

# Enquêtes et sondages

UE STA 108

## MANUEL D'EXERCICES

Sylvie Rousseau

# Table des matières

<b>I. Rappels de probabilités et de statistique inférentielle .....</b>	<b>3</b>
<i>Rappels sur les lois de probabilités</i>	5
<i>Rappels sur les intervalles de confiance</i>	7
<b>II. Sondage aléatoire simple .....</b>	<b>11</b>
<i>Rappels sur le sondage aléatoire simple</i>	16
<b>III. Plans à probabilités inégales .....</b>	<b>18</b>
<i>Rappels sur les plans à probabilités inégales</i>	20
<b>IV. TP1 : Simulations de tirage d'échantillons .....</b>	<b>21</b>
<b>V. Plans stratifiés.....</b>	<b>24</b>
<i>Rappels sur les plans stratifiés</i>	29
<b>VI. Plans par grappes .....</b>	<b>31</b>
<i>Rappels sur les plans par grappes</i>	35
<b>VII. Plans à plusieurs degrés .....</b>	<b>37</b>
<i>Rappels sur les plans à plusieurs degrés</i>	40
<b>VIII. Redressements .....</b>	<b>42</b>
<i>Rappels sur les redressements</i>	44
<b>IX. TP2 : Calage sur marges .....</b>	<b>49</b>
<b>X. TP3 : Correction de la non-réponse.....</b>	<b>49</b>
<b>XI. Compléments et révisions.....</b>	<b>49</b>

# I. Rappels de probabilités et de statistique inférentielle

## Exercice 1

### Notions d'espérance et de variance

Un passager du métro mesure son temps de trajet domicile-travail pendant 10 jours et relève successivement (en minutes) : 32 ; 25 ; 28 ; 36 ; 30 ; 26 ; 37 ; 25 ; 33 ; 28 .

Quel est en moyenne la durée du trajet ? Évaluer aussi la variabilité de cette durée.

Comparer avec un autre itinéraire emprunté par notre voyageur pendant les jours suivants et qui lui prend : 46 ; 21 ; 24 ; 38 ; 44 ; 22 ; 37 ; 20 ; 25 ; 23 minutes.

## Exercice 2

### Loi binomiale

A chaque balade qu'il effectue, un cavalier a une probabilité  $p$  d'être désarçonné.

1. Quelle est la probabilité que le cavalier ait chuté  $k$  fois au terme de  $n$  balades ? On suppose que les différentes promenades sont indépendantes les unes des autres.
2. Quelle est la loi du nombre de chutes en  $n$  balades ?
3. Donner l'espérance et la variance du nombre de chutes en  $n$  balades.

## Exercice 3

### Loi hypergéométrique

Le responsable qualité d'une usine contrôle 20 objets dans chaque lot de 1000 objets avant de le laisser partir vers le client. Il accepte seulement les lots pour lesquels il ne trouve aucun objet non conforme dans l'échantillon ; dans le cas contraire, le lot est trié unité par unité.

1. Si  $p\%$  des pièces fabriquées sont défectueuses, quelle est la probabilité d'en trouver  $k$  dans un lot donné de taille 20 ?
2. Quelle est la probabilité pour qu'un lot contenant une proportion  $p = 0,05$  d'objets non conformes soit accepté ?
3. Même question pour  $p = 0,1$ .

## Exercice 4

### La moyenne empirique

Soient  $X_1, X_2, \dots, X_n$   $n$  variables aléatoires indépendantes et identiquement distribuées (*i.i.d.*) de moyenne  $m$  et de variance  $\sigma^2$ . La moyenne empirique est :  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ . Calculer  $E(\bar{X})$  et  $V(\bar{X})$ .

## Exercice 5

### Intervalle de confiance pour une moyenne

On a mesuré le rendement de 100 parcelles de blé d'une variété donnée. On a obtenu  $\frac{1}{100} \sum_{i=1}^{100} x_i = 86$  et  $\sum_{i=1}^{100} x_i^2 = 750000$  où  $x_i$  exprime le rendement observé sur la  $i^{\text{ème}}$  parcelle (en qx/ha).

On suppose que les rendements sont mutuellement indépendants et qu'ils sont issus d'une population infinie distribuée selon une loi normale de moyenne  $m$  et de variance  $\sigma^2$ .

Construire un intervalle de confiance pour le rendement moyen au niveau de confiance 95%.

Pour préserver l'anonymat dans certaines enquêtes par sondage, le procédé suivant peut être suivi. Admettons que l'on veuille estimer la proportion de personnes qui remplissent leur déclaration fiscale de manière honnête. On demande alors à chaque personne interrogée de se retirer dans une pièce isolée, et de jouer à pile ou face.

- si elle obtient « pile » alors elle doit répondre honnêtement par « oui » ou « non » à la question « Votre déclaration fiscale est-elle honnête ? »
- si elle obtient « face », elle devra lancer la pièce une nouvelle fois et répondre par « oui » ou « non » à la question « Avez-vous obtenu « face » au deuxième tirage ? ».

Grâce à ce procédé, il est impossible à l'enquêteur de savoir à quelle question se rapporte la réponse de la personne interrogée, celle-ci peut donc fournir sans crainte une réponse sincère.

1. On note  $p$  la proportion inconnue de déclarations fiscales remplies honnêtement dans la population et  $\pi$  la proportion de réponses « oui ». Montrer que  $\pi = p/2 + 1/4$ .
2. Soit  $X$  la variable aléatoire désignant le nombre de réponses « oui » dans une enquête auprès de  $n$  personnes. Quelle est la loi de  $X$ ? Donner un estimateur de  $\pi$  et un estimateur de  $p$ . Calculer leur espérance et variance respectives.
3. En déduire un intervalle de confiance de niveau  $1 - \alpha$  pour  $p$ . On utilisera l'approximation normale de la loi binomiale.
4. Application numérique avec  $n = 1000$  et 600 réponses affirmatives. Donner une estimation de  $p$  et un intervalle de confiance pour  $p$  au niveau 95%. Quel est le prix payé pour la confidentialité ?

## Quelques rappels sur les lois de probabilité

### Variable aléatoire $X$

C'est une grandeur qui peut prendre différentes valeurs avec différentes probabilités. Elle est définie sur l'ensemble des résultats possibles (ou événements) d'une expérience aléatoire (ex : résultat d'un jeu de hasard, durée d'attente,...).

### Loi de probabilité

La loi de probabilité, ou distribution, d'une variable aléatoire  $X$  est définie par l'ensemble des valeurs prises par  $X$  ainsi que par :

- la probabilité de chaque valeur possible de  $X$  quand  $X$  est une v.a. discrète,
- la probabilité que  $X$  se réalise dans un intervalle donné quand  $X$  est une v.a. continue. La fonction de densité de  $X$ , dérivée de la fonction de répartition caractérise la loi de probabilité.

### Espérance $E(X)$

C'est la valeur que l'on peut espérer obtenir, en moyenne, en réalisant une v.a.  $X$ . On l'assimile à la moyenne de  $X$  par abus de langage.

Pour une variable aléatoire discrète,  $E(X) = \sum_k k \times P(X = k)$ .

Pour une variable aléatoire continue admettant une densité  $f(x)$ ,  $E(X) = \int_{-\infty}^{+\infty} xf(x)$

#### Propriétés :

- Pour  $c$  constante réelle,  $E(c) = c$
- $E(X + Y) = E(X) + E(Y)$  : on dit que l'espérance est un opérateur linéaire
- Si  $X$  et  $Y$  sont indépendantes alors  $E(XY) = E(X) \times E(Y)$

### Variance $Var(X)$

C'est une mesure de la variabilité des valeurs par rapport à la moyenne. Plus les valeurs de  $X$  sont « imprévisibles », plus elle est grande. Elle se définit par  $Var(X) = \sigma_X^2 = E[X - E(X)]^2 = E(X^2) - [E(X)]^2$  (« moyenne des carrés des écarts à la moyenne »)

#### Propriétés :

- La variance est toujours positive ou nulle
- $Var(X) = 0 \Leftrightarrow X$  constante
- $Var(cX) = c^2 Var(X)$  où  $c$  est une constante réelle
- $Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$ 
  - o  $Cov(X, Y) = \sigma_{XY} = E[X - E(X)] \times E[Y - E(Y)]$
  - o  $Cov(X, Y) = 0$  si  $X$  et  $Y$  sont indépendantes

### Loi de Bernoulli $B(p)$

C'est la loi de la variable  $X$  qui indique si le résultat d'une épreuve est un échec ou un succès (par exemple : jouer à pile ou face).

Loi de probabilité :  $P(X=1)=p$  et  $P(X=0)=1-p$

Espérance :  $E(X)=p$

Variance :  $Var(X)=p(1-p)$

### Loi binomiale $B(n,p)$

C'est la loi de la variable  $X$  qui compte le nombre de boules blanches obtenues à l'issue de  $n$  tirages, indépendants et avec remise, dans une urne de taille  $N$  contenant  $p$  % de boules blanches.

Loi de probabilité :  $P(X=k) = C_n^k p^k (1-p)^{n-k}$  avec  $k \in \{0, 1, \dots, n\}$

Espérance :  $E(X) = np$

Variance :  $Var(X) = np(1-p)$

N.B. : une loi binomiale de paramètres  $n$  et  $p$  est aussi la somme de  $n$  lois de Bernoulli indépendantes et de même paramètre  $p$ .

### Loi hypergéométrique $H(N, n, p)$

C'est la loi de la variable  $X$  qui compte le nombre de boules blanches sélectionnées à l'issue de  $n$  tirages sans remise dans une urne de taille  $N$  contenant des boules blanches en proportion  $p$ .

$$\text{Loi de probabilité : } P(X=k) = \frac{C_{Np}^k C_{N-Np}^{n-k}}{C_N^n} \text{ avec } \max(0, n-(N-Np)) \leq k \leq \min(n, Np)$$

$$\text{Espérance : } E(X) = np$$

$$\text{Variance : } \text{Var}(X) = np(1-p) \frac{N-n}{N-1}$$

### Convergence de la loi hypergéométrique vers la loi binomiale

Si  $N$  tend vers l'infini, la loi  $H(N, n, p)$  tend vers la loi  $B(n, p)$ , c'est-à-dire que lorsqu'on effectue un tirage dans une grande population, il importe peu que ce tirage se fasse avec ou sans remise (en pratique, on considère que la population est « grande » lorsque l'échantillon représente moins de 10% de cette population :  $n/N < 0,1$ ).

### Loi normale ou loi de Laplace-Gauss $N(m, \sigma^2)$

C'est la loi d'une variable  $X$  continue, variant de  $-\infty$  à  $+\infty$ , dont la densité de probabilité vaut :

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-m}{\sigma}\right)^2\right]$$

$$\text{Espérance : } E(X) = m$$

$$\text{Variance : } \text{Var}(X) = \sigma^2$$

### Convergence de la loi binomiale vers la loi normale

$$\text{Si } X \text{ suit une } B(n, p) \text{ et que } n \text{ tend vers l'infini alors } \frac{X - np}{\sqrt{np(1-p)}} \longrightarrow N(0,1)$$

En pratique, on considère que l'approximation est correcte dès que  $np(1-p) > 18$ , d'autant plus que  $n$  est grand et  $p$  proche de 0,5.

### Loi uniforme $U(0,1)$

Une variable  $X$  suit une loi uniforme  $U(0,1)$  si sa densité de probabilité vaut :  $f(x) = 1_{]0,1[}(x)$

$$\text{Espérance : } E(X) = 1/2$$

$$\text{Variance : } \text{Var}(X) = 1/12$$

$$F(x) = P(X \leq x) = x \text{ sur } [0,1]$$

### Loi faible des grands nombres

Si  $(X_1, X_2, \dots, X_n)$  sont des variables indépendantes et identiquement distribuées (i.i.d.) selon une loi

$$\text{quelconque de même moyenne } m, \text{ alors: } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{p} m$$

Autrement dit, la moyenne d'une variable sur un échantillon aléatoire simple tend vers la moyenne dans la population, quand la taille de l'échantillon tend vers l'infini. Par exemple, si l'on pouvait jouer indéfiniment à "pile ou face" avec une pièce bien équilibrée, le pourcentage de "pile" obtenu tendrait vers 50 %.

### Théorème central limite

Si  $(X_1, X_2, \dots, X_n)$  sont des variables i.i.d. selon une loi quelconque de moyenne  $m$  et de variance  $\sigma^2$ ,

$$\text{alors: } \sqrt{n} \frac{\bar{X}_n - m}{\sigma} \xrightarrow[n \rightarrow \infty]{\text{Loi}} N(0,1)$$

# Quelques rappels sur les intervalles de confiance

## I/ Généralités

Soient  $X$  une variable aléatoire de loi paramétrée par  $\theta$  et  $X_1, \dots, X_n$   $n$  variables i.i.d. selon la loi de  $X$ .

### 1) Principe d'un intervalle de confiance

Plutôt que d'estimer ponctuellement la vraie valeur inconnue du paramètre  $\theta$ , on recherche un intervalle recouvrant «très vraisemblablement» cette vraie valeur.

Définition : On appelle intervalle de confiance de niveau de confiance  $1-\alpha$  du paramètre  $\theta$  tout intervalle  $IC$  tel que :  $P(IC \ni \theta) = 1-\alpha$  pour  $\alpha \in [0,1]$  fixé.

Les bornes de l'intervalle de confiance  $IC$  dépendent de l'échantillon, elles sont donc aléatoires.

Par abus de langage, on note souvent  $P(\theta \in IC) = 1-\alpha$ .

Remarquons que si  $\alpha$  augmente (ou que si  $n$  augmente), l'amplitude de l'intervalle de confiance diminue.

### 2) Vocabulaire

La probabilité  $\alpha$  pour que l'intervalle de confiance ne contienne pas la vraie valeur peut être répartie différemment de part et d'autre des bornes de l'intervalle de confiance. Écrivons donc  $\alpha = \alpha_1 + \alpha_2$  où  $\alpha_1$  et  $\alpha_2$  mesurent respectivement les risques à gauche et à droite de dépasser un seuil plancher ou plafond.

- L'intervalle de confiance est dit bilatéral quand  $\alpha_1 \neq 0$  et  $\alpha_2 \neq 0$ . Si  $\alpha_1 = \alpha_2 = \frac{\alpha}{2}$ , l'intervalle est dit symétrique. Il est dissymétrique sinon.
- L'intervalle de confiance est dit unilatéral si  $\alpha_1 \alpha_2 = 0$  :
  - quand on veut assurer une valeur minimale au paramètre à estimer, on considère  $\alpha_1 = \alpha$  et  $\alpha_2 = 0$ , l'intervalle de confiance est alors de la forme :  $IC = [a, +\infty[$ .
  - quand on ne veut pas dépasser un seuil maximal, on prend  $\alpha_1 = 0$  et  $\alpha_2 = \alpha$  et on obtient alors un intervalle de confiance de la forme :  $IC = ]-\infty, b]$ .

### 3) Construction

Pour construire un intervalle de confiance, on utilise une variable aléatoire dont on connaît la distribution de probabilité.

Définition : une fonction pivotale pour le paramètre  $\theta$  est une fonction des observations  $(X_1, \dots, X_n)$  et du paramètre  $\theta$  dont la loi ne dépend pas du paramètre  $\theta$ .

On recherche dans la suite des fonctions pivotales particulières adaptées aux cas étudiés.

## II/ Intervalles de confiance pour l'espérance

On envisage deux cas :

- la variable aléatoire mesurée est normale et le nombre de réalisations est quelconque,
- la variable aléatoire mesurée n'est pas normale et le nombre de réalisations est important. Dans ce cas, la distribution de la moyenne empirique tend vers une loi normale d'après le théorème central limite. On parlera d'intervalle de confiance asymptotique.

Dans la suite on considère  $X \sim N(m, \sigma^2)$  et  $X_1, \dots, X_n$   $n$  variables i.i.d. selon la loi de  $X$ .

On définit la moyenne empirique  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  et la variance empirique modifiée

$$S_n'^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

### **1) Cas où la variance est connue**

Après centrage et réduction de la moyenne empirique, on obtient :  $\sqrt{n} \frac{\bar{X}_n - m}{\sigma} \sim N(0,1)$

On a :  $P\left(-u \leq \sqrt{n} \frac{\bar{X}_n - m}{\sigma} \leq u\right) = 1 - \alpha$  où  $u$  est le fractile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $N(0,1)$ .

Ce qui revient à :  $P\left(\bar{X}_n - u \frac{\sigma}{\sqrt{n}} \leq m \leq \bar{X}_n + u \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha.$

Quand la variance est connue, l'intervalle de confiance bilatéral symétrique pour l'espérance d'une loi normale s'écrit donc au niveau  $1 - \alpha$  sous la forme suivante :

$$\boxed{IC(m) = \left[ \bar{x}_n - u \frac{\sigma}{\sqrt{n}}, \bar{x}_n + u \frac{\sigma}{\sqrt{n}} \right]} \quad \bar{x}_n \text{ est la réalisation de } \bar{X}_n \text{ sur l'échantillon.}$$

Remarque : si  $\alpha = 5\%$ , le fractile d'ordre 0,975 de la loi normale centrée réduite correspond à 1,96.  
si  $\alpha = 10\%$ , le fractile d'ordre 0,95 de la loi normale centrée réduite vaut environ 1,64.

### **2) Cas où la variance est inconnue**

On a :  $\sqrt{n} \frac{\bar{X}_n - m}{S_n'} \sim St(n-1)$  (loi de Student à  $n-1$  degrés de libertés).

d'où  $P\left(-t \leq \sqrt{n} \frac{\bar{X}_n - m}{S_n'} \leq t\right) = 1 - \alpha$  où  $t$  est le fractile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $St(n-1)$

et donc  $P\left(\bar{X}_n - t \frac{S_n'}{\sqrt{n}} \leq m \leq \bar{X}_n + t \frac{S_n'}{\sqrt{n}}\right) = 1 - \alpha.$

Quand la variance est inconnue, l'intervalle de confiance bilatéral symétrique pour l'espérance d'une loi normale s'écrit donc au niveau  $1 - \alpha$  sous la forme suivante :

$$\boxed{IC(m) = \left[ \bar{x}_n - t \frac{s_n'}{\sqrt{n}}, \bar{x}_n + t \frac{s_n'}{\sqrt{n}} \right]} \quad \bar{x}_n \text{ et } s_n' \text{ sont les réalisations respectives de } \bar{X}_n \text{ et } S_n' \text{ sur l'échantillon.}$$

Remarque : quand  $n \rightarrow \infty$ , on approxime la loi de Student par la loi normale centrée réduite. On retrouve alors le cas précédent.

### 3) Cas particulier : intervalle de confiance pour une proportion

Soient  $X_1, \dots, X_n$  i.i.d. selon  $B(p)$  et  $X = \sum_{i=1}^n X_i \sim B(n, p)$ . Notons  $F_n = \frac{X}{n}$  estimateur sans biais de  $p$ .

- Dans le cas de grands échantillons :

En approchant une loi binomiale vers une loi normale, on a :  $\sqrt{n} \frac{F_n - p}{\sqrt{p(1-p)}} \xrightarrow[n \rightarrow \infty]{\text{loi}} N(0,1)$ .

Ce qui permet d'écrire :  $P\left(-u \leq \sqrt{n} \frac{F_n - p}{\sqrt{p(1-p)}} \leq u\right) = 1 - \alpha$  où  $u$  est le fractile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $N(0,1)$ .

Et donc l'intervalle de confiance bilatéral symétrique pour une proportion  $p$  au niveau  $1 - \alpha$

s'obtient en résolvant l'inéquation :  $\left| \sqrt{n} \frac{F_n - p}{\sqrt{p(1-p)}} \right| \leq u$

Ce qui donne en notant  $f_n$  la réalisation de  $F_n$  sur l'échantillon:

$$IC(p) = \left[ \frac{f_n + \frac{u^2}{2n} - \frac{u}{\sqrt{n}} \sqrt{\frac{u^2}{4n} + f_n(1-f_n)}}{1 + \frac{u^2}{n}}, \frac{f_n + \frac{u^2}{2n} + \frac{u}{\sqrt{n}} \sqrt{\frac{u^2}{4n} + f_n(1-f_n)}}{1 + \frac{u^2}{n}} \right]$$

Pour une taille d'échantillon importante, on considère l'approximation suivante :

$$IC(p) = \left[ f_n - u \sqrt{\frac{f_n(1-f_n)}{n}}, f_n + u \sqrt{\frac{f_n(1-f_n)}{n}} \right]$$

Cette approximation est parfaitement justifiée sur le plan théorique.

En effet, d'après le théorème de Slutsky, on a :  $\sqrt{F_n(1-F_n)} \xrightarrow{p} \sqrt{p(1-p)}$ .

On en déduit donc que :  $\sqrt{n} \frac{F_n - p}{\sqrt{F_n(1-F_n)}} \xrightarrow[n \rightarrow \infty]{\text{loi}} N(0,1)$ .

D'où :  $P\left(-u \leq \sqrt{n} \frac{F_n - p}{\sqrt{F_n(1-F_n)}} \leq u\right) = 1 - \alpha$  où  $u$  est le fractile d'ordre  $1 - \frac{\alpha}{2}$  de la loi  $N(0,1)$ .

Quand  $n$  est grand, l'intervalle de confiance bilatéral symétrique pour une proportion s'écrit donc au niveau  $1 - \alpha$  sous la forme :

$$IC(p) = \left[ f_n - u \sqrt{\frac{f_n(1-f_n)}{n}}, f_n + u \sqrt{\frac{f_n(1-f_n)}{n}} \right] \quad f_n \text{ est la réalisation de } F_n \text{ sur l'échantillon.}$$

- Sinon, construction d'intervalles de confiance « exacts » :

On construit ces intervalles en considérant la fonction de répartition de la loi binomiale. Si la probabilité de recouvrement de l'intervalle ne vaut pas exactement  $1 - \alpha$ , on prend l'intervalle ayant la plus petite probabilité de recouvrement parmi ceux ayant une probabilité de recouvrement supérieure à  $1 - \alpha$ .

### III/ Intervalles de confiance pour la variance d'une loi normale

Soient  $X \sim N(m, \sigma^2)$  et  $X_1, \dots, X_n$   $n$  variables i.i.d. selon la loi de  $X$ .

#### 1) Cas où l'espérance est connue

Soit  $S_n^{*2} = \frac{1}{n} \sum_{i=1}^n (X_i - m)^2$ . On a  $n \frac{S_n^{*2}}{\sigma^2} \sim \chi^2(n)$

D'où  $P\left(\chi^2_{\frac{\alpha_1}{2}} \leq n \frac{S_n^{*2}}{\sigma^2} \leq \chi^2_{1-\frac{\alpha_2}{2}}\right) = 1 - \alpha$  où  $\chi^2_{\alpha_1}$  est le fractile d'ordre  $\alpha_1$  de la loi  $\chi^2(n)$ ,  
et  $\chi^2_{1-\alpha_2}$  est le fractile d'ordre  $1 - \alpha_2$  de la loi  $\chi^2(n)$ .

Quand l'espérance est connue, l'intervalle de confiance bilatéral pour la variance d'une loi normale s'écrit donc au niveau  $1 - \alpha$  sous la forme suivante :

$$IC(\sigma^2) = \left[ n \frac{S_n^{*2}}{\chi^2_{1-\frac{\alpha_2}{2}}}, n \frac{S_n^{*2}}{\chi^2_{\frac{\alpha_1}{2}}} \right]$$

$S_n^*$  est la réalisation de  $S_n^*$  sur l'échantillon.

Remarque : cet intervalle n'est pas centré car la loi du khi-deux n'est pas symétrique.

#### 2) Cas où l'espérance est inconnue

On considère la variance empirique modifiée  $S_n'^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  comme fonction pivotale pour  $\sigma^2$ .

On sait que  $\frac{(n-1)S_n'^2}{\sigma^2} \sim \chi^2(n-1)$ .

On a donc  $P\left(\chi^2_{\frac{\alpha_1}{2}} \leq (n-1) \frac{S_n'^2}{\sigma^2} \leq \chi^2_{1-\frac{\alpha_2}{2}}\right) = 1 - \alpha$  où  $\chi^2_{\alpha_1}$  est le fractile d'ordre  $\alpha_1$  de la loi  $\chi^2(n-1)$   
et  $\chi^2_{1-\alpha_2}$  le fractile d'ordre  $1 - \alpha_2$  de la loi  $\chi^2(n-1)$ .

Quand l'espérance est inconnue, l'intervalle de confiance bilatéral pour la variance d'une loi normale s'écrit donc au niveau  $1 - \alpha$  sous la forme suivante :

$$IC(\sigma^2) = \left[ (n-1) \frac{S_n'^2}{\chi^2_{1-\frac{\alpha_2}{2}}}, (n-1) \frac{S_n'^2}{\chi^2_{\frac{\alpha_1}{2}}} \right]$$

$S_n'$  est la réalisation de  $S_n'$  sur l'échantillon.

## II. Sondage aléatoire simple

### Exercice 1

### Un petit exemple

L'exercice propose de retrouver sur un exemple les résultats de la théorie pour un sondage aléatoire simple sans remise de taille fixe. On considère pour cela tous les échantillons possibles de taille 2 pris dans une population de taille  $N = 5$ . On connaît par ailleurs les valeurs de la variable d'intérêt  $Y$  pour chaque unité de la population, à savoir respectivement : 8, 3, 11, 4 et 7.

1. Calculer la moyenne  $\bar{Y}$  et la dispersion  $S_Y^2$  du caractère d'intérêt sur la population.
2. Lister tous les échantillons possibles de taille 2.
3. Pour chacun de ces échantillons, calculer l'estimateur  $\hat{Y}$  de la moyenne de la variable d'intérêt ainsi que l'estimateur de sa variance  $\hat{V}(\hat{Y})$ .
4. Vérifier que  $\hat{Y}$  estime sans biais la vraie moyenne.
5. Calculer la variance  $V(\hat{Y})$ .
6. Vérifier que  $V(\hat{Y})$  coïncide avec la formule de la variance donnée par la théorie.
7. Vérifier que  $\hat{V}(\hat{Y})$  estime sans biais la vraie variance  $V(\hat{Y})$ .

### Exercice 2

### Rappels de cours

L'exercice propose de démontrer des résultats présentés dans le cours et d'insister sur des techniques de raisonnement usuelles en sondage. Considérons qu'on veuille estimer le total et la moyenne d'une grandeur  $Y$  dans une population  $U$  de taille  $N$ . Pour cela, on procède à un sondage aléatoire simple sans remise de taille  $n$  et on note  $S$  l'échantillon aléatoire obtenu.

1. Combien y a-t-il d'échantillons possibles ? Quelle est la probabilité de tirer chacun d'entre eux ?
2. On considère un individu  $k$  quelconque dans  $U$ . Combien y a-t-il d'échantillons contenant cet individu ? En déduire la probabilité de tirage de  $k$ .
3. On note  $I_k$  la variable aléatoire valant 1 si  $k$  appartient à l'échantillon et 0 sinon.
  - a. Que vaut  $E(I_k)$  ?
  - b. Comment peut-on réécrire  $\sum_{k \in S} Y_k$  à partir des  $I_k$  ?
4. En déduire que :
  - a.  $\hat{t}_y = \frac{N}{n} \sum_{k \in S} Y_k$  estime sans biais le vrai total  $t_y = \sum_{k \in U} Y_k$
  - b. et que  $\hat{Y} = \frac{1}{n} \sum_{k \in S} Y_k$  estime sans biais la vraie moyenne  $\bar{Y} = \frac{1}{N} \sum_{k \in U} Y_k$ .
5. Combien y a-t-il d'échantillons comprenant les individus identifiés  $k$  et  $l$  ? En déduire la probabilité de tirer ces deux individus conjointement. Que vaut alors  $E(I_k I_l)$  ? En déduire  $Cov(I_k, I_l)$ .

6. On note  $S_y^2 = \frac{1}{N-1} \sum_{k \in U} (Y_k - \bar{Y})^2$  et  $f = \frac{n}{N}$ . Montrer que :

a.  $Var(\hat{t}_y) = N(N-n) \frac{S_y^2}{n}$

b.  $Var(\hat{Y}) = (1-f) \frac{S_y^2}{n}$

7. Quel est l'intérêt du sondage sans remise par rapport au sondage avec remise ?

8. Montrer que  $s^2 = \frac{1}{n-1} \sum_{k \in S} (Y_k - \hat{Y})^2$  estime sans biais  $S_y^2$ .

9. En déduire des estimateurs sans biais de  $Var(\hat{t}_y)$  et de  $Var(\hat{Y})$ .

**Exercice 3** Estimation de la surface agricole utile d'un canton  
(d'après P.Ardilly et Y.Tillé, Exercices corrigés de méthode de sondage, Ellipses, 2003)

On veut estimer la surface moyenne cultivée dans les fermes d'un canton rural. Sur 2010 fermes que comprend ce canton, on en tire 100 par sondage aléatoire simple. On mesure  $Y_k$  la surface cultivée par la ferme  $k$  en hectares et on trouve :

$$\sum_{k \in S} Y_k = 2907 \text{ ha} \text{ et } \sum_{k \in S} Y_k^2 = 154\,593 \text{ ha}^2$$

1. Donner la valeur de l'estimateur sans biais classique de la moyenne  $\bar{Y} = \frac{1}{N} \sum_{k \in U} Y_k$ .
2. Donner un intervalle de confiance à 95% pour  $\bar{Y}$ .

**Exercice 4** Estimation d'une retombée touristique  
(d'après A-M. Dussaix et J-M. Grosbras, Exercices de sondage, Economica, 1992)

145 ménages de touristes séjournant en France dans une région donnée ont dépensé 830 € en moyenne par jour. L'écart type estimé de leurs dépenses s'élève à 210 €. Sachant que 50 000 ménages de touristes ont visité la région où a été effectuée l'enquête, que peut-on dire de la dépense totale journalière de l'ensemble de ces ménages ? On supposera pour cela que l'échantillon est issu d'un plan aléatoire simple à probabilités égales.

**Exercice 5** Taille d'échantillon pour un sondage d'opinion  
(d'après A-M. Dussaix et J-M. Grosbras, Exercices de sondage, Economica, 1992)

Un sondage sur la popularité d'une personnalité politique lui accorde un pourcentage  $\hat{p} = 30\%$  d'opinions favorables. En admettant qu'il s'agisse d'un sondage aléatoire simple sans remise et que la taille de l'échantillon est négligeable au regard de celle de la population, combien de personnes ont-elles été interrogées pour que l'on puisse dire avec un degré de confiance de 95% que la vraie proportion d'opinions favorables dans la population ne s'écarte pas de  $\hat{p}$  de plus de deux points ?

**Exercice 6****Taille d'échantillon pour une proportion***(d'après P.Ardilly et Y.Tillé, Exercices corrigés de méthode de sondage, Ellipses, 2003 )*

On s'intéresse à l'estimation de la proportion  $P$  d'individus atteints par une maladie professionnelle dans une entreprise de 1500 salariés. On sait par ailleurs que trois personnes sur dix sont ordinairement touchées par cette maladie dans des entreprises du même type. On se propose de sélectionner un échantillon au moyen d'un sondage aléatoire simple.

1. Quelle taille d'échantillon faut-il sélectionner pour que la longueur totale d'un intervalle de confiance avec un niveau de confiance 0,95 soit inférieure à 0,01 pour un plan simple :
  - a. avec remise ?
  - b. sans remise ?
2. Que faire dans le cas du plan sans remise si on ne connaît pas la proportion d'individus habituellement touchés par la maladie ?

**Exercice 7****Nombre d'espaces de stationnement à prévoir***(d'après A-M. Dussaix et J-M. Grosbras, Exercices de sondage, Economica, 1992 )*

Une entreprise de promotion immobilière désire estimer le nombre d'espaces de stationnement requis pour une nouvelle tour devant abriter des bureaux. Elle décide de procéder à un sondage aléatoire simple sans remise. Elle sait que le nouveau bâtiment abritera 5 000 personnes et que, dans des entreprises de même type que celles devant emménager dans les futurs locaux, la proportion de personnes se rendant à leur bureau en utilisant les moyens de transport en commun est toujours supérieure à 75%. Quelle doit être la taille de l'échantillon pris au sein des futurs occupants des bureaux pour pouvoir estimer le nombre d'espaces de stationnement à prévoir avec une marge d'erreur symétrique d'au plus 150 places au niveau de confiance 90% ?

**Exercice 8****Application au marketing direct***(d'après A-M. Dussaix et J-M. Grosbras, Exercices de sondage, Economica, 1992 )*

Les sondages sont très largement utilisés dans le marketing direct : il arrive souvent que l'on estime par sondage le rendement d'un fichier donné, ou que l'on souhaite comparer les rendements de plusieurs fichiers, ou encore, que disposant de plusieurs fichiers, on souhaite estimer par sondage le rendement global de l'ensemble de ces fichiers. Dans cet exercice, on suppose l'existence d'un fichier de  $N = 200\ 000$  adresses. On note  $p$  le rendement inconnu du fichier à une offre d'abonnement à prix réduit avec calculatrice offerte en prime ; c'est donc la proportion d'individus qui s'abonneraient si l'offre était offerte à tous les individus du fichier. Selon l'usage  $\hat{p}$  est l'estimation de  $p$  obtenue à partir d'un test fait sur un échantillon de  $n$  adresses choisies à probabilités égales et sans remise sur le fichier.

1. On sait par expérience que les rendements à ce type d'offre sur ce fichier ne dépassent pas généralement 3%. Quelle taille d'échantillon doit-on prendre pour estimer  $p$  avec une précision absolue de 0,5 point et un degré de confiance de 95% ?
2. Mêmes questions pour une précision de 0,3 point et 0,1 point.
3. Le test a porté sur 10 000 adresses et on a noté 230 abonnements. En déduire l'intervalle de confiance bilatéral à 95% pour le rendement  $p$  ainsi que le pour le nombre total d'abonnements si la même offre était faite sur l'ensemble du fichier.

**Rappel :** on appelle **précision absolue** au niveau de confiance  $1-\alpha$  la quantité  $t_{1-\frac{\alpha}{2}} \sqrt{V(\hat{p})}$  où  $t_{1-\frac{\alpha}{2}}$

est le fractile d'ordre  $1 - \frac{\alpha}{2}$  de la loi normale centrée réduite.

**Exercice 9**

**Un cas d'enquête répétée**

(d'après P.Ardilly et Y.Tillé, Exercices corrigés de méthode de sondage, Ellipses, 2003 )

On considère une population de 10 stations-services et on s'intéresse au prix du litre de supercarburant que chacune d'entre elles affiche. Plus exactement, sur deux mois consécutifs, mai et juin, les données de prix figurent dans le tableau ci-dessous :

Prix du litre de supercarburant

Station	1	2	3	4	5	6	7	8	9	10
Mai	5,82	5,33	5,76	5,98	6,20	5,89	5,68	5,55	5,69	5,81
Juin	5,89	5,34	5,92	6,05	6,20	6,00	5,79	5,63	5,78	5,84

On veut estimer l'évolution du prix moyen du litre entre mai et juin. On choisit, comme indicateur de cette évolution la différence des prix moyens. On propose deux méthodes concurrentes:

- **Méthode 1** : on échantillonne  $n$  stations ( $n < 10$ ) en mai et  $n$  stations en juin, les deux échantillons étant totalement indépendants ;
  - **Méthode 2** : on échantillonne  $n$  stations en mai, et on interroge de nouveau ces stations en juin (technique de *panel*).
1. Comparer l'efficacité des deux méthodes.
  2. Même question si on souhaite cette fois estimer un prix moyen sur la période globale mai-juin.
  3. Si on s'intéresse au prix moyen de la question 2, ne vaut-il pas mieux tirer, non pas 2 fois  $n$  relevés avec la méthode 1 ( $n$  chaque mois) mais directement  $2n$  relevés sans se soucier des mois (méthode 3) ? Aucun calcul n'est nécessaire.

**Exercice 10**

**Échantillonnages successifs**

En cours de collecte, la taille d'un échantillon s'avère parfois insuffisante pour assurer la précision attendue. Une solution naturelle est d'enquêter un échantillon complémentaire. Intéressons-nous au plan de sondage final obtenu après :

- Un premier échantillonnage simple sans remise de  $n_1$  unités parmi  $N$  à probabilités égales,
- Suivi d'un second tirage simple sans remise de  $n_2$  unités parmi  $N-n_1$  à probabilités égales

La sélection des  $n = n_1 + n_2$  unités ainsi retenues obéit-elle à un plan simple sans remise et à probabilités égales dans la population de taille  $N$ ?

**Exercice 11**

**Estimation dans un domaine**

On souhaite estimer la moyenne et le total d'une variable  $y$  sur un domaine  $U_0$  d'une population finie  $U$  de taille  $N$ . Ces quantités sont notées :

$$t_{y0} = \sum_{k \in U_0} Y_k \quad \text{et} \quad \bar{Y}_0 = \frac{t_{y0}}{N_0} = \frac{1}{N_0} \sum_{k \in U_0} Y_k$$

où  $N_D$  est la taille du domaine.

On sélectionne un échantillon  $s$  au sein de la population entière par un sondage aléatoire simple sans remise de taille  $n$ . On observe un sous-échantillon  $s_0$  de taille  $n_0$  dont les individus sont dans le domaine  $U_0$ .

On dispose des deux estimateurs suivants de la moyenne et du total de  $y$  sur le domaine  $U_0$  :

$$1. \hat{t}_{y0} = \frac{N}{n} \sum_{k \in s_0} Y_k \quad \text{et} \quad \hat{Y}_0 = \frac{\hat{t}_{y0}}{N_0} = \frac{N}{N_0 \cdot n} \sum_{k \in s_0} Y_k$$

$$2. \hat{Y}_0 = \frac{1}{n_0} \sum_{k \in s_0} Y_k \quad \text{et} \quad \hat{t}_{y0} = N_0 \cdot \hat{Y}_0 = \frac{N_0}{n_0} \sum_{k \in s_0} Y_k$$

- La taille  $n_0$  du sous-échantillon  $s_0$  est aléatoire. Calculer sa valeur moyenne .
- Montrer que les deux estimateurs du total (ou de la moyenne) sont tous deux sans biais pour le vrai total (ou la vraie moyenne) du domaine. Est-ce que l'un est préférable à l'autre ?
- Donner les expressions de variance des deux estimateurs de la moyenne. Comparer ces deux variances.
- Donner les estimateurs sans biais pour les variance de ces deux estimateurs.
- Exemple : considérons une population de  $N = 5\,793$  entreprises. Supposons connues les quantités suivantes :

$$N_0 = 984, \quad \sum_{k \in U_0} Y_k = 154814, \quad \sum_{k \in U_0} Y_k^2 = 42148912$$

où  $y$  désigne le chiffre d'affaires.

Calculer les vraies variance pour les deux estimateurs de la moyenne pour un échantillon de taille  $n = 579$ .

- On a observé sur un échantillon particulier de taille  $n = 579$

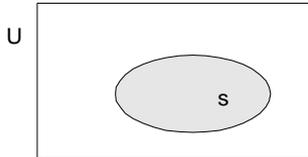
$$n_0 = 89, \quad \sum_{i \in s_0} y_i = 13782, \quad \sum_{i \in s_0} y_i^2 = 4530306$$

Donner les valeurs des deux estimateurs de la moyenne et calculer les valeurs de leur variance estimée.

# Rappels sur le sondage aléatoire simple

## I/ Définition

Tirage d'un échantillon de  $n$  unités sans remise et à probabilités égales dans une population finie composée de  $N$  unités identifiables.



## II/ Notations

1. **Dans la population (ou univers)**  $U = \{1, 2, \dots, k, \dots, N\}$

- Variable d'intérêt :  $Y$  de caractéristique individuelle  $Y_k$
- Total :  $T_Y = \sum_{k \in U} Y_k$
- Moyenne :  $\bar{Y} = \frac{T_Y}{N} = \frac{1}{N} \sum_{k \in U} Y_k$
- Variance :  $\sigma_y^2 = \frac{1}{N} \sum_{k \in U} (Y_k - \bar{Y})^2$
- Dispersion (variance modifiée) :  $S_y^2 = \frac{1}{N-1} \sum_{k \in U} (Y_k - \bar{Y})^2 = \frac{N}{N-1} \sigma_y^2$

2. **Dans l'échantillon**  $s$  : sous-ensemble de  $U$  de taille  $n(s)$

- Ensemble des échantillons possibles :  $S$
- Plan de sondage probabiliste : loi de probabilité sur  $S$

$$p(s) \geq 0, \forall s \in S, \text{ et } \sum_{s \in S} p(s) = 1.$$

- Moyenne :  $\hat{y} = \frac{1}{n} \sum_{k \in s} Y_k$
- Dispersion empirique :  $\hat{s}_y^2 = \frac{1}{n-1} \sum_{k \in s} (Y_k - \hat{y})^2$
- Probabilité d'inclusion d'ordre un de  $k$  :  $\pi_k = P(k \in s) = \sum_{s \in S / k \in s} p(s)$
- Probabilité d'inclusion ou double de  $k$  et  $l$  :  $\pi_{kl} = P(k \in s, l \in s) = \sum_{s \in S / k, l \in s} p(s)$
- $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$

### III/ Formulaire du sondage aléatoire simple

Probabilité de sélectionner l'échantillon  $s$  :  $p(s) = 1/C_N^n$

Probabilité de sélectionner l'individu  $k$  :  $\forall k \in U, \pi_k = P(k \in s) = \frac{n}{N} = f$  (taux de sondage)

Paramètre d'intérêt / Statistique	Moyenne	Proportion $p = N_0/N$	Total
Estimateur du paramètre d'intérêt	$\hat{y} = \frac{1}{n} \sum_{k \in S} Y_k = \hat{y}(s)$	$\hat{p} = \frac{1}{n} \sum_{k \in S} y_k = \frac{n_0}{n}$	$\hat{t}_y = N \times \hat{y} = \frac{N}{n} \sum_{k \in S} Y_k$
Vraie variance d'échantillonnage de cet estimateur	$Var(\hat{y}) = \left(1 - \frac{n}{N}\right) \frac{S_y^2}{n}$	$Var(\hat{p}) = \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \frac{p(1-p)}{n}$	$Var(\hat{t}_y) = N^2 \left(1 - \frac{n}{N}\right) \frac{S_y^2}{n}$
Estimateur de la variance d'échantillonnage	$\hat{Var}(\hat{y}) = \left(1 - \frac{n}{N}\right) \frac{\hat{S}_y^2}{n}$	$\hat{Var}(\hat{p}) = \left(1 - \frac{n}{N}\right) \frac{\hat{p}(1-\hat{p})}{n-1}$	$\hat{Var}(\hat{t}_y) = N^2 \left(1 - \frac{n}{N}\right) \frac{\hat{S}_y^2}{n}$

Intervalle au niveau de confiance 95% pour la moyenne :

$$IC_{95\%}(\bar{Y}) = \left[ \hat{y} - 1,96\sqrt{\hat{Var}(\hat{y})}, \hat{y} + 1,96\sqrt{\hat{Var}(\hat{y})} \right]$$

sous hypothèse que  $n$  est grand  $\frac{\hat{y} - \bar{Y}}{\sqrt{Var(\hat{y})}} \rightarrow N(0, 1)$

### III. PLANS À PROBABILITÉS INÉGALES

#### Exercice 1 Rappels de cours sur l'estimateur d'Horvitz-Thompson

On considère une population  $U$  et on s'intéresse à l'estimation du total d'une variable d'intérêt  $Y$  noté  $t_y = \sum_{k \in U} Y_k$ . Pour cela, on prélève un échantillon  $s$  avec des probabilités individuelles de sélection notées  $(\pi_k)_{k \in U}$ .

1. Rappeler l'expression de l'estimateur d'Horvitz-Thompson (ou «  $\pi$ -estimateur » ou encore « estimateur des valeurs dilatées »).
2. Étudier son espérance et sa variance.

#### Exercice 2 Application directe du cours

On considère une population  $U = \{1,2,3\}$ , sur laquelle on définit le plan de sondage suivant :

$$p(\{1,2\}) = \frac{1}{2}, p(\{1,3\}) = \frac{1}{4}, p(\{2,3\}) = \frac{1}{4}$$

$Y$  est une variable définie sur  $U$ , telle que :  $Y_1 = Y_2 = 3, Y_3 = 6$  dont on veut estimer le total  $t_y$ .

1. Calculer les probabilités d'inclusion simple  $\pi_k$  et double  $\pi_{k\ell}$ .
2. Donner la distribution de probabilité de l'estimateur de Horvitz-Thompson  $\hat{t}_{y\pi}$  du total. Calculer la variance de cet estimateur.
3. Donner la distribution de probabilité d'un estimateur de variance de  $\hat{t}_{y\pi}$  (il est conseillé de choisir l'estimateur le plus simple à calculer). On pourra vérifier que cet estimateur est sans biais.

#### Exercice 3 Volume d'archives

On désire estimer à l'échelle d'un canton le nombre de kilomètres linéaires d'archives stockées dans les mairies. Pour cela, on procède à un tirage de 4 communes parmi les 9 du canton, proportionnellement à leur population.

1. Calculer les probabilités d'inclusion de chaque communes, à partir des données suivantes :

N° de commune	Nom de la commune	Population
1	Val le Grand	1100
2	Les Gries	650
3	Les Combres	500
4	Flins	2300
5	Villers le Lac	4000
6	Fortin	5500
7	Montlebon	1900
8	Sanzeau	200
9	Aumont	150

2. Estimer le métrage total des archives du canton à partir des résultats suivants :

N° de commune	Nom de la commune	Mètres d'archives
2	Les Gries	17
4	Flins	38
5	Villers le Lac	55
6	Fortin	70

**Exercice 4****Tirage systématique d'entreprises**

On veut sélectionner un échantillon de taille 4 dans une population de 8 entreprises dont on connaît la taille, mesurée en termes d'effectif salarié. L'échantillon est tiré à probabilités proportionnelles à la taille.

Entreprise	1	2	3	4	5	6	7	8
Taille	300	300	150	100	50	50	25	25

1. Donner les probabilités d'inclusion d'ordre 1 des entreprises.
2. Sélectionner l'échantillon selon un tirage systématique en utilisant 0,27 comme nombre aléatoire ;
3. Lister les échantillons possibles que l'on peut obtenir avec un tirage systématique, et indiquer les probabilités de tirage de chacun d'eux.
4. A partir des échantillons obtenus, donner une estimation du total de l'effectif salarié des entreprises. Le résultat était-il prévisible ?
5. Calculer la matrice des probabilités d'inclusion d'ordre 2 ? Commenter.

**Exercice 5****Tirage de Poisson**

(d'après P.Ardilly et Y.Tillé, *Exercices corrigés de méthode de sondage*, Ellipses, 2003)

Lorsqu'on effectue des tirages à probabilités inégales, on utilise en général des méthodes d'échantillonnage de taille fixe. Il existe cependant des algorithmes très simples permettant des tirages à probabilités inégales mais conférant à l'échantillon une taille variable. On s'intéresse ici au tirage de Poisson dont le principe consiste à effectuer une loterie sur chaque individu de la population indépendamment d'un individu à l'autre. Ainsi, pour une population de taille  $N$  où les probabilités d'inclusion individuelles  $\pi_k$  sont connues pour tout  $k$ , on simule  $N$  aléas indépendants dans la loi uniforme sur  $[0,1]$  et on retient l'individu  $k$  si et seulement si  $u_k \leq \pi_k$

1. Vérifier que l'algorithme de tirage respecte les probabilités d'inclusion d'ordre 1 en calculant la probabilité pour que l'individu  $k$  soit sélectionné.
2. La taille de l'échantillon est une variable aléatoire notée  $n_S$ .
  - a. Écrire  $n_S$  en fonction des variables indicatrices de Cornfield.
  - b. Que vaut l'espérance et la variance de  $n_S$  ?
  - c. Quelle est la probabilité pour que l'échantillon ait une taille au moins égale à 1 ?

On supposera dans la suite que l'échantillon a une taille au moins égale à 1.

3. On utilise l'estimateur du total  $\hat{Y} = \sum_{k \in S} \frac{Y_k}{\pi_k}$  où  $S$  désigne l'échantillon aléatoire obtenu à l'issue des  $N$  loteries.
  - a. Vérifier que  $\hat{Y}$  estime le vrai total sans biais.
  - b. Quelle est la variance de  $\hat{Y}$  ? Comment peut-on l'estimer sans biais ?
  - c. Que valent les probabilités d'inclusion d'ordre 2 ?
4. Comparer à un plan général de taille fixe  $n$  de mêmes probabilités d'inclusion. Quelles sont les inconvénients d'un plan de taille non-fixe ?

## Rappels sur les plans à probabilités inégales

### I/ Intérêt

Retenir de préférence les unités les plus porteuses d'information afin d'accroître la précision.

### III/ Formulaire

**Probabilité de sélectionner l'individu k :**

- Pour un plan à probabilités proportionnelles à une variable X de taille (corrélée positivement à Y)

$$\forall k \in U, \pi_k = P(k \in S) = n \frac{X_k}{\sum_{k \in U} X_k}$$

- Pour un plan de taille fixe,  $\sum_{k \in U} \pi_k = n$

Statistique \ Paramètre d'intérêt	Moyenne	Total
<b>Estimateur d'Horvitz-Thompson du paramètre d'intérêt (<math>\pi</math>-estimateur)</b>	<p>Si la taille N est connue :</p> $\hat{\mu}_{y\pi} = \frac{1}{N} \sum_{k \in S} \frac{Y_k}{\pi_k} = \frac{\hat{t}_{y\pi}}{N}$ <p>Sinon, estimateur de Hájek :</p> $\hat{\mu}_{yH} = \frac{1}{\hat{N}_\pi} \sum_{k \in S} \frac{Y_k}{\pi_k} = \frac{\sum_{k \in S} Y_k}{\sum_{k \in S} \frac{1}{\pi_k}} = \frac{\hat{t}_{y\pi}}{\hat{N}_\pi}$	$\hat{t}_{y\pi} = \sum_{k \in S} \frac{Y_k}{\pi_k}$ <p>En particulier : <math>\hat{N}_\pi = \sum_{k \in S} \frac{1}{\pi_k}</math></p>
<b>Vraie variance d'échantillonnage de cet estimateur</b>	<p>Cas général</p> $Var(\hat{\mu}_{y\pi}) = \frac{1}{N^2} \sum_{k \in U} \sum_{l \in U} \frac{Y_k}{\pi_k} \frac{Y_l}{\pi_l} \Delta_{kl}$ <p>Si la taille de l'échantillon est fixe</p> $Var(\hat{\mu}_{y\pi}) = \frac{1}{2N^2} \sum_{k \in U} \sum_{l \in U} \left( \frac{Y_k}{\pi_k} - \frac{Y_l}{\pi_l} \right)^2 \Delta_{kl}$	<p>Cas général :</p> $Var(\hat{t}_{y\pi}) = \sum_{k \in U} \sum_{l \in U} \frac{Y_k}{\pi_k} \frac{Y_l}{\pi_l} \Delta_{kl}$ <p>Si la taille de l'échantillon est fixe</p> $Var(\hat{t}_{y\pi}) = \frac{1}{2} \sum_{k \in U} \sum_{l \in U} \left( \frac{Y_k}{\pi_k} - \frac{Y_l}{\pi_l} \right)^2 \Delta_{kl}$
<b>Estimateur de la variance d'échantillonnage</b>	<p>Cas général</p> $\hat{Var}(\hat{\mu}_{y\pi}) = \frac{1}{N^2} \sum_{k \in S} \sum_{l \in S} \frac{Y_k}{\pi_k} \frac{Y_l}{\pi_l} \frac{\Delta_{kl}}{\pi_{kl}}$ <p>Si la taille de l'échantillon est fixe</p> $\hat{Var}(\hat{\mu}_{y\pi}) = \frac{1}{2N^2} \sum_{k \in S} \sum_{l \in S} \left( \frac{Y_k}{\pi_k} - \frac{Y_l}{\pi_l} \right)^2 \frac{\Delta_{kl}}{\pi_{kl}}$	<p>Cas général</p> $\hat{Var}(\hat{t}_{y\pi}) = \sum_{k \in S} \sum_{l \in S} \frac{Y_k}{\pi_k} \frac{Y_l}{\pi_l} \frac{\Delta_{kl}}{\pi_{kl}}$ <p>Si la taille de l'échantillon est fixe</p> $\hat{Var}(\hat{t}_{y\pi}) = \frac{1}{2} \sum_{k \in S} \sum_{l \in S} \left( \frac{Y_k}{\pi_k} - \frac{Y_l}{\pi_l} \right)^2 \frac{\Delta_{kl}}{\pi_{kl}}$

Si  $n$  est grand, l'intervalle de confiance pour la moyenne au niveau de confiance  $1 - \alpha$  est :

$$IC_{1-\alpha}(\mu_y) = \left[ \hat{\mu}_{y\pi} - u_{1-\frac{\alpha}{2}} \sqrt{\hat{Var}(\hat{\mu}_{y\pi})}; \hat{\mu}_{y\pi} + u_{1-\frac{\alpha}{2}} \sqrt{\hat{Var}(\hat{\mu}_{y\pi})} \right]$$

où  $u_{1-\frac{\alpha}{2}}$  désigne le fractile d'ordre  $1 - \alpha/2$  de la loi  $N(0,1)$

## IV. TP1 : SIMULATIONS DE TIRAGE D'ÉCHANTILLONS

### Objectifs de la séance

- Utiliser différents algorithmes de tirages d'échantillons pour des plans simples sans remise et des plans à probabilités inégales ;
- Évaluer le paramètre d'intérêt et la précision de cette estimation ;
- Valider de manière empirique certaines propriétés de la théorie des sondages ;
- Comparer les méthodes d'échantillonnage.

### Données utilisées

La population étudiée est celle des 771 communes rurales d'Île-de-France recensées en 1999. On cherche à estimer le nombre total d'habitants résidant dans ces communes ainsi que le nombre moyen d'habitants par commune. Les données datent des recensements de 1999 et de 1990.

### Partie I : Tirage d'un échantillon

On cherche à échantillonner 100 communes en raisonnant successivement à probabilités égales puis à probabilités inégales, proportionnellement à la population recensée en 1990. Sélectionner un tel échantillon en utilisant les différents algorithmes suivants :

- 1) Tirage de Bernoulli ;
- 2) Méthode du tri aléatoire ;
- 3) Méthode de sélection-rejet ;
- 4) Tirage de Poisson ;
- 5) Tirage systématique ;
- 6) Algorithme de Sunter.

### Partie II : Simulations

- 1) On choisit d'abord d'échantillonner les communes selon un plan simple sans remise.
  - a. Sélectionner 100 échantillons de taille 50. Pour chaque échantillon, estimer le paramètre d'intérêt ainsi que la variance d'échantillonnage.
  - b. Vérifier empiriquement l'absence de biais de l'estimateur de la moyenne.
  - c. Tracer la distribution de l'estimateur de la moyenne et commenter.
  - d. Vérifier empiriquement l'absence de biais de l'estimateur de la variance d'échantillonnage.
- 2) On choisit maintenant de sélectionner les communes proportionnellement à leur taille, mesurée en nombre d'habitants recensés en 1990.
  - a. Sélectionner 100 échantillons de taille 50. Pour chaque échantillon, estimer le paramètre d'intérêt.
  - b. Vérifier empiriquement l'absence de biais de l'estimateur de la moyenne.
  - c. Tracer la distribution de l'estimateur de la moyenne.
- 3) Comparer les deux plans de sondage.

Le choix du logiciel est libre. A toutes fins utiles, la suite de l'énoncé propose deux modes d'emploi :

- l'un sous Excel (des macros pré-programmées sont mises à disposition),
- l'autre sous SAS qui appelle aux procédures SURVEYSELECT et SURVEYMEANS.

## Mode d'emploi sous Excel

### La base de sondage et le catalogue de macros



TP1.xls

### Entrée

La base de sondage est décrite dans l'onglet « *BS* ». Par commodité, le contenu de cette base se limite à l'identifiant, la variable d'intérêt, voire la variable auxiliaire utile au calcul de probabilités inégales proportionnelles.

### Paramètres

L'utilisateur spécifie le nombre d'échantillons à tirer ainsi que leur taille dans l'onglet « *Paramètres* » prévu à cet effet.

Dans le cas de simulations, un paramètre supplémentaire permet également de spécifier si les tirages systématiques sont à probabilités égales ou inégales.

### Algorithmes pré-programmés

Les macros mises à disposition permettent de sélectionner un ou plusieurs échantillon(s) selon différents algorithmes de tirage. Elles fournissent également les estimations de total et de moyenne de la variable d'intérêt sur l'(les) échantillon(s) obtenu(s). Dans le cas de simulations, elles dressent aussi le bilan de l'ensemble des tirages.

Les algorithmes pré-programmés sont ceux-ci :

- Méthode du tri aléatoire pour un plan simple sans remise (*macro Tri\_aléatoire*) ;
- Méthode de sélection-rejet pour un ou plusieurs plan(s) simple(s) sans remise (*macros Sélection\_rejet et Simulations\_SAS\_SR*) ;
- Tirage de Bernoulli pour un plan à probabilités égales et sans remise (*macro Bernoulli*) ;
- Algorithme de Sunter pour un plan à probabilités inégales, de taille fixe et sans remise (*macro Sunter*) ;
- Tirage systématique pour un ou plusieurs plan(s) à probabilités inégales, de taille fixe et sans remise (*macros Tirage\_systématique et Simulations\_systématique*) ;
- Tirage de Poisson pour un plan à probabilités inégales, sans remise (*macro Poisson*).

### Sorties

Les résultats de chaque macro alimentent un onglet précis. Avant lancement de chaque macro, il convient donc de vérifier la présence de la feuille vierge ad-hoc ainsi que l'absence d'un onglet portant le nom réservé aux sorties. Plus précisément, les onglets réservés par chaque méthode sont :

Algorithme	Nom de l'onglet en entrée	Nom de l'onglet en sortie
Tri aléatoire	Feuil1	Ech.Tri_Aléatoire
Systématique	Feuil2	Ech.Systématique
Sélection-Rejet	Feuil3	Ech. Sélection-Rejet
Sunter	Feuil4	Ech.Sunter
Bernoulli	Feuil5	Ech.Bernoulli
Poisson	Feuil6	Ech. Poisson
Simulation de plans simples sans remise	Feuil7	Simul_SAS_SR
Simulation de plans à probabilités inégales	Feuil8	Simul_Systématique

### Mise en œuvre

1. A l'ouverture du fichier Excel, cliquer sur *Activer les macros* ;
2. Renseigner la feuille « *BS* » en indiquant l'identifiant de chaque unité de la base de sondage en 1ère colonne, la variable d'intérêt en 2ème colonne, voire la variable auxiliaire en 3ème colonne si le plan est à probabilités inégales proportionnelles à cette donnée ;
3. Renseigner les paramètres souhaités dans la feuille « *Paramètres* » ;
4. Vérifier la disponibilité des onglets requis dans le classeur ;
5. Cliquer sur *Outils*, puis *Macro* suivi de *Macros* ;
6. Sélectionner la méthode voulue, puis cliquer sur *Exécuter* pour lancer la macro retenue ;
7. Consulter les résultats dans la feuille correspondante à la méthode choisie.

### Remarques

1. Au 1er lancement, il est conseillé de limiter le nombre de simulations afin de contrôler le temps d'exécution des macros.
2. Pour modifier le contenu des macros,
  - a. Cliquer sur *Modifier* après *Outils* > *Macro* > *Macros*
  - b. Saisir le nouveau code.NB : des commentaires permettent de comprendre le rôle de chaque action.
3. Pour tracer un histogramme, une possibilité est d'utiliser l'utilitaire d'analyse d'Excel. Pour cela, cliquer sur *Outils*, puis *Macro Complémentaire*. Cocher *Utilitaire d'analyse* et valider par *OK*. Ensuite, cliquer sur *Outils*, puis *Utilitaire d'analyse*. Choisir *histogramme* dans le menu déroulant qui s'affiche et suivre les indications.

## Mode d'emploi sous SAS

### La base de sondage



tp1.sas7bdat

### Les procédures SURVEYSELECT et SURVEYMEANS



Procédures SAS  
d'échantillonnage.pdf

## V. PLANS STRATIFIES

### Exercice 1

### Rappels de cours

Dans une population de taille  $N$  partitionnée en  $H$  strates, on sélectionne un échantillon de taille  $n$  suivant un plan stratifié. Dans chaque strate  $h$ , on tire  $n_h$  individus parmi  $N_h$  selon un sondage aléatoire simple sans remise de taille fixe.

Préalable : montrer la formule de décomposition de la variance :

$$\sigma_y^2 = \frac{1}{N} \sum_{k \in U} (Y_k - \bar{Y})^2 = \frac{1}{N} \sum_{h=1}^H N_h \sigma_{yh}^2 + \frac{1}{N} \sum_{h=1}^H N_h (\bar{Y}_h - \bar{Y})^2$$

1. Pour une variable d'intérêt  $Y$ , donner les estimateurs du total  $t_Y$  et de la moyenne.
2. Montrer que ces deux estimateurs sont sans biais et calculer leur variance.
3. On considère l'allocation proportionnelle de l'échantillon : on décide de tirer dans chaque strate  $h$  un nombre d'individus  $n_h$  tel que :

$$\frac{n_h}{N_h} = \frac{n}{N} \text{ (en supposant que } N_h \frac{n}{N} \text{ soit entier).}$$

- a. Comment s'écrivent alors les estimateurs du total et de la moyenne ?
  - b. Que vaut leur variance ?
  - c. Montrer alors, que si on suppose :  $\sigma_y^2 \approx S_y^2$  et  $\sigma_{yh}^2 \approx S_{yh}^2$  pour tout  $h$ , l'allocation proportionnelle est toujours meilleure qu'un sondage aléatoire simple.
4. Le point de vue envisagé maintenant est celui d'une allocation optimale afin de satisfaire un souci de précision. Sous la contrainte que  $\sum_{h=1}^H n_h = n$ ,
    - a. Quelle est l'allocation des  $n_h$  qui minimise la variance de l'estimateur du total ?
    - b. Que vaut alors la variance ?
    - c. Comment peut-on interpréter le choix des allocations optimales ?

### Exercice 2

### Estimation du poids des éléphants d'un cirque

(d'après P.Ardilly et Y.Tillé, Exercices corrigés de méthode de sondage, Ellipses, 2003)

Un directeur de cirque possède 100 éléphants classés en deux catégories : "mâles" et "femelles". Le directeur veut estimer le poids total de son troupeau, car il veut traverser un fleuve en bateau. Il a la possibilité de faire peser seulement 10 éléphants de son troupeau. Cependant, en 1998, ce même directeur a pu faire peser tous les éléphants de son troupeau, et il a obtenu les résultats suivants (en tonnes) :

	Effectif	Moyenne	Variance
Mâles	60	6	4,00
Femelles	40	4	2,25

1. Calculer la variance dans la population de la variable "poids de l'éléphant" en 1998.
2. Si, en 1998, le directeur avait procédé à un sondage aléatoire simple sans remise de 10 éléphants, quelle aurait été la variance de l'estimateur du poids total du troupeau ?
3. Si le directeur avait procédé à un sondage stratifié, avec SAS dans chaque strate, avec allocation proportionnelle de 10 éléphants, quelle aurait été la variance de l'estimateur du poids total du troupeau ?
4. Si le directeur avait procédé à un sondage stratifié optimal, avec SAS dans chaque strate, de 10 éléphants, quels auraient été les effectifs de l'échantillon dans les strates, et quelle aurait été la variance de l'estimateur du poids total du troupeau ?

### Exercice 3

### L'âge du personnel

Une grande entreprise veut réaliser une enquête auprès de son personnel qui comprend 10 000 personnes. Des études préliminaires ont montré :

- que les variables que l'on cherche à analyser dans l'enquête sont très contrastées selon les catégories de personnel et qu'il y a donc intérêt à stratifier selon ces catégories. Pour simplifier, on considérera qu'il y a 3 grandes catégories qui formeront les strates,
- que ces variables sont également très fortement liées à l'âge des individus.

On va donc proposer des plans d'échantillonnage comme si on voulait étudier l'âge des individus : si une stratégie est meilleure que d'autres pour estimer l'âge moyen, alors on a de bonnes raisons de penser qu'elle le sera aussi pour les variables d'intérêt. Comme on connaît l'âge des membres du personnel, on peut raisonner en faisant les comparaisons exactes.

On dispose des renseignements suivants :

Catégorie de personnel	Poids dans l'ensemble du personnel	Écart type des âges
1	20%	18,0
2	30%	12,0
3	50%	3,6
Ensemble	100%	16,0

1. Soit  $\bar{Y}$  l'âge moyen et  $\hat{Y}$  l'estimateur issu d'un échantillon aléatoire simple sans remise à probabilités égales de  $n = 100$  individus. Quelle est l'erreur type de  $\hat{Y}$  ?
2. On décide que l'échantillon de 100 individus doit être stratifié selon les catégories de personnel. Quelle est la répartition « représentative » ? Quelle est l'erreur type de l'estimateur de  $\bar{Y}$  qui en découle ? Comparer avec les résultats de la question 1.
3. Quelle serait la répartition optimale de l'échantillon ? Quelle est l'erreur type de l'estimateur de  $\bar{Y}$  qui en découle ? Comparer avec les résultats de la question 2.

### Exercice 4 proportion

### Estimation d'une

Sur les 7500 employés d'une entreprise, on souhaite connaître la proportion  $p$  d'entre eux qui possèdent au moins un véhicule. Pour chaque individu de la base de sondage, on dispose de la valeur de son revenu. On décide alors de constituer trois strates dans la population : individus de faibles revenus (strate 1), de revenus moyens (strate 2) et de revenus élevés (strate 3).

On note :

- $N_h$  la taille de la strate  $h$ ,
- $n_h$  la taille de l'échantillon dans la strate  $h$ ,
- $\hat{p}_h$  l'estimateur de la proportion d'individus possédant au moins un véhicule dans la strate  $h$ .

On obtient le résultat suivant :

	$h = 1$	$h = 2$	$h = 3$
$N_h$	3500	2000	2000
$n_h$	500	300	200
$\hat{p}_h$	0,13	0,45	0,50

1. Quel estimateur  $\hat{p}$  de  $p$  proposez-vous ? Que peut-on dire de son biais ?
2. Calculez la précision de  $\hat{p}$ , et donnez un intervalle de confiance à 95% pour  $p$ .
3. Estimez-vous que le critère de stratification est adéquat ?

**Exercice 5** **Optimalité pour une différence**  
*(d'après J-M. Grosbras, Méthodes statistiques des sondages, Economica, 1987)*

Le but de l'exercice est de montrer que si une stratégie est optimale pour estimer précisément une quantité dans l'ensemble d'une population stratifiée, elle peut ne plus l'être tout à fait si l'objectif du sondage est justement de comparer les strates entre elles. La bonne définition des objectifs à atteindre est donc essentielle au choix de la technique à employer. Considérons une population de taille  $N$  formée de deux strates, de taille  $N_1$  et  $N_2$  et intéressons-nous à la moyenne  $\bar{X}$  d'une variable  $X$ . Les moyennes de  $X$  dans les strates 1 et 2 sont notées  $\bar{X}_1$  et  $\bar{X}_2$  et leurs estimateurs  $\hat{X}_1$  et  $\hat{X}_2$ .

On dispose d'un budget  $C$  et on suppose que :

- le tirage effectué est un sondage aléatoire simple sans remise de  $n_h$  unités parmi  $N_h$  dans la strate  $h$  ( $h=1$  ou  $2$ ),
- la fonction de coût s'écrit  $C_1 n_1 + C_2 n_2$  où  $C_h$  désigne le coût unitaire dans la strate  $h$ .

1. Si on cherche à estimer précisément la moyenne  $\bar{X}$ ,
  - a. Donner l'expression de  $\hat{X}$ , estimateur sans biais de  $\bar{X}$  en fonction de  $\hat{X}_1$  et  $\hat{X}_2$ .
  - b. Calculer sa variance.
  - c. Quelle répartition  $(n_1, n_2)$  de l'échantillon donne une variance  $V(\hat{X})$  minimale ? Que vaut alors  $V(\hat{X})$  ?
  - d. Application numérique : calculer  $n_1, n_2, n$  et  $V(\hat{X})$  avec :
 

$N_1 = 10\ 000$	$N_2 = 20\ 000$	
$S_1 = 2$	$S_2 = 1$	
$C_1 = 4$	$C_2 = 9$	$C = 1\ 000$
2. Si on avait appliqué une allocation proportionnelle, c'est-à-dire :  $n_h / N_h = n / N$ ,
  - a. Qu'aurait-on trouvé pour  $n_1, n_2$  et  $n$  ?
  - b. Que vaudrait alors  $V(\hat{X})$  ?
  - c. Avec les mêmes données numériques, évaluer la perte relative de précision par rapport à l'échantillon optimal.

3. En fait, on cherche à évaluer l'écart entre les moyennes des deux groupes :  $\bar{X}_1 - \bar{X}_2$ .
- Montrer que  $\hat{X}_1 - \hat{X}_2$  est un estimateur sans biais de  $\bar{X}_1 - \bar{X}_2$ .
  - Calculer sa variance.
  - Déterminer la répartition  $(n_1, n_2)$  de l'échantillon pour que  $V(\hat{X}_1 - \hat{X}_2)$  soit minimale, toujours avec la même contrainte de budget. (on pourra éventuellement utiliser, en les adaptant, certains résultats de la question 1).
  - Calculer dans ces conditions  $V(\hat{X})$ . Comparer ce résultat avec celui de la 1<sup>ère</sup> question en écrivant la différence  $\Delta$  des variances de ces deux estimateurs.
  - Reprendre l'application numérique pour trouver les nouvelles valeurs de  $n_1, n_2, n, V(\hat{X})$  et la perte relative de précision par rapport à l'échantillon optimal.

### Exercice 6

### Choix des allocations

(d'après A-M. Dussaix et J-M. Grosbras, Exercices de sondage, Economica, 1992)

Cet exercice est une application du principe : "à chaque objectif son échantillon". Une entreprise comporte 400 exécutants et 100 cadres. La direction de l'entreprise désire évaluer un indice de satisfaction, assimilable à une variable numérique positive  $Y$ , mesurable pour chaque individu à partir d'un ensemble de questions : elle décide pour cela de faire réaliser une enquête auprès de 100 personnes employées dans l'entreprise, à l'aide d'un plan de sondage stratifié, avec un sondage aléatoire simple dans chaque strate. Le coût d'une interview est le même dans les deux strates.

On pense *a priori* que la dispersion de la variable  $Y$  doit être la même au sein de chacun des deux groupes. Comment répartir l'échantillon entre les deux groupes, selon que l'on vise l'un des objectifs suivants :

- obtenir la meilleure précision possible sur la valeur moyenne de l'indice de satisfaction dans l'entreprise ;
- obtenir la même précision sur la valeur moyenne de l'indice de satisfaction dans chacune des deux catégories ;
- obtenir la meilleure précision possible sur la différence entre les valeurs moyennes de l'indice de satisfaction dans les deux catégories.

### Exercice 7

### Estimation d'une différence

On considère une population  $U$  de taille  $N$  partitionnée en  $H$  strates notées  $U_1 \dots U_h \dots U_H$ , de tailles respectives  $N_1 \dots N_h \dots N_H$ . On note  $\bar{Y}_1 \dots \bar{Y}_h \dots \bar{Y}_H$  les moyennes d'une variable d'intérêt  $Y$  au sein de chaque strate, et  $S_1^2 \dots S_h^2 \dots S_H^2$  les dispersions.

La moyenne de  $Y$  dans la population vaut bien sûr :  $\bar{Y} = \sum_{h=1}^H \frac{N_h}{N} \bar{Y}_h = \sum_{h=1}^H w_h \bar{Y}_h$ .

On réalise un sondage stratifié, avec sondage aléatoire simple sans remise dans chaque strate, de taux de sondage  $f_h = n_h / N_h$ . La taille de l'échantillon total est  $n = \sum_{h=1}^H n_h$ .

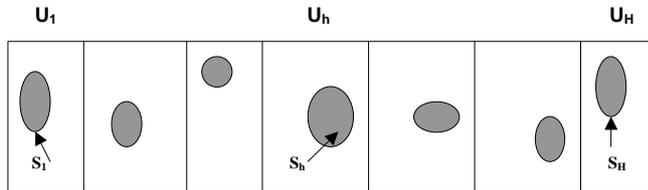
L'objectif est de comparer une strate particulière  $U_i$  à la population totale : on veut estimer  $D_i = \bar{Y}_i - \bar{Y}$

1. Donner l'expression de l'estimateur de Horvitz-Thompson de  $D_i$ , noté  $\hat{D}_i$ , ainsi que l'expression de sa variance.
2. Pour une taille d'échantillon fixée  $n$ , trouver l'allocation optimale  $n_1 \dots n_h \dots n_H$ , qui minimise la variance de  $\hat{D}_i$ . Comparer avec l'allocation optimale de Neyman.

# Rappels sur les plans stratifiés

## I/ Définition

Partition de la population en sous-groupes appelés strates selon un critère lié au paramètre d'intérêt puis tirage d'autant d'échantillons indépendants qu'il y a de strates.



Constituer des strates homogènes en intra au regard de la variable d'intérêt permet de gagner en précision.

## II/ Notations

### 1. Dans la population

- $U = \bigcup_{h=1}^H U_h$  et  $N = \sum_{h=1}^H N_h$
- Total :  $t_y = \sum_{h=1}^H t_{yh} = \sum_{h=1}^H N_h \bar{y}_h$
- Moyenne :  $\bar{y} = \frac{t_y}{N} = \sum_{h=1}^H \frac{N_h}{N} \bar{y}_h$  avec  $\bar{y}_h = \frac{1}{N_h} \sum_{k \in U_h} y_k$
- Variance :  $\sigma_y^2 = \frac{1}{N} \sum_{k \in U} (y_k - \bar{y})^2 = \sum_{h=1}^H \frac{N_h}{N} \sigma_{yh}^2 + \sum_{h=1}^H \frac{N_h}{N} (\bar{y}_h - \bar{y})^2 = \sigma_{y_{intra}}^2 + \sigma_{y_{inter}}^2 = \frac{N-1}{N} S_y^2$   
avec  $\sigma_{yh}^2 = \frac{1}{N_h} \sum_{k \in U_h} (y_k - \bar{y}_h)^2$

### 2. Dans l'échantillon

- $S = \bigcup_{h=1}^H S_h$  et  $n = \sum_{h=1}^H n_h$
- Moyenne dans  $S_h$  :  $\hat{y}_h = \frac{1}{n_h} \sum_{k \in S_h} y_k$
- Dispersion dans  $S_h$  :  $\hat{S}_{yh}^2 = \frac{1}{n_h - 1} \sum_{k \in S_h} (y_k - \hat{y}_h)^2$

III/ Formulaire du sondage stratifié

Paramètre d'intérêt / Statistique	Moyenne	Proportion	Total
<b>Estimateur du paramètre d'intérêt</b>	$\hat{y} = \sum_{h=1}^H \frac{N_h}{N} \hat{y}_h$	$\hat{p} = \sum_{h=1}^H \frac{N_h}{N} \hat{p}_h$	$\hat{t}_y = N\hat{y} = \sum_{h=1}^H \hat{t}_{yh} = \sum_{h=1}^H N_h \hat{y}_h$
<b>Vraie variance d'échantillonnage de cet estimateur</b>	$Var[\hat{y}] = Var\left[\sum_{h=1}^H \frac{N_h}{N} \hat{y}_h\right] = \sum_{h=1}^H Var\left[\frac{N_h}{N} \hat{y}_h\right]$ Si plan simple dans chaque strate : $Var[\hat{y}] = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{yh}^2}{n_h}$	Si plan simple dans chaque strate $Var[\hat{p}] = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{p_h(1-p_h)}{n_h}$	$Var[\hat{t}_y] = Var[N\hat{y}] = N^2 Var[\hat{y}]$ Si plan simple dans chaque strate : $Var[\hat{t}_y] = \sum_{h=1}^H N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{yh}^2}{n_h}$
<b>Estimateur de la variance d'échantillonnage</b>	Si plan simple dans chaque strate $\hat{V}ar[\hat{y}] = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\hat{S}_{yh}^2}{n_h}$	Si plan simple dans chaque strate $\hat{V}ar[\hat{p}] = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\hat{p}_h(1-\hat{p}_h)}{n_h-1}$	Si plan simple dans chaque strate $\hat{V}ar[\hat{t}_y] = \sum_{h=1}^H N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{\hat{S}_{yh}^2}{n_h}$

Intervalle au niveau de confiance 95% pour la moyenne :

$$IC_{95\%}(\bar{Y}) = \left[ \hat{y} - 1,96\sqrt{\hat{V}ar(\hat{y})}, \hat{y} + 1,96\sqrt{\hat{V}ar(\hat{y})} \right]$$

sous hypothèse que  $n$  est grand  $\frac{\hat{y} - \bar{Y}}{\sqrt{\hat{V}ar(\hat{y})}} \rightarrow N(0, 1)$

Choix des allocations

- Allocations proportionnelles :  $\frac{n_h}{n} = \frac{N_h}{N} \quad \forall h \in \{1, \dots, H\}$
- Allocations optimales de Neyman (sans contrainte de budget) :  $n_h = n \frac{N_h S_{yh}}{\sum_{l=1}^H N_l S_{yl}}$
- Allocations optimales sous contrainte budgétaires :  $n_h = C \frac{N_h S_{yh}}{\sqrt{C_h} \sum_{l=1}^H N_l S_{yl} \sqrt{C_l}}$

## VI. PLANS PAR GRAPPES

### Exercice 1

### Problématique d'un plan par grappes

L'objet de cet exercice est de rappeler le formulaire établi en cours et de revenir sur les notions d'effet de sondage et d'effet de grappe.

Un sondage en grappes se pratique sur une population partitionnée en groupes d'individus appelés « grappes » : il consiste à sélectionner certaines grappes, selon un plan quelconque, et à retenir tous les individus des grappes désignées dans l'échantillon final. Procéder de la sorte permet de réduire les coûts d'enquête. On s'intéresse ici au cas particulier où  $m$  grappes sont choisies par sondage aléatoire simple sans remise parmi les  $M$  grappes de taille  $N_i$  d'une population de taille  $N$ .

On cherche à estimer le total  $t_y$  et la moyenne  $\bar{y}$  sur la population d'un caractère d'intérêt  $Y$ .

#### 1. Partie 1 : généralités

- 1.1. Quelle est la probabilité pour qu'un individu appartienne à l'échantillon ?
- 1.2. Que pouvez-vous dire de la taille finale de l'échantillon ? Même question si toutes les grappes sont de même taille  $N_0$ .
- 1.3. Quels estimateurs sans biais  $\hat{t}_y$  et  $\hat{\bar{y}}$  proposez-vous ?
  - 1.3.1. Quelle est la précision de ces estimateurs ?
  - 1.3.2. Montrez que dans le cas où les grappes sont de même taille alors on obtient

$$Var(\hat{\bar{y}}) = \frac{M-m}{M-1} \frac{\sigma_{yinter}^2}{m}.$$

- 1.3.3. En déduire comment constituer les grappes pour obtenir des résultats précis.
- 1.4. Comment estimez-vous sans biais la précision des estimateurs du total et de la moyenne ?
- 1.5. Dans le cas où  $N$  est inconnue, quel estimateur de  $\bar{y}$  proposez-vous ? Cet estimateur est-il sans biais ? Approcher son espérance et son erreur quadratique moyenne.

#### 2. Partie 2 : effet de sondage

On souhaite caractériser la précision de l'échantillonnage par grappes par rapport au sondage aléatoire simple de même taille dans le cas où les grappes sont d'effectifs égaux  $N_0$ .

- 2.1. Montrez que l'effet de sondage défini par  $Deff = \frac{Var(\hat{\bar{y}})}{Var_{sas}(\hat{\bar{y}})}$  vaut  $N_0 \eta^2$  où  $\eta^2$  désigne le

$$\text{rapport de corrélation « inter-grappes » : } \eta^2 = \frac{\sum_{i=1}^M N_0 (\bar{Y}_i - \bar{Y})^2}{\sum_{i=1}^M \sum_{k=1}^{N_0} (Y_k - \bar{Y})^2} = \frac{\sigma_{yinter}^2}{\sigma_y^2}$$

- 2.2. En déduire quand le plan par grappes est plus précis que le sondage aléatoire simple.

#### 3. Partie 3 : effet de grappe

On définit le coefficient de corrélation « intra-grappes » par :

$$\rho = \frac{\sum_{i=1}^M \sum_{k=1}^{N_0} \sum_{l=1, l \neq k}^{N_0} (Y_k - \bar{Y})(Y_l - \bar{Y})}{(N_0 - 1)(N - 1)S_y^2}.$$

Ce coefficient mesure l'effet de grappe. Il se rapproche de 1 si à l'intérieur de chaque grappe, il n'y a pas de différence entre les individus ; au contraire, il est négatif si les individus sont très disparates à l'intérieur de leurs grappes.

3.1. Montrez que l'effet de grappe vaut :

$$\rho = \frac{1}{N_0 - 1} \left[ N_0 \frac{\sigma_{y^{inter}}^2}{\sigma_y^2} - 1 \right]$$

3.2. En déduire que  $Deff = 1 + \rho(N_0 - 1)$  et que  $Var(\hat{y}) = \left(1 - \frac{n}{N}\right) \frac{1}{n} S_y^2 [1 + \rho(N_0 - 1)]$ .

#### 4. Partie 4 subsidiaire: estimation de l'effet de sondage et de l'effet de grappe

On cherche à estimer l'effet de sondage et l'effet de grappe et donc à estimer sans biais  $Var_{sas}(\hat{y})$  autrement dit la dispersion  $S_y^2$ . Les grappes sont de même taille.

4.1. Montrez que la dispersion empirique observée sur l'échantillon  $s_y^2 = \frac{1}{n-1} \sum_{k \in S} (y_k - \hat{y})^2$  possède un biais sous un plan complexe de taille fixe et à probabilités égales (comme ici avec des grappes de même taille) donné par :

$$E[s_y^2] = \frac{n}{n-1} [\sigma_y^2 - Var(\hat{y})]$$

4.2. En déduire que l'expression  $\hat{Deff} = \frac{Var(\hat{y})}{\left(1 - \frac{n}{N}\right) \frac{s_y^2}{n}}$  est justifiée si  $n$  est assez grand.

#### Exercice 2

#### Nombre de signataires d'une pétition

(Extrait de Cochran, Sampling Technics)

On a collecté des signatures pour une pétition sur 676 feuilles. Sur chacune d'entre elles, il y a la place pour 42 signatures, mais beaucoup ne sont pas très remplies. Le nombre de signatures par feuille a été étudié sur un échantillon de 50 feuilles (à peu près 7% de l'ensemble donc). A partir des résultats sont consignés dans le tableau ci-contre, estimer le nombre total de signatures et donner un intervalle de confiance pour ce nombre à 95% et à 80% .

Nombre de signatures	Fréquence
42	23
41	4
36	1
32	1
29	1
27	2
23	1
19	1
16	2
15	2
14	1
11	1
10	1
9	1
7	1
6	3
5	2
4	1
3	1

**Exercice 3****Sélection d'îlots***(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)*

L'objectif est d'estimer le revenu moyen des ménages dans un arrondissement d'une ville composée de 60 îlots de maisons (un îlot est « un pâté de maison », de taille variable). Pour cela, on sélectionne 3 îlots par sondage aléatoire simple sans remise et on interroge tous les ménages qui y résident. On sait en outre que 5 000 ménages résident dans cet arrondissement. Le résultat est donné dans le tableau ci-dessous.

1. Estimez le revenu moyen et le revenu total des ménages de l'arrondissement par l'estimateur d'Horvitz-Thompson.
2. Estimez sans biais la variance de l'estimateur d'Horvitz-Thompson de la moyenne.
3. Estimez le revenu moyen des ménages de l'arrondissement par le ratio de Hájek, et comparez à l'estimation issue de 1. Le sens de variation était-il prévisible ?

Numéro de l'îlot	Nombre de ménages dans l'îlot	Revenu total des ménages de l'îlot
1	120	2100
2	100	2000
3	80	1500

**Exercice 4****Emprunts bancaires**

Une société bancaire structurée en 3 980 succursales gère 39 800 clients, à raison de 10 clients par agence. On choisit 40 succursales par sondage aléatoire simple sans remise pour lesquelles on compte le nombre de clients ayant obtenu un prêt durant une période donnée.

On note  $t_{yi}$  le nombre obtenu dans la succursale  $i$  et on observe :  $\sum_{i=1}^{40} t_{yi} = 185$  et  $\sum_{i=1}^{40} t_{yi}^2 = 1263$ .

1. Estimer le nombre total de clients de la banque qui ont obtenu un prêt durant la période de référence ainsi que leur proportion dans l'ensemble de la clientèle. On notera ces estimateurs  $\hat{t}_y$  et  $\hat{p}$ .
2. Calculer la variance des estimateurs  $\hat{t}_y$  et  $\hat{p}$ .
3. Estimer ces variances et fournir un intervalle de confiance approché à 95% pour chacune des quantités estimées.
4. Calculer l'effet de sondage défini comme le ratio mesurant la perte de variance estimée par rapport à un sondage aléatoire simple sans remise de même taille (indication : on commencera par estimer la dispersion  $S_y^2$ ). On pourra commenter le résultat en comparant les amplitudes des intervalles de confiance à 95% obtenus pour la proportion d'intérêt entre les deux plans de sondage.
5. Calculer le coefficient de corrélation intra-grappe.

**Exercice 5****Influence de la taille et du nombre de grappes échantillonnées***(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)*

Un statisticien souhaite réaliser une enquête sur la qualité des soins assurés dans les services de cardiologie des hôpitaux. Pour cela, il tire par sondage aléatoire simple 100 hôpitaux parmi les 1 000 hôpitaux répertoriés, puis, dans chacun des hôpitaux tirés, il recueille l'avis de tous les malades du service de cardiologie.

1. Comment se nomme ce plan de sondage et quelle est sa raison d'être ?

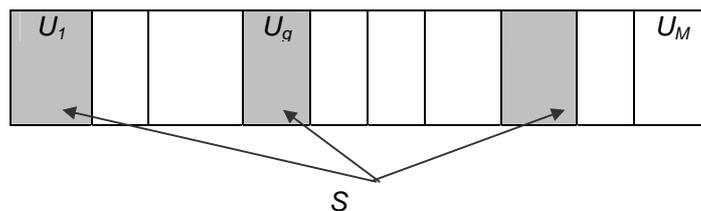
2. On considère que chaque service de cardiologie comprend exactement 50 lits et que l'intervalle de confiance à 95% sur la vraie proportion  $P$  de malades insatisfaits est :  $P \in [0,10 \pm 0,018]$ , (cela signifie en particulier que, dans l'échantillon, 10 % des malades sont insatisfaits de la qualité des soins). Comment estimez-vous l'effet de grappe (commencer par estimer  $S_y^2$ , dispersion du caractère d'intérêt sur toute la population) ?
3. Le statisticien se demande comment évoluerait la précision de son enquête de satisfaction si, d'un seul coup, il échantillonnait deux fois plus d'hôpitaux mais que dans chaque hôpital tiré, il ne collectait ses données que sur la moitié du service de cardiologie (mettons que les services soient systématiquement partagés par un couloir et que notre statisticien ne s'intéresse exclusivement qu'aux 25 lits qui se situent à droite du couloir) ?
4. Commentez ce résultat par rapport à ce que donnait le premier plan de sondage.

## Rappels sur les plans par grappes

### I/ Définition

Objectif principal : réduire les coûts d'enquête et/ou de pallier le manque d'une base de sondage.

Principe : partition de la population en sous-groupes appelés grappes, puis tirage de grappes et enfin recensement de toutes les unités qui les composent.



Règle : constituer des grappes hétérogènes en intra au regard de la variable d'intérêt.

### II/ Notations

#### 1. Dans la population $U$ constituée de $M$ grappes et $N$ individus

- $U = \bigcup_{g=1}^M U_g$  et  $N = \sum_{g=1}^M N_g$
- $t_y = \sum_{g=1}^M t_{yg} = \sum_{g=1}^M N_g \bar{y}_g$
- $\bar{y} = \frac{t_y}{N} = \sum_{g=1}^M \frac{N_g}{N} \bar{y}_g$  avec  $\bar{y}_g = \frac{1}{N_g} \sum_{k \in U_g} y_k$
- $S_G^2 = \frac{1}{M-1} \sum_{g=1}^M \left( t_{yg} - \frac{t_y}{M} \right)^2$

#### 2. Dans l'échantillon $S$ constitué de $m$ grappes et $n_s$ individus

- $S = \bigcup_{g \in S_G} U_g$  et  $n_s = \sum_{g \in S_G} N_g$

III/ Formulaire du plan par grappe dans le cas d'un plan simple de grappes

Paramètre d'intérêt / Statistique	Total	Moyenne
Estimateur du paramètre d'intérêt	$\hat{t}_y = \frac{M}{m} \sum_{g \in S_G} t_{yg}$	$\hat{y} = \frac{1}{N} \hat{t}_y = \frac{M}{Nm} \sum_{g \in S_G} N_g \bar{y}_g$
Vraie variance d'échantillonnage de cet estimateur	$Var[\hat{t}_y] = M^2 \left(1 - \frac{m}{M}\right) \frac{1}{m} \frac{1}{M-1} \sum_{g=1}^M \left(t_{yg} - \frac{t_y}{M}\right)^2$	$Var[\hat{y}] = \frac{1}{N^2} Var[\hat{t}_y]$
Estimateur de la variance d'échantillonnage	$\hat{Var}[\hat{t}_y] = M^2 \left(1 - \frac{m}{M}\right) \frac{1}{m} \frac{1}{m-1} \sum_{g \in S_G} \left(t_{yg} - \frac{\hat{t}_y}{M}\right)^2$	$\hat{Var}[\hat{y}] = \frac{1}{N^2} \hat{Var}[\hat{t}_y]$

Intervalle au niveau de confiance 95% pour la moyenne :

$$IC_{95\%}(\bar{Y}) = \left[ \hat{y} - 1,96\sqrt{\hat{Var}(\hat{y})}, \hat{y} + 1,96\sqrt{\hat{Var}(\hat{y})} \right]$$

sous hypothèse que la taille de l'échantillon est assez grande.

## VII. PLANS À PLUSIEURS DEGRÉS

### Exercice 1

### Probabilités d'inclusion et plans de sondage

(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

On considère une population  $U = \{1,2,3,4,5,6,7,8,9\}$ , sur laquelle on définit le plan de sondage suivant :

$$\begin{aligned} p(\{1,2\}) &= \frac{1}{6}, p(\{1,3\}) = \frac{1}{6}, p(\{2,3\}) = \frac{1}{6} \\ p(\{4,5\}) &= \frac{1}{12}, p(\{4,6\}) = \frac{1}{12}, p(\{5,6\}) = \frac{1}{12} \\ p(\{7,8\}) &= \frac{1}{12}, p(\{7,9\}) = \frac{1}{12}, p(\{8,9\}) = \frac{1}{12} \end{aligned}$$

1. Calculer les probabilités d'inclusion simple  $\pi_k$ .
2. Ce plan de sondage est-il simple, stratifié, en grappes, à deux degrés, ou aucun de ces plans particuliers?

### Exercice 2

### Rappels de cours

Considérons une population de taille  $N$  répartie en  $M$  unités primaires elles-mêmes quadrillées en  $N_i$  unités secondaires. Le premier degré de tirage consiste à extraire un échantillon d'unités primaires parmi lesquelles, dans un second degré de tirage, sont sélectionnées des unités secondaires. Les individus des unités secondaires désignées composent l'échantillon final. Par exemple, si les UP quadrillent le territoire selon un découpage en communes, elles-mêmes composées d'US définies à partir des îlots (ou « pâtés de maisons »), alors l'enquête sera limitée géographiquement aux communes et îlots sélectionnés.

Dans la suite, on considérera le cas où les UP sont choisies selon un sondage aléatoire simple sans remise de taille  $m$  et où les US sont tirées dans les UP retenues au 1<sup>er</sup> degré selon un plan simple sans remise de taille  $n_i$  parmi  $N_i$ . On s'intéresse au total  $t_y$  d'un caractère d'intérêt  $Y$ .

1. Quelle est l'expression de  $\hat{t}_y$  estimateur sans biais de  $t_y$  ?
2. Donner l'expression de la variance de  $\hat{t}_y$  et interpréter les différents termes de ce calcul.
3. Comment estime-t-on cette variance ?
4. Que pouvez-vous dire de la taille finale de l'échantillon ?

### Exercice 3

### Estimation d'un effectif

Un camion transporte des vis sur 500 palettes, chacune d'elles contenant 40 boîtes de vis. L'industriel réceptionnant ces palettes souhaite estimer le nombre moyen de vis par boîte. Pour cela, il tire un échantillon de 100 palettes, selon un sondage aléatoire simple sans remise, puis il tire dans chacune de ces 100 palettes un échantillon de 5 boîtes, selon un sondage aléatoire simple sans remise également, et enfin il compte le nombre de vis dans les boîtes ainsi tirées.

L'industriel, et néanmoins statisticien, calcule pour chaque palette  $i$  de son échantillon le nombre moyen de vis par boîte, et la dispersion du nombre de vis par boîte (ces deux quantités sont calculées à partir des 5 boîtes échantillonnées dans la palette).

Il calcule ensuite les moyennes, sur les 100 palettes, de ces deux quantités :

- moyenne du nombre moyen de vis par boîte = 50
- moyenne de la dispersion du nombre de vis par boîte = 455.

Il calcule aussi la dispersion des 100 estimations du nombre de vis par palette et obtient 375 000.

1. Donner un estimateur sans biais du nombre moyen de vis par boîte.
2. Donner la précision de cet estimateur.
3. Donnez un intervalle de confiance à 95% pour le nombre moyen de vis par boîte.

#### Exercice 4 Nombre de caractères par enregistrement (Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

Sur un disque dur de micro-ordinateur, on compte 400 fichiers, chacun comprenant exactement 50 enregistrements. Pour estimer le nombre moyen de caractères par enregistrement, on décide de tirer par sondage aléatoire simple 80 fichiers, puis 5 enregistrements dans chaque fichier. On note  $m = 80$  et  $\bar{n} = 5$ .

On mesure après tirage :

- la dispersion des estimateurs du nombre total de caractères par fichier, soit  $s_f^2 = 905000$ ,
- la moyenne des  $m$  dispersions  $s_i^2$  est égale à 805 où  $s_i^2$  représente la dispersion du nombre de caractères par enregistrement dans le fichier  $i$ .

1. Comment estimez-vous le nombre moyen  $\bar{Y}$  de caractères par enregistrement ?
2. Comment estimez-vous sans biais la précision de l'estimateur précédent ?
3. Donnez un intervalle de confiance à 95% pour  $\bar{Y}$ .

#### Exercice 5 Étude d'impact préalable au lancement d'un produit financier

En vue de préparer le lancement d'un nouveau produit financier, une société bancaire ayant un réseau de  $M$  succursales souhaite mener une étude approfondie auprès de particuliers possesseurs de comptes chez elle. Les variables d'intérêt de l'enquête ont trait aux caractéristiques de la clientèle et à ses motivations éventuelles. On cherche à estimer la proportion  $p$  de personnes potentiellement intéressées par la nouvelle offre. L'enquête opère selon un plan à 2 degrés : dans un premier temps, on choisit  $m$  succursales pour participer à l'opération parmi lesquelles, au second temps, on désigne des échantillons de titulaires de comptes à interroger. Le plan de sondage est le suivant :

- Au premier degré, on réalise un sondage aléatoire simple sans remise de  $m = 10$  succursales parmi  $M = 100$ . Le taux de sondage  $f_1$  vaut 0,10. La société bancaire gère  $N = 100\ 000$  titulaires de compte.
- Au second degré de tirage, le taux de sondage  $f_2$  est uniforme à 10%.

1. Donner un estimateur sans biais de  $p$  qu'on notera  $\hat{p}$ .

2. Montrer que  $V(\hat{p}) \cong \left(\frac{M}{N}\right)^2 (1 - f_1) \frac{S_T^2}{m} + \frac{1 - f_2}{N f_1 f_2} \sum_{i=1}^M \frac{N_i}{N} p_i (1 - p_i)$

3. Montrer que  $\hat{V}(\hat{p}) \cong \left(\frac{M}{N}\right)^2 (1 - f_1) \frac{\hat{s}_T^2}{m} + \frac{1 - f_2}{N f_1 f_2} \sum_{i \in S_1} \frac{N_i}{N} \hat{p}_i (1 - \hat{p}_i)$

4. Application numérique : donner un intervalle de confiance à 95% pour  $p$  avec les résultats

d'enquête suivants :  $\sum_{k \in S} y_k = 102$ ,  $s_T^2 = 1200$ ,  $\sum_{i \in S_1} \frac{N_i}{N} \hat{p}_i (1 - \hat{p}_i) = 0,01$

### Exercice 6

### Choix entre méthodes concurrentes

Une population de 1010 saucisses est partitionnée en deux unités primaires, de tailles respectives 1000 et 10. Pour estimer le nombre moyen de bouts de saucisses dans cette population, on emploie le plan de sondage suivant :

- on sélectionne une UP selon un sondage aléatoire simple,
- on sélectionne deux saucisses dans l'UP tirée selon un sondage aléatoire simple sans remise.

La première UP est sélectionnée. On observe que chacune des deux saucisses tirées dans l'UP possède deux bouts.

Le statisticien A calcule le nombre moyen de bouts sur son échantillon de deux saucisses et trouve 2. Il affirme que cette valeur est une estimation sans biais du nombre moyen de bouts dans la population.

Le statisticien B propose comme estimation sans biais de ce nombre moyen de bouts la valeur :

$$\frac{1000}{1010} \times 4 = 3.96$$

Discuter les deux méthodes d'estimation, en précisant les logiques qui les sous-tendent.

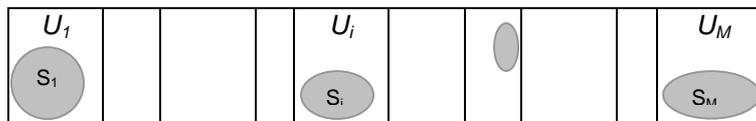
## Rappels sur les plans à deux degrés

### I/ Définition

Objectif principal : réduire les coûts d'enquête et/ou de pallier le manque d'une base de sondage.

Principe : dans une population partitionnée en sous-groupes appelés unités primaires, eux-mêmes composés d'unités secondaires :

- au 1<sup>er</sup> degré, tirage d'unités primaires
- au 2<sup>nd</sup> degré, tirage d'unités secondaires dans les unités primaires retenues au 1<sup>er</sup> degré (indépendamment d'une unité primaire à l'autre)



Règle : constituer des unités primaires hétérogènes en intra au regard de la variable d'intérêt.

### II/ Notations

#### 1. Dans la population U constituée de M unités primaires et N individus

- $U = \bigcup_{i=1}^M U_i$  et  $N = \sum_{i=1}^M N_i$
- $t_y = \sum_{i=1}^M t_{yi} = \sum_{i=1}^M N_i \bar{y}_i$
- $\bar{y} = \frac{t_y}{N} = \sum_{i=1}^M \frac{N_i}{N} \bar{y}_i$  avec  $\bar{y}_i = \frac{1}{N_i} \sum_{k \in U_i} y_k$
- $S_I^2 = \frac{1}{M-1} \sum_{i=1}^M \left( t_{yi} - \frac{t_y}{M} \right)^2$  et  $S_i^2 = \frac{1}{N_i-1} \sum_{k=1}^{N_i} (y_k - \bar{y}_i)^2$

#### 2. Dans l'échantillon S constitué de m unités primaires et $n_s$ individus

- $S = \bigcup_{i \in SUP} S_i$  et  $n_s = \sum_{i \in S_i} n_i$
- $\hat{S}_I^2 = \frac{1}{m-1} \sum_{i \in SUP} \left( \hat{t}_{yi} - \frac{\hat{t}_y}{M} \right)^2$  et  $\hat{S}_i^2 = \frac{1}{n_i-1} \sum_{k \in S_i} (y_k - \hat{y}_i)^2$

III/ Formulaire du plan à deux degrés dans le cas d'un plan simple des unités primaires et des unités secondaires

Paramètre d'intérêt / Statistique	Total	Moyenne
Estimateur du paramètre d'intérêt	$\hat{t}_y = \sum_{k \in S} \frac{y_k}{\pi_k} = \frac{M}{m} \sum_{i \in SUP} \hat{t}_{yg} = \frac{M}{m} \sum_{i \in SUP} \frac{N_i}{n_i} \sum_{k \in Si} y_k$	$\hat{y} = \frac{1}{N} \hat{t}_y = \frac{M}{Nm} \sum_{g \in SG} N_g \hat{y}_g$
Vraie variance d'échantillonnage de cet estimateur	$Var[\hat{t}_y] = M^2 \left(1 - \frac{m}{M}\right) \frac{S_t^2}{m} + \frac{M}{m} \sum_{i=1}^M N_i^2 \left(1 - \frac{n_i}{N_i}\right) \frac{S_i^2}{n_i}$	$Var[\hat{y}] = \frac{1}{N^2} Var[\hat{t}_y]$
Estimateur de la variance d'échantillonnage	$\hat{V}ar[\hat{t}_y] = M^2 \left(1 - \frac{m}{M}\right) \frac{\hat{S}_t^2}{m} + \frac{M}{m} \sum_{i \in SUP} N_i^2 \left(1 - \frac{n_i}{N_i}\right) \frac{\hat{S}_i^2}{n_i}$	$\hat{V}ar[\hat{y}] = \frac{1}{N^2} \hat{V}ar[\hat{t}_y]$

Intervalle au niveau de confiance 95% pour la moyenne :

$$IC_{95\%}(\bar{Y}) = \left[ \hat{y} - 1,96\sqrt{\hat{V}ar(\hat{y})}, \hat{y} + 1,96\sqrt{\hat{V}ar(\hat{y})} \right]$$

sous hypothèse que la taille de l'échantillon est assez grande

## X. REDRESSEMENTS

### Exercice 1

### Post-stratification

Un institut de sondage est chargé de mesurer l'audience d'un nouveau magazine. Il interroge pour cela un échantillon de taille  $n$  selon un procédé que l'on assimilera à un plan simple à probabilités égales et sans remise au sein de la population française des individus âgés de 15 ans et plus. On supposera de plus qu'il n'y a pas de non-réponse. Pour satisfaire à la demande de l'éditeur, les résultats sont ventilés selon le critère « habitant en zone urbaine » ou « habitant en zone rurale ». Les données recueillies se présentent ainsi :

	Habitant en zone rurale	Habitant en zone urbaine	Total
Lecteurs	64	476	540
Non lecteurs	576	884	1 460
Total	640	1 360	2 000

1. Estimez la proportion du lectorat du magazine dans l'ensemble de la population et proposez un intervalle de confiance à 95% de ce taux de lecture.
2. Sachant que la proportion réelle d'habitants en zone urbaine vaut 75%, proposez un nouvel estimateur de la proportion de lecteurs et donnez en un intervalle de confiance à 95%. Quel gain de précision obtient-on ?

### Exercice 2

### Chiffre d'affaires et effectif salarié

(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

Dans une population de 10 000 entreprises, on veut estimer le chiffre d'affaires moyen  $\bar{Y}$ . Pour cela, on échantillonne  $n=100$  entreprises par sondage aléatoire simple. On dispose par ailleurs de l'information auxiliaire « nombre de salariés » notée  $x$  par entreprise. Les données issues du sondage sont :

- $\bar{X}=50$  salariés (vraie moyenne sur les  $x_k$ ),
- $\hat{y}=5.2 \times 10^6$  euros (chiffre d'affaires moyen dans l'échantillon),
- $\hat{x}=45$  salariés (effectif moyen dans l'échantillon),
- $s_y^2=25 \times 10^{10}$  (dispersion corrigée des  $y_k$  calculée dans l'échantillon),
- $s_x^2=15$  (dispersion corrigée des  $x_k$  calculée dans l'échantillon),
- $\hat{\rho}=0.8$  (coefficient de corrélation linéaire entre  $x$  et  $y$  calculé dans l'échantillon).

1. Que vaut l'estimateur par le ratio ? Cet estimateur est-il biaisé ?
2. Rappelez la formule de variance « vraie » de cet estimateur.
3. Calculez une estimation de la variance vraie. L'estimateur de variance utilisé est-il biaisé ?
4. Donnez un intervalle de confiance à 95% pour  $\bar{Y}$ .

### Exercice 3

### Estimation d'une surface cultivée

(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

On considère une région agricole comprenant  $N = 2010$  fermes où on cherche à estimer la moyenne de la surface cultivée en céréales (variable  $Y$  mesurée en hectares). On possède l'information auxiliaire sur la surface agricole totale cultivée de chaque ferme. En particulier, on sait qu'il y a 1 580 fermes de moins de 160 hectares (post-strate 1) et 430 fermes d'au moins 160 hectares (post-strate 2). On réalise un sondage aléatoire simple de  $n = 100$  exploitations et on obtient (avec les indices 1 et 2 pour les deux post-strates définies) :  $n_1=70$   $n_2=30$   $\hat{y}_1=19,40$   $\hat{y}_2=51,63$   $s_{y_1}^2=312$   $s_{y_2}^2=922$ .

1.
  - a. Quel est l'estimateur post-stratifié  $\hat{Y}_{post}$  ? Est-il différent de la moyenne simple?
  - b. Quelle est la loi de  $m$  ? Que valent son espérance et sa variance?
  - c. Calculer l'estimateur sans biais de la variance de  $\hat{Y}_{post}$  et donner un intervalle de confiance à 95% pour la surface moyenne cultivée en céréales.
2. On exploite désormais la variable auxiliaire  $X$  mesurant la surface agricole totale cultivée pour construire un estimateur par le ratio. On connaît la moyenne  $\bar{X}=118,32$  ha et on obtient sur l'échantillon :  $\hat{x}=132,25$   $s_x^2=9173$   $s_y^2=708$   $\hat{\rho}=0,57$  où  $\hat{\rho}$  est l'estimateur du vrai coefficient de corrélation linéaire inconnu  $\rho$ .
  - a. Rappeler l'expression de  $\rho$ .
  - b. Comment définissez-vous  $\hat{\rho}$  ? S'agit-il d'une estimation sans biais de  $\rho$  ?
  - c. Montrez que l'estimateur par le ratio de  $\bar{Y}$  apparaît préférable à la moyenne simple si et seulement si  $\hat{\rho} > \frac{1}{2} \frac{\hat{C}V(x)}{\hat{C}V(y)}$  où les  $\hat{C}V$  estiment les coefficients de variation.  
Qu'obtient-on dans le cas présent ?
  - d. Calculez l'estimateur par le ratio  $\hat{y}_q$  de  $\bar{Y}$ .
  - e. Estimez sa précision et donnez un intervalle de confiance à 95% pour  $\bar{Y}$ .

#### Exercice 4

Taille des pieds

(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

Le directeur d'une entreprise de confection de chaussures veut estimer la longueur moyenne des pieds droits des hommes adultes d'une ville. Soient  $y$  le caractère « longueur du pied droit » (en centimètre) et  $x$  la taille de l'individu (en centimètres). Le directeur sait en outre par les résultats d'un recensement que la taille moyenne des hommes adultes de cette ville est de 168 cm. Pour estimer la longueur des pieds, le directeur effectue un sondage aléatoire simple sans remise de 100 hommes adultes. Les résultats sont les suivants :  $s_y=2, s_x=10, s_{xy}=15, \hat{x}=169, \hat{y}=24$ . Sachant que 400 000 hommes adultes vivent dans cette ville,

1. Calculez l'estimateur d'Horvitz-Thompson, l'estimateur par le quotient, l'estimateur par différence et l'estimateur par la régression.
2. Estimez les variances de ces 4 estimateurs
3. Quel estimateur conseilleriez-vous au directeur ?
4. Exprimez la différence littérale entre la variance de l'estimateur par le quotient et la variance de l'estimateur par la régression en fonction de  $\hat{x}$ ,  $\hat{y}$  et de la pente  $\hat{b}$  de la droite de régression de  $y$  sur  $x$  dans l'échantillon. Commentez.

#### Exercice 5

Comparaison d'estimateurs

(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

On se propose d'estimer la moyenne  $\bar{Y}$  d'un caractère d'intérêt au moyen d'un échantillon sélectionné selon un plan aléatoire simple sans remise de taille 1 000 dans une population de taille 1 000 000. On connaît la moyenne  $\bar{X}=15$  d'un caractère auxiliaire  $x$  et on donne, avec les notations usuelles,

$$\hat{y}=10 ; \hat{x}=14 ; s_x^2 = 25 ; s_y^2 = 20 \text{ et } s_{xy} = 15$$

1. Estimez  $\bar{Y}$  au moyen des estimateurs d'Horvitz-Thompson, par différence, par le quotient et par la régression. Estimez les variances de ces estimateurs.
2. Quel estimateur choisiriez-vous pour estimer  $\bar{Y}$  ?

## Rappels sur les redressements

### I/ Intérêt

Accroître la précision en tirant parti d'information auxiliaire liée au caractère d'intérêt.

Selon la nature de l'information auxiliaire, techniques de post-stratification, d'estimation par différence, par le ratio, par la régression, par calage généralisé.

### II/ Formulaire pour le total

En notant X l'information auxiliaire,

Méthode	Estimateur du total	Vraie erreur quadratique moyenne de cet estimateur
Estimateur d'Horvitz-Thompson	$\hat{t}_{y\pi} = \sum_{k \in S} \frac{y_k}{\pi_k} = \frac{N}{n} \sum_{k \in S} y_k$	$Var(\hat{t}_{y\pi}) = N^2 \left(1 - \frac{n}{N}\right) \frac{S_y^2}{n}$
Estimateur par la différence	$\hat{t}_{yD} = \hat{t}_{y\pi} + t_x - \hat{t}_{x\pi}$	$Var(\hat{t}_{yD}) = N^2 \left(1 - \frac{n}{N}\right) \frac{(S_y^2 + S_x^2 - 2S_{xy})}{n}$
Estimateur par le ratio (ou le quotient)	$\hat{t}_{yQ} = \hat{t}_{y\pi} \frac{t_x}{\hat{t}_{x\pi}}$	$Var(\hat{t}_{yQ}) = N^2 \left(1 - \frac{n}{N}\right) \frac{(S_y^2 + R^2 S_x^2 - 2RS_{xy})}{n}$ avec $R = \frac{t_y}{t_x} = \frac{\bar{Y}}{\bar{X}}$
Estimateur par la régression	$\hat{t}_{yD} = \hat{t}_{y\pi} + \hat{b}(t_x - \hat{t}_{x\pi})$ avec $\hat{b} = \frac{\hat{S}_{xy}}{\hat{S}_x^2}$	$Var(\hat{t}_{yQ}) = N^2 \left(1 - \frac{n}{N}\right) \frac{S_y^2}{n} (1 - \rho^2)$ avec $\rho = \frac{S_{xy}}{S_x S_y}$
Estimateur post-stratifié	$\hat{t}_{y_{post}} = \sum_{h=1}^H N_h \hat{y}_h = \sum_{h=1}^H \frac{N_h}{n_h} \sum_{k \in S_h} y_k$	$Var(\hat{t}_{y_{post}}) = \frac{N-n}{n} \sum_{h=1}^H N_h S_{yh}^2 + \frac{N-n}{N-1} \frac{N^2}{n^2} \sum_{h=1}^H \frac{N-N_h}{N} S_{yh}^2$

Estimateur par substitution de l'erreur quadratique moyenne

Intervalle au niveau de confiance 95% pour la moyenne :

$$IC_{95\%}(\bar{Y}) = \left[ \hat{y} - 1,96\sqrt{\hat{Var}(\hat{y})}, \hat{y} + 1,96\sqrt{\hat{Var}(\hat{y})} \right]$$

sous hypothèse que la taille de l'échantillon est assez grande.

## IX. TP2 : CALAGE SUR MARGES

Cette séance utilise la macro SAS %CALMAR développée par l'Insee. Elle est disponible sur le site [insee.fr](http://insee.fr), accompagnée de sa documentation.

### Exercice 1

Un institut spécialisé a réalisé une enquête auprès des salariés d'une entreprise, qui compte 230 salariés répartis sur deux établissements A (70 salariés) et B (160 salariés). L'institut a effectué un sondage aléatoire simple dans chaque établissement, de taux de sondage respectifs 1/10 (A) et 1/20 (B). Le but est d'estimer la proportion de salariés prêts à monétariser une partie de leurs congés.

Pour chaque salarié enquêté, on dispose de :

- son identifiant (variable ID), à 3 caractères : le premier indique l'établissement, les deux suivants constituent un numéro d'ordre dans l'établissement ;
- la variable SERVICE indiquant si le salarié travaille dans un service productif (1) ou administratif (2) ;
- la variable CATEG qui indique la catégorie de personnel à laquelle appartient le salarié : employés (1), ouvriers (2), autres (9) ;
- la variable SEXE ;
- la variable SALAIRE annuel brut ;
- la variable Y indiquant si l'employé est intéressé par le paiement de jours de son compte-épargne temps (oui = 1, non = 0).

Par ailleurs, la direction de l'entreprise a aimablement fourni les informations suivantes sur ses salariés : l'entreprise compte 80 employés, 90 ouvriers, 140 hommes, 100 personnes travaillent dans le secteur productif, et le salaire total annuel vaut 47 000.

On vous demande d'utiliser cette information auxiliaire, en réalisant différents calages :

- par la méthode linéaire ;
- par la méthode raking ratio ;
- par la méthode logit LO=0.5 UP=2.2.

Comment estimez-vous le paramètre d'intérêt ?

Les données sont disponibles ci-joint au format SAS :



tp2\_exercice1.sas7b  
dat

## Exercice 2

Vous disposez d'une base de sondage de 11 600 individus décrits par la région, l'âge, le niveau scolaire, la catégorie socio-professionnelle, etc. (cf. tableau ci-dessous).

Le but de l'exercice est d'en sélectionner un échantillon, puis de procéder à des estimations et des redressements, en faisant comme si l'information d'intérêt avait été collectée sur l'échantillon seulement. Les variables d'intérêt mesurent l'importance consacrée aux activités sportives et culturelles.

Ci-dessous le contenu de la base de données :

Nom	Type	Libellé et modalités
IDENTIND	C	Identifiant
TRAGE	C	Tranche d'âge 1 : de 15 à 25 ans 2 : de 25 à 29 ans 3 : de 30 à 39 ans 4 : de 40 à 49 ans 5 : de 50 à 64 ans 6 : de 65 à 69 ans 7 : plus de 70 ans
NIVSCO2	C	Niveau scolaire 1 : inférieur au baccalauréat 2 : supérieur au baccalauréat
CS	C	Catégorie socio-professionnelle 1 : agriculteurs 2 : artisans, commerçants, chefs d'entreprises, professions libérales 3 : cadres 4 : professions intermédiaires 5 : employés 6 : ouvriers 7 : retraités
REGION	C	Région
ZEAT	C	Zone d'emploi et d'aménagement du territoire 1 : Région parisienne 2 : Bassin parisien 3 : Nord 4 : Est 5 : Ouest 6 : Sud-ouest 7 : Centre-est 8 : Méditerranée
CINEMA	N	Nombre de fois où l'individu est allé au cinéma au cours des 12 derniers mois
EXPO	N	Nombre d'expositions visitées au cours des 12 derniers mois
SPORT	N	Nombre d'heures de sport pratiquées au cours de la dernière semaine
LECTURE	N	Nombre d'heures de lecture au cours de la dernière semaine
TELE	N	Nombre d'heures passées devant la télévision au cours de la dernière semaine

Les données vous sont fournies au format SAS :



tp2\_exercice2.sas7b  
dat

1 / Donner la répartition de la population par tranche d'âge et niveau scolaire.

2 / Sélectionner un échantillon de taille 1 160 selon un sondage aléatoire simple.

*Pour rappel, la syntaxe de la procédure SURVEYSELECT de SAS est la suivante :*

```
PROC SURVEYSELECT DATA = nom de la base de sondage lue en entrée
  STATS
  METHOD      = SRS pour un sondage aléatoire simple sans remise
  SEED       = germe
  SAMPsize   = taille de l'échantillon souhaitée
  OUT        = nom de la table de sortie (l'échantillon);
RUN;
```

3 / A partir de l'échantillon, estimer la répartition de la population par tranche d'âge et niveau scolaire. Évaluer également le nombre moyen d'heure par semaine consacrées à la lecture, au sport, passées devant la télévision ainsi que le nombre moyen d'expositions visitées en une année et le nombre moyen de séances de cinéma en un an.

*Pour rappel, la syntaxe de la procédure SURVEYMEANS de SAS est la suivante :*

```
PROC SURVEYMEANS DATA = nom de la table-échantillon
  N = Effectif de la population
  MEAN STDERR CLM CV = Statistiques éditées en sortie;
VAR listes de variable d'intérêt;
WEIGHT variable de pondération;
RUN;
```

4 / Caler l'échantillon sur la vraie structure par tranche d'âge et niveau scolaire. Pourquoi ces variables de calage sont-elles pertinentes ?

5/ Ré-estimer les grandeurs citées à la question 3.

## X. TP3 : CORRECTION DE LA NON-REPONSE

Le but de l'étude de cas est de corriger la non-réponse (totale et partielle) pour une enquête conduite auprès de 2 389 personnes interrogées sur leur perception de leur état de santé.

L'échantillon a été choisi par sondage aléatoire simple sans remise dans une population de 2 millions d'individus. Les informations disponibles sont les suivantes :

- l'identifiant de l'enquêté (variable « ident »),
- son poids de sondage initial (« poids »),
- son âge (« âge »),
- son sexe (« sexe »),
- son niveau de revenu (« revenu »),
- sa région d'habitat (« region »),
- son nombre de consultations chez un professionnel de la santé en un an (« visites »),
- sa consommation de tabac (« tabac »),
- sa perception de son état de santé (« sante »),
- une indicatrice de la non-réponse totale (« nrt »),
- une indicatrice de la non-réponse partielle (« nrp »).

Les modalités des caractéristiques qualitatives sont définies de la sorte :

Variable	Modalités
Age	3-4 : Junior 5-6 : Jeune adulte 7-8 : Adulte 9-11: Senior
Sexe	1 : Homme 2 : Femme
Revenu	1-4 : Bas revenus 5-6 : Moyens revenus 7-8 : Bons revenus 9-11: Hauts revenus
Tabac	1 : Fume quotidiennement 2 : Ayant fumé quotidiennement 3 : Fume occasionnellement 4 : Ayant fumé occasionnellement 5 : Jamais fumé
Sante	1-2 : Excellente 3 : Bonne 4-5 : Passable

La base de données est fournie au format SAS :



tp3.sas7bdat

1. Dresser un état des lieux sur l'importance des non-réponses.
2. Corriger la non-réponse totale. On commencera par décrire le comportement de non-réponse totale en fonction des caractéristiques disponibles pour tous les individus.
3. Corriger la non-réponse partielle pour la variable d'intérêt en envisageant diverses méthodes (imputation par la moyenne, imputation par la moyenne par classe, imputation par déduction, imputation par hot-deck, imputation par hot-deck par classe, etc.).

## XI. COMPLÉMENTS ET RÉVISIONS

### Exercice 1

### Algorithme de tirage bernoullien

On considère une population  $U$  de 1000 individus composée de trois sous-populations disjointes  $U_1, U_2, U_3$  de tailles respectives  $N_1 = 600, N_2 = 300, N_3 = 100$ . On va échantillonner dans cette population au moyen de tirage bernoullien : cette méthode consiste à choisir une probabilité d'inclusion commune  $\pi$ , puis à simuler sur la population une variable aléatoire distribuée selon une loi uniforme sur  $[0, 1[$  et à sélectionner les individus pour lesquels la réalisation de cette variable est inférieure à  $\pi$ .

1. On décide dans un premier temps de tirer un échantillon dans  $U$  en utilisant le plan de sondage suivant :
  - dans la sous-population  $U_1$ , on réalise un tirage bernoullien, tel que chaque élément  $k$  a la probabilité  $\pi_k = 0.1$  d'être sélectionné,
  - dans la sous-population  $U_2$ , on réalise un tirage bernoullien, tel que chaque élément  $k$  a la probabilité  $\pi_k = 0.2$  d'être sélectionné,
  - dans la sous-population  $U_3$ , on réalise un tirage bernoullien, tel que chaque élément  $k$  a la probabilité  $\pi_k = 0.8$  d'être sélectionné,
  - l'échantillon complet est constitué de la réunion des trois sous-échantillons ainsi obtenus.

Calculer l'espérance et la variance de la taille  $n_s$  de l'échantillon.

2. On réalise maintenant un tirage bernoullien directement dans  $U$ , tel que chaque élément a la probabilité  $\pi$  d'être sélectionné.
  - a. Déterminer  $\pi$  pour que l'espérance de la taille de l'échantillon, sous ce plan de sondage, soit égale à l'espérance de la taille de l'échantillon calculée à la question précédente.
  - b. Calculer alors la variance de la taille de l'échantillon, et comparer cette variance à celle de la question précédente.

### Exercice 2

### Tendance linéaire et tirage systématique

(d'après J-M. Grosbras, *Méthodes statistiques des sondages, Economica, 1987*)

On considère une population de taille  $N$  avec  $N = nk$  où  $n$  est la taille souhaitée de l'échantillon et  $k$  un nombre entier. On suppose que pour tout individu  $k$  de la population, on a  $Y_k = k$  pour  $k = 1$  à  $N$ .

1. On note respectivement  $\bar{Y}$  et  $S_Y^2$  la moyenne et la dispersion du caractère d'intérêt sur la population. Vérifier que  $\bar{Y} = \frac{N+1}{2}$  et  $S_Y^2 = \frac{N(N+1)}{12}$ .
2. On réalise un sondage aléatoire simple sans remise de taille  $n$ .
  - a. Quel est l'estimateur classique  $\hat{Y}$  de la moyenne ?
  - b. Montrer que sa variance vaut :  $V(\hat{Y}) = \frac{(k-1)(N+1)}{12}$ .

3. On réalise à présent un tirage systématique de taille  $n$  : on tire un nombre  $a$  au hasard entre 1 et  $k$  et on forme un échantillon de taille voulue avec les unités  $a, a + k, a + 2k, \dots, a + (n-1)k$ .  
Soit  $\hat{Y}_{sys}$  la moyenne des unités sélectionnées dans l'échantillon.

Montrer que :  $E(\hat{Y}_{sys}) = \bar{Y}$  et que  $V(\hat{Y}_{sys}) = \frac{k^2 - 1}{12}$

4. Comparer  $V(\hat{Y})$  et  $V(\hat{Y}_{sys})$  et commenter

### Exercice 3

### Algorithme du tri aléatoire

On veut estimer le poids moyen de 10 éléphants d'un cirque. Pour cela, on réalise un sondage aléatoire simple sans remise de taille 5 à l'aide d'un tri aléatoire. On simule donc une variable aléatoire uniforme  $U \sim U[0,1]$  sur la population des éléphants, puis on trie les réalisations obtenues par ordre croissant (ou décroissant) et on retient l'échantillon correspondant aux 5 plus grandes valeurs (ou plus petites). La simulation a été effectuée à partir de la fonction *ALEA()* sous Excel et a donné les réalisations ci-dessous :

N° de l'éléphant	Valeur générée
1	0,84
2	0,12
3	0,36
4	0,60
5	0,68
6	0,11
7	0,87
8	0,44
9	0,21
10	0,77

1. Quel est l'échantillon tiré ?
2. On pèse les éléphants retenus et on obtient en tonnes les poids respectifs suivants : 3,65 ; 3,17 ; 4,18 ; 3,55 et 4,26.
3. Donnez un estimateur du poids moyen des éléphants puis un intervalle de confiance à 95% de ce poids moyen.
4. Finalement, on réalise une pesée exhaustive des éléphants. On obtient un poids moyen de 3,45 tonnes. Que dire de l'intervalle de confiance précédent ? D'où peut venir le problème ?

### Exercice 4

### Algorithme de « sélection-rejet »

La méthode de sélection-rejet permet d'obtenir un échantillon de taille  $n$  en une seule lecture du fichier. L'algorithme est le suivant :

- On initialise à 0 les compteurs  $k$  et  $j$  renseignant respectivement le nombre d'unités du fichier déjà examinées et le nombre d'unités déjà sélectionnées dans l'échantillon. On se positionne sur le premier individu du fichier.
- Tant que  $j$  est strictement inférieur à la taille d'échantillon voulue, on a généré un nombre aléatoire  $u$  selon une loi uniforme sur  $[0,1[$  pour l'individu de rang  $k+1$  sur lequel on est positionné et on décide :
  - Si on obtient  $u < \frac{n-j}{N-k}$ , alors on sélectionne l'unité de rang  $k+1$ . On incrémente donc  $j$  d'une unité, puis on passe à l'individu suivant en incrémentant  $k$ .
  - Sinon, l'unité  $k+1$  n'est pas tirée et on passe à l'individu suivant en incrémentant  $k$ .

1. Montrer que le plan est de taille fixe  $n$  et qu'il suffit effectivement donc d'au plus  $N$  opérations pour sélectionner ces  $n$  unités
2. Montrer que le plan est simple. En déduire que les probabilités d'inclusion individuelles sont bien égales à :  $\pi_k = \frac{n}{N}, \forall k \in U$ .
3. Application : sélectionner un échantillon de taille 4 dans une population de taille 10 selon cette méthode en utilisant les réalisations suivantes d'une variable aléatoire  $U$  uniforme sur  $[0,1[$  :

Individu k	1	2	3	4	5	6	7	8	9	10
$u_k$	0,375	0,620	0,518	0,0454	0,633	0,246	0,927	0,326	0,646	0,178

**Exercice 5** **Non-réponse dans les enquêtes par quotas**  
*(A-M. Dussaix, J-M. Grosbras, 1992, Exercices de sondage, Economica)*

L'objet de cet exercice est de montrer l'existence de biais pouvant découler de non-réponses dans les enquêtes par quotas. On considère une enquête où sont imposés des quotas relatifs à une variable qualitative donnée. Pour fixer les idées, on supposera, par exemple, qu'il y a dans la population,  $H$  variables d'âge ou de poids en proportion  $N_h/N$  pour  $h = 1$  à  $H$ . On demande aux enquêteurs de compléter un échantillon représentatif, c'est-à-dire tel que  $n_h/n = N_h/N$ . A la fin de l'enquête, la moyenne  $\bar{Y}$  de la variable d'intérêt est estimée par la moyenne simple sur l'échantillon  $\hat{y}$ , ce qui peut encore s'écrire :

$$\hat{y} = \sum_{h=1}^H \frac{n_h}{n} \hat{y}_h = \sum_{h=1}^H \frac{N_h}{N} \hat{y}_h \quad \text{où} \quad \hat{y}_h = \frac{1}{n_h} \sum_{k \in S_h} y_k.$$

Pour étudier l'influence de la non-réponse, on fait l'hypothèse qu'il existe dans la population un partage en 2 catégories :

- La 1<sup>ère</sup> est celle des personnes accessibles et répondant volontiers à l'enquête caractérisée par les effectifs  $N_1$  et  $N_{h1}$  dans les tranches d'âge  $h$ , et les moyennes  $\bar{Y}_1$  et  $\bar{Y}_{h1}$ .
- La 2<sup>ème</sup> est celle des personnes inaccessibles ou refusant de répondre à l'enquête caractérisée par les effectifs  $N_0$  et  $N_{h0}$  dans les tranches d'âge  $h$ , et les moyennes  $\bar{Y}_0$  et  $\bar{Y}_{h0}$ .

Naturellement, les quantités  $N_1, N_{h1}, N_0, N_{h0}, \bar{Y}_1, \bar{Y}_{h1}, \bar{Y}_0$  et  $\bar{Y}_{h0}$  sont inconnues.

1. Si on fait l'hypothèse que les  $n_h$  réponses constituent un échantillon d'un plan aléatoire simple sans remise prélevé dans un ensemble d'effectif  $N_{h1}$ , montrer que  $\hat{y}$  est un estimateur biaisé pour  $\bar{Y}$ . On écrira l'expression du biais en fonction de  $N, N_{h0}$  et  $\Delta_h = \bar{Y}_{h1} - \bar{Y}_{h0}$ .
2. Commentez brièvement cette expression. Construire un exemple numérique illustrant une situation où le biais est élevé (on prendra  $H = 3$ ).

**Exercice 6** **Nombre de titulaires de comptes CODEVI à interroger**  
*(d'après A-M. Dussaix et J-M. Grosbras, Exercices de sondage, Economica, 1992)*

Une banque désire étudier par sondage (interviews par enquêteur) les caractéristiques socio-démographiques (âge, catégorie sociale,...) et les comportements financiers des titulaires de comptes CODEVI. Leur répartition en fonction des montants moyens annuels des comptes est la suivante :

Solde moyen annuel	Nombre de comptes
De 0 à 100 €	15 000
De 100 à 900 €	15 000
Plus de 900 €	30 000
Ensemble	60 000

Pour chacun des trois groupes, on veut étudier la répartition des titulaires par classe d'âge, catégorie sociale, etc. Par exemple, on s'intéresse à la proportion de titulaires ayant entre 25 et 35 ans. Quelle taille d'échantillon doit-on prévoir dans chaque strate s'il s'agit de déterminer les différentes proportions avec une précision de  $\pm 2,5\%$  au niveau de confiance 95% ?

**Exercice 7** Tirage des UP avec remise – Taille de ménages  
(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

Pour estimer le nombre moyen  $\bar{Y}$  de personnes par ménage dans un pays donné, on réalise un tirage à 2 degrés :

- 1<sup>er</sup> degré : tirage aléatoire avec remise de  $m = 4$  villages parmi  $M = 400$  proportionnellement à leur taille. La taille d'un village est le nombre de ménages qu'il contient. Ainsi, à chacun des 4 tirages indépendants, un village est sélectionné avec une probabilité proportionnelle à sa taille.
- 2<sup>ème</sup> degré : tirage aléatoire simple de  $n_i$  ménages parmi les  $N_i$  si le village  $i$  est tiré.

Le nombre total de ménages dans le pays est  $N = 10\ 000$ . Les données sont représentées dans le tableau ci-dessous ;  $\hat{Y}_i$  est le nombre moyen de personnes par ménage dans le village  $i$  d'après l'échantillon.

i	1	2	3	4
$N_i$	20	23	25	18
$\hat{Y}_i$	5.25	5.50	4.50	5

1. a. Quelle est la probabilité de tirage  $p_i$  de chacun des 4 villages sélectionnés ? (on appelle probabilité de tirage la probabilité qu'a le village d'être choisi lors de chacun des 4 tirages indépendants réalisés successivement dans les mêmes conditions).
- b. Calculer  $\Pr(i \notin S)$  en fonction de  $(1 - p_i)$ . En déduire la probabilité d'inclusion  $\pi_i = \Pr(i \in S)$  en fonction de  $p_i$ . Examiner le cas où  $p_i$  est petit.
2. Quelle est l'expression de  $\bar{Y}$  (vraie valeur) et quel est son estimateur sans biais ?
3. Estimer la variance de cet estimateur. Quel intérêt a-t-on à utiliser un tirage avec remise au 1<sup>er</sup> degré ?

**Exercice 8** Raking-ratio  
(Ardilly P, Tillé Y, 2003, Exercices corrigés de méthodes de sondage, Ellipses)

On s'intéresse à la population des 10 000 étudiants inscrits en 1<sup>ère</sup> année dans une université. On connaît le nombre total d'étudiants dont les parents ont un diplôme d'études primaires, secondaires et supérieures (respectivement 5000, 3000 et 2000). On effectue un sondage selon un plan aléatoire simple sans remise de 150 étudiants. On ventile ces 150 étudiants selon le diplôme des parents et leurs résultats (échec ou réussite) à l'examen de 1<sup>ère</sup> année et on obtient le résultat ci-dessous :

Niveau d'études	Échec	Réussite
Primaire	45	15
Secondaire	25	25
Supérieure	10	30

1. Estimez le taux de réussite des étudiants en utilisant l'estimateur de Horvitz-Thompson et donnez un estimateur de variance et un intervalle de confiance à 95% de ce taux.
2. Expliquez pourquoi il est a priori intéressant d'effectuer un redressement, et pourquoi le redressement doit diminuer la valeur de l'estimation issue de 1.
3. Estimez le taux de réussite des étudiants par l'estimateur post-stratifié et donnez un estimateur de variance et un intervalle de confiance à 95% de ce taux.
4. Estimez le taux de réussite par niveau d'études des parents en utilisant une technique de raking-ratio et sachant que dans la population totale étudiante, le taux de réussite est en réalité de 40%.

## Exercice 9

## Cas pratique dans une pisciculture

Un éleveur de poissons souhaite connaître le poids moyen de ses poissons. Il dispose de 3 bassins selon l'âge des animaux : n°1 pour ceux « de petite taille », n°2 « de taille moyenne » et n°3 « de grande taille ». Le nombre total de poissons par bassin est respectivement de 1000, 900 et 950.

Notre pisciculteur appelle un statisticien à sa rescousse pour estimer le poids moyen des poissons. Armé de son épuisette, le statisticien attrape 20 poissons dans le bassin n°1, 15 dans le n°2 et 10 dans le n°3. Ensuite, il calcule le poids moyen sur les 3 échantillons relatifs aux 3 bassins. Il trouve : 0.152 Kilo pour le bassin N°1, 0.255 Kilo pour le n°2 et 0.305 Kilo pour le n°3. Il calcule également la dispersion corrigée des poids des poissons sur les 3 échantillons et trouve respectivement:  $(0.05)^2$  Kilo<sup>2</sup>,  $(0.02)^2$  Kilo<sup>2</sup> et  $(0.01)^2$  Kilo<sup>2</sup> pour les bassins N°1, 2 et 3.

On admettra que le mode de tirage des échantillons de poissons dans chacun des trois bassins est assimilable à un sondage aléatoire simple de taille fixe.

- 1)
  - a) Proposer un estimateur sans biais du poids moyen des poissons relativement à un bassin.
  - b) Donner les 3 estimations des poids moyens relatifs aux 3 bassins puis les 3 intervalles de confiance à 95% correspondants.
  - c) Pour estimer le poids moyen relatif à l'ensemble des 3 bassins, le statisticien a mis en œuvre l'estimateur stratifié. Après avoir rappelé la forme générale de cet estimateur et précisé les strates adoptées par le statisticien, donner l'estimation recherchée et l'intervalle de confiance à 95% correspondant.
- 2)
  - a) Est-ce que l'allocation définie par le statisticien correspond à l'allocation proportionnelle?
  - b) Compte tenu des mesures effectuées sur les échantillons, expliquer (qualitativement) pourquoi l'allocation du statisticien semble être légitime.
  - c) A partir des résultats obtenus sur les trois échantillons, calculer l'allocation de Neyman pour une taille totale de l'échantillon de poissons de 45.
- 3) Le pisciculteur propose d'estimer le poids moyen des poissons sur l'ensemble des 3 bassins en faisant la moyenne arithmétique des poids des poissons sur l'ensemble des 3 échantillons.
  - a) Calculer l'estimation fournie par le pisciculteur.
  - b) Montrer que cet estimateur est en réalité biaisé (on exprimera ce biais théorique en fonction des vrais poids moyens des poissons relatifs aux bassins, des vrais effectifs de poissons et des tailles des échantillons de poissons relatifs aux bassins).
  - c) Donner une estimation de ce biais.
- 4) Le statisticien apprend par hasard, en discutant avec l'un des employés, qu'un contrôle de la taille des poissons a été réalisé récemment. Ce contrôle a été effectué dans chacun des bassins et de façon quasi-exhaustive. Il révèle que la taille moyenne des poissons par bassin est de : 25 cm pour le bassin n°1, 40 cm pour n°2 et 50 cm pour le n°3.
  - a) Expliquer pourquoi la connaissance de cette nouvelle information est intéressante par rapport au phénomène étudié.
  - b) A partir de cette nouvelle information, proposer un nouvel estimateur du poids moyen des poissons pour un bassin fixé . Donner les 3 nouvelles estimations du poids moyen relatives à chacun des bassins. On donne pour cela les tailles moyennes des poissons mesurées sur les échantillons : 23 cm (bassin n°1), 42 cm (n°2), 51 cm (n°3).
  - c) Proposer une nouvelle estimation du poids moyen pour l'ensemble des 3 bassins.