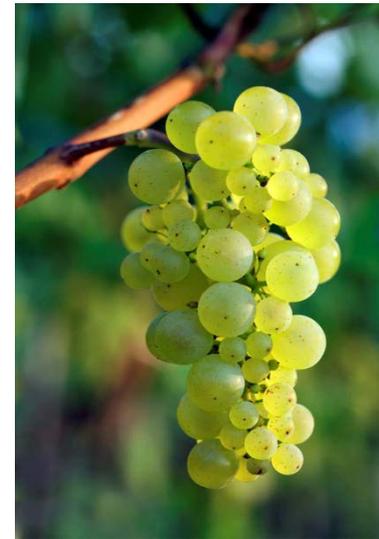


# Sondages à plusieurs degrés et par grappes



Philippe Périé, novembre 2011,  
et Gilbert Saporta, novembre 2013

# Sondages à plusieurs degrés et par grappes

- Introduction
- Sondages à plusieurs degrés
  - Tirage des unités primaires à probabilités égales sans remise (PESR)
  - Tirage des unités primaires à probabilités inégales avec remise (PIAR)
- Sondage par grappes
  - Taille des grappes connues a priori
  - Taille des grappes non connues a priori
- Grappes et stratification - mise en œuvre efficace

Sondages à plusieurs degrés et par grappes

# **INTRODUCTION**

# Sondages à plusieurs degrés

Les sondages à plusieurs degrés utilisent une succession de regroupements des unités statistiques pour tirer l'échantillon.

Par exemple :

- Tirer un échantillon de villes,
- Tirer un échantillon d'ilots
- Tirer un échantillon de ménages (logements) dans ces ilots

On a ici un exemple de sondage à 3 degrés, mais on peut généraliser à 2,3,4 degrés ...

Pour chaque degré, les méthodes déjà présentées (probabilités égales ou inégales, stratification, ...) peuvent s'appliquer.

# Dans quels cas ?

Les plans de sondage en plusieurs degrés visent le plus souvent à améliorer l'organisation de l'enquête ou sa réalisation économique

Dans la pratique, il arrive de ne pas avoir à disposition une base de sondage complète et disponible au moment où l'on planifie l'étude. Dans la plupart des cas, on peut n'avoir simplement qu'un degré (les villes, les quartiers, ...).

Ainsi, souvent ce type de sondage se rencontre quand les degrés constituent des unités géographiques et dans les enquêtes dont le mode de collecte est le face à face, car il y a un intérêt économique à limiter les déplacements autour de 'points de chute' définis.

On réalise dans ce dernier cas des économies de temps et de frais de déplacement. C'est moins vrai au téléphone, en online ou en postal : la dispersion des unités ne crée pas vraiment de coût

# Exemple

On veut interroger 5000 ménages en France métropolitaine qui en comporte 27 millions répartis sur plus de 36000 communes.

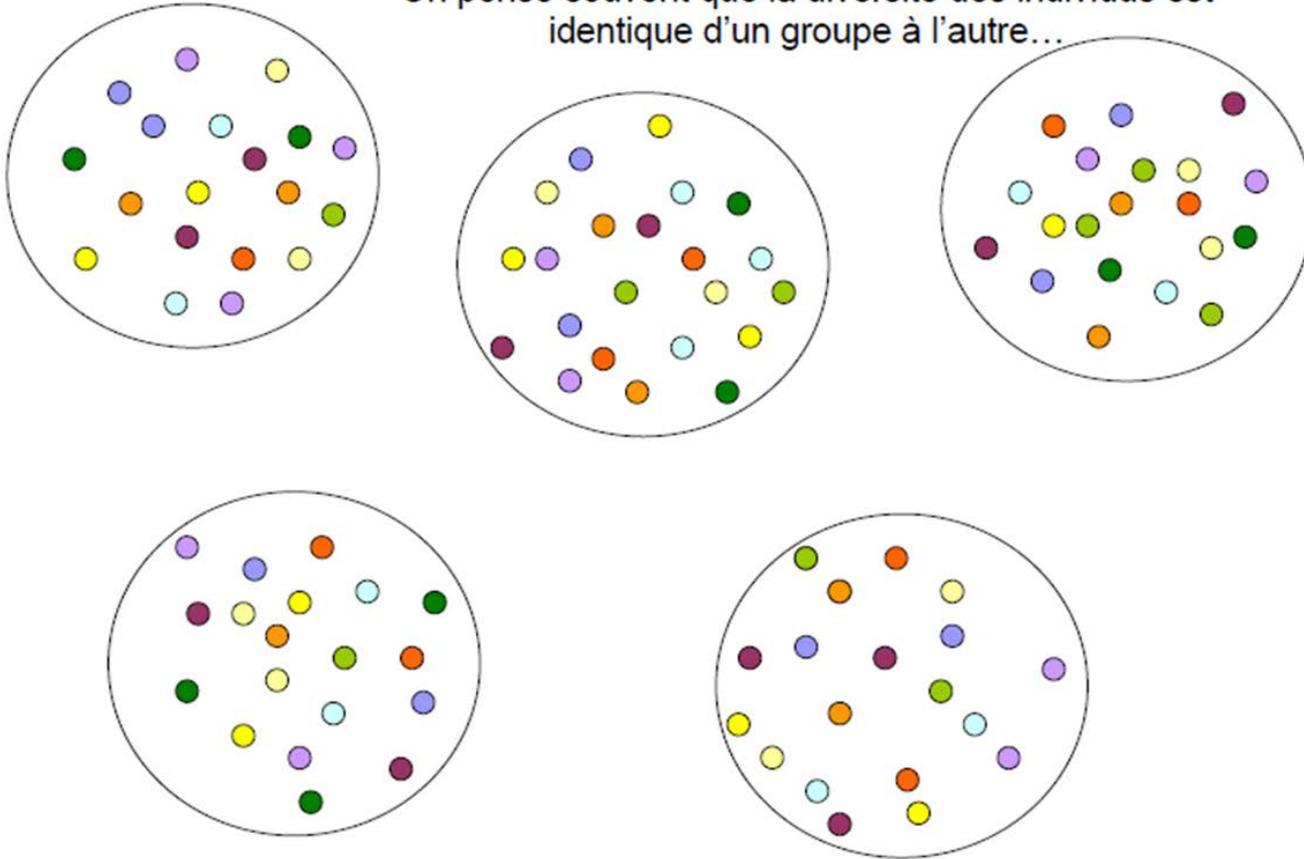
La liste des communes et des ilots, avec leurs caractéristiques est disponible à partir des enquêtes de recensement. Par contre il serait prohibitif en termes de temps et de coûts de vouloir constituer une liste exhaustive de ménages avant de lancer l'enquête, et d'y envoyer les enquêteurs au hasard.

Un sondage à plusieurs degrés permet de réaliser cette enquête avec une base de sondage exhaustive seulement au premier degré (ville) et au deuxième degré (ilot).

Une fois tiré les ilots, on envoie l'enquêteur interroger tout ou partie des logements. Si on interroge tous les ménages de chaque ilot tiré et non pas une sélection, on parle de grappes de ménages.

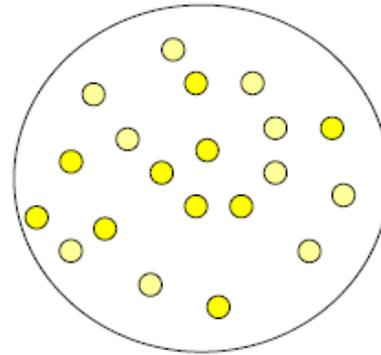
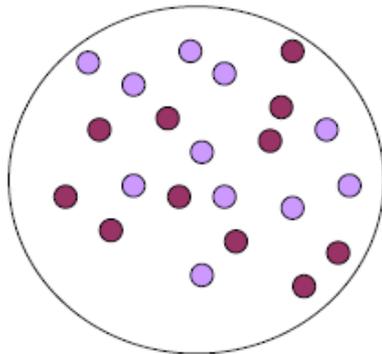
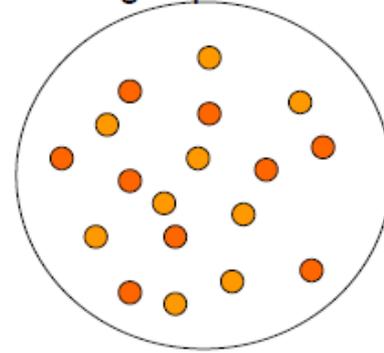
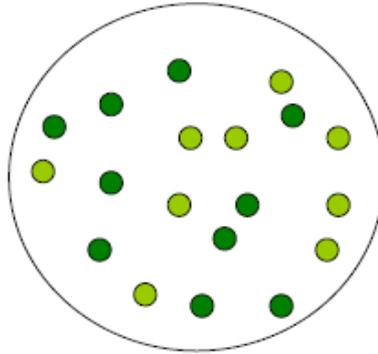
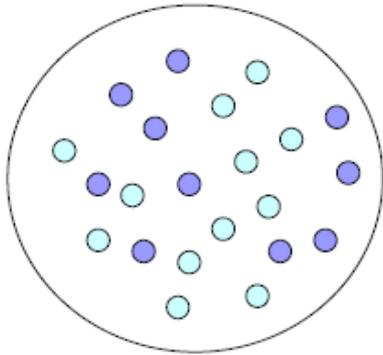
# Les limites ...

On pense souvent que la diversité des individus est identique d'un groupe à l'autre...



# Les limites ...

En réalité, les individus au sein d'un même groupe se ressemblent plus que s'ils viennent de groupes différents



# Limites ...

Un sondage à plusieurs degrés sera par contre moins précis qu'un sondage aléatoire simple par exemple pour une même taille d'échantillon : on parle d'effet de grappe (cluster effect)

Une idée intuitive est que l'on '**disperse moins**' l'échantillon : les unités regroupées dans un même groupe (une grappe) ont une certaine tendance à se ressembler (penser aux habitants d'un immeuble par exemple). Il y a donc une certaine redondance d'information : chaque unité supplémentaire d'une grappe apporte 'moins' qu'une unité tirée au hasard dans l'ensemble de la population

La plus grande partie de la variance dans le cas des tirages à plusieurs degrés vient souvent des premiers degrés. A la limite, si toutes les unités se ressemblaient parfaitement dans une grappe, alors c'est comme si l'on avait interrogé un échantillon non pas de individus mais de grappes.

# Sondages en grappes vs sondages à plusieurs degrés

Les  $N$  unités de la population sont réparties en  $M$  sous ensembles, appelées unités primaires (ou grappes)

La grappe  $\alpha$  ( $\alpha = 1, \dots, M$ ) contient  $N_\alpha$  unités de la population appelées unités secondaires

Lors d'un sondage en grappes, on prend un échantillon de  $m$  grappes, la grappe  $i$  ( $i = 1, \dots, m$ ) de l'échantillon est complètement enquêtée

# Sondages en grappes vs sondages à plusieurs degrés

Le sondage par grappes est un cas particulier de sondage à plusieurs degrés dans lequel l'ensemble des unités du dernier degré est enquêtée

# Exemples

Etudes médicales 'cas patients' : un échantillon de médecins qui donnent tout ou partie de leur patientèle, un effet de grappe médecin.

Etudes pour suivre certaines épidémies : grappes de laboratoires

INSEE : enquête emploi en continu

<http://www.insee.fr/fr/methodes/default.asp?page=sources/ope-enq-emploi-continu.htm>

*« L'échantillon est aréolaire. Les grappes ont été constituées à partir des informations collectées à l'occasion de la campagne 2006 de la taxe d'habitation. Chaque année, cette base de tirage est complétée par les logements nouveaux repérés dans les fichiers de la taxe d'habitation. La taille moyenne des grappes est de 20 logements. Au moment du tirage, on a utilisé une stratification par région et degré d'urbanisation. Chaque trimestre, environ 67 000 logements sont identifiés comme résidences principales et enquêtés. Ils sont renouvelés par sixième chaque trimestre. Au final, les fichiers d'enquête comptent environ 108 000 personnes de 15 ans ou plus répondantes chaque trimestre, réparties dans 57 000 ménages. »*

Sondages à plusieurs degrés et par grappes

# **SONDAGE À PLUSIEURS DEGRÉS**

# Notations

Pour simplifier les notations, nous développons le sondage à 2 degrés, mais toutes les notions sont généralisables à 3,4 .. Degrés

## Unités primaires:

M unités primaires dans la population  $\alpha = (1, \dots M)$

m unités tirées dans l'échantillon  $i = (1, \dots m)$

## Unités secondaires :

$N_\alpha$  dans l'unité primaire  $\alpha$  (*unités secondaires*  $\beta = 1, \dots N_\alpha$ )

$n_\alpha$  dans l'échantillon pour chaque unité primaire  $\alpha$  ( $j = 1, \dots n_i$ )

Dans chaque unité primaire  $\alpha$ , le **total**  $T_\alpha$  par unité secondaire est :

$$T_\alpha = \sum_{\beta=1}^{N_\alpha} Y_{\alpha\beta}$$

où  $Y_{\alpha\beta}$  est la valeur de la variable Y pour l'unité secondaire  $\beta$  de l'unité primaire  $\alpha$

# Notations

**Le total sur l'ensemble de la population** est donné par :

$$\sum_{\alpha=1}^M \sum_{\beta=1}^{N_{\alpha}} Y_{\alpha,\beta} = \sum_{\alpha=1}^M T_{\alpha}$$

Dans chaque unité primaire  $\alpha$ , **la moyenne  $\mu_{\alpha}$  par unité secondaire** est :

$$\mu_{\alpha} = \frac{1}{N_{\alpha}} \sum_{\beta=1}^{N_{\alpha}} Y_{\alpha,\beta}$$

où  $Y_{\alpha,\beta}$  est la valeur de la variable Y pour l'unité secondaire  $\beta$  de l'unité primaire  $\alpha$

**La moyenne sur toutes les unités secondaires** est :

$$\mu = \frac{T}{N} \sum_{\alpha=1}^M \frac{N_{\alpha}}{N} \mu_{\alpha}$$

Un estimateur de la variance des totaux est :

$$\sigma^2(T) = \frac{1}{M-1} \sum_{\alpha=1}^M (T_{\alpha} - \bar{T})^2$$

Avec  $\bar{T} = \frac{1}{M} \sum_{\alpha=1}^M T_{\alpha} = \frac{T}{M}$  le total moyen sur les unités primaires

Sondages à plusieurs degrés et par grappes

## **TIRAGE À PROBABILITÉS ÉGALES SANS REMISE (PESR) AUX DEUX DEGRÉS**

# Tirage à probabilités égales

Un plan assez classique consiste à sélectionner les unités primaires et les unités secondaires selon un plan aléatoire simple sans remise.

Les probabilités d'inclusion d'ordre 1 et 2 pour le premier tirage sont donc :

$$\pi_{1i} = \frac{m}{M} \text{ et } \pi_{1ij} = \frac{m(m-1)}{M(M-1)}$$

Pour le second tirage, la taille des échantillons au sein des unités primaires est  $n_i$ , la probabilité d'inclusion pour l'ensemble du plan de sondage vaut donc :

$$\pi_k = \frac{mn_i}{MN_i}$$

*→ Ce plan présente des inconvénients importants : la taille de l'échantillon est aléatoire -puisque cela dépend des unités tirées- et donc son coût*

# Estimateurs du total et de la moyenne

**Un estimateur du total T** d'une population U est :

$$\hat{T} = \frac{M}{m} \sum_{i=1}^m \hat{T}_i = \frac{M}{m} \sum_{i=1}^m \left( \frac{N_i}{n_i} \sum_{j=1}^{n_i} Y_{ij} \right)$$

Où  $\hat{T}_i$  est l'estimateur du total d'une unité primaire  $i$  ( $i = 1, \dots, m$ )

C'est un estimateur sans biais formé simplement par les totaux aux deux degrés de tirage

**Un estimateur de la moyenne par unité secondaire :**

$$\hat{\mu} = \frac{1}{N} \hat{T}$$
$$\hat{\mu} = \frac{1}{N} \frac{M}{m} \sum_{i=1}^m \hat{T}_i = \frac{1}{N} \frac{M}{m} \sum_{i=1}^m \left( \frac{N_i}{n_i} \sum_{j=1}^{n_i} Y_{ij} \right)$$

Cet estimateur est aussi sans biais, mais suppose que N est connu.

Plus tard nous proposons une solution dans le cas où N n'est pas connu

# Estimateurs du total et de la moyenne

**La variance de l'estimateur du total T** est donnée par :

$$\text{Var}(\hat{T}) = M^2 \left(1 - \frac{m}{M}\right) \frac{\sigma^2_1}{m} + \frac{M}{m} \sum_{\alpha=1}^M N^2_{\alpha} \left(1 - \frac{n_{\alpha}}{N_{\alpha}}\right) \frac{\sigma^2_2}{n_{\alpha}}$$

$$\text{Avec } \sigma^2_1 = \frac{1}{M-1} \sum_{\alpha=1}^M (T_{\alpha} - T)^2$$

$$\text{Et } \sigma^2_2 = \frac{1}{N_{\alpha} - 1} \sum_{\beta=1}^{N_{\alpha}} (Y_{\alpha,\beta} - \mu_{\alpha})^2$$

... **La variance de l'estimateur de la moyenne** est donnée par :

$$\text{Var}(\hat{\mu}) = M^2 \frac{1}{N^2} \left(1 - \frac{m}{M}\right) \frac{\sigma^2_1}{m} + \frac{M}{m} \frac{1}{N^2} \sum_{\alpha=1}^M N^2_{\alpha} \left(1 - \frac{n_{\alpha}}{N_{\alpha}}\right) \frac{\sigma^2_2}{n_{\alpha}}$$

Dans la formule, les deux termes correspondent aux deux degrés de tirage, et permettent donc de décomposer la variance à chacune des deux étapes. Si on augmente le nombre d'unités primaires  $m$ , les deux termes diminuent : on a intérêt au nom de la dispersion de les maximiser. Si on augmente le nombre d'unités secondaires, seul le deuxième terme diminue.

# Estimateurs du total et de la moyenne

**Un estimateur sans biais de la variance de l'estimateur du total est :**

$$\widehat{Var}(\widehat{T}) = M^2 \left(1 - \frac{m}{M}\right) \frac{\widehat{\sigma}_1^2}{m} + \frac{M}{m} \sum_{i=1}^m N_i^2 \left(1 - \frac{n_i}{N_i}\right) \frac{\widehat{\sigma}_2^2}{n_i}$$

$$\text{Avec } \widehat{\sigma}_1^2 = \frac{1}{M-1} \sum_{\alpha=1}^m \left(\widehat{T}_i - \frac{\widehat{T}}{M}\right)^2$$

$$\text{Et } \sigma_2^2 = \frac{1}{n_i-1} \sum_{\beta=1}^{n_i} (Y_{ij} - \widehat{\mu}_i)^2$$

## Dans le cas où N n'est pas connue

Nous avons vu comment estimer la moyenne à partir du total, mais parfois la taille de la population N est inconnue. Cela correspond à des cas pour lesquels nous avons la liste des unités primaires, mais pas celles des unités secondaires

Dans ce cas on l'estime par :

$$\hat{N} = M \hat{\bar{N}} = \frac{M}{m} (\sum_{i=1}^m N_i)$$

Où  $\hat{\bar{N}}$  est l'effectif moyen observé pour les unités primaires de l'échantillon

**Un estimateur de la moyenne  $\mu$  est alors :  $\hat{\mu} = \frac{\hat{T}}{\hat{N}}$**

Sondages à plusieurs degrés et par grappes

## **TIRAGE DES UNITÉS PRIMAIRES À PROBABILITÉS INÉGALES AVEC REMISE (PIAR)**

## Avec ou sans remise ?

Comme on a pu la voir les formules sans remise sont assez lourdes, celles à probabilités inégales le sont encore plus

En général les tailles d'échantillon sont assez importantes pour que l'on puisse considérer les approximations faites en l'approchant par un tirage avec remise comme acceptables

# Tirages des unités primaires avec remise

Si  $\pi_{1\alpha}$  est la probabilité de tirage de l'unité primaire  $\alpha$ , alors, **un estimateur du total est** :  $\hat{T} = \frac{1}{m} \sum_{i=1}^m \frac{\hat{T}_i}{\pi_{1\alpha}}$  où est  $\hat{T}_i$  l'estimateur du total pour l'unité primaire.

$\hat{T}$  est sans biais,  $\hat{T}_i$  est lié à la méthode de sondage utilisée au second degré de tirage

**La variance de l'estimateur du total T est donnée par :**

$$Var(\hat{T}) = \frac{1}{m} \sum_{\alpha=1}^M \pi_{1\alpha} \left( \frac{T_{\alpha}}{\pi_{1\alpha}} - T \right)^2 + \frac{1}{m} \sum_{\alpha=1}^M \frac{Var_{\alpha}}{\pi_{1\alpha}}$$

Où  $Var_{\alpha}$  est la variance de l'estimateur  $\hat{T}_{\alpha}$  du total  $T_{\alpha}$  dans l'unité primaire  $\alpha$ , qui est liée au plan de sondage du deuxième degré

# Interprétation

**Un estimateur de la variance de l'estimateur du total T est donnée par :**

$$\widehat{Var}(\widehat{T}) = \frac{1}{m} \sum_{i=1}^m \pi_{1i} \left( \frac{\widehat{T}_i}{\pi_{1i}} - \widehat{T} \right)^2 + \frac{1}{m} \sum_{i=1}^m \frac{\widehat{Var}_i}{\pi_{1i}}$$

Où  $\widehat{Var}_i$  est l'estimateur de la variance de l'estimateur  $\widehat{T}_i$  du total  $T_i$  dans l'unité primaire  $i$ , qui est liée au plan de sondage du deuxième degré

En augmentant le nombre d'unités primaires on diminue la variance des deux termes (deux termes en  $\frac{1}{m}$ ). En travaillant le plan de sondage au deuxième degré, on diminue seulement le deuxième terme. Pour cela il faut que la méthode de tirage au premier degré soit avec remise et qu'il y ait indépendance entre les degrés de tirage. La méthode de tirage au second degré est quelconque, la condition nécessaire est qu'elle permettent d'obtenir des estimateurs sans biais des totaux des unités primaires

## Remarque : le sondage à deux degrés autopondéré

Le sondage autopondéré résout le problème de taille aléatoire rencontré dans la méthode PESR. A la première étape, on sélectionne les unités primaires avec des probabilité d'inclusion proportionnelles à la taille de ces unités primaires

Les probabilités de sélection des unités primaires sont donc :  $\pi_{1i} = \frac{N_i m}{N}$

A la deuxième étape, on sélectionne des unités secondaires selon un plan aléatoire simple sans remise avec une taille d'échantillon  $n_i = n_0$  constante (quelle que soit la taille de l'unité primaire), on a donc pour chaque unité primaire :  $\pi_{k|i} = \frac{n_0}{N_i}$

La probabilité d'inclusion d'ordre 1 vaut donc :  $\pi_k = \pi_{1i} \pi_{k|i} = \frac{N_i m}{N} \frac{n_0}{N_i} = \frac{m n_0}{N}$

- **Les probabilité d'inclusion sont donc toutes constantes pour tous les individus de la population**
- **Le plan est de taille fixe, la taille de l'échantillon vaut toujours  $n = m n_0$**

# Tirage PESR autopondéré

La formule de l'estimateur du total devient dans ce cas

$$\hat{T} = \frac{1}{m} \sum_{i=1}^m \frac{N}{N_i} \left( \frac{N_i}{n_0} \sum_{j=0}^{n_0} y_{ij} \right) = \frac{N}{mn_0} \sum_{i=1}^m \sum_{j=1}^{n_0} y_{ij}$$

Chaque unité a le même coefficient d'extrapolation : le sondage est bien 'autopondéré'

Estimateur de la moyenne: N peut être inconnu

$$\hat{Y} = \bar{y}$$

Sondages à plusieurs degrés et par grappes

# **SONDAGE PAR GRAPPES**

# Sondages en grappes vs sondages à plusieurs degrés

RAPPEL : Le sondage par grappes est un cas particulier de sondage à plusieurs degrés dans lequel l'ensemble des unités du dernier degré est enquêtée

- En conséquence l'estimation de la moyenne générale sera simplement un problème d'estimation à partir d'une population de grappes, les échantillons seront constitués des quantités calculées des moyennes dans les grappes
- Dans les formules de variance, il n'y aura plus d'aléa au deuxième niveau, puisque l'on tire tous les individus dans une grappe (on effectue un 'recensement' dans chaque grappe tirée)

# Notations

## Unités primaires:

M unités primaires dans la population  $\alpha = (1, \dots, M)$

m unités tirées dans l'échantillon  $i = (1, \dots, m)$

## Unités secondaires :

$N_\alpha$  taille de la grappe  $\alpha$  (*unités secondaires*  $\beta = 1, \dots, N_\alpha$ )

$Y_{\alpha\beta}$  est la valeur de la variable étudiée pour l'unité secondaire  $\beta$  de la grappe  $\alpha$

## Par grappe :

Dans chaque grappe  $\alpha$ , le total  $T_\alpha$  est :  $T_\alpha = \sum_{\beta=1}^{N_\alpha} Y_{\alpha\beta}$

Taille moyenne des grappes :  $\bar{N} = \frac{1}{M} \sum_{\alpha=1}^M N_\alpha$

Total moyen par grappe :  $\bar{T} = \frac{1}{M} \sum_{\alpha=1}^M T_\alpha$

Moyenne à l'intérieur de chaque grappe  $\alpha$  :  $\bar{Y}_\alpha = \frac{1}{N_\alpha} \sum_{\beta=1}^{N_\alpha} Y_{\alpha\beta} = \frac{T_\alpha}{N_\alpha}$

# Notations

## Pour l'ensemble de la population

La taille de la population :  $N = \sum_{\alpha=1}^M N_{\alpha}$

Le total général :  $T = \sum_{\alpha=1}^M Y_{\alpha}$

La moyenne générale :  $\bar{Y} = \frac{1}{N} \sum_{\alpha=1}^M \sum_{\beta=1}^{N_{\alpha}} Y_{\alpha\beta} = \sum_{\alpha=1}^M \frac{N_{\alpha}}{N} \bar{Y}_{\alpha} = \frac{T}{N}$

Sondages à plusieurs degrés et par grappes

**TAILLES DES GRAPPES CONNUES A PRIORI - TIRAGE  
DES GRAPPES PESR, GRAPPES DE TAILLES ÉGALES**

# Grappes de tailles égales, probabilités égales, estimation d'une moyenne

On réalise un sondage aléatoire simple sans remise dans une population de grappes, les échantillons seront constitués des quantités calculées dans les grappes, chaque grappe apporte le même nombre d'individus

**La taille de l'échantillon est donc fixe** : nombre de grappes x nombre d'individus tirés dans chaque grappe

L'estimateur de la moyenne découle de la définition d'un SAS : moyenne arithmétique des moyennes calculées dans les grappes, la variance découle des écarts entre la moyenne globale et les moyennes calculées dans les strates

# Grappes de tailles égales, probabilités égales, estimation d'une moyenne

On note  $N_0$  la taille des grappes ( $N_\alpha = N_0 = \bar{N} \forall \alpha$ ), on a donc  $N = MN_0$

Dans ce cas simple, la moyenne générale est simplement la moyenne arithmétique des moyennes par grappe :

$$\bar{Y}_g = \sum_{\alpha=1}^M \frac{N_0}{N} \bar{Y}_\alpha = \frac{1}{M} \sum_{\alpha=1}^M \bar{Y}_\alpha$$

Si les  $m$  grappes sont tirées à probabilités égales alors un estimateur sans biais de la moyenne générale est :

$$\widehat{\bar{Y}}_g = \frac{1}{m} \sum_{i=1}^m \widehat{Y}_i$$

# Grappes de tailles égales, probabilités égales, estimation d'une moyenne

La variance de l'estimateur :

$$Var(\widehat{Y}_g) = \frac{M - m}{Mm} \frac{1}{M - 1} \sum_{\alpha=1}^M (\bar{Y}_\alpha - \bar{Y})^2$$

Son estimation :

$$\widehat{Var}(\widehat{Y}_g) = \frac{M - m}{Mm} \frac{1}{m - 1} \sum_{i=1}^m (\widehat{Y}_i - \widehat{Y})^2$$

**→ A ce stade, une première conclusion est qu'un sondage par grappes sera d'autant plus précis qu'il y a beaucoup de grappes (m est grand) qui se ressemblent en moyenne.**

# Grappes de tailles égales, probabilités égales, estimation d'une moyenne

Nous allons commencer à partir de ce cas simple à étudier les conditions qui vont rendre un sondage par grappe intéressant du point de vue de la précision

Nous allons faire la comparaison avec un plan de sondage de référence, le SAS. Pour cela il nous faudra établir une mesure du degré de similarité entre les grappes.

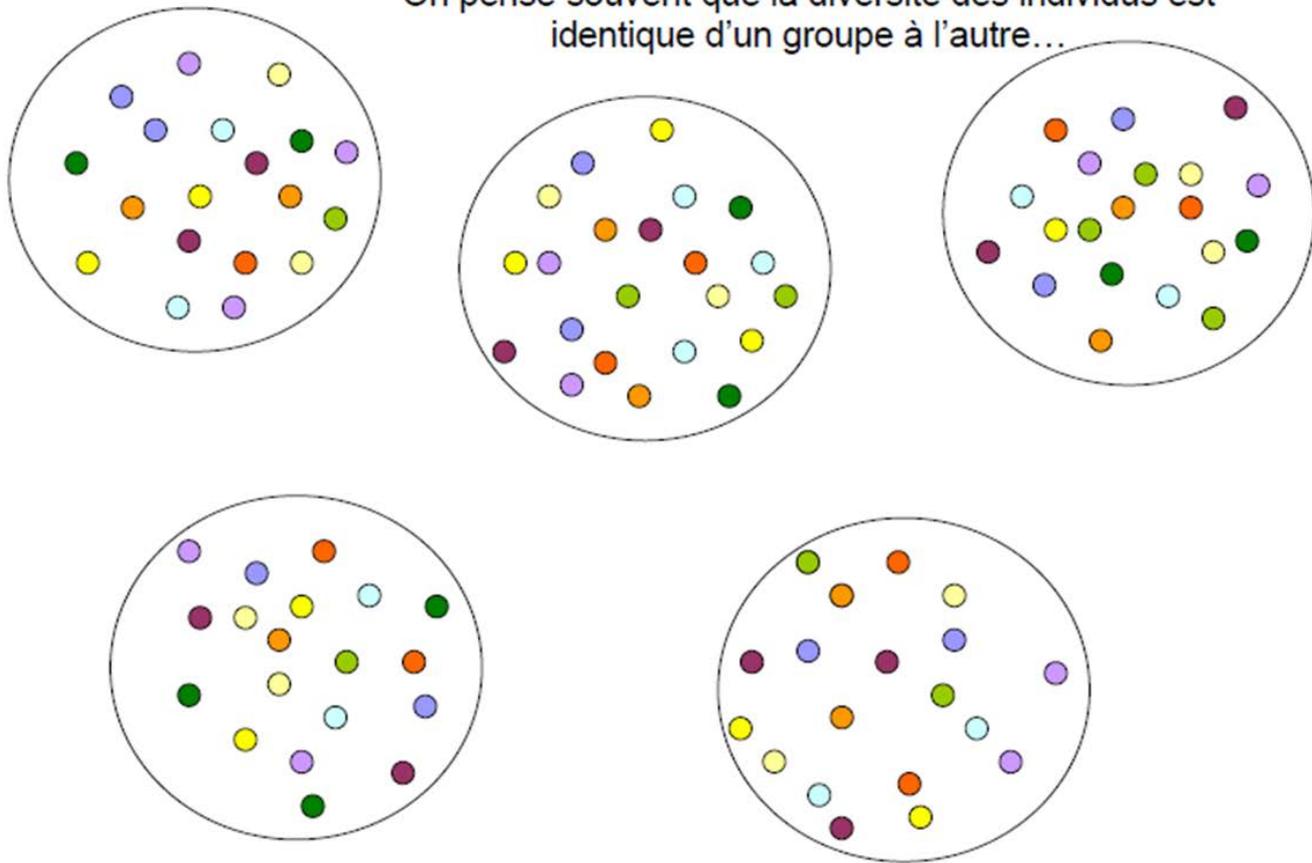
# Notion de rapport de corrélation inter-grappes

La rapport de corrélation inter-grappes est le rapport de la variance inter grappes (entre les différentes grappes) sur la variance totale

$$\rho^2 = \frac{\sum_{\alpha=1}^M N_0 (\bar{Y}_\alpha - \bar{Y})^2}{\sum_{\alpha=1}^M \sum_{\beta=1}^{N_0} (Y_{\alpha\beta} - \bar{Y})^2}$$

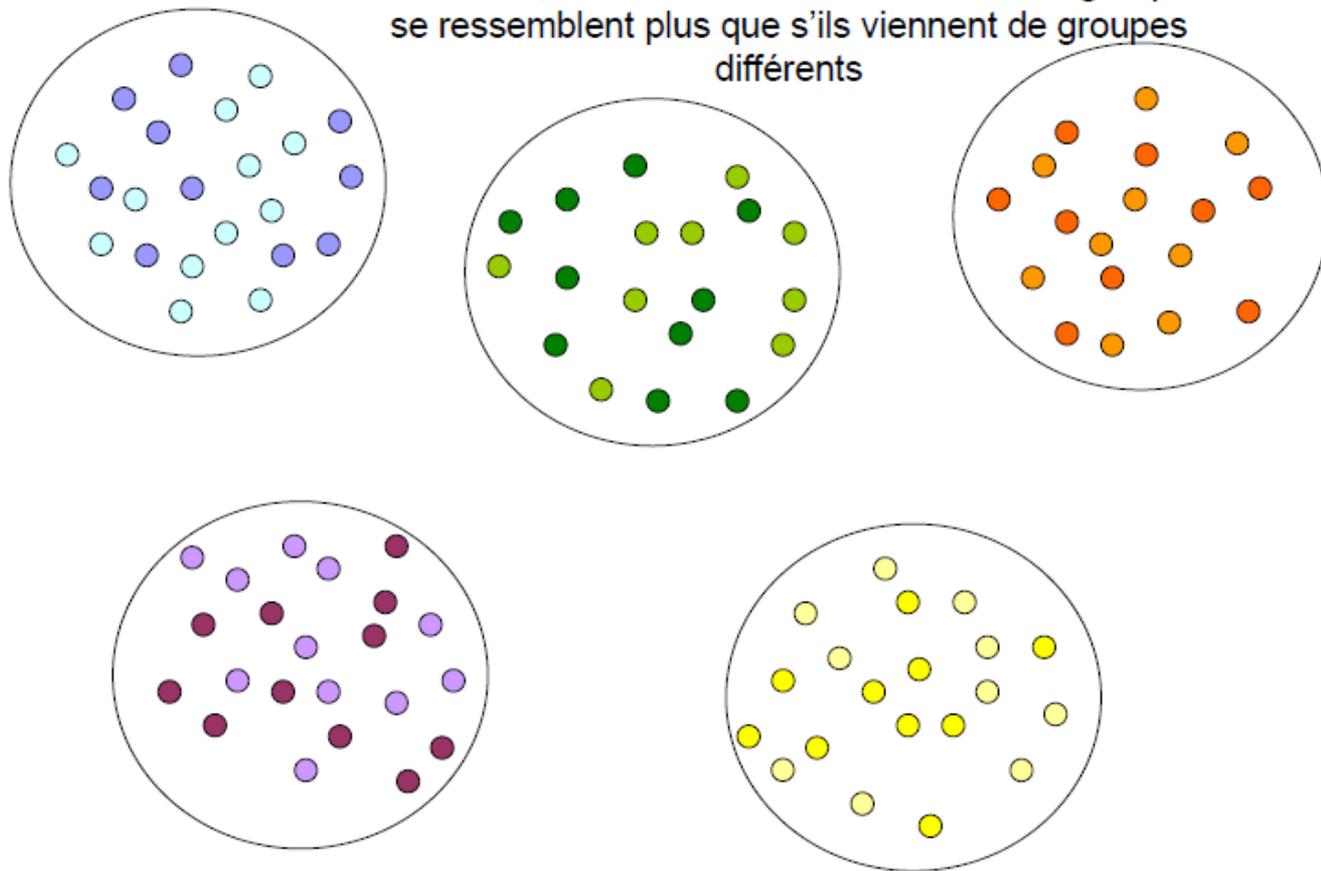
# Cas 1 : $\rho^2$ sera faible

On pense souvent que la diversité des individus est identique d'un groupe à l'autre...



## Cas 2 : $\rho^2$ sera fort

En réalité, les individus au sein d'un même groupe se ressemblent plus que s'ils viennent de groupes différents



# Comparaison avec un SAS de même taille

Comparons avec la réalisation d'un échantillon de même taille, en ignorant les grappes. On a donc :  $n = \sum_{i=1}^m N_i = mN_0$  puisque  $(N_i = N_0 = \bar{N} \forall i)$

**La moyenne générale est estimée** par :  $\hat{Y} = \frac{1}{n} \sum_{j=1}^n y_j$

**La variance de  $\hat{Y}$  est :**  $Var(\hat{Y}) = \frac{N-n}{Nn} S^2_c$

$$\text{Avec } S^2_c = \frac{1}{N-1} \sum_{\alpha=1}^M \sum_{\beta=1}^{N_0} (Y_{\alpha\beta} - \hat{Y})^2$$

Puisque  $N = MN_0$  et  $n = mN_0$ , alors :

$$Var(\hat{Y}) = \frac{1}{N_0} \frac{M-n}{Mn} S^2_c$$

# Comparaison avec un SAS de même taille

**Variance dans le cas des grappes :**

$$Var(\hat{Y}_g) = \frac{M-m}{Mm} \frac{1}{M-1} \sum_{\alpha=1}^M (\bar{Y}_\alpha - \bar{Y})^2$$

Si N et M sont grands,  $N_0$  petit (grande population constituée d'un grand nombre de grappes), alors l'approximation suivante est acceptable :

$$Var(\hat{Y}_g) \approx \frac{M-m}{Mm} \rho^2 S_c^2$$

**Dans le cas d'un SAS :**

$$Var(\hat{Y}) = \frac{1}{N_0} \frac{M-n}{Mn} S_c^2$$

➔ **Le sondage par grappes est intéressant si le rapport de corrélation inter grappes est faible**, inférieur à  $\frac{1}{N_0}$  à l'approximation près

# Conclusions

Il est souhaitable que les moyennes des grappes soient les plus semblables possible. Il ne faut pas que la taille des grappes soit trop élevée

La sondage par grappes dans une population de grappes de tailles égales est d'autant plus efficace que la dispersion totale est essentiellement constituée par l'hétérogénéité des individus au sein des classes.

➔ Le sondage par grappes est efficace s'il y a beaucoup de petites grappes, les plus ressemblantes possibles

# Application au tirage systématique

- Très utilisé à la place d'un tirage aléatoire à probabilités égales
- Soit  $N$  multiple de  $n$ . Par exemple on veut tirer 10 individus parmi 1000 : on commence par tirer au hasard un nombre entier entre 1 et 100, si ce nombre est 27, le premier individu sera le n°27, le deuxième le n°127 etc. jusqu'au n°927.
- De façon générale si on a tiré un entier  $h$ , les individus sélectionnés ont les numéros :  $h, h+M, h+2M, \dots, h+(n-1)M$ .
- **Tirage d'une seule grappe** parmi  $M=N/n$  grappes.

- L'estimateur de la moyenne est simplement la moyenne de la grappe sélectionnée et sa variance est

$$V\left(\hat{\bar{Y}}\right) = M \sum_{i=1}^M \left( \frac{\bar{Y}_i N_i}{N} - \frac{\bar{Y}}{M} \right)^2$$

- Lorsque le fichier se trouve être trié selon un ordre proche de Y, la variance peut être notablement plus faible que pour le tirage aléatoire simple. Exemple  $Y_i = i$
- Mais la variance n'est pas estimable .

– Voir formule  $s_1^2 = \frac{1}{m-1} \sum_{i=1}^m \left( \hat{T}_i - \frac{\hat{T}}{M} \right)^2$

- Il est incorrect d'utiliser la variance de l'estimateur du tirage aléatoire simple sauf si la base de sondage a été triée préalablement au hasard.

## Tirage systématique: un exemple théorique

$Y_i=i$  Population triée par ordre croissant  $N=Kn$

$$\bar{Y} = \frac{N+1}{2} \quad S^2 = \frac{(N+1)^2}{12}$$

• Tirage équiprobable sans remise :

$$V(\bar{y}_{sr}) = \left(1 - \frac{n}{N}\right) \frac{(N+1)^2}{12n} = \left(1 - \frac{1}{K}\right) \frac{(Kn+1)^2}{12n}$$

une grappe :  $h, h+K, h+2K, \dots, h+(n-1)K$

• Moyenne

$$\bar{Y}_h = h + \frac{n-1}{2} K$$

$$E(\bar{Y}_h) = E(h) + \frac{n-1}{2} K = \frac{K+1}{2} + \frac{n-1}{2} K = \frac{nK+1}{2} = \frac{N+1}{2}$$

• Variance

$$V(\hat{Y}_{syst}) = V\left(h + \frac{n-1}{2} K\right) = V(h) = V(h) = \frac{K^2 - 1}{12}$$

$$V(\hat{Y}_{syst}) < V(\bar{y}_{sr})$$

Exemple  $N=20$   $n=4$        $V(\hat{Y}_{syst}) = 1.33$        $V(\bar{y}_{sr}) = 7.35$

Sondages à plusieurs degrés et par grappes

**TAILLES DES GRAPPES CONNUES A PRIORI - TIRAGE  
DES GRAPPES PESR, GRAPPES DE TAILLES INÉGALES**

# Grappes de tailles inégales, probabilités égales, estimation d'une moyenne

On réalise un sondage aléatoire simple sans remise dans une population de grappes, les échantillons seront constitués des quantités calculées des moyennes dans les grappes, chaque grappe apporte un nombre différent d'individus. La taille de l'échantillon n'est plus fixe : même si on décide à priori un nombre à tirer dans chaque strate de la population, elle dépendra des grappes choisies finalement

L'estimateur est maintenant la moyenne pondérée par les tailles relatives des grappes des moyennes calculées dans les grappes, la variance des écarts entre moyenne globale et moyennes calculées dans les strates

# Grappes de tailles inégales, probabilités égales, estimation d'une moyenne

→ Les formules sont analogues à celles du cas simple introductif des tailles égales, on intègre le facteur des tailles relatives de grappes dans l'échantillon

Dans ce cas , la moyenne générale est simplement la moyenne pondérée des moyennes par grappe :

$$\overline{Yg} = \frac{1}{M} \sum_{\alpha=1}^M \frac{N_{\alpha}}{\overline{N}} \overline{Y}_{\alpha} \text{ avec } N = \frac{N}{M}$$

Si les m grappes sont tirées à probabilités égales alors un estimateur sans biais de la moyenne générale est :

$$\widehat{\overline{Yg}} = \frac{1}{m} \sum_{i=1}^m \frac{N_i}{\overline{N}} \widehat{Y}_i$$

# Grappes de tailles inégales, probabilités égales, estimation d'une moyenne

La variances de l'estimateur :

$$Var(\widehat{Y}_g) = \frac{M - m}{Mm} \frac{1}{M - 1} \sum_{\alpha=1}^M \left( \overline{Y}_\alpha \frac{N_\alpha}{N} - \bar{Y} \right)^2$$

Son estimation :

$$\widehat{Var}(\widehat{Y}_g) = \frac{M - m}{Mm} \frac{1}{m - 1} \sum_{i=1}^m \left( \widehat{Y}_i \frac{N_i}{N} - \widehat{Y} \right)^2$$

**→ A ce stade, une première conclusion est qu'un sondage par grappes sera d'autant plus précis qu'il y a beaucoup de grappes (m est grand) qui se ressemblent en moyenne.**

# Comparaison avec un SAS

Le nombre d'unités statistiques de l'échantillon est aléatoire car il dépend des grappes choisies. On réalise donc la comparaison avec un SAS de taille  $m\bar{N}$  l'espérance mathématique de la taille de l'échantillon

Les conclusions sont identiques à ce que l'on a vu précédemment : il est préférable d'avoir beaucoup de grappes, dont la taille moyenne est faible et dont les moyennes ne soient pas trop dissemblables

Sondages à plusieurs degrés et par grappes

## **TAILLES DES GRAPPES CONNUES A PRIORI - TIRAGE DES GRAPPES PIAR**

# Tirage des grappes à probabilités inégales avec remise, estimation d'une moyenne

On réalise un sondage à probabilités inégales avec remise dans une population de grappes, les échantillons seront constitués des quantités calculées des moyennes dans les grappes.

A chaque tirage, la grappe  $\alpha$  est retenue avec la probabilité  $P_\alpha$

L'estimateur du total  $T$  est maintenant :

$$\hat{T} = \frac{1}{m} \sum_{i=1}^m \frac{Y_i}{\pi_i}$$

L'estimateur de la moyenne est :

$$\hat{Y} = \frac{1}{m} \sum_{i=1}^m \frac{Y_i}{\pi_i}$$

# Tirage des grappes à probabilités inégales avec remise, estimation d'une moyenne

La variance de l'estimateur de la moyenne est :

$$Var(\hat{Y}) = \frac{1}{N^2} Var(\hat{T}) = \frac{1}{N^2} \frac{1}{M} \sum_{\alpha=1}^M \pi_{\alpha} \left( \frac{\hat{T}_{\alpha}}{\pi_{\alpha}} - \hat{T} \right)^2$$

Son estimation :

$$\widehat{Var}(\hat{Y}) = \frac{1}{N^2} \widehat{Var}(\hat{T}) = \frac{1}{N^2} \frac{1}{M} \sum_{i=1}^m \pi_i \left( \frac{\hat{T}_i}{\pi_i} - \hat{T} \right)^2$$

## Cas particulier : probabilités proportionnelles aux tailles.

Si les moyennes ou les totaux par strates sont corrélés avec le nombre d'unités qu'elles contiennent, on sera efficace de choisir les

probabilités proportionnelles à ces tailles :  $\pi_{\alpha} = \frac{N_{\alpha}}{N}$

Sondages à plusieurs degrés et par grappes

**TAILLES DES GRAPPES INCONNUES A PRIORI -  
TIRAGE DES GRAPPES PESR, GRAPPES DE TAILLES  
INÉGALES**

# Tailles des grappes inconnues a priori, mais population totale connue

A défaut d'information complémentaire, les tirages se feront PESR. **Si la taille globale de la population est connue**, mais pas celles des grappes, alors on est dans le cas d'un sondage à probabilités inégales sans remise.

L'estimateur du total  $T$  est donc :  $\hat{T} = \frac{M}{m} \sum_{i=1}^m T_i = M\hat{Y}$

Sa variance est :

$$Var(\hat{T}) = M^2 \frac{M-m}{Mm} S_c^2 \text{ avec } S_c^2 = \frac{1}{M-1} \sum_{\alpha=1}^M (\bar{T}_\alpha - \bar{T})^2$$

Un estimateur de sa variance est:

$$\widehat{Var}(\hat{T}) = M^2 \frac{M-m}{Mm} \widehat{S}_c^2 \text{ avec } \widehat{S}_c^2 = \frac{1}{m-1} \sum_{i=1}^m (\widehat{T}_i - \hat{T})^2$$

Pour la moyenne, on a :  $\hat{Y} = \frac{1}{N} \hat{T}$  et  $\widehat{Var}(\hat{Y}) = \frac{1}{N^2} \widehat{Var}(\hat{T})$

# Tailles des grappes inconnues a priori, mais population totale inconnue

Si la taille globale  $N$  de la population est inconnue, il faut l'estimer, et on se retrouve dans le cas de l'estimation d'un ratio (quantités aléatoires au numérateur et au dénominateur)

Si on note  $\widehat{N}$  l'estimation de  $N$  alors  $\widehat{Y} = \frac{1}{\widehat{N}} \widehat{T}$

Sondages à plusieurs degrés et par grappes

# **GRAPPES ET STRATIFICATION, MISE EN ŒUVRE EFFICACE**

# Mise en œuvre pratique

Pour avoir un rapport de corrélation inter grappes le plus petit possible, nous avons vu qu'il faut un grand nombre de grappes dont les moyennes sont peut être différentes les unes des autres, ce qui n'est pas réalisé dans les conditions concrètes (on voudrait que chaque grappe constitue une 'mini population', on contredit la notion même de grappe ...)

Par contre, cette condition peut être approchée si l'on constitue des sous ensembles de grappes : des strates

C'est ce que l'on fait en pratique pour conjuguer les effets bénéfiques de la stratification sur la précision et des grappes sur l'économie des moyens

# Mise en œuvre pratique

Le lien avec le principe de la stratification est facile. : les strates doivent être les plus contrastées possible pour bien prendre en compte la variabilité du phénomène étudié. Mais à l'intérieur d'une strate, les grappes doivent se ressembler le plus possible

La répartition de l'échantillon dans les strates doit aussi intégrer la variabilité interne aux strates : si dans une strate, les grappes sont très ressemblantes, on pourra en sélectionner moins que dans les strates où les grappes sont plus différentes les unes des autres (application du principe de l'allocation optimale de Neyman)

## Quelques cas

Etudes de satisfaction des passagers de compagnies aérienne : stratification selon le type de vol (les périodes, les horaires sont plus ou moins loisir vs business) et les faisceaux (Asie, Europe, ...)  
Une fois cette stratification opérée, les vols sont des grappes de passagers.

Etudes de marché : en général, stratification région x catégorie d'agglomération puis tirage des unités secondaires (iris/ilot, ...) proportionnel à la taille. Les instituts privés font à la différence de l'INSEE (du fait de l'absence de base de sondage) la dernière étape par quotas : de 10 personnes par 'point de chute' à partir d'une feuille de quotas.

Exemple: étude du taux de contact en  
face à face (BVA, Paris 1991)

---

## SOMMAIRE

	PAGES
<b>I - METHODOLOGIE .....</b>	<b>3</b>
• ECHANTILLON .....	5
• LES ILOTS .....	7
• LE TERRAIN .....	9
• METHODE D'ENQUETE .....	10
<b>II - LES RESULTATS .....</b>	<b>11</b>
• RESULTATS SUR L'ENSEMBLE DES ILOTS	12
• SONDAGES ANTERIEURS .....	13
• STRUCTURE SOCIO-PROFESSIONNELLE	15

## ECHANTILLON

→ Population étudiée :

Ménages ordinaires habitant Paris intra-muros (résidence principale)

→ Base de sondage (source APUR) :

- ensemble des îlots parisiens (recensement 1990)
- structure socio-démographique (recensement 1982)



4 337 îlots (sur total 5 133) habités en 1982 et ayant au moins 5 résidences principales en 1990

→ Mode de tirage :

- sondage en grappes (îlots)
- probabilités inégales après stratification



→ Stratification :

- analyse factorielle des correspondances du tableau de contingence juxtaposé croisant les îlots et les caractéristiques socio-démographiques
- classification mixte sur la totalité des facteurs de l'analyse des correspondances



8 types d'îlots

→ Tirage des îlots :

Pour chaque type d'îlot : proportionnellement au nombre de résidences principales dans chaque îlot.



## LES ILOTS

20 îlots retenus :

• 4ème arrondissement .....	3
• 6ème arrondissement .....	1
• 9ème arrondissement .....	2
• 12ème arrondissement .....	1
• 13ème arrondissement .....	1
• 14ème arrondissement .....	2
• 15ème arrondissement .....	1
• 16ème arrondissement .....	3
• 17ème arrondissement .....	2
• 18ème arrondissement .....	2
• 19ème arrondissement .....	1
• 20ème arrondissement .....	1

ILOT LE PLUS IMPORTANT : 720 résidences principales

ILOT LE MOINS IMPORTANT : 22 résidences principales

La structure du plan de sondage est donc la suivante :

	POPULATION		ECHANTILLON	
TYPE	Nbre total	Nbre de résidences principales	Nbre d'îlots	Nbre résidences principales
1	679	206863	3	1162
2	302	52685	1	434
3	353	105182	2	748
4	505	132620	2	692
5	695	220505	3	1296
6	421	82791	2	809
7	608	107896	3	561
8	774	141363	4	1177
TOTAL	4337	1049905	20	6879

On trouvera plus loin le tableau donnant les caractéristiques des 20 îlots échantillonnés.

La répartition géographique des 20 îlots couvre 12 arrondissements (4°, 6°, 9°, 12°, 13°, 14°, 15°, 16°, 17°, 18°, 19°, 20°)



191, Avenue du Général Leclerc  
78220 VIROFLAY

ETUDE : F 1 / 3 / 7 /

MAI 1991

ENQUETEUR : ..... N° ENQ. :  / / / / / / / /

Ce questionnaire témoigne d'un contact avec une personne occupant un logement de l'immeuble. L'acceptation de l'interview ou son refus ont lieu en tête à tête, à la porte de l'appartement ou à l'intérieur, jamais dans l'entrée ou l'escalier, etc. ou en utilisant l'interphone.

En cas de refus d'interview, vous notez les informations prévues ci-dessous.

- A - Arrondissement .....  / / /      • B - Numéro d'ilot .....  / / / / / / / /  
• C - Rue .....      • D - Numéro dans la rue ...  / / / /

Si à ce numéro de rue, on se trouve en présence de plusieurs immeubles, bâtiments, escaliers ou ascenseurs, chaque escalier ou ascenseur définit un immeuble.

- E - A ce numéro de rue, l'immeuble a un seul accès aux appartements ..... 1  
- A ce numéro de rue, on trouve un immeuble à plusieurs entrées ou un ensemble d'immeuble. Il est identifié par : (ex. lettre, chiffre, nom, etc.) : ..... 2  
.....  
.....
- IMPORTANT : Reporter ici le numéro de la fiche immeuble .....  / / / / /
- F - Le contact est pris avec une personne qui habite l'étage numéro  / / /
- G - Date :  / / / / / / / 9 / 1 /      • H - 1ère visite ..... 1  
2ème visite ..... 2

Je suis ..... de la société BVA, institut de sondages

**Si la personne refuse l'interview prenez congé et notez en dehors de l'appartement: R1 à R3**

- R1 - Refus non motivé ..... 1  
- Trop souvent interrogé ..... 2  
- Manque de temps ..... 3  
- Autres (*préciser*) ..... 4  
.....

- R2 - Sexe :  
- Homme ..... 1  
- Femme ..... 2

- R3 - Age approximatif :  
- Moins de 15 ans ..... 1  
- De 15 à 17 ans ..... 2  
- De 18 à 24 ans ..... 3  
- De 25 à 34 ans ..... 4  
- De 35 à 49 ans ..... 5  
- De 50 à 64 ans ..... 6  
- 65 ans et plus ..... 7

- R4 - Etait-elle, selon vous ou selon ce qu'elle a déclaré?  
• Le chef de famille ..... 1  
• La maîtresse de maison ..... 2  
• Un autre adulte ..... 3  
• Une femme de ménage ..... 4

**Vous posez tout de suite la question 1 :**

1 - Aimez-vous habiter dans ce quartier ?

- OUI ..... 1  
• NON ..... 2  
• NSP ..... 3

2 - Vous sentez-vous ici, plus, autant ou moins en sécurité que dans d'autres quartiers de Paris ?

- Plus ..... 1  
• Autant ..... 2  
• Moins ..... 3  
• NSP ..... 4

3 - Si c'était à refaire, préféreriez-vous habiter ailleurs que dans Paris ?

- OUI ..... 1      --> 4  
• NON ..... 2      --> 5  
• NSP ..... 3      --> 5

4 - Où préféreriez-vous habiter ?  
(Plusieurs réponses possibles)

- En proche banlieue .. 1  
• En grande banlieue .. 2  
• En province ..... 3  
• A l'étranger ..... 4  
• NSP ..... 5

**A TOUS**

5 - D'habitude, tous les combien, personnellement, lisez-vous ou feuilletez-vous chez vous ou ailleurs, un numéro de ... (**Tendre carte 5**)

	France Soir	Le Figaro	Le Monde	Le Parisien	Libération
• Tous les jours .....	1	1	1	1	1
• 3 à 5 fois par semaine .....	2	2	2	2	2
• 1 à 2 fois par semaine .....	3	3	3	3	3
• Moins souvent .....	4	4	4	4	4
• Jamais .....	5	5	5	5	5

6 - Avez-vous personnellement déjà été interrogé au cours d'un sondage ?

- OUI ..... 1      -> **7a**
- NON ..... 2      -> **S1**
- NSP ..... 3      -> **S1**

7a - Combien de fois avez-vous répondu à un sondage à votre domicile sans compter aujourd'hui ?

- ... 1 fois ..... 1
- ... 2 ou 3 fois..... 2
- ... plus de 3 fois..... 3
- ... jamais..... 4

7b - Dans la rue ou dans un lieu public ...

- ... 1 fois ..... 1
- ... 2 ou 3 fois..... 2
- ... plus de 3 fois..... 3
- ... jamais..... 4

7c - Au téléphone à votre domicile ...

- ... 1 fois ..... 1
- ... 2 ou 3 fois..... 2
- ... plus de 3 fois..... 3
- ... jamais..... 4

7d - Par correspondance ...

- ... 1 fois ..... 1
- ... 2 ou 3 fois..... 2
- ... plus de 3 fois..... 3
- ... jamais..... 4

7e - Sur votre lieu de travail que ce soit au téléphone ou par la visite d'un enquêteur ...

- ... 1 fois ..... 1
- ... 2 ou 3 fois..... 2
- ... plus de 3 fois..... 3
- ... jamais..... 4

8a - Quand avez-vous répondu à un sondage pour la dernière fois, était-ce ...

(**Enq. : Enumérer**)

- ... il y a moins de 8 jours ..... 1
- ... entre 8 et 15 jours ..... 2
- ... entre 15 jours et 1 mois .... 3
- ... entre 1 et 3 mois..... 4
- ... il y a plus de 3 mois ..... 5
- ... il y a plus d'1 an..... 6
- ... NSP ..... 4

8b - Avez-vous cette dernière fois été interrogé ... (**Enq. : Enumérer**)

- ... à votre domicile ..... 1
- ... au téléphone à votre domicile ..... 2
- ... dans une camionnette ou dans un salon d'hotel..... 3
- ... dans la rue ou dans un lieu public 4
- ... sur votre lieu de travail ..... 5
- ... par correspondance ..... 6
- ... autre réponse, **laquelle ?** ..... 7

**Préciser :**

**A TOUS**

S1 - Sexe :

- Homme ..... 1
- Femme ..... 2

S2 - Age en clair : (**le demander**)    /    /    ans

S3 - Quelle est votre profession ?

S4 - Etes-vous le chef de famille ?

- OUI ..... 1      -> **S.6**
- NON ..... 2      -> **S.5**

S5 - Quelle est la profession du chef de famille ?

**S.6 Coder seulement la colonne "chef de famille" s'il s'agit de l'interviewé.**

	Chef de famille	Interviewé
• Artisan, petit commerçant (moins de 10 salariés) .....	01	01
• Industriel, gros commerçant (10 salariés et plus) .....	02	02
• Prof. libérale, cadre sup. profes. intellectuelle sup. ....	03	03
• Cadre moyen, agent technique	04	04
• Employé .....	05	05
• Ouvrier .....	06	06
• Inactif : <b>préciser :</b>		
• Etudiant .....	07	07
• Ménagère .....	08	08
• Retraité .....	09	09
• Autre .....	10	10



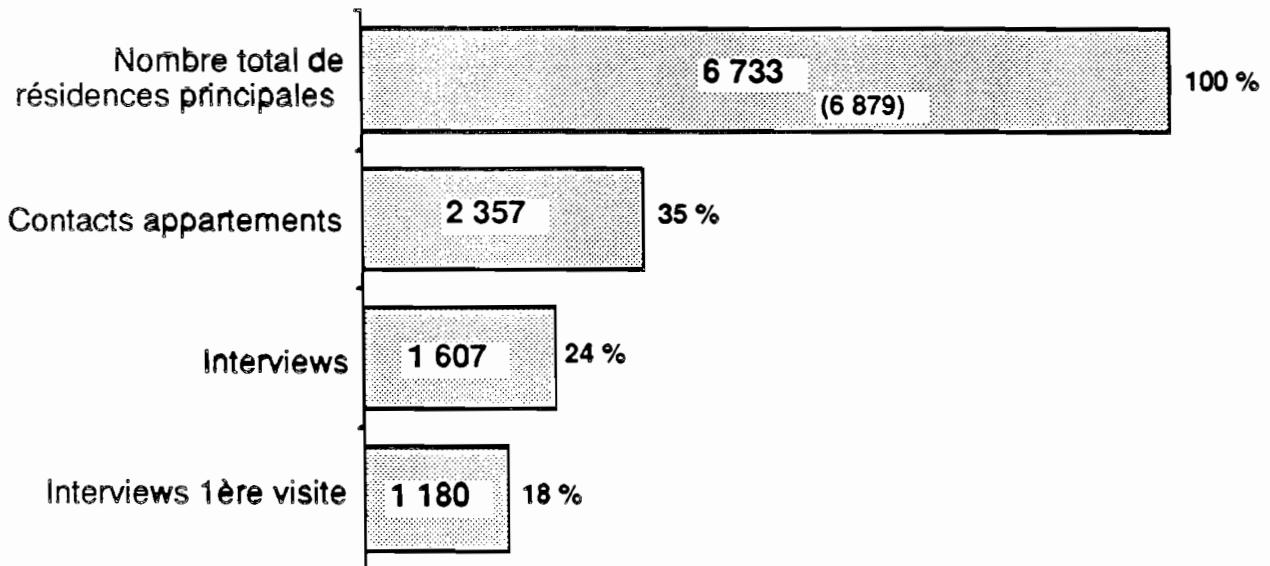
## STRUCTURE SOCIO-PROFESSIONNELLE

	ENS. PARIS %	INTERVIEWS (1 607) %
CSP Chef de famille		
• Artisans, commerçants, chefs d'entreprise.....	5,5	5,0
• Cadres, professions intellectuelles supérieures.....	17,8	29,3
• Professions intermédiaires.....	13,1	11,6
• Employés.....	17,2	14,4
• Ouvriers.....	13,0	11,6
• Retraités.....	<del>28,5</del> 22,5	18,9
• Autres.....	10,9	9,2

### Nombre de personnes au foyer

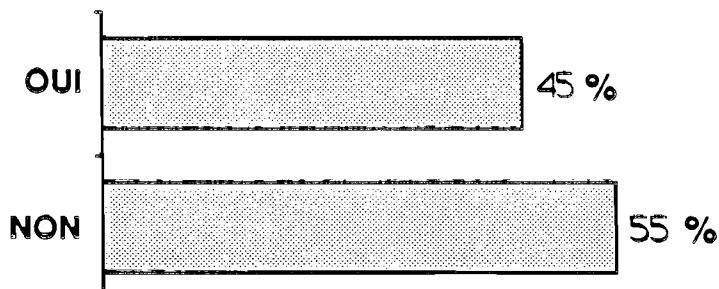
• 1.....	48,2	34,7
• 2.....	27,8	30,2
• 3.....	12,1	13,3
• 4.....	7,9	12,2
• 5.....	2,7	5,7
• 6 et plus.....	1,3	3,9

**RESULTATS SUR L'ENSEMBLE  
DES ILOTS**

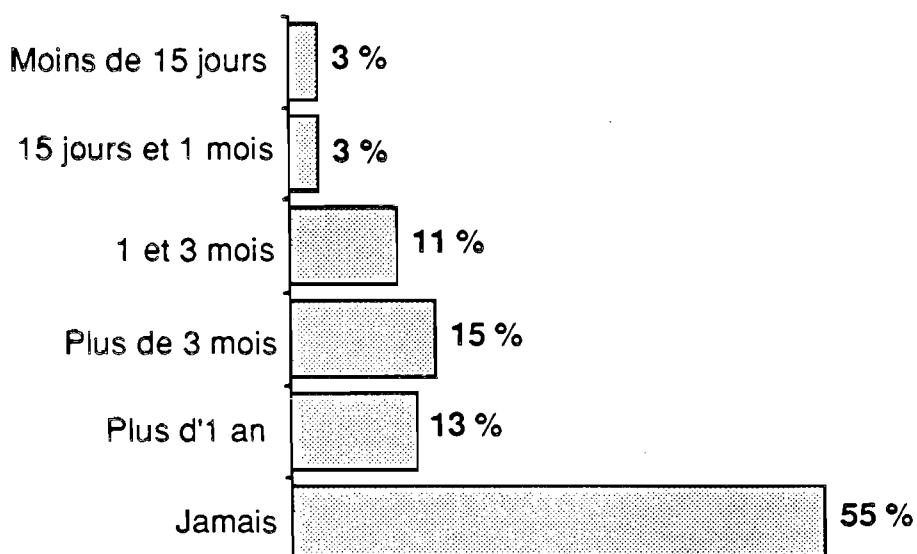


## SONDAGES ANTERIEURS

→ Déjà interrogé au cours d'un sondage



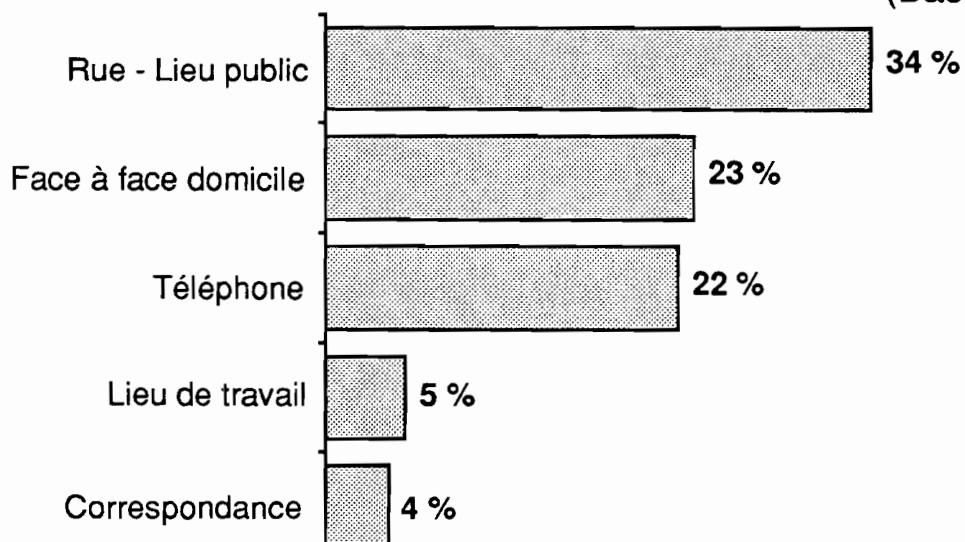
→ Date du dernier sondage



## SONDAGES ANTERIEURS

### → Méthodes d'interviews rencontrées

(Base : 1 607)



### → Face à face à domicile

