

Jeudi, 26 Septembre 2024

Chapitre 01: CONSIDERATIONS PRATIQUES DANS L'ELABORATION D'UN PLAN DE SONDAGE.

1- Introduction

Le Sondage statistique est une technique qui consiste à enquêter sur un phénomène auprès d'un échantillon ^{individuel} sélectionné selon certaines règles scientifiques pour représenter toute la population dont il est issu. L'exemple le plus connu est le sondage d'opinion développés aux Etats-Unis dès le début du XIX^e siècle.

On distingue le sondage probabiliste et le sondage non probabiliste. Dans le premier cas la probabilité de chaque individu est connue d'avance. Dans le second cas, on ne peut calculer la probabilité de sélection ou d'inclusion. Par ailleurs, la décision d'un statisticien ou d'un chercheur de réaliser une enquête par sondage doit soigneusement être pensée. Cette décision doit tenir compte non seulement des objectifs de l'étude mais aussi d'autres facteurs tel que les contraintes budgétaires, la qualité de sondage.

Ce chapitre est consacré aux différentes considérations à prendre en compte dans le choix de réaliser une enquête par sondage aléatoire.

2 - Contexte et Justification d'une enquête

Le contexte doit permettre d'appréhender la situation en terme d'étude et de production de données statistiques pour des prises de décisions dans les domaines sélectionnés. Il doit en outre faire ressortir le besoin à combler en réalisant une collecte de données. Cette analyse situationnelle devra finalement déboucher sur la justification de réaliser une enquête par rapport aux besoins d'indicateurs recherchés.

3 - Formulation des objectifs

L'élaboration de l'énoncé des objectifs est processus itératif qui engage les producteurs et les utilisateurs des données statistiques. Les étapes du processus sont les suivantes :

- les besoins d'informations
- les utilisateurs et les utilisations des données
- les principaux concepts et les données opérationnelles
- le choix de l'enquête

→ le plan d'analyse

Exemple : On a besoin de suivre le pouvoir d'achat de la population, en calculant la dépense par tête et sa répartition selon les fonctions de consommation.

La dépense par tête est obtenue en divisant la dépense totale de consommation par l'effectif total de la population. L'enquête doit alors mesurer la composition des ménages, les conditions de logement et leur dépense de consommation.

On distinguera ainsi, l'objectif générale de l'enquête qui consiste à suivre le pouvoir d'achat de la population et les objectifs spécifiques qui consistent à mesurer :

- a - les dépenses monétaires des consommateurs des ménages
- b - l'auto consommation
- c - les transferts reçus en nature
- d - les caractéristiques socio-économiques et démographiques de la population

4 - Concepts et indicateurs

Les principaux concepts de l'enquête, doivent être bien définis. En tenant compte de l'exemple précédant, il faudra

definis le ménage, la dépense monétaire, l'auto-consommation, les sources des revenus des ménages, le cycle de vie, la situation d'emploi de la population.

Les concepts sont généralement abstraits. Pour les opérationnaliser, il est requis la définition des indicateurs de suivi et des variables qui contribueront aux calculs des indicateurs et variables ou à l'explication du phénomène qui est suivi. L'inventaire des indicateurs à calculer, et des données à collecter est indispensable pour l'élaboration des questionnaires et des outils de collecte.

Le tableau suivant illustre l'opérationnalisation d'un concept par la mesure d'un indicateur.

Tableau 1 : Opérationnalisation d'un concept par la mesure d'un indicateur

Concepts	Indicateurs	Variables requises
dépense de consommation par tête ; dépense effective par tête ; en moyenne pour un ménage pour satisfaire ses besoins de consommation finale	dépense de consommation totale effective de la population	<ul style="list-style-type: none"> Noms des biens et services consommés Dépense de consommation calculée pour l'ensemble des ménages. Taille des ménages : ou unités de consommation

5- Champs géographique et social de l'enquête

Le champ social d'une enquête est l'ensemble des personnes cibles de l'enquête. Dans l'exemple précédent, la population cible peut être l'ensemble des ménages toute catégories confondues.

La définition du champ géographique implique le choix d'une enquête. Il y a lieu de savoir, si les résultats attendus sont significatifs au niveau national ou infra-national.

La définition du champ géographique et du champ social de l'enquête induit des questions suivantes pour la méthodologie de collecte :

- Qui enquêter
- Où enquêter
- Comment enquêter
- Quand enquêter

6- Etape de conception d'une enquête

- Concept d'
- Contexte et justification
 - Définition des objectifs
 - Concept et indicateurs
 - Champs de l'étude
 - Plan de sondage

- Conception du questionnaire
- Collecte de données
- Traitement de données
- Analyse et diffusion des données
- Documentation de l'enquête

02-10-2024

7- Détermination de la taille de l'échantillon.

La plus grande question à laquelle doit répondre un statisticien d'enquête, est avant tout, la détermination de la taille de l'échantillon des unités statistiques à enquêter. Dans la détermination de la taille de l'échantillon, il faut tenir compte de ces préoccupations à savoir :

- la précision des estimations
- les contraintes de mise en œuvre du plan de sondage
- l'efficacité du plan de sondage

a- Précision des estimations

Le recours à une base de données d'enquêtes antérieures ayant déjà permis de calculer les indicateurs recherchés est très important. Il permet de tirer les enseignements pour la réalisation d'une nouvelle enquête. Les enseignements peuvent porter sur la taille de l'échantillon, le

type de plan de sondage, les outils utilisés pour la collecte des données. Les difficultés rencontrées aussi bien pendant la collecte des données qu'à la phase du traitement et d'analyse des résultats (taux de réponses, erreurs d'observation, etc), les estimations des principaux indicateurs et leurs précisions.

Toutes ces informations doivent aider le statisticien à effectuer des simulations en vue de la détermination de la taille de l'échantillon requise pour la réalisation de l'enquête prévue. Par ailleurs, la taille de l'échantillon doit être déterminée en fonction de :

- la fréquence d'apparition des événements du phénomène étudiés
- le niveau des statistiques d'analyse

Toutes ces formules qui seront présentées dans les chapitres suivants, peuvent être utilisées pour le calcul de la taille d'échantillon lors d'une enquête.

b- Contraintes de mise en œuvre du plan de sondage

Plusieurs contraintes peuvent avoir des incidences sur la détermination de la taille de l'échantillon à enquêter. La plus importante est la contrainte budgétaire. Le budget d'une enquête doit être réparti en coût fixe et coût variable.

Les coûts fixes sont généralement indépendants de la taille de l'échantillon à enquêter (coordonnées, fonctionnement, organisation des ateliers de formation des agents etc.)

Les coûts variables sont étroitement liés à la taille de l'échantillon (nombre de questionnaires, déplacement et salaire des agents de terrain et exploitation des données etc.)

c- efficacité du plan de sondage

Le plan de sondage efficace est celui qui offre la plus grande précision pour une taille d'échantillon prenant en compte les contraintes de mise en œuvre. Par exemple, si la population concernée par un phénomène à étudier est de petite taille, il est préférable de réaliser une enquête par sondage aléatoire simple à condition que il y ait une bonne maîtrise des coûts de déplacement.

Le statisticien peut aussi décider de réaliser un sondage stratifié pour améliorer la précision des résultats à condition de bien choisir les critères de stratification et de ne pas avoir plusieurs strates qui n'apportent rien de gain d'amélioration des estimations.

8- base de sondage

Une autre étape à franchir est le choix de la base de sondage de bonne qualité :

→ Couverture exhaustive

→ Pas de double enregistrement

→ Pas d'unité morte

→ Non basé essentiellement sur des informations anciennes type

Comme exemple de base de sondage, on peut citer une base de sondage aérolienne constituée des images satellitaires et de la géolocalisation.

et un exemple plus typique (plus utilisé)

Le terme base de sondage désigne généralement une liste concrète d'unités ayant un lien avec la population à étudier. Il s'agit de la description d'éléments déjà existants sous forme de liste, annuaire à partir desquels on peut constituer des unités et sélectionner un ensemble d'unités à enquêter.

La base de sondage doit comprendre toute les informations auxiliaires (mesure de taille, données démographiques etc.) nécessaires pour la mise en œuvre technique spéciale de sondage tel que la stratification ou le type de tirage etc.

La construction d'une base de sondage

doit tenir compte de :

→ la nature de la population à utiliser
exple : individu, famille, ménage, exploi-
tation

→ la répartition géographique de la po-
pulation (elle est due permise à une lo-
calité, une région ou est-elle répartie sur
l'ensemble du pays)

→ la nature des opérations de terrain
(enquête par téléphone, par correspondance
ou interview direct par des enquêteurs)

Une base de sondage peut convenir,
pour une enquête donnée et ne pas être
adaptée pour un autre type d'enquête.
Les unités statistiques dans la base doi-
vent être de taille assez homogène. Le
nombre d'unité dans la base doit être
connu pour permettre les extrapolations

9. Principaux paramètres à estimer par une enquête

Les principales caractéristiques d'une
variable d'intérêt sont souvent estimées
lors d'une enquête par sondage :

→ le total

→ la moyenne (la proportion est une
moyenne calculée pour une variable qua-
litative)

→ le ratio

→ la variance et le coefficient de

Variation des estimateurs

→ l'effet de sondage

10. Type d'erreur rencontrés dans une enquête par sondage

a. erreur d'échantillonnage

Elle est due au plan de sondage et est
mesurable par le biais, la variance ou
l'écart quadratique moyen. L'estimateur de
l'erreur de l'échantillonnage suit une loi de
probabilité.

b. erreur d'observation ou erreur de mesure

C'est une erreur liée aux dispositifs et
au support de collecte des données. Elle sur-
vient surtout dans la formulation des ques-
tions ou lors de la manipulation des
instruments de mesure, elle est très difficile
à quantifier à moins de retourner sur le ter-
rain.

c. erreur due au défaut de couverture ou à la non réponse

Elle est due d'une part, à l'utili-
sation d'une base de sondage incomple-
te et d'autre part à la non réponse
complète ou partielle à l'enquête par cer-
tains individus. Les non réponses peuvent
être régularisées ou imputées.
Par contre il est difficile de corriger à pos

postérieur l'erreur due au défaut de non couverture

11 - Choix et justification du plan de sondage

Les différents types de sondage que nous étudierons sont :

- le sondage aléatoire simple
- le sondage à probabilité inégale
- le sondage stratifié
- le sondage en grappes
- le sondage à plusieurs degrés

Le choix du plan de sondage doit être justifié. Par exemple, lorsqu'il s'agit de réaliser l'étude d'un phénomène sur une population de petite taille et bien connue, le statisticien peut décider de réaliser un sondage aléatoire simple ou avec stratification si des variables de stratification existent. Par contre, la réalisation d'une étude qui concerne une grande population pourra utiliser un plan de sondage par grappes ou à plusieurs degrés avec les possibilités de stratification.

12 - Notations

a - Sur l'univers (Population totale)

→ Unité statistique $X = 1, \dots, N$

→ Moyenne de la variable Y

$$\mu = \frac{1}{N} \sum_{\alpha=1}^N Y_{\alpha}$$

→ Variance de $Y : V(Y) = \frac{1}{N} \sum_{\alpha=1}^N (Y_{\alpha} - \mu)^2$

On utilise souvent la notation σ^2 pour la variance. $V(Y) = \sigma^2$. Par ailleurs, on définit aussi la variance corrigée $S^2 = \frac{1}{N-1} \sum_{\alpha=1}^N (Y_{\alpha} - \mu)^2$

b - Sur l'échantillon

→ Unité statistique $i = 1, \dots, n$

→ $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ moyenne de la variable Y calculée sur l'échantillon. \bar{y} est une variable aléatoire, $E(\bar{y})$ est l'espérance de la variable aléatoire \bar{y} , $V(\bar{y})$ est la variance de variable aléatoire \bar{y} qui est rappelons le, la variance de l'estimateur \bar{y} en fonction de la variance de la variable Y calculée sur les unités de l'échantillon.

c - Taux de sondage

On le note $f = \frac{n}{N}$

Marsred, 03 Octobre 2024

Chapitre 02 :

I - DEFINITIONS

Le sondage aléatoire simple est le modèle d'échantillonnage en apparence le plus simple que l'on puisse imaginer. Il consiste à considérer que dans une population d'effectif (N) , tous les échantillons de n unités sont possibles avec la même probabilité.

Un plan de sondage est dit avec remise si un même individu peut apparaître plusieurs fois dans l'échantillon et si l'ordre dans lequel apparaissent les individus compte.

exple : $P = \{\text{bleu, blanc, rouge, violet, vert}\}$

$n = 3$

l'échantillon $\{\text{blanc, blanc, rouge}\}$ est différent de l'échantillon $\{\text{blanc, rouge, blanc}\}$

Dans le cas d'un plan avec remise, il y a N^n échantillon possible.

Un plan de sondage est dit sans remise si un individu ne peut apparaître qu'une seule fois dans l'échantillon.

Dans l'exple précédent, l'échantillon $\{\text{blanc, blanc, rouge}\}$ n'est pas possible. Dans le cas d'un plan sans remise, il y a

$$\frac{N!}{n!(N-n)!} \text{ échantillon possible}$$

La plus part du temps nous nous intéressons au plan sans remise. Ces échantillons de deux fois le même individu n'apporte pas d'information supplémentaire. Cependant il n'est pas intéressant de considérer des plans avec remise, nécessaires que pour servir d'élément de comparaison et de référence.

Un plan de sondage aléatoire est dit simple ou à probabilités égales si chaque échantillon a la même probabilité qu'un autre d'être tiré au sort.

Exple : Dans le cas simple d'un plan sans remise, un échantillon de taille fixe n a donc une probabilité

$$\frac{1}{\frac{N!}{n!(N-n)!}}$$

d'être tiré au sort, puis qu'il est com
Si $N=5$ et $n=2$, cette probabilité est

$$\text{donc } \pi_k = \frac{2 \times 3 \times 2}{5 \times 4 \times 3 \times 2} = \frac{1}{10}$$

Remarque : les données concernant la population toute entière (X_i, μ, T, P, \dots etc.) sont inconnues et déterministes, puisque l'on a pas accès aux informations concernant toute la population. En revanche, les valeurs obtenues à partir de l'échantillon sont connues et aléatoires, elles dépendent en effet du hasard puisqu'elles varient d'un échantillon aléatoire à un autre et elles sont connues puisqu'on dispose des informations nécessaires pour les calculer sur l'échantillon.

II - PLAN SIMPLE SANS REMISE.

1- Plan de sondage et probabilité d'inclusion.

Soit un échantillon aléatoire S , la variable aléatoire $1\{k \in S\}$, $k \in U$ nous sera très utile. Notons que c'est bien une variable aléatoire puisque S est aléatoire (variable indicatrice de Confield).
 $I_k = (1 \text{ quand } k \in S \text{ et } 0 \text{ quand } k \notin S)$

Definition

La probabilité d'inclusion de la k -ième

unité notée π_k correspond à la probabilité que cette k -ième unité = $\pi_k = P(k \in S)$
 $= \sum_{s \in \mathcal{S}} P(s) 1\{k \in s\}$

$$= \sum_{s \in \mathcal{S}} P(s) k \in U$$

Notons également que par définition, π_k c'est $E(1\{k \in S\})$.

De même nous pouvons définir les probabilités d'ordre supérieur

La probabilité d'inclusion d'ordre 2 est la probabilité que deux (ou) unités distinctes appartiennent simultanément à un échantillon. C'est à dire :

$$\pi_{kl} = P(k \in S, l \in S) = \sum_{s \in \mathcal{S}} P(s) k, l \in U, k \neq l$$

Notons comme précédemment que :

$$\pi_k = E(1\{k \in S\}), \text{ on a :}$$

$$\text{Var}(1\{k \in S\}) = E(1\{k \in S\}^2) - E(1\{k \in S\})^2 = \pi_k(1 - \pi_k)$$

$$\text{et } \text{Cov}(1\{k \in S\}, 1\{l \in S\}) = E(1\{k \in S\} 1\{l \in S\}) - E(1\{k \in S\}) E(1\{l \in S\}), k, l \in U, k \neq l$$

Dans la suite on notera :

$$\Delta_{kl} = \begin{cases} \text{Cov}(1\{k \in S\}, 1\{l \in S\}), k \neq l \\ \text{Var}(1\{k \in S\}), k = l \end{cases}$$

Il est souvent utile de connaître les probabilités d'inclusion afin d'établir des estimateurs et sa variance. Les probabilités d'inclusion se calculent com-

$$m_{\text{Holt}} = T_k = \sum_{s \in R} P(s) = \frac{\binom{n-1}{N-1} \binom{N-1}{n} - 1}{(n-1)!(N-n)!} \cdot \frac{n!(N-n)!}{N!} = \frac{n}{N}$$

nbre d'échantillon contenant \$k\$

$$T_{k \neq l} = \sum_{s \in R, k \neq l} P(s) = \frac{\binom{n-2}{N-2} \binom{N-2}{n} - 1}{(n-2)!(N-n)!} \cdot \frac{n!(N-n)!}{N!} = \frac{n(n-1)}{N(N-1)}$$

De ces deux (ou) expressions, on déduit :

$$\Delta_{k \neq l} = \begin{cases} \frac{n(n-1)}{N(N-1)} - \frac{n^2}{N^2} = -\frac{n(N-n)}{N^2(N-1)}, & k \neq l \\ \frac{n}{N} \left(1 - \frac{n}{N}\right) = \frac{n(N-n)}{N^2}, & k = l \end{cases}$$

2. Le \$\pi\$-estimateur

Les estimateurs de Horvitz-Thompson sont des outils essentiels pour les statisticiens, permettant d'obtenir des estimations possibles précises, même dans des contextes de sélection inégale. La maîtrise de leur utilisation peut permettre aux chercheurs d'améliorer la qualité de leurs analyses et de leurs résultats.

Dans de nombreux cas, les unités d'une population n'ont pas toute la même probabilité d'être sélectionnées dans un échantillon. Par exemple dans un sondage sur les

ménages, il peut être plus probable de sélectionner certains types de ménage (comme ceux avec un revenu plus élevé) en raison de la structure du sondage.

Les estimateurs de Horvitz-Thompson corrigent cette inégalité en ajustant les estimations.

* Estimation d'un total et d'une moyenne

Horvitz-Thompson (1952) ont introduit un estimateur linéaire sans biais d'un total \$y\$ pour tout plan de sondage.

$$\hat{y}_{HT} = \sum_{k \in S} \frac{y_k}{\pi_k}$$

Cette estimateur est appelé

le \$\pi\$-estimateur, l'estimateur d'Horvitz-Thompson ou encore l'estimateur des inverses de probabilités.

Théorème :

Si \$\pi_k > 0 \forall k \in U\$, alors \$\hat{y}_{HT}\$ estime \$y\$ sans biais.

Nous avons introduit le \$\pi\$-estimateur pour estimer le total \$y\$, mais nous pouvons également l'utiliser pour estimer la moyenne \$\bar{y}\$ par \$N - \hat{y}_{HT} = \frac{1}{N} \sum_{k \in S} \frac{y_k}{\pi_k}\$.

Néanmoins, toute fois que pour utiliser cet estimateur, il faut que la taille de la population soit connue, ce qui n'est pas toujours

raisonnement pas toujours à cas. Cependant plus que $N = \sum_{k \in U} 1$ on peut estimer N par Horvitz-Thompson c-à-d $\hat{N}_H = \sum_{k \in U} \frac{1}{\pi_k}$

* Variance du π -estimateur

Il est également possible de connaître la variance du π -estimateur.

Théorème :

Soit \hat{y}_π le π -estimateur d'un total y . Soit $\pi_k > 0$ pour tout $k \in U$ alors :

$$\text{Var}(\hat{y}_\pi) = \sum_{k, l \in U} \frac{y_k y_l}{\pi_k \pi_l} \Delta_{kl}$$

3 - Variance pour les plans de taille fixe

Dans le cas des plans de taille fixe, on peut réécrire la variance du π -estimateur sous une forme différente.

Théorème :

Soit \hat{y}_π d'un total y , si le plan est de taille fixe n et que

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

4 - Estimation de la Variance de π -estimateur

L'idée de base du π -estimateur, peut

être naturellement étendue au contexte des fonctions de deux (ou) variables $f(\cdot, \cdot)$.

Théorème

Soit $f(\cdot, \cdot)$ une fonction de deux variables quelconque.

Si $\pi_{kl} > 0$ pour tous $k, l \in U, k \neq l$ alors

$$\sum_{k, l \in U, k \neq l} \frac{f(k, l)}{\pi_{kl}} \cdot 1(k \in s, l \in s) \text{ est un estimateur sans biais de } \sum_{k, l \in U, k \neq l} \frac{f(k, l)}{\pi_{kl}}$$

Avec les probabilités d'inclusion calculées plus haut, nous pouvons donner une version plus explicite du π -estimateur. Le π -estimateur d'une moyenne μ_y sera :

$$\hat{\mu}_{y\pi} = \frac{1}{N} \sum_{k \in U} \frac{y_k}{\pi_k} 1(k \in s) = \frac{1}{N} \sum_{k \in U} y_k \frac{1(k \in s)}{\pi_k}$$

et le π -estimateur du total y est :

$$\hat{y}_\pi = N \hat{\mu}_{y\pi} \text{ avec } \hat{\mu}_{y\pi} = \frac{1}{n} \sum_{k \in U} y_k \frac{1(k \in s)}{\pi_k}$$

Puis que le plan est de taille fixe, et que les probabilités d'inclusion des deux premiers ordres sont strictement positives, on peut utiliser la formule de la variance de la $\mu_{y\pi}$ de Sen - Jais - Grundy c'est à dire $\text{Var}(\hat{\mu}_{y\pi}) = - \frac{1}{2N^2} \sum_{k, l \in U, k \neq l} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \Delta_{kl}$

$$= \frac{1}{2N^2} \times \frac{N-1}{2(N-1)} \sum_{k \neq l} (y_k - y_l)^2$$

$$= \frac{N-1}{N} \frac{\sum y_k^2}{n}$$

avec $s_y^2 = \frac{1}{N-1} \sum_{k \in U} (y_k - \mu_y)^2$

$$= \frac{1}{2N(N-1)} \sum_{k \neq l} (y_k - y_l)^2$$

Remarque :

La variance de l'estimateur du total peut également s'écrire :

$$\text{Var}(\hat{t}_{yT}) = N(N-1) \frac{s_y^2}{n}$$

Théorème

Pour un plan de taille fixe n simple et sans remise, la variance corrigée de la population s_y^2 est estimée sans biais par :

$$s_y^2 = \frac{1}{n-1} \sum_{k \in U} (y_k - \bar{y})^2 \quad \text{et} \quad \text{Var}(s_y^2) = \frac{1}{n-1} \sum_{k \in U} (y_k - \bar{y})^2$$

En fait, on peut estimer sans biais la variance de \hat{t}_{yT} par ces plans particuliers par $\text{Var}(\hat{t}_{yT}) = \frac{(N-1)}{N} \frac{s_y^2}{n}$ et pour

le π -estimateur du total \hat{t}_y par :

$$\text{Var}(\hat{t}_{yT}) = N(N-1) \frac{s_y^2}{n}$$

16-10-2024 II - PLANS SIMPLES AVEC REMISE

Le plan de taille fixe n simple et avec remise correspond au cadre de la statistique inférentielle. En effet, le plan de sondage consiste à sélectionner une unité aléatoire avec probabilité égale à $\frac{1}{N}$ et de recommencer l'opération n fois indépendamment. On est donc ramené au cas des variables aléatoires indépendantes et identiquement distribuées de moyenne μ_y .

$$\mu_y = \frac{1}{N} \sum_{k \in U} y_k \quad \text{et de variance}$$

$$s_y^2 = \frac{1}{N} \sum_{k \in U} (y_k - \mu_y)^2$$

La moyenne sur la population μ_y est estimée sans biais par $\hat{\mu}_y = \frac{1}{n} \sum_{k \in U} y_k$ et $\text{E}(\hat{\mu}_y) = \mu_y$. En effet $\text{E}(\hat{\mu}_y) = \frac{1}{n} \sum_{k \in U} \frac{y_k}{N} = \mu_y$.

De plus, puisque les y_k de l'échantillon sont sélectionnés indépendamment et sont de même loi, $\text{Var}(\hat{\mu}_y) = \frac{s_y^2}{n}$.

Théorème :

Pour un plan de taille fixe n simple et sans remise, la variance non corrigée de la population, $\hat{s}_y^2 = \frac{1}{n} \sum_{k \in U} (y_k - \hat{\mu}_y)^2$ est estimée sans biais par :

$$\sigma_y^2 = \frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2$$

Au final, la variance de \bar{y}_g pour un plan simple avec remise est estimée sans biais par $\text{Var}(\bar{y}_g) = \frac{\sigma_y^2}{n}$

TABIEAU RECAPITULATIF DES RESULTATS

	Sans Remise	Avec Remise
Estimateur de la moyenne	$\frac{1}{n} \sum_{k \in U} y_k$	$\frac{1}{n} \sum_{k \in U} y_k$
Variance de l'estimateur de la moyenne	$\frac{N-n}{N} \frac{\sigma_y^2}{n}$	$\frac{\sigma_y^2}{n}$
Estimateur de la variance de l'estimateur de la moyenne	$\frac{N-n}{N} \frac{s_y^2}{n}$	$\frac{s_y^2}{n}$

La plus part du temps, l'on s'intéresse essentiellement au plan sans remise car interroger deux fois le même individu n'apporte pas d'information supplémentaire de plus, le sondage simple et sans remise est toujours préférable à celui avec remise. En effet, si l'on appelle \bar{y}_{π} et $\bar{y}_{\pi r}$, les π -estimateurs de la moyenne avec et sans remise, alors $\text{Var}(\bar{y}_{\pi}) \leq \text{Var}(\bar{y}_{\pi r})$

$$\frac{\text{Var}(\bar{y}_{\pi})}{\text{Var}(\bar{y}_{\pi r})} = \frac{(N-n)}{N} \times \frac{s_y^2}{\sigma_y^2} =$$

$$= \frac{(N-n)}{N} \times \frac{n}{N-1} = \frac{N-n}{N-1} < 1$$

Variance du plan sans remise est plus petit qu'une variance du plan avec remise. Ainsi la précision est meilleure dans un plan sans remise qu'avec remise.

II - PLANS SIMPLES SANS REMISE ET FONCTION D'INDICET

* Estimation d'une proportion

Il est fréquent qu'une étude porte sur une estimation d'une proportion p . Estimer une proportion revient à compter le nombre d'unité y_k avec $k \in U$, possédant une caractéristique $y_k = 1$ si y_k possède la caractéristique $y_k = 1$ si y_k possède la caractéristique $y_k = 1$ si non $y_k = 0$.

Estimer une proportion c'est estimer une moyenne. Toute fois pour des proportions, les expressions pour la variance se voit considérablement simplifiée du fait que $y_k^2 = y_k$ $\forall k \in U$. Ainsi nous aurons pour un plan simple sans remise nous aurons :

$$\sigma_y^2 = \frac{1}{n-1} \sum_{k \in U} (y_k - p)^2 = \frac{n}{n-1} p(1-p)$$

$$\text{Var}(\bar{y}) = \frac{N-n}{N} \frac{s_y^2}{n} = \frac{N-n}{N-1} \frac{p(1-p)}{n}$$

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

V - DETERMINATION DE LA TAILLE DE L'ECHANTILLON

Avant de commencer un sondage, il est toujours souhaitable de se poser la question des incertitudes liées à nos estimations. Généralement, les limites budgétaires fixent la taille de l'échantillon et on se contente alors de répondre si le budget alloué est suffisante pour la précision donnée.

Pour précision donnée, nous attendons que le paramètre d'intérêt θ soit contenu dans un intervalle de confiance centré en $\hat{\theta}$ avec une probabilité d'au moins $1-\alpha$, c-à-d. d'avoir $P\{\theta \in [\hat{\theta}-e, \hat{\theta}+e]\} \geq 1-\alpha$ tel que $P\{\theta \in [\hat{\theta}-e, \hat{\theta}+e]\} \geq 1-\alpha$.

En supposant que notre estimateur $\hat{\theta}$ suit approximativement une loi normale (ce qui sera souvent le cas), on sait que $P\{\theta \in [\hat{\theta}-e, \hat{\theta}+e]\} = P\{\hat{\theta} - z_{1-\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\theta})} \leq \theta \leq \hat{\theta} + z_{1-\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\theta})}\} = 1-\alpha$

où $z_{1-\frac{\alpha}{2}}$ est le quantile d'une loi normale centrée réduite au non dépassement $1-\frac{\alpha}{2}$, c-à-d. $P\{Z \leq z_{1-\frac{\alpha}{2}}\} = 1-\frac{\alpha}{2}$, $Z \sim N(0,1)$.

Besoin : Puis que $\text{Var} \hat{\theta}$ dépend

de la taille de l'échantillon n , on cherchera la taille minimale n_0 induisant la précision requise.

Dans le cas de l'estimation de la moyenne μ_y pour un plan simple sans remise, on aura donc :

$$P\left\{\mu_y \in \left[\bar{y} - z_{1-\frac{\alpha}{2}} \sqrt{\frac{N-n}{N} \frac{S_y^2}{n}}, \bar{y} + z_{1-\frac{\alpha}{2}} \sqrt{\frac{N-n}{N} \frac{S_y^2}{n}}\right]\right\} = 1-\alpha$$

Exercice d'application

Soit une caractéristique X définie sur une population $N=4$ unités.

Individus : 1, 2, 3, 4, 5

Valeurs de X : 11, 10, 8, 11

1. Calculez la valeurs des paramètres suivants de la population : la moyenne, la variance et la variance corrigée notées respectivement μ , σ^2 et s^2 .

2. On tire un échantillon sans remise de taille $n=2$ à probabilité égale

- Combien d'échantillon peut-on tirer ?
- Pour chaque échantillon possible, calculez pour l'échantillon, la variance et la moyenne (s^2 et \bar{y}).
- Calculez $E(\bar{y})$, $\text{Var}(\bar{y})$ et $E(s^2)$.

Corrigé de la fiche de Td.

Exercice 1

$(j_k - \bar{y})^2$	1,36	2,06	13,36	6,76	0,16	41,2
	1	2	3	4	5	
	8	3	11	4	7	
						$N=5$

1- Calculons

- la moyenne \bar{y}

$$\bar{y} = \frac{1}{N} \sum j_k = \frac{8+13+11+4+7}{5}$$

$$\bar{y} = 6,6$$

- la variance corrigée S_y^2

$$S_y^2 = \frac{1}{N-1} \sum (j_k - \bar{y})^2$$

$$= \frac{41,2}{4} = 10,3$$

2- Listons tous les échantillons possibles

$(1,2)$ $(1,3)$ $(1,4)$ $(1,5)$

$(2,3)$ $(2,4)$ $(2,5)$

$(3,4)$ $(3,5)$ $(4,5)$

3-

j_k	$\frac{1}{j}$	$(j - \bar{y})^2$	$V(\frac{1}{j})$
$(1,2)$	5,5	1,21	4,25
$(1,3)$	3,5	8,41	2,25
$(1,4)$	6	0,36	4
$(1,5)$	7,5	0,81	0,25
$(2,3)$	7	0,16	16
$(2,4)$	3,5	9,61	0,25
$(2,5)$	5	2,56	4
$(3,4)$	7,5	0,81	12,25
$(3,5)$	3	5,76	4
$(4,5)$	5,5	1,21	2,25
		30,7	

$$E(\frac{1}{j}) = \frac{1}{10} \sum \frac{1}{j_k}$$

$$= \frac{5,5 + 3,5 + 6 + 7,5 + 7 + 3,5 + 5 + 7,5 + 3 + 5,5}{10}$$

$$= \frac{66}{10} = 6,6 = \bar{y}$$

d'où \bar{y} est un estimateur sans biais de \bar{y}

$$V(\frac{1}{j}) = \frac{1}{10} \sum (j_k - \bar{y})^2 = \frac{31,1}{10} = 3,11$$

$$4) V(\bar{y}) = (1 - \frac{n}{N}) \frac{S_y^2}{n}$$

$$= (1 - \frac{2}{5}) \frac{10,3}{2}$$

$$= \frac{3}{5} \times \frac{10,3}{2}$$

$$= 3,09$$

5) Estimation $\hat{y}(\bar{y}) = \frac{1}{n} \sum (y_k - \bar{y})^2$

Exercice 2

1- $n=2$, $P(\{1,2\}) = \frac{1}{2}$, $P(\{1,3\}) = \frac{1}{4}$

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

* Non.

2- $P(A \in S) = \sum_{S \ni 1} P(S) = \frac{1}{2} + \frac{1}{4} = \frac{3}{4}$

$P(B \in S) = \frac{1}{2} + \frac{1}{4} = \frac{3}{4}$

$P(C \in S) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$

3.

Éch.	μ
$\{1,2\}$	1,5
$\{1,3\}$	2
$\{2,3\}$	2,5

4- $E(\hat{\mu}) = \sum p_i \mu_i$

$= \frac{1}{2} \times 1,5 + \frac{1}{4} \times 2 + \frac{1}{4} \times 2,5$

$E(\hat{\mu}) = 1,875$

ou $\mu = \frac{1+2+3}{3} = 2$

d'où $\hat{\mu}$ est biaisé.

Exercice 4

Ind	1	2	3	4	Total
y_k	102,1	102,16	108,5	100,3	430,1
$(y_k - \bar{y})^2$	97,050	24,255	0,15	163,203	266,285

1- Calculons

$\bar{y} = \frac{1}{4} \sum_{k=1}^4 y_k = \frac{430,1}{4} = 107,525$

$\sigma_y^2 = \frac{266,285}{4} = 66,571$

2- On peut extraire $C_4^3 = 4$ échantillons de taille $n=3$

	$\hat{\mu}$	$\hat{\sigma}^2$	P_i
$S_1 = \{1,2,3\}$	107,2	116,228	1/4
$S_2 = \{1,2,4\}$	107,2	88,34	1/4
$S_3 = \{1,3,4\}$	107,166	77,88	1/4
$S_4 = \{2,3,4\}$	110,466	54,148	1/4
Total	430,038		

$E(\hat{\mu}) = \frac{1}{n} \sum p_k \hat{\mu}_k$

$= \frac{430,038}{4} = 107,5095$

$B(\hat{\mu}) = E(\hat{\mu}) - \mu$

$= -0,0005$

$$EOM = R^2(\hat{\mu}) + V(\hat{\mu})$$

$$V(\hat{\mu}) = \left(1 - \frac{n}{N}\right) \frac{\hat{S}^2}{n}$$

$$= \left(1 - \frac{12}{N}\right) \cdot \frac{1}{n} \left(\frac{45^2}{N-1} \right)$$

A.N

$$V(\hat{\mu}) = \left(1 - \frac{3}{4}\right) \times \frac{1}{3} \times \left(\frac{4 \times 66,571}{3} \right)$$

$$= \frac{1}{4} \times \frac{1}{3} \times \frac{4 \times 66,571}{3}$$

$$= 7,396$$

$$EOM = (-0,0005)^2 + 7,396$$

$$= 7,39600025$$

3 - Supposons que la femelle est l'individu 1.

	P_k
Q_1	115
Q_2	115
Q_3	115
N	215

$$E(\mu) = \sum P_k Q_k$$

$$= \frac{1}{5} (103,266 + 107,2 + 103,166) + \frac{2}{5} \times 1,10,46$$

$$= 10,8128$$

Exo 6

$$\sum_{k=1}^{100} u_k = 29,07$$

$$\sum_{k=1}^{100} u_k^2 = 154,553$$

1 - 1^{re} estimation

$$\hat{\mu} - \bar{u} = \frac{1}{n} \sum_{k=1}^n u_k$$

$$= \frac{1}{100} \sum_{k=1}^{100} u_k$$

$$= 29,07$$

2.

$$IC(\hat{\mu}) = \hat{\mu} \pm z_{1-\frac{\alpha}{2}} \sqrt{V(\hat{\mu})}$$

$$= \hat{\mu} \pm z_{1-\frac{\alpha}{2}} \sqrt{\left(1 - \frac{n}{N}\right) \frac{\hat{S}^2}{n}}$$

$$\text{on } \alpha = 0,05 \Rightarrow 1 - \frac{\alpha}{2} = 0,975$$

$$z_{0,975} = 1,96$$

$$s_k^2 = \frac{1}{n} \sum u_k^2 - \bar{u}^2$$

$$\Rightarrow \hat{S}^2 = \frac{n}{n-1} \left[\frac{1}{n} \sum u_k^2 - \bar{u}^2 \right]$$

$$z_{1-\frac{\alpha}{2}} \sqrt{V(\hat{\mu})}$$

$$= 1,96 \times \sqrt{\left(1 - \frac{n}{N}\right) \times \frac{1}{n} \times \frac{n}{n-1} \left[\frac{1}{n} \sum u_k^2 - \bar{u}^2 \right]}$$

$$= 1,96 \times \sqrt{\left(1 - \frac{100}{2010}\right) \times \frac{1}{100} \times (154,553 - (29,07)^2)}$$

$$\frac{R^2}{n} = \frac{1}{n} s^2$$

Mars 2021, 30 octobre 2021.

Chapitre 03 : SONDAGE A PROBABILITES INEGALES

I. CARACTERE AUXILIAIRE ET PROBABILITES D'INCLUSION.

Soit $X_k, k \in U$ les valeurs prises par le caractère auxiliaire. Notons que cela implique sa connaissance sur toute la population. Notre étude porte toujours sur une fonction d'intérêt telle que la moyenne ou le total d'un caractère Y . Le principe d'un plan à probabilités inégales, consiste à définir les probabilités d'inclusion du 1^{er} ordre proportionnelles aux X_k .

Rappelons que pour un plan de taille X , la variance du π -estimateur, du total Y est :

$$\text{Var}(\hat{Y}) = \frac{1}{2} \sum_{\substack{k, l \in U \\ k < l}} \left(\frac{X_k}{\pi_k} - \frac{X_l}{\pi_l} \right)^2 \Delta_{kl} \quad (1)$$

Si nous souhaitons minimiser (1), en jouant seulement sur les probabilités d'inclusion de 1^{er} ordre π_k , prendre $\pi_k = \frac{X_k}{\sum_{l \in U} X_l}$

(proportionnel à) $X_k, k \in U$, est un choix ~~raisonnable~~ judicieux puisque $\text{Var}(\hat{Y})$ est alors ~~minimale~~. Bien évidemment, cette approche est impossible puisqu'elle suppose que l'on connaît les valeurs prises par le caractère Y sur toute la population U (il serait alors inutile de faire un sondage).

En revanche, si nous disposons d'un caractère auxiliaire X connu sur toute la population et dont on pense qu'il est approximativement proportionnel à Y , alors on pourra à définir les probabilités d'inclusion du 1^{er} ordre proportionnelles aux X_k .

Remarque : Si au contraire, le caractère X n'est pas du tout proportionnel à Y , le plan de sondage selon alors ~~est inapproprié~~ et il sera préférable de prendre un plan simple.

Pour un plan de taille fixe $n, \sum_{k \in U} \pi_k = 1$. Puisque, on aura obtenu les probabilités d'inclusion proportionnelles aux X_k c-à-d $\pi_k = c X_k$ avec $c = \frac{n}{\sum_{k \in U} X_k} = \frac{n}{T_X}$

Toutefois, il n'y a aucune garantie que le $\pi_k \in [0, 1]$, il sera fréquent que

certaines probabilités de inclusion sont > 1 . Pour que de telles situations, on sélectionnera d'office les unités correspondantes, c-à-d, $\pi_k = 1$, et l'on recommencera la procédure avec les unités restantes en prenant soin de diminuer la taille de l'echantillon n dans \textcircled{a} .

Ex 2: Considérons la population $U = \{1, 2, 3, 4, 5, 6\}$ avec une variable auxiliaire x telle que $x_1 = 1, x_2 = 9, x_3 = 10, x_4 = 70, x_5 = 90, x_6 = 120$.

On a donc $tx = 360$. Si l'on souhaite obtenir un plan de taille fixe $n = 3$ alors les probabilités d'inclusion temporaire

$$\pi_1 = \frac{3 \times 1}{360} = 0,01$$

$$\pi_4 = \frac{3 \times 70}{360} = 0,7$$

$$\pi_2 = \frac{3 \times 9}{360} = 0,03$$

$$\pi_5 = \frac{3 \times 90}{360} = 0,9$$

$$\pi_3 = \frac{3 \times 10}{360} = 0,01$$

$$\pi_6 = \frac{3 \times 120}{360} = 1,2 > 1$$

L'unité 6 est alors sélectionnée d'office. Le total sur la 6^{ème} unité est $\sum_{k \in U} x_k = tx = 360 = 120$ et les probabilités de inclusion $n = 3 - 1 = 2$ deviennent:

$$\pi_1 = \frac{2 \times 1}{120} = 0,011$$

$$\pi_3 = \frac{2 \times 10}{120} =$$

$$\pi_2 = \frac{2 \times 9}{120} =$$

$$\pi_4 = \frac{2 \times 70}{120} =$$

2e pl

$$\pi_5 = \frac{2 \times 90}{120} = 1$$

On arrête ici

$$\pi_1 = \frac{1}{90}, \pi_2 = \frac{1}{10}, \pi_3 = \frac{1}{9}, \pi_4 = \frac{7}{9}$$

$\pi_5 = \pi_6 = 1$ les unités 5 et 6 sont donc sélectionnées d'office et il restera à choisir une unité entre 1, 2, 3 et 4

$$\text{Notons que } \sum_{k=1}^6 \pi_k = \frac{1+9+10+70}{90} + 2 = 3$$

comme souhaité.

Rappelons qu'un plan de sondage est défini par les $P(S)$ et non par les π_k . Pour avoir un plan à probabilité égale, il faut donc définir un plan de sondage $P(S)$, telle que $\forall k \in U, \sum_{S \ni k} P(S) = \pi_k$

$$\mathcal{P}_n = \{S \subset U \text{ tel que } |S| = n\}$$

Remarque: Il existe une infinité de plans de sondage vérifiant ces conditions. Nous allons donc introduire quelques plans de sondage à probabilité égale à taille fixe n couramment utilisés.

I - PLAN DE POISSON

Pour mettre en œuvre le plan de sondage aléatoire de Poisson,

→ On considère n probabilités π_1, \dots, π_n
 → On genère grand N nombre X_1, \dots, X_n (indépendamment des uns des autres) suivant la loi uniforme $U([0, 1])$.

→ $V_i \in 1, \dots, N$, l'individu V_i s'il vérifie $U_i \leq \pi_i$

→ les individus sélectionnés constituent l'échantillon.

Le Plan de Poisson a de très bonne qualité mais également un gros défaut, la taille de l'échantillon ne pouvant être supérieure à N .

Puisque les unités sont sélectionnées indépendamment, on a :

$$\pi_{kl} = \Pr(k \in S, l \in S) = \Pr(k \in S) \Pr(l \in S) = \pi_k \pi_l$$

et donc $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l = 0 \quad \forall k \neq l$

Il s'ensuit

Le plan de sondage est donc donné

$$P(S) = \prod_{k \in S} \pi_k \times \prod_{k \notin S} (1 - \pi_k)$$

Probabilité de sélectionner les unités choisies
 Probabilité de ne pas sélectionner les unités non choisies

Puis que $\Delta_{kl} = 0, \forall k \neq l$, la variance du π -estimateur du total T_y est

$$\text{Var}(T_y) = \sum_{k \in U} \frac{y_k^2 \pi_k}{\pi_k^2} \Delta_{kk} = \sum_{k \in U} \frac{y_k^2}{\pi_k} \pi_k (1 - \pi_k)$$

$$= \sum_{k \in U} \frac{y_k^2 (1 - \pi_k)}{\pi_k} \text{ ne peut qu'être estimé}$$

par $\hat{\text{Var}}(T_y) = \sum_{k \in U} \frac{y_k^2 (1 - \pi_k)}{\pi_k^2} \pi_k$

Obtenir le plan de Poisson est intéressant car il est simple.

Nous introduisons maintenant une mesure du désordre

Définition :

On appelle entropie d'un plan $P(\cdot)$ la quantité $H(P) = - \sum_{S \subset U} P(S) \log P(S)$ avec la convention que $0 \log 0 = 0$

L'entropie est toujours positive. De plus, comme mesure du désordre, plus $H(P)$ sera grand, plus le plan $P(\cdot)$ sera aléatoire.

Pour des probabilités d'inclusion fixes, on cherchera donc le plan le plus aléatoire donné, il est celui maximisant l'entropie

différence : $\sum_{S \subset U, k \in U} \pi_k = \pi_k (1 + \sum_{S \subset U, k \notin U} 1)$

B - Sondage systématique à probabilité inégale.

Ce plan de sondage a été introduit vers 1950 et est toujours largement utilisé puisqu'il a le mérite d'être simple et exact. Contrairement au plan de Poisson, il est également de taille fixe.

On désire tirer des échantillons dont les probabilités d'inclusion d'ordre 1 sont fixées à priori et telle que $0 < \pi_i < 1$

$$k \in U \text{ et } \sum_{k \in U} \pi_k = n.$$

Définissons les probabilités d'inclusions cumulées $G_k = \sum_{l=1}^k \pi_l$, $k \in U$, $G_0 = 0$.

Considérons l'approche à generer u et de sélection des unités à partir de cette unique réalisation. La première unité sélectionnée notée k_1 sera telle que $G_{k_1-1} \leq u \leq G_{k_1}$.

La deuxième unité sélectionnée notée k_2 sera cette fois ci telle que $G_{k_2-1} \leq 1+u \leq G_{k_2}$.

Et ainsi de suite. De manière générale, la j -ième unité sélectionnée notée k_j sera alors $G_{k_j-1} \leq j-1+u \leq G_{k_j}$.

Exemple : Prenons la situation où $N=6$
 $n=3$; $\pi_1 = 0,2$; $\pi_2 = 0,7$; $\pi_3 = 0,8$; $\pi_4 = 0,5$
 $\pi_5 = \pi_6 = 0,4$.
 Déterminer l'échantillon sélectionné sachant que $u = 0,3658$.

$$\begin{aligned} G_0 &= 0 & G_0 \leq u < G_1 \\ G_1 &= \pi_1 = 0,2 & \text{1}^\circ \quad 0 \leq 0,3658 < 0,2 \quad \times \\ & & \text{2}^\circ \quad 0,2 \leq 0,3658 < 0,9 & \checkmark \\ & & \text{3}^\circ \quad 0,9 \leq 1+0,3658 < 1,7 & \checkmark \\ & & \text{4}^\circ \quad 1,7 \leq 2,3658 < 2,2 & \times \end{aligned}$$

$$\begin{aligned}
 1,3 &\leq 2,3658 < 2,2 \quad \times \\
 1,3 &\leq 2,3658 < 2,2 \quad \times \\
 5 &: 6,9 \leq 2,3658 < 2,6 \\
 2,2 &\leq 2,3658 < 2,6 \quad \checkmark \\
 6 &: 2,6 \leq 2,3658 < 3 \quad \times
 \end{aligned}$$

Mardi, 06 Novembre 2024

Chapitre 04 : SONDAGE SIMPLIFIÉ

I - PRINCIPES ET JUSTIFICATIONS

Dans un sondage aléatoire simple, tout les échantillons de population de taille N sont possibles avec la même probabilité. On imagine que certains d'entre eux puissent savoir la a priori indésirable. Prenons par exemple le cas où nous disposons de cinq (5) 'jetons' : 1, 2, 4, 10 et 20 dont nous souhaitons évaluer la moyenne ($\mu = 7$) à l'aide d'un échantillon de taille 2. Parmi les échantillons de deux (2) unités, on trouve des cas extrêmes $\{1, 2\}$ et $\{10, 20\}$, qui sont particulièrement mauvais.

Plus concrètement, dans l'étude du lancement d'un nouveau produit financier, on peut supposer des différences de comportement entre les petits et les gros clients de la banque. Il serait malheureux que les hasards de l'échantillon conduisent à n'interroger que les clients appartenant à une seule de ces catégories, ou simplement que l'échantillon soit trop déséquilibré en faveur de l'une d'elles. S'il existe dans

de sondage une information auxiliaire qui permet de distinguer, à priori, les catégories de petits et gros clients, on peut tout à fait utiliser cette information pour répartir l'échantillon dans chaque sous-population. C'est la principale de la stratification : découper la population en sous-ensembles appelés strates et réaliser un sondage dans chacune d'elles.

L'intérêt de cette méthode, en comparaison des plans simples, est qu'elle permet d'augmenter la précision des estimateurs. Elle nécessite l'utilisation d'une information auxiliaire connue pour l'ensemble de la population.

Exemple : Nous souhaitons estimer l'âge moyen de toutes les personnes évoluant sur le site UJI. La base de sondage est composée de l'ensemble des personnes de l'UJI. Supposons que nous disposons de la répartition des éléments de la base suivant les catégories :

- Étudiants
- Enseignants
- Personnel administratif.

Autrement dit, nous connaissons la répartition des personnes de l'UJI suivant ces trois (ou) catégories. Il y'a

fort ap à penser que la variable âge ne se comporte pas de la même manière dans les trois classes (en moyenne, on peut penser en effet que la population enseignante ou personnel administrative est plus âgée que la population étudiante). Il paraît dès lors pertinent que c de prendre en compte cette information dans le plan de sondage.

La répartition des personnes de l'UJI fournit une information auxiliaire en outre problématique.

L'objectif principal consiste donc à mettre à profit cette information pour obtenir des résultats précis. L'information auxiliaire peut être utilisée à deux moments :

- À l'étape de la conception du plan de sondage.
- À l'étape de l'estimation des paramètres.

Dans ce chapitre, nous utiliserons cette information pour bâtir le plan de sondage.

II - POPULATION ET STRATES

Supposons que la population U soit partitionnée en H sous-ensembles U_1, \dots, U_H ap-

insérer petite strate et telle que l'union des U_h $U = \bigcup_{h=1}^H U_h$, $U_h \cap U_k = \emptyset$ $\forall h \neq k$

Chaque strate U_h admet une taille N_h et l'on a $\sum_{h=1}^H N_h = N$ où N est la taille de la population U .

Remarque :

Les tailles de strates N_h sont ici supposées connues et constituent l'information auxiliaire.

Notre but étant toujours d'estimer un total ou une moyenne, remarquons que le total (respectivement la moyenne) s'écrit à l'aide des strates $t_y = \sum_{k \in U} y_k = \sum_{h=1}^H \sum_{k \in U_h} y_k$
 $= \sum_{h=1}^H t_{y,h}$

Où $t_{y,h}$ est le total des valeurs prises par le caractère y sur la strate U_h
 c-à-d $t_{y,h} = \sum_{k \in U_h} y_k$

De même la moyenne sur la population s'écrit $\bar{y} = \frac{1}{N} \sum_{k \in U} y_k$
 $= \frac{1}{N} \sum_{h=1}^H \sum_{k \in U_h} y_k = \frac{1}{N} \sum_{h=1}^H N_h \bar{y}_{y,h}$

où $\bar{y}_{y,h}$ est la moyenne des valeurs prises par le caractère y sur la strate

$$U_h \text{ c-à-d } \sigma_{y,h}^2 = \frac{1}{N_h} \sum_{k \in U_h} (y_k - \bar{y}_{y,h})^2$$

$$\bar{y}_{y,h} = \frac{1}{N_h} \sum_{k \in U_h} y_k$$

$$\sigma_{y,h}^2 = \frac{1}{N_h - 1} \sum_{k \in U_h} (y_k - \bar{y}_{y,h})^2$$

Remarque :

La variance sur la population (total), σ_y^2 s'écrit :

$$\sigma_y^2 = \frac{1}{N} \sum_{k \in U} y_k^2$$

$$= \frac{1}{N} \sum_{h=1}^H \sum_{k \in U_h} [(y_k - \bar{y}_{y,h}) + (\bar{y}_{y,h} - \bar{y})]^2$$

$$= \frac{1}{N} \sum_{h=1}^H \left\{ \sum_{k \in U_h} (y_k - \bar{y}_{y,h})^2 + 2(\bar{y}_{y,h} - \bar{y}) \sum_{k \in U_h} (y_k - \bar{y}_{y,h}) + N_h (\bar{y}_{y,h} - \bar{y})^2 \right\}$$

$$= \frac{1}{N} \sum_{h=1}^H N_h \sigma_{y,h}^2 + \frac{1}{N} \sum_{h=1}^H N_h (\bar{y}_{y,h} - \bar{y})^2$$

$$= \sigma_y^2 \text{ intra} + \sigma_y^2 \text{ inter}$$

où $\sigma_y^2 \text{ intra}$ est la variance intra strate (moyenne des variances des strates) et $\sigma_y^2 \text{ inter}$ la variance inter strate (due aux différences entre strates).

13-11-2022

III - ECHANTILLONNAGE, PROBABILITE D'INCLUSION ET INFORMATION

Un sondage est dit stratifié, si, pour chaque strate, on tire un échantillon. Selon le sondage électronique simple sans remise de taille fixe n_h et que les tirages au sein de chaque strate sont mutuellement indépendants.

Soit S_h l'échantillon aléatoire tiré dans la strate U_h à l'aide d'un plan de sondage $P_h(\cdot)$. L'échantillon aléatoire S obtenu au final est donc $S = \bigcup_{h=1}^H S_h$.

Le plan de sondage associé π n'est rien d'autre que $P(S) = \prod_{h=1}^H P_h(S_h)$, $S = \bigcup_{h=1}^H S_h$ et la taille de l'échantillon S est $n = \sum_{h=1}^H n_h$.

Le calcul des probabilités d'inclusion pour un sondage stratifié n'est pas difficile mais il faut tout de même faire attention. Pour les probabilités d'inclusion d'ordre 1 et si l'unité $k \in U_h$ alors $\pi_k = \frac{n_h}{N_h}$ puisqu'on a effectué un plan simple sans remise de taille n_h pour cette strate.

Pour les probabilités d'inclusion d'ordre 2, c'est un peu plus difficile et

le résultat dépend du fait que les unités $k, l \in U_h$ à la même strate ou non.

- si $k, l \in U_h$ à la même strate alors $\pi_{kl} = \frac{n_h(n_h-1)}{N_h(N_h-1)}$

- si $k, l \in U_h$ à deux strates différentes U_{h_1} et U_{h_2} alors l'indépendance entre les strates $\pi_{kl} = \pi_k \pi_l = \frac{n_{h_1} n_{h_2}}{N_{h_1} N_{h_2}}$

En conséquence on a :

$$\Delta_{kl} = \begin{cases} \frac{n_{h_1}}{N_{h_1}} \left(1 - \frac{n_{h_1}}{N_{h_1}}\right), & k = l, k \in U_{h_1} \\ -\frac{n_{h_1}}{N_{h_1}} \frac{(N_{h_1} - n_{h_1})}{(N_{h_1} - 1)}, & k \neq l, k, l \in U_{h_1} \\ 0, & k \in U_{h_1}, l \notin U_{h_1} \end{cases}$$

Soit donc, les T-estimateurs du total T_y et la moyenne μ_y sont :

$$\hat{T}_{y, \text{strat}} = \sum_{k \in S} \frac{y_k}{\pi_k}$$

$$= \sum_{h=1}^H \sum_{k \in S_h} \frac{n_h y_k}{n_h}$$

$$= \sum_{h=1}^H \frac{1}{n_h} \sum_{k \in S_h} y_k$$

$$\hat{\mu}_{y, \text{strat}} = \frac{1}{N} \sum_{h=1}^H \frac{N_h}{n_h} \sum_{k \in S_h} y_k$$

$$= \frac{1}{N} \sum_{h=1}^H N_h \bar{y}_h$$

de \bar{y}_h est l'estimateur du total pour la strate h .
 c-à-d $\hat{t}_h = \frac{N_h}{n_h} \sum_{k=1}^{n_h} y_k$

et \bar{y}_h est la moyenne de l'échantillon prise à la strate h c-à-d
 $\bar{y}_h = \frac{1}{n_h} \sum_{k=1}^{n_h} y_k$

Exemple: les résultats du sondage sont donnés dans le tableau suivant.

Strate	1	2	3	4	5	6	7	8	9	10
Age	25	30	35	40	45	50	55	60	65	70

1- Calculer la moyenne des âges des individus de l'échantillon par strate.

2- Déterminer la moyenne μ , calculer la variance de la moyenne, calculer la variance du total et calculer la variance pour $\hat{\mu}$.

Solution

1- Calcul de la moyenne des âges de l'échantillon par strate.

$$\bar{y}_h = \frac{1}{n_h} \sum_{k=1}^{n_h} y_k$$

Pour $h=1$, $\bar{y}_1 = \frac{1}{n_1} \sum_{k=1}^{n_1} y_k$ avec n_1

$n_1=5$, $\bar{y}_1 = \frac{1}{5} (20 + 25 + 23 + 22 + 26) = 23,2$

Pour $h=2$, $\bar{y}_2 = \frac{1}{n_2} \sum_{k=1}^{n_2} y_k$

$\bar{y}_2 = \frac{1}{3} (30 + 35 + 38) = 34,3$

Pour $h=3$, $\bar{y}_3 = \frac{1}{n_3} \sum_{k=1}^{n_3} y_k$

$\bar{y}_3 = \frac{1}{2} (42 + 44) = 43$

2- Calcul de la moyenne μ

$$\mu_{\text{strate}} = \frac{1}{N} \sum_{h=1}^H N_h \bar{y}_h$$

$\Rightarrow \mu = \frac{1}{10} (5 \times 23,2 + 3 \times 34,3 + 2 \times 43) = 30,5$

$$\text{Var}(\hat{\mu}) = \frac{1}{N^2} \sum_{h=1}^H N_h \frac{N_h - n_h}{n_h} \sigma_h^2$$

$$\text{Var}(\hat{\mu}) = \sum_{h=1}^H \text{Var}(\hat{y}_h)$$

$$\text{Var}(\hat{y}_h) = \sum_{k=1}^{n_h} N_h (N_h - n_h) \frac{y_k^2}{n_h}$$

$$\hat{y}_h = \frac{1}{n_h - 1} \sum_{k=1}^{n_h} (y_k - \bar{y}_h)^2, h=1, 2, \dots, H$$

IV - PLAN STRATIFIÉ ET ALLOCATION PROPORTIONNELLE

1 - Définition

Un plan stratifié est dit à allocation proportionnelle si $\frac{n_h}{N_h} = \frac{n}{N} \quad \forall h = 1, \dots, H$

C'est à dire que, les strates de tailles importantes ontent pour plus d'unités dans l'échantillon que celles de tailles plus petites.

Remarque: Généralement, la taille de l'échantillon par chaque strate $n_h = \frac{n N_h}{N}$ ne sera pas entière. Mais afin de simplifier les développements théoriques qui viennent, nous allons tout de même le supposer.

les E.T. - estimateurs du total et de la moyenne sont alors :

$$\begin{aligned} \hat{t}_{y, \text{stat prop}} &= \sum_{h=1}^H \hat{t}_{y, h} \\ &= \frac{N}{n} \sum_{h=1}^H \hat{t}_{y, h} \end{aligned}$$

$$\hat{\bar{y}}_{y, \text{stat prop}} = \frac{1}{n} \sum_{h=1}^H \hat{t}_{y, h}$$

la variance de la moyenne est alors :

$$\begin{aligned} \text{Var}(\hat{y}_{y, \text{stat prop}}) &= \frac{1}{n^2} \left(1 - \frac{n}{N}\right) \frac{1}{n} \sum_{h=1}^H N_h S_{y, h}^2 \\ \text{Var}(\hat{t}_{y, \text{stat prop}}) &= \sum_{h=1}^H N_h (N_h - n) \frac{S_{y, h}^2}{n} \\ &= \sum_{h=1}^H N_h \left(\frac{N}{n} - 1\right) S_{y, h}^2 \\ &= \frac{N-n}{n} \sum_{h=1}^H N_h S_{y, h}^2 \end{aligned}$$

Remarque: lorsque les tailles de strates N_h sont suffisamment grande, alors $S_{y, h}^2 \approx \sigma_{y, h}^2$ et donc: $\text{Var}(\hat{t}_{y, \text{stat prop}}) \approx N-n$

$$\begin{aligned} \text{Var}(\hat{t}_{y, \text{stat prop}}) &= \frac{N-n}{n} \sum_{h=1}^H N_h \sigma_{y, h}^2 \\ &= N(N-n) \frac{\sigma_y^2 \text{ intra}}{n} \end{aligned}$$

Alors que la variance de l'estimateur du total par un plan simple sans remise vaut: $\text{Var}(\hat{t}_{y, \pi}) \approx N(N-n) \frac{\sigma_y^2}{n}$. les deux expressions sont quasiment identiques, mais plus que $\sigma_y^2 = \sigma_y^2 \text{ intra} + \sigma_y^2 \text{ inter}$, la 1^{re} expression est plus petite c-à-d que l'on obtient de meilleurs résultats avec un plan stratifié avec l'allocation proportionnelle qu'avec un plan simple sans remise.

Ceci sera d'autant plus vrai que la variance inter strates sera grande, car est le

Les caractéristiques d'intérêt dépendent fortement du caractère servant à la classification, les tailles Nh .

On minimise sans biais cette variance pour $Var(hg, n) \leq N(N-n) \frac{S_y^2}{n}$

$$Var(hg, n) = \frac{N-n}{n} \sum_{h=1}^H N_h^2 \frac{S_y^2}{N_h}$$

$$S_y^2 = \frac{1}{N-1} \sum_{h=1}^H (J_h - \bar{J}_h)^2 \quad J_h = J_1, J_2, \dots, J_H$$

Ex. 2

On donne dans le tableau suivant pour chaque faculté de l'université de Guel :

- Répartition par âge
- Les catégories : 1 = étudiants, 2 = enseignants
- 3 = personnel administratif
- Les valeurs de cheveux : a = brun ; b = blond ; c = châtain

Pour simplifier, on considère une population de 20 individus.

Age	Catégorie	cheveux
24	1	C
52	2	A
42	3	B
19	1	C
33	3	A
26	1	B
45	2	C
23	1	A
39	2	A
24	1	B

Age	Catégorie	cheveux
22	1	C
48	2	A
24	1	A
38	3	A
26	1	B
36	3	B
46	2	B
23	1	C
39	2	A
18	1	C

$n = 10$

1- On veut estimer la moyenne d'âge μ à l'aide d'un plan simple. Quelle est la variance de l'estimateur ?

2- On désire stratifier la population suivant la catégorie. Quelle est la variance de l'estimateur μ pour un tel plan ?

3- On choisit maintenant de stratifier suivant la couleur des cheveux. Quelle est la variance de l'estimateur pour un tel plan.

Solution

20-11-2021

$$1 - Var(\mu) = \left(1 - \frac{n}{N}\right) = \frac{52}{n}$$

$$= \left(1 - \frac{10}{20}\right) = \frac{10}{20}$$

$$= 5,77$$

$$2 - S_{y,h}^2 = \frac{1}{N_h - 1} \sum_{k \in U_h} (y_k - \bar{y}_{h,h})^2$$

$$\begin{array}{ll} N_1 = 22,7 & S_{y,1}^2 = 6,99 \\ N_2 = 41,93 & S_{y,2}^2 = 26,17 \\ N_3 = 35,5 & S_{y,3}^2 = 6,33 \end{array}$$

* On cherche la variance de l'estimateur

$$Var(\hat{\mu}) = \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N} \sum_{h=1}^H N_h S_{y,h}^2$$

$$= \frac{1}{10} \left(1 - \frac{10}{20}\right) \frac{1}{20} [10 \times 6,99 + 6 \times 26,17 + 4 \times 6,33]$$

$$= 0,13$$

3 - Variance de l'estimateur est de 4,86.

V - PLAN OPTIMAL DU LI TOTALE

Si votre intérêt est d'estimer un total ou une moyenne, alors il existe une taille optimale pour les strates. On cherche donc les tailles d'échantillon n_1, \dots, n_H minimisant la variance de l'estimateur du total y_T pour une

taille d'échantillon fixée n , c-à-d minimiser la variance.

$$Var(y_{T, optimal}) = \sum_{h=1}^H N_h (N_h - n_h) \frac{S_{y,h}^2}{n_h}$$

par rapport à n_h et sous contrainte de

$$\sum_{h=1}^H n_h = n$$

Le Lagrangien de ce problème de minimisation est

$$L(n_1, \dots, n_H, \lambda) = \sum_{h=1}^H N_h (N_h - n_h) \frac{S_{y,h}^2}{n_h} + \lambda \left(\sum_{h=1}^H n_h - n \right)$$

On a donc

$$\frac{\partial L}{\partial n_h} = 0 \Rightarrow -\frac{N_h^2}{n_h^2} S_{y,h}^2 + \lambda = 0$$

$$\Rightarrow n_h = \frac{N_h S_{y,h}}{\sqrt{\lambda}}$$

Mais plus que $\sum n_h = n$, on a donc

$$\lambda \frac{1}{2} N_h S_{y,h} = n \quad \lambda \frac{1}{2} \sum_{h=1}^H N_h S_{y,h} = n$$

$$\text{et il vient : } n_h = n \frac{N_h S_{y,h}}{\sum_{j=1}^H N_j S_{y,j}}$$

$$h = 1, \dots, H$$

Remarque: la taille optimale pour une strate h est donc proportionnelle au produit de la taille de cette strate et de l'écart type du caractère y de cette strate.

Bien entendu en pratique on ne connaît pas $S_{y,h}$ et donc la formule précédente n'est pas d'un grand intérêt. Elle est cependant assez instructive et intuitive. Instructive puisqu'elle indique qu'il faut surreprésenter les strates qui ont une forte variabilité.

Remarque: En pratique les tailles $n_h \notin \mathbb{N}$ et on arrondira les résultats. De plus il peut arriver également que $n_h > N_h$ pour au un $h \in 1, \dots, H$. Pour de tel cas on pose alors $n_h = N_h$ et on déterminera les tailles optimales sur les strates restantes en procédant par itération si nécessaire.

Supposons que nos tailles optimales soient des entiers et telles que $n_h \leq N_h \forall h$. Alors la variance du \bar{Y} estimée est donc :

$$\begin{aligned} \text{Var}(\hat{y}_{\text{strat}}) &= \sum_{h=1}^H N_h (N_h - n_h) \frac{S_{y,h}^2}{n_h} \\ &= \sum_{h=1}^H N_h^2 \sum_{h=1}^H \frac{N_h S_{y,h}^2}{N_h N_h S_{y,h}^2} S_{y,h}^2 - \sum_{h=1}^H N_h S_{y,h}^2 \\ &= \left(\sum_{h=1}^H \frac{N_h S_{y,h}^2}{n_h} \right) \sum_{h=1}^H N_h S_{y,h}^2 - \sum_{h=1}^H N_h S_{y,h}^2 \\ &= \frac{1}{n} \left(\sum_{h=1}^H N_h S_{y,h}^2 \right)^2 - \sum_{h=1}^H N_h S_{y,h}^2 \end{aligned}$$