

# Estimation de la variance inconditionnelle dans le cadre d'enquêtes complexes

Yves G. Berger



12 février 2026

Basé sur Berger (2025) "*International Statistical Review*"

# Introduction: exemples

- Données de population :  $y_i \sim N(\mu, \sigma)$   $i = 1, \dots, N$ , i.i.d.
- $N = 1000$ , population finie
- Échantillon  $S$  de taille  $n = 200 \implies \frac{n}{N} = 0.2$
- Nous souhaitons estimer la moyenne  $\mu$ .
- Considérons différents plans d'échantillonnage.

# Échantillonnage aléatoire simple sans remise

- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

# Échantillonnage aléatoire simple sans remise

- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{\mathbb{V}}(\bar{y}) := \left(1 - \frac{n}{N}\right) \frac{\hat{\sigma}^2}{n}, \quad \text{avec } \hat{\sigma}^2 := \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

# Échantillonnage aléatoire simple sans remise

- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{\mathbb{V}}(\bar{y}) := \left(1 - \frac{n}{N}\right) \frac{\hat{\sigma}^2}{n}, \quad \text{avec } \hat{\sigma}^2 := \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

- $\hat{\mathbb{V}}(\bar{y})$  **sans biais?**

# Échantillonnage aléatoire simple sans remise

- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{V}(\bar{y}) := \left(1 - \frac{n}{N}\right) \frac{\hat{\sigma}^2}{n}, \quad \text{avec } \hat{\sigma}^2 := \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

- $\hat{V}(\bar{y})$  **sans biais?**

$\Rightarrow \hat{V}(\bar{y})$  est **biaisé!**

Relative bias = -20%

# Échantillonnage aléatoire simple avec remise

- Échantillon  $\equiv n$  tirages indépendants
- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{V}(\bar{y}) := \frac{\hat{\sigma}^2}{n}, \quad \text{avec } \hat{\sigma}^2 := \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

# Échantillonnage aléatoire simple avec remise

- Échantillon  $\equiv n$  tirages indépendants
- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{\mathbb{V}}(\bar{y}) := \frac{\hat{\sigma}^2}{n}, \quad \text{avec } \hat{\sigma}^2 := \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

- $\hat{\mathbb{V}}(\bar{y})$  **sans biais?**



## Échantillonnage aléatoire simple avec remise

- Échantillon  $\equiv n$  tirages indépendants
- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{V}(\bar{y}) := \frac{\hat{\sigma}^2}{n}, \quad \text{avec } \hat{\sigma}^2 := \frac{1}{n-1} \sum_{i \in S} (y_i - \bar{y})^2$$

- $\hat{V}(\bar{y})$  **sans biais?**

$\Rightarrow \hat{V}(\bar{y})$  est **biaisé!**

# Échantillonnage de Poisson

- Échantillon  $\equiv$  sélectionner chaque unité avec une probabilité  $n/N$  indépendamment.
- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{\mathbb{V}}(\bar{y}) := \left(1 - \frac{n}{N}\right) \frac{1}{n^2} \sum_{i \in S} y_i^2$$

# Échantillonnage de Poisson

- Échantillon  $\equiv$  sélectionner chaque unité avec une probabilité  $n/N$  indépendamment.
- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{V}(\bar{y}) := \left(1 - \frac{n}{N}\right) \frac{1}{n^2} \sum_{i \in S} y_i^2$$

- $\hat{V}(\bar{y})$  sans biais?

# Échantillonnage de Poisson

- Échantillon  $\equiv$  sélectionner chaque unité avec une probabilité  $n/N$  indépendamment.
- Estimateur

$$\bar{y} := \frac{1}{n} \sum_{i \in S} y_i$$

- Estimateur de variance

$$\hat{V}(\bar{y}) := \left(1 - \frac{n}{N}\right) \frac{1}{n^2} \sum_{i \in S} y_i^2$$

- $\hat{V}(\bar{y})$  sans biais?

$\Rightarrow \hat{V}(\bar{y})$  est **biaisé!**

# Échantillonnage aléatoire simple sans remise

- Ignorons la “*Correction pour population finie*”:

⇒ Estimateur de variance

$$\hat{V}(\bar{y}) := \frac{\hat{\sigma}^2}{n}, \quad \text{c'est-à-dire, même estimateur qu'avec remise}$$

- **Is  $\hat{V}(\bar{y})$  sans biais?**

# Échantillonnage aléatoire simple sans remise

- Ignorons la “*Correction pour population finie*”:

⇒ Estimateur de variance

$$\hat{V}(\bar{y}) := \frac{\hat{\sigma}^2}{n}, \quad \text{c'est-à-dire, même estimateur qu'avec remise}$$

- **Is  $\hat{V}(\bar{y})$  sans biais?**

⇒  $\hat{V}(\bar{y})$  est **sans biais!**

# Échantillonnage aléatoire simple sans remise

- Ignorons la “*Correction pour population finie*” :

⇒ Estimateur de variance

$$\hat{V}(\bar{y}) := \frac{\hat{\sigma}^2}{n}, \quad \text{c'est-à-dire, même estimateur qu'avec remise}$$

- **Is  $\hat{V}(\bar{y})$  sans biais?**

⇒  $\hat{V}(\bar{y})$  est **sans biais!**

- Malgré que
  - ▶ Échantillonnage sans remise
  - ▶ Nous avons une fraction d'échantillonnage importante ( $n/N = 0.2$ )
  - ▶ La population est finie ( $N = 1000$ ). Cependant, pas de “*correction de population finie*” !

Je reviendrai sur ces exemples plus tard.

Considérons une “*approche inconditionnelle*”.



# Approche inconditionnelle

- Une **distribution**  $\xi$  qui génère les données de population:

$$\mathbf{Y} := (\mathbf{y}_1, \dots, \mathbf{y}_i, \dots, \mathbf{y}_N)^\top \sim \text{distribution } \xi$$

- $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  ne sont pas nécessairement i.i.d.
- On pourrait qualifier  $\xi$  de “*modèle*” au sens général.
- Le **plan d'échantillonnage**  $P(s)$  spécifie la distribution de l'échantillon  $S$ , qui pourrait être informatif
- **Deux processus aléatoires:**

Une distribution  $\xi$  qui génère  $\mathbf{Y}$

Une sélection aléatoire d'un échantillon  $S \subset U = \{1, \dots, N\}$

- Approche hybride basée sur l'échantillonnage et un “*modèle*”

# Approche inconditionnelle

- **Deux paramètres cibles possibles :**

- ▶ spécifié par la distribution  $\xi$ , par exemple  $\theta_0$  est l'espérance de la distribution  $\xi$

- ▶  $\theta_N$  une fonction de  $\mathbf{Y}$ , par exemple  $\theta_N = \frac{1}{N} \sum_{i \in \mathcal{U}} y_i$

- **Estimation ponctuelle:** Un estimateur  $\hat{\theta}$  sans biais (sous le plan d'échantillonnage) de  $\theta_N$  est généralement sans biais pour  $\theta_0$ :

$$\mathbb{E}_{\xi} \mathbb{E}_P[\hat{\theta} \mid \mathbf{Y}] = \mathbb{E}_{\xi}[\theta_N] = \theta_0$$

- **Pour l'estimation de variance**, la distinction entre  $\theta_0$  et  $\theta_N$  est cruciale

## Approche inconditionnelle: $\theta_N$ est la cible

- Soit  $\hat{\theta}$  un estimateur sans biais (sous le plan d'échantillonnage) de  $\theta_N$
- La **variance inconditionnelle**:

$$\begin{aligned}\mathbb{V}(\hat{\theta} - \theta_N) &= \mathbb{E}_{\xi} \mathbb{V}_P(\hat{\theta} - \theta_N \mid \mathbf{Y}) + \mathbb{V}_{\xi} \mathbb{E}_P(\hat{\theta} - \theta_N \mid \mathbf{Y}) \\ &= \mathbb{E}_{\xi} \mathbb{V}_P(\hat{\theta} - \theta_N \mid \mathbf{Y}) + 0 \\ &= \mathbb{E}_{\xi} \mathbb{V}_P(\hat{\theta} \mid \mathbf{Y})\end{aligned}$$

- **Tout estimateurs de variance sans biais  $\hat{\mathbb{V}}_P(\hat{\theta} \mid \mathbf{Y})$  basé sur le plan d'échantillonnage, est un estimateur sans biais de la variance inconditionnelle**

$$\mathbb{E}_{\xi} \mathbb{E}_P[\hat{\mathbb{V}}_P(\hat{\theta} \mid \mathbf{Y}) \mid \mathbf{Y}] = \mathbb{E}_{\xi} \mathbb{V}_P(\hat{\theta} \mid \mathbf{Y}) = \mathbb{V}(\hat{\theta} - \theta_N)$$

## Approche inconditionnelle: $\theta_N$ est la cible

- Une approche basée sur le plan d'échantillonnage est valable
- On constate que la **distribution/modèle  $\xi$  est ignorable**, puisque  $\xi$  ne joue aucun rôle dans l'estimation de la variance.
- Notez que les  $y_i$  sont des variables aléatoires !

- ▶ La caractéristique principale des approches basées sur le plan d'échantillonnage est le fait que  $\theta_N$  est la cible
- ▶ Le fait que les données  $\mathbf{Y}$  (les  $y_i$ ) soient fixes (pas aléatoires) n'est pas la caractéristique principale des approches basées sur le plan d'échantillonnage !

- Puisque nous estimons  $\mathbb{V}_P(\hat{\theta} \mid \mathbf{Y})$ , les données  $\mathbf{Y}$  **peuvent être considérées comme fixes**, du fait du conditionnement. Cela ne signifie pas pour autant que les données  $\mathbf{Y}$  sont fixes (et non aléatoires) !

## Exemple: quantiles

- L'estimation de la variance d'un quantile  $\alpha$  implique une linéarisation

$$\widehat{V}_P(\widehat{Y}_\alpha) \approx \widehat{V}_P\left(\frac{1}{N} \sum_{i \in S} \frac{\widehat{z}_i}{\pi_i}\right),$$

$$\text{où } \widehat{z}_i = \frac{1}{N \widehat{f}(\widehat{Y}_\alpha)} \left\{ \delta(y_i \leq \widehat{Y}_\alpha - \alpha) \right\}$$

- $\widehat{f}(\cdot)$  est la **densité** de la distribution de  $y_i$ .

Densité  $\implies$  Les  $y_i$  sont aléatoires et non constantes.

# Approche inconditionnelle: échantillonnage non-informatif

- Échantillonnage non-informatif:  $\mathbf{Y} \perp\!\!\!\perp S$ , c'est à dire  $\xi$  et  $P(s)$  sont des processus indépendants:
- Soit  $\hat{\theta}$  un estimateur sans biais sous le modèle de  $\tilde{\theta}$ , qui pourrait être  $\theta_N$  ou  $\theta_0$ .
- $\xi$  et  $P$  peuvent être inversés:

$$\begin{aligned}\mathbb{V}(\hat{\theta} - \tilde{\theta}) &= \mathbb{E}_P \mathbb{V}_{\xi}(\hat{\theta} - \tilde{\theta} \mid S) + \mathbb{V}_P \mathbb{E}_{\xi}(\hat{\theta} - \tilde{\theta} \mid S) \\ &= \mathbb{E}_P \mathbb{V}_{\xi}(\hat{\theta} - \tilde{\theta} \mid S) + 0\end{aligned}$$

- **Tout estimateurs de variance sans biais  $\hat{\mathbb{V}}_{\xi}(\hat{\theta} - \tilde{\theta} \mid S)$  basé sur un modèle, est un estimateur sans biais de la variance inconditionnelle**

$$\mathbb{E}_P \mathbb{E}_{\xi}[\hat{\mathbb{V}}_{\xi}(\hat{\theta} - \tilde{\theta} \mid S) \mid \mathbf{Y}] = \mathbb{E}_P \mathbb{V}_{\xi}(\hat{\theta} - \tilde{\theta} \mid \mathbf{Y}) = \mathbb{V}(\hat{\theta} - \tilde{\theta})$$

# Approche inconditionnelle: échantillonnage non-informatif

- Une approche basée sur un modèle est valable
- On constate que **le plan d'échantillonnage  $P(s)$  est ignorable**, car  $P(s)$  n'intervient pas dans l'estimation de la variance.
- Notez que l'échantillon reste une variable aléatoire !

- ▶ La **caractéristique principale des approches basées sur des modèles est l'échantillonnage non informatif**.
- ▶ Le fait que les données  $S$  ne soient pas aléatoires n'est pas la caractéristique clef des approches basées sur des modèles !

- Puisque nous estimons  $\mathbb{V}_\xi(\hat{\theta} - \tilde{\theta} \mid S)$ , l'échantillon  $S$  peut être **considéré comme fixe**, du fait du conditionnement. Cela ne signifie pas pour autant que l'échantillon  $S$  est fixes (et non aléatoires) !

# Approche inconditionnelle: échantillonnage non-informatif

- $\xi$  peut impliquer un modèle de régression pour les données  $\mathbf{Y}$  afin d'estimer  $\mathbb{V}_{\xi}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}} \mid S)$ , par exemple pour l'estimation de petits domaines ; mais le modèle régression n'est pas la caractéristique principale.
- La caractéristique principale de l'approche basée sur un modèle n'est pas la modélisation par régression, mais l'échantillonnage non informatif.
- Les approches basées sur des modèles devraient être appelées “*échantillonnage non informatives*” pour mettre en évidence la caractéristique clef



# Approche inconditionnelle: échantillonnage non-informatif

- Sous un échantillonnage non informatif, nous pouvons estimer  $\theta_N$  (prédiction) ou  $\theta_0$ .
- Lorsque  $\theta_0$  est la cible, cela ne signifie pas que nous devons envisager des approches basées sur un modèle (non informatives) !

# Approche inconditionnelle: échantillonnage informatif

- **Nous ne pouvons pas échanger  $\xi$  et  $P$** , c'est-à-dire que nous utilisons  $\mathbb{V} \equiv \mathbb{E}_\xi \mathbb{V}_P + \mathbb{E}_P \mathbb{V}_\xi$ , et pas  $\mathbb{E}_P \mathbb{V}_\xi + \mathbb{E}_P \mathbb{V}_\xi$
- **Lorsque la cible est  $\theta_N$** :  $\implies$  approches basées sur le plan d'échantillonnage

$$\begin{aligned}\mathbb{V}(\hat{\theta} - \theta_N) &= \mathbb{E}_\xi \mathbb{V}_P(\hat{\theta} - \theta_N \mid \mathbf{Y}) + \mathbb{V}_\xi \mathbb{E}_P(\hat{\theta} - \theta_N \mid \mathbf{Y}) \\ &= \mathbb{E}_\xi \mathbb{V}_P(\hat{\theta} - \theta_N \mid \mathbf{Y}) + 0\end{aligned}$$

Problème résolu lorsque la cible est  $\theta_N$ , inutile d'en discuter davantage.

- Cibler  $\theta_N$  est plus approprié lorsque l'on s'intéresse à une inférence descriptive.
- Inférence analytique : que signifie  $\theta_N$  ?  $\theta_0$  est plus logique.

# Approche inconditionnelle: échantillonnage informatif

- **Lorsque  $\theta_0$  est la cible:**

$$\mathbb{V}(\hat{\theta} - \theta_0) = \mathbb{E}_{\xi} \mathbb{V}_P(\hat{\theta} - \theta_0 \mid \mathbf{Y}) + \mathbb{V}_{\xi} \mathbb{E}_P(\hat{\theta} - \theta_0 \mid \mathbf{Y})$$

Maintenant  $\mathbb{V}_{\xi} \mathbb{E}_P(\hat{\theta} - \theta_0 \mid \mathbf{Y}) \neq 0$ ! Les estimateurs de variance traditionnels basés sur le plan d'échantillonnage sont biaisés !

Contribution : Estimation de la variance sous échantillonnage informatif, lorsque  $\theta_0$  est la cible.

- Il existe de nombreuses situations où la cible devrait être  $\theta_0$ :
  - ▶  $\theta_0$  est un paramètre de régression (inférence analytique)
  - ▶ Estimation de petits domaines
  - ▶ Estimation de quantiles?
  - ▶ Lorsque  $\theta_0 = \mathbb{E}_{\xi} \mathbb{E}_P(\hat{\theta})$ . Indice des prix à la consommation?

# Exemples de l'introduction

- $y_i \sim N(\mu, \sigma)$
- La cible était  $\theta_0 = \mu$ , et non  $\theta_N = \frac{1}{N} \sum_{i \in \mathcal{U}} y_i$

$\implies$  l'estimateur de **variance sous le plan d'échantillonnage sont biaisés**

- Sous un **échantillonnage aléatoire simple sans remise**:

$$\hat{\theta} := \frac{1}{n} \sum_{i \in S} y_i \quad \text{estimateur sans bias de } \theta_0 = \mu$$

$$\hat{\mathbb{V}}(\hat{\theta}) := \frac{\hat{\sigma}^2}{n}, \quad \text{estimateur **sans biais** de } \mathbb{V}(\hat{\theta})$$

$$\mathbb{V}(\hat{\theta}) := \mathbb{E}_{\xi} [\mathbb{V}_P(\hat{\theta} \mid \mathbf{Y})] + \mathbb{V}_{\xi} [\mathbb{E}_P(\hat{\theta} \mid \mathbf{Y})].$$

- **$\mathbb{V}(\hat{\theta})$  a 2 termes et  $\hat{\mathbb{V}}(\hat{\theta})$  a 1 terme !?!?**

# Échantillonnage aléatoire simple sans remise

- $y_i \sim N(\mu, \sigma)$  i.i.d
- $\mathbb{V}_P(\hat{\theta} \mid \mathbf{Y}) = \left(1 - \frac{n}{N}\right) \frac{S^2}{n}$ , où  $S^2 = \frac{1}{N-1} \sum_{i \in \mathcal{U}} (y_i - \theta_N)^2$
- $\mathbb{E}_\xi \mathbb{V}_P(\hat{\theta} \mid \mathbf{Y}) = \left(1 - \frac{n}{N}\right) \frac{\sigma^2}{n} = \frac{\sigma^2}{n} - \frac{\sigma^2}{N}$ , car  $\theta_0 = \mu$  est la cible
- $\mathbb{V}_\xi \mathbb{E}_P(\hat{\theta} \mid \mathbf{Y}) = \mathbb{V}_\xi(\theta_N) = \frac{\sigma^2}{N}$

$$\begin{aligned}\mathbb{V}(\hat{\theta}) &:= \mathbb{E}_\xi[\mathbb{V}_P(\hat{\theta} \mid \mathbf{Y})] + \mathbb{V}_\xi[\mathbb{E}_P(\hat{\theta} \mid \mathbf{Y})] \\ &= \frac{\sigma^2}{n} - \frac{\sigma^2}{N} + \frac{\sigma^2}{N} = \frac{\sigma^2}{n}\end{aligned}$$

- La **correction pour population finie** disparaît à cause de  $\mathbb{V}_\xi \mathbb{E}_P(\hat{\theta} \mid \mathbf{Y})$

# Échantillonnage aléatoire simple sans remise

- Le terme  $(1 - n/N)$  est la principale raison du biais de l'estimateur de variance basé sur le plan
- Le biais est  $-\sigma^2 N^{-1} = -\mathbb{V}_\xi \mathbb{E}_P(\hat{\theta})$  à cause du terme  $(1 - n/N)$
- **L'interprétation correcte de  $(1 - n/N)$  est**

Le rôle du terme  $(1 - n/N)$  est de réduire la variance de  $\mathbb{V}_\xi \mathbb{E}_P(\hat{\theta})$ , afin de compenser le fait que la quantité d'intérêt est  $\theta_N$ , plutôt que  $\theta_0$

## Correction de population finie?

- Cette correction n'est pas due au fait que  $N$  soit fini, car  $(1 - n/N)$  doit être remplacé par 1, lorsque  $\theta_0$  est le paramètre cible, même lorsque  $N$  est fini.
- Cette correction est nécessaire lorsque  $\theta_N$  est le paramètre cible.
- Qualifier  $(1 - n/N)$  de **“correction de population finie” est trompeur**, car cela n'a rien à voir avec le fait que  $N$  soit fini.

## Correction de population finie?

- Qualifier  $(1 - n/N)$  de **“correction de population finie”** est **trompeur**

Nous devrions l'appeler *“Correction de population fixe”*. Le terme *“fixe”* est utilisé pour souligner que la correction doit être appliquée lorsque  $\theta_N$  est la cible et que nous devons utiliser une variance conditionnelle étant donné  $\mathbf{Y}$ , en traitant les données  $\mathbf{Y}$  comme fixes.

- *“Correction pour grande fraction d'échantillonnage”* est aussi adéquat
- *“L'inférence en population finie”* doit être comprise comme le fait que  $\theta_N$  est la cible et que  $\mathbf{Y}$  est considéré comme fixe (conditionnement sur  $\mathbf{Y}$ ). L'expression *“inférence en population fixe”* est plus appropriée



# Généralisation

- Les principes développés jusqu'à présent s'appliquent à des **plans d'échantillonnage plus élaborés**, tels que
  - ▶ Échantillonnage à probabilité inégale sans remise
  - ▶ Échantillonnage stratifiée
  - ▶ Échantillonnage à plusieurs degré
  - ▶ Échantillonnage systématique ordonné
- Nous supposons un **échantillonnage informatif**
- **Inférence analytique**, car nous nous intéressons à  $\theta_0$
- Extension aux **variables auxiliaires** : voir l'article
- Échantillonnage de Poisson informatif : voir l'article

# Échantillonnage à probabilité inégale sans remise

- $(y_i, \pi_i)^\top \sim$  distribution bivariée  $\xi$
- $\theta_0$  est la solution (supposée unique) de

$$\mathbb{E}_\xi[g(y, \theta)] = 0 \quad \text{iff } \theta = \theta_0.$$

- $\hat{\theta}$  solution (supposée unique) de l'équivalent empirique

$$\hat{\Gamma}(\theta) := \sum_{i \in S} \frac{1}{\pi_i} g(y_i, \theta) = 0.$$

- Soit  $\varrho_i := g(y_i, \theta_0)$ .
- $\pi_i$  et  $\varrho_i$  **peuvent être dépendant** (échantillonnage informatif).

# Hypothèses

- $N/n < \infty$  est borné, lorsque  $n \rightarrow \infty$
- $d := \sum_{i \in \mathcal{U}} \pi_i (1 - \pi_i) \rightarrow \infty$
- $\pi_{ij} = \pi_i \left[ \pi_j \left\{ 1 - [1 + o(1)](1 - \pi_i)(1 - \pi_j)d^{-1} \right\} \right]^{1 - \delta_{ij}}$ ,  
où  $\delta_{ij}$  est le delta de Kronecker

Entropie élevée. Expression la plus courante pour  $\pi_{ij}$  (Hájek, 1964)

- **Indépendance**  $\varrho_i \perp\!\!\!\perp \varrho_j$  et  $\pi_i \perp\!\!\!\perp \pi_j$ , pour  $i \neq j$ , où  $\varrho_i := g(y_i, \theta_0)$ .

**Indépendance entre les  $\varrho_i$ , pas entre les  $y_i$  !**

- $\frac{1}{N} \sum_{i \in \mathcal{U}} \mathbb{V}_\xi(\varrho_i) < \infty$ .
- $\pi_{ij} \geq \pi_i \pi_j$ , a.s.  $\forall i, j$ .
- Hypothèses asymptotiques standard dans un contexte de population finie

## Estimation de la variance

- Considérons **l'estimateur de variance de Hansen & Hurwitz (1943)**:

$$\widehat{\mathbb{V}}[\widehat{\Gamma}(\theta)] = \sum_{i \in S} \frac{1}{\pi_i^2} g(y_i, \theta)^2 - \frac{1}{n} \left\{ \sum_{i \in S} \frac{1}{\pi_i} g(y_i, \theta) \right\}^2.$$

- Asymptotiquement **sans biais**:

$$\frac{n}{N^2} \left\{ \mathbb{E}_{\xi} \mathbb{E}_P \widehat{\mathbb{V}}[\widehat{\Gamma}(\theta_0)] - \mathbb{V}[\widehat{\Gamma}(\theta_0)] \right\} = o(1)$$

- Il n'est pas nécessaire d'estimer séparément  $\mathbb{E}_{\xi} \mathbb{V}_P[\widehat{\Gamma}(\theta) \mid \mathbf{Y}]$  et  $\mathbb{V}_{\xi} \mathbb{E}_P[\widehat{\Gamma}(\theta) \mid \mathbf{Y}]$

$\implies$  Théorème de Taylor:

$$\widehat{\mathbb{V}}(\widehat{\theta}) := \left\{ \frac{\partial \widehat{\Gamma}(\widehat{\theta})}{\partial \widehat{\theta}} \right\}^{-2} \widehat{\mathbb{V}}[\widehat{\Gamma}(\widehat{\theta})],$$

# L'estimateur de variance de Hansen & Hurwitz (1943)

- L'estimateur de variance de Hansen & Hurwitz (1943) est sans biais sous le plan, lorsque l'échantillon est sélectionné avec remise et que  $\theta_N$  est la cible.
- Cependant, lorsque  $\theta_0$  est la cible, cet estimateur est **biaisé** sous un plan d'échantillonnage avec remise
- Cet estimateur n'est pas limité à l'échantillonnage avec remise
- **Il est asymptotiquement sans biais sous échantillonnage sans remise avec de grandes fractions d'échantillonnage, lorsque  $\theta_0$  est la cible**

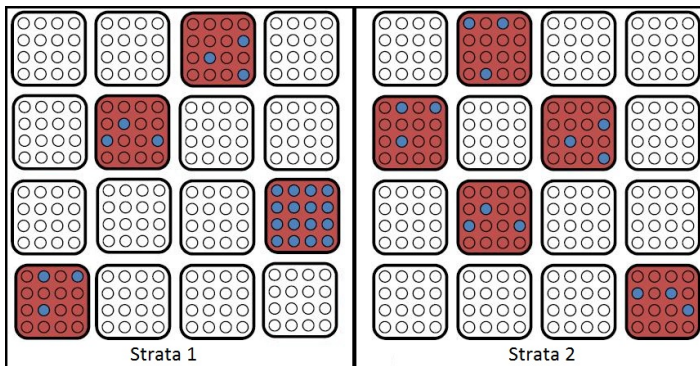
# L'estimateur de variance de Hansen & Hurwitz (1943)

- L'absence de FPC et de  $\pi_{ij}$  n'est pas due au fait que  $n/N$  est asymptotiquement négligeable, ni à une hypothèse "avec remise"
- La variance est asymptotiquement sans biais car  $\theta_0$  est la cible
- **Interprétation de  $\pi_{ij}$ :**

Les  $\pi_{ij}$  sont nécessaires lorsque  $\theta_N$  est la cible, et non parce que la population est finie

# Échantillonnage stratifié à plusieurs degrés

- Un échantillon  $S$  de  $n$  grappes est sélectionné sans remise avec des probabilités inégales  $\pi_i$
- Au sein de chaque grappe  $\mathcal{U}_i$  échantillonnée ( $i \in S$ ), un échantillon  $S_i$  de  $m_i$  unités est sélectionné



# Échantillonnage stratifié à plusieurs degrés

- $\theta_0 \in \mathbb{R}^{d_\theta}$  est défini par une “condition de moment multivariée”, c'est-à-dire

$$\mathbb{E}_\xi \{ \mathbf{g}(\mathbf{y}, \theta) \} = \mathbf{0}_{d_g}, \quad \text{iff } \theta = \theta_0,$$

- $\mathbf{g}(\mathbf{y}, \theta) \in \mathbb{R}^{d_g}$  et  $\mathbf{0}_d$  est le  $d$ -vecteur de 0.
- Le modèle peut être sur-spécifié, c'est-à-dire  $d_g \geq d_\theta$ , afin d'inclure des contraintes de calibration (voir l'article).
- L'estimateur ponctuel  $\hat{\theta}$  est la solution de

$$\hat{\Gamma}(\theta) := \sum_{i \in S} \frac{1}{\pi_i} \hat{\rho}_i(\theta) = \mathbf{0}_{d_g},$$

avec

$$\hat{\rho}_i(\theta) := \sum_{j \in S_i} \frac{1}{\pi_{j|i}} \mathbf{g}_{j|i}(\theta), \quad \mathbf{g}_{j|i}(\theta) := \mathbf{g}(\mathbf{y}_j, \theta), \quad \text{pour } j \in \mathcal{U}_i,$$



## Estimateur de variance stratifié “*grappe ultime*” de Hansen & Hurwitz (1943)

$$\widehat{\mathbb{V}}\{\widehat{\Gamma}(\theta)\} := \sum_{h=1}^H \widehat{\mathbb{V}}\{\widehat{\Gamma}_h(\theta)\}, \quad \text{avec } \widehat{\Gamma}_h(\theta) := \sum_{i \in S_h} \frac{1}{\pi_i} \widehat{\rho}_i(\theta),$$

où

$$\widehat{\mathbb{V}}\{\widehat{\Gamma}_h(\theta)\} := \sum_{i \in S_h} \frac{1}{\pi_i^2} \widehat{\rho}_i(\theta) \widehat{\rho}_i(\theta)^\top - \frac{1}{n_h} \widehat{\Gamma}_h(\theta) \widehat{\Gamma}_h(\theta)^\top,$$

- Asymptotiquement sans biais  $\implies$  Théorème de Taylor:

$$\widehat{\mathbb{V}}(\widehat{\theta}) := \left\{ \frac{\partial \widehat{\Gamma}(\theta)}{\partial \theta} \Big|_{\theta=\widehat{\theta}} \right\}^{-1} \widehat{\mathbb{V}}\{\widehat{\Gamma}(\widehat{\theta})\} \left\{ \frac{\partial \widehat{\Gamma}(\theta)}{\partial \theta} \Big|_{\theta=\widehat{\theta}} \right\}^{-1\top}.$$

- Inutile d'estimer les variances a l'intérieur des grappes** (Gustave (Chevalier & Richer, 2023)), même lorsque  $n/N$  est grand

# Hypothèse clef pour l'absence de biais asymptotique

- **Indépendance entre les grappes**, c'est-à-dire

$$\mathbf{g}_{j|i}(\boldsymbol{\theta}_0) \perp\!\!\!\perp \mathbf{g}_{\ell|k}(\boldsymbol{\theta}_0) \text{ et } \pi_i \perp\!\!\!\perp \pi_k, \forall j \in \mathcal{U}_i, \ell \in \mathcal{U}_k \text{ et } i \neq k.$$

avec  $\mathbf{g}_{j|i}(\boldsymbol{\theta}) := \mathbf{g}(\mathbf{y}_j, \boldsymbol{\theta})$

- Nous pourrions avoir une **dépendance au sein des grappes**
- Entropie élevée (Hájek, 1964)
- Il n'est pas nécessaire d'utiliser des effets aléatoires ni d'estimer la variance au sein des grappes

# Ordered systematic sampling

- S'il y a une tendance, alors les  $y_i$  ne peuvent pas être i.i.d.
- L'hypothèse d'entropie élevée est erronée, car certains des  $\pi_{ij}$  sont nuls
- Une solution consiste à sur-spécifier la fonction d'estimation en se basant sur un modèle de tendance, tel que :

$$y_i = \alpha_0 + \ell_i \beta_0 + e_i,$$

où  $\ell_i$  est le label à l'intérieur de  $U$  de l'unité  $i \in S$

# Échantillonnage systématique ordonné

- Une **fonction d'estimation sur-spécifiée** peut être utilisée, c'est-à-dire

$$\mathbf{g}(y_i, \boldsymbol{\theta}) = \{\varepsilon(y_i, \boldsymbol{\theta}) - \mu + \alpha + \mu_\ell \beta, \varepsilon(y_i, \boldsymbol{\theta}), \ell_i \varepsilon(y_i, \boldsymbol{\theta})\}^\top,$$

avec  $\varepsilon(y_i, \boldsymbol{\theta}) := y_i - \alpha - \ell_i \beta,$

$$\boldsymbol{\theta} := (\mu, \alpha, \beta)^\top,$$

$$\mu_\ell := (N + 1)/2, \quad \text{la moyenne de la population des } \ell_i,$$

$$\mu_0 := \mathbb{E}_\xi(y), \quad \text{le paramètre cible,}$$

$$(\alpha_0, \beta_0)^\top \quad \text{un paramètre de nuisance}$$

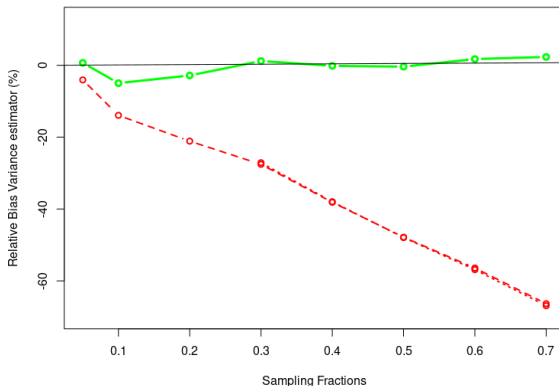
- **$\mathbf{g}(y_i, \boldsymbol{\theta}_0)$  sont i.i.d., malgré le fait que les  $y_i$  ne soient pas i.i.d.**
- L'échantillonnage systématique aléatoire et l'échantillonnage systématique ordonné sont équivalents, car i.i.d.  
 $\implies$  Entropie large  $\implies$  variance de Hansen & Hurwitz (1943)

# Numerical example

- $y_i \sim N(10, sd = 2) \Rightarrow \theta_0 = 10$
- $p(S) = \text{échantillonnage de Chao}$
- $n = 100, \quad 0.05 \leq f = \frac{n}{N} \leq 0.7$
- $0 \leq \text{cor}(y_i, \pi_i) \leq 0.9$
- Population finie  $N = n/f$

# Biais relatifs

$$\text{cor}(y_i, \pi_i) = 0.4$$

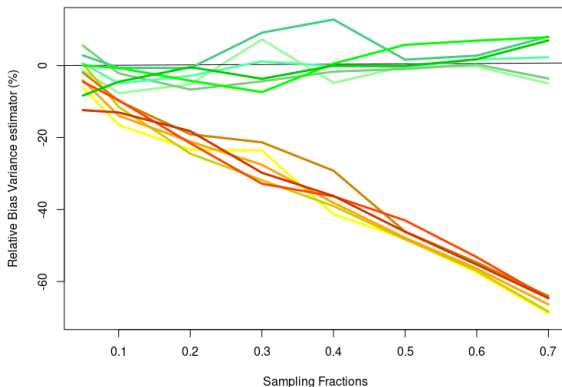


Variance de Hansen & Hurwitz (1943)

Variance de Horvitz & Thompson (1952):  $\widehat{V}_P(\widehat{\theta})$

# Biais relatifs pour différentes corrélations

La corrélation n'a aucun effet sur les biais.



Variance de Hansen & Hurwitz (1943)

Variance de Horvitz & Thompson (1952):  $\widehat{V}_P(\widehat{\theta})$

# Échantillonnage à deux degrés

- Données générées à partir d'un modèle à effets aléatoires

$$\mathcal{Y}_{j|i} = \alpha_0 + \beta_0 x_{j|i} + \epsilon_{j|i}, \quad \text{avec } j \in \mathcal{U}_i \text{ et } \alpha_0 = \beta_0 = 1, \\ \epsilon_{j|i} = u_i + e_{ji}.$$

- Supposons que nous souhaitions tester  $H_0 : \theta_0 = \tilde{\theta}$  contre  $H_a : \theta_0 \neq \tilde{\theta}$ , où  $\tilde{\theta} = (1, 1)^\top$ . Niveau  $\alpha = 5\%$
- Statistique de pivot:

$$r(\tilde{\theta}) := \hat{\Gamma}(\tilde{\theta})^\top \hat{\mathbb{V}}\{\hat{\Gamma}(\tilde{\theta})\}^{-1} \hat{\Gamma}(\tilde{\theta})$$

tend en distribution vers une distribution  $\chi^2$  avec  $d_\theta$  degrés de liberté, sous  $H_0$ .

- Plus simple que les ajustements de Rao & Scott (1987)



# Échantillonnage à deux degrés. Tailles observées (%)

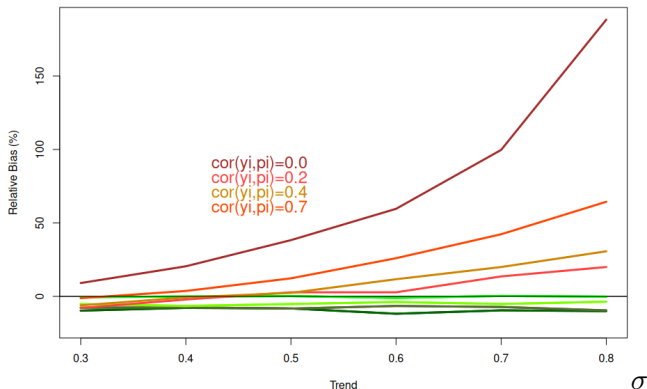
		$f = n/N$							
Statistique de test	$\rho$	0.05	0.1	0.2	0.3	0.4	0.5	0.6	0.7
Proposé	0.0	6.1	6.3	5.5	5.9	4.5	5.5	6.1	5.8
	0.2	5.6	5.3	5.1	5.2	4.4	5.9	5.5	5.4
	0.4	6.6 <sup>†</sup>	4.3	6.3	5.3	4.9	5.6	5.0	5.1
	0.6	5.1	5.5	6.2	6.3	6.0	6.1	5.0	5.1
	0.8	5.6	4.3	5.7	5.6	5.6	5.6	5.4	4.8
	0.9	4.5	4.6	4.3	6.8 <sup>†</sup>	6.2	5.6	4.9	5.5
Rao & Scott 1	0.0	5.6	6.5 <sup>†</sup>	5.2	5.5	4.4	5.5	5.6	5.0
	0.2	5.6	5.3	5.3	4.0	5.9	5.1	6.0	5.4
	0.4	6.6 <sup>†</sup>	5.9	7.1 <sup>†</sup>	7.0 <sup>†</sup>	5.5	8.3 <sup>†</sup>	8.2 <sup>†</sup>	5.7
	0.6	6.5 <sup>†</sup>	7.1 <sup>†</sup>	8.0 <sup>†</sup>	8.3 <sup>†</sup>	7.1 <sup>†</sup>	8.6 <sup>†</sup>	6.9 <sup>†</sup>	9.8 <sup>†</sup>
	0.8	8.1 <sup>†</sup>	6.4 <sup>†</sup>	8.3 <sup>†</sup>	9.2 <sup>†</sup>	7.8 <sup>†</sup>	8.1 <sup>†</sup>	9.6 <sup>†</sup>	8.0 <sup>†</sup>
	0.9	6.1	8.2 <sup>†</sup>	7.4 <sup>†</sup>	8.5 <sup>†</sup>	8.5 <sup>†</sup>	7.8 <sup>†</sup>	8.3 <sup>†</sup>	8.5 <sup>†</sup>
Rao & Scott 2	0.0	5.6	6.5 <sup>†</sup>	5.2	5.5	4.4	5.5	5.6	5.0
	0.2	5.6	5.3	3.6 <sup>†</sup>	3.2 <sup>†</sup>	5.9	5.1	6.0	5.4
	0.4	6.6 <sup>†</sup>	3.5 <sup>†</sup>	7.1 <sup>†</sup>	7.0 <sup>†</sup>	3.1 <sup>†</sup>	8.3 <sup>†</sup>	8.2 <sup>†</sup>	3.4 <sup>†</sup>
	0.6	6.5 <sup>†</sup>	7.1 <sup>†</sup>	8.0 <sup>†</sup>	8.3 <sup>†</sup>	7.1 <sup>†</sup>	8.6 <sup>†</sup>	6.9 <sup>†</sup>	9.8 <sup>†</sup>
	0.8	8.1 <sup>†</sup>	2.4 <sup>†</sup>	8.3 <sup>†</sup>	9.2 <sup>†</sup>	7.8 <sup>†</sup>	8.1 <sup>†</sup>	9.6 <sup>†</sup>	8.0 <sup>†</sup>
	0.9	6.1	8.2 <sup>†</sup>	3.2 <sup>†</sup>	8.5 <sup>†</sup>	8.5 <sup>†</sup>	7.8 <sup>†</sup>	8.3 <sup>†</sup>	8.5 <sup>†</sup>

<sup>†</sup> Taille significativement différente de 5%: p-valeur < 0.05.

# Échantillonnage systématique ordonné

- $y_i = 5 + (i - 1)(N - 1)^{-1} + \varepsilon_i$ , où  $\varepsilon_i \sim N(0, \sigma)$  i.i.d.
- C'est-à-dire  $y_i$  suit une tendance linéaire.
- La valeur différente de  $\sigma$  permet de contrôler la force de la tendance
- Les corrélations  $\rho_{y\pi}$  entre  $\pi_i$  et  $y_i$  sont de 0,2, 0,4 et 0,7.
- $n = 500$
- Population de taille  $N = \text{round}(n/f)$ , où  $f = 0.05, 0.2$  et  $0.4$ .

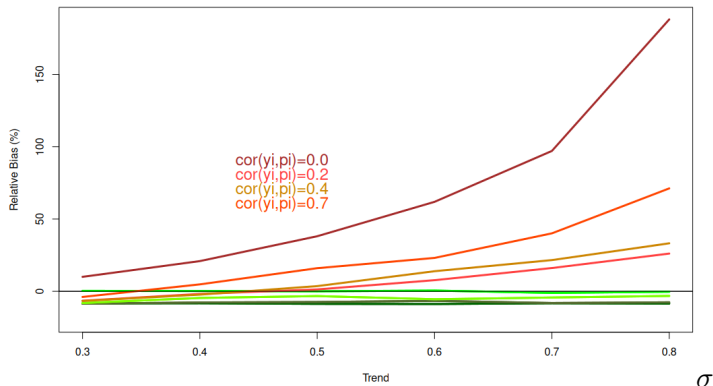
## Biais relatifs pour différentes corrélations: $n/N = 0.05$



Variance de Hansen & Hurwitz (1943) surspécifié

Variance de Hansen & Hurwitz (1943)

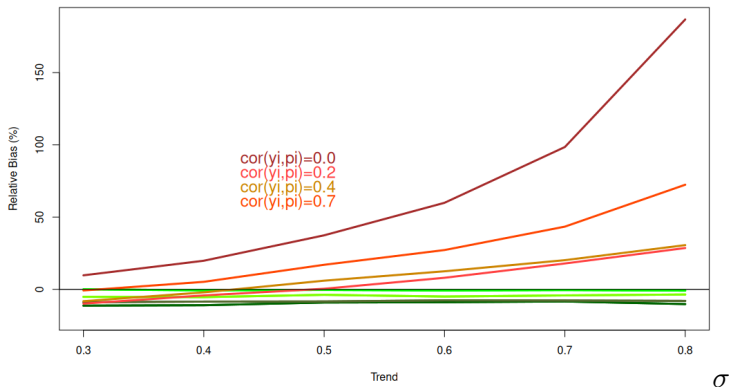
## Biais relatifs pour différentes corrélations: $n/N = 0.2$



Variance de Hansen & Hurwitz (1943) surspécifié

Variance de Hansen & Hurwitz (1943)

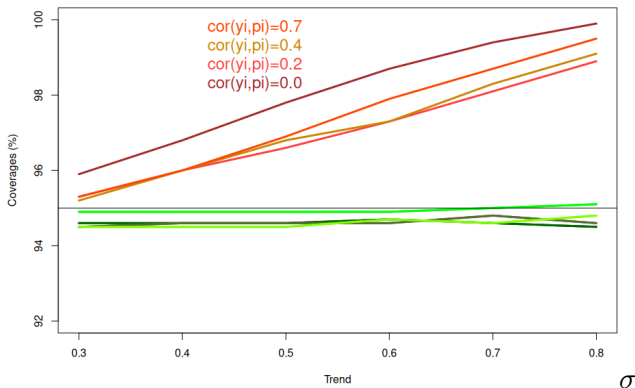
## Biais relatifs pour différentes corrélations: $n/N = 0.4$



Variance de Hansen & Hurwitz (1943) surspécifié

Variance de Hansen & Hurwitz (1943)

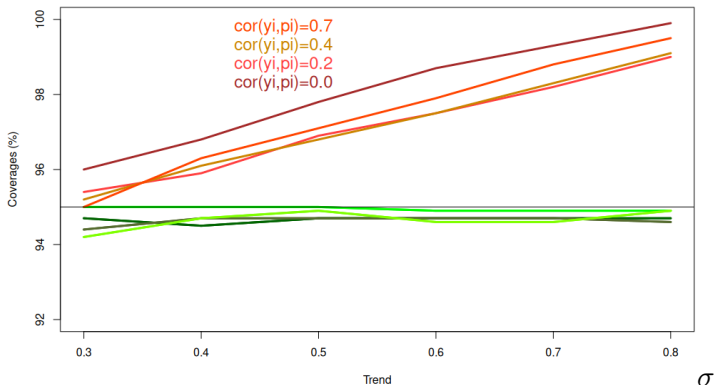
## Couvertures pour différentes corrélations: $n/N = 0.05$



Variance de Hansen & Hurwitz (1943) surspécifié

Variance de Hansen & Hurwitz (1943)

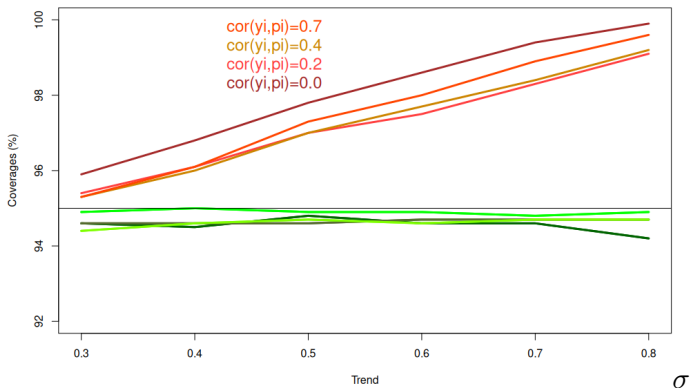
## Couvertures pour différentes corrélations: $n/N = 0.2$



Variance de Hansen & Hurwitz (1943) surspécifié

Variance de Hansen & Hurwitz (1943)

## Couvertures pour différentes corrélations: $n/N = 0.4$



Variance de Hansen & Hurwitz (1943) surspécifié  
Variance de Hansen & Hurwitz (1943)



# Conclusion

- Le choix entre les approches basées sur le plan et les approches basées sur les modèles ne doit pas être guidé par
  - ▶ Le fait que l'on suppose les  $y_i$  fixes ou aléatoires
  - ▶ Le fait que l'on traite l'échantillon comme fixe ou aléatoire
- Ces questions sont hors de propos

# Conclusion

Deux questions clefs:

- ▶ L'échantillonnage est-il informatif ?
- ▶ Quel est le paramètre cible ?  $\theta_N$  ou  $\theta_0$  ?

# Conclusion

Deux questions clefs:

- ▶ L'échantillonnage est-il informatif ?
- ▶ Quel est le paramètre cible ?  $\theta_N$  ou  $\theta_0$  ?

- **Si l'échantillonnage n'est pas informatif**: variance basée sur le plan, si la cible est  $\theta_N$  ou variance basée sur le modèle (avec des hypothèses de distribution supplémentaires) si la cible est  $\theta_N$  ou  $\theta_0$

# Conclusion

Deux questions clefs:

- ▶ L'échantillonnage est-il informatif ?
- ▶ Quel est le paramètre cible ?  $\theta_N$  ou  $\theta_0$  ?

- **Si l'échantillonnage n'est pas informatif**: variance basée sur le plan, si la cible est  $\theta_N$  ou variance basée sur le modèle (avec des hypothèses de distribution supplémentaires) si la cible est  $\theta_N$  ou  $\theta_0$
- **Si l'échantillonnage est informatif et la cible est  $\theta_N$**   
⇒ Variance traditionnels basés sur le plan (Gustave (Chevalier & Richer, 2023))

# Conclusion

Deux questions clefs:

- ▶ L'échantillonnage est-il informatif ?
- ▶ Quel est le paramètre cible ?  $\theta_N$  ou  $\theta_0$  ?

- **Si l'échantillonnage n'est pas informatif**: variance basée sur le plan, si la cible est  $\theta_N$  ou variance basée sur le modèle (avec des hypothèses de distribution supplémentaires) si la cible est  $\theta_N$  ou  $\theta_0$
- **Si l'échantillonnage est informatif et la cible est  $\theta_N$**   
⇒ Variance traditionnels basés sur le plan (Gustave (Chevalier & Richer, 2023))
- **Si l'échantillonnage est informatif et la cible est  $\theta_0$**   
⇒ **Variance de Hansen & Hurwitz (1943)**  
⇒ Variance basé sur le plan n'a aucun sens et est beaucoup plus compliquée

# Conclusion

Le rôle du terme  $(1 - n/N)$  est de réduire la variance de  $\mathbb{V}_\xi \mathbb{E}_P(\hat{\theta})$ , afin de compenser le fait que la quantité d'intérêt est  $\theta_N$ , plutôt que  $\theta_0$

- Qualifier  $(1 - n/N)$  de “*correction pour population finie*” est trompeur, car cela n'a rien à voir avec le fait que la population est finie.
- Nous devrions l'appeler “**correction de population fixe**”. Le mot “*Fixe*” est utilisé pour souligner que la correction doit être utilisée lorsque  $\theta_N$  est la cible et que nous devons utiliser une variance conditionnelle étant donné  $\mathbf{Y}$ , en traitant les données  $\mathbf{Y}$  et  $\theta_N$  comme étant fixes
- Le fait que  $N$  soit fini est sans importance

# Autres sujets abordés dans l'article

- Variables auxiliaires
- Variance de l'estimateur de régression tenant compte de l'estimation du paramètre de régression
- Échantillonnage informatif de Poisson

# References



**Berger, Y. G. (2025) Unconditional variance estimation under complex surveys. *International Statistical Review*, 26pp.**

<https://onlinelibrary.wiley.com/doi/10.1111/insr.70001>.

Chevalier, M. & Richer, J.-S. (2023) *gustave: User-Oriented Statistical Toolkit for Analytical Variance Estimation*. URL <https://CRAN.R-project.org/package=gustave>. R package version 0.3.0.

Hájek, J. (1964) Asymptotic theory of rejective sampling with varying probabilities from a finite population. *The Annals of Mathematical Statistics*, **35**, 1491–1523. URL <https://doi.org/10.1214/aoms/1177700375>.

Hansen, M. H. & Hurwitz, W. N. (1943) On the theory of sampling from finite populations. *The Annals of Mathematical Statistics*, **14**, pp. 333–362.

Horvitz, D. G. & Thompson, D. J. (1952) A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, **47**, 663–685. URL <https://doi.org/10.1080/01621459.1952.10483446>.

Rao, J. N. K. & Scott, A. J. (1987) On simple adjustments to chi-square tests with sample survey data. *The Annals of Statistics*, **15**, 385–397.