

T. D. n° 7

Estimation par différence et par régression

Corrigé

Exercice 1. *Stratification et estimateur par la différence.* D'après l'examen de Janvier 2006, M2-Statistique.

1. L'estimateur est sans biais. En effet, puisque nous avons :

$$\widehat{Y}_\pi = \sum_{h=1}^H \frac{N_h}{N} \widehat{Y}_h,$$

où \widehat{Y}_h désigne la moyenne simple des y_k dans l'échantillon de la strate h ,

$$\begin{aligned} \mathbb{E}(\widehat{Y}_D) &= \mathbb{E}(\widehat{Y}_\pi + \bar{X} - \widehat{X}_\pi) \\ &= \mathbb{E}(\widehat{Y}_\pi) + \bar{X} - \mathbb{E}(\widehat{X}_\pi) \\ &= \bar{Y} + \bar{X} - \bar{X} \\ &= \bar{Y}. \end{aligned}$$

2. Nous posons :

$$z_k = y_k - x_k.$$

Nous avons :

$$\widehat{Y}_D = \bar{X} + \widehat{Z}_\pi.$$

Donc nous pouvons calculer la variance de l'estimateur en se servant de cette décomposition :

$$\begin{aligned} \text{Var}(\widehat{Y}_D) &= \text{Var}(\widehat{Z}_\pi) \\ &= \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{zh}^2}{n_h}, \end{aligned}$$

où S_{zh}^2 se définit par :

$$\begin{aligned} S_{zh}^2 &= \frac{1}{N_h - 1} \sum_{k \in U_h} (z_k - \widehat{Z}_h)^2 \\ &= S_{yh}^2 + S_{xh}^2 - 2S_{xyh}, \end{aligned}$$

et

$$S_{xyh} = \frac{1}{N_h - 1} \sum_{k \in U_h} (x_k - \widehat{X}_h)(y_k - \widehat{Y}_h).$$

3. En posant $z_k = y_k - x_k$, le problème revient à minimiser $\text{Var}(\widehat{Z}_\pi)$ sous la seule contrainte de taille fixe, qui s'écrit ici $\sum_{h=1}^H n_h = n$. En effet, le coût unitaire est le même dans toutes les strates, ce qui donne :

$$n_h = \frac{N_h S_{zh}}{\sum_{l=1}^H N_l S_{zl}} n.$$

En pratique, nous estimons a priori les S_{zh} , nous arrondissons n_h à l'entier le plus proche, après avoir fixé n en fonction du budget global dont nous disposons. Il peut arriver que nous obtenons $n_h > N_h$ pour certains h : dans ce cas, nous imposons $n_h = N_h$ et nous reprenons l'ensemble du calcul avec les strates restantes.

4. Puisque la variance de l'estimateur s'écrit :

$$\text{Var}(\widehat{Y}_\pi) = \sum_{h=1}^H \left(\frac{N_h}{N}\right)^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{yh}^2}{n_h},$$

et que les deux estimateurs sont sans biais, \widehat{Y}_D est indiscutablement préférable à \widehat{Y}_π lorsque, pour tout h :

$$S_{yh}^2 > S_{zh}^2,$$

soit, pour tout h :

$$\frac{S_{xyh}}{S_{xh}^2} > \frac{1}{2}.$$

Cette condition revient à obtenir une droite de régression de y sur x qui, dans chaque strate, ait une pente supérieure à $1/2$. C'est en particulier le cas si nous posons $y = x$ (pente égale à 1) : ce résultat est naturel, car alors :

$$\widehat{X}_D = \bar{X}$$

quel que soit l'échantillon tiré. Nous disons que l'estimateur \widehat{Y}_D est « calé » sur \bar{X} .

Exercice 2. Comparaison de plusieurs estimateurs.

- 1.(i) **Estimation par la moyenne à probabilités égales avec remise.**

L'estimateur de la moyenne \bar{X} est égal à :

$$\bar{x}_{SASPEAR} = \frac{1}{n} \sum_{i=1}^n x_i$$

et sa variance est égale à :

$$\text{Var}(\bar{x}_{SASPEAR}) = \frac{\sigma_{pop}^2}{n}.$$

(ii) **Estimation par la moyenne à probabilités égales sans remise.**

L'estimateur de la moyenne \bar{X} est égal à :

$$\bar{x}_{SASPESR} = \frac{1}{n} \sum_{i=1}^n x_i$$

et sa variance est égale à :

$$\text{Var} [\bar{x}_{SASPESR}] = \frac{N-n}{N-1} \frac{\sigma_{pop}^2}{n}.$$

(iii) **Estimation par le quotient.**

L'estimateur de la moyenne \bar{X} est égal à :

$$\bar{x}_Q = \bar{Y} \frac{\bar{x}}{\bar{y}}.$$

Cet estimateur est biaisé et par conséquent ce n'est pas la variance qu'il faut étudier mais l'erreur quadratique de cet estimateur.

2. Construisons l'estimateur par la différence \bar{x}_D :

$$\bar{x}_D = \bar{Y} + \bar{x}_{SAS} - \bar{y}_{SAS},$$

où \bar{Y} dénote la moyenne des y pour la population entière.

Justification de l'estimateur. Il y a deux façons de justifier intuitivement cet estimateur.

Première justification. Nous pouvons décomposer \bar{X} en deux parties, l'une connue, l'autre pas :

$$\bar{X} = \bar{Y} + (\bar{X} - \bar{Y}).$$

La première partie, \bar{Y} , est connue, et nous n'avons pas besoin de l'estimer. La deuxième partie, $(\bar{X} - \bar{Y})$, la différence entre la moyenne des x et celle des y , n'est pas connue et doit être estimée. Nous l'estimons, naturellement, par la différence entre les deux moyennes échantillonnées, $\bar{x} - \bar{y}$. C'est ce qui donne l'estimateur par la différence.

Deuxième justification. L'estimateur naturel de \bar{X} est \bar{x} et a priori c'est l'estimateur privilégié. Dans l'estimateur par la différence, écrit comme :

$$\bar{x}_D = \bar{x} + (\bar{Y} - \bar{y}),$$

l'ajout du terme $(\bar{Y} - \bar{y})$ peut s'interpréter comme un ajustement à l'estimateur \bar{x} . Grâce à notre information sur la variable y , nous pouvons deviner si, en l'occurrence, l'estimateur \bar{x} a surestimé ou sous-estimé la moyenne \bar{X} .

3.(i) Calculons l'espérance de l'estimateur \bar{x}_D :

$$\begin{aligned} \mathbb{E}(\bar{x}_D) &= \mathbb{E}(\bar{x}_{SAS} + \bar{Y} - \bar{y}_{SAS}) \\ &= \mathbb{E}(\bar{x}_{SAS}) + \bar{Y} - \mathbb{E}(\bar{y}_{SAS}) \\ &= \bar{X} + \bar{Y} - \bar{Y} \\ &= \bar{X}. \end{aligned}$$

Nous en concluons que l'estimateur de la différence \bar{x}_D est un estimateur sans biais de \bar{X} .

(ii) Calculons la variance de l'estimateur \bar{x}_D :

$$\begin{aligned}\text{Var}(\bar{x}_D) &= \text{Var}(\bar{x}_{SAS} + \bar{Y} - \bar{y}_{SAS}) \\ &= \text{Var}(\bar{x}_{SAS}) - 2\text{Cov}(\bar{x}_{SAS}, \bar{y}_{SAS}) + \text{Var}(\bar{y}_{SAS}).\end{aligned}$$

À partir de cette égalité, il est d'usage de se demander comment nous devons tirer l'échantillon : à probabilités égales avec ou sans remise.

4.(i) **Estimation par la moyenne à probabilités égales avec remise.**

$$22072,63.$$

(ii) **Estimation par la moyenne à probabilités égales sans remise.**

$$22072,63.$$

(iii) **Estimation par le quotient.**

$$32039,66 \times \left(\frac{22072,63}{21585,26} \right) = 32763,08.$$

(iv) **Estimateur par la différence.**

$$32039,66 + (22072,63 - 21585,26) = 32527,03.$$

Nous notons une grande variabilité dans les estimations. Pour choisir la meilleure estimation et par conséquent le meilleur estimateur, il faut savoir quelle est la plus petite variance associée à chaque estimateur. Pour cela, nous allons calculer les différents intervalles de confiance à 95%.

(i) **Intervalle de confiance pour la moyenne \bar{X} (PEAR).**

$$\begin{aligned}22072,63 \pm 1,96 \times \sqrt{\frac{4\,131\,789\,466}{35}} \\ 22072,63 \pm 1,96 \times 10865,13 \\ 22072,63 \pm 21295,67\end{aligned}$$

où 1,96 est le quantile de la loi normale centrée réduite à 95%. Or nous faisons une large approximation en utilisant la loi normale centrée réduite. Il serait conseillé d'utiliser plutôt un quantile de Student à 95%.

(ii) **Intervalle de confiance pour la moyenne \bar{X} (PESR).**

$$\begin{aligned}22072,63 \pm 1,96 \times \sqrt{\left(1 - \frac{35}{180}\right) \frac{4\,131\,789\,466}{35}} \\ 22072,63 \pm 1,96 \times 9751,76 \\ 22072,63 \pm 19113,44\end{aligned}$$

où 1,96 est le quantile de la loi normale centrée réduite à 95%. Or nous faisons une large approximation en utilisant une loi normale centrée réduite. Il serait conseillé d'utiliser plutôt un quantile de Student à 95%.

- (iii) **Intervalle de confiance pour la moyenne \bar{X} (par le quotient).**
 Estimons maintenant la variance de \bar{x}_Q :

$$\begin{aligned}\widehat{\text{Var}}(\bar{x}_Q) &= (1-f) \frac{s_x^2 - 2\hat{R}s_{xy} + \hat{R}^2 s_y^2}{n} \\ &= \left(1 - \frac{35}{180}\right) \frac{4\,131\,789\,466 - 2,04(4\,059\,448\,772) + (1,02)^2 3\,989\,656\,072}{35} \\ &= 32\,918,77.\end{aligned}$$

D'où l'intervalle de confiance pour la moyenne \bar{X} est égal à :

$$\begin{aligned}32763,08 \pm 1,96 \times 181,43 \\ 32763,08 \pm 355,61.\end{aligned}$$

- (iv) **Intervalle de confiance pour la moyenne \bar{X} (par la différence).**
 Estimons maintenant la variance de \bar{x}_D :

$$\begin{aligned}\widehat{\text{Var}}(\bar{x}_D) &= (1-f) \frac{s_x^2 - 2s_{xy} + s_y^2}{n} \\ &= \left(1 - \frac{35}{180}\right) \frac{4\,131\,789\,466 - 2(4\,059\,448\,772) + 3\,989\,656\,072}{35} \\ &= 58\,644,31.\end{aligned}$$

D'où l'intervalle de confiance pour la moyenne \bar{X} est égal à :

$$\begin{aligned}32527,03 \pm 1,96 \times 242,17 \\ 32527,03 \pm 474,64.\end{aligned}$$

Conclusion : Les deux derniers intervalles de confiance pour la moyenne \bar{X} sont les plus intéressants pour nous car les deux écart-types sont très petits. D'autre part, il faut noter qu'il y a peu de différence entre l'estimateur par le quotient et l'estimateur par la différence.

5. Nous constatons que les deux estimateurs qui utilisent la variable auxiliaire sont très nettement meilleurs. Les estimations elles-mêmes sont très différentes : de l'ordre de 22 000 pour l'estimateur par la moyenne avec ou sans remise ; et de l'ordre de 32 000 pour les deux autres estimateurs. S'il fallait décider lequel des estimateurs à employer, nous dirions 32 000 plutôt que 22 000. Mais il n'est pas question, en pratique, de calculer plusieurs estimateurs et puis leur écart-type. Il faudrait pouvoir faire un choix d'avance. Est-ce que nous aurions pu prévoir la supériorité des deux derniers estimateurs, et s'en tenir à l'un ou l'autre de ces deux, sans même considérer le premier ? Dans le cas présent, nous aurions pu le prévoir. En général, les estimateurs qui font appel à une variable auxiliaire sont avantageux dans la mesure où la variable auxiliaire est corrélée positivement avec la variable d'intérêt. Il est évident que le nombre d'habitants en 1996 est corrélé avec le nombre d'habitants en 2001. C'est donc une information pertinente, et les estimateurs par la différence et par le quotient en tirent profit. La leçon importante qui se dégage est « si la variable auxiliaire est fortement liée à la variable d'intérêt, il vaut mieux utiliser l'estimateur par la différence ou l'estimateur par le quotient ».