

Sondage aléatoire simple à probabilités égales

Myriam Maumy¹

¹IRMA, Université Louis Pasteur
Strasbourg, France

Master 1ère Année 06-03-2006

Ce chapitre s'appuie essentiellement sur trois livres :

- « Éléments de statistiques »,
de Jean-Jacques Dreesbeke,
Université de Bruxelles, 2001.
- « Les techniques de sondage »
de Pascal Ardilly,
Éditions Technip, 2006.
- « Exercices corrigés de méthodes de sondage »
de Pascal Ardilly et de Yves Tillé,
Éditions Ellipses, 2003.

Sommaire

- 1 Introduction
- 2 Sondage aléatoire simple à probabilités égales avec remise
- 3 Sondage aléatoire simple à probabilités égales sans remise
- 4 Comparaison des prélèvements PEAR et PESR

Définition

Un sondage aléatoire est simple (SAS) si tous les échantillons de taille n fixée a priori, prélevés au sein d'une population U d'effectif N , sont réalisables avec la même probabilité.

Remarque

Dans ce cas, les individus de la population U ont tous la même probabilité d'être choisis pour faire partie de l'échantillon S : leur **probabilité d'inclusion** est une constante.

Remarque

Qu'appelle-t-on une **probabilité d'inclusion** ?

Réponse

Vous pouvez trouver une définition p: 51 dans le livre de Ardilly.

Remarque

Si l'on reprend le choix d'*une seule observation*, chaque individu de la population U a une probabilité égale à $1/N$ d'être prélevé dans la population U afin de constituer l'échantillon \mathcal{S} .

Il y a deux méthodes pour sélectionner des individus pour constituer un échantillon \mathcal{S} .

La première méthode

Elle consiste à replacer chaque valeur observée dans la population U avant le tirage suivant et cela n fois de suite.

⇒ Prélèvement **avec remise**. On parle d'un sondage à probabilités égales avec remise (PEAR).

La deuxième méthode

Elle consiste à ne pas remettre l'individu dans la population U à chaque tirage.

⇒ Prélèvement **sans remise**. On parle d'un sondage à probabilités égales sans remise (PESR).

Sommaire

- 1 Introduction
- 2 Sondage aléatoire simple à probabilités égales avec remise**
- 3 Sondage aléatoire simple à probabilités égales sans remise
- 4 Comparaison des prélèvements PEAR et PESR

Propriété

Dans ce cas, il y a N^n échantillons S possibles.

Remarque

Un même individu peut-être sélectionné plusieurs fois !

Remarque

À chaque tirage, on a toujours *la même population U .*

Chaque valeur observée est prise *indépendamment* des autres.

Propriété

L'échantillon S est alors considéré comme une suite de variables aléatoires **indépendantes et équadistribuées** $\{X_1, \dots, X_n\}$, où X_i est la valeur observée pour le i -ème individu sélectionné, telles que

$$\forall i = 1, \dots, n \quad \mathbb{E}[X_i] = \mu \quad \text{et} \quad \text{Var}[X_i] = \sigma^2,$$

où μ est la moyenne de la population U et σ^2 la variance de la population U .

Définition

Un estimateur classique de la moyenne μ d'une population U se définit par

$$\hat{\mu} = \frac{1}{n} \sum_{i \in S} x_i.$$

Propriété

On montre, par calculs directs, que

$$\mathbb{E}[\hat{\mu}] = \mu \quad \text{et} \quad \text{Var}[\hat{\mu}] = \frac{\sigma^2}{n}.$$

Remarque

L'avant dernière égalité de la dernière propriété implique que $\hat{\mu}$ est un estimateur sans biais de la moyenne μ de la population U .

Remarque

Dans l'expression de la variance de $\hat{\mu}$, on remarque que le terme de la variance σ^2 de la population U intervient. Or, dans la plupart des cas, on ne connaît pas la variance σ^2 de la population U . On sera donc amené à construire un estimateur de la variance de $\hat{\mu}$.

Définition

Un estimateur de la variance de $\hat{\mu}$ se définit par

$$\widehat{\text{Var}}[\hat{\mu}] = \frac{s_c^2}{n},$$

où s_c^2 désigne la variance corrigée de l'échantillon S .

Propriété

On montre, par calcul direct, que

$$\mathbb{E} \left[\widehat{\text{Var}}[\hat{\mu}] \right] = \frac{\sigma^2}{n}.$$

Remarque

On rappelle que la variance corrigée s_c^2 de l'échantillon S se définit par

$$s_c^2 = \frac{1}{n-1} \sum_{i \in S} (x_i - \hat{\mu})^2$$

et que s_c^2 est un estimateur sans biais de la variance σ^2 de la population U .

Remarque

De cette dernière propriété, on en déduit que $\frac{s_c^2}{n}$ est un estimateur sans biais de la variance de $\hat{\mu}$.

Définition

Un estimateur classique du total T d'une population U se définit par

$$\hat{T} = N\hat{\mu} = \frac{N}{n} \sum_{i \in \mathcal{S}} x_i.$$

Propriété

On montre, par calculs directs, que

$$\mathbb{E} \left[\hat{T} \right] = T \quad \text{et} \quad \text{Var} \left[\hat{T} \right] = N^2 \frac{\sigma^2}{n}.$$

Remarque

L'avant dernière égalité de la dernière propriété implique que \hat{T} est un estimateur sans biais du total T de la population U .

Remarque

Dans l'expression de la variance de \hat{T} , on remarque que le terme de la variance σ^2 de la population U intervient. Or, dans la plupart des cas, on ne connaît pas la variance σ^2 de la population U . On sera donc amené à construire un estimateur de la variance de \hat{T} .

Définition

Un estimateur de la variance de \hat{T} se définit par

$$\widehat{\text{Var}} \left[\hat{T} \right] = N^2 \frac{s_c^2}{n},$$

où s_c^2 désigne la variance corrigée de l'échantillon S .

Propriété

On montre, par calcul direct, que

$$\mathbb{E} \left[\widehat{\text{Var}} \left[\hat{T} \right] \right] = N^2 \frac{\sigma^2}{n}.$$

Remarque

De cette dernière propriété, on en déduit que $N^2 \frac{s_c^2}{n}$ est un estimateur sans biais de la variance de \hat{T} .

On rappelle que :

Définition

Un estimateur classique de la variance σ^2 d'une population U se définit par

$$s_c^2 = \frac{1}{n-1} \sum_{i \in S} (x_i - \hat{\mu})^2.$$

Propriété

On montre, par calculs directs, que

$$\mathbb{E} [s_c^2] = \sigma^2 \quad \text{et} \quad \text{Var} [s_c^2] = \frac{1}{n(n-1)} \left[(n-1)\mu_4 - (n-3)\sigma^4 \right].$$

Remarque

L'avant dernière égalité de la dernière propriété implique que s_c^2 est un estimateur sans biais de la variance σ^2 de la population U .

Remarque

Dans l'expression de la variance de s_c^2 , on remarque que le terme σ^4 , qui est le carré de la variance de la population U , intervient ainsi que le moment d'ordre 4, μ_4 . Or, dans la plupart des cas, on ne connaît ni σ^4 , ni μ_4 . On sera donc amené à construire un estimateur de la variance corrigée s_c^2 , si besoin est.

Le **prélèvement avec remise** est susceptible de fournir plusieurs fois un individu de la population. Deux situations se présentent.

Les n tirages fournissent n individus distincts.

Dans ce cas, \mathcal{S} correspond à un sous-ensemble de U de taille n .

Les définitions de $\hat{\mu}$, \hat{T} et s_c^2 sont équivalentes si on renumérote les individus de la population U de telle sorte que

$$\mathcal{S} = \{1, \dots, n\}.$$

Les n tirages fournissent m individus, où $m < n$.

Dans ce cas, deux comportements sont à envisager.

- Le premier consiste à prendre en compte les observations autant de fois qu'elles ont été recueillies.
- Le second consiste de prendre la moyenne des m valeurs distinctes observées dont l'ensemble est désigné par \mathcal{S}_m :

$$\hat{\mu}_m = \sum_{k \in \mathcal{S}_m} x_k.$$

Il est clair que dans ce cas, la taille de n de l'échantillon n'est plus une constante mais devient elle-même une v.a., fonction du processus de prélèvement.

On montre que, en moyenne, $\hat{\mu}_m$ est encore égal à μ .

Sommaire

- 1 Introduction
- 2 Sondage aléatoire simple à probabilités égales avec remise
- 3 Sondage aléatoire simple à probabilités égales sans remise**
- 4 Comparaison des prélèvements PEAR et PESR

Définition

*Un sondage aléatoire simple est **sans remise** si l'observation prélevée au i -ème tirage n'est pas replacée dans la population avant les prélèvements suivants. On parle alors d'un sondage à probabilités égales sans remise (PESR)*

Remarque

Un individu est choisi au plus une fois, chaque tirage fait décroître la population U d'une unité.

⇒ Les observations ne sont plus des variables aléatoires indépendantes les unes des autres.

Définition

Un estimateur classique de la moyenne μ d'une population U se définit par

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Propriété

On montre, par calculs (Ardilly, p :196-198), que

$$\mathbb{E}[\hat{\mu}] = \mu,$$

et

$$\text{Var}[\hat{\mu}] = \frac{N-n}{N-1} \frac{\sigma^2}{n} = (1-f) \frac{N}{N-1} \frac{\sigma^2}{n} = (1-f) \frac{S_c^2}{n}.$$

Remarque

L'avant dernière égalité de la dernière propriété implique que $\hat{\mu}$ est un estimateur sans biais de la moyenne μ de la population.

Remarque

Si la taille N de la population U est grande, la variance de $\hat{\mu}$ vaut :

$$\text{Var}[\hat{\mu}] \approx (1 - f) \frac{\sigma^2}{n}.$$

Remarque

Dans l'expression de la variance de $\hat{\mu}$, on remarque que le terme de la variance corrigée S_c^2 de la population U intervient. On sera donc amené à construire un estimateur de la variance de $\hat{\mu}$.

Remarque

On rappelle que la variance corrigée s_c^2 de l'échantillon S se définit par

$$s_c^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

et que s_c^2 est un estimateur sans biais de la variance corrigée S_c^2 de la population U .

Remarque

De cette dernière propriété, on en déduit que $(1-f) \frac{s_c^2}{n}$ est un estimateur sans biais de la variance de $\hat{\mu}$.

Définition

Un estimateur classique du total T d'une population U se définit par

$$\hat{T} = N\hat{\mu} = \frac{N}{n} \sum_{i=1}^n x_i.$$

Propriété

On montre, par calculs (Ardilly p :196-198), que

$$\mathbb{E}[\hat{T}] = T \quad \text{et} \quad \text{Var}[\hat{T}] = N^2(1-f)\frac{S_c^2}{n}.$$

Remarque

L'avant dernière égalité de la dernière propriété implique que \hat{T} est un estimateur sans biais du total T de la population U .

Remarque

Dans l'expression de la variance de \hat{T} , on remarque que le terme de la variance corrigée S_c^2 de la population U intervient. Or, dans la plupart des cas, on ne connaît pas la variance corrigée S_c^2 de la population U . On sera donc amené à construire un estimateur de la variance de \hat{T} .

Définition

Un estimateur de la variance de \widehat{T} se définit par

$$\widehat{\text{Var}} \left[\widehat{T} \right] = N^2(1 - f) \frac{s_c^2}{n},$$

où s_c^2 désigne la variance corrigée de l'échantillon S .

Propriété

On montre, par calcul direct, que

$$\mathbb{E} \left[\widehat{\text{Var}} \left[\widehat{T} \right] \right] = N^2(1 - f) \frac{S_c^2}{n}.$$

Remarque

De cette dernière propriété, on en déduit que $N^2(1 - f)\frac{s_c^2}{n}$ est un estimateur sans biais de la variance de \hat{T} .

Définition

Un estimateur de la variance σ^2 d'une population U , dans le cas d'un sondage aléatoire simple à probabilités égales sans remise, se définit par

$$\widehat{\sigma^2} = \frac{N-1}{N} s_c^2 = \frac{N-1}{N} \frac{1}{n-1} \sum_{i \in S} (x_i - \widehat{\mu})^2.$$

Propriété

On montre, par calculs (Ardilly et Tillé, p :43-49), que

$$\mathbb{E} \left[\widehat{\sigma^2} \right] = \mathbb{E} \left[\frac{N-1}{N} s_c^2 \right] = \sigma^2$$

et

$$\begin{aligned} \text{Var} \left[\widehat{\sigma^2} \right] &= \frac{(N-n)}{n(n-1)N(N-2)(N-3)} \\ &\quad \times \left\{ \mu_4(N-1) [N(n-1) - (n+1)] \right. \\ &\quad \left. \sigma^4 \left[N^2(n-3) + 6N - 3(n+1) \right] \right\}. \end{aligned}$$

Remarque

L'avant dernière égalité de la dernière propriété implique que $\widehat{\sigma^2}$ est un estimateur sans biais de la variance σ^2 de la population U .

Remarque

Dans l'expression de la variance de $\widehat{\sigma^2}$, on remarque que le terme σ^4 , qui est le carré de la variance de la population U , intervient ainsi que le moment d'ordre 4, μ_4 . Or, dans la plupart des cas, on ne connaît ni σ^4 , ni μ_4 . On sera donc amené à construire un estimateur de la variance de $\widehat{\sigma^2}$, si besoin est.

Sommaire

- 1 Introduction
- 2 Sondage aléatoire simple à probabilités égales avec remise
- 3 Sondage aléatoire simple à probabilités égales sans remise
- 4 **Comparaison des prélèvements PEAR et PESR**

Remarque

Les deux méthodes conduisent toutes les deux à des estimateurs $\hat{\mu}$ qui sont, **en moyenne** égaux au paramètre μ de la population.

Remarque

Par contre les variances de $\hat{\mu}$ ne sont pas égales !

Problème :

Qui est le meilleur estimateur de la moyenne μ de la population parmi ces deux estimateurs ?

Pour répondre à cette question, on va utiliser une méthode.

Remarque

En général, les estimateurs que l'on doit comparer sont en moyenne égaux au paramètre à estimer. Ils ne diffèrent que par leur variance. (La variance est un paramètre de précision de l'estimateur.)

Proposition

*Pour comparer deux estimateurs ou deux méthodes qui produisent des estimateurs différents, on utilise l'**effet de sondage**.*

Définition

L'effet de sondage *est défini par* :

$$D(\hat{\theta}^*|\hat{\theta}) = \frac{\text{Var}[\hat{\theta}^*]}{\text{Var}[\hat{\theta}]}.$$

Remarque

Si $D(\hat{\theta}^*|\hat{\theta}) < 1$, alors $\hat{\theta}^*$ sera plus précis que $\hat{\theta}$.

On rappelle que

Propriété

$$\text{Var} [\hat{\mu}_{PEAR}] = \frac{\sigma^2}{n} \quad \text{et} \quad \text{Var} [\bar{x}_{PESR}] = \frac{\sigma^2}{n} \frac{N-n}{N-1}.$$

Il s'en suit que

$$D(PESR|PEAR) = \frac{N-n}{N-1}.$$

Si

$$n > 1,$$

alors

$$N-n < N-1.$$

Par conséquent

$$D(PESR|PEAR) < 1.$$

Conclusion

La précision de l'estimateur est donc meilleure si l'on utilise un échantillon aléatoire simple PESR qu'un échantillon aléatoire simple PEAR.

Remarque

Ce dernier résultat est intuitif car il y a une perte d'information dès que certains individus sont observés plus d'une fois, ce qui est impossible lors d'un tirage sans remise.

Remarques

- 1 Si la taille de la population est grande, l'effet de sondage est tel que

$$D(\text{PESR}|\text{PEAR}) = \frac{N-n}{N-1} \approx \frac{N-n}{N} = 1-f,$$

où f est le taux de sondage. L'amélioration de la précision est d'autant meilleure que f est grand.

- 2 La différence entre les deux procédures faiblit quand la taille de l'échantillon est petite par rapport à celle de la population, i.e. quand f est faible ! Dans ce cas l'effet de sondage est proche de 1, les deux méthodes fournissent des estimateurs de précision analogue.