

Théorie des sondages

Semestre 6

Table des matières

I	1
1	1
2 Les mesures des erreurs d'échantillonnage	2
2.1 Le biais	2
2.2 La variance	2
2.3 L'erreur quadratique moyenne (EQM) (MSE)	3
3 Les bases de sondage	5
3.1 Les propriétés des bases de sondages	5
3.2 Absence de base de sondage	5
3.3 Les différents types de bases	6
4 Les différents types d'erreurs rencontrées dans les enquêtes par sondage	6
4.1 L'erreur d'échantillonnage	6
4.2 Les erreurs d'observation	6
4.3 Le défaut de couverture et les non réponses	7
II Le sondage aléatoire simple	7

Première partie

1

Une fois l'échantillon sélectionné on dispose de l'information suivante : $Y_{i1}, Y_{i2}, \dots, Y_{in}$. Cherchons toujours θ , il convient de combiner ces petites n valeurs recueillies sur l'échantillon pour obtenir une expression dont la valeur numérique est proche de celle de θ . La formule qui agrège les petites n valeurs s'appelle l'estimateur de θ , on le note : $\hat{\theta}$ il s'agit d'une fonction calculée à partir des données de l'échantillon, et la procédure qui permet de sauter le pas, c'est à dire de passer des données recueillies sur les échantillons à la vraie valeur

inconnue dans la population s'appelle l'inférence statistique (cœur même de la théorie des sondages).

2 Les mesures des erreurs d'échantillonnage

Une fois que la fonction qui agrège $\hat{\theta} = g(Y_{i1}, Y_{i2}, \dots, Y_{in})$ a été choisie il est nécessaire d'évaluer sa pertinence en d'autres termes il convient de répondre à la question suivante : Sommes nous proche de la valeur θ lorsque l'on calcule g à partir de l'échantillon sélectionné ?

En fait il n'est pas possible de répondre précisément à cette question pour la simple et bonne raison que l'on ne connaît pas θ , mais la statistique permet d'apporter des éléments de réponse en exploitant l'aspect probabiliste des choses. Un des points fondamentaux qu'il faut à présent bien comprendre, réside dans la nature de l'aléa introduit dans notre problème, l'aléa se situe exclusivement au niveau des identifiants des individus de l'échantillon.

Ce qui est aléatoire sont les $i1, i2, \dots, in$ et non les $Y_{i1}, Y_{i2}, \dots, Y_{in}$, ainsi si l'on réalise une enquête sur les revenus, le sondeur considère dans son enquête que chaque individu dans la population est capable de fournir son revenu, en revanche l'échantillon s d'individus interrogés sur leur revenu aura une composition aléatoire. Notre estimateur g (ou $\hat{\theta}$) est donc aléatoire fonction de l'échantillon s . La réponse à notre question initiale fait appel à trois notions : le biais, la variance, et l'erreur quadratique moyenne.

2.1 Le biais

Le biais permet de détecter la présence éventuelle d'erreur systématique, il correspond donc à la différence entre l'espérance mathématique de l'estimateur et le paramètre lui-même.

$$\text{Biais}(\hat{\theta}) = E(\hat{\theta}) - \theta$$

Le biais constitue une première mesure d'erreur d'échantillonnage.

2.2 La variance

On présente bien au travers de l'exemple précédent que la notion de moyenne ne suffit pas à mesurer la qualité d'un échantillonnage et qu'il faut une autre grandeur d'avantage liée à la dispersion des valeurs de $\hat{\theta}$. On calcule la variance des estimateurs $\hat{\theta}_s$ lorsque l'aléa est l'échantillon, s .

$$\begin{aligned} V(\hat{\theta}) &= E \left[\left(\hat{\theta} - E(\hat{\theta}) \right)^2 \right] \\ &= \sum_{i=1}^s P(i) \left[\left(\hat{\theta} - E(\hat{\theta}) \right)^2 \right] \end{aligned}$$

La présence du terme quadratique fait toute la différence avec le biais puisque les écarts positifs et négatifs ne se compensent plus. En terme de sondage $\sigma_{\hat{\theta}}$ et $V(\hat{\theta})$ mesurent la précision et constituent après le biais une seconde mesure de l'erreur d'échantillonnage, plus ils sont grands, moins le plan de sondage est bon.

Si $\sigma_{\hat{\theta}}$ et $V(\hat{\theta})$ sont grands il faut agir sur l'expression de $\hat{\theta}$ soit sur les probabilités de tirage ($P(s)$). Le meilleur estimateur $\hat{\theta}$ est celui qui a la plus petite variance compte tenu du budget dont on dispose.

Pour terminer, il faut insister une fois encore sur le fait que ni l'espérance, ni la variance ne peuvent renseigner sur l'écart exact entre la valeur de $\hat{\theta}$ obtenue et la vraie valeur de θ inconnue.

2.3 L'erreur quadratique moyenne (EQM) (MSE)

On peut également construire un indicateur de précision qui englobe les notions de biais et de variance pour se faire, il suffit de calculer l'erreur quadratique moyenne tel que :

$$\text{EQM} = E \left[\hat{\theta} - \theta \right]^2$$

Cette erreur quadratique moyenne est un indicateur synthétique de qualité permettant de répondre à la question complexe suivante : est-il préférable d'avoir un biais fort et une variance faible ou un biais faible et une variance forte ?

Il n'y a évidemment pas de bonne réponse à cette question, mais la pratique montre que l'on cherche en général à éviter en priorité les biais forts.

Pour conclure, l'erreur d'échantillonnage est donc le fait que les résultats numériques publiés à la suite d'un sondage dépendent des individus qui composent l'échantillon. Elle est présente, et même une caractéristique des enquêtes par sondage. Cette erreur d'échantillonnage se mesure par le biais, la variance ou l'erreur quadratique moyenne, elle ne peut pas être évitée.

Exemple 1. Considérons une population composée de 4 entreprises $N = 4$ et l'on s'intéresse au chiffre d'affaire mensuel moyen de cette population \overline{CA} mensuel. Le CA mensuel des entreprises est le suivant :

$$CA_1 = 6000\text{€} \quad CA_2 = 12000\text{€} \quad CA_3 = 8000\text{€} \quad CA_4 = 6000\text{€}$$

Supposons que pour raison de contrainte de budget, je ne peux interroger que 2 entreprises parmi les 4, $n = 2$.

$$\theta = \frac{CA_1 + CA_2 + CA_3 + CA_4}{N} = 8000$$

Nombre d'échantillons possibles :

$$C_4^2 = \frac{4!}{2!(4-2)!} = 6$$

Ces 6 échantillons sont les suivants :

$s_1 = (1, 2)$	$P(s_1) = 0.25$
$s_2 = (1, 3)$	$P(s_2) = 0.25$
$s_3 = (1, 4)$	$P(s_3) = 0.2$
$s_4 = (2, 3)$	$P(s_4) = 0.1$
$s_5 = (2, 4)$	$P(s_5) = 0.1$
$s_6 = (3, 4)$	$P(s_6) = 0.1$

Parce que l'on juge que l'entreprise 1 est particulièrement coopérative sur ce sujet, on veut lui donner une probabilité de tirage supérieure aux trois autres, si bien que les trois échantillons s_1, s_2, s_3 sont un peu plus probables que les autres.

Et par soucis de simplicité on choisit comme estimateur la moyenne simple dans l'échantillon. Autrement dit, on suppose que : $\hat{\overline{CA}} = \bar{y}$ La moyenne simple calculée pour les entreprises des échantillons.

s_1	$\overline{CA}(s_1) = \frac{6000 + 12000}{2} = 9000$
s_2	$\overline{CA}(s_2) = \frac{6000 + 8000}{2} = 7000$
s_3	$\overline{CA}(s_3) = \frac{6000 + 6000}{2} = 6000$
s_4	$\overline{CA}(s_4) = \frac{12000 + 8000}{2} = 10000$
s_5	$\overline{CA}(s_5) = \frac{12000 + 6000}{2} = 9000$
s_6	$\overline{CA}(s_6) = \frac{8000 + 6000}{2} = 7000$

Calculons le biais : $\theta = \overline{CA} = 8000$

$$\begin{aligned} E(\theta) &= \sum_{s=1}^6 P(s) \cdot \hat{\theta}(s) \\ &= 0.25 * 9000 + 0.25 * 7000 + 0.2 * 6000 + 0.1 * 10000 + 0.1 * 9000 + 0.1 * 7000 = 7800 \end{aligned}$$

En moyenne l'estimateur est de 7800, on peut donc calculer le biais :

$$\text{Biais} = 7800 - 8000 = -200$$

Variance :

$$0.25(9000 - 7800)^2 = 360000$$

3 Les bases de sondage

3.1 Les propriétés des bases de sondages

Pour pouvoir réaliser un tirage probabiliste, c'est à dire un tirage pour lequel par définition chaque individu de la population a une probabilité connue et fixée à l'avance de faire partie de l'échantillon à enquêter, il est absolument nécessaire de disposer d'une liste de toutes les unités d'échantillonnage faisant partie du champ de l'enquête, cette liste est appelée base de sondage, doit avoir trois qualités essentielles.

1. Elle doit permettre de repérer l'unité sans aucune ambiguïté, une bonne base de sondage est une bonne base d'identifiant.
2. Elle doit être exhaustive, cela signifie que chaque unité faisant partie du champ de l'enquête doit être nécessairement présente dans la liste des identifiants. Si ce n'est pas le cas, on parle de base de sondage incomplète ou de défaut de couverture.
3. Elle doit être sans double compte, aucun individu ne doit être présent deux fois dans la base de sondage, même sous deux identifiants différents.

De façon générale il est extrêmement difficile en pratique de s'affranchir du manque d'exhaustivité (la plupart des bases de sondage est incomplète) et de la présence de double compte. L'important est de juger de leur impact et de ne conserver que les bases faiblement imparfaites. A ces trois propriétés on en rajoute souvent une quatrième qui sans être indispensable, peut s'avérer très bénéfique, c'est la suivante : une information auxiliaire de bonne qualité (information supplémentaire dont on va pouvoir se servir pour améliorer soit la qualité des estimateurs ou la méthode de tirage). En fait on appelle information auxiliaire, toute variable quantitative ou qualitative autre que la variable d'intérêt Y et autre que les variables nécessaires et suffisantes à l'identification des individus de la population.

3.2 Absence de base de sondage

Lorsque l'on s'intéresse à une population, il est tout à fait possible que on ne puisse pas trouver de base de ce sondage reproduisant cette population ou qu'une telle base existe mais soit jugée de mauvaise qualité, il peut également arriver que l'on renonce à l'utilisation de base de bonne qualité pour des raisons d'ordre pratique, parmi ces raisons :

- La base existe mais ne peut pas être donnée
- La base est trop volumineuse et le matériel informatique dont on dispose ne permet pas son traitement, cet obstacle ne doit pas être négligé, car on oublie trop souvent que si les échantillons sont de petite taille, les bases elles sont souvent très grande.

Si pour une raison quelconque on ne dispose pas de base de sondage acceptable permettant de réaliser un échantillonnage de bonne qualité, trois solutions sont envisageables :

1. Ne pas utiliser du tout de base de sondage, sortir du cadre probabiliste, réaliser un sondage empirique.
2. On peut rechercher des bases de sondage, non plus d'individus directement susceptibles de fournir l'information, mais accepter de passer par

un niveau intermédiaire de groupe d'individus, par un tirage à plusieurs degrés, et en réalisant des recensements intermédiaires dans les groupes sélectionnés, on peut échantillonner rigoureusement des unités de la population qui nous intéresse. Cette façon de faire conduit fréquemment sur ce que l'on appelle les sondages de type aréolaire, pour lesquels on échantillonne en premier lieu des aires géographiques.

3. Si le sujet auquel on s'intéresse s'y prête, on peut recueillir l'information en passant par une population intermédiaire, qui conduit aux unités d'observation, population intermédiaire d'une autre nature que celle des individus de la population sur laquelle portera l'inférence.

3.3 Les différents types de bases

Conformément au type d'enquêtes que l'on est amené à réaliser, on peut opter pour une base de liste, c'est à dire une base constituée par une liste d'identifiants d'individus, ou pour une base aréolaire, constituée par des aires géographiques, bien délimitées et bien identifiées, la base aréolaire permet d'effectuer des sondages en grappe, où une grappe est constituée par l'ensemble des individus de l'aire. On attribue trois avantages comparatifs principaux aux sondages aréolaires par rapport aux sondages réalisés à partir d'une base de liste.

Avantage des bases aréolaires

1. L'aire est une entité relativement stable qui permet de mieux prendre en compte l'évolution de la structure réelle de la population qu'elle contient.
2. Le regroupement géographique est par définition : maximal, ce qui permet de limiter les coûts de déplacement
3. Les taux de réponse sont en général meilleurs que dans les sondages par liste

4 Les différents types d'erreurs rencontrées dans les enquêtes par sondage

4.1 L'erreur d'échantillonnage

cf sec 1

4.2 Les erreurs d'observation

Dans la pratique, il existe une seconde famille d'erreurs, que l'on appelle les erreurs d'observation ou les erreurs de mesures, qui tiennent au fait que la valeur que l'on recueille lors de l'enquête peut être une valeur Y_i^* différente de la vraie valeur Y_i . Ce type d'erreur survient dans les questions sensibles. Sur de tels sujets, l'erreur d'observation est volontairement introduite par l'enquêté, mais il existe parallèlement à cela une longue liste de sources d'erreurs d'observations qui n'ont pas d'origine volontaire, parmi lesquelles on peut citer les erreurs de bonne foi de l'enquêté. Ensuite il y a les erreurs introduites par l'enquêteur qui interprète les questions et souffle potentiellement les réponses, il s'agit là d'un défaut de formation de l'enquêteur qui n'a pas à influencer la réponse de l'enquêté

4.3 Le défaut de couverture et les non réponses

On peut distinguer une troisième famille d'erreurs liée à l'existence d'une base de sondage incomplète et à la non réponse de certains individus aux questions posées. En effet, une base de sondage incomplète est une situation qui donne lieu à un défaut de couverture de la population, dans ce cas il y a dès l'origine un biais de l'estimateur que l'on ne peut pas mesurer. Il convient néanmoins de distinguer le défaut d'une base comprenant des doubles compte du défaut de couverture, dans le premier cas (a double compte), toute l'information est disponible dans le fichier, mais elle est mal utilisée alors que dans le défaut de couverture, l'information n'existe pas du tout. Concernant la non exhaustivité d'une base est un voisin de la non réponse dans la mesure où ces deux sources d'erreur sont dues à l'existence d'individus pour lequel on ne peut pas tirer d'information sur la valeur de Y . Néanmoins la non réponse se différencie sur deux points par rapport aux erreurs liées au défaut de couverture. L'ampleur du phénomène de non réponse est mesurable à partir de l'échantillon sélectionné alors que l'étendue d'un défaut de couverture ne l'est pas nécessairement. La taille de l'échantillon qui importe dans la mesure de l'échantillon avec une base incomplète et la taille totale de l'échantillon (taille qui peut être fixée à l'avance.) Alors que la taille de l'échantillon qui importe pour le calcul de la précision lorsqu'il y a non réponse est l'effectif des répondants, qui est aléatoire. Pour ce qui concerne la non réponse on peut avoir à faire en une réponse partielle ou une non réponse complète. Finalement l'erreur totale que l'on connaît

$$\begin{aligned} \text{Erreur totale} = & \boxed{\text{Erreur d'échantillonnage}} \\ & + \text{Erreur de mesure/observation} \\ & + \text{Erreurs dues au défaut de couverture} \\ & + \boxed{\text{Erreurs dues à la non réponse}} \end{aligned}$$

En général on a tendance naturellement à faire porter l'effort sur l'erreur d'échantillonnage et sur l'erreur de non réponse. En général les défauts de modélisation et de mesure portent sur l'erreur d'échantillonnage et sur la non réponse, malheureusement on ne sait que très peu de choses sur les erreurs d'observation et les défauts de couverture si ce n'est qu'ils sont de la nature de biais. Ce constat n'est donc pas très réjouissant car dans la mesure au contraire aux erreurs d'échantillonnage, ces erreurs ne diminuent pas avec la taille de l'échantillon.

Deuxième partie

Le sondage aléatoire simple

Il s'agit d'une méthode de tirage qui consiste à tirer dans la population de taille N un échantillon de taille fixée n sans remise à partir des seuls identifiants des individus de manière à ce que, chaque individu de la population ait la même probabilité d'inclusion et cela sans manipulation au préalable de la population. C'est exactement ce que l'on réalise lorsqu'on sélectionne dans une urne des boules sans remise. Le sondage aléatoire simple attribue à chaque échantillon s qui peut être formé, la même probabilité de sortie $p(s)$ qui est égale à l'inverse du

nombre d'échantillons distincts que l'on peut constituer dans la population, cette propriété remarquable le caractérise il ne nécessite pas non plus d'informations auxiliaires lors de sa mise en oeuvre.