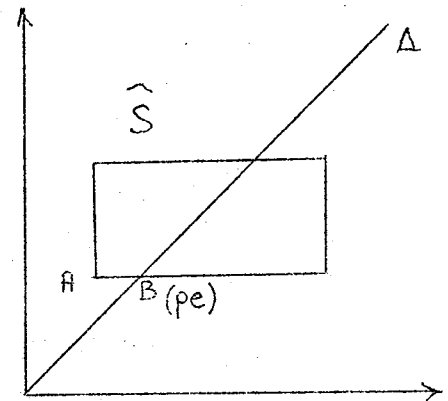
$\phi$  est  $\mu$ -admissible, si  $M_\phi \in AB$ , mais n'est aussi admissible que si  $M_\phi = A$ .

Minimiser le risque maximum, c'est minimiser l'abscisse de  $M_\phi$ , si  $M_\phi$  est en dessous de la 1ère bissectrice  $\Delta$  et l'ordonnée de  $M_\phi$ , si  $M_\phi$  est en dessus. Par conséquent, le minimum du risque maximum est le nombre  $\rho = \{\sup\{\alpha ; Q[\alpha e] \cap \hat{S} = \emptyset\}$ , où  $e$  est le vecteur de  $\mathbb{R}^k$ , dont toutes les coordonnées valent 1 ; si  $M_\phi \in Q[\rho e] \cap \hat{S}$ , alors  $\phi$  est une stratégie minimax. En particulier, si  $S$  est compact, il existe une stratégie minimax.

Exemples :  $k = 2$



$\phi$  est minimax, de risque constant et admissible.



$\phi$  est minimax, si  $M_\phi \in AB$ .  
Mais  $\phi$  n'est admissible que si  $M_\phi = A$ , et de risque constant que si  $M_\phi = B$ .

Pour permettre de résoudre des problèmes simples dans le cas fini nous devons préciser la structure de  $\hat{S}$ , lorsque  $A$  est fini.

Préliminaires : Soient  $S_1$  et  $S_2$  2 convexes disjoints de  $\mathbb{R}^k$ . Il existe une forme linéaire  $u$  sur  $\mathbb{R}^k$  et un nombre  $c \in \mathbb{R}$ , tels que  $u(x) \geq c$  pour  $x \in S_1$  et  $u(x) \leq c$  pour  $x \in S_2$ . Autrement dit, il existe un hyperplan qui sépare  $S_1$  et  $S_2$ . Nous ne montrerons pas cette propriété, qui est une forme du théorème de Hahn Banach, qu'on peut montrer élémentairement sur  $\mathbb{R}^k$ . On en déduit le :

lemme : Soit  $S$  un convexe de  $\mathbb{R}^k$ ; soit  $Z$  un vecteur aléatoire de dimension  $k$ , tel que  $Z \in S$  presque sûrement et que  $E(Z)$  existe; alors  $E(Z) \in S$ .

Démonstration - Le théorème est vrai si  $k = 1$ , car un convexe est ici un intervalle. Supposons le vrai pour  $k - 1$ . Si  $E(Z)$  n'est pas dans  $S$ , il existe une forme linéaire  $u$  telle que  $u(E(Z)) \leq u(x)$  pour tout  $x \in S$ . Donc :

$$u[E(Z)] = E[u(Z)] \leq u(Z) \quad \text{p.s.}$$

cela implique que  $u(Z)$  vaut  $E[u(Z)]$  (p.s.).

Soit  $S' = \{x; u(x) = E[u(Z)]\} \cap S$ ; c'est un convexe d'un espace euclidien à  $k - 1$  dimensions et  $Z \in S'$ .

D'après l'hypothèse de récurrence,  $E(Z) \in S' \subset S$ .

Proposition : Si  $A$  est fini,  $\hat{S}$  est l'enveloppe convexe de  $S$  (le plus petit convexe contenant  $S$ ).

Démonstration :  $\hat{S}$  est convexe ; il suffit donc de montrer que  $\hat{S}$  est contenu dans l'enveloppe convexe de  $S_0$  de  $S$ .

Soit  $A = (a_1, a_2, \dots, a_n)$ . Soit  $\nu$  une probabilité sur  $A$ . Munissons  $[0, 1]$  de la mesure de Lebesgue  $P$  (soit  $E$  l'espérance par rapport à  $P$ ). Il existe une fonction  $Z_\nu$  de  $[0, 1]$  dans  $A$ , de loi  $\nu$  ; il suffit en effet de diviser  $[0, 1]$  en  $m$  intervalles disjoints de longueurs respectives  $\nu(a_i)$  ( $1 \leq i \leq k$ ) et de poser  $Z_\nu = a_i$  sur l'intervalle de longueur  $\nu(a_i)$ . Soit alors une stratégie  $\phi \in \mathcal{F}$  ; à chaque  $x \in E$  correspond une mesure  $\phi(x)$  sur  $A$  donc une fonction  $Z_{\phi(x)}$  de  $[0, 1]$  dans  $A$  -

$$\begin{aligned} R(\theta, \phi) &= \int dP_0(x) \int d\phi(x) (a) W(\theta, a) = \int dP_0(x) E \left[ W(\theta, Z_{\phi(x)}) \right] \\ &= E \left[ \int dP_0(x) W(\theta, Z_{\phi(x)}) \right] \end{aligned}$$

Pour chaque  $\omega \in [0, 1]$ ,  $x \rightarrow Z_{\phi(x)}(\omega)$  est une règle de décision notée  $Z_\phi(\omega)$ .

$$R(\theta, \phi) = \int dP(\omega) R(\theta, Z_\phi(\omega))$$

Pour chaque  $\omega \in [0, 1]$ ,  $R(\theta, Z_\phi(\omega))$  est dans  $S$  donc dans  $S_0$ . Par suite  $R(\theta, \phi) \in S_0$  et  $\hat{S} \subset S_0$ .

Conséquence : Si  $\theta$  et  $A$  sont finis et si  $S$  est compact, alors  $\hat{S}$  est convexe compact ; c'est en effet l'image de  $[0, 1]^k \times S^k$  dans  $\mathbb{R}^k$  par l'application continue

$$((\lambda_i)_{1 \leq i \leq k}, (x_i)_{1 \leq i \leq k}) \longmapsto \sum_{i=1}^k \lambda_i x_i.$$

En particulier, si  $\Omega$  est fini,  $S$  est fini et  $\hat{S}$  est convexe compact.

Classes complètes de stratégies

Une classe complète de stratégies est une partie  $C \subset \mathcal{F}$  telle que si  $\phi \in \mathcal{F} \setminus C$ , il existe une stratégie dans  $C$  meilleure que  $\phi$ . Une classe essentiellement complète de stratégies est une partie  $C \subset \mathcal{F}$ , telle que si  $\phi \in \mathcal{F} \setminus C$ , il existe une stratégie dans  $\mathcal{F} \setminus C$  aussi bonne que  $\phi$ . Une classe complète (resp. essentiellement complète)  $C$  est minimale, si il n'y a pas de sous ensemble strict de  $C$  qui soit une classe complète (resp. essentiellement complète).

Le théorème suivant permettra, dans les problèmes d'estimation, de ne considérer que les fonctions de décision (cf. Chapitre II).

Théorème

Supposons que  $A$  est un convexe de  $\mathbb{R}^d$  et que pour chaque  $\theta \in \Theta$ ,  $a \mapsto W(\theta, a)$  est une fonction convexe continue. Si  $\phi \in \mathcal{F}$  et si pour chaque  $x$ ,  $\int |a| \phi(x, da)$  est fini, alors la fonction de décision  $\psi(x) = \int a \phi(x, da)$  est aussi bonne que  $\phi$ . Par suite si  $\int |a| \phi(x, da)$  est finie pour toute  $\phi \in \mathcal{F}$ , la classe  $\mathcal{S}$  est essentiellement complète.

Remarque : Une fonction d'un convexe  $C$  dans  $\mathbb{R}$  est convexe si :

$$\alpha \in ]0, 1[ , x \in C \text{ et } y \in C \implies f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha) f(y)$$

En particulier si  $d = 1$ , et si  $f$  est 2 fois dérivable, avec une dérivée seconde positive,  $f$  est convexe. Si  $d$  est quelconque, si  $|\cdot|$  est la valeur absolue dans  $\mathbb{R}^d$  et si  $0 \in \mathbb{R}^d$ ,  $x \mapsto |0-x|$  est convexe ; si  $\phi$  est une fonction convexe de  $\mathbb{R}_+$  dans  $\mathbb{R}$ ,  $\phi(|0-x|)$  est convexe ; en particulier  $|0-x|^k$  est convexe pour  $k \geq 1$ .

Cas particulier : Soit  $\alpha$  une fonction de  $\Theta$  dans un ensemble convexe  $A$  de  $\mathbb{R}^d$ . On estime  $\alpha(0)$  à l'aide de la fonction de perte quadratique  $W(\theta, a) = |\alpha(0) - a|^2$ . Alors pour  $\phi \in \mathcal{F}$ ,  $\int |\alpha(0) - a|^2 \phi(x, da)$  est fini et  $\int |a| \phi(x, da) < \infty$ .

Donc pour ces problèmes, on pourra se limiter aux fonctions de décision (l'introduction de stratégies aléatoires n'ajoute rien) : dans la théorie de l'estimation une fonction de décision est appelée estimateur.

La remarque précédente serait encore valable pour  $W(0, a) = |a(0) - a|^k$  avec  $k \geq 1$ . Par contre, dans les problèmes de test nous aurons besoin de stratégies aléatoires.

Pour démontrer le théorème nous avons besoin du lemme suivant :

Lemme : Inégalité de Jensen

Soit  $C$  un ensemble convexe non vide de  $\mathbb{R}^d$ . Soit  $f$  une fonction convexe continue de  $C$  dans  $\mathbb{R}$ . Soit  $Z$  un vecteur aléatoire de dimension  $d$ , sur un espace de probabilité  $(\Omega, \mathcal{A}, P)$ . On suppose que  $Z \in C$ , P p.s. et que  $Z$  est intégrable. Alors :

- a)  $E(Z) \in C$
- b)  $f(E(Z)) \leq E f(Z)$

Démonstration - Nous avons déjà vu que  $E(Z) \in C$ . Soit alors  $C_1$  le convexe de  $\mathbb{R}^{d+1}$  défini par :

$$C_1 = \{(x_1, \dots, x_{d+1}) ; (x_1, \dots, x_d) \in C, x_{d+1} > f(x_1, \dots, x_d)\}$$

Le point  $\{E(Z), f(E(Z))\}$  est un point de  $\mathbb{R}^{d+1}$  qui n'est pas dans  $C_1$ . Il existe donc un vecteur  $(\alpha_i)_{1 \leq i \leq d+1}$  de  $\mathbb{R}^{d+1}$  tel que pour tout  $(x_1, \dots, x_{d+1}) \in C_1$

$$\sum_{i=1}^{d+1} \alpha_i x_i \geq \sum_{i=1}^d \alpha_i E(Z_i) + \alpha_{d+1} f(E(Z))$$

(où  $(Z_i)_{1 \leq i \leq d}$  sont les composantes de  $Z$ ).

On peut prendre  $x_{d+1}$  aussi grand que l'on veut, donc  $\alpha_{d+1} \geq 0$ .

Mais le point  $\{(x_i)_{1 \leq i \leq d} = x, f(x)\}$  appartient à l'adhérence de  $C_1$  pour tout  $x \in C$  ; donc pour  $x = (x_i)_{1 \leq i \leq d} \in C$

$$\alpha_{d+1} f(x) \geq \sum_{i=1}^d \alpha_i [E(Z_i) - x_i] + \alpha_{d+1} f(E(Z))$$

Par suite, p.s. :

$$\alpha_{d+1} f(Z) \geq \sum_{i=1}^d \alpha_i (E(Z_i) - Z_i) + \alpha_{d+1} f(E(Z))$$

\* Supposons  $\alpha_{d+1} > 0$ . Alors en intégrant on obtient :

$$E(f(Z)) \geq f(E(Z))$$



\* Si  $\alpha_{d+1} = 0$ , on obtient  $\sum_{i=1}^d \alpha_i [E(Z_i) - Z_i] \leq 0$ , mais l'espérance de cette v.a. étant nulle, on a  $\sum_{i=1}^d \alpha_i (E(Z_i) - Z_i) = 0$ . Si  $d = 1$ , cela veut dire que  $Z = E(Z)$  p.s. et  $f(Z) = f(E(Z)) = E(f(Z))$  p.s. Sinon,  $Z$  est un vecteur aléatoire à valeurs dans le convexe  $S' = S \cap \{(x_i)_{1 \leq i \leq d} ; \sum_{i=1}^d \alpha_i (E(Z_i) - x_i) = 0\}$ , qui est un convexe d'un espace euclidien de dimension  $d-1$ . Le théorème est donc prouvé par récurrence sur la dimension  $d$ .

### Démonstration du théorème

Avec les hypothèses faites, on a pour tout  $x \in E$  :

$$W(0, \psi(x)) = W\left[0, \int a \phi(x, da)\right] \leq \int W(0, a) \phi(x, da)$$

On le voit en écrivant l'inégalité de Jensen avec la probabilité  $\phi(x, \cdot)$ .

D'où  $R(0, \psi) \leq R(0, \phi)$ .

### EXERCICES

- ① Dans le cadre de la théorie des jeux, soit  $\mathcal{Y} = A = ]0 + \infty[$  et  $W(0, a) = e^{-\theta a}$  quels que soient le paramètre  $\theta$  élément de  $\mathcal{Y}$ , et l'action  $a$  choisie par le statisticien dans  $A$ , c'est à dire que le statisticien a intérêt à choisir la plus grande valeur possible pour  $a$ . Soit  $F$  la fonction de répartition d'une variable aléatoire  $Z$  à valeurs dans  $]0 + \infty[$ , telle que  $EZ = +\infty$ .  $F$  peut être considéré comme un élément de  $\hat{\mathcal{Y}}$ . Montrer qu'aucun élément de  $\mathcal{Y}$ , ensemble des stratégies pures, ne peut être aussi bon que la stratégie  $F$ .

### 4. DONNEES - ECHANTILLON - ECHANTILLONNAGE - LOIS CLASSIQUES - NOTATIONS CLASSIQUES.

En statistique descriptive nous avons vu quelques manières de décrire une population considérée dans son entier. Le modèle statistique est alors fini. En statistique mathématique ou inductive, on cherche à partir d'un certain nombre de données que l'on fabrique avec l'aide de l'expérimentateur ou soi-même, à obtenir des renseignements sur le modèle statistique.

Un échantillon est un sous-ensemble de l'ensemble population.

## I - Echantillonnage d'une population finie

### a) Description

Considérons une population finie de  $N$  individus parmi lesquels on distingue  $r$  types distincts ; pour  $1 \leq i \leq r$ , soit  $N_i$  le nombre d'individus de type  $i$ . Pour déterminer les valeurs des  $N_i$ , on peut faire un recensement, c'est-à-dire observer toute la population. Mais on préfère souvent, faute de temps ou d'argent, effectuer un "sondage", c'est-à-dire prélever un échantillon de  $n$  individus de la population et déduire de cet échantillon des estimations sur les paramètres  $N_1, N_2, \dots, N_r$ .

Les résultats d'un sondage seront donc une suite de  $n$  éléments de  $\{1, 2, \dots, r\}$ , autrement dit un point  $(x_i)_{1 \leq i \leq n} \in \{1, 2, \dots, r\}^n = E$ . Le point  $(x_i)_{1 \leq i \leq n}$  représente le sondage, où le  $i^{\text{ème}}$  individu examiné est de type  $r$ . A chaque sondage  $x = (x_i)_{1 \leq i \leq n}$ , on peut associer la suite  $\{n_1(x), n_2(x), \dots, n_r(x)\}$  où pour  $1 \leq j \leq r$ ,  $n_j(x)$  est le nombre des individus examinés de type  $j$  :

$$n_j(x) = \sum_{i=1}^n 1_{\{j\}}(x_i)$$

On peut envisager 2 types de sondages :

### b) Sondage avec remise

A chacun des  $n$  tirages, la probabilité de tirer un individu de type  $i$  est la fréquence de ce type  $p_i = \frac{N_i}{N}$ . Les  $n$  tirages sont indépendants entre eux. Donc  $E$  est muni d'une probabilité  $P$  définie par :

$$P(x) = P[(x_1, \dots, x_n)] = p_1^{n_1(x)} \cdot p_2^{n_2(x)} \cdot \dots \cdot p_r^{n_r(x)}$$

Cette probabilité dépend du paramètre  $(p_1, p_2, \dots, p_r)$ .

On obtient le modèle statistique

$$\{E, \mathcal{P}(E), (P_{p_1, p_2, \dots, p_r})_{(p_1, p_2, \dots, p_r) \in \Pi}\}$$

où  $\Pi$  est le sous ensemble des points de  $[0, 1]^r$  dont la somme des coordonnées vaut 1.

Les  $n$  fonctions coordonnées  $(X_i)_{1 \leq i \leq n}$  sont des v.a. indépendantes de loi  $v_{p_1, \dots, p_r} = \sum_{i=1}^r p_i \delta_i$ .  $X_i$  représente la valeur obtenue au

$i^{\text{ème}}$  tirage. On a un  $n$ -échantillon de la famille de lois

$$\{v_{p_1, p_2, \dots, p_r} \}_{(p_1, \dots, p_r) \in \Pi}$$

Le vecteur  $x \mapsto \bar{n}(x) = \{n_i(x)\}_{1 \leq i \leq r}$  est une statistique très

importante. Si  $\{i_1, \dots, i_r\}$  est une suite de  $r$  nombres de

$\{0, 1, \dots, n\}$  dont la somme vaut  $n$ , les sondages  $x$  tels que  $\bar{n}(x) = \{i_1, \dots, i_r\}$  sont équiprobables ; il y a  $\frac{n!}{\prod_{k=1}^r i_k!}$  suites de ce type.

On obtient donc :

$$\begin{aligned} P_{p_1, \dots, p_r} & (n_1 = i_1, \dots, n_r = i_r) \\ &= \frac{n!}{\prod_{k=1}^r i_k!} \prod_{k=1}^r (p_k)^{i_k} \end{aligned}$$

si  $\sum_{k=1}^r i_k = n$  ; sinon la probabilité que  $\bar{n} = (i_k)_{1 \leq k \leq r}$  est nulle.

Supposons  $r = 2$  ; alors on remplace le type "2" par la notation "0"

Le modèle est alors :

$$\left\{ \{0, 1\}^n = \Omega, \mathcal{P}(\Omega), (P_p)_{p \in [0, 1]} \right\}$$

Les v.a.  $(X_i)_{1 \leq i \leq n}$  sont indépendantes et leur loi est la loi de

bernoulli  $b(1, p) = (1-p) \delta_0 + p \delta_1$ . La v.a.  $n_1(x)$  a la loi binomiale  $b(n, p) = \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \delta_k$ .

### c) Sondage sans remise

Contrairement au cas précédent un sondage sans remise revient à extraire de la population un bloc non ordonné de  $n$  individus.

Il y a  $\binom{N}{n}$  tels blocs. La donnée du résultat d'un sondage est la donnée des nombres  $(n_i)_{1 \leq i \leq r}$  de tirages du sondage  $x$  où on obtient  $i$ . Il y a  $\prod_{i=1}^r \binom{N_i}{n_i}$  tels blocs non ordonnés de  $n$  individus.

Le modèle statistique est alors  $\{E, \mathcal{D}(E), (P_{N_1, \dots, N_r})_{N_1, \dots, N_r} \in \Pi\}$  où :

$$E = \{(n_i)_{1 \leq i \leq r} \mid 0 \leq n_i \leq n, \sum_{i=1}^r n_i = n\}$$

$$\Pi = \{N_1, \dots, N_r \mid 0 \leq N_i \leq N \text{ et } \sum_{i=1}^r N_i = N\}$$

$$P_{N_1, \dots, N_r}(n_1, n_2, \dots, n_r) = \frac{\binom{N_1}{n_1} \binom{N_2}{n_2} \dots \binom{N_r}{n_r}}{\binom{N}{n}}$$

Le loi  $P_{N_1, \dots, N_r}$  est une loi "hypergéométrique".

Remarque : Lorsque  $N$  est très grand, les premiers tirages ne modifient à peu près pas la composition de l'urne. Le tirage sans remise sera alors considéré comme un tirage avec remise (de loi plus simple).

## II - Echantillonnage dans une population infinie. Forme classique du modèle statistique.

En général, dans l'étude d'un phénomène expérimental, on ne suppose rien sur la population du point de vue mathématique. Mais concrètement les données sont très souvent considérées comme une réalisation de  $n$  v.a. indépendantes et de même loi  $F_\theta$ . Si ces variables sont réelles, on notera

$$(\mathbb{R}, \mathcal{E}, F_\theta)_{\theta \in \Theta}^{\otimes n} \quad \text{le modèle}$$

le modèle statistique ainsi obtenu. L'échantillon est donné  $S = (X_1, \dots, X_n)$ , et on dit aussi échantillon pour réalisation de l'échantillon  $(X_1(\omega), \dots, X_n(\omega))$ . Les  $X_i$  sont donc les applications coordonnées de ce modèle.

On verra à la fin du chapitre II sur l'estimation d'autres types de fabrication d'échantillons dans des problèmes à objet

bien spécifique, et pour certains types de population. (Sondages stratifiés).

$$\text{On note } \bar{X} = \frac{X_1 + \dots + X_n}{n}$$

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

Si les  $X_i$  sont réels, on note  $X_{(1)} \dots X_{(n)}$  les réarrangements par ordre croissant des  $X_i$ .

On notera  $N(m, \sigma^2)$  la loi normale de moyenne  $m$ , de variance  $\sigma^2$ ,  $N_r(m, K)$  la loi normale de moyenne  $m$  dans  $\mathbb{R}^r$ , de covariance  $K$ . Si  $L$  est une loi  $X \sim L$  signifie que la v.a.  $X$  suit la loi  $L$ .

Il sera fréquent de considérer des modèles images de  $(\mathbb{R}, \mathbb{D}, F_\theta)_{\theta \in \Theta}^{\otimes n}$  par des statistiques à 1 ou de dimension, et de raisonner uniquement sur ces modèles images.

Fréquemment dans les modèles paramétriques ou non paramétriques on regardera une famille de loi du type

$$F\left(\frac{x - \theta}{\sigma}\right)$$

s'appelle alors paramètre de position (centrage, translation et paramètre d'échelle).

$$\text{Ex. } \frac{1}{\sqrt{2\pi}} \exp \frac{1}{2} \frac{(x - \theta)^2}{\sigma} \frac{dx}{\sigma}, \quad \frac{1}{\pi} \frac{1}{1 + \left(\frac{x - \theta}{\sigma}\right)^2} \frac{dx}{\sigma}$$

## CHAPITRE II

### THEORIE ELEMENTAIRE DE L'ESTIMATION ET SONDAGES

Soit  $(X_1 \dots X_n) = X$ , un  $n$ -échantillon d'une v.a. de loi  $F$ . Sauf à la fin du chapitre, le modèle statistique est toujours  $(\mathbb{R}, \mathcal{B}, F)_{F \in \mathcal{F}}^{\otimes n}$ , où  $\mathcal{F}$  est un ensemble de lois de probabilités sur  $\mathbb{R}$  (nous verrons sur des exemples des v.a. à valeurs dans  $\mathbb{R}^K$ ,  $K \geq 1$ ).

Soit  $\theta$  un paramètre réel associé à  $F$ , par exemple : moyenne, variance, médiane, fonctions de ces quantités. Notant  $\Theta$  l'ensemble des valeurs de  $\theta$ , nous remarquerons qu'en général, l'application  $F \rightarrow \theta(F)$  n'est pas biunivoque. Faire une estimation, c'est au vu de l'échantillon  $X$  prendre une décision  $\hat{\theta}(X)$  au sujet de la valeur de  $\theta$ .

Définition 1. Un estimateur  $\hat{\theta}$  de  $\theta$  est une application  $\hat{\theta} : X \rightarrow \Theta$ .  $\hat{\theta} \circ X(\omega)$  s'appelle l'estimation de  $\theta$ .

La fonction de perte dans ce problème de décision, est la perte quadratique définie par

$$w(\theta, \hat{\theta}(X)) = |\hat{\theta}(X) - \theta|^2,$$

on suppose donc toujours  $E_F |\hat{\theta}(X)|^2 < \infty$ . L'estimation de  $\theta$  par  $\hat{\theta}$  comporte donc le risque :

$$R_{\hat{\theta}}(F) = E_F |\hat{\theta}(X) - \theta|^2.$$

Le risque peut aussi être considéré comme une mesure de la précision de  $\hat{\theta}$ .

Lorsqu'il existe une application biunivoque  $\mathcal{F} \longleftrightarrow \Theta$  (modèles paramétriques) et on a alors :

$$R_{\hat{\theta}}(\theta) = E_{\theta} |\hat{\theta}(X) - \theta|^2.$$

La règle générale de comparaison des règles de décision donne ici :

$\hat{\theta}_1$  est meilleur que  $\hat{\theta}_2$  si

$$R_{\hat{\theta}_1}^*(F) \leq R_{\hat{\theta}_2}^*(F) \text{ pour tout } F.$$

Exemple 1. Soit  $\theta$  la moyenne de  $F$ ,

$$\begin{aligned} \mathcal{F} &= \{F, F \text{ probabilité sur } \mathbb{R}\} \\ &= \pi(\mathbb{R}) \\ \bar{X} &= \frac{X_1 + \dots + X_n}{n} \text{ est un estimateur de } \theta, \end{aligned}$$

de risque  $E_F |\bar{X} - \theta|^2 = \sigma^2(F)$ .

Remarquons que l'on a  $E_F \bar{X} = \theta(F)$ . Cette propriété de centrage d'un estimateur nous amène à poser la définition suivante :

Définition 2. Un estimateur de  $\theta$  est sans biais si pour tout  $F \in \mathcal{F}$

$$E_F \hat{\theta} = \theta(F).$$

Exemple 2. Soit à estimer  $\sigma^2(F)$ .

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} \text{ est un estimateur de } \sigma^2. \text{ On a } E_F s^2 = \sigma^2(F).$$

$s^2$  est sans biais.

Exemple 3. Soit  $X_1 \dots X_n$  un  $n$ -échantillon de  $N(0, \sigma^2)$ ,  $\sigma^2 \in \mathbb{R}^+$ .

Soit  $T = \sum_{i=1}^n X_i^2$ ,  $U = \frac{T}{n}$  est un estimateur sans biais de  $\sigma^2$ ,  $E_\sigma U = \sigma^2$ .

On a  $E_\sigma X_i^4 = 3\sigma^4$  (calculs aisés) et donc

$$\begin{aligned} E_\sigma T^2 &= E_\sigma \left( \sum_{i=1}^n X_i^4 + \sum_{i \neq j} X_i^2 X_j^2 \right) \\ &= n(n+2)\sigma^4 \end{aligned}$$

et donc  $R_U(\sigma) = \frac{n+2}{n}\sigma^4 - \sigma^4 = \frac{2}{n}\sigma^4$ .

Posons  $U_c = cnU$ .

Des calculs élémentaires donnent

$$\begin{aligned} E_\sigma U_c &= cn\sigma^2 \\ E_\sigma (U_c - \sigma^2)^2 &= \sigma^4 [n(n+2)c^2 - 2nc + 1] \end{aligned}$$

qui est minimum pour  $c = \frac{1}{n+2}$  et vaut

$$\left[ 1 - \left( \frac{n}{n+2} \right) \right] \sigma^4 = \frac{2}{n+2} \sigma^4.$$

On voit sur cet exemple qu'un estimateur sans biais peut être moins bon qu'un estimateur avec biais.

L'intérêt essentiel d'un estimateur sans biais est le suivant :

si l'on répète des estimations indépendantes du même paramètre, à l'aide du même estimateur sans biais : c.à.d. que si :  $X^k = (X_1^k \dots X_n^k)$ ,  $k \in K$  est une suite d'échantillons indépendants, de même loi, si  $\hat{\theta}_k$  ( $k \in K$ ) est la suite des estimateurs sans biais de  $\theta$ , associés aux  $X^k$ , indépendants et de même loi, alors :

le risque associé à  $\frac{\hat{\theta}_1 + \hat{\theta}_2 + \dots + \hat{\theta}_k}{K}$  valant alors  $\frac{\text{var}_{\mathbb{F}}(\hat{\theta}_1)}{K}$ , ce qui est une bonne propriété que n'ont pas les estimateurs biaisés.

Dans le cas où  $n$  est petit ou moyen, la recherche d'estimateurs sans biais n'a pas toujours un sens. Un ensemble extrême est le suivant :

Exemple 4. (téléphoniste).

Un téléphoniste ne connaît pas son standard. Il écoute une heure et enregistre  $x$  appels. Le nombre d'appels suit une loi de Poisson de paramètre inconnu  $\theta$ . Pour aller voir son amie, il lui faut 2 heures. Quelle est la probabilité pour qu'il n'y ait pas d'appel pendant ces deux heures.

La probabilité pour qu'il n'y ait pas d'appel pendant ces deux heures, est égale à :

$$P_{\theta}(X_1 + X_2 = 0) = e^{-2\theta}.$$

On a donc à estimer  $e^{-2\theta} = f(\theta)$ .

Soit  $\hat{f}$  un estimateur sans biais de  $e^{-2\theta}$ . On doit avoir

$$E_{\theta}(d) = e^{-2\theta} \text{ pour tout } \theta \geq 0,$$

or

$$E_{\theta}(d) = \sum_{n=0}^{\infty} d(n) e^{-\theta} \frac{\theta^n}{n!} = \sum_{n=0}^{\infty} e^{-\theta} (-1)^n \frac{\theta^n}{n!} \text{ pour tout } \theta \geq 0,$$

donc nécessairement  $\hat{f}(n) = (-1)^n$ , ce qui ne présente aucun intérêt !

Certains modes sont dit non paramétriques. Essentiellement ce sont ceux où  $\mathcal{F}$  est très grand et en particulier ceux où il n'existe pas d'application biunivoque  $\mathcal{F} \longleftrightarrow \theta$  (il n'y a pas d'accord général sur cette définition). C'est le cas des exemples 1, 2 et de l'exemple 5 ci-dessous par opposition au cas paramétrique (ex. 3, 6, 7).

Exemple 5. Soit  $X_{(1)}, X_{(2)} \dots X_{(2n+1)}$  le réarrangement par ordre crois-



sant d'un  $(2n+1)$  échantillon  $X_1, \dots, X_{2n+1}$ . Supposons que  $\mathcal{F} = \{F$  probabilités diffuses sur  $\mathbb{R}\}$ . Soit  $\mu$  la médiane de  $F$ , ( $F(\mu) = 1/2$ ). Alors  $\hat{\mu} = X_{(n+1)}$  est un estimateur de  $\mu$  (nous reviendrons sur ce cas au chapitre VI).

Exemple 6. Soit une loi de Bernoulli de paramètre  $p$ .  $\bar{X}$  est un estimateur de  $p$  de risque  $p-p^2$ .

Exemple 7. Soit  $N(m, \sigma^2)$  une loi normale  $\theta = (m, \sigma^2)$ ,  $\hat{\theta} = (\bar{X}, s^2)$  est un estimateur de  $\theta$ , le risque n'est pas défini car  $\theta \in \mathbb{R}^2$ . On peut prendre pour fonction de perte  $\|\hat{\theta} - \theta\|^2$  pour une certaine norme euclidienne sur  $\mathbb{R}^2$ .

Supposons maintenant avoir un modèle paramétrique  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta}$ , tel que l'on ait une idée a priori sur la valeur de  $\theta$  traduite par une distribution de probabilité  $\mu$  sur  $\Theta$ , supposé doté d'une structure d'espace mesuré  $(\Theta, \mathcal{E})$ .  $\mu$  est appelée distribution a priori du paramètre  $\theta$ , et le cadre des problèmes que l'on va poser est dit bayésien. La probabilité  $\mu$  sera souvent choisie pour de simples raisons de commodité mathématique, que nous exposerons plus loin, ou bien dans les meilleurs cas, en tenant compte d'expériences précédentes.

Un estimateur, noté ici  $d$ , de  $f(\theta)$  sera dit estimateur de Bayes si le problème de théorie de la décision que l'on pose est le suivant, la perte est toujours  $|d(X) - f(\theta)|^2$ , si  $X$  est l'échantillon, mais le risque est le risque bayésien,  $\mu$ -moyenne du risque ordinaire, soit

$$\begin{aligned} R_\mu(d) &= \int_\Theta R_\theta(d) \mu(d\theta) \\ &= \int_\Theta E_\theta |d - f(\theta)|^2 \mu(d\theta). \end{aligned}$$

Un estimateur  $d$  est dit  $\mu$ -admissible, si  $R_\mu(d) \leq R_\mu(d')$ , pour tout autre estimateur  $d'$ .

Exemple 8.a. Soit  $\Theta = \{\theta_1, \dots, \theta_k\}$ ,  $\mu = (p_1, \dots, p_k)$ ,  $p_i > 0$  pour tout  $i = 1 \dots k$ .

Alors si  $d$  est  $\mu$ -admissible,  $d$  est admissible au sens ordinaire.

En effet  $R_{\mu}(d) = \sum_{i=1}^k p_i E_{\theta} (d - f(\theta_i))^2$ .

Soit  $d'$  un autre estimateur, par hypothèse  $R_{\mu}(d) \leq R_{\mu}(d')$ .

Donc il existe  $i_0$  tel que

$$E_{\theta_{i_0}} [d - f(\theta_{i_0})]^2 \leq E_{\theta_{i_0}} [d' - f(\theta_{i_0})]^2.$$

Deux cas seulement sont alors possibles : les estimateurs  $d$  et  $d'$  ne sont pas comparables, ou alors  $d$  est meilleur que  $d'$ .  $d$  est donc admissible au sens ordinaire

Exemple 8.b. Soit  $\theta = \mathbb{R}$ . On suppose que  $R_d(\theta) = E_{\theta} |d - \theta|^2$  est une fonction continue de  $\theta$  (c'est un cas très fréquent), pour tout estimateur.

Soit  $d_0$  un estimateur  $\mu$ -admissible. Alors  $d_0$  est admissible si  $\mu$  a pour support  $\mathbb{R}$  (le support de  $\mu$  est le plus petit fermé  $F$  tel que  $\mu(F^c) = 0$ ).

Supposons  $d_0$  non-admissible. Alors il existe  $\theta_0$  et  $d$  tels que :

$$E_{\theta_0} |d - \theta_0|^2 < E_{\theta_0} |d_0 - \theta_0|^2 \quad (1)$$

et  $E_{\theta} |d - \theta|^2 \leq E_{\theta} |d_0 - \theta|^2$ .

Comme les fonctions  $R_d(\theta)$  et  $R_{d_0}(\theta)$  sont continues, l'inégalité (1) a lieu pour tous les  $\theta$  d'un intervalle  $I$  qui est donc chargé par  $\mu$  puisque  $\mu$  est de support  $\mathbb{R}$ .

$$\int E_{\theta} |d_0 - \theta|^2 d\mu(\theta) > \int E_{\theta} |d - \theta|^2 d\mu(\theta),$$

$d_0$  n'est donc pas  $\mu$ -admissible, contraire à l'hypothèse.

On voit sur ces deux exemples que la méthode de Bayes est très intéressante ; bien entendu, pratiquement son caractère est très discuté, mais il faut noter qu'elle donne souvent, par des calculs simples, de très bons estimateurs.

L'interprétation suivante, quoique restrictive, a été souvent utilisée pour interpréter le point de vue bayésien, bien qu'elle repose sur des fondements expérimentaux difficiles, souvent, à justifier, comme nous l'avons vu.

On suppose  $\theta \in \mathbb{R}$ ,  $d\mu(\theta) = g(\theta) d\theta$ , et que le modèle statistique est :

$$(\mathbb{R}, \mathcal{B}, f_{\theta}(x) dx)_{\theta \in \Theta}^{\otimes n}.$$

On considère  $\theta$  comme une v.a. et plus précisément on considère que le couple  $(X, \theta)$  admet sur  $(\mathbb{R}^{n+1}, \mathcal{B}_{n+1})$  la densité  $g(\theta) f_{\theta}(x)$ .

$g(\theta)$  est la densité a priori de  $\theta$ . C'est la densité marginale de  $\theta$ , celle qu'on attribue à  $\theta$  avant toute expérience.

La densité a posteriori de  $\theta$  est la densité de  $\theta$  connaissant le résultat  $X$  de l'expérience, c.à.d. :

$$g_x(\theta) = \frac{f_{\theta}(x)g(\theta)}{\int_{\mathbb{R}} f_{\theta}(x)g(\theta)d\theta}.$$

La recherche d'un estimateur bayésien de  $k(\theta)$  a alors une interprétation simple : il s'agit de trouver une v.a.  $X$ -mesurable telle que :

$$\int_{\mathbb{R}} E_{\theta} |d - k(\theta)|^2 g(\theta) d\theta \quad \text{soit minimale c.à.d.}$$

$$\int_{\mathbb{R}^{n+1}} |d(x) - k(\theta)|^2 f_{\theta}(x)g(\theta) dx d\theta$$

soit minimale,  $k(\theta)$  étant donné.

Dans  $L^2(\mathbb{R}_n \times \bar{\mathbb{R}}, \mathcal{B}_n \otimes \bar{\mathcal{B}}, f_{\theta}(x) g(\theta) d\theta \times dx)$ , il s'agit de trouver une fonction  $X$ -mesurable, dont la distance à  $k(\theta)$  soit minimale.

Nous savons, d'après le C 4 de probabilités, qu'il existe une solution et une seule à ce problème,  $E^X[k(\theta)]$ . L'estimateur de Bayes existe, est unique et égal à :

$$d(x) = \frac{\int k(\theta) f_{\theta}(x) g(\theta) d\theta}{\int f_{\theta}(x) g(\theta) d\theta},$$

d'après la formule donnant l'espérance conditionnelle dans le cas d'une densité.

Exemple 8.c. On cherche à estimer le paramètre d'une loi uniforme

$\frac{1}{\theta} 1_{[0, \theta]}(x)$ , avec un 1-échantillon. On cherche un estimateur de Bayes associé à  $g(\theta) = \theta e^{-\theta}$ .

L'estimateur est donc, sur  $R^+$ ,

$$d(x) = \frac{\int_0^\theta e^{-\theta} 1_{[0,\theta]}(x) d\theta}{\int_x^\infty e^{-\theta} 1_{[0,\theta]}(x) d\theta} = x+1 = \frac{\int_x^\infty \theta e^{-\theta} d\theta}{\int_x^\infty e^{-\theta} dx} = \frac{x e^{-x} + e^{-x}}{e^{-x}}$$

De plus  $R_d(\theta) = \int_0^\theta d^2(x) \frac{dx}{\theta} - 2 \int_0^\theta d(x) + \theta^2$  est une fonction continue de  $\theta$ .

Exemple 8.d. Soit à estimer le paramètre  $\theta$  d'une loi de Bernoulli à l'aide d'un  $n$ -échantillon. Soit  $S$  le nombre de face (+),  $n-S$  le nombre de pile (0).

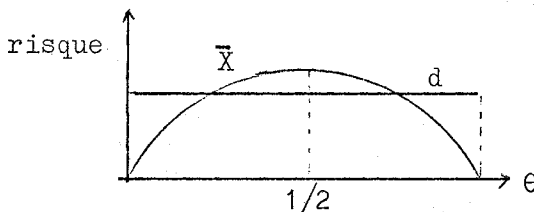
Considérons les estimateurs  $\frac{S+a}{n+b}$ ,  $E_\theta S = n\theta$ ,  $\text{var}_\theta S = n\theta(1-\theta)$

$$E\left|\frac{S+a}{n+b} - \theta\right|^2 = \frac{1}{(n+b)^2} E(S - n\theta + a - b\theta)^2 = \frac{n\theta(1-\theta) + (a-b\theta)^2}{(n+b)^2}$$

Choisissons  $a = \frac{\sqrt{n}}{2}$ ,  $b = \sqrt{n}$   $E\left(\frac{S + \frac{\sqrt{n}}{2}}{n + \sqrt{n}} - \theta\right)^2 = \frac{n\theta(1-\theta) + n\left(\frac{1}{2} - \theta\right)^2}{(n + \sqrt{n})^2} = \frac{1}{4(\sqrt{n} + 1)^2}$

L'estimateur  $d = \frac{S + \sqrt{n}/2}{n + \sqrt{n}}$  a un risque constant et égal à

$$\frac{1}{4n(1 + 1/\sqrt{n})^2}$$



L'estimateur  $\bar{X} = \frac{S}{n}$  a un risque  $\frac{\theta(1-\theta)}{n}$ .

$d$  est meilleur que  $\bar{X}$  sur l'intervalle défini par  $\theta(1-\theta) \leq \frac{1}{4(1 + (1/\sqrt{n}))^2}$

dont la taille est de l'ordre de  $1/\sqrt{n}$ . On peut voir que  $d$  est un estimateur de Bayes.

Considérons les densités définies pour  $a, b > 0$  par

$$f_{x,y}(\theta) = \frac{1}{B(x,y)} \theta^{x-1} (1-\theta)^{y-1}$$

où  $B(x,y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)} = \frac{(x-1)!(y-1)!}{(x+y-1)!}$ , si  $x, y \in \mathbb{N}$ .

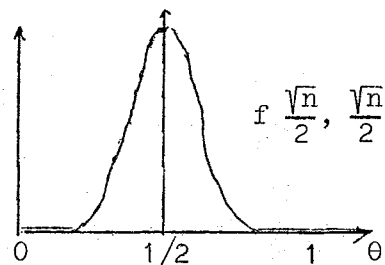
On a alors, comme estimateur de Bayes pour la densité  $f_{x,y}(\theta)$

$$d_{x,y} = \frac{\int_0^1 \theta^{S+x} (1-\theta)^{n-S+y-1} d\theta}{\int_0^1 \theta^{S+x-1} (1-\theta)^{n-S+y-1} d\theta}$$

$$\begin{aligned}
&= \frac{\Gamma(S+x+1)\Gamma(n-S+y)}{\Gamma(n+x+y+1)} \frac{\Gamma(n+x+y)}{\Gamma(S+x)\Gamma(n-S+y)} \\
&= \frac{S+x}{n+x+y} .
\end{aligned}$$

Pour  $x = \frac{\sqrt{n}}{2}$ ,  $y = \frac{\sqrt{n}}{2}$ , il vient

$$\begin{aligned}
d_{x,y} &= \frac{S+\sqrt{n}/2}{S+\sqrt{n}} \\
&= d
\end{aligned}$$



$d$  est donc de Bayes.

Comme sa mesure charge tout  $(0,1)$   $d$  est admissible, et comme le risque de  $d$  est constant, tout autre estimateur a un risque supérieur à cette constante, au moins en un point. Donc  $d$  est minimax. Il est intuitif que le maximum du risque est atteint pour  $1/2$  (maximum d'incertitude) et donc que, la loi  $f \frac{\sqrt{n}}{2}, \frac{\sqrt{n}}{2}$  qui "concentre" tout le risque de Bayes autour du point  $1/2$  (car très pointue) est très bonne pour les valeurs voisines de  $1/2$ . Si l'on prend comme mesure de Bayes la mesure de Lebesgue ( $x = y = 1$ ), l'estimateur est  $\frac{S+1}{n+2}$  dont le risque vaut  $\frac{np(1-p)+(1-p)^2}{(n+2)^2}$  qui vaut  $\frac{1}{4(n+2)}$  pour  $p = \frac{1}{2}$ .

Cet estimateur est également admissible et biaisé. On a

$$\frac{np(1-p)+(1-p)^2}{(n+2)^2} \leq \frac{p(1-p)}{n} .$$

Si  $n(1-p) \leq (2n+4)p$

soit  $\frac{n}{3n+4} \leq p$ , donc pour  $n$  grand, pour  $p \geq \frac{1}{3}$ . Cet estimateur est meilleur que  $\bar{X}$ , mais il est mauvais au voisinage de 0 (si  $S = 0$ , il donne  $\frac{1}{n+2}$ , si  $S = n$  il donne  $\frac{n+1}{n+2}$ , alors que  $\bar{X}$  donne les vraies valeurs. Nous n'avons pas montré que  $\bar{X}$  était admissible !.

On a donc un moyen de fabrication d'estimateurs minimax : on cherche une mesure de Bayes telle que l'estimateur associé soit de risque constant.

Cet exemple est très intéressant en ce qu'il conduit à un résultat peu intuitif. (On pourrait croire que  $\frac{S}{n}$  est le meilleur dans tous les cas possibles par raison de symétrie, notamment non minimax).

Le lecteur pourra cependant se persuader à l'aide de ces quelques exemples que le problème du choix entre plusieurs estimateurs admissibles n'est pas simple, pour  $n$  petit, même si le modèle est simple (lois de Bernoulli ou de Gauss), et que la pratique du statisticien est un élément important du choix.

Nous verrons plus tard une théorie plus élaborée de l'estimation en introduisant certaines classes d'estimateurs et aussi les problèmes relatifs aux échantillons de taille  $n$  grand et au comportement asymptotique.

E cet effet on se limitera souvent à rechercher un estimateur dans une classe particulière, par exemple

Définition 3. Un estimateur  $\hat{\theta}$  sans biais de  $\theta$  est dit de variance minimum si il est de risque minimum parmi les estimateurs sans biais ; autrement dit si  $E_{\theta} \hat{\theta} = \theta$  implique

$$E_{\theta} |\hat{\theta}_1 - \theta|^2 \geq E_{\theta} |\hat{\theta} - \theta|^2 .$$

Il arrive fréquemment comme nous le verrons plus tard qu'il existe des estimateurs de variance minimum.

Remarquons enfin, que l'on peut avoir des stratégies d'estimation de type aléatoire, mais elles sont peu utilisées.

2. ESTIMATEURS CONSISTANTS. LA METHODE DU MAXIMUM DE VRAISEMBLANCE.

a) Le problème est d'étudier des estimateurs ayant un bon comportement asymptotique.

Soit donc une suite  $(T_n)$  d'estimateurs de  $g(\theta)$  dans un modèle  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta}$ .

En fait, on considère une suite de modèles  $\otimes^n(\mathbb{R}, \mathcal{B}, F_\theta)_{\theta \in \Theta}$ , associés à des échantillons de taille croissante et pour chaque  $n$  on a un estimateur de  $g(\theta)$  mais pour pouvoir écrire la définition, il vaut mieux considérer tous les estimateurs définis sur le même espace.

Définition. Une suite  $T_n$  d'estimateurs est dite consistante si

$$\lim_{n \rightarrow \infty} T_n = g(\theta) \quad P_\theta \text{ - p.s.} \quad \text{pour tout } \theta \in \Theta.$$

Exemple 1. Soit  $\theta$  la moyenne de  $F_\theta$ . D'après la loi p.s. des grands nombres

$$\frac{X_1 + \dots + X_n}{n} = T_n \xrightarrow{P_\theta \text{ p.s.}} \theta$$

$(T_n)$  est donc une suite consistante d'estimateurs.

De même  $\frac{\sum_{i=1}^n (X_i - T_n)^2}{n} \xrightarrow{\text{p.s.}} \sigma^2(\theta)$  et donc  $S_n = \frac{\sum (X_i - T_n)^2}{n}$  est une suite d'estimateurs consistants de la variance.

Mais cette définition mathématique maniée sans précaution, a peu de sens concret. Si on pose

$$\begin{aligned} T'_n &= 0 & \text{si } n < 10^6 \\ T'_n &= T_n & \text{si } n \geq 10^6 \end{aligned}$$

alors  $T'_n$  est également une suite d'estimateurs consistants de la moyenne mais son impact concret est nul (car concrètement on a  $n$  grand mais  $n \ll 10^6$ ).

Il faudrait donc si l'on veut donner un intérêt autre que mathématique, préciser la définition de manière à pouvoir choisir raisonnablement entre plusieurs suites d'estimateurs consistants. Nous ne le ferons pas, nous contentant d'étudier des méthodes dont on démontre par ailleurs la qualité.

La méthode du maximum de vraisemblance.

Soit un modèle régulier sur  $(\mathbb{R}^n, \mathfrak{B}_n)$  défini par des probabilités du type  $(p_\theta dF)_{\theta \in \Theta}$ , où  $p_\theta$  sont les densités de la loi  $F_\theta$  par rapport à une mesure  $dF$ . On a donc

$$dP_\theta(X_1 \dots X_n) = p_\theta(x_1) \dots p_\theta(x_n) dF(X_1 \dots X_n).$$

Soit  $x = (x_1 \dots x_n)$  la valeur observée, estimer  $\theta$  par la méthode du maximum de vraisemblance, c'est calculer  $\theta$  tel que  $g_{n,\theta}(x) = p_\theta(x_1) \dots p_\theta(x_n)$  soit maximum. Donc  $\hat{\theta}$  est l'estimateur défini (s'il existe) par :

$$g_{n,\hat{\theta}_n}(x) = \sup_{\theta \in \Theta} g_{n,\theta}(x)$$

Dans la pratique, on utilise souvent cette méthode même pour des échantillons de taille  $n \neq 20$ , et concrètement les résultats ne sont pas mauvais. Pour chaque  $n$ , on a un estimateur  $\hat{\theta}_n$ . On démontre que sous des hypothèses très larges

$$\hat{\theta}_n \xrightarrow{\text{p.s.}} \theta \quad \text{si } \theta \text{ est une partie de } \mathbb{R}^d.$$

Remarquons que si  $\Theta$  est un ouvert de  $\mathbb{R}^d$  et si  $g_{n,\theta}$  est dérivable on a

$$\frac{d}{d\theta} g_{n,\hat{\theta}_n}(x) = 0$$

soit aussi

$$\frac{d}{d\theta} \log g_{n,\hat{\theta}_n}(x) = 0$$

ou

$$\sum_{j=1}^n \frac{d}{d\theta} \log p_{\hat{\theta}_n}(x_j) = 0$$

équation dite du maximum de vraisemblance, qui dans ce cas est une condition nécessaire pour  $\hat{\theta}_n$ .

Exemple 1. Famille exponentielle. On a

$$\log p_\theta(x) = \langle \theta, T(x) \rangle - \Phi(\theta) = \sum_{\ell} \theta_\ell T_\ell(x) - \Phi(\theta)$$

où  $\theta \in \mathbb{R}^k$ ,  $T: \mathbb{R} \rightarrow \mathbb{R}^k$ ,  $T \circ X$  est donc un vecteur aléatoire à  $K$  dimension, et  $\Phi$  une fonction normative, supposée convexe et définie dans un voisinage ouvert de  $0$ , et

$$\exp \Phi(\theta) = \int \exp \langle \theta, T(x) \rangle p_\theta(x) dF(x)$$

(ces densités seront étudiées en détail plus tard).



et donc l'équation de vraisemblance s'écrit

$$\sum_{i=1}^n T_{\ell}(x_i) = n \frac{\partial \Phi}{\partial \theta_{\ell}} \quad \ell = 1 \dots d,$$

puisque  $\theta$  est un ouvert d'intérieur non vide de  $\mathbb{R}^d$ . Soit encore

$$\begin{aligned} \Phi'(\theta) &= \frac{\sum_{i=1}^n T(x_i)}{n} \\ &= S_n(x) \end{aligned}$$

Comme  $\Phi$  est convexe par rapport à chaque variable (cf. définition de  $\Phi$ ) il existe en général une solution  $\hat{\theta}$  qui correspond effectivement à un maximum de  $g_{n,\theta}(x)$ .

Remarquons que comme

$$E_{\theta} T(x_i) = \Phi'(\theta)$$

on a  $\lim S_n(x) = \Phi'(\theta)$

et donc, puisque  $\Phi'$  est monotone on a  $\hat{\theta}_n(x) \rightarrow \theta$ , donc l'estimateur est consistant.

Donc, pour une famille exponentielle, s'il existe un estimateur du maximum de vraisemblance  $\hat{\theta}_n$  pour  $\theta$ , on a

$$\Phi'(\hat{\theta}_n) = \frac{\sum_{i=1}^n T(x_i)}{n}$$

et de plus cet estimateur est alors consistant.

Exemple 1.a. Soit une loi de Bernoulli  $B_{\theta}$ . Par rapport à la mesure de masse 1 sur chaque point de  $\{0,1\}^n$

$$g_{n,\theta}(x) = \theta^{\sum_{j=1}^n x_j} (1-\theta)^{n - \sum_{j=1}^n x_j}.$$

Soit  $T = \sum X_i$

$$\frac{d}{d\theta} g_{n,\theta}(x) = \frac{t}{\theta} - \frac{n-t}{1-\theta}$$

l'équation du maximum de vraisemblance donne donc

$$\hat{\theta} = \frac{t}{n},$$

et on vérifie sans difficulté avec la dérivée seconde qu'il s'agit bien d'un maximum. On a donc ici  $\hat{\theta} = \bar{x}$ .

Exemple 1.b. Si  $F_{\theta} = N(\mu, \sigma^2)$

$$\log g_{n,\theta}(x) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2$$

L'équation du maximum de vraisemblance donne ( $d = 2$ ,  $\theta$  et à 2 dimensions)

$$\frac{\partial \log g_{n,\theta}(x)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{j=1}^n (x_j - \mu) = 0$$

$$\frac{\partial \log g_{n,\theta}(x)}{\partial \sigma} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{1}{\sigma^3} \sum_{j=1}^n (x_j - \mu)^2 = 0$$

d'où

$$\mu = \frac{\sum_{j=1}^n x_j}{n} = \bar{x}$$

$$\sigma^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 .$$

On vérifie, en utilisant l'inégalité  $\log x \leq \frac{x^2-1}{2}$  pour  $x > 0$  que

$\hat{\theta} = (\bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2)$  donne bien un maximum de  $\log g_{n,\theta}(x)$ .

Exemple 1.c. Si  $X$  a une loi uniforme sur  $[a, b]$ ,  $\theta = \{a, b\}$

$$g_{n,\theta}(x) = \pi_{i=1}^n 1_{[a,b]}(x_i)$$

alors  $p_{\theta}$  n'est pas dérivable mais il est facile d'étudier directement le maximum qui vaut  $\frac{1}{(b-a)^n}$  si  $a = \inf_j x_j$ ,  $b = \sup_j x_j$ .

Nous allons maintenant étudier un problème qui concerne un autre type d'échantillonnage.

b. Estimation de la moyenne dans une population non homogène (sondage). Echantillonnage en classes (stratifié).

La population est divisée en  $M$  classes sociales. Soit  $N$  la population totale, on connaît  $N_j$  nombre d'individus de la classe  $j$ ,  $j = 1 \dots M$ . On veut estimer la moyenne d'un caractère défini sur la population ; par exemple le % du revenu dépensé par tête pour l'alimentation.

On peut imaginer que dans certaines classes la variance de cette variable est plus petite que dans d'autres. En particulier que certaines classes nombreuses sont particulièrement homogènes.

Dans ces conditions si on a une idée de ces variances, on a intérêt à faire un échantillonnage d'un type particulier.

La classe  $j$  contient  $N_j$  individus,  $N_j/N = p_j$ . Si on fait un échantillonnage de taille  $n_j$  dans la classe, la moyenne de l'échantillon  $(X_j^1, \dots, X_j^{n_j})$  est  $\hat{\mu}_j = \frac{X_j^1 + \dots + X_j^{n_j}}{n_j}$ , c'est un estimateur de la moyenne  $m$  du caractère  $X$  (% du revenu...) pour la classe  $j$ , c'est-à-dire de la moyenne de la loi conditionnelle de  $X/X \in j$ , soit  $\mu_j$ .

Supposons maintenant faire les 2 types d'échantillonnage suivants :

a) tirer  $n$  individus au sort dans la population entière ( $N$  supposé infini) soit  $X_1 \dots X_n$  les v.a. % du revenu ... des  $n$  individus tirés. On en déduit une estimation classique

$$\text{de } m, \bar{X} = \frac{X_1 + \dots + X_n}{n}.$$

b) tirons  $n_j$  individus dans la classe  $j$  (supposée de taille infinie), les  $n_j$  étant soumis à la contrainte

$$\sum_{j=1}^n n_j = n$$

A cet échantillon, on associe un nouvel estimateur de  $m$ ,

$$\hat{\mu} = \sum_{i=1}^M p_j \hat{\mu}_j$$

$$\text{On a } E \hat{\mu} = \sum_{i=1}^M p_j E \hat{\mu}_j$$

$$= \sum_{i=1}^M p_j \mu_j$$

Or on a, pour une v.a.  $X_\ell$  quelconque tirée de la population totale

$$\begin{aligned}
 EX_{\ell} &= \sum_{j=1}^n E (X_{\ell} \cap X_{\ell} \in j) \\
 &= \sum_{j=1}^n E (X_{\ell} / X_{\ell} \in j) P (X_{\ell} \in j) \\
 &= \sum_{j=1}^n p_j \mu_j
 \end{aligned}$$

Donc  $\hat{\mu}$  estimateur de  $m$ , constitué de la moyenne pondérée des  $\hat{\mu}_j$  est un estimateur sans biais.

c) Comparons maintenant les variances de ces estimateurs.

$\sigma^2 (\bar{X}) = \frac{1}{n} \sigma^2$  si  $\sigma^2$  est la variance de  $X_{\ell}$ , élément quelconque de la population.

$$\sigma^2 (\hat{\mu}) = \sum_{j=1}^M p_j^2 \sigma^2 (\hat{\mu}_j)$$

(variance de la somme de v.a. indépendantes)

$$= \sum_{j=1}^M p_j^2 \frac{\sigma_j^2}{n_j}$$

si  $\sigma_j^2$  est la variance d'un élément de la classe  $j$ .

Comparons ces variances. On a

$$\begin{aligned}
 E (X_{\ell} - m)^2 &= \sum_{j=1}^M E [(X_{\ell} - m)^2 \cap (X_{\ell} \in j)] \\
 &= \sum_{j=1}^M E [(X_{\ell} - m)^2 / X_{\ell} \in j] p_j \\
 &= \sum_{j=1}^M E (X_{\ell} - \mu_j + \mu_j - m / X_{\ell} \in j)^2 p_j \\
 &= \sum_{j=1}^M E [(X_{\ell}^j - \mu_j) + (\mu_j - m)]^2 p_j
 \end{aligned}$$

Si  $X_\ell^j$  est élément quelconque de la classe  $j$ ,  
 $E (X_\ell^j - \mu_j) = 0$ ,

$\mu_j - m$  est une constante, donc  $E (X_\ell^j - \mu_j) (\mu_j - m) = 0$   
 et par suite on a

$$\begin{aligned} \sigma^2 &= E (X_\ell - m)^2 \\ &= \sum_{j=1}^M p_j \sigma_j^2 + p_j (\mu_j - m)^2 \end{aligned}$$

Interprétation :

Le premier terme correspond à la partie de la variance  $\sigma^2$  due à la variabilité  $\sigma_j$  à l'intérieur du groupe  $j$ , le 2ème terme à la variabilité  $(\mu_j - m)^2$  entre le groupe  $j$  et la population générale (carré de la "distance" du groupe  $j$  à la population générale). On a donc :

$$\sigma^2 (\bar{X}) = \frac{1}{n} \left[ \sum_{j=1}^M p_j \sigma_j^2 + p_j (\mu_j - m)^2 \right]$$

On vérifiera aisément, si  $M = 2$  que

$$\frac{1}{n_1} p_1^2 \sigma_1^2 + \frac{1}{n_2} p_2^2 \sigma_2^2 \leq \frac{1}{n_1 + n_2} (p_1 \sigma_1^2 + p_2 \sigma_2^2)$$

implique que  $\hat{\mu}$  est meilleur que  $\bar{X}$ . Montrons que l'on peut toujours choisir les  $n_j$  pour qu'il en soit ainsi et déterminons le choix optimal des  $n_j$ . Remarquons d'abord que si l'on prend l'estimateur "naturel" qui consiste à prendre les  $n_j$  proportion-

nels aux  $p_j$  soit  $\frac{n_j}{n} \neq p_j$ ,  $n_j = [np_j]$  éventuellement corrigé d'une unité pour que la somme fasse 1, ( $[np_j]$  signifiant le plus grand entier  $\leq np_j$ ) on a

$$\begin{aligned} \sigma^2 (\hat{\mu}_1) &\neq \sum_{j=1}^M p_j^2 \frac{\sigma_j^2}{p_j} n \\ &\neq \sum_{j=1}^M p_j \sigma_j^2 \end{aligned}$$

et donc

$$\sigma^2(\bar{X}) \neq \sigma^2(\hat{\mu}_1) + \sum_{j=1}^M p_j (m - \nu_j)^2$$

ce qui montre déjà que  $\hat{\mu}_1$  est meilleur que  $\bar{X}$ .

Pour avoir le choix optimal des  $n_j$  il faut minimiser, par le meilleur choix possible des  $n_j$ ,

$$\sum_{j=1}^M p_j^2 \frac{\sigma_j^2}{n_j} \text{ sous la contrainte } \sum_{j=1}^M n_j = n.$$

Nous laissons faire le calcul pour  $M = 2$  par une méthode directe. Pour  $M$  quelconque employons la méthode de Lagrange.

Soit à minimiser

$$L(n_1 \dots n_M) = \sum_{j=1}^M p_j^2 \frac{\sigma_j^2}{n_j} + \lambda^2 \left( \sum_{j=1}^M n_j - n \right)$$

$$\text{Ecrivons } \left\{ \frac{\partial L}{\partial n_j} = 0 \right. \quad j = 1 \dots M$$

$$\text{Soit } \left\{ \frac{p_j^2 \sigma_j^2}{n_j^2} = \lambda^2 \right. \quad j = 1 \dots M$$

$$\text{ou } \left\{ \frac{p_j \sigma_j}{\lambda} = n_j \right. \quad j = 1 \dots M$$

et en additionnant,

$$n = \frac{1}{\lambda} \sum_{j=1}^M p_j \sigma_j$$

d'où le choix optimal des  $n_j$

$$n_j = \frac{p_j \sigma_j}{\sum_{j=1}^M p_j \sigma_j} n \quad j = 1 \dots M$$

(on vérifie aisément que cet extremum de  $L$  est un minimum car lorsque  $n_j \rightarrow 0$ ,  $L \rightarrow \infty$ ).

Le choix optimal (en nombres entiers) est donc de prendre

$$n_j = \left[ \frac{p_j \sigma_j n}{\sum_j p_j \sigma_j} \right]$$

c'est-à-dire de prendre les  $n_j$  proportionnels aux  $p_j \sigma_j$ , proportions de la classe  $j$  pondérée par sa variance. Par exemple, et c'est le cas pratique, on augmentera les proportions dans le sondage des individus d'une classe peu représentée dans la population mais à forte variance (i.e. classes aisées dans ce problème socio-économique) et on diminuera au  $X$  d'une classe nombreuse à revenu peu dispersé. On augmentera ainsi la précision du résultat obtenu sur la population entière.

### CHAPITRE III

#### INTRODUCTION A LA THEORIE DES TESTS ET DES REGIONS DE CONFIANCE

##### 1. Position du problème.

Soit  $X = (X_1 \dots X_n)$  un  $n$ -échantillon d'un modèle statistique  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta}$ .

Soit  $H_0$  une hypothèse faite sur le modèle et  $H_1$  l'alternative de  $H_0$ . Soit  $\theta_0 = \{\theta \in \Theta, "H_0 \text{ est vraie dans } (\Omega, \mathcal{A}, P_\theta)\}$ .

$\theta_1 = \{\theta \in \Theta, "H_1 \text{ est vraie dans } (\Omega, \mathcal{A}, P_\theta)\}$ . Donc  $\Theta = \theta_0 \cup \theta_1$ .

Faire un test de l'hypothèse  $H_0$  contre l'alternative  $H_1$  c'est, ayant étudié une réalisation  $X_1(\omega), \dots, X_n(\omega)$  de l'échantillon, énoncer la

conclusion "H<sub>0</sub> est vraie" (c.à.d.  $\theta \in \theta_0$ )

ou "H<sub>1</sub> est vraie" (c.à.d.  $\theta \in \theta_1$ ).

(Remarquons qu'en fait, dans beaucoup de problèmes concrets, on n'a pas  $\Theta = \theta_0 \cup \theta_1$ , on teste  $H_0$  contre  $H_1$ , c'est-à-dire on regarde qui de  $H_0$  ou de  $H_1$  a le "plus de chances" d'être vraie, mais  $\theta_0 \cup \theta_1$  est simplement une partie de  $\Theta$ ,  $H_0$  et  $H_1$  étant deux hypothèses s'excluant mais  $H_1$  n'est pas l'alternative à  $H_0$  et on étudie souvent plusieurs hypothèses  $H_1$ ).

##### Exemple 1. Test sur le paramètre d'une loi de Bernoulli.

Soit  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta} = (\{0,1\}, \mathcal{P}\{0,1\}, B_\theta)_{\theta \in ]0,1[}^{\otimes n}$  où  $B_\theta$  est la loi de Bernoulli

$$P(X_1 = 0) = 1 - \theta, P(X_1 = 1) = \theta.$$

a) On peut tester  $\theta = \frac{1}{2}$  contre  $\theta \neq \frac{1}{2}$

b)  $\theta > \frac{1}{2}$  contre  $\theta < \frac{1}{2}$

c)  $\theta = \frac{1}{2}$  contre  $\theta > \frac{1}{2}$  (voir la remarque).



Exemple 2. Test sur la variance d'une loi  $N(0, \sigma^2)$ .

$$(\Omega, \mathcal{A}, P_\sigma)_{\sigma^2 \in \mathbb{R}^+} = (\mathbb{R}, \mathcal{B}, N(0, \sigma^2))^{\otimes n}_{\sigma^2 \in \mathbb{R}^+}$$

- a) On peut tester  $\theta = 1$  contre  $\theta \neq 1$
- b)  $\theta \geq 2$  contre  $\theta < 2$
- c)  $\theta = 1$  contre  $\theta > 2$
- d)  $\theta \in ]2, 4[$  contre  $\theta \geq 4$ .

Exemple 3. Test (non paramétrique) sur la médiane.

$$(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta} = (\mathbb{R}, \mathcal{B}, F)_{\theta \in \pi_c(\mathbb{R})}^{\otimes n}$$

où  $\pi_c(\mathbb{R})$  est l'ensemble des lois de probabilités continues sur  $\mathbb{R}$ . On peut tester

- a)  $\mu \geq \frac{1}{2}$  contre  $\mu < \frac{1}{2}$
- b)  $\mu = \frac{1}{2}$  contre  $\mu \neq \frac{1}{2}$ .

Exemple 4. Test (non paramétrique) d'homogénéité.

Soient  $(X_1 \dots X_n)$  et  $(Y_1 \dots Y_m)$  deux échantillons indépendants définis respectivement  $(\mathbb{R}, \mathcal{B}, F)_{F \in \pi(\mathbb{R})}^{\otimes n}$  et  $(\mathbb{R}, \mathcal{B}, F_1)_{F_1 \in \pi(\mathbb{R})}^{\otimes m}$ . Le modèle global est  $(\mathbb{R}^{n+m}, \mathcal{B}_{n+m}, (\bigotimes_1^n F) \otimes (\bigotimes_1^m F_1))_{F \times F_1 \in \pi^2(\mathbb{R})}$   $F$  et  $F_1$  lois de probabilités sur  $\mathbb{R}$ ; donc  $\theta = \{F, F_1\} \in \pi^2(\mathbb{R})$ , ensemble des couples de lois de probabilité sur  $\mathbb{R}$ . Un test d'homogénéité est un test portant sur l'hypothèse : les deux échantillons proviennent de la même loi, soit  $H_0 : F = F_1$  contre  $H_1 : F \neq F_1$  (donc  $\theta_0$  est la diagonale de  $\pi^2(\mathbb{R})$ ).

Exemple 5. Test non paramétrique d'indépendance.

Avec les mêmes notations que dans le cas précédent, on veut tester si les échantillons  $(X_1 \dots X_n)$  et  $(Y_1 \dots Y_n)$  proviennent de v.a. indépendantes,  $X_k$  et  $Y_k$  étant mesurées simultanément. Le modèle est donc

$$(\mathbb{R}^{2n}, \mathcal{B}_{2n}, \bigotimes_1^n F)$$

On teste l'hypothèse

$$"F = F_1 \otimes F_2" \text{ contre } "F \text{ n'est pas du type } F_1 \otimes F_2",$$

où  $F_1$  et  $F_2 \in \pi(\mathbb{R})$ .

Exemple 6. Test du caractère strictement aléatoire.

Soit  $(X_1(\omega), \dots, X_n(\omega))$  une réalisation de  $n$  v.a.r. dont on ne sait pas si elles sont indépendantes. Problème : répondre à la question :

il existe une loi de probabilité  $F$  telle que  $(X_1 \dots X_n)$  puissent être considérées comme la réalisation d'un  $n$ -échantillon constitué de v.a. indépendantes) de là  $F$ .

Le modèle à prendre est donc :

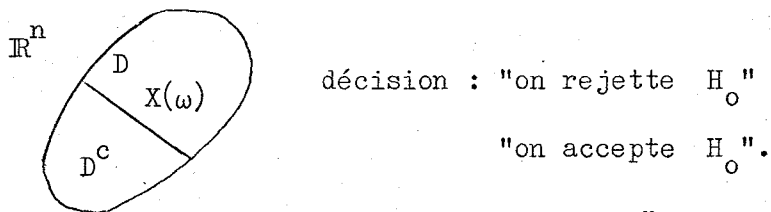
$(\mathbb{R}^n, \mathcal{B}_n, F)_{F \in \pi(\mathbb{R}^n)}$ , où  $\pi(\mathbb{R}^n)$  est l'ensemble des probabilités sur  $\mathbb{R}^n$ . On dispose donc d'un échantillon de taille 1 pour répondre à la question  $H_0 : F = \bigotimes_1^n G$ , où  $G$  est une loi de probabilité sur  $\mathbb{R}$ . Dans certains problèmes on doit préciser  $H_0 : F = \bigotimes_1^n \lambda$  où  $\lambda$  est la mesure de Lebesgue sur  $[0,1]$ . Par exemple, si on a fabriqué à l'aide d'un ordinateur  $n$  nombres et que l'on veuille vérifier que ses nombres peuvent être considérés comme  $n$  nombres "réels" tirés au sort suivant la loi uniforme sur  $[0,1]$  (cette spécification de  $F$  s'appelle test d'ajustement) et indépendamment.

On a l'habitude de classer les tests en 2 catégories. On parle de test paramétrique lorsque  $\Theta$  est un ensemble de  $\mathbb{R}$  ou de  $\mathbb{R}^k$  les exemples 1 et 2 sont des tests paramétriques. On parle de test non paramétrique dans les autres cas. Bien entendu, cette division n'est pas claire mathématiquement puisque un test paramétrique apparaîtrait comme un cas particulier de test non paramétrique et que souvent on peut, à l'aide d'une bijection "peu naturelle" identifier  $\Theta$  à un sous-ensemble de  $\mathbb{R}$  ou de  $\mathbb{R}^n$ .

Intuitivement les tests sont non-paramétriques si l'ensemble  $\Theta$  est une partie de  $\pi(\mathbb{R})$  ensemble des probabilités sur  $\mathbb{R}$  (ou  $\pi(\mathbb{R}^n)$ ) qui ne se représente pas "naturellement" comme un ensemble simple de  $\mathbb{R}$  ou  $\mathbb{R}^k$ . C'est le cas des tests 3, 4, 5, 6.

## 2. Forme de la règle de décision.

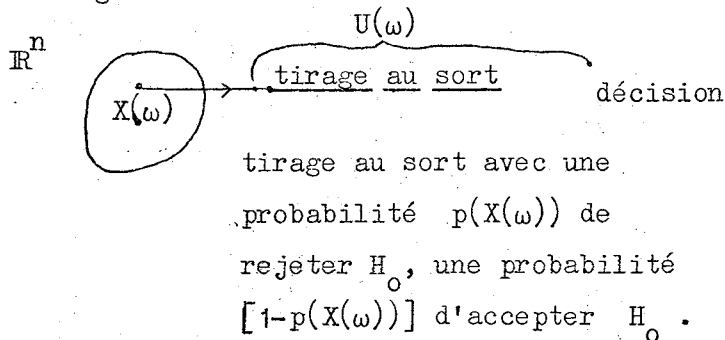
Ayant observé  $(X_1(\omega), \dots, X_n(\omega))$  on doit répondre " $H_0$  est vraie" ou " $H_0$  est fausse" statistiquement, c'est-à-dire avec un risque d'erreur dans la réponse qui va être étudiée plus bas. Ceci étant, les règles de décision utilisées seront dans les bons cas du type suivant :



Une règle décision est une partition  $(D, D^c)$  de  $\mathbb{R}^n$ . Si  $X(\omega) \in D$  on rejette  $H_0$ , si  $X(\omega) \in D^c$  on accepte  $H_0$ .

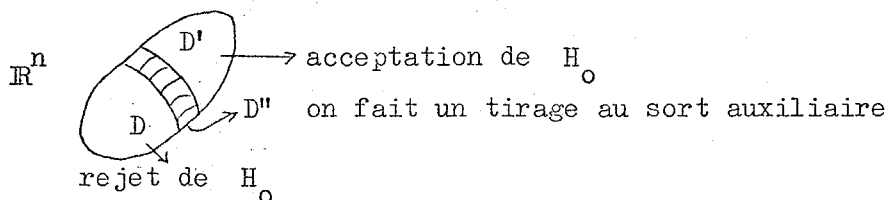
Cependant, on ne peut pas toujours se contenter de ce type de règle, car nous verrons que même dans certains cas simples, il existe un type de règles plus compliquées mais meilleures, comme nous l'avons vu au chapitre I.

Le schéma général sera le suivant :



A chaque valeur observée  $X(\omega)$ , on associe une probabilité  $p$  et on effectue un tirage au sort, (c'est-à-dire on réalise une v.a. de Bernoulli  $U$  de toute manière indépendante de  $X$  une fois choisie  $p$ ) et suivant la valeur de  $U$ , on choisira  $H_0$  ou on rejettera  $H_0$ .  $p$  s'appelle la fonction critique du test, c'est donc une fonction  $\mathbb{R}^n \rightarrow [0,1]$ .

Concrètement la situation sera très souvent la suivante



$D$  et  $D'$  sont en général des ouverts,  $D''$  leur frontière commune.

3. Pertes et risques dans un problème de test.

La fonction de perte que l'on prendra est la suivante : soit  $d$  la décision prise,  $d$  est à valeurs dans l'ensemble à 2 éléments  $(H_0, H_1)$ , la perte vaut 1 si la décision prise est mauvaise, 0 si la décision prise est la bonne.

Si  $(\Omega, \mathcal{A}, P_\theta)_{\theta \in \Theta}$  est le modèle, si  $H_0 = \{\theta \in \Theta_0\}$  et si le test est défini par une région de rejet  $D$  et une région d'acceptation  $D^c$ , on a une fonction de perte :

$$W : \Omega \times \Theta \rightarrow \{0, 1\} \text{ définie par}$$

$$w(\omega, \theta) = 1_D(X(\omega)) \text{ si } \theta \in \Theta_0$$

$$w(\omega, \theta) = 1_{D^c}(X(\omega)) \text{ si } \theta \in \Theta_1 .$$

Le risque est comme d'habitude l'espérance de la perte, soit :

$$R_D(\theta) = E_\theta w(\omega, \theta)$$

$$= P_\theta(D) \text{ si } \theta \in \Theta_0$$

$$= P_\theta(D^c) \text{ si } \theta \in \Theta_1 .$$

Remarque. Le choix de la fonction de perte très spéciale  $1_D$  sur  $\Theta_0$ ,  $1_{D^c}$  sur  $\Theta_1$ , peut paraître arbitraire. On pourrait pénaliser les erreurs grossières, c'est-à-dire lorsque  $\theta_0$  est par exemple un bout de  $\mathbb{R}^d$  dire qu'accepter à tort  $H_0$  est peu grave si il est voisin de  $\theta_0$ , et est très grave si  $\theta$  est loin de  $\theta_0$  et donc réintroduire une fonction de perte de type quadratique.

En fait, de même que dans la théorie de l'estimation, le fait de travailler avec une fonction de perte quadratique simplifiait l'exposé de la théorie sans perdre grand-chose sur la puissance d'application, le fait de se limiter ici à cette fonction de perte simplifie l'exposition et il n'y a pas grande difficulté à changer de point de vue. De plus le paragraphe suivant justifiera beaucoup plus profondément le choix de cette fonction de perte, qui aura alors un sens expérimental très précis.

La relation  $(D, D^c) > (D', D'^c)$  si et seulement si

$$R_D(\theta) \leq R_{D'}(\theta) \text{ pour tout } \theta \in \Theta$$

définit une relation d'ordre partiel sur les tests. On a donc une notion de test admissible et une notion de test optimal (cf. Chap. I) et on pourrait dérouler la théorie des tests exactement comme celle de l'estimation. Mais nous allons voir par la suite que ce point de vue n'est pas en général, celui que l'on choisit.

Etudions maintenant perte et risque dans le cas général d'un test avec tirage au sort auxiliaire  $U$

$$\begin{cases} w(\theta, \omega, \omega') = 1 & \text{si } U(\omega') = H_0 \text{ et si } \theta \notin \theta_0 \\ w(\theta, \omega, \omega') = 1 & \text{si } U(\omega') = H_1 \text{ et si } \theta \in \theta_0 \\ w(\theta, \omega, \omega') = 0 & \text{dans les autres cas (c.à.d. quand il n'y a pas d'erreur de faite).} \end{cases}$$

On a donc

$$\begin{aligned} R_d(\theta) &= E w(\theta, \omega, \omega') \\ &= E E^X(\theta, \omega, \omega') \\ &= \int_{\Omega} p(X(\omega)) dP_{\theta}(\omega) \end{aligned}$$

si  $\theta \in \theta_0$ , et

$$R_d(\theta) = \int [1 - p(X(\omega))] dP_{\theta}(\omega) \quad \text{si } \theta \notin \theta_0.$$

Exemple. Soit à tester  $\theta = \frac{1}{4}$  contre  $\theta = \frac{3}{4}$  pour une loi de Bernoulli à l'aide d'un  $(2n+1)$  échantillon. Le modèle est donc  $(\{0,1\}^{2n+1}, (\{0,1\}^{2n+1}, \otimes_{\theta} B_{\theta})_{\theta \in \{\frac{1}{4}, \frac{3}{4}\}}$  où  $B_{\theta}$  est la loi de Bernoulli de paramètre  $\theta$ .

Pour des raisons de symétrie (justifiées au chap. VII), on peut se limiter aux régions  $D$  de  $\{0,1\}^{2n+1}$  définies par  $X_1(\omega) + \dots + X_{2n+1}(\omega) \in E$  où  $E$  est une partie de  $[0, 2n+1] \subset \mathbb{N}$ . La décision sera donc

$$\begin{cases} \frac{1}{4} & \text{si } S \in E^c, \quad S = \sum_{i=1}^{2n+1} X_i \\ \frac{3}{4} & \text{si } S \in E. \end{cases}$$

On a donc

$$\begin{aligned} R\left(\frac{1}{4}\right) &= \sum_{k \in E} \binom{2n+1}{k} \left(\frac{3}{4}\right)^{2n+1-k} \left(\frac{1}{4}\right)^k \\ R\left(\frac{3}{4}\right) &= \sum_{k \in E^c} \binom{2n+1}{k} \left(\frac{1}{4}\right)^{2n+1-k} \left(\frac{3}{4}\right)^k \end{aligned}$$

Pour  $k < n$  on a  $\left(\frac{3}{4}\right)^{2n+1-2k} > \left(\frac{1}{4}\right)^{2n+1-2k}$  et donc pour minimiser le risque il faut prendre  $[0, n-1] \in E$

$$[n+1, 2n+1] \in E^c.$$

Mais que se passe-t-il pour la valeur  $n$  ?

Si on choisit par exemple  $\frac{3}{4}$ , on augmente le risque de  $\left(\frac{2n+1}{n}\right) \left(\frac{1}{4}\right)^{n+1} \left(\frac{3}{4}\right)^n$  dans le cas où  $\theta = \frac{1}{4}$  et on ne l'augmente pas dans le cas où  $\theta = \frac{3}{4}$ . Par contre si l'on fait un tirage au sort auxiliaire, par exemple avec  $p = \frac{1}{2}$ ,

on augmente les deux risques, mais on diminue le risque maximum, comme on le vérifiera par un calcul simple.

Les deux règles de décision ne sont donc pas comparables.

#### 4. Désymétrisation du risque et formulation classique de la théorie des tests.

Dans les faits le problème se pose souvent de la manière suivante. Soit à tester l'innocuité d'une pilule. Il est beaucoup plus grave de rejeter à tort l'hypothèse que la pilule a des effets dangereux que d'accepter à tort cette hypothèse. De manière générale, ceci étant particulièrement vrai dans le domaine des sciences expérimentales. La démarche est la suivante. Soit une expérience associée à un modèle à tester contre l'alternative (ou une autre hypothèse exclusive de  $H_0$ )  $H_1$ . Le risque de rejeter à tort  $H_0$  est considéré comme plus grave que celui d'accepter à tort  $H_0$ , la deuxième conclusion n'amenant que la poursuite du travail de recherche, alors que la première amène à l'application de ce travail (voir exemple ci-dessus). Dans beaucoup de problèmes, il est un type de risque maximal que l'on accepte de courir dans une certaine situation : un test sur la solidité d'un pont doit donner une certitude très forte, autrement dit on ne peut prendre de forts risques en rejetant l'hypothèse "le pont n'est pas solide".

Bien entendu, la notion de risque dépend beaucoup du domaine considéré. La quantification d'évènement négligeable est très variable. Par exemple chacun sait, dans notre monde, que la probabilité qu'une chaise soit dans un état non solide est de l'ordre de  $10^{-34}$  et ceci n'empêche pas les physiciens de s'asseoir.

Le risque admissible à courir pour une hypothèse donnée n'est pas le même pour la ruine de l'assureur, le joueur en bourse, le joueur de belote ou le biologiste qui cherche à prouver une certaine loi. Ce problème de quantification du risque acceptable relève donc du domaine concret, et de ce fait ce choix est quelquefois assez subjectif.

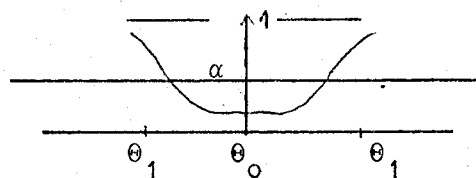
Définition. Soit le test d'une hypothèse  $H_0$  contre une hypothèse  $H_1$  effectué sur un modèle statistique  $(\mathbb{R}^n, \mathcal{B}_n, (P_\theta)_{\theta \in \Theta})$ ,  $H_0$  étant associée à  $\theta_0$ ,  $H_1$  à  $\theta_1$ . Le test est effectué à l'aide d'une région de rejet  $D$  et d'une région d'acceptation  $D^c$ .

On appelle risque de 1<sup>ère</sup> espèce  $P_\theta(D)$  pour  $\theta \in \theta_0$  : c'est la probabilité de rejeter à tort l'hypothèse  $H_0$ . On appelle risque de 2<sup>e</sup> espèce  $P_\theta(D^c)$  pour  $\theta \in \theta_1$ . C'est la probabilité d'accepter à tort  $H_0$  (nous verrons que l'on utilise plutôt  $P_\theta(D) = 1 - P_\theta(D^c)$  que le risque de 2<sup>e</sup> espèce).

Le niveau d'un test est  $\alpha = \sup_{\theta \in \theta_0} P_\theta(D)$ , c'est le maximum du risque de 1<sup>ère</sup> espèce.  $\alpha$  est le risque maximum que l'on accepte de courir dans le problème de test; on dit aussi que le test est significatif au niveau  $\alpha$ .

La puissance du test est la fonction

$$\begin{aligned} \theta &\rightarrow [0, 1] \\ \theta &\rightarrow P_\theta(D) \end{aligned}$$



(fin de la définition).

test de niveau  $\alpha$  avec fonction puissance.

On introduit alors une nouvelle relation d'ordre sur les tests.

Cette relation de pré-ordre ne compare que des tests de niveau inférieur ou égal à un niveau donné.

Définition. Soient  $D$  et  $D'$  deux tests (de région de rejet  $D$  et  $D'$  respectivement) de la même hypothèse  $H_0$ , contre la même alternative et de niveau  $\leq \alpha$ ,  $\alpha$  donné.  $D$  est dit plus puissant que  $D'$  si  $P_\theta(D) \geq P_\theta(D')$  pour tout  $\theta \in \theta_1$ . Il s'agit donc de la relation d'ordre partiel adapté au point de vue symétrisé.

(Etant sûr de courir un risque acceptable de 1<sup>ère</sup> espèce, l'expérimentation peut comparer deux tests).

Dans une classe donnée de tests de niveau  $\leq \alpha$ , un test maximal pour cet ordre est dit test uniformément plus puissant (U.M.P. en abrégé d'anglais) plutôt qu'optimal. Nous noterons cette propriété par U.P.P..

Nous pouvons étendre facilement toutes ces notions au cas de test avec tirage au sort auxiliaire.

Soit  $X$  l'échantillon et  $p(X)$  la probabilité de choisir  $D$  dans le tirage au sort auxiliaire après observation de  $X$ . Le niveau du test est :

$$\alpha = \sup_{\theta \in \theta_0} \int_{\mathbb{R}^n} p(X(\omega)) dP_{\theta}(\omega)$$

et la puissance la fonction

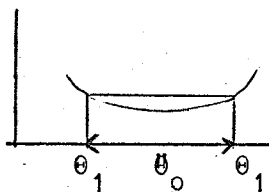
$$\theta \rightarrow \int_{\mathbb{R}^n} p(X(\omega)) dP_{\theta}(\omega).$$

Mathématiquement un des intérêts des tests avec tirage au sort auxiliaire sera de donner de bons théorèmes d'existence (et de construire) des tests à niveau  $\alpha$  donné (et non à niveau  $\leq \alpha$ ).

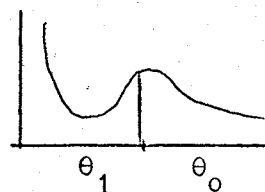
5. Test sans biais.

On dit qu'un test est sans biais si

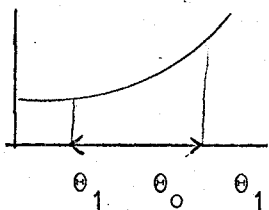
$$P_{\theta_0}(D) \leq P_{\theta_1}(D) \text{ pour tout } \theta_0 \in \theta_0, \theta_1 \in \theta_1$$



test sans biais



test biaisé



test biaisé

Si le test est à tirage au sort auxiliaire de fonction critique  $p$ , le test est sans biais si

$$E_{\theta_0} p \leq E_{\theta_1} p \text{ pour tout } \theta_0 \in \theta_0, \theta_1 \in \theta_1.$$

Il est pratiquement exclu de s'intéresser concrètement aux tests biaisés.

6. Test et intervalles de confiance. Dualité pour les tests ( $\theta = a$ ) contre ( $\theta \neq a$ ).

A) REGIONS DE CONFIANCE.

Il s'agit d'une troisième méthode de formulation des décisions en statistiques. Dans la théorie des tests la réponse est oui ou non. En estimation c'est  $\theta = \theta_0$ . En région de confiance c'est  $\theta \in S$  où  $S$  est un



sous-ensemble de  $\theta$ .

Soit  $\mathcal{X}$  les valeurs de l'échantillon (en général  $\mathcal{X} = \mathbb{R}^n$ ). On dit qu'une famille  $(S_x)_{x \in \mathcal{X}}$  de parties de  $\theta$  est une famille de région de confiance au niveau  $1-\alpha$  pour la fonction  $g(\theta)$  si l'on a

$$\inf_{\theta \in \theta} P_{\theta} \{x, g(\theta) \in S_x\} \geq 1-\alpha$$

autrement dit : à chaque  $x$ , on associe  $S_x$ , tel que pour  $P_{\theta}$ , c'est-à-dire si  $\theta$  est la valeur "vraie", la probabilité que  $g(\theta)$  soit dans la région  $S_x$  déterminée par la valeur  $x$  de l'échantillon soit au moins  $1-\alpha$ .

On peut alors définir un pré-ordre sur les régions de confiance. Une famille  $(S_x)$  est dite plus courte ou meilleure qu'une famille  $S'_x$  (au niveau  $1-\alpha$ ) si on a :

pour tout  $\theta_1, \theta_2, g(\theta_1) \neq g(\theta_2)$

$$P_{\theta_2} \{(x, g(\theta_1)) \in S(x)\} \leq P_{\theta_2} \{x, g(\theta_1) \in S'_x\}.$$

C'est une définition très naturelle qui signifie que  $S_x$  sépare mieux les valeurs distinctes  $g(\theta_1)$  et  $g(\theta_2)$  que  $S'_x$ . Cette définition dissymétrique est évidemment choisie en rapport avec la théorie des tests.

Enfin par analogie avec la théorie des tests, on dit que des régions de confiance sont sans biais si et seulement si pour tout  $\theta_1, \theta_2$  tels que  $g(\theta_1) \neq g(\theta_2)$  on a :

$$P_{\theta_2} \{x, g(\theta_1) \in S_x\} \leq 1-\alpha.$$

#### B) FAMILLES DE TESTS ET FAMILLES DE REGIONS DE CONFIANCE.

Considérons une famille de tests  $T_a$  de  $(\theta = a)$  contre  $(\theta \neq a)$  (famille de tests bilatères) et soit  $D_a$  leur région de rejet. Posons  $S_x = \{a, x \in D_a^c\}$ , (on suppose l'application  $(a, x) \rightarrow 1_{D_a}(x)$  mesurable). Inversement soit  $S_x$  une famille de régions de confiance. Définissons une famille de tests  $D_a$  par  $D_a = \{x, a \in S_x^c\}$ . Les deux définitions sont duales : l'espace  $\theta$  du paramètre et l'espace de l'échantillon jouent un rôle symétrique, puisque

$$S_x = \{a, x \in D_a^c\} \iff D_a = \{x, a \in S_x^c\}.$$

Donc il y a correspondance biunivoque entre famille de tests bilatères et familles de régions de confiance.

Si la famille de tests est au niveau  $\alpha$ , on a  $P_a(D_a) = \alpha$ . Le niveau de la famille de régions de confiance  $S_x$  est défini par

$$\inf_{\theta} P_a\{x, a \in S_x\} = \inf_{\theta} P_a(D_a^c) = 1 - \alpha.$$

Montrons enfin que les relations de pré-ordre sur les familles de régions de confiance et de tests sont les mêmes.

Dire que  $D_a$  est meilleure que  $D'_a$  pour tout  $a$  ( $D'_a$  étant une autre famille de tests au niveau  $\alpha$ ) c'est dire que  $P_b(D_a) > P_b(D'_a)$  pour tout  $b \neq a$  ( $D_a$  est plus puissant que  $D'_a$ ), donc

$$P_b\{x, a \in S_x^c\} \geq P_b\{x, a \in S'_x{}^c\}$$

ou

$$P_b\{x, a \in S_x\} \leq P_b\{x, a \in S'_x\}$$

qui est la relation de pré-ordre sur les familles  $S_x$  et  $S'_x$ . La réciproque est immédiate.

Exemple. Considérons la famille de tests de  $\{\theta = \theta_0\}$  contre  $\{\theta \neq \theta_0\}$ , où  $\theta$  est la moyenne d'une loi normale de variance 1, tests au niveau  $\alpha$ .

On verra qu'il existe alors un test sans biais de la forme

$$D_{\theta_0} = \{T < c_1 \cup T > c_2\},$$

où  $T = \sum_1^n X_i$ ,  $X_1 \dots X_n$  étant l'échantillon, avec de plus

$P_{\theta_0}(T < c_1 \cup T > c_2) = \alpha$ . Cette équation ne suffisant pas à déterminer  $c_1$  et  $c_2$ , choisissons les symétriques par rapport à  $n\theta_0$ , soit

$$c_1 = n\theta_0 - \sqrt{n} b(\theta_0), \quad c_2 = n\theta_0 + \sqrt{n} b(\theta_0).$$

La région de test s'écrit donc :

$$\left\{ \left| \sum_{i=1}^n X_i - n\theta_0 \right| > \sqrt{n} b(\theta_0) \right\} = D_{\theta_0},$$

où  $b$  est déterminé par

$$P_{\theta_0} \left\{ \left| \sum_{i=1}^n X_i - n\theta_0 \right| > \sqrt{n} b(\theta_0) \right\} = \alpha.$$

La loi de  $\sum X_i$  est  $N(n\theta_0, n)$ , d'où

$$\frac{1}{\sqrt{n} \sqrt{2\pi}} \int_{n\theta_0 - b\sqrt{n}}^{n\theta_0 + b\sqrt{n}} \exp \frac{(x - n\theta_0)^2}{2n} dx = 1 - \alpha$$

$$\text{soit } \frac{1}{\sqrt{2\pi}} \int_{-b}^b e^{-\frac{x^2}{2}} dx = 1-\alpha .$$

$b(\theta_0)$  est donc une constante qui se lit simplement sur la table de la loi normale. La région de test est un intervalle de l'espace des valeurs de l'échantillon.

Etudions maintenant la famille des régions de confiance associées à ces régions de tests.

On ne considèrera pas la valeur  $(x_1 \dots x_n)$  de l'échantillon mais seulement la valeur  $\sum_{i=1}^n x_i = t$  de la statistique  $\sum_{i=1}^n X_i$ .

$$\text{On a } S_x = \{\theta, x \in D_\theta^c\}$$

$$S_t = \{\theta, t \in D_\theta^c\}$$

(puisque  $D_\theta$  ne "porte" que sur  $T$ ).

$$\text{Donc } S_t = \{\theta, |t-n\theta| \leq \sqrt{n} b\}$$

$\{S_t\}_{t \in \mathbb{R}}$  est donc une famille d'intervalles de confiance. Si  $\alpha$  est donné,

$b$  est calculé par  $\frac{2}{\sqrt{2\pi}} \int_0^b e^{-\frac{x^2}{2}} dx = 1-\alpha$ , si  $t$  est la valeur observée, l'intervalle de confiance au niveau  $1-\alpha$  est donc  $[\frac{t}{n} - \frac{b}{\sqrt{n}}, \frac{t}{n} + \frac{b}{\sqrt{n}}]$ . On voit que la taille de l'intervalle décroît à la vitesse  $\frac{1}{\sqrt{n}}$ .

## CHAPITRE IV

### LA STATISTIQUE NON PARAMETRIQUE

#### I. Introduction.

Les statistiques descriptives utilisent souvent le terme de paramètre dans un contexte où il n'y a pas de modèle. Les paramètres comme la médiane servent à décrire l'échantillon ou les données. Il en est de même en statistique mathématique (ou inductive). Très souvent, on n'a pas d'idées précises sur le modèle (on : que ce soit l'expérimentateur ou le statisticien). Fixer arbitrairement un modèle de manière à pouvoir employer une méthode mathématique (test ou estimation) bien connue, relève souvent plus d'une démarche idéologique ("couverture" rationnelle d'une théorie mal justifiée) que d'une démarche scientifique. On désigne par statistique non-paramétrique deux sortes de techniques a priori assez différentes.

D'abord toutes celles qui opèrent sur un modèle très général et qui en particulier dans le cas des tests n'ont pas d'équivalent dans le domaine classique (ex : test du caractère "aléatoire" de l'échantillon). D'autre part, l'utilisation de statistiques dont la loi est indépendante du modèle sous des hypothèses très larges (free-statistics en anglais, les traductions françaises ne sont pas très sûres : libres ou universelles), ce sont essentiellement des statistiques fondées sur la notion d'ordre, comme nous le verrons ci-dessous. Les hypothèses sur le modèle sont très générales (loi continue ou quelquefois symétrique).

Parmi les avantages de la statistique non paramétrique, citons son caractère général (indépendant du modèle), le fait que très souvent, elle procède plus du classement des données que de leur mesure (on dit en américain

qu'elle est bien adaptée aux données sales c.à.d. entachées d'erreurs de mesures), ces caractères sont dits caractères de robustesse.

Par contre, il n'est pas possible de calculer des puissances en un sens simple. Disons cependant que de manière très générale, surtout pour les petits échantillons, quand on emploie une méthode non paramétrique sur un échantillon de source paramétrique (simulation) ou que l'on calcule mathématiquement la puissance d'un test non-paramétrique dans un modèle classique (gaussien, exponentiel etc...), les méthodes non-paramétriques, si elles sont moins bonnes que les paramétriques, leur sont cependant comparables (puissance souvent supérieure à la moitié de la paramétrique). Concluons en disant que trop souvent ce sont la tradition de modélisation abusive, le goût d'utilisation de "belles techniques mathématiques", le caractère ennuyeux des théories non-paramétriques qui font utiliser à tort des théories paramétriques. De ce point de vue, nous mettons en garde le lecteur. La place réservée dans ce cours aux techniques non paramétriques est insuffisante par rapport à leur importance, surtout dans le type d'application où la modélisation est très arbitraire (par exemple : les Sciences Humaines). Pour terminer, remarquons que nous avons classé à part le test du  $\chi^2$ , bien qu'il soit non-paramétrique. Ceci essentiellement parce qu'il n'est pas fondé sur les statistiques d'ordre et que sa justification n'est qu'asymptotique.

## 2. Statistiques d'ordre.

Soit  $X_1 \dots X_n$  un  $n$ -échantillon de v.a. de loi  $F$  sur  $\mathbb{R}$ . On notera encore par  $F$  la fonction de répartition des  $X$  et l'on supposera  $F$  continue. L'échantillon ordonné  $X_{(1)} \dots X_{(n)}$  est l'ensemble des  $X$  réordonné par ordre croissant, c.à.d.  $X_{(1)} = \inf_{1 \leq i \leq n} X_i$ ,  $X_{(2)} = 2^{\text{e}} X_i$  par ordre croissant,  $X_{(n)} =$  plus grand  $X_i$ .

Une statistique  $T$  est dite universelle ou libre si sa loi est indépen-

dante de la loi  $F$  de l'échantillon. L'exemple fondamental de statistique universelle  $T$  à valeurs dans  $\mathbb{R}^n$ , est celui de la suite  $U_1, \dots, U_n$ ; où  $U_1 = F(X_1), \dots, U_n = F(X_n)$ . Les  $U_i$  sont alors des v.a. indépendantes, de loi uniforme sur  $[0, 1]$  :

$$\begin{aligned} P(a < F(X_i) < b) &= P(F^{-1}(a) < X_i < F^{-1}(b)) \\ &= b-a. \end{aligned}$$

Donnons quelques interprétations évidentes des  $X_{(i)}$ .

- a)  $X_{(n)}$  intervient dans les problèmes de sécurité.
- b)  $X_{(1)}$  intervient dans les problèmes de sécurité, ou de fiabilité
- c) La médiane  $(X_{(\frac{n+1}{2})})$  si  $n$  est impair, n'importe quel nombre entre  $X_{(\frac{n}{2})}$  et  $X_{(\frac{n}{2}+1)}$  si  $n$  est pair, est une mesure de position ou de tendance moyenne (tendance dite centrale ou principale).

d)  $\frac{X_{(1)} + X_{(n)}}{2}$  est aussi une mesure de la tendance centrale.

e) La plage  $X_{(n)} - X_{(1)}$  est une mesure de dispersion.

f) La transformation  $(X_1 \dots X_n) \rightarrow (X_{(1)} \dots X_{(n)}) = Y$  n'est évidemment pas biunivoque. Elle l'est sur chaque  $\mathbb{R}_\sigma^n \subset \mathbb{R}^n$  où  $Y = \sigma X$ , (c.à.d.  $(X_{(1)} \dots X_{(n)}) = (X_{\sigma(1)} \dots X_{\sigma(n)})$ ),  $\sigma$  permutation donnée. Dans ce domaine le jacobien de  $Y = \sigma X$  est 1. On en déduit aisément que la densité de  $(X_{(1)} \dots X_{(n)})$  est  $n! \prod_{\substack{i=1 \\ y_1 < y_2 < \dots < y_n}}^n f_X(y_i)$ ;  $f$  densité de  $X$  (puis-que il y a  $n!$  permutations  $\sigma$ ).

Les lois des  $X_{(r)}$  (marginales) sont immédiates à calculer, par un argument combinatoire, par exemple si  $F$  a une densité  $f$  :

$$f_r(y_r) = \frac{n!}{(r-1)!(n-r)!} F^{r-1}(y_r) [1-F(y_r)]^{n-r} f(y_r).$$

#### Exercices.

- a) Montrer que la densité de  $X_{(r)}, X_{(s)}$  est donnée par

$$f_{r,s}(x,y) = \frac{n!}{(r-1)!(s-r-1)!(n-s)!} \\ \times [F(x)^{r-1}(F(y)-F(x))^{s-r-1}(1-F(y))^{n-s}] \\ \times f(x)f(y), \quad x < y, \quad r < s.$$

b) Trouver la densité de la médiane  $\mu$  :

$$\text{si } n = 2m+1 \text{ on a } f_{\mu}(x) = f_m(x).$$

$$\text{si } n = 2m \text{ on prendra } \mu = \frac{X_{(m)} + X_{(n+1)}}{2} \text{ et ensuite}$$

partir de  $f_{m,m+1}(x,y)$  et faire le changement de variable  $u = \frac{x+y}{2}$ ,  $v = y$

$$\text{il vient } f_{\mu}(u) = \frac{(2m)!2}{(m-1)!^2} \int_u^{\infty} F(2u-v)^{m-1} \\ \times (1-F(v))^{m-1} f(2u-v)f(v)dv.$$

c) Calculer la densité de la plage (range)  $X_{(n)} - X_{(1)}$ .

$$\text{d) Montrer que } E|U_r|^k = \frac{n!(r+k-1)!}{(n+k)!(r-1)!}$$

$$\text{cov}(U_r, U_s) = \frac{r(n-s+1)}{(n+1)^2(n+2)} \quad r < s$$

$$\text{corrélation } (U_{(1)}, U_{(n)}) = \frac{1}{n}.$$

Citons un résultat de convergence en loi dû à Smirnov, trop compliqué pour être démontré ici mais qui donne une idée des approximations utilisables lorsque  $n$  est grand.

Théorème. Supposons que  $n \rightarrow \infty$ ,  $\frac{r(n)}{n} = p$  et soit  $\pi$  défini par  $F(\pi) = p$ . Alors

$$\sqrt{\frac{n}{p(1-p)}} f(\pi) (X_{(r)} - \pi) \xrightarrow[\text{loi}]{} N(0,1), \text{ où } f(\pi) \text{ est une constante.}$$

### 3. Runs.

Soit un échantillon tiré d'une population à 2 caractères A et B. On note par un mot la réalisation de cet échantillon de taille  $n$ , par exemple,  $X_1(\omega) \dots X_n(\omega) = A A B A B B B \dots A B A$ .

Un run est une suite de mêmes symboles

$$A A B B B A B A = (AA)(BBB)(A)(B)(A)$$

intervenant successivement dans le mot et de longueur maximale parmi les

suites ayant cette propriété (précédée et suivie soit d'un symbole différent, soit de rien). Nombres de runs et longueur des runs sont des statistiques exploitables, notamment pour tester le caractère aléatoire d'un échantillon donné, c'est-à-dire pour vérifier que l'échantillon peut être considéré comme tiré au sort dans une population infinie, chacun des individus représentant la réalisation d'une v.a.  $X_i$  dans une suite  $X_1 \dots X_n$  de v.a. indépendantes et équidistribuées. L'hypothèse alternative n'étant souvent pas détaillée, on peut perfectionner la technique des runs dans le cas des caractères quantitatifs en considérant les runs croissants et décroissants (cf. plus loin).

a) Nombre total de runs

$n$	taille échantillon	
$n_1$	nombre de A	} fixés
$n_2$	nombre de B	
$r_1$	nombre de A-runs	
$r_2$	nombre de B-runs	
$R = r_1 + r_2$		



Problème : trouver la distribution de  $R$  lorsque l'hypothèse  $H_0$  :

"on a tiré au hasard et indépendamment  $n$  individus dans une population infinie" est vérifiée.

Lemme 1. Le nombre de configurations distinctes de  $n$  objets indiscernables dans  $r$  boîtes discernables est  $\binom{n-1}{r-1}$ .

Lemme 2. Soient  $R_1, R_2$  les v.a. nombres de A-runs et B-runs, la loi du couple  $R_1, R_2$  est donnée par

$$P(R_1 = r_1 \cap R_2 = r_2) = \frac{c \binom{n_1-1}{r_1-1} \binom{n_2-1}{r_2-1}}{\binom{n_1+n_2}{n_1}}$$

lorsque  $r_1 = 1, 2, \dots, n_1$ ,  $r_2 = 1, 2, \dots, n_2$

$$r_1 = r_2, \text{ ou } r_1 = r_2 \pm 1$$



$$C = 2 \text{ si } r_1 = r_2, \quad C = 1 \text{ si } r_1 = r_2 \pm 1.$$

Lemme 3.

$$P(R_1 = r_1) = \frac{\binom{n_1-1}{r_1-1} \binom{n_2+1}{r_1}}{\binom{n}{n_1}}$$

$$P(R_1 + R_2 = r) = \frac{2 \binom{n_1-1}{r/2-1} \binom{n_2-1}{r/2-1}}{\binom{n}{n_1}}, \quad r \text{ pair}$$

$$= \frac{\binom{n_1-1}{(r-1)/2} \binom{n_2-1}{(r-3)/2} + \binom{n_1-1}{r-3/2} \binom{n_2-1}{r-1/2}}{\binom{n}{n_1}} \quad \text{si } r \text{ impair}$$

et  $2 \leq r \leq n$ .

Ex :  $n_1 = 5, n_2 = 4$

$$P(R = 9) = 0.008 \qquad P(R = 8) = 0.063$$

$$P(R = 2) = 0.0016 \qquad P(R = 3) = 0.056.$$

Si on rejette pour  $R \in \{2, 9\}$ , on trouve comme niveau significatif 0.024  
 $R \in \{2, 3, 8, 9\}$ ,  $\alpha = 0.143$ . (cf. Tables).

Etude asymptotique.

On a, si  $\lambda = \frac{n_1}{n}$ , et si  $n \rightarrow \infty$  avec  $\lambda$  constant,

$$\lim E \frac{R}{n} = 2\lambda(1-\lambda)$$

$$\lim \text{Var} \frac{R}{\sqrt{n}} = 4\lambda^2(1-\lambda)^2$$

et un calcul pénible montre que

$$\frac{R - 2 \frac{\lambda(1-\lambda)}{\lambda(1-\lambda)}}{2\sqrt{n} \lambda(1-\lambda)} \rightarrow N(0, 1).$$

Remarques. 1. On ne sait évidemment rien sur la puissance, l'alternative à  $H_0$  étant trop large !

2. Si on rejette  $H_0$  on a intérêt à regarder s'il n'y a pas une tendance unilatérale par exemple "les A plus grands que les B" et à faire alors un test unilatéral car cette tendance indique une tendance à un nombre anormalement petit de runs (les grandes valeurs étant sans signification).

3. On peut aussi remplacer ce test par un test fondé sur la distribution du run le plus long (un run trop long est suspect, un plus long

(un run trop long est suspect, un plus long run trop court peut indiquer une périodicité).

4. Dans le cas où les v.a. observées sont quantitatives c.à.d. réelles, on utilise souvent la méthode dit du run up and down. Le run démarre chaque fois que la v.a. observée est plus grande que la précédente (run up ou croissant) si la suite précédente était du type down (décroissante) et réciproquement

Ex. : 7 8 2 3 4 5 s'écrit

+ - + + +  $\Rightarrow$  3 runs (1,+) (1,-) (3,+).

Des calculs compliqués donnent la distribution du nombre de runs up and down sous  $H_0$  et le comportement asymptotique,  $\frac{R-(2n-1)/3}{[(16n-29)/90]^{1/2}} \sim N(0,1)$ .

Le test de  $H_0$  a pour région de rejet les trop petites valeurs de  $R$ .

#### 4. Tests d'ajustement.

L'hypothèse  $H_0$  dans un tel test est la suivante, l'échantillon est la réalisation de  $n$  v.a. indépendantes équidistribuées, de loi  $F_0$  fixée. On ne discute ni de l'indépendance ni de l'équidistribution, c.à.d. que l'on admet que l'échantillon est du type "classique". On discute  $F_0$

$$H_0 : F = F_0$$

$$H_1 : F \neq F_0 .$$

Remarque importante : très souvent le test d'ajustement consiste à prendre pour  $F_0$  une loi, par exemple, gaussienne dont on a estimé les paramètres  $m$  et  $\sigma^2$ . Si l'estimation s'est faite à partir de l'échantillon, les méthodes de tests d'ajustement à  $F(m, \sigma^2)$  exposées ci-dessous ne sont plus valables. On doit utiliser des résultats plus élaborés.

Remarque. Comme nous l'avons déjà dit, le test du  $\chi^2$  est reporté à un chapitre ultérieur, c'est typiquement un test d'ajustement.

a) Distribution empirique.

$$\begin{aligned} \text{On pose } \hat{F}_n(x) &= \frac{k}{n} \text{ si } X_{(k)} \leq x < X_{(k+1)} \\ &= 0 \text{ si } x < X_{(1)} \\ &= 1 \text{ si } x > X_{(n)} \end{aligned}$$

(où  $X_{(j)}$  est l'échantillon réordonné).

On a donc

$$P(\hat{F}_n(x) = \frac{k}{n}) = \binom{n}{j} (F(x))^k (1-F(x))^{n-k}$$

pour  $k = 0 \dots n$ , si  $F$  est la fonction (inconnue) de répartition de  $X$ .

En particulier  $E \hat{F}_n(x) = F(x)$

$$\text{Var } \hat{F}_n(x) = \frac{F(x)(1-F(x))}{n}.$$

La loi large des grands nombres montre que  $\hat{F}_n(x) \rightarrow F(x)$  en probabilité.

Un résultat beaucoup plus fort est le suivant :

Théorème de Glivenko-Cantelli.

$\hat{F}_n(x) \rightarrow F(x)$  uniformément p.s. c'est-à-dire

$$\forall \varepsilon > 0, P(\sup |\hat{F}_n(x) - F(x)| > \varepsilon) \rightarrow 0.$$

Nous laisserons la démonstration de côté. Par contre une application immédiate du théorème de limite centrale donne

$$\sqrt{n} \frac{(\hat{F}_n(x) - F(x))}{\sqrt{F(x)(1-F(x))}} \rightarrow N(0,1).$$

b) Les tests de Kolmogorov-Smirnov.

On pose, pour  $F$  continue,

$$\begin{aligned} D_n &= \sup_x |\hat{F}_n(x) - F(x)| \\ D_n^+ &= \sup_x \hat{F}_n(x) - F(x) \\ D_n^- &= \sup_x F(x) - \hat{F}_n(x); \\ D_n &= \max_{0 \leq i \leq n} \sup_{X(i) \leq x < X_{i+1}} \left| \frac{i}{n} - F(x) \right| \end{aligned}$$

en posant  $X_{(0)} = -\infty$ ,  $X_{(n+1)} = \infty$

$$\begin{aligned} \text{d'où } D_n &= \max_i \left| \frac{i}{n} - \inf_{X(i) \leq x < X_{i+1}} F(x) \right| \\ &= \max_{1 \leq i \leq n} \left| \frac{i}{n} - F(X_{(i)}) \right| \\ &= \max_{1 \leq i \leq n} \left| \frac{i}{n} - U_{(i)} \right|. \end{aligned}$$

Donc la loi de  $D_n$  est indépendante de  $F$ . De même pour  $D_n^+$  et  $D_n^-$ . Le

calcul de  $D_n$  est pénible. On a

$$P(D_n < \frac{1}{2n} + v) = \begin{cases} 0 & \text{si } v \leq 0 \\ \int_{1/2n-v}^{1/2n+v} \dots \int_{2n-1/2n-v}^{2n-1/2n+v} f(u_1 \dots u_n) du_1 \dots du_n & \text{si } 0 < v < \frac{2n-1}{2n} \\ 1 & \text{si } v \geq \frac{2n-1}{2n} \end{cases}$$

où  $f(u_1 \dots u_n) = n!$  si  $0 < u_1 < \dots < u_n < 1$   
 $= 0$  sinon

est la densité de  $U_{(1)} \dots U_{(n)}$ .

Démonstration :  $0 < \frac{1}{2n} + v < 1$  ou  $-\frac{1}{2n} < v < \frac{2n-1}{2n}$  car  $0 \leq D_n \leq 1$ .

$$P(D_n < \frac{1}{2n} + v) = P\{\sup_x |\hat{F}_n(x) - x| < \frac{1}{2n} + v\}$$

$$= P\left(\bigcap_{i=0}^n \left[ \left| \frac{i}{n} - x \right| < \frac{1}{2n} + v, U_{(i)} < x < U_{(i+1)} \right]\right).$$

Soit  $A_i = \left\{ \frac{2i-1}{2n} - v \leq U_{(i)} < U_{(i+1)} < \frac{2i+1}{2n} + v \right\}$ .

$$P(D_n < \frac{1}{2n} + v) = P\left(\bigcap_{i=0}^n A_i\right)$$

$$= P\left(\bigcap_{i=0}^{n-1} A_i \cap A_{i+1}\right)$$

$$= P\left[\bigcap_{i=0}^{n-1} \left(\frac{2i+1}{2n} - v \leq X_{i+1} < \frac{2i+1}{2n} + v\right)\right]$$

d'où le résultat.

Des tables donnent  $P(D_n > D_{n,\alpha}) = \alpha$ , on en déduit un test d'ajustement.

Théorème limite de Kolmogorov et Smirnov.

Si  $F$  est continue,

$$\lim_{n \rightarrow \infty} P(D_n < \frac{z}{\sqrt{n}}) = 1 - 2 \sum_{k=1}^{\infty} (-1)^{k-1} e^{-2k^2 z^2}.$$

Cette fonction est bien tabulée. L'approximation est très bonne pour  $n \geq 35$   
 (de l'ordre de 0,01).

Théorème. On a, pour  $F$  continue

$$P(D_n^+ < t) = \begin{cases} \int_{1-t}^1 \int_{\frac{n-1}{n}-t}^u \dots \int_{\frac{2}{n}-t}^u \int_{1/n-t}^u f(u) du_1 \dots du_n & \text{pour } 0 < t < 1 \\ 1 & \text{si } t \geq 1. \end{cases}$$

Démonstration.  $P(D_n^+ < t) = P\left[\bigcap_{1 \leq i \leq n} \left(\frac{i}{n} - X_{(i)} < t\right)\right]$   
 $= P\left[\bigcap_{i=1}^n (X_{(i)} > \frac{i}{n} - t)\right]$  C.Q.F.D.

Propriété limite :

$$\lim_{n \rightarrow \infty} P(D_n^+ < \frac{z}{\sqrt{n}}) = 1 - e^{-2z^2}$$

et  $4nD_n^{+2} \xrightarrow{\text{loi}} \chi^2(2)$  (cf. chapitre sur le  $\chi^2$ ).

### 5. Statistiques fondées sur les rangs des statistiques d'ordre.

#### A) Définition des rangs.

Soit  $\sigma(\omega)$  la permutation définie par

$$X_{(j)}(\omega) = X_{\sigma(\omega)(k)}(\omega)$$

pour tout  $j$ ,  $j = 1 \dots n$ .

$j$  s'appelle le rang de  $X_k$  dans l'échantillon  $X_1 \dots X_n$ . On notera donc  $\sigma(k)$  ce rang, qui est une v.a. à valeurs dans  $\{1 \dots n\}$ , on peut encore la définir par :

$$\sigma(k) = \sum_{i \neq k} 1_{(X_k > X_i)} + 1.$$

La distribution de  $\sigma(k)$  est uniforme sur  $\{1 \dots n\}$ . On a :

$$\begin{aligned} P[(X_k < x) \cap \sigma(k) = j] &= \frac{1}{n} P(X_k < x / \sigma(k) = j) \\ &= \frac{1}{n} P(X_{(j)} < x). \end{aligned}$$

Supposons que  $X$  ait une loi continue  $F$ . On a :

$$\begin{aligned} E X_k \sigma(X_k) &= \sum_{j=1}^n \frac{j E X_{(j)}}{n} \\ &= \sum_{j=1}^n \frac{jn!}{(j-1)!(n-j)!} \int_{\mathbb{R}} x F(x)^{j-1} (1-F(x))^{n-j} dF(x), \end{aligned}$$

d'après la forme vue plus haut de la loi de  $X_{(j)}$ .

Des calculs élémentaires donnent :

$$\begin{aligned} \sum_{j=1}^n j E X_{(j)} &= n(n-1) E X_j F(X_j) + n E X_j \\ E \sigma(k) &= \frac{n+1}{2}, \quad \text{var } \sigma(k) = \frac{n^2-1}{12} \end{aligned}$$

et  $\rho(X_k, \sigma(k)) = \sqrt{\frac{12(N-1)}{N+1}} \frac{E X F(X) - 1/2 E X}{\sqrt{\text{Var } X}}$ .

#### Remarques.

1. Cette corrélation est assez variable et montre que les rangs ne

représentent pas très fidèlement l'échantillon, mais son sens n'est pas très clair. Nous y reviendrons.

2. On dit qu'il y a des noeuds lorsque l'on observe plusieurs valeurs égales des variables. Ceci est fréquent lorsque l'on utilise des rangs pour des lois  $F$  discrètes, ce que l'on fait couramment. Pour traiter les noeuds on peut :

a) faire un tirage au sort pour déterminer le rang entre valeurs égales. Cette méthode a l'avantage d'étendre la validité de la théorie aux lois discrètes.

b) Affecter à toutes les variables concernées le rang moyen. Ceci fausse les lois des statistiques concernées (par exemple : cela diminue la variance de  $\sigma(k)$ ) ; il existe d'autres procédures.

#### B) Tests de la médiane.

La médiane joue, en non paramétrique, le rôle que joue la moyenne pour certaines lois paramétriques comme la gaussienne (où la moyenne est la médiane) et d'autres.

Les tests non paramétriques que nous allons définir seront utilement comparés aux tests paramétriques. Rappelons que la médiane  $\mu$  est définie comme l'un des nombres tels que

$$P(X \geq \mu) = P(X \leq \mu) = 1/2 .$$

Le problème est de tester  $H_0 : "\mu = \mu_0"$  contre  $H_1 : "\mu \neq \mu_0"$  (test bilatère) ou contre  $H'_1 : "\mu > \mu_0"$  test unilatère.

On supposera  $P(X = \mu_0) = 0$ .

#### a) Le test du signe.

Si  $H_0$  est vrai, la v.a.  $S_n = \sum_{i=1}^n 1_{(X_i > \mu)}$  suit une loi binômiale  $B(n, 1/2)$ . Les grandes et les petites valeurs de  $S_n$  conduiront donc au rejet de l'hypothèse  $H_0$  (dans le cas bilatère). De manière précise si  $\alpha$  est donné, on peut déterminer des entiers  $k_1(\alpha)$  et  $k_2(\alpha)$  tels que

$$k_1(\alpha) = \inf\{h, P(S_n \geq h) \leq \frac{\alpha}{2}\}$$

$$k_2(\alpha) = \sup\{h, P(S_n \leq h) \leq \frac{\alpha}{2}\}$$

et la région de rejet sera

$$(S_n \geq k_1(\alpha)) \cup (S_n \leq k_2(\alpha)).$$

De même dans le cas unilatère la région de rejet sera  $(S_n \geq k(\alpha))$  avec

$$k(\alpha) = \inf\{h, \sum_{k=h}^n \frac{1}{2^n} \binom{n}{k} \leq \alpha\}.$$

On sait que pour  $n$  grand

$$\frac{S_n - n/2}{1/2\sqrt{n}} \sim N(0,1),$$

on peut donc approximer la loi de  $S_n$  par une loi normale.

Cependant, comme on prend des v.a. discrètes, on a intérêt pour tel unilatère par exemple, à choisir  $k(\alpha)$  de la manière suivante:

$$k(\alpha) = \left[ \frac{1}{2} + \frac{1}{2\sqrt{n}} z_\alpha + \frac{1}{2} n \right]$$

où  $z_\alpha$  est défini par :  $\frac{1}{\sqrt{2\pi}} \int_{z_\alpha}^{\infty} e^{-x^2/2} dx = \alpha$  (la correction  $1/2$  venant du caractère discret de la binômiale, on laisse au lecteur le soin de la justifier).

On peut calculer la puissance de ce test sur des modèles paramétriques. Si  $H_0$  est rejeté, et si on appelle  $p$  la probabilité que  $X > \mu$ , on a, par exemple pour le test unilatère :

$$\pi(p) = \sum_{k=k(\alpha)}^n \binom{n}{k} p^k (1-p)^{n-k}.$$

Exemple numérique :  $\alpha = 0,05$ . Alors il n'existe pas de test de niveau significatif  $\alpha$ , la valeur la plus proche, par exemple pour  $n = 16$  est  $k = 12$ , qui donne

$$\sum_{k=12}^{16} \binom{16}{k} \frac{1}{2^{16}} = 0,038.$$

Si on veut  $\alpha = 0,05$ , il faut, comme à l'habitude pour les lois discrètes, faire un tirage au sort auxiliaire. Si on prend  $\alpha = 0,038$ , la loi gaussienne  $N(29,04 ; 1)$  on peut tester  $H_0 : \mu_0 = 28$  contre  $H_1 : \mu_0 > 28$ . On a alors  $p = 0,85$  (c.à.d. on est très loin de  $H_0$ ).

La puissance du test binomial  $\pi(0,85)$  est 0,92. La puissance du test de la moyenne (cf. chapitre sur les gaussiennes) est

$$\begin{aligned}\pi'(29,04) &= P(\bar{X} > 28,44) \\ &= 0,99.\end{aligned}$$

La différence est considérable mais on est dans un cas vraiment défavorable car l'écart à la moyenne est de l'ordre de  $\sigma$ .

b) Intervalle de confiance et test d'ordre sur la médiane.

On a, sous  $H_0$ ,

$$\begin{aligned}P(X_{(k_1)} \leq \mu_0 \leq X_{(k_2)}) &= P(U_{(k_1)} < \frac{1}{2} < U_{(k_2)}) \\ &= \int_{u < \frac{1}{2} < v} g_{k_1, k_2, n}(u, v) du dv = c(n, k_1, k_2).\end{aligned}$$

Si  $g_{k_1, k_2, n}$  est la densité des statistiques d'ordre numéro  $k_1$  et  $k_2$  d'un échantillon de loi uniforme. Si  $\alpha$  est donné, il existe des choix (non uniques) de  $k_1$  et  $k_2$  tels que

$$c(n, k_1, k_2) \geq 1 - \alpha$$

et aussi voisin que possible de  $1 - \alpha$ . On en déduit :

a) un intervalle de confiance au niveau  $\alpha$  pour  $\mu$ , l'intervalle

$$[X_{(k_1)}, X_{(k_2)}].$$

b) un test :  $H_0$  est rejeté si

$$\mu \notin [X_{(k_1)}, X_{(k_2)}].$$

Pour  $n$  grand, on utilise l'approximation normale, pour

$$S_n = \sum_{j=1}^n 1_{(-\infty, \mu_0)}(X_j)$$

qui est telle, sous  $H_0$ , que

$$\frac{S_n - n/2}{1/2\sqrt{n}} \sim N(0, 1).$$

Ayant calculé  $z_\alpha$  comme précédemment, on calcule les solutions entières de l'équation approchée,

$$\frac{2y - n}{\sqrt{n}} = \pm z_\alpha,$$

soient  $y_1, y_2$  puis  $m$  tel que  $X_{(m)} \leq \mu_0 \leq X_{(m+1)}$  et si  $m \notin [y_1, y_2]$



on rejette  $H_0$ . Le lecteur détaillera aisément cette procédure.

c) Cas de 2 échantillons apariés.

Soient 2 échantillons apariés de v.a.  $X$  et  $Y$ , non nécessairement indépendants.  $(X_1 \dots X_n)$  et  $(Y_1 \dots Y_n)$ , sous la forme  $(X_1, Y_1) \dots (X_n, Y_n)$ .

Soit  $Z_i = X_i - Y_i$ , la loi de  $Z_i$  étant supposée continue. On teste ici, la médiane des différences est 0 (et non la différence des médianes, si  $X$  et  $Y$  sont des v.a. indépendantes ou simplement symétriques, l'égalité des médianes équivaut à une médiane des différences, nulle). On est ramené au problème du test de signe déjà vu.

d) Test de Wilcoxon ou test signe et rang.

Soit  $\mu_0$  la médiane des  $X_i$ . N'utiliser que les signes est appauvrir beaucoup l'information donnée par l'échantillon. Soit  $D_i = X_i - \mu_0$ . On ordonne les  $|D_i|$  par ordre croissant et on note  $T^+$  (resp.  $T^-$ ) la somme des rangs des  $D_i \geq 0$  (resp. des  $D_i < 0$ ),  $T^+ + T^- = \frac{n(n+1)}{2}$ . On s'intéresse à la loi de  $T^+$  (qui fait intervenir les valeurs des  $D_i$ , pas seulement leurs signes). On suppose, et il est important de le retenir, que la loi  $F$  des  $X$  est symétrique par rapport à la médiane  $\mu_0$ , c'est-à-dire que :  $F(\mu_0 - x) = 1 - F(\mu_0 + x)$ ,  $x > 0$  et donc que  $P(D_i > a) = P(D_i < a)$ , la loi des  $D_i$  étant symétrique au sens ordinaire.

Comme nous l'avons déjà dit, la médiane, dans le cas symétrique doit être considérée comme le paramètre naturel de centrage et un test sur la valeur de la médiane est un test de centrage. Notant

$$Z_{(i)} = 1_{(D_{(i)} > 0)}$$

on a

$$T^+ = \sum_{i=1}^n i Z_{(i)}.$$

Il est clair que sous  $H_0$ ,  $P(Z_{(i)} = 1) = \frac{1}{2}$ . On en déduit :  $E_{H_0} T^+ = \frac{n(n+1)}{4}$

et  $\text{Var}_{H_0} T^+ = \frac{n(n+1)(2n+1)}{24}$ .

Quand  $H_0$  est fautive, un calcul simple montre que

$$P(Z_{(i)} = 1) = n \binom{n-1}{i-1} \int_0^{\infty} [F(u+\mu) - F(-u-\mu)]^{i-1} \\ \times [1-F(u+\mu) + F(-u-\mu)]^{n-i} dF(u+\mu).$$

On a alors  $\text{Var } T^+ = 4 \sum_{i=1}^n \frac{i^2 p_i (1-p_i)}{n^2 (n+1)^2}$  où  $p_i = P(Z_{(i)} = 1)$ , et comme

$p_i (1-p_i) \leq \frac{1}{4}$  on a  $\text{Var } \frac{T^+}{n(n+1)} \rightarrow 0$  dans tous les cas.

Si  $\mu > \mu_0$ , il est facile de voir que  $p_i > \frac{1}{2}$  et si  $\mu \neq \mu_0$  on a  $p_i \neq \frac{1}{2}$  en général.

Le test unilatère sera donc défini par une région de rejet du type

$$\frac{2T^+}{N(N+1)} - \frac{1}{2} \geq k$$

et le test bilatère par

$$\left\{ \frac{2T^+}{N(N+1)} - \frac{1}{2} > k_1 \right\} \left\{ \frac{2T^+}{N(N+1)} - \frac{1}{2} < k_2 \right\}.$$

Les valeurs de la loi de  $T^+$  sont tabulées.

On vérifie que pour  $n$  grand ( $> 20$ )

$$\frac{4T^+ - n(n+1)}{\sqrt{\frac{2}{3} n(n+1)(2n+1)}} \sim N(0,1).$$

On peut évidemment employer cette technique pour le cas de 2 échantillons appariés et aussi pour tester la symétrie d'une loi (mais ce n'est pas très bon).

Conclusion : l'intérêt des tests de signe est qu'ils ne nécessitent pas des mesures très précises. Le test de signe est bien adapté à des v.a. qualitatives (0,1).

L'efficacité asymptotique relative de ces tests par rapport au test de Student pour une loi normale est de 0,64 pour le test de signe et de 0,96 pour le test signe et rang ce qui est bon.

### c) Comparaison de 2 échantillons.

a) On ne suppose plus les 2 échantillons appariés (de manière naturelle ou artificielle) et en particulier on ne suppose plus leurs tailles égales. De plus, on suppose toujours ici les 2 échantillons indépendants :

soit  $X_1 \dots X_m$  et  $Y_1 \dots Y_n$ . On teste  $H_0 : F_X = F_Y$ , c'est-à-dire les échantillons sont issus de populations identiques en loi (et que l'on peut donc considérer comme identiques tout court). Si les populations sont supposées normales, on sait bien tester " $m_X = m_Y$ " ou " $\sigma_X = \sigma_Y$ ", assez mal " $m_X = m_Y$  et  $\sigma_X = \sigma_Y$ ". Dans les problèmes non paramétriques, utilisables dès que l'on a quelques doutes, sur les modèles paramétriques, on fait l'hypothèse que les lois  $F$  sont continues (quoique on puisse appliquer, avec quelques précautions les résultats à des lois discrètes).

Les alternatives classiques à  $H_0$  sont :

$$H_1 : F_Y \neq F_X$$

$$H_2 : F_Y(x) \geq F_X(x) \text{ pour tout } x \text{ où } Y \text{ est plus dispersé vers } +\infty \text{ que } X.$$

Nous étudierons plus systématiquement comme dans le cas paramétrique, les hypothèses sur le centrage (en anglais : location) et sur l'échelle, à savoir :

$$H_{C_1} : "F_Y(x) = F(x-\theta_0) \text{ pour un certain } \theta_0 \neq 0 "$$

$$H_{C_2} : "il existe \theta, F_Y(x) = F(x-\theta)"$$

et de même

$$H_{E_1} : "F_Y(x) = F_X(\sigma_0 x)" \text{ pour un certain } \sigma_0 > 0$$

$$H_{E_2} : "il existe \sigma \text{ tel que } F_Y(x) = F_X(\sigma x)".$$

Une autre hypothèse, dite de Lehman qui peut intervenir dans des problèmes de fiabilité, est

$$H_L : "F_Y(x) = F_X(x)^k, \text{ } k \text{ entier}"$$

( $Y$  est un maximum de  $k$  v.a.  $X$ ).

Nous allons d'abord passer en revue l'application des tests déjà vus.

a) Test de runs (v.a. nulles)

On classe les  $X$  et les  $Y$  par ordre croissant pour obtenir un écrit symboliquement sous la forme

X X Y X X X Y Y X Y Y X ...

On compte le nombre (total) de runs  $R$ , s'il est trop petit, cela indique souvent que  $X$  a tendance à être plus grand que  $Y$ , ou réciproquement; la région de rejet de  $H_0$  contre  $H_1$  est donc  $R < c_\alpha$ ,  $R$  est tabulé. (cf. paragraphes précédents, échantillons de tailles  $m, n$ ).

b) Test de Kolmogorov et Smirnov.

Soient  $\hat{F}_X$  et  $\hat{F}_Y$  les distributions empiriques de  $X$  et  $Y$ .

Si  $H_0$  est vraie,

$$D_{m,n} = \max |\hat{F}_X - \hat{F}_Y|$$

suit une loi indépendante de  $F$ , (immédiat en passant aux variables uniformes  $U_i$  et  $V_j$  associées aux  $X$  et aux  $Y$ ).

Le calcul de la loi de  $D_{m,n}$  est assez finaud et utilise profondément les propriétés des marches aléatoires arrêtées sur certaines barrières (cf. Cours de Probabilités). Un théorème de Smirnov donne la distribution asymptotique de

$$\lim_{\substack{m, n \rightarrow \infty \\ \frac{m}{n} = \text{cte}}} P \left[ \sqrt{\frac{mn}{m+n}} D_{m,n} < z \right].$$

On trouve que c'est celle du  $\sqrt{n} D_n$  déjà vu.

On peut aussi étudier  $D_{m,n}^+ = \max(\hat{F}_X - F_Y)$  et tester  $H_0$  contre  $H_2$

Ce sont des tests très bons, applicables moyennant les précautions habituelles aux cas discrets.

c) Le test de la médiane.

Soit à tester  $H_0 : "F_X = F_Y"$  contre une alternative  $H_1$  à préciser, où  $F_X$  (resp.  $F_Y$ ) est la loi de  $X$  (resp.  $Y$ ). On a un  $m$ -échantillon  $X_1 \dots X_m$  de  $X$ , un  $n$ -échantillon  $Y_1 \dots Y_n$  de  $Y$

$$N = n+m, \quad Z = (Z_1 \dots Z_{m+n}) = (X_1 \dots X_m, Y_1 \dots Y_n).$$

Soit  $\hat{\mu}$  la médiane empirique de  $Z$ ,

$$\hat{\mu} = Z_{\left(\frac{N+1}{2}\right)} \quad \text{si } N \text{ est impair}$$

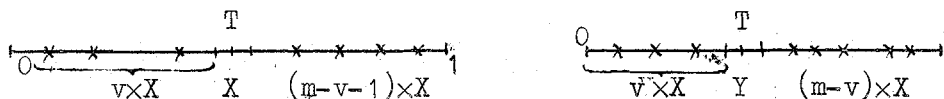
$$\hat{\mu} \in \left[ Z_{\left(\frac{N}{2}\right)}, Z_{\left(\frac{N+2}{2}\right)} \right], \quad N \text{ pair,}$$

(pour se fixer les idées on peut prendre le milieu de l'intervalle).

Soit  $V = \sum_{i=1}^n 1(X_i < T)$ . Sous  $H_0$ , et si  $F$  est continue, posant  $U = F(Z)$ , on montre aisément que la loi de  $V$  est indépendante de  $F$ , c'est celle que l'on obtiendrait pour des variables de la loi uniforme sur  $(0, 1)$ .

Les schémas ci-dessous décrivent la situation

$$\begin{aligned}
 P(V = v) &= \left[ m \binom{m-1}{v} \binom{n}{\frac{N-1}{2}-v} + n \binom{m}{v} \binom{n-1}{\frac{N-1}{2}-v} \right] \\
 &\times \int_0^1 x^{\frac{N-1}{2}-v} (1-x)^{\frac{N-1}{2}} dx \\
 &= \frac{\binom{m}{v} \binom{n}{\frac{N-1}{2}-v}}{\binom{N}{\frac{N-1}{2}}}
 \end{aligned}$$



Dans le cas où  $N$  est pair, on fait le même calcul en remplaçant dans les formules  $\frac{N-1}{2}$  par  $\frac{N}{2}$  (cf. interprétation plus haut).

Remarquons que sous  $H_0$ ,  $\frac{V}{m} \rightarrow \frac{1}{2}$  ( $m$  grand).

Le test contre l'alternative générale  $H_1$  " $F_X = F_Y$ " a donc pour région de rejet

$$\{V > c_1(\alpha) \cup V < c_2(\alpha)\}$$

$c_1, c_2$  convenablement choisis.

Le test devient très bon quand on prend pour alternative,

$$H_1' : F_Y(x) = F_X(x-\theta) \text{ pour un } \theta \neq 0.$$

On a alors :

$$H_1'' : \theta > 0 \Rightarrow \text{région de rejet } V > c_\alpha$$

$$H_1'' : \theta < 0 \Rightarrow \text{région de rejet } V < c_\alpha$$

On a 
$$E_{H_0} V = m \frac{\frac{N-1}{2}}{N} \quad (\text{cas impair})$$

et 
$$\text{Var}_{H_0} V = \frac{mn(N+1)}{N^2}.$$

Si  $N \rightarrow \infty$ ,  $m \rightarrow \infty$  et  $\frac{m}{N} = \text{cte} = \lambda$ , la loi hypergéométrique de  $V$ , se comporte comme une binômiale  $(N, \lambda)$  et est donc à peu près normale une fois normée.

On utilise la statistique

$$V' = \frac{V - m/2}{(mn/N^3)^{1/2}} \sim N(0, 1).$$

Une remarque importante est la suivante :

$$\begin{aligned} NV &\sim m \frac{N-1}{2} \\ V' &\sim \frac{V - m/2}{\sqrt{mn \left(\frac{N-1}{2}\right)^2 / N}} \\ &\sim \frac{V/m - W/n}{\sqrt{\left(\frac{V+W}{N}\right) \frac{1-V+W}{N} \frac{N}{nm}}} \end{aligned}$$

où

$$W = \sum_1^T 1_{(Y_i < T)},$$

soit

$$V' \sim \frac{V/m - W/n}{\sqrt{\hat{p}(1-\hat{p})(1/m + 1/n)}}$$

où  $\hat{p}$  est l'estimateur du nombre de v.a.  $X < T$ , donc  $\hat{p} \neq 1/2$  sous  $H_0$ . Sous cette forme, on a une forte analogie avec le test comparant 2 proportions supposées normales et indépendantes. On peut aussi, évidemment obtenir par cette méthode un intervalle de confiance pour  $\theta$ . On suppose que  $\theta$  représente la médiane de  $Y$  (et donc que celle de  $X$  est nulle). On choisit  $c_1(\alpha)$  et  $c_2(\alpha)$  tels que

$$P(c_1(\alpha) < U < c_2(\alpha)) \geq 1-\alpha \quad (c_1, c_2 \text{ entiers}).$$

On calcule la médiane  $T$  de  $Z$ . On utilise la dualité tests-intervalles de confiance (cf. Chap. II).

Soit  $H_\theta$  l'hypothèse "la loi de  $Y-\theta$  est celle de  $X$ ".

$H_\theta$  est acceptée si et seulement si

$$X_{(c_1+1)} < Y_{\left(\frac{N-1}{2} - c_1\right)^+ + \theta}$$

et

$$X_{(c_2)} > Y_{\left(\frac{N-1}{2} - c_2 + 1\right) - \theta}$$

l'intervalle de confiance est donc par dualité

$$\left[ Y_{\left(\frac{N-1}{2} - c_2 + 1\right)} - X_{(c_2)}, Y_{\left(\frac{N-1}{2} - c_1\right)} - X_{(c_1 + 1)} \right].$$

d) Le test U (Mann-Whitney).

Les notations sont les mêmes qu'en c). Soit

$$D_{ij} = 1_{(Y_j > X_i)}, \quad i = 1 \dots m, \quad j = 1 \dots n$$

$$U = \sum_{i=1}^m \sum_{j=1}^n D_{ij}.$$

On pose  $\pi = P(Y < X)$ .

Si l'hypothèse  $H_0$  est vraie, on a  $\pi = 1/2$ . On a  $\pi > 1/2$  si Y a "tendance" à être plus petite que X.

Il est donc clair que sous  $H_0$  la loi de U est indépendante de celle de X. Elle est assez difficile à calculer et est tabulée. L'alternative générale  $H_1$  " $F_X \neq F_Y$ " donne une région de rejet du type :

$$\left| U - \frac{mn}{2} \right| > C,$$

fondée sur l'idée que  $E U = \frac{mn}{2}$  sous  $H_0$ . On a de plus  $E U = mn \pi$ , ce qui explique la force du test et  $\text{var} \frac{U}{mn} \rightarrow 0$  pour toute hypothèse sur (X, Y).

Si on teste  $H_0$  contre  $H_2$   $F_X(x) > F_Y(x)$  ( $\pi > \frac{1}{2}$ ) et la région de rejet est du type  $U > \frac{mn}{2} + C$ .

Sous  $H_0$ ,  $\frac{U - mn/2}{\sqrt{mn(N+1)/12}} \sim N(0, 1)$

ce qui donne une bonne approximation (dès que  $n, m \geq 6$ ).

Ce test présente l'avantage d'être utilisable pour des v.a. discrètes, en présence de noeuds. On pose alors

$$D_{ij} = \begin{cases} 1 & X_i > Y_j \\ 0 & X_i = Y_j \\ -1 & X_i < Y_j \end{cases} \quad \pi^+ = P(X > Y)$$

$$E U = mn(\pi^+ - \pi^-) \quad \text{où} \quad \pi^- = P(X < Y)$$

et  $E U = 0$  sous  $H_0$ .

Des calculs simples permettent d'obtenir  $\text{Var} U$  en fonction du  $H_0$ .

nombre de noeuds.

U permet aussi d'obtenir un intervalle de confiance dans une estimation où  $F_Y$  est choisie du type  $F_X(x - \theta)$ . On calcule les différences  $y_j - x_i = d_{ij}$ . On cherche le nombre U de différences  $< \theta$  ( $H_{\theta=0}$  est rejeté (au niveau  $\alpha/2$ ) si le nombre de telles différences est  $< k(\alpha/2)$ ). On ordonne alors ces différences. En appliquant le principe de dualité test-intervalle de confiance, on montrera qu'un intervalle de confiance pour  $\theta$  est tel que  $\theta \geq (k+1)^{\text{ème}}$  (différence  $d_{ij}$ ) (les  $d_{ij}$  ordonnées par ordre croissant).

Même raisonnement pour la borne supérieure de l'intervalle de confiance qui est donc

$$(d_{ij})_{(k(\frac{\alpha}{2}))}, (d_{ij})_{(N-k(\frac{\alpha}{2}))}$$

Dans ces conditions

$$(F_Y(x) = F_X(x - \theta))$$

le test U est très bon. Comparé au test de Student, sur une population quelconque, son efficacité asymptotique est  $> 0,83$  et sur une population normale  $> 0,96$ .

e) Le test de Wilcoxon.

Il s'agit d'une présentation (légèrement) différente de d). Nous gardons les notations précédentes. Posons

$$\begin{aligned} R_i &= 1 \text{ si } Z_{(i)} \text{ est un } X \\ &= 0 \text{ si } Z_{(i)} \text{ est un } Y \end{aligned}$$

analogue en statistiques non paramétriques de ce que nous appellerons contraste dans le cas des modèles gaussiens d'analyse de variance, est constitué des combinaisons linéaires  $\sum_{i=1}^N a_i R_i = a(R)$ . Si l'on pose en particulier  $W = \sum_{i=1}^N i R_i$  on obtient le test de Wilcoxon.

Si U désigne la statistique introduite précédemment au d), nous laissons au lecteur le soin de vérifier que  $W = U + \frac{m}{2} (m+1)$ ; en particu-



lier  $W$  est symétrique autour de  $E_{H_0} W = \frac{m(N+1)}{2}$

$$\min W = \frac{m(m+1)}{2}, \quad \max W = \frac{m(2N-m+1)}{2}.$$

Les valeurs extrêmes de  $W$  constituent une région de rejet.

f) Tests sur l'échelle.

Supposons maintenant avoir à tester " $H_0$  contre  $H_1$ " : il existe  $\theta \neq 1$ ,  $F_Y(x) = F_X(\theta x)$ , où  $\theta$  doit être interprété par exemple, comme  $\frac{\sigma_Y}{\sigma_X}$  rapport des dispersions de  $Y$  et de  $X$ .

On interprète donc  $\theta$  comme une variable d'échelle (cf. chapitre I).

Le test gaussien analogue est le test  $F$  pour des moyennes inconnues mais égales. Le test de Mood consiste à choisir comme statistique

$$M = \sum_{i=1}^N \left(i - \frac{N+1}{2}\right)^2 R_i$$

qui pondère les grandes valeurs (absolues) par de fortes masses et caractérise donc la dispersion. Si  $M$  est trop grand, on rejettera  $H_0$  ( $X$  plus dispersé que  $Y$ ), si  $M$  est trop petit on rejettera  $H_0$  ( $X$  moins dispersé que  $Y$ ).

Le test  $S$  de Siegel-Tukey utilise

$$S = \sum_{i=1}^N a_i Z_i \quad \text{où} \quad a_i = 2i, \quad 1 \leq i \leq N/2$$

$$= 2(N-i) + 2, \quad \frac{N}{2} < i \leq N \quad \text{si } N \text{ est pair}$$

et

$$a_i = 2i-1, \quad 1 \leq i \leq N/2, \quad N \text{ impair}$$

$$= 2(N-i) + 1, \quad N/2 < i < N.$$

Son avantage est d'avoir la même distribution que la statistique  $W$  du test de Wilcoxon.

## CHAPITRE V

### ESTIMATIONS ET TESTS EN VARIABLE GAUSSIENNE

#### I. Lois de certaines variables aléatoires définies à partir de variables gaussiennes.

a)  $\chi^2$  et  $\chi'^2$ .

Définition. Une v.a.  $Y$  est une variable  $\chi'^2(n, \lambda)$  si

$$Y = \sum_{i=1}^n X_i^2$$

où les  $X_i$  sont des v.a. gaussiennes indépendantes de moyenne  $\mu_i$  et de variance 1 :

$$X_i \sim N(\mu_i, 1)$$

avec  $\sum \mu_i^2 = \lambda$ .  $\lambda$  est appelé paramètre de non-centralité. On note

$$\chi^2(n) \sim \chi'^2(n, 0).$$

Proposition. La loi de  $Y$  ne dépend en effet que de  $\lambda$ .

Soient  $X_i = N_i + \mu_i$

$Z_i = T_i + v_i$  où  $N_i$  et  $T_i$  sont indépendantes de loi  $N(0, 1)$  et  $\sum \mu_i^2 = \sum v_i^2$ .

Soit  $G$  une transformation orthogonale de  $\mathbb{R}^n$  telle que  $G(\mu) = v$

(il en existe une puisque  $\|\mu\| = \|v\|$ ) alors :

$$\|GX\|_{\mathbb{R}^n}^2 = \|X\|_{\mathbb{R}^n}^2 \quad \text{et} \quad \|G(N+\mu)\|^2 = \|G(N)\|^2 + 2 \langle GN, G\mu \rangle + \|G(\mu)\|^2$$

si on pose :  $T' = GN$  ; les  $T'_i$  sont des v.a. indépendantes de loi  $N(0, 1)$

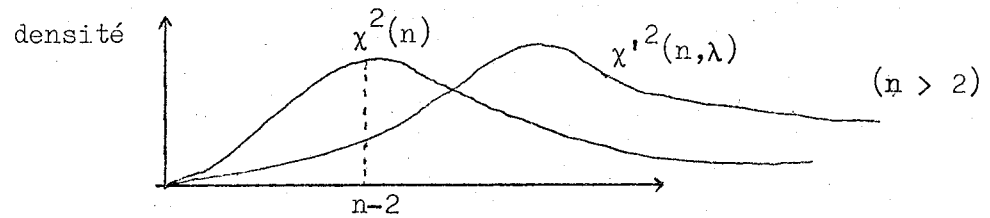
car  $G^t G = I$ , ( $G$  orthogonale cf. Cours C 4) et

$$\langle GN, G\mu \rangle = \langle N, \mu \rangle$$

donc  $\sum X_i^2$  a même loi que  $\sum Y_i^2$ , puisque  $(T_1 \dots T_n)$  et  $(T'_1 \dots T'_n)$  ont même loi, et  $\sum X_i^2 = \|T'\|^2 + \langle T', v \rangle + \|v\|^2$ ,  $\sum Y_i^2 = \|T\|^2 + \langle T, v \rangle + \|v\|^2$ .

Remarque. On peut aussi calculer la fonction caractéristique de  $Y$  et constater qu'elle ne dépend que de  $n$  et  $\lambda$  :

$$\Phi(t) = \frac{1}{(\sqrt{1-2it})^n} e^{\frac{it}{1-2it} \sum_{j=1}^n \mu_j^2}$$

Distribution d'un  $\chi'^2$ .

On calcule :

$$EY = \sum EX_i^2 = n + \lambda$$

$$\text{Var } Y = 2n + 4\lambda$$

on démontre que :

$P(Y < \rho)$  pour un  $\rho$  fixé est une fonction strictement croissante de  $\lambda$  (paramètre de non centralité).

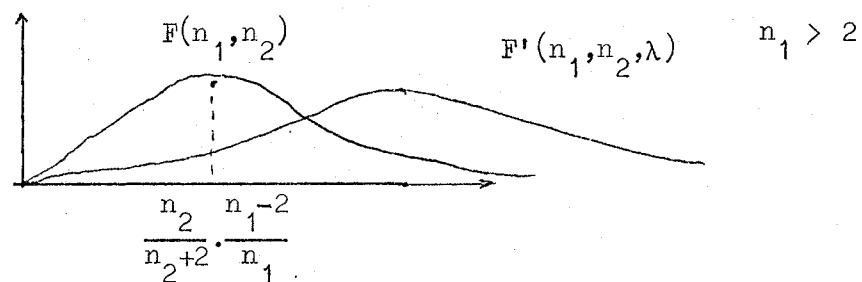
b) Variables F et F'.

Si  $U_1$  et  $U_2$  sont des v.a. indépendantes de loi  $\chi'^2(n_1, \lambda)$  et  $\chi'^2(n_2)$

$$\frac{U_1/n_1}{U_2/n_2} = F'(n_1, n_2, \lambda)$$

est une variable  $F'$  de Fischer de paramètres  $n_1, n_2$  et  $\lambda$ .

Si  $\lambda = 0$  on a une variable  $F(n_1, n_2)$ .

Densités

On vérifie, comme plus haut :

$$P(F'_{n_1, n_2, \lambda} < \rho)$$

est une fonction strictement croissante de  $\lambda$ .

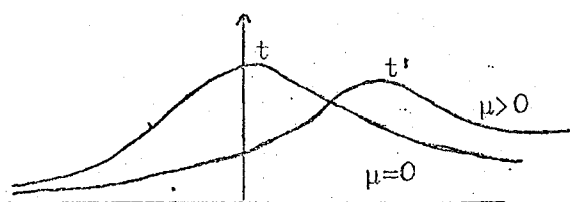
c) Variable t de Student.

Si  $X$  et  $U$  sont des v.a. indépendantes de loi respectivement  $N(\mu, 1)$  et  $\chi^2(n)$ , la variable  $\frac{X}{\sqrt{U/n}}$  a une distribution  $t'$  de Student :

$$t'(n, \mu) = \frac{X}{\sqrt{U/n}}$$

si  $\mu = 0$

$$t'(n, 0) = t(n).$$

DensitésII. THEOREME DE COCHRAN.

Théorème. Soient  $Y_1, \dots, Y_i, \dots, Y_n$  des v.a. gaussiennes  $N(\eta_i, \gamma)$  indépendantes et  $Q_1, Q_2, \dots, Q_s$  des formes quadratiques des  $(Y_i)_{i \in 1 \dots n}$  telles que :

$$\sum_1^n Y_i^2 = Q_1 + Q_2 + \dots + Q_s$$

soit  $n_j = \text{rang } Q_j$ . Si  $\sum_1^s n_j = n$  alors les  $Q_j$  sont des v.a. indépendantes de loi  $\gamma \chi^2(n_j, \lambda_j)$  où

$$\lambda_j = Q_j(\eta_1 \dots \eta_n).$$

Application. Estimation de l'espérance et de la variance d'une loi gaussienne.

Soit  $\{Y_i, i = 1 \dots n\}$  un échantillon de variables  $N(\mu, \sigma^2)$  alors :

$$\sum Y_i^2 = n\bar{Y}^2 + \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad \text{où } \bar{Y} = \frac{1}{n} \sum_1^n Y_i$$

les rangs des deux formes quadratiques étant respectivement 1 et  $n-1$  le théorème de Cochran s'applique :

$\sum (Y_i - \bar{Y})^2$  est indépendante de  $\bar{Y}$  et distribuée comme un  $\gamma \chi_{n-1}^2$  ;

en effet  $E(Y_i - \bar{Y}) = 0$  le 2<sup>e</sup> paramètre est donc nul.

Comme  $E \chi_{n-1}^2 = n-1$

$$S = \frac{1}{n-1} \sum (Y_i - \bar{Y})^2 \text{ est estimateur sans biais de } \gamma,$$

comme on l'a vu au chapitre II.

Pour la démonstration du théorème démontrons la

Proposition (algèbre).

Si  $\sum_{i=1}^n y_i^2 = Q_1 + \dots + Q_s$  où  $Q_j$  est une forme quadratique en  $(y_i)_{i=1 \dots n}$  de rang  $n_j$ , la condition

$$n = \sum_{j=1}^s n_j$$

est nécessaire et suffisante pour qu'il existe une transformation orthogonale  $A$  telle que

$$z = Ay, \quad z = (z_1 \dots z_n)$$

et que

$$Q_1 = \sum_{i=1}^{n_1} z_i^2$$

$$Q_2 = \sum_{i=n_1+1}^{n_1+n_2} z_i^2$$

$$\dots$$

$$Q_s = \sum_{i=n_1+\dots+n_{s-1}+1}^{n_1+\dots+n_s} z_i^2$$

Démonstration. Condition nécessaire. Si une telle transformation existe

alors  $\sum_{i=1}^n y_i^2 = \sum_{i=1}^{n_1+\dots+n_s} z_i^2$  dans un espace euclidien, le rang d'une forme quadra-

tique étant invariant par changement de base, on a

$$n = n_1 + \dots + n_s.$$

Condition suffisante. Le rang de  $Q_j$  est  $n_j$ , il existe donc des formes linéaires  $z_\alpha^j$  en  $y_i$  ( $i = 1 \dots n$ ) telle que

$$Q_j = \sum_{\alpha} \delta_{\alpha}^j z_{\alpha}^j \quad \text{où } \delta_{\alpha}^j = +1 \text{ ou } -1$$

et  $\alpha = n_1 + \dots + n_{j-1} + 1, \dots, n_1 + \dots + n_j$ .

Si  $\sum_{j=1}^s n_j = n$  on a donc  $n$  formes linéaires  $z_{\alpha}^j$  définissant une matrice  $A$  dans la base de  $y_i$ , soit :

$$[z] = A[y] \quad ([z], [y] \text{ vecteurs colonnes}).$$

Soit  $D$  la matrice diagonale des  $\delta_{\alpha}^j$ ,  $\alpha = 1 \dots n$ . Alors

$$\sum_{j=1}^s Q_j = \sum_{j=1}^s \sum_{\alpha} \delta_{\alpha}^j z_{\alpha}^j{}^2 = {}^t[z] D [z] = {}^t[y] {}^t A D A [y]$$

mais on a aussi

$$\sum_{j=1}^s Q_j = \sum y_i^2 = {}^t[y][y],$$

la matrice  ${}^t ADA$  est symétrique donc l'égalité des 2 formes quadratiques entraîne

$${}^t ADA = I$$

donc  $A$  est régulière ;

démontrons que  $D = I$  en effet si  $\delta_{\beta}^{j_0} = -1$  comme  $[y] = A^{-1}[z]$  en prenant

$z_{\alpha}^j = 0$  pour  $(\alpha, j) \neq (\beta, j_0)$  on a

$$\begin{aligned} \sum_{i=1}^n y_i^2 &= \sum_{j=1}^s \sum_{\alpha=1}^{n_j} \delta_{\alpha}^j z_{\alpha}^j{}^2 \\ &= \delta_{\beta}^{j_0} \end{aligned}$$

$= -1$  ce qui est impossible.

Donc  $D = I$  et  ${}^t AA = I$  soit  $A$  est orthogonale

cqfd.

Remarque. La condition  $\sum_{j=1}^s n_j = n$  implique en particulier que les formes  $Q_j$  sont positives.

Démonstration du théorème de Cochran.

Condition nécessaire. Si les  $Q_j$  sont des v.a.  $\sim \chi^2$  indépendantes de dimension  $n_j$ , il est immédiat que  $\sum_{j=1}^s Q_j \sim \chi^2$  de dimension  $\sum_{j=1}^s n_j$  et comme  $\sum_{j=1}^s Q_j = \sum_{i=1}^n y_i^2$  on a  $n = \sum_{j=1}^s n_j$ .

Condition suffisante. Supposons  $n = \sum_{j=1}^s n_j$ . Soit  $A$  la transformation orthogonale définie dans le lemme précédent. Les v.a.  $Z_{\alpha}^j = z_{\alpha}^j(Y_1 \dots Y_n)$  sont des v.a. gaussiennes indépendantes. On en déduit que les  $Q_j$  sont des v.a.  $\sim \gamma \chi^2$  indépendantes. En effet si  $[Z] = (Z_{\alpha}^j)_{\alpha=1 \dots n_j, j=1 \dots s}$  On a  $\text{cov}[Z] = E(( [Z] - E[Z] ) \otimes ( [Z] - E[Z] )) = {}^t A \text{cov} Y A = {}^t A \gamma I A = \gamma I$  (cf. cours probabilité chap. 4).

On a le coefficient de non centralité de  $Q_j$  par  $\lambda_j = \sum_{\alpha} (E Z_{\alpha}^j)^2$   
 or  $Q_j(Y_1 \dots Y_n) = \sum_{\alpha=1}^{n_j} Z_{\alpha}^j{}^2$ .  $Z_{\alpha}^j$  est la  $(\alpha, j)$ <sup>ème</sup> composante de  $A[Y]$   
 donc  $E Z_{\alpha}^j = E(A[Y])_{(\alpha, j)}$

$$\sum_{\alpha=1}^{n_j} (E Z_{\alpha}^j)^2 = \sum_{\alpha=1}^{n_j} [(E A[Y])_{(\alpha, j)}]^2$$

$$\sum_{\alpha=1}^{n_j} (A[y])_{(\alpha, j)}^2 = Q_j(y_1 \dots y_n)$$

donc  $\sum_{\alpha=1}^{n_j} (E Z_j^\alpha)^2 = Q_j (E Y_1 \dots E Y_n)$  puisque E et A commutent. cqfd.

Remarque. Le théorème est encore valable si on remplace  $\sum_{j=1}^s n_j = n$  par  $\sum_{j=1}^s n_j \leq n$ .

### III. EXEMPLES DE TESTS CONSTRUITS AVEC DES V.A. GAUSSIENNES.

#### 1. Comparaison des moyennes de 2 v.a. gaussiennes indépendantes de même variance.

On se propose de comparer les moyennes de 2 v.a. X et Y avec  $X \sim N(\mu, \sigma^2)$  et  $Y \sim N(\nu, \sigma^2)$ .

On veut tester  $H_0 : \{\mu = \nu\}$  contre  $H_1 : \{\mu \neq \nu\}$ .

Si on a un n-échantillon du premier type de v.a. soit  $X_1 \dots X_n$  et un m-échantillon du 2<sup>e</sup> type de modèle est

$$(\mathbb{R}, \mathcal{B}, \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{(x-\mu)^2}{2\sigma^2})^{\otimes n} \otimes (\mathbb{R}, \mathcal{B}, \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{(y-\nu)^2}{2\sigma^2})^{\otimes m}, (\mu, \nu) \in \mathbb{R}^2$$

et il s'agit d'un test particulier d'homogénéité à savoir est-ce que les v.a. ont même loi ?

Si  $m = n$  (on dit que l'échantillonnage est équilibré), en groupant arbitrairement les v.a.  $X_i$  et  $Y_j$  on peut prendre comme modèle

$$\bigotimes_1^n (\mathbb{R}, \mathcal{B}, \frac{1}{\sqrt{2\pi}\sqrt{2\sigma}} \exp \frac{(z-\theta)^2}{4\sigma^2}) \quad \text{où } \theta = \mu - \nu.$$

a) On connaît  $\sigma^2$ .

Prenons d'abord  $m = n = 1$  et étudions complètement le problème dans ce cas.

$$\begin{aligned} U &= X - Y \\ &\sim N(\mu - \nu, 2\sigma^2) \\ &\sim N(\theta, 2\sigma^2) \end{aligned}$$

donc

$$\begin{aligned} V &= \frac{U}{2\sigma^2} \\ &\sim N\left(\frac{\theta}{2\sigma^2}, 1\right). \end{aligned}$$

On a un test défini par la région de rejet  $D = \{\omega, |V(\omega)| > \rho\}$ .

Si  $\alpha$  est le niveau désiré, alors  $\alpha$  est déterminé par

$$P_{\theta=0}(|V| > \rho) = 2(1 - F(\rho))$$

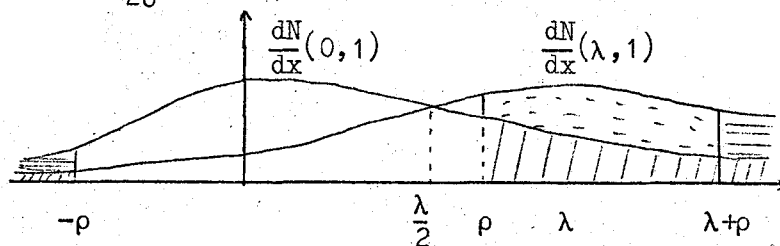
où F est la fonction de répartition tabulée d'une loi  $N(0, 1)$ .

Etudions la puissance de ce test. On a

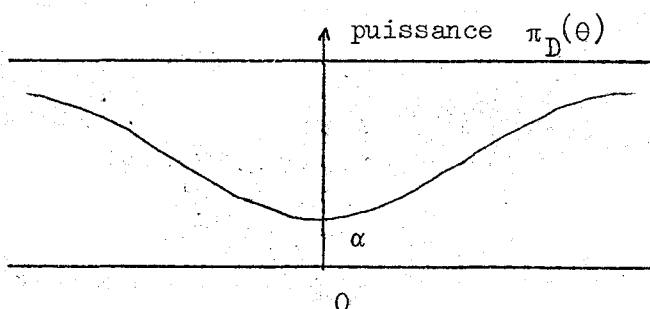
$$P_{\theta}(|V| > \rho) = P\left(|N + \frac{\theta}{2\sigma^2}| > \rho\right)$$

où  $N \sim N(0,1)$ .

Posant  $\lambda = \frac{\theta}{2\sigma^2}$ , on voit que la puissance vaut  $2 - F(\rho - \lambda) - F(\rho + \lambda)$



La puissance est représentée par l'aire en pointillé, c'est une fonction paire, croissante de  $|\theta|$



Remarquons que le test est sans biais.

Rappel : Définition. Un test est dit sans biais si  $P_{\theta}(D) > P_{\theta'}(D)$  pour tout  $\theta \in \theta_1$  ( $H_0$  fausse),  $\theta' \in \theta_0$  ( $H_0$  vraie)

Ici  $\theta_0 = \{0\}$ .

Supposons maintenant disposer pour  $X$  d'un  $n$ -échantillon  $(X_1 \dots X_n)$  et pour  $Y$  d'un  $m$ -échantillon  $(Y_1 \dots Y_m)$ . Posons comme d'habitude

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}, \quad \bar{Y} = \frac{Y_1 + \dots + Y_m}{m}$$

$$U = \bar{X} - \bar{Y} \sim N\left(\mu - \nu, \frac{\sigma^2}{n} + \frac{\sigma^2}{m}\right)$$

$$V = \frac{\bar{X} - \bar{Y}}{\sigma(\frac{1}{n} + \frac{1}{m})^{1/2}} \sim N(\lambda, 1)$$

avec  $\lambda = \frac{\mu - \nu}{\sigma(\frac{1}{n} + \frac{1}{m})^{1/2}}$ . La courbe de puissance précédente reste valable.

On voit comment la puissance  $\pi_D(\theta)$  croît quand  $n$  et  $m$  croissent ce qui est évidemment normal.



b) Supposons maintenant  $\sigma^2$  inconnu.

Dans le test à faire, on ne s'intéresse pas à  $\sigma^2$ , puisque celui-ci ne figure pas dans l'hypothèse. Il faut donc faire un test valable pour tout  $\sigma^2$ .

Si on pose  $s = \frac{1}{n+m-2} \left[ \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{j=1}^m (X_j - \bar{Y})^2 \right]$ , on a

$E_{(\mu, \nu, \sigma^2)} s = \sigma^2$ .  $s$  est donc un estimateur sans biais de  $\sigma^2$ .

$\bar{X} - \bar{Y}$  n'a pas une loi indépendante de  $\sigma^2$  puisque

$\bar{X} - \bar{Y} \sim N(\mu - \nu, (\frac{1}{n} + \frac{1}{m})\sigma^2)$ , mais

$$Z = \frac{\bar{X} - \bar{Y}}{\sqrt{s(\frac{1}{n} + \frac{1}{m})}} \quad \text{a une loi qui est}$$

la loi d'un quotient de 2 v.a. indépendantes d'après Cochran,  $\bar{X} - \bar{Y}$  et  $s$ .

Alors  $Z$  a une loi  $t'(n+m-2, \lambda)$ , soit  $t(n+m-2, \lambda)$  qui est tabulée.

Le test est alors construit.

Si  $H_0$  est vraie, soit  $\lambda = 0$ ,  $Z$  est centrée. Intuitivement, il y a donc une grande probabilité de trouver  $Z'$  près de 0 si  $H_0$  est vraie.

On prend tout pour  $D$  la région ( $|Z'| > \rho$ ).

Ce test est sans biais.  $\rho$  est déterminé par  $\alpha = P_{\theta=0}(|Z'| > \rho)$ .

(Si  $\lambda = 0$ , on a vu la loi de  $Z$  qui est  $t(n+m-2)$ ).

$Z^2$  est une variable  $F'_{1, n+m-2, (\lambda)^2}$

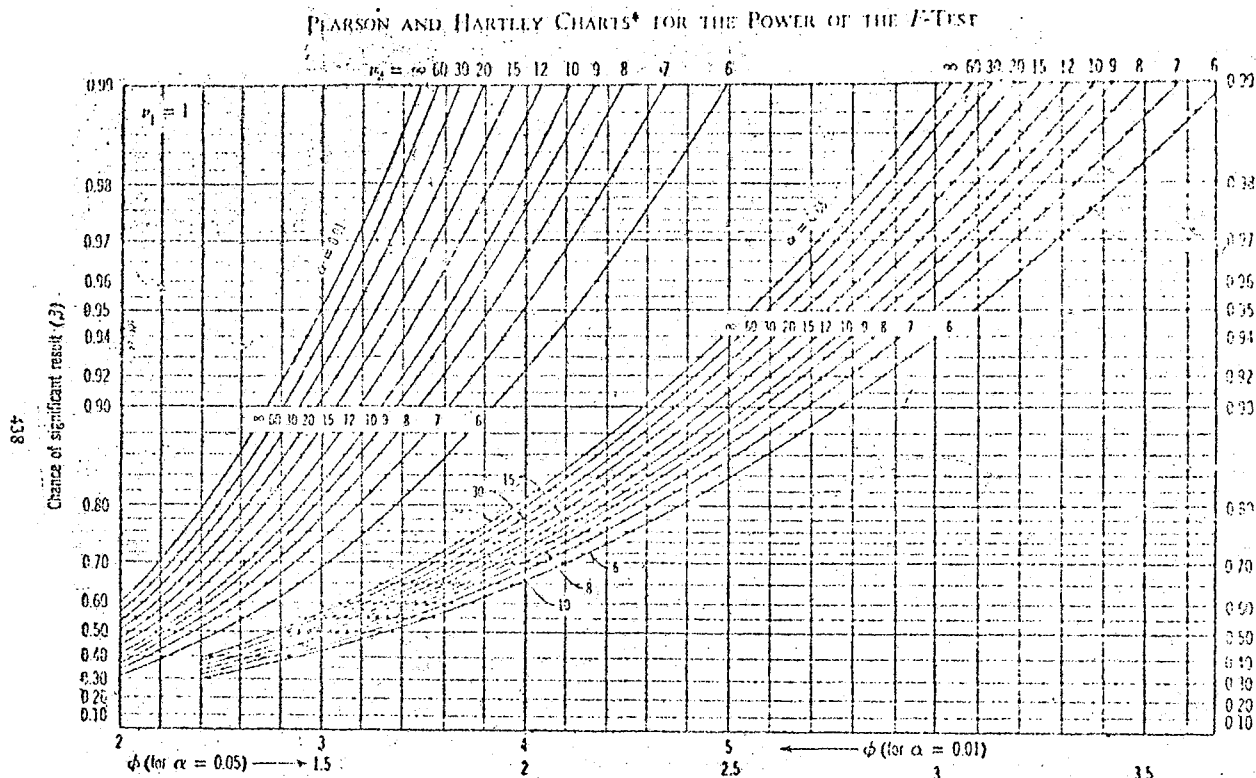
(quotient de deux  $\chi^2$  (un  $\chi^2(1, (\lambda)^2)$  est un  $\chi^2(n+m-2)$  centré), et

comme  $P(Z^2 > \rho^2)$  est une fonction strictement croissante de  $(\lambda)^2$  (voir plus haut), on a bien

$$P_0(Z^2 > \rho^2) < P_{(\mu-\nu)^2}(Z^2 > \rho^2) \quad \text{pour } \nu \neq \mu.$$

Remarque. Nous verrons plus loin les qualités de ce test (quant à l'optimalité).

Etude de la puissance du test ainsi obtenu



\* By E. S. Pearson and H. O. Hartley in *Biometrika*, Vol. 38, pp. 115-122 (1951). Reproduced with the kind permission of the authors and the editor.

Ces courbes donnent la puissance d'un test  $F$ ,  $P(F'_{1,v_2,\lambda} > \rho^2)$  pour des niveaux  $\alpha = 0,01$  et  $\alpha = 0,05$ , c.à.d. que  $\rho^2$  est calculé pour que le niveau soit le bon.

On a porté en abscisse  $\phi = \frac{\sqrt{\lambda}}{\sqrt{2}}$ , ( $v_2 = n+m-2$ ,  $\lambda = (\mu-v)^2$ ).

Supposons que la valeur de l'échantillon donne une valeur supérieure à  $\rho^2$  (pour  $Z'^2$ ) c.à.d. qu'il faudrait rejeter  $H_0$  : il est important de s'assurer que dans la région où on l'applique le test est assez puissant.

Exemple. On sait (par une estimation, par exemple, ou quelquefois par des expériences antérieures) que  $\Phi$  est de l'ordre de 2,5. Si  $v_2 = 6$  la puissance du test de niveau 0,01 est 0,5 c'est-à-dire que la probabilité d'avoir conclu juste est 0,5, ce n'est pas très sérieux. Il faut continuer l'étude. La situation n'est pas très probante. Le lecteur réfléchira à cette

situation courante qui explique l'importance de l'étude de la puissance.

Remarque. Si le degré de liberté du  $\chi^2$  du dénominateur augmente la puissance du test aussi : d'où un intérêt d'avoir de grands échantillons.

Le cas  $v_2 = \infty$  porté sur la figure correspond au carré d'une variable gaussienne, en effet  $\lim_n \frac{1}{n} \sum X_i^2 = 1$  (loi des grands nombres) on a donc aussi la courbe de puissance du test a).

Autre exemple de test :

Comparaison des variances de deux lois normales.

On dispose de deux échantillons indépendants

$(X_i)_{i=1 \dots n}$  de loi  $N(\mu, \gamma_1)$

$(Y_i)_{i=1 \dots m}$  de loi  $N(v, \gamma_2)$

alors

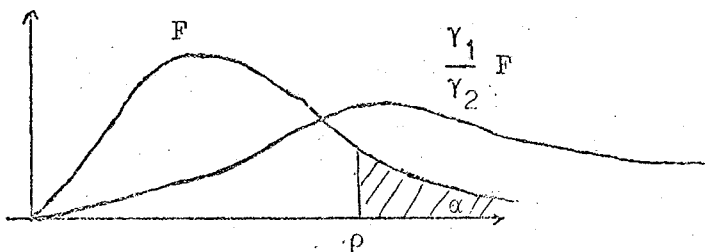
$$s_X = \sum_{i=1}^n (X_i - \bar{X})^2 \sim \gamma_1 \chi_{n-1}^2$$

$$s_Y = \sum_{j=1}^m (Y_j - \bar{Y})^2 \sim \gamma_2 \chi_{m-1}^2$$

donc

$$U = \frac{s_X/n-1}{s_Y/m-1} \sim \frac{\gamma_1}{\gamma_2} F_{n-1, m-1}$$

a) Supposons qu'on veuille tester  $\gamma_1 = \gamma_2$  contre  $\gamma_1 > \gamma_2$ . Les densités de  $F$  et  $\frac{\gamma_1}{\gamma_2} F$  ont l'allure suivante :



on propose donc le test

$$D = \text{rejet de } H_0 \text{ si } U > \rho$$

où  $\rho$  est calculé par  $P[F_{n-1, m-1} > \rho] = \alpha$  (niveau désiré).

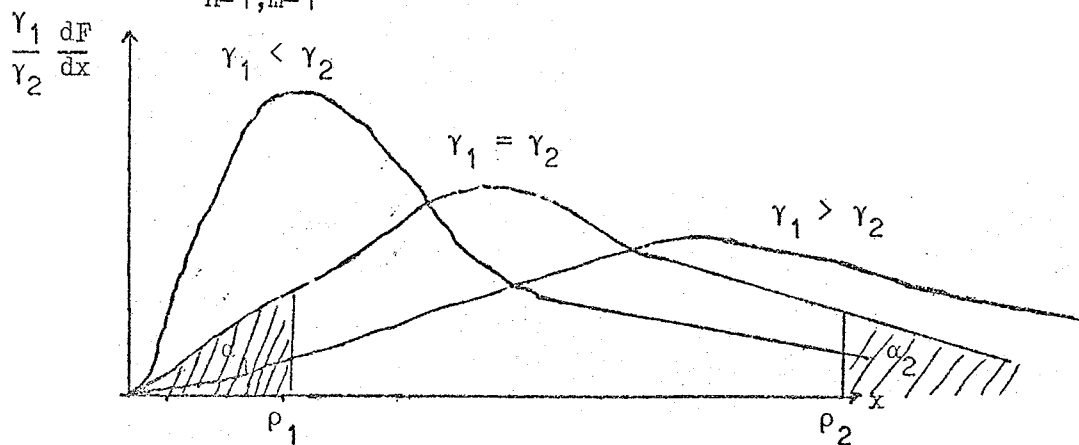
Puissance du test.

$$P[U > \rho] = P\left[\frac{\gamma_1}{\gamma_2} F_{n-1, m-1} > \rho\right] = P\left[F > \frac{\gamma_2}{\gamma_1} \rho\right]$$

si  $\frac{\gamma_2}{\gamma_1} < 1$  c.à.d.  $H_0$  fautive, la puissance est supérieure à  $\alpha$  et pour avoir

une idée de sa valeur on utilise une table donnant la fonction de répartition d'une variable  $F$ .

b) Supposons maintenant qu'on veuille tester  $\gamma_1 = \gamma_2$  contre  $\gamma_1 \neq \gamma_2$  il faut donc exclure les valeurs de  $U$  où les variables  $\frac{\gamma_1}{\gamma_2} F$  pour  $\gamma_1 \neq \gamma_2$  sont trop probables, c.à.d. les petites et les grandes valeurs de  $U$ . (nous écrivons  $F$  pour  $F_{n-1, m-1}$  dans la suite)



on va exclure  $H_0$  si  
 $U < \rho_1$  ou  $U > \rho_2$   
 $D = [U < \rho_1] \cup [U > \rho_2]$ .

Pour obtenir le niveau  $\alpha$  il faut que  
 $P[F < \rho_1] + P[F > \rho_2] = \alpha$

mais il reste une indétermination puisqu'il y a deux bornes :  
 on peut par exemple choisir

$$P[F < \rho_1] = \alpha/2 \text{ et } P[F > \rho_2] = \alpha/2$$

on peut aussi chercher un test sans biais c.à.d. tel que

$$P\left[\frac{1}{\lambda} F < \rho_1\right] + P\left[\frac{1}{\lambda} F > \rho_2\right]$$

soit minimum pour  $\lambda = 1$ . Si  $g$  est la densité d'une variable  $F$

$$P[F < \lambda \rho_2] + P[F > \lambda \rho_2] = \int_0^{\lambda \rho_2} g(x) dx + \int_{\lambda \rho_2}^{\infty} g(x) dx$$

la condition "sans biais" s'écrit :

$$\rho_1 g(\rho_1) - \rho_2 g(\rho_2) = 0$$

or

$$g(x) = \frac{\frac{p}{q} \cdot \left(\frac{p}{q} x\right)^{\frac{p}{2}-1}}{\beta\left(\frac{p}{2}, \frac{q}{2}\right) \left(1 + \frac{p}{q} x\right)^{\frac{p+q}{2}}}$$

pour une  $F_{p, q}$

d'où une résolution possible en  $\rho_1, \rho_2$ .

En fait ce test est intéressant si on garde l'hypothèse  $H_0$  et dans ce cas le caractère "sans biais" est secondaire. En effet, pour conclure au rejet, il faut encore s'assurer qu'on est dans une région où la puissance est raisonnable : c.à.d. que, ayant une idée du rapport  $\frac{\gamma_1}{\gamma_2}$  on calcule la puissance en ce point, mais dans ce cas le test unilatère précédent est plus justifié (puisque on a déjà une estimation de  $\frac{\gamma_1}{\gamma_2}$ ).

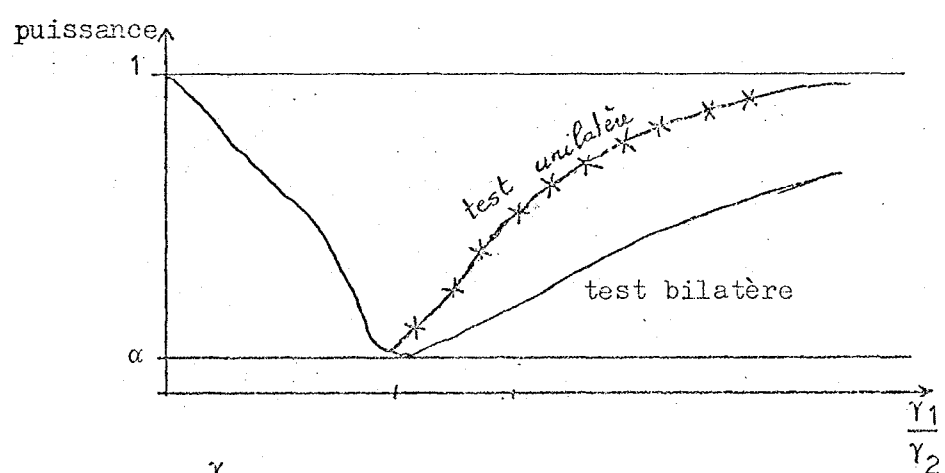
Remarque pour le calcul de  $\rho_1$  et  $\rho_2$ . Les tables donnent  $\rho$  tel que

$$P[F_{n,m} > \rho] = \alpha$$

pour des niveaux  $\alpha$  classiques. L'autre borne s'obtient en remarquant que

$$P[F_{n,m} < \rho] = P[F_{m,n} > \frac{1}{\rho}]$$

Nous n'avons fait qu'amorcer ici la discussion. On voit que le choix entre le test unilatère ( $\gamma_1 = \gamma_2$  contre  $\gamma_1 > \gamma_2$ ) et le test initial ( $\gamma_1 = \gamma_2$  contre  $\gamma_1 \neq \gamma_2$ ) est très technique et dépend du crédit que l'on accorde à l'estimation de  $\frac{\gamma_1}{\gamma_2}$ .



Pour les valeurs  $\frac{\gamma_1}{\gamma_2}$  un peu plus grande que 1 le choix est délicat et nous reviendrons dessus plus loin du point de vue théorique (concrètement bien que démarrant avec l'idée de faire un test bilatère il peut être mauvais de rester sur ce point de vue).

Les moyennes  $\mu$  et  $\nu$  interviennent ici comme paramètres fantômes, elles sont été éliminées d'entrée.

Nous n'aborderons pas dans le cours l'aspect théorique de l'élimination de ces paramètres fantômes, nous contentant de voir sur divers exemples comment D est fabriqué de manière à ce qu'ils ne figurent pas.

Intervalles de confiance.

Nous ne détaillerons pas ici la forme des régions de confiance possible. A l'aide des statistiques introduites pour les tests, le lecteur le fera aisément. Indiquons simplement deux petits résultats numériques fort utiles. Soit un  $n$ -échantillon de  $N(\mu, \sigma^2)$ ,  $\Phi$  la fonction de répartition de  $N(0,1)$ . On a  $\Phi^{-1}(0,025) = -1,96 \approx -2$

$$\Phi^{-1}(0,015) = -3.$$

Par exemple si  $\sigma$  est connu les intervalles de confiance pour  $\mu$  seront

a) bilatère

$$\bar{X} - \frac{\sigma}{\sqrt{n}} |\Phi^{-1}(\frac{\alpha}{2})|, \bar{X} + \frac{\sigma}{\sqrt{n}} |\Phi^{-1}(\frac{\alpha}{2})|$$

b) unilatère à droite

$$[\bar{X} + \frac{\sigma}{\sqrt{n}} |\Phi^{-1}(\alpha)|, \infty [.$$

APPENDICE ET RESUME

A) Notes sur quelques lois intervenant souvent à propos des problèmes sur les gaussiennes notamment.

La loi Gamma  $\Gamma(a, \lambda)$  a pour densité

$$\gamma_{a, \lambda}(x) = \frac{1}{\Gamma(a)} \lambda^a e^{-\lambda x} x^{a-1} 1_{\mathbb{R}^+}(x)$$

où

$$\Gamma(a) = \int_0^{\infty} e^{-x} x^{a-1} dx$$

$$\Gamma(a+1) = \Gamma(a), \quad \Gamma(n+1) = n!$$

Pour  $a = 1$ , on obtient la loi exponentielle ordinaire.

$$\text{On a} \quad \int_0^{\infty} e^{-tx} \gamma_{a, \lambda}(x) dx = \left(1 + \frac{t}{\lambda}\right)^{-a}$$

et

$$\Gamma(a, \lambda) * \Gamma(b, \lambda) = \Gamma(a+b, \lambda)$$

$$\Gamma\left(\frac{n}{2}, 1/2\right) = \chi^2(n).$$

Soient  $X$  et  $Y$  deux variables indépendantes  $X \sim \Gamma(a, \lambda)$ ,  $Y \sim \Gamma(b, \lambda)$ .

Alors  $X/Y$  a pour loi  $\beta(a, b)$  de densité

$$g_{a, b}(x) = \frac{\Gamma(a+b)}{\Gamma(a) \Gamma(b)} \frac{x^{a-1}}{(1+x)^{a+b}} 1_{\mathbb{R}^+}(x).$$

Si  $U \sim \beta\left(\frac{n_1}{2}, \frac{n_2}{2}\right)$ ,  $\frac{n_2}{n_1} U \sim F(n_1, n_2)$ .

B) RESUME. Statistique à utiliser :

sur la moyenne  $\mu$

a)  $\sigma$  connu  $\sqrt{n} \frac{\bar{X} - \mu}{\sigma} \sim N(0,1)$

b)  $\sigma$  inconnu  $\sqrt{n} \frac{\bar{X} - \mu}{s} \sim T(n-1)$ .

Statistique à utiliser sur  $\sigma$

a)  $\mu$  connu  $\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \sim \chi^2(n)$

b)  $\mu$  inconnu  $\frac{(n-1)s^2}{\sigma^2} \sim \chi^2(n-1)$ .

Pour 2 échantillons on a

$$\begin{aligned} \bar{X} - \bar{Y} &\sim N(\mu - \mu', \frac{\sigma^2}{n} + \frac{\sigma'^2}{n'}) \\ \frac{(n-1)s_X^2}{\sigma^2} - \frac{(n'-1)s_Y^2}{\sigma'^2} &\sim \chi^2(n+n'-2) \\ \frac{\sigma'^2}{\sigma^2} \cdot \frac{s_X^2}{s_Y^2} &\sim F(n-1, n'-1) \end{aligned}$$

et si  $\sigma^2 = \sigma'^2$

$$\frac{\bar{X} - \bar{Y} - (\mu - \mu')}{[(n-1)s_X^2 + (n'-1)s_Y^2]^{1/2}} \frac{\sqrt{n+n'-2}}{\sqrt{\frac{1}{n} + \frac{1}{n'}}} \sim T(n+n'-2).$$

CHAPITRE VI

LE TEST DU  $\chi^2$

Il s'agit d'une méthode non paramétrique. Nous l'avons présentée à part car elle n'est pas fondée sur des considérations portant sur les statistiques d'ordre, elle n'a de justification qu'asymptotique. Ceci étant c'est sûrement le test statistique le plus employé et son emploi correct présente de nombreuses difficultés. Nous expliquons dans ce chapitre le fondement du test (loi multinomiale) et les conditions d'applications les plus communes, sans démontrer les théorèmes intéressants, concernant les rapports test-estimation, quant au choix convenable du nombre de degrés de liberté.

1. PRELIMINAIRES. Convergence d'une loi multinomiale vers un  $\chi^2$ .

Soient  $p_1 \dots p_r$ ,  $r$  nombres strictement positifs tels que  $\sum_{i=1}^r p_i = 1$ . Soient  $X_j$  des v.a. indépendantes telles que  $P(X_j = i) = p_i$ ,  $j = 1 \dots n$ ;  $i = 1 \dots r$ . On pose

$$Z_i^n = \sum_{j=1}^n 1_{\{i\}}(X_j), \quad Z^n = (Z_1^n, \dots, Z_r^n).$$

On a :  $P(Z^n = (n_1, \dots, n_r)) = \frac{n!}{n_1! \dots n_r!} p_1^{n_1} \dots p_r^{n_r}$ , où  $n_1, \dots, n_r \in \mathbb{N}$ ,

$\sum_{i=1}^r n_i = n$ , par un raisonnement combinatoire élémentaire.  $Z^n$  est donc une v.a. portée par hyperplan  $\sum_{i=1}^r n_i = n$  de  $\mathbb{N}^r$ . La loi de  $Z^n$ , notée  $M_d$

$(n, p_1, \dots, p_r)$  est appelée loi multinomiale dégénérée. Si  $M_d(n_1, p_1, \dots, p_r)$

et  $M_d(n_2, p_1, \dots, p_r)$  sont les lois de 2 v.a. indépendantes  $Z_1^{n_1}$  et  $Z_2^{n_2}$

alors  $M_d(n_1+n_2, p_1, \dots, p_r)$  est la loi de  $Z_1^{n_1} + Z_2^{n_2}$ .

$$\begin{aligned} \text{Posons } Y^n &= (Y_1^n, \dots, Y_r^n)^t \\ &= \left( \frac{Z_1^n - np_1}{\sqrt{np_1}}, \dots, \frac{Z_r^n - np_r}{\sqrt{np_r}} \right) \end{aligned}$$

et  $\sqrt{P} = (\sqrt{p_1} \dots \sqrt{p_r})^t$ .



Posons  $V_j = \left( \frac{1}{\sqrt{p_i}} 1_{\{i\}}(X_j) \right)_{i=1 \dots r}$

Les  $V_j$  forment une suite de  $n$  vecteurs équadistribués indépendants. Comme

$$E\left(\frac{1}{\sqrt{p_h}} 1_{\{h\}}(X_j) - \sqrt{p_h}\right) \left(\frac{1}{\sqrt{p_k}} 1_{\{k\}}(X_j) - \sqrt{p_k}\right) = \delta_{hk} - \sqrt{p_h p_k},$$

on a  $\text{cov } V_j = I_r - {}^t \sqrt{p} \sqrt{p}$  où  $I_r$  est la matrice identité sur  $\mathbb{R}^r$ .

D'après le théorème de convergence vers une loi normale dans  $\mathbb{R}^r$  (cf. cours de Probabilités, chapitre 6) on a :

$$\frac{\sum_{j=1}^n (V_j - \sqrt{p})}{\sqrt{n}} \xrightarrow{\text{loi}} N(0, I_r - {}^t \sqrt{p} \sqrt{p})$$

et comme l'application  $\mathbb{R}^n \rightarrow \mathbb{R}$  définie par  $x \rightarrow \|x\|$  est continue, on a

$$\left\| \frac{\sum_{j=1}^n V_j - \sqrt{p}}{\sqrt{n}} \right\|^2 \xrightarrow{\text{loi}} \text{loi} \|N(0, I_r - {}^t \sqrt{p} \sqrt{p})\|^2$$

(si  $T$  est une fonction  $\mathbb{R}^n \rightarrow \mathbb{R}^k$  continue et si  $X \xrightarrow{\text{loi}} Y$  alors

$T(X) \xrightarrow{\text{loi}} T(Y)$  puisque  $\int f \circ (T \circ X) = \int (f \circ T)(X)$  et que  $f \circ T \in C(\mathbb{R}^n)$  si  $f$  est continue, bornée de  $\mathbb{R}^k$  dans  $\mathbb{R}$ ).

Il reste à calculer la loi du carré de la norme d'une v.a.  $Z \sim N(0, I_r - {}^t \sqrt{p} \sqrt{p})$ . La norme est invariante par un changement de base orthonormée.

Soit  $B$  une telle matrice orthonormée du type  $(\sqrt{p}/c)$ , ( $\sqrt{p}$  est acceptable pour une première colonne puisque  $\|\sqrt{p}\|^2 = 1$ ). On a  $BZ \sim N(0, B(I_r - {}^t \sqrt{p} \sqrt{p}) {}^t B)$

et  $B(I_r - {}^t \sqrt{p} \sqrt{p}) {}^t B = B {}^t B - (\sqrt{p}/c) {}^t \sqrt{p} \sqrt{p} {}^t (\sqrt{p}/c)$  or  $(\sqrt{p}/c) \begin{pmatrix} \sqrt{p} \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & I_{r-1} \end{pmatrix}$  puisque  $B$  est orthogonale. Soit  $c {}^t \sqrt{p} = 0$  et donc  $(\sqrt{p}/c) \begin{pmatrix} \sqrt{p} \\ 0 \\ \vdots \\ 0 \end{pmatrix} = (\sqrt{p}/c)$

$$\begin{aligned} \text{et donc } B(I_r - {}^t \sqrt{p} \sqrt{p}) {}^t B &= \begin{pmatrix} 1 & 0 \\ 0 & I_{r-1} \end{pmatrix} - \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \\ &= \begin{pmatrix} 0 & 0 \\ 0 & I_{r-1} \end{pmatrix} \end{aligned}$$

donc la loi de  $\|Z\|^2$  est un  $\chi_{r-1}^2$ .

## 2. Mesure de la distance de deux lois. (Cf. aussi Chap. IV test de Kolmogorov et Smirnov).

Commençons par rappeler la démonstration du lemme suivant :

Soit un  $n$ -échantillon  $X_1, \dots, X_n$  de v.a. à valeurs dans  $\mathbb{R}$  (ou  $\mathbb{R}^n$ ).

Si  $A$  est un borélien poscons  $F_n(A) = \frac{n_A}{n}$  où  $n_A = \text{nombre } (X_j \in A, 1 \leq j \leq n)$ .

Alors  $F_n \rightarrow F$ , où  $F$  est la loi de  $X$ .

$F_n$  est en effet (c'est trivial) une probabilité sur  $\mathbb{R}$  (dite répartition ou loi empirique).

$$\text{On a } n_A = \sum_{j=1}^n 1_A(X_j).$$

Les v.a.  $1_A(X_j)$  sont indépendantes, de même loi, et  $E 1_A(X_j) = F(A)$ .

La loi forte des grands nombres assure alors que  $\frac{n_A}{n} \rightarrow F(A)$  p.s. d'où le lemme.

Il y a de très nombreuses manières de mesurer la distance de 2 lois  $F_n$  et  $F$ . Les mesures théoriques classiques sont intéressantes. Dans ce chapitre nous définirons des "distances"  $d(F_n, F)$  avec l'idée suivante. Les lois  $F_n$  que nous considérons convergent sous une hypothèse convenable  $H_0$  vers la loi  $F$ . Dans un problème d'estimation on choisira  $F$  (loi à estimer) de manière que  $d(F_n, F)$  soit minimale,  $F_n$  étant la loi empirique associée à l'échantillon. Dans un problème de test,  $F$  sera fixée, si  $H_0$  hypothèse à tester est vraie, alors on aura  $d(F_n, F)$  petite, si  $H_0$  est fausse on aura  $d(F_n, F)$  grande. Ceci étant dit, les "distances" élémentaires (qui ne sont pas des vraies distances) sont fondées sur la loi multinomiale. Soit  $(p_i)_{i=1 \dots r}$  une loi à  $r$  valeurs  $a_i$ , et soit  $X_1 \dots X_n$  un échantillon prenant ces mêmes  $r$  valeurs  $a_i$ . Soit  $n_i = \text{nombre } (X_j = a_i)$ . La loi empirique est  $(\frac{n_i}{n})_{i=1 \dots r}$ . On mesure sa distance à la loi  $(p_i)_{i=1 \dots r}$  par

$$\begin{aligned} \chi^2 &= \sum_{i=1}^r \frac{(\frac{n_i}{n} - p_i)^2 n}{p_i} \\ &= \sum_{i=1}^r \frac{(n_i - np_i)^2}{np_i} \end{aligned}$$

ou par

$$\begin{aligned} d_{\text{Hellinger}} &= \cos^{-1} \sum_{i=1}^r \frac{n_i}{n} p_i \\ d_{\text{Kullback}} &= \sum p_i \log \left[ \frac{p_i}{\frac{n_i}{n}} \right] \text{ etc... } \end{aligned}$$

Pour estimer les nombres  $p_i$  qui caractérisent la vraie loi, on décide (arbitrairement mais on peut justifier asymptotiquement) de minimiser

ces distances. Pour ce faire, on suppose que  $p_i$  est une fonction bien dérivable d'un paramètre  $\theta$  variant dans un bon "ouvert", du type

$$(\theta < x_i < 1, \sum_{i=1}^r x_i = 1).$$

Par exemple dans le cas du  $\chi^2$ , on a en dérivant par rapport à  $p_i$ ,

$$\frac{\partial \chi^2}{\partial p_i} = 0 \implies \hat{p}_i = \frac{n_i}{n}.$$

L'estimateur obtenu  $\frac{n_i}{n} \xrightarrow{(n \rightarrow \infty)} p_i$  vraie valeur, la suite d'estimateurs fabriquée ainsi est donc consistante. Nous laissons au lecteur le soin d'étudier les autres cas en précisant les hypothèses.

Si on a des lois  $F$  quelconques on se ramène aux techniques précédentes en approximant  $F$  par une loi à  $r$  valeurs, en choisissant une partition  $A_1 \dots A_r$  de  $\mathbb{R}$ . On estime  $p_i = \int_{A_i} dF_i$  en le comparant à  $\frac{n_i}{n}$ , où  $n_i = \sum_{j=1}^n 1_{A_i}(X_j)$ .

### 3. Les tests du $\chi^2$ .

#### a) Test sur la loi multinomiale.

Soit  $X_1 \dots X_n$  un  $n$ -échantillon ( $n$  grand) d'une v.a. à  $r$ -valeurs de loi  $(\pi_1(\theta), \dots, \pi_r(\theta))$ , où  $\theta \in \Theta$ . Soit  $p_1 = \pi_1(\theta_0), \dots, p_r = \pi_r(\theta_0)$ .  $0 < p_i < 1$ . Soit à tester  $H_0$  :

$$" \pi_1(\theta) = p_1, \dots, \pi_r(\theta) = p_r " \text{ contre } H_1 :$$

$$" \pi_1(\theta) \neq p_1 \text{ ou } \dots \text{ ou } \pi_r(\theta) \neq p_r ".$$

On sait, en reprenant les notations du § 1, que d'après la loi p.s.

des grands nombres, on a

$$\frac{Z^n}{n} \rightarrow (\pi_1(\theta), \dots, \pi_r(\theta)) P_\theta \text{ p.s.}$$

si  $P_\theta = \otimes^n (\pi_1(\theta), \dots, \pi_r(\theta))$ .

Donc  $\frac{\|Z^n - n.p\|^2}{n} = n \left\| \frac{Z^n}{n} - p \right\|^2$

$$\sim c^2 n \text{ avec } c^2 = \sum |\pi_1(\theta) - p_i|^2 \\ = \|\pi(\theta) - p\|^2 \text{ si } \pi(\theta) \neq p.$$

Posons toujours  $Y^n = (Y_1^n \dots Y_r^n)$

avec 
$$Y_i^n = \frac{Z_i^n - np_i}{\sqrt{n} \sqrt{p_i}} .$$

Le raisonnement ci-dessus montre que si  $\pi(\theta) \neq p$ , donc  $\|Y^n\|^2 \sim C'n \rightarrow \infty$  si  $H_0$  est fautive et on a vu, que si  $\pi(\theta) = p$ , on avait par contre  $\|Y^n\|^2 \rightarrow \chi^2(r-1)$ .

On en déduit un test. On commence par assimiler pour  $n$  grand la distribution de  $Y^n$  (dans le cas où  $H_0$  est vraie) à celle d'un  $\chi^2(r-1)$ .

Soit  $\alpha$  le niveau donné. Définissons  $d$ , par  $P(\chi^2(r-1) \geq d) = \alpha$ . Le test est défini par  $D = (\|Y^n\|^2 \geq d)$ . La puissance d'un tel test est évidemment très compliquée à calculer puisqu'elle dépend de  $\pi_1(\theta), \dots, \pi_r(\theta)$  et aussi un peu de  $n$ . Pour cette raison, on est amené à introduire dans l'étude des tests asymptotiques de ce type (c.à.d. de test où la loi de la statistique utilisée est confondue avec la loi limite pour  $n = \infty$  dans le cas où l'hypothèse est vraie) d'autres notions plus adaptées comme celle d'efficacité.

Nous ne développerons pas dans ce cours ces notions, nous contentant d'exposer le pourquoi et la méthode des tests les plus importants, sans rechercher leurs qualités dans l'ensemble des tests asymptotiques. Les  $\chi^2$  sont tabulés.

#### b) Tests d'ajustements.

Soit  $F$  une loi quelconque sur un ensemble  $E$ , on a un  $n$ -échantillon  $X_1, \dots, X_n$ .

Problème : tester  $H_0$  "la loi de  $X$  est  $F$ " contre  $H_1$  "la loi de  $X$  n'est pas  $F$ ".

Un tel problème est dit test d'ajustement (on essaie de voir si la loi  $F$  s'ajuste bien à la répartition des  $X_j$ ).

A cet effet, on fait une partition de  $E$  en  $A_1, \dots, A_r$  et on con-

sidère les variables  $X_{j,i} = 1_{\{i\}} X_j$ ,  $Z_i^n = \sum_{j=1}^n X_{j,i}$ ,  $Z^n = (Z_1^n, \dots, Z_r^n)$ .

$Z^n$  a une loi multinomiale  $M_d(p_1 \dots p_r)$  avec  $p_i = P(X_j \in A_i)$ . Le principe du test d'ajustement est alors le suivant : remplacer l'hypothèse  $H_0$  par l'hypothèse  $H'_0$  (plus faible),  $Z^n$  a une loi multinomiale  $M_d(p_1 \dots p_r)$ .

(Si  $H_0$  est vraie,  $H'_0$  est vrai a fortiori, l'implication contraire étant fausse). On teste alors  $H'_0$  comme au paragraphe précédent contre  $H_1 : Z^n$  a une loi  $M_d(p'_1 \dots p'_r)$ ,  $p' \neq p$ .

Ce test d'ajustement très simple dépend évidemment beaucoup du partage  $A_1 \dots A_r$  (on prend en général  $r \geq 7$  quand c'est possible).

c) Comparaison de plusieurs échantillons de lois multinomiales.

Tests d'homogénéité.

Soient  $m$  échantillons  $(X_{1,i}, \dots, X_{n_i,i})$   $i = 1 \dots m$ , de lois multinomiales associées à  $m$  populations. Problème : ces échantillons proviennent-ils d'une même population ?

Nous allons étudier le problème pour  $m = 2$  sur l'exemple de groupes sanguins relevant de 2 populations (d'après Rao). On a les fréquences des groupes O, A, B, AB.

Soit  $p, q$  la fréquence des gènes A et B pour la première population,  $p', q'$  les mêmes fréquences pour la deuxième population. A-t-on  $p = p'$  ? On a

	O	A	B	AB	TOTAL
1 <sup>ère</sup> population	121	120	79	33	353
2 <sup>ème</sup> population	118	95	121	30	364

Les classes sont toutes faites. Le problème se pose donc directement comme problème sur les multinomiales.

Posant  $\pi_{i,k}$ ,  $i = 1, 2$ ,  $k = 1, 2, 3, 4$

les probabilités pour un individu de la population  $i$  d'appartenir au groupe  $k$ , il faut tester si  $\pi_{i,k} = \pi_{i',k}$  pour tout  $k$ . On estime  $\pi_{i,k}$  par  $\frac{n_{i,k}}{n_i}$  où  $n_i$  est le nombre total d'individus de la population  $i$ , et  $n_{i,k}$

le nombre de ceux du groupe  $k$ .

Supposons vérifiée l'hypothèse  $H_0$  : "homogénéité des populations".

On a alors pour la population globale  $n_k = n_{1,k} + n_{2,k}$  individus dans la classe  $k$  et l'on estime la probabilité  $\pi_k$  d'être dans la classe  $k$  par  $\frac{n_k}{n}$ . On calcule alors la somme des distances  $\chi^2$  (des 2 populations), distances à la répartition estimée sous  $H_0$ , soit

$$\begin{aligned}\chi^2 &= \sum_{i=1}^2 \sum_k \frac{\left(\frac{n_{i,k}}{n_i} - \frac{n_k}{n}\right)^2 n_i}{\frac{n_k}{n}} \\ &= \sum_i \sum_k \frac{(n_{i,k} - \frac{n_i n_k}{n})^2}{n_i n_k}\end{aligned}$$

Si l'on avait pas fait d'estimation c'est-à-dire si la loi était une loi a priori on aurait la somme de deux  $\chi^2$  ( $k-1$ ) indépendants soit un  $\chi^2$  ( $2k-2$ ). Mais, on démontre que chaque estimation d'un paramètre fait perdre un degré de liberté du  $\chi^2$  (en fait puisque l'on estime la loi a priori à partir de l'échantillon, on se place intuitivement près de la loi multinomiale la plus voisine). Ceci est valable à condition que les paramètres estimés soient fonctionnellement indépendants, c'est-à-dire qu'il n'existe une fonction (implicite) des paramètres identiquement nulle. Ici  $\sum_k \pi_k = 1$  et on a donc estimé  $k-1$  paramètres, il reste donc  $2(k-1) - (k-1)$  degrés de liberté. Si l'on avait comparé entre elles  $h$  populations le nombre de degré de liberté serait donc  $h(k-1) - (k-1) = (h-1)(k-1)$ . Dans l'exemple concret  $\chi^2 = 11,73$ , ce qui pour 3 degrés de liberté amène au rejet de l'hypothèse au niveau 5 %, paramètre  $\pi_k$ . Le degré de liberté du  $\chi^2$  est donc  $2(k-1) - k-1 = k-1$ . Dans l'exemple étudié, on a un  $\chi^2(3)$  qui vaut 11,73 et qui amène à rejeter l'hypothèse d'homogénéité des populations au niveau 5 %.

Remarquons que dans certains problèmes on considère une des populations comme population test, et donc connue, non interprétée comme échantillon. On regarde alors si une autre population s'ajuste à celle-ci et l'on obtient un test sur une loi fixée, sans estimation donc à  $(k-1)$  degrés de

liberté pour 2 populations et à  $(h-1)(k-1)$  degré de liberté pour  $h$  population dont une fixe ! Mais alors l'hypothèse testée ne compare pas entre elles les  $(h-1)$  populations, elle les compare séparément à  $H_0$ .

d) Test d'indépendance et tables de contingence.

On se propose de tester l'indépendance de 2 caractères constatés ou mesurés sur une même population par exemple la forme du pollen et la couleur de certaines fleurs obtenues par croisement (Bateson).

Soient  $F$  et  $C$  ces deux caractères  $i$  ( $i = 1, 2$ ) ils peuvent prendre les valeurs  $k$  ( $k = 1, 2$  long et rond) et  $l$  ( $l = 1, 2$  violet, rouge). On a (table de contingence).

Forme pollen	violet	couleur rouge	TOTAL
long	296	27	323
rond	19	85	104
Total	315	112	427

Problème : tester l'hypothèse  $H_0$  :  $F$  et  $C$  sont indépendants.

Notons  $n_{i,k}$  le nombre d'individus ayant les valeurs  $k, l$  de caractères.

On pose  $n_l = \sum_k n_{l,k}$ ,  $n_k = \sum_l n_{l,k}$ . On a  $\sum_l n_l = n$ ,  $\sum_k n_k = n$ .

Supposons que la probabilité inconnue d'être dans la classe  $(l,k)$  soit  $\pi_{l,k}$ . La probabilité d'obtenir l'échantillonnage donné (loi multinomiale) vaut

$$n! \prod_l \prod_k \frac{\pi_{l,k}^{n_{l,k}}}{n_{l,k}!} = p_{l,k}$$

L'indépendance se traduit par  $\pi_{l,k} = \pi_l \pi_k$  où  $\pi_l$  est la probabilité de l'état  $l$ ,

$$\pi_l = \sum_k \pi_{l,k}, \quad \pi_k = \sum_l \pi_{l,k}$$

$$\sum_l \pi_l = 1, \quad \sum_k \pi_k = 1.$$

On peut écrire  $p_{l,k}$  sous la forme

$$p_{l,k} = n! \prod_l \frac{\pi_l^{n_l}}{n_l!} \prod_k \frac{\pi_k^{n_k}}{n_k!} \frac{\prod_l n_l! n_k!}{n!} \\ \times \prod \frac{1}{n_{l,k}!} \left( \frac{\pi_{l,k}}{\pi_l \pi_k} \right)^{n_{l,k}}$$

S'il y a indépendance, on obtient simplement l'expression

$$p_{l,k} = n! \prod_l \frac{\pi_l^{n_l}}{n_l!} \prod_k \frac{\pi_k^{n_k}}{n_k!} \frac{\prod_l n_l! n_k! / n_{l,k}!}{n!}$$

La quantité  $\frac{\prod n_{\ell}! n_k!}{n! n_{\ell k}!}$  représente exactement (et pour toute loi multinomiale) la loi conditionnelle des  $n_{\ell k}$  donnée les lois marginales  $n_{\ell}$  et  $n_k$ .

Si l'on connaît à l'avance les  $\pi_k$  et les  $\pi_{\ell}$  on fera le test d'indépendance suivant. On calcule

$$\chi^2 = \sum_k \sum_{\ell} \frac{(n_{k\ell} - n\pi_k\pi_{\ell})^2}{n\pi_k\pi_{\ell}}$$

et on fait simplement un test d'ajustement sur la loi  $\pi_k\pi_{\ell}$ . Donc si  $k$  peut prendre les valeurs  $1 \dots m$  et  $\ell$  les valeurs  $1 \dots m'$ , le nombre de degré de liberté à choisir est  $mm' - 1$  pour faire le test. Mais on peut améliorer la méthode de la manière suivante. On peut d'abord isoler les  $\chi^2$  qui correspondent aux déviations de l'échantillon par rapport aux marginales.

Soit

$$\chi_1^2 = \sum_{\ell} \frac{(n_{\ell} - n\pi_{\ell})^2}{n\pi_{\ell}}, \quad \chi_2^2 = \sum_k \frac{(n_k - n\pi_k)^2}{n\pi_k}$$

et pour  $\chi_3^2 = \chi^2 - \chi_1^2 - \chi_2^2$ . Comme  $\chi_3^2$  a  $mm' - 1$  degré de liberté

$\chi_1^2$  a  $m - 1$  et  $\chi_2^2$  a  $m' - 1$ ,  $\chi_3^2$  en a  $mm' - 1 - m + 1 - m' + 1 = (m - 1)(m' - 1)$ .

$\chi_3^2$  mesure la déviation (la distance) de l'échantillon à l'indépendance, une fois corrigées les déviations normales aux marginales et c'est donc  $\chi_3^2$  que l'on utilisera pour faire le test.

Remarque : on raisonne quelquefois en disant que  $\sum \pi_i = 1$ ,  $\sum \pi_k = 1$  et donc que le nombre de  $(\pi_i, \pi_k)$  intervenant dans la loi multinomiale à tester est  $(r - 1)(s - 1)$ , mais ce raisonnement n'est pas très clair.

Dans l'exemple précédent si il y a indépendance des caractères et si les rapports dans chaque classe valent 3 on obtient

$$\begin{aligned} \chi_1^2 &= 0,0945 && 1^{\text{e}} \text{ degré de liberté} \\ \chi_2^2 &= 0,3443 && \text{ " " " } \\ \chi_3^2 &= 221,6833 && \text{ pour 1 degré de liberté.} \end{aligned}$$

Donc en fait seul  $\chi_3^2$  importe.

Supposons maintenant que les marginales  $\pi_k$ ,  $\pi_{\ell}$  ne soient pas connues. On veut les estimer, par exemple par la méthode de maximum de vraisemblance

et l'on trouve  $\hat{\pi}_{\ell} = \frac{n_{\ell}}{n}$ ,  $\hat{\pi}_k = \frac{n_k}{n}$

$$\chi^2 = \sum_k \sum_{\ell} \frac{(n_{i,j} - n \frac{n_{\ell} n_k}{n})^2}{\frac{n_{\ell} n_k}{n}}$$



Quel degré de liberté faut-il choisir pour faire le test. Nous avons estimé  $(m-1)+(m'-1)$  paramètres  $\pi_{\ell}, \pi_k$ . Le nombre total de paramètres est  $mm'-1$ . L'estimation faite conduit à diminuer d'autant le degré de liberté du  $\chi^2$ . Intuitivement parce que les paramètres inconnus  $\pi_{k,\ell}$  sont liés par les relations  $\sum_{\ell} \pi_{k\ell} = \hat{\pi}_k$  et  $\sum_k \pi_{k\ell} = \hat{\pi}_{\ell}$  une fois l'estimation faite et donc que le nombre de paramètres à ajuster après estimation est  $mm'-1 - (m-1) - (m'-1)$ , (la dernière des relations  $\sum_{\ell} \pi_{k\ell} = \hat{\pi}_k$  se déduit des autres puisque  $\sum_k \sum_{\ell} \pi_{k\ell} = \sum_k \hat{\pi}_k = 1$ ).

Ce raisonnement intuitif peut être rendu rigoureux en utilisant le résultat énoncé plus haut. Nous laisserons cette démonstration de côté, en insistant sur le fait que dans tous les cas un test d'indépendance amène un  $\chi^2$  à  $(m-1)(m'-1)$  degrés de liberté.

Dans le cas numérique étudié, on trouve sur  $\pi_k, \pi_{\ell}$  ne sont pas spécifiés  $\chi^2 = 218,8722$ , ce qui est à peu près la valeur obtenue en utilisant des marginales connues à l'avance.

De toute manière, il faut bien voir la valeur relative de ces tests. Si la taille du nombre d'individus dans une des cases  $\pi_{ik}$  est petite cette case peut fausser le  $\chi^2$  et il faut faire appel à des méthodes plus fines mais encore basées sur le  $\chi^2$ .



