

Section 2. Méthodes d'échantillonnage

Pierre Duchesne

September 6, 2019

- ▶ On dispose d'une base de sondage de taille finie N ;
- ▶ Notation: $U = \{1, 2, \dots, k, \dots, N\}$;
- ▶ Habituellement on va tenter d'utiliser l'indice k pour l'identificateur de l'unité k ;
- ▶ La variable d'intérêt sera notée y .
 1. Exemple, variable comme le salaire; Salaire $y_k = 50500\$$ indique le salaire de l'unité k ;
 2. Autre exemple, variable dichotomique: $y_k = 1$ si l'unité k est chômeur, et $y_k = 0$ sinon.

Paramètres de la population

- ▶ Ces quantités reposent sur la connaissance de l'ensemble des unités dans la population.
- ▶ Total: $t_y = \sum_{k=1}^N y_k$;
- ▶ Moyenne: $\bar{y}_U = \frac{1}{N} \sum_{k=1}^N y_k$;
- ▶ Variance: $S_{yU}^2 = \frac{1}{N-1} \sum_{k=1}^N (y_k - \bar{y}_U)^2$;
- ▶ D'autres exemples sont possibles; on pourrait penser à un paramètre comme la médiane, ou encore un quantile de population.

Convention pour le symbole \sum

- ▶ Soit A in sous-ensemble d'unités, $A \subseteq U$;
- ▶ Pour le total de la variable y sur A :

$$\sum_A y_k = \sum_{k \in A} y_k;$$

- ▶ Avec cette notation, on écrit:

$$t_y = \sum_U y_k;$$

et

$$\bar{y}_U = \frac{1}{N} \sum_U y_k;$$

- ▶ Un échantillon s est un sous-ensemble de U , on aura donc $s \subseteq U$.
- ▶ On note que U est de cardinalité finie, N .
- ▶ Il existe un nombre fini d'échantillons, mais le nombre possible peut en fait être très grand.
- ▶ On peut dénombrer le nombre total d'échantillons:

$$\binom{N}{0} + \binom{N}{1} + \binom{N}{2} + \dots + \binom{N}{N} = 2^N.$$

Plan d'échantillonnage

- ▶ Dans un sondage, on sélectionne un échantillon s et on observe la variable d'intérêt y pour les unités $k \in s$.
- ▶ On observe alors:

$$\{y_k; k \in s\}$$

- ▶ Échantillon probabiliste: c'est un échantillon obtenu selon un mécanisme aléatoire que l'on appelle **plan d'échantillonnage**.
- ▶ Plan d'échantillonnage: c'est un ensemble de règles strictes, qui une fois mise en oeuvre, nous donne le sous-ensemble de U , c'est-à-dire l'échantillon s .
- ▶ On requiert que la probabilité d'être inclus dans l'échantillon s est strictement positive, et ce $\forall k \in U$.

Échantillon représentatif

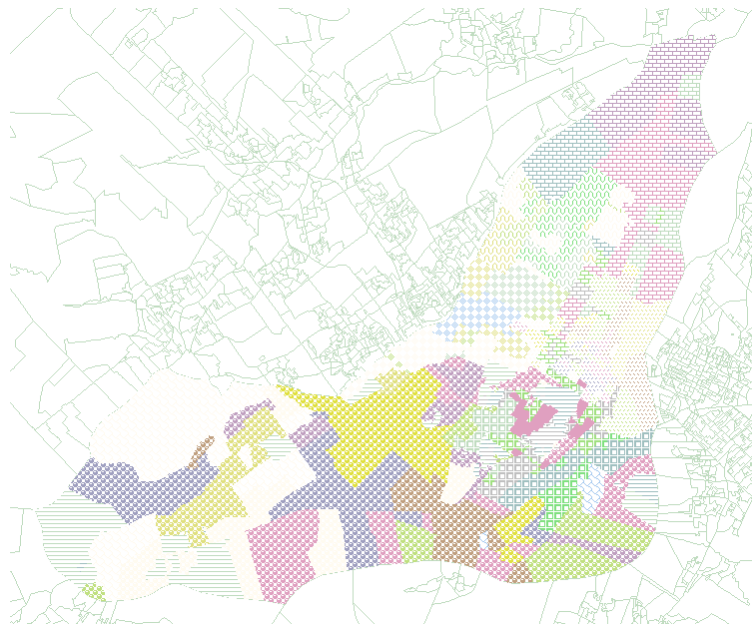
- ▶ Un échantillon a pour but de représenter la population, donc d'être représentatif.
- ▶ En quelque sorte, l'échantillon est un modèle de la population.
- ▶ Il n'est pas possible de déterminer si un échantillon est *représentatif* dans le sens qu'il est un *modèle réduit* de la population.
- ▶ Un bon plan d'échantillonnage nous assure que l'échantillon est obtenu selon une méthodologie que le méthodologiste peut tenir compte.
- ▶ Dans un sens, un plan d'échantillon performant contribue à éliminer des échantillons non représentatifs, donc la sens que les échantillons reflètent des caractéristiques utiles de la population.

Exemples de plan de sondage:

- ▶ Tirage aléatoire simple avec remise;
- ▶ Tirage aléatoire simple sans remise (chaque échantillon de taille n , n fournie, possède les mêmes chances de survenir);
- ▶ Tirage de Bernoulli; Tirage de Poisson;
- ▶ Tirage stratifié; tirage stratifié simple;
- ▶ Tirage systématique;
- ▶ Tirage en grappes;
- ▶ Tirage à plusieurs degrés.

Exemples de sondages probabilistes

- ▶ Dans une étude de marketing, on désire sonder les ménages de la ville de Montréal.
- ▶ On suppose que l'on dispose d'une liste de 1 à M des unités géographiques sur une carte; on parle sondages aréolaires.
- ▶ On sélectionne un échantillon au hasard d'unités géographiques.
- ▶ Dans chaque unité géographique, on sélectionne et observe tous ou une partie des ménages.



- ▶ Sondages de grappes et sondages à deux degrés.
- ▶ Sondages de grappes: On observe tous les ménages dans les unités géographiques sélectionnées.
- ▶ Sondage à deux degrés: On observe un sous-ensemble de ménages dans les unités géographiques sélectionnées. Il y aura échantillonnage (premier degré: les grappes) suivi de sous-échantillonnage (second degré: les ménages).

- ▶ C'est une liste de N unités de la population;
- ▶ Dans l'échantillonnage probabiliste, chaque unité a une chance positive d'être dans l'échantillon;
- ▶ Idéalement on possède une base de sondage ou encore on en construit une;
- ▶ Sinon, on peut construire des grappes avec l'aide d'une carte géographique et observer tous les ménages dans les grappes choisies: l'échantillonnage de grappes permet de contourner (du moins en théorie) les problèmes lorsqu'il n'y a pas de base de sondage.

Exemple d'une base de sondage: population MU284 (SSW, Appendice B, p. 652)

LABEL	P85	P75	RMT85	CS82	SS82	S82	ME84	REV84	REG	CL
1	?	27	288	13	24	49	2135	2836	1	1
2	?	15	139	14	12	41	957	2035	1	1
3	?	20	196	12	14	41	1530	6030	1	1
4	?	15	159	12	19	41	1059	4704	1	1
5	?	52	536	20	27	61	3951	5183	1	1
6	?	15	134	16	12	41	918	2157	1	2
7	?	62	623	18	27	61	4367	7072	1	2
8	?	54	517	15	32	61	4345	5246	1	2
9	?	12	96	10	12	31	754	951	1	2
10	?	50	467	14	29	61	3902	6067	1	2
11	?	29	277	14	20	45	1993	3264	1	3
12	?	14	155	10	21	41	1312	1899	1	3
13	?	40	386	24	13	51	2780	5931	1	3
14	?	27	241	24	8	45	1649	3877	1	3
15	?	43	422	19	18	51	2983	4968	1	3
16	?	671	6263	34	41	101	45324	59877	1	4
17	?	78	612	14	31	61	5331	7027	1	4
18	?	54	532	23	23	61	3994	6529	1	4
19	?	28	250	9	22	41	1616	2208	1	4
20	?	55	412	20	27	61	3240	3976	1	4
...										
282	?	27	226	7	28	49	1682	2898	8	50
283	?	9	63	5	19	41	604	594	8	50
284	?	31	233	5	27	45	1788	2366	8	50

- ▶ Dans l'exemple précédent, la variable P85 est inconnue et on veut faire un sondage afin d'estimer le total de la population;
- ▶ Les autres variables sont des variables auxiliaires qui sont utiles au niveau de la:
 1. Conception d'un sondage; construction d'estimateurs.
- ▶ Variables auxiliaires quantitatives: P75, RMT85, CS82, . . . ,
- ▶ Variables auxiliaires qualitatives: REG, CL.

Pour un échantillon s , on introduit

- ▶ Taille de l'échantillon: n_s ;
- ▶ Moyenne échantillonnale:

$$\bar{y}_s = \frac{1}{n_s} \sum_s y_k$$

- ▶ Variance échantillonnale:

$$S_{ys}^2 = \frac{1}{n_s - 1} \sum_s (y_k - \bar{y}_s)^2$$

On considère $\bar{y}_s = \frac{1}{n_s} \sum_s y_k$.

- ▶ Les propriétés statistiques de cette quantité dépendent du plan d'échantillonnage;
- ▶ Pour le tirage aléatoire simple sans remise, on sait que:

$$E_{SI}(\bar{y}_s) = \bar{y}_U.$$

- ▶ Ainsi, la moyenne échantillonnale possède la propriété d'absence de biais lorsque le plan de sondage est le tirage aléatoire simple sans remise.

- ▶ Pour le tirage aléatoire simple sans remise, la moyenne de \bar{y}_s sur tous les s égale la moyenne de la population \bar{y}_U ;
- ▶ Le biais d'un estimateur est formellement défini ainsi:

$$E_p(\hat{\theta}) - \theta$$

- ▶ Biais positif: en moyenne, les estimateurs excèdent la véritable valeur du paramètre;
- ▶ Biais négatif: en moyenne, les estimateurs sont inférieurs à la véritable valeur du paramètre.