

Introduction à la théorie des sondages - Cours 2

Thomas Merly-Alpa
thomas.merly-alpa@insee.fr

INSEE, département des méthodes statistiques

23 janvier 2017



Organisation

- 8 cours, 4 TD en demi-groupes
- 1/3 de la note : devoir maison à rendre **le 27 février**
- 2/3 de la note : examen final le 20 mars
- 2 intervenants :
 - Thomas Merly-Alpa - thomas.merly-alpa@insee.fr
 - Martin Chevalier - martin.chevalier@insee.fr
- Les slides et TD du cours sont à l'adresse
<http://nc233.com/teaching>

Principe du sondage

Objectif Construire un estimateur $\Phi(Y)$ d'une variable Y à partir d'un échantillon s de taille n tiré dans une population \mathcal{U} de taille N .

Plan de sondage On définit un plan de sondage p comme une loi de probabilité sur l'ensemble des échantillons possibles \mathcal{S} .

Exemple : $\mathcal{U} = \{1, 2, 3\}$. On définit le plan de sondage p_1 par :

$$\begin{aligned} p_1(\{1\}) &= p_1(\{2\}) = p_1(\{3\}) = 0 \\ p_1(\{1, 2\}) &= 0,5 \quad p_1(\{1, 3\}) = 0,2 \quad p_1(\{2, 3\}) = 0,2 \\ p_1(\{1, 2, 3\}) &= 0,1 \end{aligned}$$

Remarque p_1 n'est pas un plan de sondage de taille fixe.

Probabilités d'inclusion

Le plan de sondage permet de déterminer des probabilités d'inclusion pour chaque unité de la population.

Probabilité d'inclusion simple $\pi_k = \sum_{s \in \mathcal{S}} \delta_k p(s)$

Probabilité d'inclusion double $\pi_{k,l} = \sum_{s \in \mathcal{S}} \delta_k \delta_l p(s)$

avec $\delta_k(s) = \mathbf{1}(k \in s)$

Enfin, on note $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$.

Sommaire

- 1 Stratégie d'estimation
 - L'estimateur plugin
 - L'estimateur d'Horvitz-Thompson
 - L'estimateur de Hájek
 - Recherche d'un estimateur optimal
- 2 Sondage aléatoire simple
 - Définitions
 - Réaliser un tirage
- 3 Estimation dans un SAS
 - Estimation d'un total
 - Estimation d'une proportion
 - Échantillonnage dans le temps
 - Estimation sur un domaine
 - Estimation d'un ratio

Chapitre 1

Stratégie d'estimation

Partie 1

L'estimateur plugin

L'estimation "plugin"

Pour l'estimation du total et de la moyenne d'une variable Y ,
l'estimateur plugin s'écrit :

$$\hat{T}(Y)_{plugin} = \sum_{k \in s} y_k$$
$$\hat{y}_{plugin} = \frac{1}{n} \sum_{k \in s} y_k$$

L'estimation "plugin"

En général, l'estimation plugin est biaisée :

$$\mathbb{E}(\hat{\Phi}_{plugin}) = \sum_s p(s) \cdot \hat{\Phi}(s) \\ \neq \Phi$$

$\mathbb{E}(\hat{\Phi})$ est la valeur moyenne de $\hat{\Phi}$ obtenue avec le plan de sondage considéré **sur tous les échantillons possibles**.

Exemple d'estimation plug-in

Un exemple : étude du temps quotidien passé sur Internet :

Père			15 minutes
Mère			30 minutes
Enfant 1			215 minutes
Enfant 2			240 minutes

Vraie moyenne : 125 minutes.

Exemple d'estimation plug-in

On interroge les deux parents, et un des enfants au hasard.

Père	Dans l'échantillon	Poids = 1	15 minutes
Mère	Dans l'échantillon	Poids = 1	30 minutes
Enfant 1	Dans l'échantillon	Poids = 2	215 minutes
Enfant 2	/	/	? minutes

Estimateur plug-in : $15 + 30 + 215 / 3 \approx 87$ minutes

Exemple d'estimation plug-in

On interroge les deux parents, et un des enfants au hasard.

Père	Dans l'échantillon	Poids = 1	15 minutes
Mère	Dans l'échantillon	Poids = 1	30 minutes
Enfant 1	Dans l'échantillon	Poids = 2	215 minutes
Enfant 2	/	/	? minutes

Estimateur plug-in : $15 + 30 + 215 / 3 \approx 87$ minutes

Estimateur pondéré : $15 + 30 + 2*215 / 4 = 118,75$ minutes

Exemple d'estimation plug-in

On interroge les deux parents, et un des enfants au hasard.

Père	Dans l'échantillon	Poids = 1	15 minutes
Mère	Dans l'échantillon	Poids = 1	30 minutes
Enfant 1	/	/	? minutes
Enfant 2	Dans l'échantillon	Poids = 2	240 minutes

Estimateur plug-in : $15 + 30 + 240 / 3 = 95$ minutes

Exemple d'estimation plug-in

On interroge les deux parents, et un des enfants au hasard.

Père	Dans l'échantillon	Poids = 1	15 minutes
Mère	Dans l'échantillon	Poids = 1	30 minutes
Enfant 1	/	/	? minutes
Enfant 2	Dans l'échantillon	Poids = 2	240 minutes

Estimateur plug-in : $15 + 30 + 240 / 3 = 95$ minutes

Estimateur pondéré : $15 + 30 + 2*240 / 4 = 131,25$ minutes

Partie 2

L'estimateur d'Horvitz-Thompson

Définition

Définition

L'estimateur d'Horvitz-Thompson (ou π -estimateur) est défini :

$$\text{pour un total : } \hat{T}_{y\pi} = \sum_{k \in s} \frac{y_k}{\pi_k}$$

$$\text{pour une moyenne : } \hat{y}_{\pi} = \frac{1}{N} \sum_{k \in s} \frac{y_k}{\pi_k}$$

*C'est donc un **estimateur pondéré** utilisant les poids $w_k = \frac{1}{\pi_k}$*

Estimation sans biais

Théorème

Si $\forall k \in \mathcal{U}, \pi_k > 0$, alors l'estimateur d'Horvitz-Thompson est **sans biais** pour le total et la moyenne.

La condition signifie que toutes les unités de la population ont une chance non nulle d'être dans l'échantillon.

Estimation sans biais

Démonstration.

$$\begin{aligned}\mathbb{E}[\hat{T}_{y\pi}] &= \mathbb{E}\left[\sum_{k \in s} \frac{y_k}{\pi_k}\right] \\ &= \mathbb{E}\left[\sum_{k \in \mathcal{U}} \frac{y_k \delta_k}{\pi_k}\right] \\ &= \sum_{k \in \mathcal{U}} \frac{y_k \mathbb{E}[\delta_k]}{\pi_k} \\ &= \sum_{k \in \mathcal{U}} y_k \\ &= T(y)\end{aligned}$$

Rappel : Variance / Précision

$$\text{Var}(\hat{\Phi}) = \sum_s p(s) \cdot \left[\mathbb{E}(\hat{\Phi}) - \hat{\Phi}(s) \right]^2$$

C'est une mesure de la dispersion des valeurs $\hat{\Phi}(s)$ autour de leur moyenne.

Variance de l'estimateur de Horvitz-Thompson

Propriété

La variance de l'estimateur de Horvitz-Thompson s'écrit :

$$\text{Var}[\hat{T}_{y\pi}] = \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}} \frac{y_k y_l}{\pi_k \pi_l} \Delta_{kl}$$

(où : $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$)

Variance de l'estimateur de Horvitz-Thompson

Démonstration.

$$\begin{aligned}\text{Var}(\hat{t}_{y\pi}) &= \text{Var}\left(\sum_{k \in U} \frac{y_k}{\pi_k} \delta_k\right) \\ &= \sum_{k \in U} \frac{y_k^2}{\pi_k^2} \text{Var}(\delta_k) + \sum_{k \in U} \sum_{l \in U, l \neq k} \frac{y_k}{\pi_k} \frac{y_l}{\pi_l} \text{Cov}(\delta_k, \delta_l) \\ &= \sum_{k \in U} \frac{y_k^2}{\pi_k^2} \pi_k \cdot (1 - \pi_k) + \sum_{k \in U} \sum_{l \in U, l \neq k} \frac{y_k}{\pi_k} \frac{y_l}{\pi_l} (\pi_{kl} - \pi_k \pi_l) \\ &= \sum_{k, l \in U} \frac{y_k}{\pi_k} \frac{y_l}{\pi_l} \Delta_{kl}\end{aligned}$$



Variance pour un plan de taille fixe

Propriété

Si le plan de sondage est de taille fixe (formule de Yates-Grundy) :

$$\text{Var}(\hat{t}_{y\pi}) = -\frac{1}{2} \sum_{k \in U} \sum_{l \in U, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \Delta_{kl}$$

Variance de l'estimateur de Horvitz-Thompson

Démonstration.

Découle de la formule de Horvitz-Thompson quand le plan de sondage est de taille fixe. Pour démontrer la formule, il vaut mieux procéder à rebours :

$$\begin{aligned}
 & -\frac{1}{2} \sum_{k \in U} \sum_{l \in U, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \Delta_{kl} \\
 &= \frac{1}{2} \sum_{k \in U} \sum_{l \in U, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 (\pi_k \pi_l - \pi_{kl}) \\
 &= \frac{1}{2} \sum_{k \in U} \sum_{l \in U, l \neq k} \left(\frac{y_k^2}{\pi_k^2} + \frac{y_l^2}{\pi_l^2} - 2 \frac{y_k y_l}{\pi_k \pi_l} \right) (\pi_k \pi_l - \pi_{kl}) \\
 &= \sum_{k \in U} \sum_{l \in U, l \neq k} \frac{y_k^2}{\pi_k^2} (\pi_k \pi_l - \pi_{kl}) - \sum_{k \in U} \sum_{l \in U, l \neq k} y_k y_l \left(1 - \frac{\pi_{kl}}{\pi_k \pi_l} \right) \\
 &= \sum_{k \in U} \frac{y_k^2}{\pi_k} \left(\sum_{l \in U, l \neq k} \pi_l - \frac{1}{\pi_k} \frac{y_k^2}{\pi_k} \pi_{kl} \right) - \sum_{k \in U} \sum_{l \in U, l \neq k} y_k y_l \left(1 - \frac{\pi_{kl}}{\pi_k \pi_l} \right)
 \end{aligned}$$

Variance de l'estimateur de Horvitz-Thompson

Démonstration.

...

Or, d'après le cours 1, on a dans le cas taille fixe : $\sum_{k \in U} \pi_k = n$ et $\sum_{l \in U, l \neq k} \pi_{kl} = \pi_k(n-1)$, cela donne :

$$\begin{aligned}
 & -\frac{1}{2} \sum_{k \in U} \sum_{l \in U, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \Delta_{kl} \\
 &= \sum_{k \in U} \frac{y_k^2}{\pi_k} \left(n - \pi_k - \frac{\pi_k(n-1)}{\pi_k} \right) - \sum_{k \in U} \sum_{l \in U, l \neq k} y_k y_l \left(1 - \frac{\pi_{kl}}{\pi_k \pi_l} \right) \\
 &= \sum_{k \in U} \frac{y_k^2}{\pi_k^2} \pi_k (1 - \pi_k) - \sum_{k \in U} \sum_{l \in U, l \neq k} \frac{y_k y_l}{\pi_k \pi_l} (\pi_k \pi_l - \pi_{kl}) \\
 &= \sum_{k, l \in U} \frac{y_k}{\pi_k} \frac{y_l}{\pi_l} \Delta_{kl}
 \end{aligned}$$

Et on retombe bien sur la formule d'Horvitz-Thompson.

Estimation de variance

Les quantités précédentes sont les **vraies variances**. On peut utiliser les estimateurs suivants, qui sont sans biais dès lors que $\forall k, l, \pi_{kl} > 0$:

$$\hat{\text{Var}}(\hat{t}_{y\pi}) = \sum_{k \in s} \frac{y_k^2}{\pi_k^2} (1 - \pi_k) - \sum_{k \in s} \sum_{l \in s, l \neq k} \frac{y_k y_l}{\pi_k \pi_l \pi_{kl}} (\pi_k \pi_l - \pi_{kl})$$

Pour un plan de taille fixe :

$$\hat{\text{Var}}(\hat{t}_{y\pi}) = -\frac{1}{2} \sum_{k \in s} \sum_{l \in s, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \frac{\pi_k \pi_l - \pi_{kl}}{\pi_{kl}}$$

Remarque

Remarque : si le plan de sondage ne vérifie pas :

$$\forall k \neq l \in \mathcal{U}, \pi_{kl} - \pi_k \pi_l \geq 0$$

(condition de Sen-Yates-Grundy), ces estimateurs de variance peuvent prendre des valeurs négatives.

Construction d'un intervalle de confiance

On fait l'**hypothèse** : $\hat{\Phi}(s) \sim \mathcal{N}(\Phi, \text{Var}(\Phi))$

L'intervalle de confiance à 95% est défini par :

$$IC_{95\%} = \left[\hat{\Phi} - 2\sigma(\hat{\Phi}); \hat{\Phi} + 2\sigma(\hat{\Phi}) \right]$$

L'intervalle de confiance **estimé** est défini par :

$$\hat{IC}_{95\%} = \left[\hat{\Phi} - 2\hat{\sigma}(\hat{\Phi}); \hat{\Phi} + 2\hat{\sigma}(\hat{\Phi}) \right]$$

Construction d'un intervalle de confiance

Pour l'estimateur de variance de l'estimateur d'Horvitz-Thompson pour un plan à taille fixe, cela revient à calculer :

La borne inférieure de l'intervalle de confiance

$$\sum_{k \in s} \frac{y_k}{\pi_k} - 2 \sqrt{\frac{1}{2} \sum_{k \in s} \sum_{l \in s, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \left(1 - \frac{\pi_k \pi_l}{\pi_{kl}} \right)}$$

et la borne supérieure de l'intervalle de confiance

$$\sum_{k \in s} \frac{y_k}{\pi_k} + 2 \sqrt{\frac{1}{2} \sum_{k \in s} \sum_{l \in s, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \left(1 - \frac{\pi_k \pi_l}{\pi_{kl}} \right)}$$

Partie 3

L'estimateur de Hájek

Définition

L'estimateur de Horvitz-Thompson de la moyenne nécessite la connaissance de N , la taille de la population. Si on ne la connaît pas, on peut utiliser dans ce cas l'estimateur de Hájek :

$$\hat{y}_H = \frac{\sum_{k \in s} \frac{y_k}{\pi_k}}{\sum_{k \in s} \frac{1}{\pi_k}}$$

Propriété

L'estimateur de Hájek est biaisé, mais en général, le biais est négligeable.

Estimateur de Hájek du total

L'estimateur de Hájek peut être utilisé pour estimer un total :

$$T(\hat{Y})_H = N \cdot \frac{\sum_{k \in s} \frac{y_k}{\pi_k}}{\sum_{k \in s} \frac{1}{\pi_k}}$$

... mais cela impose de connaître N .

Partie 4

Recherche d'un estimateur optimal

Estimateur de Horvitz-Thompson

L'estimateur de Horvitz-Thompson constitue le fondement de l'estimation par sondage (même si d'autres estimateurs peuvent être utilisés, la logique de construction découle souvent de celle de Horvitz-Thompson, voir cours suivants)

Estimateur de Horvitz-Thompson

L'estimateur de Horvitz-Thompson n'est pas le seul estimateur sans biais.

Recherche d'optimalité

Existe-t-il un estimateur optimal en sondages ?

Question centrale pour les théoriciens des sondages dans les années 1950 à 1970 : Godambe, Hanurav, Basu, etc.

Difficile à répondre car cela dépend de la population, de la taille d'échantillon, des concepts mesurés. . .

Chapitre 2

Sondage aléatoire simple

Partie 1

Définitions

Définition

Sondage aléatoire simple sans remise (SAS) de taille n : plan de sondage sans remise de taille fixe n tel que tous les échantillons de taille n ont la même probabilité d'être tirés. Cette probabilité vaut :

$$p(s) = \frac{1}{\binom{N}{n}} \quad \text{si } |s| = n$$
$$= 0 \quad \text{sinon.}$$

On note le taux de sondage : $f = \frac{n}{N}$

Un petit rappel

Combien vaut $\binom{N}{n}$? On rappelle que cette notation, n parmi N , signifie "le nombre de façons de choisir n éléments parmi N ", noté aussi C_N^n . On a ainsi :

$$\binom{N}{n} = \frac{N!}{n!(N-n)!}$$

où $n! = 1 \times 2 \times 3 \times \dots \times n$.

Probabilités d'inclusion

$$\forall k \in \mathcal{U}, \pi_k = \mathbb{P}(k \in s) = \frac{n}{N} = f$$

$$\forall k \neq l \in \mathcal{U}, \pi_{k,l} = \mathbb{P}(k \wedge l \in s) = \frac{n(n-1)}{N(N-1)}$$

Notations

On note, **dans la population** :

$$\text{Total : } T(Y) = \sum_{k \in \mathcal{U}} Y_k$$

$$\text{Moyenne : } \bar{Y} = \frac{1}{N} \sum_{k \in \mathcal{U}} Y_k$$

$$\text{Variance empirique (dispersion) : } S^2 = \frac{1}{N-1} \sum_{k \in \mathcal{U}} (Y_k - \bar{Y})^2$$

Notations

On note, **dans l'échantillon s** :

$$\text{Total : } n\bar{y} = \sum_{k \in s} y_k$$

$$\text{Moyenne : } \bar{y} = \frac{1}{n} \sum_{k \in s} y_k$$

$$\text{Variance empirique (dispersion) : } s^2 = \frac{1}{n-1} \sum_{k \in s} (y_k - \bar{y})^2$$

Partie 2

Réaliser un tirage

Tirage aléatoire simple

Comment procéder en pratique pour tirer un échantillon ? Il y a plusieurs possibilités.

- En R, par exemple, on peut utiliser la fonction `sample` qui réalise un sondage aléatoire simple.
- Sinon, le moyen le plus simple consiste à trier la population complètement au hasard, et choisir les n premiers individus.

Tirage aléatoire simple

Comment trier aléatoirement une population ?

A		
B		
C		
D		
E		
F		
G		
H		
I		
J		
K		

Tirage aléatoire simple

On génère pour chaque individu une variable aléatoire uniforme, entre 0 et 1.

A	0.123	
B	0.245	
C	0.654	
D	0.987	
E	0.015	
F	0.975	
G	0.126	
H	0.745	
I	0.811	
J	0.626	
K	0.413	

Tirage aléatoire simple

On trie la population sur cette variable, et on prend les $n = 4$ premiers (par exemple)

E	0.015	Sélection
A	0.123	Sélection
G	0.126	Sélection
B	0.245	Sélection
K	0.413	
J	0.626	
C	0.654	
H	0.745	
I	0.811	
F	0.975	
D	0.987	

Chapitre 3

Estimation dans un SAS

Partie 1

Estimation d'un total

Estimateur d'Horvitz-Thompson

L'estimateur d'Horvitz-Thompson pour le total et la moyenne s'écrit :

$$T(\hat{Y}) = \sum_{k \in s} \frac{1}{\pi_k} y_k = \frac{N}{n} \sum_{k \in s} y_k = N\bar{y}$$
$$\hat{Y} = \frac{1}{N} \sum_{k \in s} \frac{1}{\pi_k} y_k = \bar{y}$$

Poids de sondage

Les poids pour l'estimation par Horvitz-Thompson sont :

$$w_k = \frac{1}{\pi_k} = \frac{N}{n}$$

On peut dire que l'individu k "représente" $w_k = \frac{N}{n}$ individus de la population \mathcal{U} .

Attention, w_k n'est pas un effectif (en particulier, w_k n'est pas forcément entier !)

Précision

Théorème

*En utilisant la formule de Yates-Grundy, la **vraie** variance des estimateurs d'Horvitz-Thompson s'écrit :*

$$\text{Var}(\bar{y}) = (1 - f) \frac{S^2}{n}$$
$$\text{Var}(T(\hat{Y})) = N^2(1 - f) \frac{S^2}{n}$$

Précision

Démonstration.

$$\begin{aligned}\text{Var}[\hat{Y}] &= \frac{1}{N^2} \text{Var}[T(\hat{Y})] \\ &= \frac{-1}{2N^2} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}, l \neq k} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \Delta_{kl} \\ &= \frac{1}{2N^2} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}, l \neq k} \left(\frac{y_k N}{n} - \frac{y_l N}{n} \right)^2 \frac{n(N-n)}{N^2(N-1)} \\ &= \frac{N-n}{nN} \frac{1}{2N(N-1)} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}, l \neq k} (y_k - y_l)^2 \\ &= \frac{N-n}{nN} S^2 \\ &= (1-f) \frac{S^2}{n}\end{aligned}$$

Estimation de la précision

Théorème

La variance empirique (ou dispersion) dans l'échantillon

$s^2 = \frac{1}{n-1} \sum_{k \in s} (y_k - \bar{y})^2$ est un estimateur sans biais de

$$S^2 = \frac{1}{N-1} \sum_{k \in \mathcal{U}} (Y_k - \bar{Y})^2$$

Estimation de la précision

Démonstration.

$$\begin{aligned}\mathbb{E}[s^2] &= \mathbb{E} \left[\frac{1}{n-1} \sum_{k \in s} (y_k - \bar{y})^2 \right] \\ &= \mathbb{E} \left[\frac{1}{2n(n-1)} \sum_{k \in s} \sum_{l \in s, l \neq k} (y_k - y_l)^2 \right] \\ &= \frac{1}{2n(n-1)} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}, l \neq k} (y_k - y_l)^2 \mathbb{E}(\delta_k \delta_l) \\ &= \frac{1}{2n(n-1)} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}, l \neq k} (y_k - y_l)^2 \frac{n(n-1)}{N(N-1)} \\ &= \frac{1}{2N(N-1)} \sum_{k \in \mathcal{U}} \sum_{l \in \mathcal{U}, l \neq k} (y_k - y_l)^2 \\ &= S^2\end{aligned}$$

Estimation de la précision

On peut estimer sans biais la variance de l'estimateur d'Horvitz-Thompson par :

$$\hat{\text{Var}}(\bar{y}) = (1 - f) \frac{s^2}{n}$$
$$\hat{\text{Var}}(T(\hat{Y})) = N^2(1 - f) \frac{s^2}{n}$$

Partie 2

Estimation d'une proportion

Estimation d'une proportion

On cherche à estimer P la proportion d'individus portant une caractéristique dans la population \mathcal{U} .

p , la proportion dans s d'individus portant la caractéristique, est un estimateur sans biais de P .

Variance

Sa *vraie* variance vaut :

$$\text{Var}(p) = (1 - f) \frac{N}{N - 1} \frac{P(1 - P)}{n}$$

On l'estime par :

$$\hat{\text{Var}}(p) = (1 - f) \frac{p(1 - p)}{n - 1}$$

Précision

Demi-longueur de l'intervalle de confiance :

$$L = 2\sqrt{(1-f)\frac{p(1-p)}{n-1}}$$

Coefficient de variation estimé :

$$\begin{aligned}\hat{C}V(p) &= \frac{\sqrt{\hat{V}\text{ar}(p)}}{p} \\ &= \sqrt{(1-f)\frac{1}{n-1}\frac{1-p}{p}}\end{aligned}$$

Taille pour une précision absolue donnée

On fixe L ("précision absolue"). Si $f \approx 0$, on a :

$$n \approx \frac{4p(1-p)}{L^2}$$

C'est souvent le cas lorsque qu'on s'intéresse à une grande population.

Taille pour une précision absolue donnée

Cas général (f pas forcément petit) : on note z le quantile d'ordre $1 - \frac{\alpha}{2}$ de la loi $\mathcal{N}(0, 1)$:

$$n = \frac{1 + n_0}{1 + \frac{n_0}{N}}$$

$$\text{avec : } n_0 = \frac{z^2 p(1 - p)}{L^2}$$

Souvent $z = 2$, quantile à 95% de la normale centrée réduite.

Taille pour une précision relative donnée

La précision relative δ est définie par le rapport de la demi-longueur de l'intervalle de confiance à l'estimation :

$$\delta = \frac{2\hat{\sigma}}{p}$$

On peut se ramener au coefficient de variation simplement :

$$\delta = 2\hat{C}V(p) = \sqrt{(1-f) \frac{1-p}{p(n-1)}}$$

Taille pour une précision relative donnée

De manière équivalente à la précision absolue L , on peut fixer le coefficient de variation $\hat{C}\hat{V}(p)$. Dans ce cas, et si $f \approx 0$:

$$n \approx \frac{1 - p}{p(\hat{C}\hat{V}(p))^2}$$

Taille pour une précision relative donnée

Taille de l'échantillon pour une précision relative de $\pm\delta\%$ selon la valeur de la proportion recherchée :

	0,05	0,10	0,20	0,30	0,40	0,50
1 %	760000	360000	160000	93333	60000	40000
2 %	190000	90000	40000	23333	15000	10000
3 %	84444	40000	17778	10370	6667	4444
4 %	47500	22500	10000	5833	3750	2500
5 %	30400	14400	6400	3733	2400	1600
10 %	7600	3600	1600	933	600	400

Exemple

Exemple d'application : la législation sur la méthode des quotas, en France.

- `http://www.commission-des-sondages.fr/oblig/instituts.htm`
- `http://www.ipsos.fr/faq`

Partie 3

Échantillonnage dans le temps

Problème

On veut estimer l'évolution de la moyenne d'une variable Y entre deux dates 1 et 2 : $\Delta Y = \bar{Y}_1 - \bar{Y}_2$

Méthode 1

Méthode 1 : On tire deux échantillons indépendants aux dates 1 et 2, selon un sondage aléatoire simple.

On a alors : $\Delta\hat{Y} = \bar{y}_2 - \bar{y}_1$ un estimateur sans biais de ΔY , de variance :

$$\text{Var}(\Delta\hat{Y}) = \text{Var}(\bar{y}_1) + \text{Var}(\bar{y}_2)$$

Méthode 2 : panel

Méthode 2 : On utilise un panel, c'est-à-dire que l'on tire un échantillon en date 1, et on le réinterroge à la date 2. On a alors : $\hat{\Delta Y} = \bar{y}_2 - \bar{y}_1$ un estimateur sans biais de ΔY , de variance :

$$\text{Var}(\hat{\Delta Y}) = \text{Var}(\bar{y}_1) + \text{Var}(\bar{y}_2) - 2\text{Cov}(\bar{y}_1, \bar{y}_2)$$

$$\text{où : } \text{Cov}(\bar{y}_1, \bar{y}_2) = (1 - f) \frac{S_{12}}{n}$$

$$\text{et : } S_{12} = \frac{1}{N-1} \sum_{k \in \mathcal{U}} (Y_{1k} - \bar{Y}_1)(Y_{2k} - \bar{Y}_2)$$

Méthode 2 : panel

Dans les bons cas, on a : $S_{12} > 0$, d'où :

$$\text{Var}(\hat{\Delta Y}) < \text{Var}(\bar{y}_1) + \text{Var}(\bar{y}_2)$$

Exemple : enquête emploi à l'INSEE

Année	Trimestre	Sous-échantillons					
2016	T1	6	5	4	3	2	1 →
	T2	→ 7	6	5	4	3	2
	T3	8	7	6	5	4	3
	T4	9	8	7	6	5	4
2017	T1	10	9	8	7	6	5
	T2	11	10	9	8	7	6
	T3	12	11	10	9	8	7
	T4	13	12	11	10	9	8

Partie 4

Estimation sur un domaine

Notations

$\mathcal{U}_d \subset \mathcal{U} =$ sous-population d'intérêt

$N_d =$ taille de \mathcal{U}_d (connue ou inconnue)

$P_d = \frac{N_d}{N} =$ taille relative de \mathcal{U}_d

$Q_d = 1 - P_d$

$s_d = s \cap \mathcal{U}_d$

$n_d =$ taille de s_d

$p_d = \frac{n_d}{n} =$ taille relative de s_d

$q_d = 1 - p_d$

Estimation de la taille d'un domaine

On définit sur \mathcal{U} la variable Z indicatrice d'appartenance au domaine :

$$Z_k = 1 \text{ si } k \in \mathcal{U}_d$$

$$Z_k = 0 \text{ sinon}$$

Estimation de la taille d'un domaine

Alors :

$$T(Z) = \sum_{k \in \mathcal{U}} Z_k = N_d$$

$$\bar{Z} = \frac{N_d}{N} = P_d$$

$$\bar{z} = \frac{1}{n} \sum_{k \in s} z_k = p_d$$

$$S^2 = \frac{N}{N-1} P_d Q_d$$

$$s^2 = \frac{n}{n-1} p_d q_d$$

Estimation de la taille d'un domaine

Théorème

$$\hat{N}_d = N \cdot p_d = N \cdot \frac{n_d}{n} \text{ est un estimateur sans biais de } N_d$$
$$\hat{P}_d = p_d \text{ est un estimateur sans biais de } P_d$$

Estimation de la taille d'un domaine

Démonstration.

Toutes ces quantités s'écrivent sous la forme d'un total (via Z) et correspondent à l'estimateur d'Horvitz-Thompson, qui est sans biais. □

Estimation de la taille d'un domaine

On a aussi :

$$\text{Var}(\hat{N}_d) = N^2(1-f) \frac{N}{N-1} \frac{P_d Q_d}{n}$$

$$\text{Var}(\hat{P}_d) = \text{Var}(p_d) = (1-f) \frac{N}{N-1} \frac{P_d Q_d}{n}$$

$$\hat{\text{Var}}(\hat{N}_d) = N^2(1-f) \frac{p_d q_d}{n-1}$$

$$\hat{\text{Var}}(\hat{P}_d) = \hat{\text{Var}}(p_d) = (1-f) \frac{p_d q_d}{n-1}$$

Estimation d'un total sur un domaine

On veut estimer le total $T_{\mathcal{U}_d}(Y)$ d'une variable Y sur le domaine \mathcal{U}_d . On définit sur \mathcal{U} la variable Y^d par :

$$Y_k^d = Y_k \text{ si } k \in \mathcal{U}_d$$

$$Y_k^d = 0 \text{ sinon}$$

Alors le total à estimer s'écrit :

$$T(Y^d) = \sum_{k \in \mathcal{U}} Y_k^d = \sum_{k \in \mathcal{U}_d} Y_k = T_{\mathcal{U}_d}(Y)$$

Estimation d'un total sur un domaine

Un estimateur sans biais de $T_{\mathcal{U}_d}(Y)$ est :

$$\hat{T}_{\mathcal{U}_d}(Y) = \frac{n_d}{n} N \bar{y}_d$$

où : $\bar{y}_d = \frac{1}{n_d} \sum_{k \in \mathcal{S}_d} y_k$

Estimation d'un total sur un domaine

Et pour ce qui est de la précision :

$$\text{Var}(\hat{T}_{\mathcal{U}_d}(Y)) = N^2(1-f) \frac{S_{Y^d}^2}{n} \text{ avec } S_{Y^d}^2 = \frac{1}{N-1} \sum_{k \in \mathcal{U}} (Y_k^d - \bar{Y}^d)^2$$

$$\hat{\text{Var}}(\hat{T}_{\mathcal{U}_d}(Y)) = N^2(1-f) \frac{s_{Y^d}^2}{n} \text{ avec } s_{Y^d}^2 = \frac{1}{n-1} \sum_{k \in s} (y_k^d - \bar{y}^d)^2$$

avec :

\bar{Y}^d = moyenne de Y^d sur \mathcal{U}

\bar{y}^d = moyenne de Y^d sur s

Estimation d'un total sur un domaine

Remarque sur la précision : Si on pose :

$$\bar{Y}_d = \frac{1}{N_d} \sum_{k \in \mathcal{U}_d} Y_k = \text{moyenne de } Y \text{ sur } \mathcal{U}_d$$

$$\bar{S}_d^2 = \frac{1}{N_d - 1} \sum_{k \in \mathcal{U}_d} (Y_k - \bar{Y}_d)^2 = \text{dispersion de } Y \text{ sur } \mathcal{U}_d$$

alors on a :

$$\text{Var}(\hat{T}_{\mathcal{U}_d}(Y)) \sim N_d^2 \left(\frac{1}{\mathbb{E}(n_d)} - \frac{1}{N_d} \right) \left[\frac{1 - \frac{1}{N_d}}{1 - \frac{1}{N}} S_d^2 + \frac{N - N_d}{N - 1} \bar{Y}_d^2 \right]$$

C'est donc la taille (attendue) de l'échantillon dans le domaine qui est déterminante et non n .

Estimateur alternatif pour le total

Si on connaît la taille du domaine N_d , un autre estimateur "naturel" de $T_{\mathcal{U}_d}(Y)$ est :

$$\hat{T}_{\mathcal{U}_d}^{alt}(Y) = N_d \bar{y}_d$$

C'est-à-dire que l'on remplace un estimateur sans biais de N_d :

$\hat{N}_d = \frac{n_d}{n} N$ par N_d . En général, $\hat{T}_{\mathcal{U}_d}^{alt}(Y)$ est préférable à $\hat{T}_{\mathcal{U}_d}(Y)$.

Estimation de la moyenne sur un domaine

On veut estimer : $\bar{Y}_d = \frac{T_{U_d}(Y)}{N_d}$. On peut utiliser :

$$\hat{Y}_d = \frac{\hat{T}_{U_d}(Y)}{N_d} \text{ si on connaît } N_d$$

$$\hat{Y}_d^{alt} = \frac{\hat{T}_{U_d}^{alt}(Y)}{N_d} = \bar{y}_d \text{ que l'on connaisse } N_d \text{ ou non !}$$

Ce dernier estimateur est assez intuitif (plugin !), et est en général meilleur que le premier.

Partie 5

Estimation d'un ratio

Estimation d'un ratio

On cherche à estimer le rapport des totaux (ou des moyennes) de deux variables X et Y :

$$R = \frac{T(X)}{T(Y)} = \frac{\bar{X}}{\bar{Y}}$$

Attention ! L'estimateur d'Horvitz-Thompson est sans biais quand on estime un total ou une moyenne.

Estimation d'un ratio

On peut utiliser l'estimateur :

$$\hat{R} = \frac{T(\hat{X})}{T(\hat{Y})} = \frac{\hat{X}}{\hat{Y}} = \frac{\bar{x}}{\bar{y}}$$

Son biais s'écrit :

$$B(\hat{R}) \approx -\frac{1}{\bar{X}^2}(1-f)\frac{S_{XY} - RS_X^2}{n}$$

où : $S_{XY} = \frac{1}{N-1} \sum_{k \in \mathcal{U}} (y_k - \bar{y})(x_k - \bar{x})$

Précision de l'estimateur du ratio

Son écart quadratique moyen et l'EQM estimé s'écrivent :

$$EQM(\hat{R}) = \frac{1-f}{n\bar{X}^2} (S_Y^2 + R^2 S_X^2 - 2RS_{XY})$$

$$E\hat{Q}M(\hat{R}) = \frac{1-f}{n\bar{X}^2} (s_Y^2 + \hat{R}^2 s_X^2 - 2\hat{R}s_{XY})$$

Conclusion sur le SAS

- Les estimateurs ont une forme simple
- Ne nécessite aucune information sur les individus de la base de sondage
- Est essentiel pour comprendre les plans de sondage plus complexes
- Peut permettre d'approximer les plans de sondage plus complexes