

MAT 3777
Échantillonnage et sondages

Chapitre 8
Sujets choisis

P. Boily (uOttawa)

Session d'hiver – 2022

P. Boily (uOttawa)

Aperçu

8.1 – Échantillonnage avec probabilité proportionnelle à la taille (p.3)

- Méthodes de sélection d'un PPT avec remise (p.5)
- Estimation des paramètres (p.9)

8.2 – Échantillonnage à plusieurs degrés (p.16)

- Échantillonnage aléatoire simple à deux degrés (p.18)
- Estimation des paramètres (p.20)

8.3 – Échantillonnage à plusieurs phases (p.28)

- Échantillonnage aléatoire simple à deux phases (p.29)
- Estimation de la variance d'échantillonnage (p.34)

8.4 – Méli-mélo (p.38)

- Effet de plan (p.39)
- Ajustement pour la non-réponse (p.41)
- Estimation de la taille d'une population (p.46)
- Réponse aléatoire (p.53)
- Échantillonnage de Bernoulli (p.60)

8.1 – Échantillonnage avec probabilité proportionnelle à la taille

En pratique, la **taille** (que cela soit une caractéristique physique ou non) des unités d'échantillonnage est souvent très **variable** – un EAS n'est pas toujours efficace puisqu'il ne tient pas compte de l' de la population.

On peut parfois mettre à profit des afin de sélectionner un échantillon donnant un estimateur plus précis des paramètres d'intérêt.

Une façon possible de s'y prendre:

Exemple: en général, plus la superficie d'un pays est élevé, plus sa population l'est aussi ($\rho =$).

Si on cherche à estimer la population de la planète, il pourrait être souhaitable d'adopter un système d'échantillonnage dans lequel la probabilité de sélection d'un pays est **proportionnelle à sa superficie** – dans un EAS, il est fort probable que ni la , ni l' ne soient sélectionnées, ce qui entraîne une sous-estimation du total recherché.

Si la variable d'intérêt est liée (plus ou moins) à la taille de l'unité, on peut assigner une

.

Dans un PPT, les unités prélevées au préalable peuvent être **remises** dans la population, permettant la .

8.1.1 – Méthodes de sélection d'un PPT avec remise

Nous allons considérer deux méthodes de sélection d'un échantillon PPT:

- n fois, et
- .

Dans les deux cas, la procédure de sélection de l'échantillon PPT consiste à associer à chaque unité une n (ce sont souvent des n , mais ce n'est pas nécessaire), liée à la **taille de l'unité**, et à prélever les unités qui correspondent à des nombres choisis dans l'ensemble de nombres associés à la population .

Avec la **méthode des totaux cumulés**, la **taille** de la i -ème unité (dans une population contenant N unités) est dénotée par x_i , $1 \leq i \leq N$.

On associe ensuite une **étendue** à chaque unité de la manière suivante:

Unité	Étendue
1	à
2	à
3	à
⋮	⋮
$N - 1$	à
N	à

Finalement, on prélève un échantillon PPT en choisissant n entiers entre 1 et $X = x_1 + \cdots + x_{N-1} + x_N$ () et en sélectionnant les unités **associées à ces entiers**.

Exemple: Dans un village, il y a 8 vergers, contenant respectivement un certain nombre de pommiers. Un échantillon de $n = 3$ vergers est prélevé (avec remise), de manière proportionnelle au nombre de pommiers.

# série i	Taille x_i	Taille cumulée	Étendue associée
1	50		
2	30		
3	25		
4	40		
5	26		
6	44		
7	20		
8	35		

On choisit $n = 3$ entiers au hasard entre 1 et : , , et , par exemple. Les unités associées sont la i ème, la i ème, et la i ème.

Avec la **méthode de Lahiri**, on dénote toujours la taille d'une unité par x_i , $1 \leq i \leq N$, mais sans avoir
(ce qui peut s'avérer fastidieux, même avec un ordinateur).

La méthode consiste à sélectionner un couple (i, j) d'entiers au hasard, où $1 \leq i \leq N$ et $1 \leq j \leq M = \max\{x_i | 1 \leq i \leq N\}$.

Si $j \leq x_i$, la i -ème unité est ajoutée à l'échantillon. Sinon, on rejette la paire (i, j) et on continue jusqu'à ce que n unités aient été choisies.

 Il y a d'autre façon de s'y prendre; ce qui importe, c'est d'avoir un **mécanisme pour sélectionner un échantillon PPT**.

 Il est préférable de prélever sans remise, mais l'échantillonnage avec remise offre une approximation raisonnable si est “ ”.

8.1.2 – Estimation des paramètres

Revisitons l'exemple des vergers, dans lequel u_i représente le rendement de tous les pommiers du i -ème verger.

# série i	# pommiers x_i	π_i	Rendement
1	50	$50/270$	$u_1 = 2250$
2	30	$30/270$	$u_2 = 1080$
3	25	$25/270$	$u_3 = 1300$
4	40	$40/270$	$u_4 = 1400$
5	26	$26/270$	$u_5 = 1196$
6	44	$44/270$	$u_6 = 1716$
7	20	$20/270$	$u_7 = 820$
8	35	$35/270$	$u_8 = 1680$

On s'intéresse à la production **totale** de pommes du village, $\tau =$.

Puisqu'**en principe**, un verger qui contient plus de pommiers devrait produire plus de pommes, on prélève un échantillon PPT (avec remise) de $n = 3$ unités, où le nombre de pommiers dans verger représente sa taille.

Dans ce qui suit, nous illustrerons les concepts à l'aide de l'échantillon

$$\{y_1 = \quad, y_2 = \quad, y_3 = \quad\}.$$

Si l'échantillon \mathcal{Y} , avec $|\mathcal{Y}| = n$, est prélevé de la population \mathcal{U} à partir d'un plan d'échantillonnage PPT, les unités y_1, \dots, y_n sont et distribuées selon

$$\begin{array}{c|ccccc} y_i & u_1 & \cdots & u_j & \cdots & u_N \\ \hline p(y_i) & \pi_1 & \cdots & \pi_j & \cdots & \pi_N \end{array}$$

où pour tout $1 \leq j \leq N$ et .

Pour tout $1 \leq i \leq n$, posons $w_i = \frac{u_j}{\pi_j}$, si $y_i = u_j$ pour un $1 \leq j \leq N$. Les **poids** w_i sont également et distribués selon

$$P(y_i = u_j) = \pi_j, \quad 1 \leq i \leq n, 1 \leq j \leq N.$$

On remarque que, pour tout $1 \leq i \leq n$, l'**espérance des poids** équivaut à

$$E(w_i) = \sum_{j=1}^N w_j P(w_i = w_j) = \sum_{j=1}^N \frac{u_j}{\pi_j} \pi_j = \sum_{j=1}^N u_j = \tau.$$

C'est donc dire que

$$\hat{\tau}_{\text{ppt}} = \bar{w} = \frac{1}{n} \sum_{i=1}^n w_i = \frac{1}{n} \sum_{i=1}^n \frac{y_i}{\pi_j} = \tau.$$

offre un **estimateur sans biais** du total τ .

La **variance d'échantillonnage** se calcule comme suit:

$$V(\hat{\tau}_{\text{ppt}}) = V\left(\frac{1}{n} \sum_{i=1}^n w_i\right) = \frac{1}{n^2} \underbrace{\sum_{i=1}^n V(w_i)} =$$

=

=

.

En pratique, on ne connaît pas τ , alors on utilise l'estimateur **non-biaisé**

$$\hat{V}(\hat{\tau}_{\text{ppt}}) = \quad .$$

Théorème de la limite centrée – PPT

Si n et $N - n$ sont suffisamment élevés, alors

$$\hat{\tau}_{\text{ppt}} \sim_{\text{approx.}} \quad .$$

La **marge d'erreur sur l'estimation** et l'**intervalle de confiance** de τ à **environ 95%** sont ainsi

$$\hat{B}_{\tau;\text{ppt}} = \quad \text{et} \quad \text{IC}_{\text{ppt}}(\tau; 0.95) = \quad .$$

Exemple: Dans l'exemple des vergers, nous avons

$$\hat{\tau}_{\text{ppt}} = \frac{1}{3} \left[\underbrace{\quad}_{w_1} + \underbrace{\quad}_{w_2} + \underbrace{\quad}_{w_3} \right] = \quad ;$$

$$\hat{V}(\hat{\tau}_{\text{ppt}}) = \frac{1}{3(2)} \left[\left(\underbrace{\quad}_{w_1} \right)^2 + \left(\underbrace{\quad}_{w_2} \right)^2 + \left(\underbrace{\quad}_{w_3} \right)^2 - 3 \cdot \underbrace{\quad}_{\hat{\tau}_{\text{ppt}}^2} \right]$$

$$= \quad .$$

Conséquemment, l'intervalle de confiance pour le rendement total des pommiers du village à environ 95% est

$$IC_{\text{ppt}}(\tau; 0.95) = \quad .$$

La valeur réelle $\tau =$ **ne se retrouve pas** dans l'intervalle de confiance – pourquoi est-ce le cas? Est-ce problématique?

En général, $V(\hat{\tau}_{\text{ppt}}) \leq V(\hat{\tau}_{\text{EAS}})$. Dans l'exemple des pommiers, on peut montrer que

$$V(\hat{\tau}_{\text{EAS}}) \approx \quad , \quad \text{et}$$

$$V(\hat{\tau}_{\text{ppt}}) \approx \quad .$$

On peut également donner un estimé de la **moyenne** de population μ à l'aide de

$$\hat{\mu}_{\text{ppt}} = \quad , \quad \hat{V}(\hat{\mu}_{\text{ppt}}) = \quad , \quad \text{IC}_{\text{ppt}}(\mu; 0.95) = \quad .$$

8.2 – Échantillonnage à plusieurs degrés

En séparant l'échantillonnage en plusieurs étapes, on peut et

Dans un **échantillonnage à plusieurs degrés** (EnD), on prélève un échantillon d'unités de grande taille (), puis des sous-unités de ces grandes unités (), etc.

Exemple: l'échantillonnage d'une province peut se faire en trois étapes:

1. échantillon de municipalités (**unités primaires**),
2. échantillon de quartiers par municipalités (**unités secondaires**), et
3. échantillon de ménages par quartiers (**unités tertiaires**).

Dans un *EnD*, l'échantillon est concentré autour de plusieurs **pivots**: dans les études sur le terrain, par exemple, cela à l'avantage de réduire considérablement la surface d'enquête, ce qui aide à (en plus de).

De plus, il arrive souvent que l'on dispose d'informations détaillées pour des **groupes** d'unités d'échantillonnage, mais par pour des unités **individuelles**: il n'est donc pas nécessaire d'obtenir une base de sondage (pour **toutes** les unités d'échantillonnage), mais seulement pour celles appartenant aux unités primaires sélectionnés lors du premier tour, par exemple.

On peut utiliser n'importe quelle méthode d'échantillonnage probabiliste à chaque stade, et elles peuvent changer d'un stade à l'autre (un EAS de municipalités, un EAS de quartiers, un SYS de ménages, par exemple).

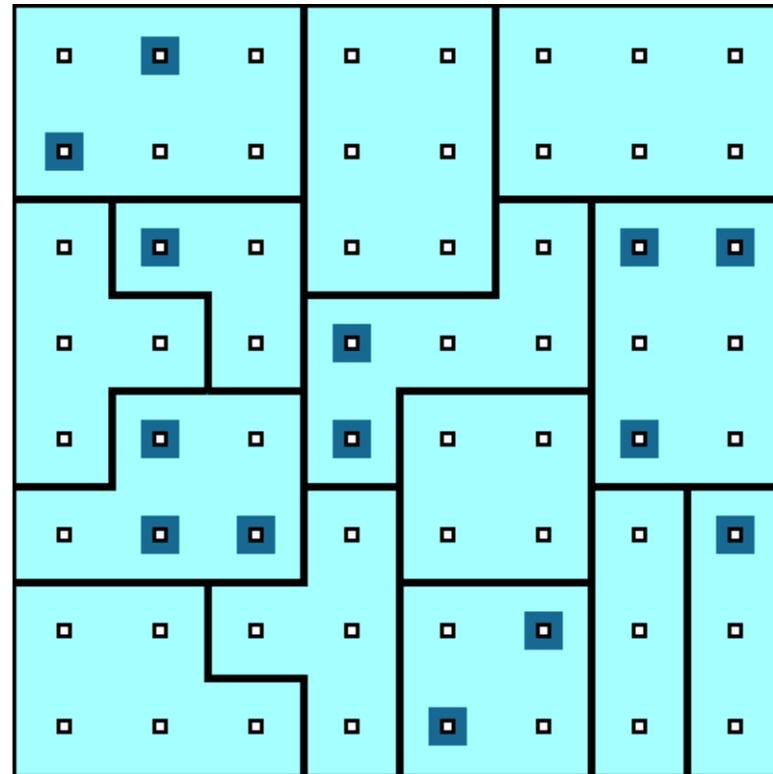
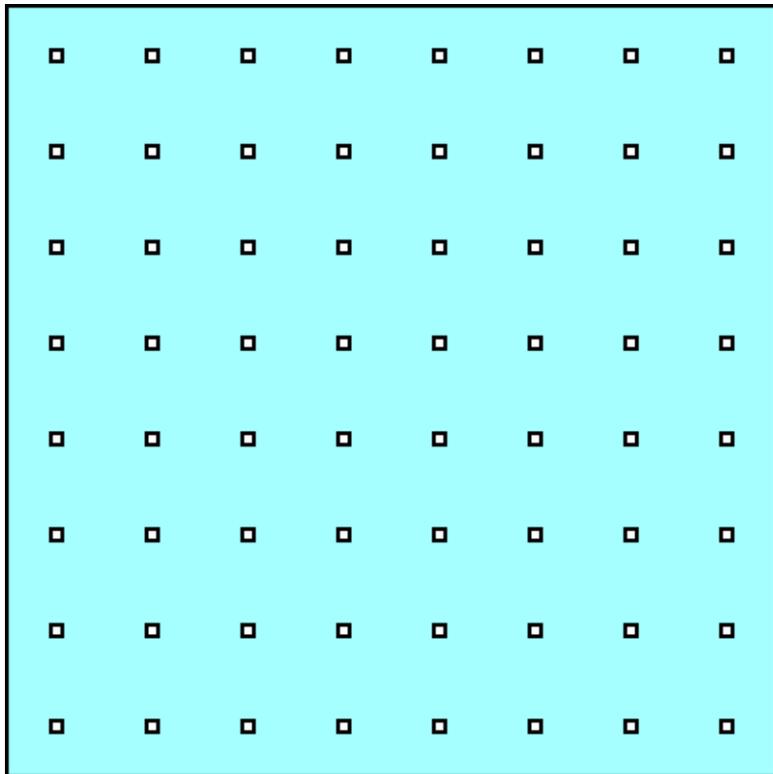
8.2.1 – Échantillonnage aléatoire simple à deux degrés

Si les deux étapes de la sélection se font par EAS, la méthode prend le nom d'**échantillonnage aléatoire simple à deux degrés** (EAS2D).

Exemple: on peut estimer la biomasse d'une espèce de plante dans une superficie forestière composée de 40 compartiments (unités primaires) en prélevant un EAS de $m = 8$ compartiments des $M = 40$ compartiments composant la population à l'étude.

Pour chacun de ces compartiments $1 \leq i \leq m$, on prélève ensuite un EAS de n_i parcelles, et on mesure la biomasse en question.

On peut calculer les estimations de la quantité moyenne ou totale de biomasse dans la superficie forestière à l'aide des formules appropriées.



8.2.2 – Estimation des paramètres

Soient une population constituée de M unités primaires, et possédant N_ℓ unités secondaires dans la ℓ -ème unité primaire.

Notons par $u_{i,j}$ la valeur de la variable réponse de la j -ième unité du second degré dans la i -ième unité du premier degré.

La **moyenne de la population** est

$$\mu = \cdot$$

Supposons que l'on prélève un EAS de m unités primaires, et un EAS de n_i unités secondaires dans la i -ème unité primaire.

L'échantillon est donc de taille $n = n_1 + \dots + n_m$.

On obtient un estimateur non biaisé de μ grâce à l'équation

$$\bar{y}_{\text{EAS2D}} = \frac{1}{n} \sum_{i=1}^m n_i \bar{y}_i, \quad ,$$

où

$$\bar{N} = \frac{1}{m} \sum_{i=1}^m N_i .$$

La variance d'échantillonnage est composée de deux éléments:

- une mesure de la variation σ^2 , et
- une mesure de la variation $\frac{1}{m} \sum_{i=1}^m \frac{N_i^2}{N}$.

Lorsque $n_i = N_i$, pour tout $1 \leq i \leq N_i$, on fait affaire à un EnD et la variance est donnée $\sigma^2 \frac{1}{m} \sum_{i=1}^m \frac{N_i^2}{N}$ (cf. chapitre 6).

Dans le cas où $m = M$, on fait affaire à un STR (cf. chapitre 3) et la variance est donnée σ^2 .

Quand $m \neq M$ et $n_i \neq N_i$ pour au moins un i , la variance est une combinaison de ces deux extrêmes: dans ce cas, le second terme représente **la contribution du sous-échantillonnage** (un autre nom pour un EnD).

On peut se servir du **théorème de la variance totale** afin d'estimer la variance d'échantillonnage:

$$\begin{aligned}
 V(\bar{y}_{EAS2D}) &= \\
 &= \\
 &\approx \quad ,
 \end{aligned}$$

où

$$s_T^2 = \quad , \quad s_i^2 = \quad .$$

Exemple:

On mesure la biomasse d'une espèce de plante (kg) dans des parcelles de 0.025 ha (unités secondaires) sélectionnées dans $m = 8$ compartiments (unités primaires) choisis au hasard parmi les $M = 40$ compartiments d'une étendue forestière.

Le sommaire des résultats se retrouve dans le tableau suivant:

Comp.	1	2	3	4	5	6	7	8
\bar{y}_i	118	107	109	110	120	95	93	90
s_i^2	436	516	586	456	412	497	755	496
N_i	1760	1975	1615	1785	1775	2050	1680	1865
n_i	9	10	8	9	9	10	8	9

Déterminer des intervalles de confiance (à environ 95%) de la biomasse moyenne par parcelle et par compartiment, et de son total dans la forêt.

Solution: Puisque l'on ne connaît pas \bar{N} , on l'approxime à l'aide de la moyenne

$$\bar{N} \approx \frac{1}{8}(1760 + \dots + 1865) = \quad .$$

On calcule ensuite les totaux dans les unités primaires sélectionnées:

Comp.	1	2	3	4	5	6	7	8
$N_i \bar{y}_i (\times 10^5)$								

et les estimateurs EAS2D de la moyenne μ , de la moyenne des totaux dans les compartiments, et du total sont:

$$\bar{y}_{EAS2D} = \frac{1}{8(\quad)} (\quad + \dots + \quad) \times 10^5 = \quad ;$$

$$\bar{N} \bar{y}_{EAS2D} = \quad . \quad = \quad ; \quad \tau_{EAS2D} = M \cdot \bar{N} \bar{y}_{EAS2D} = \quad .$$

La variance entre les compartiments (unités primaires) est ainsi:

$$s_T^2 = \frac{1}{8-1} \sum_{i=1}^8 (N_i \bar{y}_i - \quad)^2 =$$

Finalement, on calcule la variance à même les compartiments:

Comp.	1	2	3	4	5	6	7	8
$\frac{N_i^2}{N^2} \cdot \frac{s_i^2}{n_i} \left(1 - \frac{n_i}{N_i}\right)$								

La variance d'échantillonnage devient alors

$$\hat{V}(\bar{y}_{EAS2D}) = \frac{1}{8 \left(\quad \right)^2} \left(1 - \frac{8}{40}\right) + \frac{1}{8(40)} \left(\quad + \dots + \quad \right) =$$

Les variances des deux autres estimateurs se calculent aisément:

$$\hat{V}(\bar{N}\bar{y}_{EAS2D}) = \bar{N}^2 \hat{V}(\bar{y}_{EAS2D}) = (\quad)^2 \cdot \quad = \quad ;$$

$$\hat{V}(\tau_{EAS2D}) = M^2 \bar{N}^2 \hat{V}(\bar{y}_{EAS2D}) = (40)^2 \cdot (\quad)^2 \cdot \quad = \quad ;$$

tout comme les intervalles de confiance:

$$IC_{EAS2D}(\mu; 0.95) :$$

$$IC_{EAS2D}\left(\frac{N_0}{M}\mu; 0.95\right) :$$

$$IC_{EAS2D}(\tau; 0.95) :$$

... en supposant bien sûr que le théorème de la limite centrée demeure valide dans le contexte d'un EAS2D.

8.3 – Échantillonnage à plusieurs phases

L'**échantillonnage à plusieurs phases** (EnP) joue un rôle crucial dans plusieurs types d'enquêtes, incluant entre autre les enquêtes à distance, telle que celles menées par **téledétection**.

Lors de la première phase, on prélève un nombre n_1 d'unités, mais on ne capte qu'un k_1 nombre de caractéristiques pour chaque unité.

Dans chaque phase successive, on mesure un plus k_2, k_3, \dots nombre de caractéristique sur un plus n_2, n_3, \dots échantillon d'unités.

De cette façon, on arrive à estimer le paramètre visé avec **plus de précision** et à un **plus faible coût**, en étudiant la relation entre les caractéristiques mesurées lors des différentes phases.

8.3.1 – Échantillonnage aléatoire simple à deux phases

Un EnP qui ne comporte que deux phases prend le nom d'**échantillonnage à deux phases (E2P)**, ou **échantillonnage double**.

Les E2P sont particulièrement utiles dans une situation où l'énumération de la τ est dispendieuse (en \$\$\$ ou en main d'oeuvre), mais dans laquelle on peut aisément observer une μ .

Il est ainsi parfois préférable de prélever, en **première phase**, un EAS de n_1 afin d'analyser la variable auxiliaire, ce qui mène à des estimations précises (tout du moins, c'est ce que l'on espère) de τ ou de μ pour cette variable auxiliaire.

Lors de la seconde phase, on choisit un n_2 échantillon, généralement un n_2 , dans lequel on mesure et la **caractéristique principal** et la **variable auxiliaire**.

On obtient ensuite des estimations de la caractéristique principale à l'aide des renseignements obtenus lors de la **première phase**, en utilisant la τ_{v_2} ou la τ_{c_2} .

On peut augmenter la précision des estimations finales en incluant **plusieurs variables auxiliaires corrélées**, au lieu d'une seule.

Exemple: Si l'on cherche à estimer le volume total τ d'une forêt, on commence par mesurer la circonférence c_i et la hauteur h_i des arbres i dans un échantillon, puis le volume v_{i_k} des arbres i_k dans un sous-échantillon. On détermine ensuite la relation statistique entre τ_v , τ_c , et τ_h , et voilà!

Le mode d'échantillonnage E_nP aide à réduire le \dots et à accroître la \dots .

On peut aussi s'en servir afin de \dots une population: un premier échantillon est prélevé en se fondant sur la caractéristique auxiliaire, que l'on utilise pour subdiviser la population en strates dans lesquelles la caractéristique principale est plus ou moins \dots .

Tant que les deux caractéristiques sont **corrélées**, on obtient ainsi des estimations précises de la caractéristique principale à partir d'un deuxième échantillon relativement petit.

On peut aussi jumeler le mode E_2P avec le mode E_2D , par exemple (ou avec n'importe quel mode d'échantillonnage).

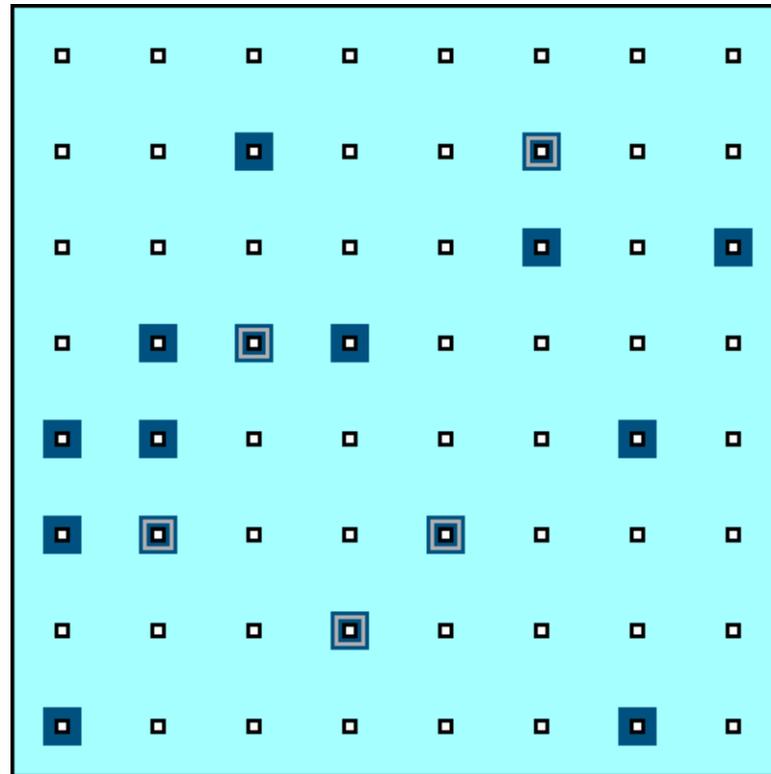
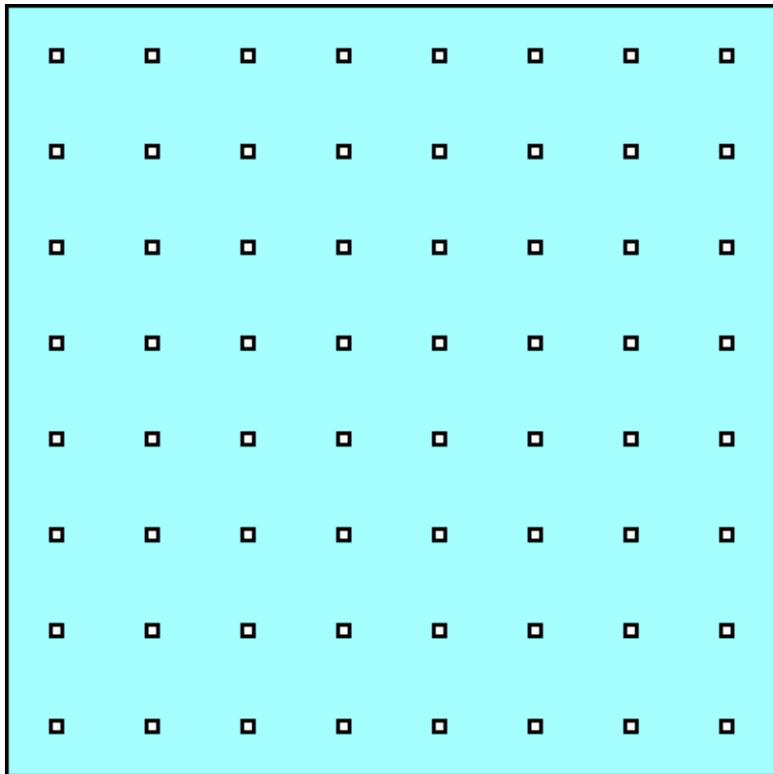
Si les deux étapes de sélection se font par EAS, la méthode prend le nom d'**échantillonnage aléatoire simple à deux phases (EAS2P)**.

Lors de la première phase, la population est divisée en unités d'échantillonnage bien définies et on y prélève un EAS \mathcal{Y}_1 de taille n_1 ; on mesure la **variable auxiliaire** x sur toutes les unités de \mathcal{Y}_1 .

Ensuite, on prélève un sous-EAS $\mathcal{Y}_2 \subseteq \mathcal{Y}_1$ de taille n_2 ; on mesure la **caractéristique principale** y sur toutes les unités de \mathcal{Y}_2 .

On évalue les paramètres $r_{\mathcal{Y}_2}$ ou $b_{\mathcal{Y}_2}$ à partir de \mathcal{Y}_2 (à l'aide de la méthode du quotient ou de la méthode de la régression), ce qui donne

$$\hat{\mu}_{Y;R;EAS2P} = \quad \text{ou} \quad \hat{\mu}_{Y;L;EAS2P} = \quad .$$



8.3.2 – Estimation de la variance d'échantillonnage

En raison du **double échantillonnage**, on retrouve deux termes qui contribuent à chacune des variances d'échantillonnage des estimateurs (la première lorsque l'on passe de \mathcal{U} à \mathcal{Y}_1 , et la seconde de \mathcal{Y}_1 à \mathcal{Y}_2):

$$\hat{V}(\hat{\mu}_{Y;R;EAS2P}) =$$

$$\hat{V}(\hat{\mu}_{Y;L;EAS2P}) =$$

où $s_{\mathcal{Y}_2}^2$, s_{XY} et s_X^2 représentent les quantités habituelles (dans \mathcal{Y}_2),

$$r_{\mathcal{Y}_2} = \quad , \quad b_{\mathcal{Y}_2} = \quad , \quad \text{et} \quad s_{XY;L}^2 = \quad .$$

Exemple:

On s'intéresse à la biomasse d'une plante quelconque dans une région, divisée en parcelles de 0.025 ha chacune.

En premier lieu, on mesure le nombre de bosquets x par unité dans un EAS \mathcal{Y}_1 de $n_1 = 200$ parcelles. Ensuite, on calcule la biomasse y de la plante en question dans chaque unité d'un sous-EAS \mathcal{Y}_2 de $n_2 = 40$ parcelles:

$$\bar{x}_{\mathcal{Y}_1} = 374.4; \quad \sum_{i=1}^{40} x_i = 15,419; \quad \sum_{i=1}^{40} y_i = 2104;$$

$$\sum_{i=1}^{40} x_i^2 = 7,744,481; \quad \sum_{i=1}^{40} x_i y_i = 960,320; \quad \sum_{i=1}^{40} y_i^2 = 125,346.$$

Donner un I.C. de la biomasse moyenne par parcelle, à environ 95%.

Solution: on calcule les quantités intermédiaires requises:

$$\bar{x}_{y_2} = \quad ; \quad \bar{y}_{y_2} = \quad ; \quad r_{y_2} = \frac{\bar{y}_{y_2}}{\bar{x}_{y_2}} = \quad ;$$

$$s_X^2 = \quad ; \quad s_Y^2 =$$

$$s_{XY} = \quad ; \quad b_{y_2} = \frac{s_{XY}}{s_X^2} = \quad ;$$

$$s_{XY;L}^2 = \quad ;$$

ce qui donne

$$\hat{\mu}_{Y;R;EAS2P} = \quad ; \quad \hat{\mu}_{Y;L;EAS2P} =$$

et

$$\hat{V}(\hat{\mu}_Y; R; \text{EAS2P}) =$$

;

$$\hat{V}(\hat{\mu}_Y; L; \text{EAS2P}) =$$

;

d'où

$$\text{IC}_{R; \text{EAS2P}}(\mu_Y; 0.95) =$$

$$\text{IC}_{L; \text{EAS2P}}(\mu_Y; 0.95) =$$

.

8.4 – Méli-mélo

Nous terminons le cours en discutant brièvement de quelques notions qui n'ont pas trouvé de place naturelle dans les sections précédentes:

- les effets de plan;
- l'ajustement pour la non-réponse;
- l'estimation de la taille d'une population,
- la méthode de la réponse aléatoire, et
- l'échantillonnage de Bernoulli.

8.4.1 – Effet de plan

[Adapté de *Méthodes et pratiques d'enquête*, Statistique Canada]

L'**effet de plan** compare la variance des estimateurs entre un plan d'échantillonnage et un EAS. Il s'agit du rapport entre la

, et la

Cette mesure est souvent appliquée pour comparer l'**efficience** des estimateurs de divers plans d'échantillonnage. Si le rapport < 1 , le plan d'échantillonnage est plus efficace que l'EAS; s'il est > 1 , il est moins efficace que l'EAS.

Nous avons comparé directement les variances théoriques de plusieurs plan d'échantillonnage aux sections 3.4, 4.5, et 6.4 – typiquement, on calcule l'effet de plan à l'aide des échantillons réalisés.

Les effets du plan d'échantillonnage aident aussi à obtenir des estimations approximatives de la variance pour des plans d'**échantillonnage complexes**.

Si une estimation de l'effet du plan d'échantillonnage est disponible dans une enquête précédente qui a utilisé le même plan d'échantillonnage, elle peut servir à déterminer la **taille de l'échantillon nécessaire de l'enquête**.

8.4.2 – Ajustement pour la non-réponse

[*Ibid.*] Les non-réponses représentent un problème dans **toutes** les enquêtes.

La **non-réponse totale** (lorsque toutes les données ou presque d'une unité échantillonnée sont manquantes) survient lorsque:

- une unité de l'échantillon _____ au sondage;
- il est impossible d' _____ ;
- l'unité ne peut être _____ , ou encore
- si l'information obtenue est _____ .

La façon la plus simple de traiter ces non-réponses est de les ignorer; dans certaines circonstances **exceptionnelles**, des proportions ou des moyennes estimées sans ajustement pour les non-réponses totales sont **plus ou moins identiques** à celles produites en appliquant un ajustement pour les non-réponses.

Si l'on néglige de **compenser** pour les unités non répondantes, les (e.g. la taille d'une population, le total des revenus ou le total d'acres récoltés).

La façon la plus commune de traiter la non-réponse totale est d' en supposant que les unités répondantes représentent les unités répondantes et non répondantes.

(Est-ce un hypothèse vraisemblable, en pratique?)

Si les **non-répondants sont équivalents aux répondants** pour les caractéristiques mesurées dans l'enquête, c'est une approche raisonnable.

Les poids de base pour les non-répondants sont ensuite redistribués entre les répondants, à l'aide d'un **facteur d'ajustement pour les non-réponses** qui est multiplié par la poids de base, afin d'obtenir une pondération ajustée.

Par exemple, si on prélève un EAS de taille $n = 25$ d'une strate de taille $N = 1000$, la **probabilité d'inclusion** de chacune de ces unités et le **poids de base** correspondant sont

$$\pi = \frac{n}{N} =$$
$$w = \frac{1}{\pi} = .$$

Chaque unité sélectionnée représente 40 unités dans la strate.

Si nous n'obtenons une réponse que de $n_r = 20$ des $n = 25$ unités sélectionnées, le **facteur d'ajustement pour les non-réponses** et la **pondération ajustée** (pour la non-réponse) deviennent:

$$\text{FANR} = \frac{n}{n_r} =$$

$$w_{nr} = w \cdot \text{FANR} = \quad ;$$

chaque unité répondante représente alors 50 unités dans la strate. C'est avec cette pondération ajustée que l'on travaillerait.

Il va de soit que la pondération ajustée peut varier d'une strate à l'autre, en fonction du plan d'échantillonnage et de la taille de l'échantillon.

Lorsque l'on cherche à déterminer la taille de l'échantillon et sa répartition dans diverses strates, on obtient en pratique la taille de l' (on suppose ici que les populations cible et à l'étude coïncident). On peut avoir alors recours à l' de la taille de l'échantillon.

Exemple: on détermine que la répartition d'un STR de taille $n = 29$ est $(17, 9, 3)$. Lors d'une étude préalable, on a déterminé que les taux de non-réponse par strate sont de $(16.2\%, 20.8\%, 31.2\%)$. Quelle répartition optimise les chances d'obtenir la répartition visée?

Solution:

8.4.3 – Estimation de la taille d'une population

Comment s'y prend-on si la taille N de la population \mathcal{U} est inconnue? Lorsque la population est suffisamment large, on peut toujours utiliser l'approximation $N \approx \infty$ dans les formules de variance d'échantillonnage.

Mais c'est parfois le paramètre N qui représente la quantité d'intérêt.

Exemple: combien de billets de 5 dollars, N , y a-t-il en circulation?

On donne un estimé de N à l'aide de la méthode de :

1. on capture n_1 billets au hasard (sans remise) dans la population;
2. on les marque et on les remet en circulation;
3. à un moment ultérieur, on capture n_2 billets au hasard (sans remise) dans la population;
4. on compte le nombre X de billets marqués, $0 < X \leq n_2$.

Si on attend assez longtemps (question de laisser les billets marqués se propager dans la population), on obtient

$$\frac{n_1}{N} \approx \frac{X}{n_2}, \quad \text{d'où } \hat{N} = \frac{n_1 n_2}{X},$$

où $X \sim \text{loi binomiale}(n_2, \frac{n_1}{N})$ dont les paramètres sont $n_1, N - n_1, n_2$, et

$$P(X = x) = \binom{n_2}{x} \left(\frac{n_1}{N}\right)^x \left(\frac{N - n_1}{N}\right)^{n_2 - x}, \quad 0 \leq x \leq n_2$$

$$\mu_X = \text{E}[X] = n_2 \frac{n_1}{N}, \quad \sigma_X^2 = \text{V}[X] = n_2 \frac{n_1}{N} \left(\frac{N - n_1}{N}\right).$$

Si $\frac{n_2}{N} < 0.05$, on peut ignorer le FCPF dans la variance:

$$\sigma_X^2 = V[X] \approx \dots$$

On peut maintenant développer des expressions pour $E[\hat{N}]$ et $V[\hat{N}]$, en se servant de la **série de Taylor d'ordre 2 près de** $X \approx \mu_X = n_2 p$:

$$f(X) \approx \dots$$

Si $\hat{N} = f(X) = \frac{n_1 n_2}{X}$, alors

$$\begin{aligned} \hat{N} &\approx \\ &= \dots \end{aligned}$$

Conséquemment,

$$E[\hat{N}] =$$

$$=$$
$$=$$
$$=$$
$$=$$

Puisque $\frac{1-p}{n_2 p} > 0$, $E[\hat{N}] \neq N$ et l'estimateur \hat{N} est **asymptotiquement non-biaisé** lorsque la taille n_2 du second échantillon augmente.

On obtient un estimateur de la variance en utilisant de la **série de Taylor d'ordre 1 près de** $X \approx \mu_X = n_2 p$:

$$\hat{N} \approx \dots = \dots = \dots .$$

Dans ce cas,

$$V[\hat{N}] \approx \dots = \dots = \dots \approx \dots = \dots .$$

En pratique, on en connaît pas p ; on utilise alors

$$\hat{V}[\hat{N}] = \frac{\hat{p}(1-\hat{p})}{n_2}, \quad \text{où } \hat{p} = \frac{\hat{N}}{N}.$$

Théorème de la limite centrée – taille de la population N

Si n_2 et N sont suffisamment élevés, alors

$$\hat{N} \underset{\text{approx.}}{\sim} N, \quad \text{,}$$

et l'intervalle de confiance de N à environ **95%** est ainsi

$$\text{IC}(N; 0.95) : \quad \left[\hat{N} - 1.96 \sqrt{\hat{V}[\hat{N}]}, \hat{N} + 1.96 \sqrt{\hat{V}[\hat{N}]} \right].$$

Exemple: supposons que $n_1 = 500$ billets aient été capturés et marqués initialement; des $n_2 = 300$ billets recapturés à la 2e étape, $X = 127$ étaient marqués. Donner un intervalle de confiance du nombre total de billets de 5\$ à environ 95%.

Solution:

8.4.4 – Réponse aléatoire

Avez-vous déjà triché lors d'un contrôle durant la pandémie?

Avec un “**Oui**”, on peut vraisemblablement conclure que c'est la vérité.

Mais puisqu'il y a un **coût social** associé à une telle réponse, on peut s'attendre à ce que certains tricheurs répondent “**Non**”.

Comment peut-on s'y prendre afin de réduire l'erreur de mesure pour les **questions délicates**?

Première approche: avec de telles questions, la compétence de l'enquêteur.e joue un rôle crucial – il ne faut pas négliger ce volet.

Seconde approche: la technique de la **réponse aléatoire** nécessite l'utilisation de deux questions:

- la question Q_1 , et
- un question Q_2 ,

et d'un θ à **paramètres connus** (pile ou face, etc.).

Le principe est le suivant: la répondante tire à pile ou face (sans annoncer le résultat à l'enquêteur), et elle répond honnêtement à une des 2 questions:

- **“face”**: “Avez-vous déjà triché lors d'un contrôle?”;
- **“pile”**: “Êtes-vous née en janvier?”;

Puisque l'enquêteur ne connaît pas le résultat du tirage au sort, il ne sait pas si la répondante répond à la question délicate ou à la question innocente.

En théorie, l'anonymat assuré par la réponse aléatoire libère les répondants (le coût social est) – conséquemment, on peut s'attendre à une réponse honnête, quelle que soit la question.

 Cette approche ne peut porter fruit que si l'on connaît les probabilités:

- θ d'observer une réponse positive à la question innocente;
- ρ de poser la question délicate, et
- ϕ d'observer une réponse positive, quelle que soit la question.

Soit p la **proportion de réponses positives à la question délicate** ().

Ainsi,

$$\phi = P(\text{réponse positive})$$

=

,

=

d'où

$$p = \cdot$$

Si $\hat{\phi}$ représente la proportion de réponses positives dans l'échantillon réalisé, on peut construire l'**estimateur**

$$\hat{p}_{ra} = \frac{\hat{\phi}^\theta}{\hat{\phi}^\theta + (1 - \hat{\phi})^\rho}, \quad \theta, \rho \text{ des constantes,}$$

dont la variance est

$$V(\hat{p}_{ra}) = \frac{\hat{\phi}^\theta (1 - \hat{\phi})^\rho}{n} \left[\frac{\theta}{\hat{\phi}^{\theta+1}} + \frac{\rho}{(1 - \hat{\phi})^{\rho+1}} \right].$$

Puisque $\hat{\phi}$ est l'estimateur d'une proportion dans une population \mathcal{U} de taille N , obtenu à l'aide d'un EAS de taille n , sa **variance d'échantillonnage** est

$$V(\hat{\phi}) = \frac{\hat{\phi}(1 - \hat{\phi})}{n},$$

d'où

$$V(\hat{p}_{ra}) = \frac{\phi(1-\phi)}{n\rho^2}$$

Puisque ϕ est inconnu en général, on utilise l'estimateur (non-biaisé)

$$\hat{V}(\hat{p}_{ra}) = \frac{\hat{p}_{ra}(1-\hat{p}_{ra})}{n\rho^2},$$

et on construit un **intervalle de confiance de p à environ 95%** avec

$$IC_{ra}(p; 0.95) : \left[\hat{p}_{ra} \pm 1.96 \sqrt{\hat{V}(\hat{p}_{ra})} \right]$$

Le facteur $1/\rho^2$ vient **pénaliser l'incertitude** apportée par la réponse aléatoire – plus ρ est élevé, plus $\hat{V}(\hat{p}_{ra})$ est faible. Mais si ρ est trop élevé, l'anonymat conféré par l'approche s'évapore...

Exemple: on cherche à déterminer l'incidence de tricherie chez les étudiants ($N = 442$) du département de mathématiques et de statistique lors des cours en ligne offerts pendant la pandémie, à l'aide d'un EAS ($n = 65$). On se sert du stratagème décrit dans cette section avec $\rho = 1/2$, et on observe $\theta = \frac{52}{442}$ et $\hat{\phi} = \frac{21}{65}$. Déterminer un intervalle de confiance de la proportion des étudiants qui ont triché pendant la pandémie.

Solution:

8.4.5 – Échantillonnage de Bernoulli

[Adapté des notes de cours de D. Haziza]

L'**échantillonnage de Bernoulli** (BE) est un plan de sondage à taille
– il est impossible de fixer la **taille de l'échantillon a priori**.

On assigne à chaque unité de la population $\mathcal{U} = \{u_1, \dots, u_N\}$ la même probabilité d'inclusion dans l'échantillon \mathcal{Y} : $\pi_j = \pi \in (0, 1)$, pour tout j .

On dénote la **taille de l'échantillon obtenu** par n_a .

Le plan BE consiste à effectuer N épreuves de Bernoulli indépendantes, chacune avec probabilité de succès π (succès: ; échec:).

La probabilité d'obtenir un échantillon \mathcal{Y} de taille n_a devient alors

$$P(|\mathcal{Y}| = n_a) = \binom{N}{n_a} 2^{-N}.$$

Il y a 2^N échantillons possibles, dont la taille varie de $n_a = 0$ à $n_a = N$.

La taille de l'échantillon suit une **loi** $n_a \sim$:

$$P(n_a = n) = \binom{N}{n} 2^{-N}, \quad E[n_a] = N/2, \quad V[n_a] = N/4.$$

Lorsque N est suffisamment élevé, cette loi est **approximativement normale**; on peut alors construire un **intervalle de confiance de la taille de l'échantillon à environ 95%** à l'aide de

$$\text{IC}(n_a; 0.95) : \left[n_a - 1.96 \sqrt{N/4}, n_a + 1.96 \sqrt{N/4} \right].$$

Soit $\pi_{j,k}$ la probabilité d'inclusion des unités $j \neq k$ dans l'échantillon \mathcal{Y} . Puisque les épreuves de Bernoulli sont indépendantes les unes des autres,

$$\pi_{j,k} = P(\{u_j, u_k\} \in \mathcal{Y}) = \dots$$

L'estimateur

$$\hat{\tau}_{\text{BE}} = \frac{1}{\pi} \sum_{i=1}^{n_a} y_i$$

est un **estimateur sans biais du total** τ : en effet

$$E[\hat{\tau}_{\text{BE}}] = \dots,$$

puisque n_a et \bar{y} sont indépendants.

Dans le même ordre d'idée, la **variance d'échantillonnage de l'estimateur** $\hat{\tau}_{BE}$ peut être approchée par

$$\hat{V}[\hat{\tau}_{BE}] = \dots$$

Si N et n_a sont suffisamment élevés, le théorème de la limite centrée entre de nouveau en jeu, et on peut construire un **intervalle de confiance de τ à environ 95%** à l'aide de

$$IC_{BE}(\tau; 0.95) : \dots$$

Les estimateurs correspondants pour la moyenne \bar{y}_{BE} et la proportion \hat{p}_{BE} s'obtiennent de la manière habituelle.

Exemple: Une professeure doit corriger 600 copies d'examen. Pour chaque copie, elle lance un dé ne le corrige que s'il montre un 6. À la fin du processus, elle a corrigé 90 copies, desquelles 60 obtiennent une note de passage. Déterminer un IC à 95% du nombre total de succès dans la classe.

Solution: