

Article

Calage et poids restreints

par Alain Théberge

Juin 2000



Calage et poids restreints

Alain Théberge¹

Résumé

Pour mieux comprendre l'impact de l'imposition d'une région de restriction sur les poids de calage, on examine le comportement asymptotique de ceux-ci. On donne des conditions nécessaires et suffisantes pour l'existence d'une solution à l'équation de calage avec des poids à l'intérieur d'intervalles donnés. Une formulation plus générale du problème de calage permet de faire un compromis entre le besoin de satisfaire l'équation de calage, et le désir d'obtenir des poids qui sont près des poids de Horvitz-Thompson. Si on relâche les exigences vis-à-vis l'équation de calage, alors diverses méthodes d'estimation avec poids restreints peuvent être utilisées. Les estimateurs présentés ont habituellement les mêmes propriétés asymptotiques que l'estimateur par calage sans restrictions sur les poids, et certains ont des poids qui peuvent être calculés explicitement, sans procédure itérative. On montre comment ces estimateurs peuvent être adaptés pour tirer parti d'un estimateur synthétique. Une stratégie semblable à celle utilisée pour restreindre les poids est appliquée aux données aberrantes.

Mots clés : Petits domaines; inverse de Moore-Penrose; solutions d'inéquations; propriétés asymptotiques; données aberrantes.

1. Introduction

L'estimateur par calage possède de bonnes propriétés asymptotiques. Mais pour de petites tailles d'échantillon, ou si le calage se fait au niveau de domaines dont certains ont peu d'observations, les poids de cet estimateur peuvent prendre des valeurs extrêmes. Une façon de pallier ce problème consiste à utiliser la méthode de calage avec des mesures de distance qui font que les poids des observations sont restreints à certains intervalles autour des poids de sondage. Cette approche a été développée par Deville et Särndal (1992). D'autres méthodes qui visent à obtenir des estimations robustes qui respectent l'équation de calage sont données dans Duchesne (1999). Cet article contient une bonne bibliographie sur les estimateurs robustes. On n'est toutefois pas assuré de l'existence d'une solution à l'équation de calage avec des poids restreints. Même si de tels poids existent, il se peut que le statisticien préfère résoudre le problème de poids extrêmes en relâchant quelque peu ses exigences vis-à-vis l'équation de calage, plutôt qu'en resserrant les contraintes sur les poids en utilisant une mesure de distance plus « contraignante ». On présentera ici, une formulation du problème de calage qui offre plus de flexibilité au statisticien. Il s'agit en fait d'un problème de minimisation du même ordre que celui rencontré lors d'une régression avec coefficients de coûts (« ridge regression »). Bardsley et Chambers (1984) ont rencontré le même problème de minimisation dans leur recherche d'estimateurs basés sur des modèles. Cette formulation du problème de calage peut être utilisée pour restreindre les poids sans utiliser de mesure de distance spéciale entre les poids calés et les poids de Horvitz-Thompson. Rao et Singh (1997) ont eux combiné cette approche avec des méthodes itératives

utilisant des mesures de distance. D'autres façons de restreindre les poids seront aussi examinées.

Dans la section suivante on présente la méthode de calage en l'absence de bornes sur les valeurs des poids. Le problème de calage qui est posé ne présume pas de l'existence d'une solution à l'équation de calage. Les propriétés asymptotiques des poids calés sont discutées. Ces propriétés sont d'intérêt pour le comportement asymptotique des estimateurs dont les poids sont restreints. La section 3 donne des conditions nécessaires et suffisantes pour l'existence de poids restreints qui satisfont l'équation de calage. À la section 4, on discute d'une façon de poser le problème d'estimation en dosant l'importance qu'on accorde à l'équation de calage. La section 5 donne différentes façons de restreindre les poids qui ne reposent pas sur l'utilisation d'une distance particulière. On propose à la section 6, un estimateur avec poids restreints qui est utile pour de petits domaines. Finalement, la section 7 aborde les données aberrantes en développant une méthode semblable à celle utilisée pour traiter les poids extrêmes.

2. Calage

Soient $Y \in \mathbb{R}^{N \times d}$ une matrice de d variables d'intérêt pour une population de taille N , et $c \in \mathbb{R}^d$ un vecteur de constantes connues, on tire un échantillon s de taille n et on notera à l'aide de l'indice s les sous-vecteurs ou les sous-matrices qui correspondent à l'échantillon. Il s'agit d'estimer $Y'c$ par $Y'_s w_s$, où $w_s \in \mathbb{R}^n$ est un vecteur de poids pour les unités échantillonnées. Pour un vecteur v et une matrice diagonale positive F de même dimension, on définit $\|v\|_F^2 = v'Fv$. Pour une matrice d'information

1. Alain Théberge, Division des méthodes d'enquêtes sociales, Statistique Canada, Ottawa (Ontario) K1A 0T6 Canada.

auxiliaire $X \in \mathbb{R}^{N \times p}$, $A \in \mathbb{R}^{N \times N}$ la matrice diagonale des poids de sondage, des matrices diagonales positives données $U_s \in \mathbb{R}^{n \times n}$ et $T \in \mathbb{R}^{p \times p}$, on cherche parmi les vecteurs de poids $\mathbf{w}_s \in \mathbb{R}^n$ qui minimisent $\|X'_s \mathbf{w}_s - X'c\|_T^2$, celui qui minimise $D_s(\mathbf{w}_s) = \|\mathbf{w}_s - A_s c_s\|_{U_s}^2$. Cette formulation du problème qui ne présume pas de l'existence de poids satisfaisant l'équation de calage, $X'_s \mathbf{w}_s = X'c$, est donnée dans Théberge (1999). La solution recherchée constitue le vecteur des poids calés \mathbf{w}_{cal} . On a

$$\mathbf{w}_{\text{cal}} = A_s c_s + U_s^{-1} X_s T^{1/2} (T^{1/2} X'_s U_s^{-1} X_s T^{1/2})^\dagger T^{1/2} (X'c - X'_s A_s c_s), \quad (1)$$

où F^\dagger dénote l'inverse de Moore-Penrose de la matrice F .

Afin de mieux étudier les propriétés asymptotiques d'estimateurs par calage avec poids restreints, on examinera maintenant le comportement de \mathbf{w}_{cal} lorsque $n \rightarrow \infty$. On suppose l'existence d'un cadre asymptotique où la taille de la population et la taille de l'échantillon tendent vers l'infini, voir par exemple Isaki et Fuller (1982), et pour lequel on a

$$\begin{aligned} Y'c &= O_p(N^\gamma) (\gamma \geq 0) \\ X'c - X'_s A_s c_s &= O_p(n^{-1/2} N^\gamma) \\ T^{1/2} X'_s U_s^{-1} X_s T^{1/2} &= O_p(n). \end{aligned} \quad (2)$$

Il s'ensuit que $(T^{1/2} X'_s U_s^{-1} X_s T^{1/2})^\dagger = O_p(n^{-1})$, puisqu'une des propriétés de l'inverse de Moore-Penrose d'une matrice F est $F^\dagger F F^\dagger = F^\dagger$. Habituellement, on peut s'attendre à avoir $\gamma = 1$ lorsque chaque élément du vecteur c vaut 1 (estimation d'un total), et $\gamma = 0$ lorsque chaque élément de c vaut $1/N$ (estimation d'une moyenne). Sous les conditions (2) on a donc,

$$\begin{aligned} \mathbf{w}_{\text{cal}} - A_s c_s &= U_s^{-1} X_s T^{1/2} (T^{1/2} X'_s U_s^{-1} X_s T^{1/2})^\dagger \\ &\quad T^{1/2} (X'c - X'_s A_s c_s) \\ &= O_p(n^{-1}) O_p(n^{-1/2} N^\gamma) \\ &= O_p(n^{-3/2} N^\gamma). \end{aligned} \quad (3)$$

Donc $\mathbf{w}_{\text{cal}} - A_s c_s$ converge en probabilité vers $\mathbf{0}$, si

$$\lim_{n, N \rightarrow \infty} n^{-3/2} N^\gamma = 0.$$

Pour un cadre asymptotique tel celui de Brewer (1979) où la fraction de sondage n/N est constante, ou tout cadre pour lequel la fraction de sondage converge vers un nombre positif, cette condition est vérifiée si $\gamma < 3/2$.

Si on écrit $\mathbf{w}_{\text{cal}} = A_s c_s + U_s^{-1} X_s T^{1/2} H_s^\dagger T^{1/2} (X'c - X'_s A_s c_s)$, où $H_s = T^{1/2} X'_s U_s^{-1} X_s T^{1/2}$, on a

$$\begin{aligned} D_s(\mathbf{w}_{\text{cal}}) &= (X'c - X'_s A_s c_s)' T^{1/2} H_s^\dagger H_s H_s^\dagger \\ &\quad T^{1/2} (X'c - X'_s A_s c_s) \\ &= (X'c - X'_s A_s c_s)' T^{1/2} H_s^\dagger T^{1/2} (X'c - X'_s A_s c_s) \\ &= O_p(n^{-1/2} N^\gamma) O_p(n^{-1}) O_p(n^{-1/2} N^\gamma) \\ &= O_p(n^{-2} N^{2\gamma}). \end{aligned} \quad (4)$$

Toujours pour un cadre asymptotique où la fraction de sondage converge vers un nombre positif, on a $D_s(\mathbf{w}_{\text{cal}})$ converge en probabilité vers 0, si $\gamma < 1$. Il y a donc des cas, notamment pour l'estimation d'un total, où $\mathbf{w}_{\text{cal}} - A_s c_s$ converge en probabilité vers $\mathbf{0}$, mais où $D_s(\mathbf{w}_{\text{cal}}) = \|\mathbf{w}_{\text{cal}} - A_s c_s\|_{U_s}^2$ ne converge pas vers 0.

3. Solutions à l'équation de calage et poids restreints

Même lorsqu'il n'y a pas de restrictions sur les poids, il est possible qu'il n'y ait pas de solution à l'équation de calage. En appliquant Graybill (1983, 113) au problème de calage, on obtient que l'équation de calage $X'_s \mathbf{w}_s = X'c$ a une solution si et seulement si $(X'_s X_s)^\dagger X'c = X'c$. Si une solution existe, il est possible que les poids calés soient négatifs ou exceptionnellement grands. Deville et Särndal (1992) ont proposé d'utiliser diverses mesures de distance autres qu'une somme pondérée de carrés pour mesurer la distance entre les poids de Horvitz-Thompson et les poids calés, afin de restreindre les poids à certains intervalles tout en satisfaisant l'équation de calage. Cette approche ne fonctionnera que s'il existe des poids dans ces intervalles qui satisfont l'équation de calage. Le but principal de cette section est de trouver des conditions nécessaires et suffisantes pour l'existence d'un vecteur de poids \mathbf{w}_s à l'intérieur de bornes données, et tel que les estimations des totaux pour les variables auxiliaires sont également bornées. C'est-à-dire qu'on cherche des conditions pour l'existence d'un vecteur \mathbf{w}_s tel que $\mathbf{w}^{(L)} \leq \mathbf{w}_s \leq \mathbf{w}^{(H)}$ et $\mathbf{t}^{(L)} \leq X'_s \mathbf{w}_s \leq \mathbf{t}^{(H)}$, où $\mathbf{w}^{(L)}$, $\mathbf{w}^{(H)}$, $\mathbf{t}^{(L)}$ et $\mathbf{t}^{(H)}$ sont donnés. En particulier, en posant $\mathbf{t}^{(L)} = \mathbf{t}^{(H)} = X'c$, on obtiendrait des conditions pour l'existence de poids restreints aux intervalles $\mathbf{w}^{(L)} \leq \mathbf{w}_s \leq \mathbf{w}^{(H)}$, qui satisfont l'équation de calage.

Une première étape est fournie par le théorème suivant de Fan (1956). Il est énoncé ici pour une matrice M de dimension finie, quoique la preuve donnée par Fan tient également pour une matrice de dimension infinie. Le théorème fait appel au noyau de M' , $N(M')$, défini comme l'ensemble des vecteurs α tels que $M'\alpha = \mathbf{0}$.

Théorème : Soient $M \in \mathbb{R}^{m \times n}$ et $l \in \mathbb{R}^m$, $\exists \mathbf{w} \in \mathbb{R}^n$ tel que $M\mathbf{w} \geq l$ si et seulement si pour tout $\lambda \geq \mathbf{0}$ dans $N(M')$, on a $l'\lambda \leq 0$.

Corollaire : Soient $M \in \mathbb{R}^{m \times n}$ et $l, h \in \mathbb{R}^m$, $\exists w \in \mathbb{R}^n$ tel que $l \leq Mw \leq h$ si et seulement si premièrement $l \leq h$ et deuxièmement $\lambda \in N(M') \Rightarrow -l' \lambda_- \leq h' \lambda_+$, où $\lambda_+ = \max(\lambda, 0)$ et $\lambda_- = \min(\lambda, 0)$ les extrema étant pris élément par élément.

Le corollaire est obtenu en posant

$$M = \begin{pmatrix} M \\ -M \end{pmatrix}, l = \begin{pmatrix} l \\ -h \end{pmatrix} \text{ et } \lambda = \begin{pmatrix} -\lambda_- \\ \lambda_+ \end{pmatrix}$$

dans le théorème.

On note p la dimension de $N(M')$. Si p est égal à zéro, alors $\lambda \in N(M')$ implique $\lambda = 0$, et la condition du théorème (ou la condition similaire du corollaire) est évidemment satisfaite. Si p est égal à un, alors $\lambda \in N(M')$ implique que λ est un multiple d'un vecteur z , et il est suffisant de vérifier la condition pour $\lambda = z$ et $\lambda = -z$. En utilisant la propriété $(-\lambda)_- = -(\lambda)_+$, le problème posé au début de la section peut maintenant être résolu si X_s est un vecteur. Le corollaire avec

$$M = \begin{pmatrix} I_n \\ X'_s \end{pmatrix}, l = \begin{pmatrix} w^{(L)} \\ t^{(L)} \end{pmatrix}, h = \begin{pmatrix} w^{(H)} \\ t^{(H)} \end{pmatrix},$$

et le fait que

$$z = \begin{pmatrix} -X_s \\ 1 \end{pmatrix}$$

engendre $N(M')$, donne les conditions nécessaires et suffisantes

$$\begin{aligned} w^{(L)} &\leq w^{(H)} \\ t^{(L)} &\leq t^{(H)} \\ (X_s)_+ w^{(L)} + (X_s)_- w^{(H)} &\leq t^{(H)} \\ t^{(L)} &\leq (X_s)_+ w^{(H)} + (X_s)_- w^{(L)}. \end{aligned} \quad (5)$$

La troisième inégalité dans (5) affirme que le total pondéré de la variable auxiliaire ne doit pas être supérieur à $t^{(H)}$, lorsque le plus petit poids possible $w^{(L)}$ est donné aux unités où la variable auxiliaire prend une valeur positive, et lorsque le plus grand poids possible $w^{(H)}$ est donné aux unités où la variable auxiliaire prend une valeur négative. La quatrième inégalité dans (5) affirme que le total pondéré de la variable auxiliaire ne doit pas être inférieur à $t^{(L)}$, lorsque le plus grand poids possible est donné aux unités où la variable auxiliaire prend une valeur positive, et lorsque le plus petit poids possible est donné aux unités où la variable auxiliaire prend une valeur négative.

Même lorsque $p > 1$, il est suffisant de vérifier la condition du corollaire pour un nombre fini de valeurs de λ . Soit $V \in \mathbb{R}^{m \times p}$, $2 \leq p \leq m$ une matrice dont les colonnes forment une base pour $N(M')$. Il est toujours possible de construire V de sorte que p de ses lignes,

v_1, v_2, \dots, v_m , soient les vecteurs unités de \mathbb{R}^p , et on supposera que V est de cette forme. On démontre dans l'annexe A qu'il est suffisant de vérifier la condition du corollaire pour les vecteurs $\lambda = V\phi$ et $\lambda = -V\phi$, où $\phi = (\phi_1, \dots, \phi_p)'$ est un vecteur non nul et satisfait $v_i' \phi = 0$ pour un sous-ensemble de $(p-1)$ vecteurs v_i linéairement indépendants. On doit donc vérifier la condition pour au plus $\binom{m}{p-1}$ vecteurs ϕ , donc au plus $2\binom{m}{p-1}$ valeurs de λ .

En utilisant le corollaire avec

$$M = \begin{pmatrix} I_n \\ X'_s \end{pmatrix}, l = \begin{pmatrix} w^{(L)} \\ t^{(L)} \end{pmatrix}, h = \begin{pmatrix} w^{(H)} \\ t^{(H)} \end{pmatrix},$$

et en notant que les colonnes de

$$V = \begin{pmatrix} -X_s \\ I_p \end{pmatrix}$$

forment une base pour $N(M')$, on obtient les conditions nécessaires et suffisantes suivantes pour l'existence d'une solution au problème mentionné au début de cette section lorsque $X_s \in \mathbb{R}^{n \times p}$ avec $p > 1$. On doit avoir $w^{(L)} \leq w^{(H)}$, $t^{(L)} \leq t^{(H)}$, et pour chaque sous-ensemble de $(p-1)$ lignes de

$$V = \begin{pmatrix} -X_s \\ I_p \end{pmatrix}$$

linéairement indépendantes,

$$(X_s \phi)_+ w^{(L)} - \phi_- t^{(L)} \leq -(X_s \phi)_- w^{(H)} + \phi_+ t^{(H)}$$

$$-(X_s \phi)_- w^{(L)} + \phi_- t^{(L)} \leq (X_s \phi)_+ w^{(H)} - \phi_+ t^{(H)} \quad (6)$$

pour un vecteur non nul $\phi \in \mathbb{R}^p$ de direction perpendiculaire à chaque ligne du sous-ensemble. La deuxième inéquation de (6) est obtenue de la première en changeant le signe de ϕ .

Si $V_{\text{sub}} \in \mathbb{R}^{p \times p}$ est une matrice non singulière dont les lignes sont des lignes de V , alors chaque colonne de V_{sub}^{-1} est un vecteur perpendiculaire à $(p-1)$ lignes de V linéairement indépendantes. D'où le résultat suivant :

Il existe un vecteur de poids w_s tels que $w^{(L)} \leq w_s \leq w^{(H)}$ et $t^{(L)} \leq X'_s w_s \leq t^{(H)}$ si et seulement si $w^{(L)} \leq w^{(H)}$, $t^{(L)} \leq t^{(H)}$ et

$$\begin{aligned} (X_s V_{\text{sub}}^{-1})'_+ w^{(L)} - (V_{\text{sub}}^{-1})'_- t^{(L)} &\leq \\ &-(X_s V_{\text{sub}}^{-1})'_- w^{(H)} + (V_{\text{sub}}^{-1})'_+ t^{(H)} \end{aligned} \quad (7)$$

$$-(X_s V_{\text{sub}}^{-1})'_- w^{(L)} + (V_{\text{sub}}^{-1})'_+ t^{(L)} \leq$$

$$(X_s V_{\text{sub}}^{-1})'_+ w^{(H)} - (V_{\text{sub}}^{-1})'_- t^{(H)}$$

pour toutes les matrices non singulières $V_{\text{sub}} \in \mathbb{R}^{p \times p}$ dont les lignes sont des lignes de

$$V = \begin{pmatrix} -X_s \\ I_p \end{pmatrix}.$$

Ces conditions sont dans une certaine mesure redondantes. Par exemple, si les inégalités (7) sont satisfaites pour $V_{\text{sub}} = V_1$, alors elles sont nécessairement satisfaites pour toute matrice V_2 obtenue de V_1 par une permutation de lignes.

Un autre exemple est fourni par la pondération d'observations dans un tableau de contingence. Soient $\hat{N}_{ij} = n_{ij}w_{ij}$ ($i=1, 2, \dots, R; j=1, 2, \dots, C$), où n_{ij} est le nombre d'observations dans la cellule (i, j) du tableau de contingence et w_{ij} est le poids de chacune de ces observations, on veut savoir s'il existe des poids w_{ij} tels que \hat{N}_{ij} satisfait certaines contraintes. Par exemple, motivé par le problème de la convergence de la procédure du quotient réitéré (raking ratio), Bacharach (1965) donne des conditions nécessaires et suffisantes pour l'existence de poids w_{ij} tels que $\hat{N}_{ij} \geq 0, \sum_{i=1}^R \hat{N}_{ij} = N_j (j=1, \dots, C), \sum_{j=1}^C \hat{N}_{ij} = N_i (i=1, \dots, R)$, où les valeurs de N_j et N_i sont données. Le résultat suivant qui est démontré à l'annexe B, est plus général.

Pour des constantes arbitraires $N_{ij}^{(L)}, N_{ij}^{(H)}, N_{.j}^{(L)}, N_{.j}^{(H)}, N_{i.}^{(L)}, N_{i.}^{(H)}, N_{..}^{(L)}$, et $N_{..}^{(H)}$, il existe des \hat{N}_{ij} tels que

$$\begin{aligned} N_{ij}^{(L)} &\leq \hat{N}_{ij} \leq N_{ij}^{(H)} & i=1, \dots, R; j=1, \dots, C; \\ N_{.j}^{(L)} &\leq \sum_{i=1}^R \hat{N}_{ij} \leq N_{.j}^{(H)} & j=1, \dots, C; \\ N_{i.}^{(L)} &\leq \sum_{j=1}^C \hat{N}_{ij} \leq N_{i.}^{(H)} & i=1, \dots, R; \\ N_{..}^{(L)} &\leq \sum_{i=1}^R \sum_{j=1}^C \hat{N}_{ij} \leq N_{..}^{(H)}, \end{aligned} \tag{8}$$

si et seulement si

$$\begin{aligned} N_{ij}^{(L)} &\leq N_{ij}^{(H)} & i=1, \dots, R; j=1, \dots, C; \\ N_{.j}^{(L)} &\leq N_{.j}^{(H)} & j=1, \dots, C; \\ N_{i.}^{(L)} &\leq N_{i.}^{(H)} & i=1, \dots, R; \\ N_{..}^{(L)} &\leq N_{..}^{(H)} \end{aligned}$$

$$\begin{aligned} &\sum_{j \in T} \left(N_{.j}^{(L)} - \sum_{i \in S} N_{ij}^{(H)} \right) & (9) \\ &\leq \sum_{i \in S} \left(N_{i.}^{(H)} - \sum_{j \in T} N_{ij}^{(L)} \right) \\ &\sum_{i \in S} \left(N_{i.}^{(L)} - \sum_{j \in T} N_{ij}^{(H)} \right) \\ &\leq \sum_{j \in T} \left(N_{.j}^{(H)} - \sum_{i \in T} N_{ij}^{(L)} \right) \\ &N_{..}^{(L)} + \sum_{j \in T} \left(N_{.j}^{(H)} - \sum_{i \in S} N_{ij}^{(H)} \right) \\ &\leq \sum_{i \in S} \left(N_{i.}^{(H)} - \sum_{j \in T} N_{ij}^{(L)} \right) + \sum_{j=1}^J N_{.j}^{(H)} \\ &\sum_{i=1}^I N_{i.}^{(L)} + \sum_{j \in T} \left(N_{.j}^{(L)} - \sum_{i \in S} N_{ij}^{(H)} \right) \\ &\leq \sum_{i \in S} \left(N_{i.}^{(L)} - \sum_{j \in T} N_{ij}^{(L)} \right) + N_{..}^{(H)} \end{aligned}$$

pour tout $S \subseteq \{1, 2, \dots, R\}, T \subseteq \{1, 2, \dots, C\}$

On peut réduire le nombre d'inéquations à vérifier. Par exemple, plutôt que de vérifier

$$\sum_{j \in T} \left(N_{.j}^{(L)} - \sum_{i \in S} N_{ij}^{(H)} \right) \leq \sum_{i \in S} \left(N_{i.}^{(H)} - \sum_{j \in T} N_{ij}^{(L)} \right)$$

pour tout $S \subseteq \{1, 2, \dots, R\}$, et $T \subseteq \{1, 2, \dots, C\}$, on montre facilement qu'il est équivalent de vérifier que

$$\sum_{j \in T} N_{.j}^{(L)} \leq \sum_{i=1}^R \min \left(\left(N_{i.}^{(H)} - \sum_{j \in T} N_{ij}^{(L)} \right), \sum_{j \in T} N_{ij}^{(H)} \right)$$

pour tout $T \subseteq \{1, 2, \dots, C\}$.

4. Calage mitigé

On peut ne pas être satisfait avec l'approche en deux temps de l'estimation par calage, où on cherche d'abord les vecteurs de poids qui satisfont au mieux l'équation de calage, puis parmi cet ensemble de vecteurs, on cherche celui qui se rapproche le plus des poids de Horvitz-Thompson. Pour de petits échantillons, cette méthode peut conduire à des poids que le statisticien considère trop éloignés des poids de Horvitz-Thompson.

On pourrait préférer doser l'importance qu'on accorde à l'équation de calage par rapport à la norme de $w_s - A_s c_s$. Ainsi, on peut désirer trouver un vecteur de poids w_s qui minimise

$$\left\| \begin{pmatrix} w_s - A_s c_s \\ X'_s w_s - X'_s c \end{pmatrix} \right\|_V^2,$$

où

$$V = \begin{pmatrix} U_s & \mathbf{0} \\ \mathbf{0} & \alpha T \end{pmatrix}$$

et $\alpha \geq 0$. On minimise alors

$$\|w_s - A_s c_s\|_{U_s}^2 + \alpha \|X'_s w_s - X'c\|_T^2 = D_s(w_s) + \alpha \|X'_s w_s - X'c\|_T^2.$$

Un problème de minimisation semblable est rencontré lors d'une régression avec coefficients de coûts (ridge regression). Pour $\alpha = 0$ la solution est donnée par les poids de Horvitz-Thompson $w_s = A_s c_s$. Tandis que pour $\alpha > 0$, on cherche $w_s(\alpha)$ qui minimise $\|K(w_s - A_s c_s) - b\|_V^2$, où $K = (I_n, X'_s)'$, $b = (\mathbf{0}_{1 \times n}, (X'c - X'_s A_s c_s)')'$ et $\mathbf{0}_{1 \times n} \in \mathbb{R}^n$ est un vecteur-ligne de 0. Ben-Israel et Greville (1980) donne

$$w_s(\alpha) - A_s c_s = (K'VK)^{-1} K'Vb. \quad (10)$$

Donc on obtient en substituant les valeurs de K , V , et b

$$w_s(\alpha) = A_s c_s + \alpha (U_s + \alpha X_s T X'_s)^{-1} X_s T (X'c - X'_s A_s c_s). \quad (11)$$

On montre facilement que

$$\alpha (U_s + \alpha X_s T X'_s)^{-1} X_s T = U_s^{-1} X_s (\alpha^{-1} T^{-1} + X'_s U_s^{-1} X_s)^{-1},$$

d'où

$$w_s(\alpha) = A_s c_s + U_s^{-1} X_s (\alpha^{-1} T^{-1} + X'_s U_s^{-1} X_s)^{-1} (X'c - X'_s A_s c_s). \quad (12)$$

L'estimateur $Y'_s w_s(\alpha)$ prend donc la forme $\hat{Y}'c + (Y_s - \hat{Y}'_s)' A_s c_s$, où $\hat{Y} = X \hat{\beta}_s(\alpha)$ et

$$\hat{\beta}_s(\alpha) = (X'_s U_s^{-1} X_s + \alpha^{-1} T^{-1})^{-1} X'_s U_s^{-1} Y_s.$$

Le vecteur des coefficients de régression est donc celui qu'on obtient lors d'une régression avec coefficients de coûts (ridge regression). Tout comme la méthode d'estimation par calage et la méthode d'estimation par régression généralisée telle que décrite dans Särndal, Swensson et Wretman (1992), conduisent aux mêmes estimateurs, un parallèle semblable peut être fait entre le calage mitigé et la régression avec coefficients de coûts.

À partir de l'équation (12), on peut également utiliser Ben-Israel et Greville (1980), et le fait que $F^\dagger = F'(FF')^\dagger$ avec $F = T^{1/2} X'_s U_s^{-1/2}$, pour montrer que

$$\lim_{\alpha \rightarrow \infty} w_s(\alpha) = w_{cal}.$$

C'est le résultat auquel on devait s'attendre, puisque trouver le vecteur $w_s(\alpha)$ qui minimise $D_s(w_s) + \alpha \|X'_s w_s - X'c\|_T^2$ lorsque $\alpha \rightarrow \infty$, revient à trouver le vecteur de poids qui minimise $D_s(w_s)$ parmi ceux qui minimisent $\|X'_s w_s - X'c\|_T^2$.

Rao et Singh (1997) ont défini des tolérances pour chacune des p contraintes de l'équation de calage, et ils ont établi une relation entre ces tolérances et la matrice αT .

Pour $\alpha \in [0, \infty[$ la fonction $w_s(\alpha)$ décrit une courbe dans \mathbb{R}^n qui relie le point $A_s c_s$ au point w_{cal} . Si $p = 1$, c'est-à-dire si X est un vecteur, cette courbe est en fait un segment de droite. En effet, dans ce cas la matrice $(\alpha^{-1} T^{-1} + X'_s U_s^{-1} X_s)^{-1}$ et le vecteur $X'c - X'_s A_s c_s$ sont des scalaires, les poids $w_s(\alpha)$ donnés par (12) sont donc égaux aux poids de Horvitz-Thompson plus un multiple du vecteur $U_s^{-1} X_s$. Puis toujours pour $p = 1$, on a

$$\lim_{\alpha \rightarrow \infty} w_s(\alpha) = w_{cal} = A_s c_s + [(X'c - X'_s A_s c_s) / (X'_s U_s^{-1} X_s)] U_s^{-1} X_s$$

ce qui donne l'estimateur

$$Y'_s w_{cal} = Y'_s A_s c_s + [(Y'_s U_s^{-1} X_s) / (X'_s U_s^{-1} X_s)] (X'c - X'_s A_s c_s)$$

Si on prend $U = A^{-1} \text{diag}(X)$, on obtient l'estimateur par le quotient

$$Y'_s A_s c_s + [(Y'_s A_s \mathbf{1}_{n \times 1}) / (X'_s A_s \mathbf{1}_{n \times 1})] (X'c - X'_s A_s c_s),$$

où $\mathbf{1}_{a \times b} \in \mathbb{R}^{a \times b}$ est une matrice de 1.

Ben-Israel et Greville (1980, 111, exercice 15) montrent que $D_s(w_s(\alpha))$ est une fonction monotone croissante de α . Il faut cependant noter que pour une unité $k \in s$, $|w_k(\alpha) - a_k c_k|$ n'est pas nécessairement une fonction monotone de α . Lorsque α augmente, le vecteur de poids $w_s(\alpha)$ s'éloigne du vecteur des poids de Horvitz-Thompson, mais ce n'est pas nécessairement le cas pour chaque coordonnée de ce vecteur.

Dans cet article on utilisera le calage mitigé pour restreindre les poids, donc lorsque la taille d'échantillon est relativement petite. Il est cependant facile de montrer que pour un cadre asymptotique qui satisfait (2) et pour lequel $\hat{\beta}_s(\alpha) - \beta(\alpha) \rightarrow \mathbf{0}$ en probabilité, avec

$$\beta(\alpha) = (X'U^{-1}X + \alpha^{-1}T^{-1})^{-1} X'U^{-1}Y,$$

on a $Y'_s w_s(\alpha)$ est un estimateur asymptotiquement non biaisé dont la variance asymptotique est

$$(Y - Y^*)' \text{diag}(c) (A \Pi A - \mathbf{1}_{N \times N}) \text{diag}(c) (Y - Y^*),$$

où $Y^* = X\beta(\alpha)$, Π est la matrice des probabilités d'inclusion de second ordre, et $\text{diag}(c)$ est la matrice diagonale formée à partir du vecteur c .

5. Méthodes d'estimation avec poids restreints

Afin d'éviter d'obtenir des poids ayant des valeurs extrêmes, on peut vouloir forcer le vecteur de poids à être à l'intérieur d'une région déterminée. On supposera que cette région de restriction est convexe et fermée, et que $A_s c_s$ est un point de cette région. Par exemple, si $w^{(L)} < A_s c_s < w^{(H)}$, on pourrait vouloir restreindre les poids à la région $R_w = \{w_s : w^{(L)} \leq w_s \leq w^{(H)}\}$. On supposera que

$$\lim_{n \rightarrow \infty} w^{(L)} - A_s c_s < 0 \text{ et } \lim_{n \rightarrow \infty} w^{(H)} - A_s c_s > 0.$$

L'approche mentionnée à la section 3 consiste à choisir une mesure de distance entre les poids calés et les poids de Horvitz-Thompson qui donnera des poids satisfaisant l'équation de calage qui sont dans la région de restriction, ceci pourvu que de tels poids existent. L'approche qui sera étudiée dans cette section consiste à tempérer l'exigence de satisfaire l'équation de calage lorsque le vecteur des poids de calage w_{cal} est à l'extérieur de la région de restriction. Différentes façons de tempérer cette exigence conduisent à différentes méthodes de pondération.

Lorsque w_{cal} est à l'extérieur de la région de restriction, on pourrait par exemple, chercher les points de la courbe $w_s(\alpha)$ paramétrisée par $\alpha \geq 0$ qui sont sur la frontière de cette région. Ces points ont la propriété d'être une solution du problème de minimisation décrit dans la section 4 pour les valeurs de α correspondantes, donc à travers la matrice T , on peut pondérer l'importance de chaque équation de calage. Avec l'exemple de la région de restriction donné plus haut, si

$$w_{\text{cal}} = \lim_{\alpha \rightarrow \infty} w_s(\alpha)$$

est à l'intérieur de cette région, alors on peut prendre comme vecteur de poids restreints $w_{\text{res1}} = w_{\text{cal}}$, sinon on peut prendre $w_{\text{res1}} = w_s(\alpha)$ pour $\alpha < \infty$ tel que $w_s(\alpha)$ est à la frontière de la région de restriction. Si le cadre asymptotique est tel que les conditions (2) sont respectées avec $\gamma < 3/2$ alors pour n suffisamment grand, la probabilité que w_{cal} soit à l'intérieur de la région de restriction est égale à un. En effet, on a $w_{\text{cal}} - A_s c_s$ converge en probabilité vers 0 . Les propriétés asymptotiques de l'estimateur qui utilise les poids restreints, w_{res1} , sont donc les mêmes que celles de l'estimateur par calage. Il est important de noter que puisque $|w_k - a_k c_k|$ n'est pas nécessairement une fonction monotone de α , il est possible que $w_s(\alpha)$ soit à la frontière de la région de restriction pour plusieurs valeurs de α , même si la région de restriction est convexe. Il n'est pas

nécessairement simple de trouver toutes ces valeurs, et il faut décider laquelle utiliser.

Une autre option pour restreindre les poids consisterait à se donner comme région de restriction, les poids w_s qui satisfont $D_s(w_s) \leq l$ pour une borne $l > 0$. On prend ensuite comme vecteur de poids restreints $w_{\text{res2}} = w_{\text{cal}}$, si w_{cal} est dans la région de restriction, sinon on cherche $\alpha > 0$ tel que $D_s(w_s(\alpha)) = l$. Cette valeur de α est unique et peut être trouvée de façon itérative. On calcule ensuite les poids $w_{\text{res2}} = w_s(\alpha)$ qui correspondent à cette valeur de α à l'aide de l'équation (12). Si le cadre asymptotique est tel que les conditions (2) sont respectées avec $\gamma < 1$, et si l ne varie pas avec n , alors pour n suffisamment grand, la probabilité que w_{cal} soit à l'intérieur de la région de restriction est égale à un. En effet, on a $D_s(w_{\text{cal}})$ converge en probabilité vers 0. Les propriétés asymptotiques de l'estimateur qui utilise les poids restreints, w_{res2} , sont alors les mêmes que celles de l'estimateur par calage. Malheureusement, lorsqu'on estime un total il faut s'attendre à avoir $\gamma = 1$. Pour pallier cet inconvénient, on peut utiliser $l\sqrt{n}$ comme borne, plutôt que l . On peut justifier cette borne par le fait que la longueur de la diagonale principale d'un hypercube de \mathbb{R}^n est égale au diamètre de la boule qui circonscrit cet hypercube, par contre, le diamètre de la boule inscrite dans ce même hypercube est plus petit par un facteur de \sqrt{n} . Il reste que le statisticien peut être inconfortable avec l'utilisation d'un cadre asymptotique où la borne croît avec la taille de l'échantillon. Il y a aussi le fait que cette approche ne permette pas de limiter individuellement les poids des observations. Seule la distance entre le vecteur des poids restreints et le vecteur des poids de Horvitz-Thompson est contrôlée.

Avec les méthodes décrites plus haut, on cherche les points de la courbe $w_s(\alpha)$ qui sont sur la frontière de la région de restriction. La valeur de α pour laquelle $w_s(\alpha)$ est à la frontière de la région de restriction doit souvent être trouvée de façon itérative. On pourrait plus simplement remplacer la courbe $w_s(\alpha)$ par le segment de droite reliant $A_s c_s$ à w_{cal} . Pour la région de restriction R_w , on aurait comme vecteur de poids restreints $w_{\text{res3}} = w_{\text{cal}}$, si w_{cal} est dans la région de restriction, et sinon w_{res3} serait égal au point où le segment de droite traverse la frontière de la région de restriction, c'est-à-dire

$$w_{\text{res3}} = A_s c_s + \xi(w_{\text{cal}} - A_s c_s),$$

où

$$\xi = \min_k \{ \max \{ (w^{(L)} - A_s c_s) / (w_{\text{cal}} - A_s c_s), (w^{(H)} - A_s c_s) / (w_{\text{cal}} - A_s c_s) \} \},$$

la division des vecteurs se faisant élément par élément, le maximum des deux vecteurs étant pris élément par élément, et \min donnant l'élément minimum. On pourrait aussi considérer le vecteur de poids de la région de restriction, w_{res4} ,

qui est le plus près de \mathbf{w}_{cal} . Toujours pour la région de restriction $R_{\mathbf{w}}$, on aurait

$$\mathbf{w}_{\text{res}4} = \min[\max(\mathbf{w}_{\text{cal}}, \mathbf{w}^{(L)}), \mathbf{w}^{(H)}].$$

Les propriétés asymptotiques des estimateurs qui utilisent les poids restreints $\mathbf{w}_{\text{res}3}$ ou $\mathbf{w}_{\text{res}4}$ sont les mêmes que celles de l'estimateur par calage, pourvu que $\mathbf{w}_{\text{cal}} - \mathbf{A}_s \mathbf{c}_s$ converge en probabilité vers $\mathbf{0}$, ce qui est normalement le cas.

Une propriété intéressante de toutes les approches discutées dans cette section est que quelle que soit la région de restriction, l'existence de poids restreints est assurée. Ce n'est pas toujours le cas avec une approche basée sur des mesures de distance. On va maintenant présenter un exemple simple pour comparer quelques-unes des méthodes.

On veut estimer un total à partir d'un échantillon aléatoire simple de taille 2 d'une population de taille 20. C'est-à-dire $\mathbf{c} = \mathbf{1}_{20 \times 1}$ et $\mathbf{a} = 10(\mathbf{1}_{20 \times 1})$. On utilise le vecteur d'information auxiliaire $\mathbf{X} = (1, 2, 3, \dots, 20)'$, on suppose que l'échantillon obtenu est $s = \{2, 12\}$ et on choisit de prendre \mathbf{U} une matrice diagonale avec $u_{kk} = x_k = k$. On se donne une région de restriction rectangulaire à l'aide des points $\mathbf{w}^{(L)} = (0, 0)'$ et $\mathbf{w}^{(H)} = (20, 13)'$. C'est-à-dire que le poids de la première unité de l'échantillon doit être supérieur à 0 et inférieur à 20, tandis que le poids de la seconde unité de l'échantillon doit être supérieur à 0 et inférieur à 13.

Sous ces conditions, les poids calés $\mathbf{w}_{\text{cal}} = (15, 15)'$ sont à l'extérieur de la région de restriction. Puisque $p = 1$, les poids $\mathbf{w}_s(\alpha)$ sont sur le segment de droite qui relie $\mathbf{A}_s \mathbf{c}_s = (10, 10)'$ à \mathbf{w}_{cal} . On a donc $\mathbf{w}_{\text{res}1} = \mathbf{w}_{\text{res}3}$, c'est-à-dire que les deux méthodes donnent le même résultat. Dans ce cas-ci on a $\mathbf{w}_{\text{res}1} = \mathbf{w}_{\text{res}3} = (13, 13)'$. La méthode qui consiste à choisir le point de la région de restriction le plus près des poids calés donne $\mathbf{w}_{\text{res}4} = (15, 13)'$. Par contre, si on cherche $\mathbf{w}_{\text{res}5}$, les poids restreints obtenus en continuant d'exiger que l'équation de calage soit satisfaite et en utilisant une mesure de distance qui prend une valeur infinie à l'extérieur de la région de restriction, alors il n'y a pas de solution. En effet, pour tout poids dans la région de restriction $\mathbf{X}'_s \mathbf{w}_s \leq 196$, tandis que $\mathbf{X}'\mathbf{c} = 210$. Si on avait, disons $\mathbf{w}^{(H)} = (30, 13)'$, alors avec $D_s(\mathbf{w}_s)$ comme mesure de distance à l'intérieur de la région de restriction on aurait $\mathbf{w}_{\text{res}5} = (27, 13)'$. Ces poids sont plutôt éloignés de $\mathbf{w}_{\text{cal}} = (15, 15)'$ et de $\mathbf{A}_s \mathbf{c}_s = (10, 10)'$. C'est le prix qu'il faut payer si on tient à avoir des poids qui respectent l'équation de calage.

6. Estimateurs pour domaines avec une composante synthétique

On utilise des poids restreints à cause des propriétés de l'estimateur par calage pour de petites tailles d'échantillon. Pour de grandes tailles d'échantillon, on sait qu'on a normalement $\mathbf{w}_{\text{cal}} - \mathbf{A}_s \mathbf{c}_s$ qui converge en probabilité vers $\mathbf{0}$, donc des poids qui ne sont pas problématiques. Un

statisticien qui est confronté au problème de poids extrêmes doit donc vraisemblablement faire face à un autre problème lié aux petites tailles d'échantillon, à savoir l'estimation pour de petits domaines. On présentera dans cette section, un estimateur dont les propriétés asymptotiques sont celles de l'estimateur par calage, mais qui utilise des poids restreints et tire profit d'un estimateur synthétique.

Soit $\tilde{\mathbf{Y}} = \mathbf{X}\tilde{\boldsymbol{\beta}}_s$ une estimation synthétique pour \mathbf{Y} , on a

$$\begin{aligned} \tilde{\mathbf{Y}}'\mathbf{w}_s &= (\mathbf{X}_s\tilde{\boldsymbol{\beta}}_s)'\mathbf{w}_s \\ &= \tilde{\boldsymbol{\beta}}_s'\mathbf{X}'_s\mathbf{w}_s \\ &\approx \tilde{\boldsymbol{\beta}}_s'\mathbf{X}'\mathbf{c} \\ &= (\mathbf{X}\tilde{\boldsymbol{\beta}}_s)'\mathbf{c} \\ &= \tilde{\mathbf{Y}}'\mathbf{c} \end{aligned} \quad (13)$$

avec égalité pour la troisième expression de droite si les poids satisfont l'équation de calage $\mathbf{X}'_s\mathbf{w}_s = \mathbf{X}'\mathbf{c}$. Les poids \mathbf{w}_{cal} donnés par (1) minimisent $\|\mathbf{X}'_s\mathbf{w}_s - \mathbf{X}'\mathbf{c}\|_T^2$. On peut donc estimer $\mathbf{Y}'\mathbf{c}$ par

$$\hat{\boldsymbol{\tau}} = (\mathbf{Y}_s - \tilde{\mathbf{Y}}_s)'\mathbf{w}_{\text{res}} + \tilde{\mathbf{Y}}'\mathbf{c}. \quad (14)$$

Il y aura égalité entre cet estimateur et l'estimateur $\mathbf{Y}'_s\mathbf{w}_{\text{cal}}$ dès que l'échantillon sera suffisamment grand pour que l'équation de calage soit satisfaite et pour que \mathbf{w}_{cal} soit dans la région de restriction, c'est-à-dire dès que $\mathbf{w}_{\text{res}} = \mathbf{w}_{\text{cal}}$. Les propriétés asymptotiques de ces deux estimateurs sont donc les mêmes sous certaines conditions discutées dans la section précédente. L'avantage de l'estimateur $\hat{\boldsymbol{\tau}}$ est qu'il donne une estimation synthétique lorsque des colonnes de \mathbf{Y}_s et $\tilde{\mathbf{Y}}_s$ sont nulles.

7. Données aberrantes

On pourrait traiter les données aberrantes d'une façon similaire aux poids extrêmes. La stratégie est la suivante : on se donne une région de restriction pour $\mathbf{Y}'_s\mathbf{w}_{\text{cal}}$, on montre que pour n suffisamment grand $\mathbf{Y}'_s\mathbf{w}_{\text{cal}}$ est à l'intérieur de cette région de restriction, et on se donne un estimateur plus « raisonnable » pour remplacer $\mathbf{Y}'_s\mathbf{w}_{\text{cal}}$ dans les cas où $\mathbf{Y}'_s\mathbf{w}_{\text{cal}}$ est à l'extérieur de la région de restriction. Dans le cas d'un échantillon stratifié on aurait normalement une région de restriction par strate.

On a montré à la section 2 que sous certaines conditions sur le cadre asymptotique, $\mathbf{w}_{\text{cal}} - \mathbf{A}_s \mathbf{c}_s = O_p(n^{-3/2}N^\gamma)$. On a donc $\mathbf{Y}'_s\mathbf{w}_{\text{cal}} - \mathbf{Y}'_s\mathbf{A}_s \mathbf{c}_s = O_p(n^{-1/2}N^\gamma)$, et si on suppose que

$$\mathbf{Y}'_s\mathbf{A}_s \mathbf{c}_s - \mathbf{Y}'\mathbf{c} = O_p(n^{-1/2}N^\gamma), \quad (15)$$

alors $\mathbf{Y}'_s\mathbf{w}_{\text{cal}} - \mathbf{Y}'\mathbf{c} = O_p(n^{-1/2}N^\gamma)$. Un expert (ou un groupe d'experts) peut déterminer à partir d'informations indépendantes des données d'enquête, qu'il ne serait pas

raisonnable d'avoir $Y'_s w_{\text{cal}}$ à l'extérieur d'une certaine région. Si $Y'c$ est dans la région de restriction (c'est-à-dire si l'expert ne juge pas déraisonnable une estimation du paramètre qui serait égale à la vraie valeur, $Y'c$, du paramètre), si $\gamma = 0$, et si la région de restriction ne varie pas avec n ou N (ou si $\gamma = 1$, et la région de restriction varie proportionnellement à N), alors pour n suffisamment grand, la probabilité que $Y'_s w_{\text{cal}}$ soit à l'intérieur de la région de restriction est égale à un. Dans les cas où $Y'_s w_{\text{cal}}$ est à l'extérieur de la région de restriction, on pourrait utiliser comme estimation le point de la région de restriction qui est le plus près de $Y'_s w_{\text{cal}}$ ou on pourrait poser égal à un le poids des quelques observations qui sont jugées aberrantes, et répartir leurs poids originaux (moins le nombre d'observations aberrantes) sur les observations non aberrantes. Les propriétés asymptotiques de cet estimateur modifié pour traiter les valeurs aberrantes sont alors les mêmes que celles de l'estimateur non modifié.

Dans le cas d'un échantillon non stratifié cette méthode est relativement facile à appliquer. Si par contre l'échantillon est stratifié, et si des contraintes sont imposées aux estimations de chaque strate, alors on doit faire face à deux problèmes additionnels. Premièrement, si le cadre asymptotique est tel que le nombre de strates augmente de façon proportionnelle à la taille d'échantillon, alors la supposition donnée en (15) ne tient pas, puisque la taille moyenne d'échantillon par strate reste constante alors que $n \rightarrow \infty$. Il s'agit de savoir s'il est raisonnable de se donner un cadre asymptotique où le nombre de strates est constant (ou croît moins vite que n). Un tel cadre asymptotique est moins plausible si le nombre d'observations par strate est petit. Le deuxième problème est la difficulté pour l'expert d'imposer des contraintes aux estimations pour chacune des strates. Plus il y a de strates, plus le risque est grand que $Y'c$ ne soit pas dans la région de restriction définie par l'expert. En fait, dans le cas d'un échantillon stratifié, il est préférable que l'expert utilise les informations indépendantes des données d'enquêtes, afin de s'assurer de l'homogénéité des strates, avant que la stratification soit finalisée. En d'autres mots il vaut mieux utiliser l'information qui est disponible avant l'enquête, pour prévenir les données aberrantes, plutôt que pour les corriger. Si l'information a été utilisée de sorte qu'avant l'enquête, on n'a aucune raison de croire qu'il y a dans quelque strate une observation non représentative, alors on n'a aucune justification pour prétendre le contraire après la collecte des données.

8. Conclusion

Si pour de grandes tailles d'échantillons, les poids calés demeurent à l'intérieur d'une région de restriction, alors les propriétés asymptotiques de l'estimateur avec poids restreints sont bien sûr identiques à celles de l'estimateur par calage. Pour un cadre asymptotique donné, on peut habituellement s'attendre à avoir $w_{\text{cal}} - A_s c_s$ qui converge en probabilité vers 0 , donc pour des tailles d'échantillons

suffisamment grandes, les poids calés w_{cal} demeureront à l'intérieur de la région de restriction R_w si $A_s c_s$ est à l'intérieur de R_w . Cependant, on a vu que pour l'estimation d'un total, on n'a pas nécessairement convergence vers 0 de $D_s(w_{\text{cal}})$. Une région de restriction définie par $\|w_s - A_s c_s\|_{U_s}^2 \leq l$ est donc à éviter, à tout le moins si on estime un total plutôt qu'une moyenne.

On a donné des conditions nécessaires et suffisantes pour l'existence de poids restreints à des intervalles qui satisfont l'équation de calage. Si de tels poids n'existent pas, il faut alors abandonner l'idée de satisfaire exactement l'équation de calage. On peut reformuler le problème de calage avec poids restreints, de manière à ce qu'une solution soit toujours possible. Quelques-unes des approches présentées dans cet article permettent d'obtenir une solution sans avoir recours à des méthodes itératives. Il s'agit de méthodes simples, faciles à interpréter. Les propriétés asymptotiques de ces estimateurs sont habituellement identiques à celles de l'estimateur par calage sans restrictions sur les poids.

Le problème de poids extrêmes survient avec des tailles d'échantillon qui sont petites, donc le problème d'estimation pour de petits domaines devrait être envisagé en même temps. Il est possible de tirer profit d'estimateurs synthétiques tout en ayant un estimateur avec des poids restreints qui a de bonnes propriétés asymptotiques.

On peut aussi modifier l'estimateur par calage, ou tout autre estimateur asymptotiquement convergent, pour traiter les données aberrantes. Les conditions pour que cet estimateur modifié ait les mêmes propriétés asymptotiques que l'estimateur non modifié ne sont pas facilement vérifiables, tout comme il est difficile de vérifier qu'une donnée aberrante est effectivement aberrante (non représentative). Cependant, ces conditions permettent d'identifier les facteurs qui font qu'un estimateur corrigé pour données aberrantes peut être statistiquement valable.

Remerciements

L'auteur tient à remercier un éditeur associé et un arbitre pour leurs remarques constructives qui ont permis d'améliorer cet article.

Annexe A

On veut vérifier que $\Omega(\phi) = l'(V\phi)_- - h'(V\phi)_+$ est inférieur ou égal à zéro. Premièrement, il est facile de montrer que ceci est vrai pour un vecteur ϕ , si et seulement si c'est vrai pour un vecteur $k\phi$ avec $k > 0$ arbitraire. Seulement la direction de ϕ importe. Il est donc suffisant de vérifier la condition pour ϕ de norme égale à un. Pour la preuve, on utilisera la norme 1_1 de $\|\phi\|_{1_1} = \sum_{i=1}^p |\phi_i|$. Les vecteurs ϕ avec $\|\phi\|_{1_1} = 1$ sont situés dans des hyperplans dont les intersections sont à des points perpendiculaires aux vecteurs unité, c'est-à-dire des points dont au

moins une des coordonnées est zéro. La fonction Ω varie linéairement sauf à des points $\boldsymbol{\varphi}$ perpendiculaires à une ou plusieurs rangées de V . Même lorsque le domaine de Ω est restreint aux vecteurs $\boldsymbol{\varphi}$ avec $\|\boldsymbol{\varphi}\|_{l_i}=1$ qui sont perpendiculaires à $0 \leq j < (p-1)$ rangées de V linéairement indépendantes, la fonction Ω varie encore linéairement sauf à des points perpendiculaires à d'autres rangées de V ou perpendiculaires à des vecteurs unité (lesquels sont également des rangées de V). Le maximum de Ω pour $\|\boldsymbol{\varphi}\|_{l_i}=1$ est donc atteint en un point $\boldsymbol{\varphi}$ perpendiculaire à $(p-1)$ rangées de V linéairement indépendantes. Il est donc suffisant de vérifier la condition pour deux vecteurs de direction opposée qui sont perpendiculaires à $(p-1)$ rangées de V linéairement indépendantes et ce, pour chaque sous-ensemble de $(p-1)$ rangées de V linéairement indépendantes.

Annexe B

On note $\text{vec}(F)$ le vecteur obtenu en empilant les colonnes successives de la matrice $F \in \mathbb{R}^{a \times b}$ avec la première colonne au-dessus, et on définit le produit de Kronecker de deux matrices F et G comme

$$F \otimes G = \begin{pmatrix} f_{11}G & \cdots & f_{1n}G \\ \vdots & & \vdots \\ f_{m1}G & \cdots & f_{mn}G \end{pmatrix}. \quad (B1)$$

Le résultat découle du corollaire de la section 3 avec

$$M = \begin{pmatrix} I_{RC} \\ I_R \otimes \mathbf{1}_{1 \times C} \\ \mathbf{1}_{1 \times R} \otimes I_C \\ \mathbf{1}_{1 \times RC} \end{pmatrix}, \quad \mathbf{w} = \text{vec}((\hat{N}_{ij})'),$$

$$l = \begin{pmatrix} \text{vec}((N_{ij}^{(L)})') \\ N_{1.}^{(L)} \\ \vdots \\ N_{R.}^{(L)} \\ N_{.1}^{(L)} \\ \vdots \\ N_{.C}^{(L)} \\ N_{..}^{(L)} \end{pmatrix}, \quad \mathbf{h} = \begin{pmatrix} \text{vec}((N_{ij}^{(H)})') \\ N_{1.}^{(H)} \\ \vdots \\ N_{R.}^{(H)} \\ N_{.1}^{(H)} \\ \vdots \\ N_{.C}^{(H)} \\ N_{..}^{(H)} \end{pmatrix}. \quad (B2)$$

Seul un ensemble fini de conditions doit être vérifié, premièrement en notant que les colonnes de

$$V = \begin{pmatrix} -I_R \otimes \mathbf{1}_{C \times 1} & -\mathbf{1}_{R \times 1} \otimes I_C & -\mathbf{1}_{RC \times 1} \\ I_R & \mathbf{0}_{R \times C} & \mathbf{0}_{R \times 1} \\ \mathbf{0}_{C \times R} & I_C & \mathbf{0}_{C \times 1} \\ \mathbf{0}_{1 \times R} & \mathbf{0}_{1 \times C} & 1 \end{pmatrix} \quad (B3)$$

forment une base pour $N(M')$. C'est-à-dire que $M'V = \mathbf{0}$, les colonnes de V sont linéairement indépendantes, et $N(M')$ est de dimension $R+C+1$. En notant également que les dernières $R+C+1$ lignes de V sont les vecteurs unité. Et finalement, en vérifiant les conditions du corollaire pour tous les vecteurs $\boldsymbol{\lambda} = V\boldsymbol{\varphi}$ et $\boldsymbol{\lambda} = -V\boldsymbol{\varphi}$, où $\boldsymbol{\varphi}$ est perpendiculaire à $R+C$ lignes de V linéairement indépendantes. Cette dernière étape est élaborée plus en détails dans le paragraphe suivant.

Un sous-ensemble arbitraire de $R+C$ lignes de V linéairement indépendantes qui inclut la dernière ligne de V sera noté L , et le sous-ensemble de toutes les lignes de V qui sont des combinaisons linéaires de lignes de L sera noté L^+ . Si L^+ inclut la ligne $RC+i$ ($i=1, \dots, R$) si et seulement si $i \notin S \subseteq \{1, 2, \dots, R\}$, et inclut la ligne $RC+R+j$ ($j=1, \dots, C$) si et seulement si $j \notin T \subseteq \{1, 2, \dots, C\}$, alors on pose $\boldsymbol{\varphi} = (\boldsymbol{\varphi}'_S, -\boldsymbol{\varphi}'_T, 0)'$, où le $i^{\text{ième}}$ élément de $\boldsymbol{\varphi}'_S \in \mathbb{R}^R$ est égal à 1 si $i \in S$ et à zéro sinon, et le $j^{\text{ième}}$ élément de $\boldsymbol{\varphi}'_T \in \mathbb{R}^C$ est égal à un si $j \in T$ et à zéro sinon. Alors

$$V\boldsymbol{\varphi} = ((-\boldsymbol{\varphi}'_S \otimes \mathbf{1}_{C \times 1} + \mathbf{1}_{R \times 1} \otimes \boldsymbol{\varphi}'_T), \boldsymbol{\varphi}'_S, -\boldsymbol{\varphi}'_T, 0)',$$

donc $\boldsymbol{\varphi}$ est perpendiculaire à toutes les lignes de L^+ , et à plus forte raison, $\boldsymbol{\varphi}$ est perpendiculaire à toutes les lignes de L . De même, le vecteur $\boldsymbol{\varphi}^* = (\boldsymbol{\varphi}'_S, \boldsymbol{\varphi}'_T, -1)'$ est perpendiculaire à toutes les lignes d'un sous-ensemble de $R+C$ lignes de V linéairement indépendantes qui inclut la ligne $RC+i$ ($i=1, \dots, R$) si et seulement si $i \notin S$, et inclut la ligne $RC+R+j$ ($j=1, \dots, C$) si et seulement si $j \notin T$, mais qui n'inclut pas la dernière ligne de V . La condition $-l'\boldsymbol{\lambda}_- \leq h'\boldsymbol{\lambda}_+$ avec $\boldsymbol{\lambda} = V\boldsymbol{\varphi}$ donne le cinquième ensemble d'inéquations dans (9). De même, en posant $\boldsymbol{\lambda}$ égal à $-V\boldsymbol{\varphi}$, $V\boldsymbol{\varphi}^*$ et $-V\boldsymbol{\varphi}^*$ on obtient les trois derniers ensembles d'inéquations dans (9).

Bibliographie

- Bacharach, M. (1965). Estimating nonnegative matrices from marginal data. *International Economic Review*, 6, 294-310.
- Bardsley, P., et Chambers, R.L. (1984). Multipurpose estimation from unbalanced samples. *Applied Statistics*, 33, 290-299.
- Ben-Israel, A., et Greville, T.N.E. (1980). *Generalized Inverses: Theory and Applications*. Huntington, New York : Robert E. Krieger Publishing Company.

- Brewer, K.R.W. (1979). A class of robust sampling designs for large scale surveys. *Journal of the American Statistical Association*, 74, 911-915.
- Deville, J.-C., et Särndal, C.-E. (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association*, 87, 376-382.
- Duchesne, P. (1999). Estimateurs de calage robustes. *Techniques d'enquête*, 25, 47-60.
- Fan, K. (1956). On systems of linear inequalities. *Annals of Mathematics Studies*, (Éds. H. W. Kuhn, et A. W. Tucker), 38, 99-156.
- Graybill, F.A. (1983). *Matrices with Applications in Statistics*, (Deuxième édition). Belmont, California: Wadsworth Publishing.
- Isaki, C.T., et Fuller, W.A. (1982). Survey design under the regression superpopulation model. *Journal of the American Statistical Association*, 77, 89-96.
- Rao, J.N.K., et Singh, A.C. (1997). A ridge-shrinkage method for range-restricted weight calibration in survey sampling. *Proceedings of the Section on Survey Research Methods, American Statistical Association*.
- Särndal, C.-E., Swensson, B. et Wretman, J. (1992). *Model Assisted Survey Sampling*. New York : Springer-Verlag.
- Théberge, A. (1999). Extensions of calibration estimators in survey sampling. *Journal of the American Statistical Association*, 94, 635-644.