

ESTIMATEURS DE LA VARIANCE PAR LINÉARISATION POUR DES DONNÉES D'ENQUÊTE AVEC DES RÉPONSES MANQUANTES

Abdellatif DEMNATI ^(*), J. N. K. RAO ^(**)

^(*) Statistique Canada, Division des méthodes d'enquêtes sociales

^(**) Université Carlton, Étude des Mathématiques et des Statistiques

Introduction

La linéarisation de Taylor est une méthode d'estimation de la variance fréquemment utilisée pour des statistiques complexes, comme les estimateurs par quotient ou par régression, ainsi que les estimateurs des coefficients de régression logistique. Elle est généralement applicable à tout plan d'échantillonnage qui permet d'obtenir une estimation non biaisée de la variance d'estimateurs linéaires, et les calculs sont plus simples que ceux des méthodes de rééchantillonnage, comme le jackknife. Cependant, elle peut donner plusieurs estimateurs de la variance asymptotiquement non biaisés par rapport au plan d'échantillonnage en cas d'échantillonnage répété. Par conséquent, pour choisir l'estimateur de la variance, il faut tenir compte d'autres facteurs, dont (i) l'absence approximative du biais de la variance de l'estimateur sous le modèle, pour un modèle supposé et (ii) la validité dans des conditions d'échantillonnage répété conditionnel. Par exemple, dans le contexte de l'échantillonnage aléatoire simple et de l'estimateur par quotient, $\hat{Y}_R = (\bar{y}/\bar{x})X$, du total de la population Y , Royall et Cumberland [5] ont montré qu'un estimateur par linéarisation utilisé couramment, $v_L = N^2(n^{-1} - N^{-1})s_z^2$, ne permet pas de suivre la variance conditionnelle de \hat{Y}_R étant donné \bar{x} , contrairement à l'estimateur jackknife de la variance, v_J . Ici, \bar{y} et \bar{x} sont les moyennes d'échantillon, X est le total connu de la population d'une variable auxiliaire x , s_z^2 est la variance dans l'échantillon des résidus $z_i = y_i - (\bar{y}/\bar{x})x_i$ et (n, N) représente les tailles de l'échantillon et de la population. Si nous linéarisons l'estimateur jackknife de la variance, v_J , nous obtenons un estimateur par linéarisation de la variance différent, $v_{JL} = (\bar{X}/\bar{x})^2 v_L$, qui suit la variance conditionnelle ainsi que la variance non conditionnelle, où $\bar{X} = X/N$ est la moyenne de x . Par conséquent, on pourrait préférer v_{JL} ou v_J à v_L . Yung et Rao [9] ont considéré des estimateurs par régression généralisée et des estimateurs ajustés par quotient stratifiés a posteriori sous un plan stratifié à plusieurs degrés et ont obtenu un estimateur par linéarisation jackknife de la variance, v_{JL} , par la linéarisation de v_J . Valliant [8] a aussi obtenu v_{JL} pour l'estimateur stratifié a posteriori et a réalisé une étude de simulation pour démontrer que v_J et v_{JL} possèdent tous deux de bonnes propriétés

conditionnelles, étant donné les estimations des tailles des strates. Särndal, Swensson et Wretman [6] ont montré que v_{JL} est à la fois asymptotiquement non biaisé par rapport au plan de sondage et asymptotiquement non biaisé par rapport au modèle au sens où $E_m(v_{JL}) = Var_m(\hat{Y}_R)$, où E_m représente l'espérance par rapport au modèle et $Var_m(\hat{Y}_R)$ est la variance de \hat{Y}_R par rapport au modèle sous un modèle par quotient? $E_m(y_i) = \mathbf{b} x_i$; $i = 1, \dots, N$ et les y_i sont indépendants avec comme variance du modèle $Var_m(y_i) = \mathbf{s}^2 x_i$, $\mathbf{s}^2 > 0$. Donc, v_{JL} est un bon choix tant du point de vue du plan d'échantillonnage que du point de vue du modèle. Demnati et Rao [2] ont proposé une nouvelle méthode d'estimation de la variance qui est justifiable théoriquement et qui donne directement un estimateur de la variance de type v_{JL} pour des plans d'échantillonnage généraux. Cette méthode est présentée à la deuxième section.

En cas de réponses manquantes, on procède souvent à l'ajustement des poids pour compenser la non-réponse complète, tandis que l'on a recourt couramment à l'imputation pour compenser la non-réponse partielle afin d'obtenir des données complètes pour calculer des estimations. Cependant, traiter les poids ajustés comme des poids de sondage et les valeurs imputées comme des valeurs réelles, et appliquer des formules standard d'estimation de la variance peut produire une sous-estimation considérable si le taux de non-réponse est important. Ces dernières années, on a proposé plusieurs méthodes pour estimer correctement la variance d'un estimateur en cas d'imputation. Rao [4], Shao et Steel [7] et d'autres ont étudié l'estimation de la variance sous imputation par quotient, cependant le problème d'estimation de la variance en cas d'ajustement de la pondération et d'imputation reste ouvert.

L'objectif principal du présent article est d'étendre la méthode de Demnati et Rao [2] au cas des réponses manquantes, lorsqu'on procède à la correction pour la non-réponse totale et à l'imputation basée sur des fonctions lisses des valeurs observées, en particulier l'imputation par quotient. À la deuxième section, nous décrivons brièvement la méthode dans le cas de la réponse complète. À la troisième section, nous décrivons la méthode de Shao et Steel [7], tandis qu'à la quatrième section, nous présentons l'extension de la méthode de Demnati-Rao.

1. Réponse complète

Pour motiver l'application de la méthode de Demnati-Rao [2] en cas de réponse complète, supposons qu'un estimateur $\hat{\mathbf{q}}$ d'un paramètre \mathbf{q} s'exprime sous forme d'une fonction dérivable $g(\underline{Y})$ des totaux estimés $\underline{Y} = (\hat{Y}_1, \dots, \hat{Y}_m)^T$, où $\hat{Y}_j = \sum_{i \in U} d_i(s) y_{ij}$ est un estimateur du total de la population Y_j , $j = 1, \dots, m$ et $d_i(s) = 0$ si l'unité i ne fait pas partie de l'échantillon s , U est l'ensemble des unités de la population et $\mathbf{q} = g(\underline{Y})$, avec $\underline{Y} = (Y_1, \dots, Y_m)^T$. Nous pouvons écrire $\hat{\mathbf{q}}$ sous la forme $\hat{\mathbf{q}} = f(\underline{d}(s), \underline{A}_y)$ et $\mathbf{q} = f(\underline{1}, \underline{A}_y)$, où \underline{A}_y est une matrice $m \times N$ dont la j^e colonne est $\underline{y}_j = (y_{1j}, \dots, y_{mj})^T$, $j = 1, \dots, m$, $\underline{d}(s) = (d_1(s), \dots, d_N(s))^T$ et $\underline{1}$ est le vecteur de taille N dont toutes les coordonnées valent 1. Par exemple, si $\hat{\mathbf{q}}$ représente l'estimateur par quotient $\hat{Y}_R = [\sum_{i \in U} d_i(s) y_i / \sum_{i \in U} d_i(s) x_i] X$, alors $m = 2$, $y_{1i} = y_i$, $y_{2i} = x_i$ et $f(\underline{1}, \underline{A}_y)$ se réduit au total Y , en remarquant que $(Y/X)X = Y$. Notons que \hat{Y}_R est une fonction de $\underline{d}(s)$, \underline{y} et \underline{x} , et du total connu X , mais que, par souci de simplicité, nous laissons tomber X et écrivons $\hat{Y}_R = f(\underline{d}(s), \underline{y}, \underline{x})$. Si nous utilisons les coefficients de pondération d'Horvitz-Thompson,

alors $d_i(s) = 1/p_i$ pour tout $i \in s$, où p_i est la probabilité de sélectionner l'unité i dans l'échantillon s .

Soit $\hat{Y} = \sum b_i y_i$ pour des nombres réels arbitraires $\underline{b} = (b_1, \dots, b_N)^T$, et $f(\underline{b}, A_y) = f(\underline{b})$. Demnati et Rao [2] ont montré que la linéarisation de Taylor de $\hat{\mathbf{q}} - \mathbf{q}$, c'est-à-dire

$$\hat{\mathbf{q}} - \mathbf{q} = g(\hat{Y}) - g(Y) \approx \left(\frac{\partial g(\underline{a})}{\partial \underline{a}} \right)^T \Big|_{\underline{a}=Y} (\hat{Y} - Y),$$

est équivalente à

$$\begin{aligned} \hat{\mathbf{q}} - \mathbf{q} &\approx \sum_{k=1}^N \left(\frac{\partial f(\underline{b})}{\partial b_k} \right) \Big|_{b=1} (d_k(s) - 1) \\ &= \tilde{\mathbf{z}}^T (\underline{d}(s) - \underline{1}) \end{aligned} \quad (1.1)$$

où $\frac{\partial g(\underline{a})}{\partial \underline{a}} = (\frac{\partial g(\underline{a})}{\partial a_1}, \dots, \frac{\partial g(\underline{a})}{\partial a_m})^T$ et $\tilde{\mathbf{z}} = (\tilde{z}_1, \dots, \tilde{z}_N)^T$ avec $\tilde{z}_k = \frac{\partial f(\underline{b})}{\partial b_k} \Big|_{b=1}$. Il découle de (1.1) qu'un estimateur de la variance de $\hat{\mathbf{q}}$ est donné approximativement par l'estimateur de la variance du total estimé $\sum d_i(s) \tilde{z}_i = \hat{Y}(\tilde{\mathbf{z}})$; c'est-à-dire, $\text{var}(\hat{\mathbf{q}}) \approx v(\tilde{\mathbf{z}})$, où $v(y)$ représente l'estimateur de la variance de $\hat{Y} = \hat{Y}(y)$ en notation opérationnelle. Maintenant, nous remplaçons \tilde{z}_k par $z_k = \frac{\partial f(\underline{b})}{\partial b_k} \Big|_{b=d(s)}$, puisque les \tilde{z}_k sont inconnus, pour obtenir un estimateur par linéarisation de la variance

$$v_L(\hat{\mathbf{q}}) = v(z). \quad (1.2)$$

Notons que $v_L(\hat{\mathbf{q}})$ donné par (1.2) s'obtient simplement à partir de la formule de $v(y)$ pour $\hat{Y} = \hat{Y}(y)$ en remplaçant y_i par z_i pour $i \in s$. Soulignons que nous ne commençons pas par évaluer les dérivées partielles $\frac{\partial f(\underline{b})}{\partial b_k} \Big|_{b=1}$ afin d'obtenir $\tilde{\mathbf{z}}$ pour ensuite substituer les estimations aux composantes inconnues de $\tilde{\mathbf{z}}$. Par conséquent, notre méthode est similaire dans l'esprit à l'approche de Binder [1]. L'estimateur de la variance $v_L(\hat{\mathbf{q}})$ est valide, car z_i est un estimateur convergent de \tilde{z}_i .

Supposons que $\hat{\mathbf{q}}$ soit l'estimateur par quotient $\hat{Y}_R = X[(\sum d_i(s) y_i) / \sum d_i(s) x_i]$, où \sum représente la sommation sur $i \in U$. Alors $f(\underline{b}) = X[(\sum b_i y_i) / (\sum b_i x_i)] = X[\hat{Y}(\underline{b}) / \hat{X}(\underline{b})]$ et

$$z_k = \frac{\partial f(\underline{b})}{\partial b_k} \Big|_{b=d(s)} = \frac{X}{\hat{X}} (y_k - \hat{R}x_k).$$

Dans le cas de l'échantillonnage aléatoire simple, $v_L(\hat{Y}_R) = v(z)$ coïncide avec $v_{JL} = (\bar{X} / \bar{x})^2 v_L$.

Demnati et Rao [2] ont appliqué la méthode à divers problèmes, qui englobent les estimateurs par régression et par calage d'un total Y et d'autres estimateurs définis explicitement ou implicitement comme étant des solutions d'équations d'estimation. Ils ont obtenu un nouvel estimateur de la variance pour une classe générale d'estimateurs

par calage qui inclut les estimateurs par la méthode itérative du quotient (raking ratio) généralisée et les estimateurs par régression généralisée. Ils ont étendu également la méthode à l'échantillonnage à deux degrés et ont obtenu un estimateur de la variance qui utilise mieux les données du premier degré d'échantillonnage que les estimateurs de la variance par linéarisation classique.

2. Non-réponse partielle

En s'inspirant de Fay [3], Shao et Steel [7] ont proposé une méthode de calcul des estimateurs de la variance de l'estimateur du total de type Horvitz-Thompson, \hat{Y}^\bullet , lorsque des valeurs sont imputées pour tenir compte de la non-réponse partielle. Ils ont supposé que l'on peut exprimer le total estimé \hat{Y}^\bullet sous forme d'une fonction lisse des totaux $\hat{Y}^\bullet = \mathbf{y}(\hat{T}_o)$, où $\hat{T}_o = \sum d_i(s) \text{diag}(o_i) \mathbf{t}_i$, t_{ki} est la valeur de y_i ou celle d'une autre variable utilisée pour imputer une valeur à y_i , et $o_i = (o_{i1}, \dots, o_{ip})^T$ est le vecteur des variables indicatrices de réponse. Par exemple, considérons l'imputation par quotient lorsqu'on dispose des valeurs de la variable auxiliaire x_i pour tous les i de s . Une valeur de y_i manquante est alors imputée par $\hat{y}_i = \hat{R}_o x_i$, où $\hat{R}_o = (\sum d_i(s) o_i y_i) / (\sum d_i(s) o_i x_i)$ et o_i est l'indicateur de réponse pour y_i , c.-à-d. $o_i = 1$ si la valeur de y_i est observée et $o_i = 0$ si elle est manquante. L'estimateur imputé \hat{Y}^\bullet est donné par

$$\begin{aligned} \hat{Y}^\bullet &= \sum d_i(s) o_i y_i + \sum d_i(s) (1 - o_i) \hat{R}_o x_i \\ &= \sum d_i(s) o_{1i} y_i (1 + \sum d_i(s) o_{2i} x_i / \sum d_i(s) o_{1i} x_i) \end{aligned} \quad (2.1)$$

où $o_{i1} = o_i$ et $o_{i2} = 1 - o_i$. Il découle de (2.1) que \hat{Y}^\bullet est de la forme $\mathbf{y}(\hat{T}_o)$ avec $o_i = (o_{i1}, o_{i1}, o_{i2})^T$ et $\mathbf{t}_{-i} = (y_i, x_i, x_i)^T$

Nous supposons que l'imputation est déterministe. Nous avons $\text{Var}(\hat{Y}^\bullet - Y) = V_1 + V_2$, où $V_1 = E_o(\text{Var}_s(\hat{Y}^\bullet - Y))$, $V_2 = \text{Var}_o(E_s(\hat{Y}^\bullet - Y))$, E_o et Var_o représentent l'espérance et la variance sous le mécanisme de réponse, et E_s et Var_s représentent l'espérance et la variance sous le mécanisme d'échantillonnage selon un plan d'échantillonnage donné. Shao et Steel [7] ont obtenu un estimateur de la variance v_1 de V_1 en se servant d'un estimateur par linéarisation de la variance de $\mathbf{y}(\hat{T}_o)$ étant donné les o_i .

Ils ont obtenu aussi un estimateur, v_2 de $V_2 = \left[\nabla \mathbf{j}(E_o \mathbf{T}_o) \right]^T \mathbf{C} \left[\nabla \mathbf{j}(E_o \mathbf{T}_o) \right]$, où $\mathbf{j}(\mathbf{T}_o) = \mathbf{y}(E_s \hat{T}_o) - Y$, en calculant \mathbf{C} dont le kl^e élément $c_{kl} = \text{cov}_o(\sum o_{ki} t_{ki}, \sum o_{li} t_{li})$ et en remplaçant les quantités inconnues par des estimateurs. Pour l'échantillonnage aléatoire simple et l'imputation par quotient, Shao et Steel [7] ont obtenu v_1 sous la forme

$$v_1 = N^2 \frac{(1 - n/N)}{n(n-1)} \left[\left(\frac{\bar{x}}{\bar{x}_o} \right)^2 \frac{s_d^2}{n_o} + 2 \frac{\bar{x}}{\bar{x}_o} \frac{\hat{R}_o s_{dx}}{n} + \frac{\hat{R}_o^2 s_x^2}{n} \right], \quad (2.2)$$

où \bar{x} et s_x^2 sont la moyenne et la variance d'échantillon des x_i , $\bar{x}_o = \sum_{i \in s} o_i x_i / n_o$ est la moyenne des x_i pour les répondants, n_o est le nombre de répondants, $s_d^2 = \sum_{i \in s} (y_i - \hat{R}_o x_i)^2 / (n_o - 1)$, et $s_{dx} = \sum_{i \in s} o_i x_i (y_i - \hat{R}_o x_i)^2 / (n_o - 1)$.

Aussi, sous l'hypothèse de réponse uniforme (c.-à-d. si les o_i sont indépendamment et identiquement distribués en ayant une moyenne p_y et une variance $p_y(1 - p_y)$), Shao et Steel [7] ont obtenu v_2 sous la forme

$$v_2 = (X / X_o)^2 \hat{p}_y (1 - \hat{p}_y) N s_d^2, \quad (2.3)$$

où $\hat{p}_y = \sum o_i d_i(s) / \sum d_i(s)$. La somme de (2.2) et (2.3) donne l'estimateur de la variance de \hat{Y}^\bullet .

La méthode de Shao et Steel [7] est fondée sur la méthode de linéarisation classique qui consiste à i) exprimer l'estimateur en fonction de composantes élémentaires, ii) calculer les dérivées partielles au niveau de la population et iii) estimer les paramètres inconnus dans la formule. Par conséquent, l'estimateur correspondant de la variance n'est pas nécessairement unique. Notre méthode permet d'éviter d'exprimer l'estimateur en fonction de composantes élémentaires, donc produit directement un estimateur unique de la variance ayant les propriétés souhaitées. Nous présentons notre méthode pour l'imputation par quotient à la sous-section 3.1, et nous examinons le cas de l'estimation de la variance sous ajustement de la pondération pour tenir compte de la non-réponse complète et imputation par quotient pour tenir compte de la non-réponse partielle à la sous-section 3.2.

3. Nouvelle méthode : réponses manquantes

Après ajustement de la pondération pour tenir compte de la non-réponse complète et imputation pour tenir compte de la non-réponse partielle, nous estimons le total de la population Y au moyen d'un total pondéré à partir de l'échantillon

$$\hat{Y}^\bullet = \sum \tilde{w}_i(s) o_i y_i + \sum \tilde{w}_i(s) (1 - o_i) \hat{y}_i^\bullet, \quad (3.1)$$

où $\tilde{w}_i(s)$ est le poids ajusté et \hat{y}_i^\bullet représente la valeur imputée pour l'unité i . Nous pouvons réécrire l'estimateur (3.1) sous la forme

$$\hat{Y}^\bullet = \sum \tilde{w}_i(s) \hat{y}_i = \hat{y}^T \tilde{w}(s),$$

(3.2)

où $\hat{y} = (\hat{y}_1, \dots, \hat{y}_N)^T$ et $\hat{y}_i = o_i y_i + (1 - o_i) \hat{y}_i^\bullet$. À la sous-section 3.1, nous étudions le cas de la non-réponse partielle uniquement (c.-à-d. $\tilde{w}_i(s) = d_i(s)$) en supposant que nous pouvons exprimer \hat{Y}^\bullet sous la forme d'une fonction lisse de totaux $\sum d_i(s) \text{diag}(o_i) t_i$, où t_{ki} est la valeur de y_i ou celle d'une autre variable utilisée pour imputer une valeur à y_i . À la sous-section 3.2, nous traitons le cas plus général de l'ajustement de la pondération pour tenir compte de la non-réponse complète et de l'imputation pour tenir compte de la non-réponse partielle.

3.1 Imputation pour tenir compte de la non-réponse partielle

Nous supposons que l'estimateur avec données imputées \hat{Y}^\bullet est une fonction lisse de totaux, $\sum d_i(s) \text{diag}(o_i) t_i$, à l'instar de Shao et Steel [7]. Dans ce cas, nous pouvons exprimer \hat{Y}^\bullet sous la forme $f(\underline{A}_w, \underline{A}_y)$, où $\underline{A}_w = \text{diag}(d(s)) \underline{A}_o$ est une matrice $m \times N$ dont la j^e colonne est $\underline{w}_j(s) = (w_{1j}(s), \dots, w_{mj}(s))^T$, $j = 1, \dots, N$. Le vecteur $\underline{w}_j(s)$ est défini comme étant

$$\begin{aligned}\underline{w}_j(s) &= (w_{1j}(s), \dots, w_{mj}(s))^T \\ &= (o_{1j}d_j(s), \dots, o_{mj}d_j(s))^T = o_j d_j(s),\end{aligned}$$

où $o_j = (o_{1j}, \dots, o_{mj})^T$ est le vecteur des variables indicatrices correspondant au vecteur $y_j = (y_{1j}, \dots, y_{mj})^T$.

Par souci de simplicité, nous laissons tomber A_y et écrivons $\hat{Y}^\bullet = f(A_w)$. En cas d'imputation par quotient, nous avons $m = 2$, $y_{1j} = x_j$, $y_{2j} = y_j$, $o_{1j} = 1$, $o_{2j} = o_j$, $\underline{w}_j(s) = (w_{1j}(s), w_{2j}(s))^T = (d_j(s), o_j d_j(s))^T$ et

$$\hat{Y}^\bullet = \sum w_{2i}(s)(y_i - \hat{R}_o x_i) + \sum w_{1i}(s)\hat{R}_o x_i,$$

où $\hat{R}_o = (\sum w_{2i}(s)y_i) / (\sum w_{2i}(s)x_i)$.

Comme l'estimateur $\hat{Y}^\bullet = f(A_w)$ est une fonction de totaux, nous pouvons appliquer la méthode de Demnati et Rao [2] pour approximer sa variance au moyen de la variance d'une fonction linéaire

$$Var(\hat{Y}^\bullet) \approx Var(\hat{Y}_L^\bullet)$$

avec

$$\hat{Y}_L^\bullet = \sum (o_i d_i(s))^T \tilde{z}_i = \sum w_i^T(s) \tilde{z}_i,$$

où \tilde{z}_i est le vecteur des dérivées de $f(A_b)$ par rapport à b_k évaluées à $A_b = E(A_w)$, où A_b est une matrice $m \times N$ de nombres réels arbitraires, $f(A_b)$ est obtenue par remplacement de A_w par A_b dans la formule de \hat{Y}^\bullet et b_k est un vecteur colonne de A_b . Nous pouvons estimer la variance totale de \hat{Y}_L^\bullet selon

$$v(\hat{Y}_L^\bullet) = v_s(o^T \underline{z}) + v_o(\underline{z}), \quad (3.4)$$

où \underline{z}_k est le vecteur des dérivées de l'estimateur $f(A_b)$ par rapport à b_k évaluées à $A_b = A_w$, et $v_o(\underline{z})$ est un estimateur de $Var(\sum diag(o_i) \underline{z}_i)$. Sous un mécanisme de réponse indépendante,

$$v_o(\underline{z}) = \sum \underline{z}_i^T cov_o(o_i) \underline{z}_i, \quad (3.5)$$

où $cov_o(o_i)$ est un estimateur (approximativement) non biaisé de $E(o_i o_i^T) - E(o_i)E(o_i^T)$.

Sous imputation par quotient, nous avons

$$\begin{aligned}\underline{z}_k &= \frac{\partial}{\partial b_k} \left(\sum b_{2i}(y_i - \hat{R}_o(A_b)x_i) + \sum b_{1i}\hat{R}_o(A_b)x_i \right) \Big|_{A_b=A_w} \\ &= \left(\hat{R}_o x_k, (\hat{X} / \hat{X}_o)(y_k - \hat{R}_o x_k) \right)^T\end{aligned} \quad (3.6)$$

Il découle de (3.6) que

$$z_k = o_k (\hat{X} / \hat{X}_o)(y_k - \hat{R}_o x_k) + \hat{R}_o x_k. \quad (3.7)$$

Par conséquent, $v_s(o^T \underline{z})$ est égal à $v_1 = v(z)$. Dans les conditions d'échantillonnage aléatoire simple, $v(z)$ avec z_k donné par (3.7) concorde avec l'expression (2.2) de Shao et Steel [7]. Aussi,

$$\text{cov}_o(o_i) = d_i(s) \begin{pmatrix} 0 & 0 \\ 0 & o_i(1 - \hat{\mathbf{x}}_{io}) \end{pmatrix},$$

où $\hat{\mathbf{x}}_{io}$ est un estimateur de la probabilité de réponse de l'unité i . Par conséquent, $v_o(\underline{z})$, donné par (3.5), se réduit à

$$v_o(\underline{z}) = (\hat{X} / \hat{X}_o)^2 \sum d_i(s) o_i (1 - \hat{\mathbf{x}}_{io}) (y_i - \hat{R}_o x_i)^2. \quad (3.8)$$

Dans le cas d'un échantillonnage aléatoire simple et d'un mécanisme de réponse uniforme, (3.8) se réduit à

$$v_o(\underline{z}) = \frac{N}{n} (\hat{X} / \hat{X}_o)^2 (1 - n_o / n) \sum o_i (y_i - \hat{R}_o x_i)^2 \quad (3.9)$$

qui est l'estimateur v_2 de Shao et Steel [7] donné par (2.3).

4.2 Ajustement de la pondération et imputation pour tenir compte de la non-réponse partielle

Soit r_i , la variable indicatrice de réponse partielle pour la i^e unité, c.-à-d. $r_i = 0$ si la non-réponse est complète et $r_i = 1$ si la réponse est partielle. La variable indicatrice de réponse partielle r_i est reliée aux variables indicatrices de réponse à une question o_p , $p = 1..m$ par

$$r_i = 1 - \prod_{p=1}^m (1 - o_{ip}). \quad (3.10)$$

Nous avons

$$\text{Cov}(r_i, o_{ip}) = E(r_i o_{ip}) - E(r_i)E(o_{ip}),$$

pour toute variable indicatrice de réponse o_{ip} . Notant que $r_i o_{ip} = o_{ip}$ pour tout o_{ip} ,

$$r_i o_{ip} = [1 - (1 - o_{ip}) \prod_{q \neq p} (1 - o_{iq})] o_{ip} = o_{ip}.$$

Donc,

$$\text{Cov}(r_i, o_{ip}) = E(o_{ip}) - E(r_i)E(o_{ip}) = E(o_{ip})(1 - E(r_i)).$$

Nous pouvons donner un estimateur de $\text{Cov}(r_i, o_{ip})$ sous la forme

$$\text{cov}(r_i, o_{ip}) = o_{ip}(1 - \hat{\mathbf{x}}_{ir})$$

où $\hat{\mathbf{x}}_{ir} = \hat{E}(r_i)$ et $\hat{E}(\cdot)$ représente un estimateur de $E(\cdot)$.

Une méthode très répandue de correction de la non-réponse complète consiste à employer un nouvel ensemble de poids, $\tilde{w}_i(s)$, dont le i^e élément est égal à

$$\tilde{w}_i(s) = d_i(s) r_i g_i(\underline{d}(s), r_i, \underline{A}_c), \quad (3.11)$$

où les $g_i(\underline{d}(s), r_i, \underline{A}_c)$ sont connus sont le nom de poids g dans le contexte de l'estimateur par régression et \underline{A}_c

est une matrice de variables auxiliaires dont les valeurs sont connues pour toutes les unités figurant dans l'échantillon. L'estimateur par quotient est un cas particulier de (4.11) pour lequel les poids g se réduisent à

$$g_i(\underline{d}(s), r_i, \underline{A}_c) = \frac{\sum d_i(s) \mathbf{c}_i}{\sum d_i(s) r_i \mathbf{c}_i} = \frac{\hat{\mathbf{c}}}{\hat{\mathbf{c}}_r}, \quad (3.12)$$

où $\hat{\mathbf{c}} = \sum d_i(s) \mathbf{c}_i$ et $\hat{\mathbf{c}}_r = \sum d_i(s) r_i \mathbf{c}_i$. L'ajustement des poids au moyen du quotient (3.12) est un cas particulier de la classe de poids de calage obtenus en utilisant l'estimateur par régression. On utilise également les poids obtenus par ajustement itératif du quotient généralisé pour compenser pour la non-réponse complète. Un autre moyen de tenir compte de la non-réponse complète consiste à pondérer chaque observation par la probabilité inverse de réponse, auquel cas

$$g_i(\underline{d}(s), r_i, \underline{A}_c) = \hat{\mathbf{x}}_{ir}^{-1},$$

où

$$\hat{\mathbf{x}}_{ir} = \mathbf{x}_{ir}(r_i, \underline{d}(s)) = \Pr(r_i = 1 | \underline{d}(s), \underline{A}_c),$$

est l'estimateur de la probabilité de réponse défini comme étant la solution d'une équation d'estimation de la forme

$$\hat{U}(\hat{\mathbf{x}}_{ir}) = \sum d_i(s) u_i(r_i, \mathbf{c}_i, \hat{\mathbf{x}}_{ir}) = 0.$$

Dans le cas logistique, nous avons

$$\hat{U}(\hat{\mathbf{x}}_{ir}) = \sum d_i(s) (r_i - \hat{\mathbf{x}}_{ir}) \mathbf{c}_i = 0,$$

où

$$u_i(r_i, \mathbf{c}_i, \hat{\mathbf{x}}_{ir}) = (r_i - \hat{\mathbf{x}}_{ir}) \mathbf{c}_i,$$

$$\mathbf{x}_{ir} = \exp(\mathbf{c}_i^T \underline{\mathbf{b}}) / (1 + \exp(\mathbf{c}_i^T \underline{\mathbf{b}})) = \Pr(r_i = 1 | \mathbf{c}_i, \underline{\mathbf{b}}),$$

et $\underline{\mathbf{b}}$ est le vecteur des variables explicatives.

Avec les méthodes d'ajustement de la pondération susmentionnées, nous pouvons obtenir la variance en suivant la méthode de Demnati et Rao [2] en exprimant \hat{Y}^* sous la forme $f(\underline{A}_w)$, puis en dérivant $f(\underline{A}_b)$ par rapport à b_k .

Nous omettons les détails par souci de simplicité, mais nous illustrons le calcul pour l'estimateur (3.1) sous ajustement de la pondération par quotient (3.12) et imputation par quotient, c'est-à-dire

$$\hat{Y}^* = \sum_i o_i \tilde{w}_i(s) y_i + \sum_j (1 - o_j) \tilde{w}_j(s) \hat{R}_o x_j,$$

avec

$$\hat{R}_o = \frac{\sum \tilde{w}_i(s) o_i y_i}{\sum \tilde{w}_i(s) o_i x_i},$$

et

$$\tilde{w}_i(s) = d_i(s) r_i \hat{\mathbf{c}} / \hat{\mathbf{c}}_r.$$

Nous avons

$$w_i(s) = (\mathbf{1}, o_i, r_i)^T d_i(s),$$

$$\mathbf{z}_k = \left(x_k \hat{R}_o (\hat{\mathbf{X}}_r / \hat{\mathbf{c}}_r), (\hat{\mathbf{c}} / \hat{\mathbf{c}}_r) (\hat{\mathbf{X}}_r / \hat{\mathbf{X}}_o) (y_k - \hat{R}_o x_k), (\hat{\mathbf{c}} / \hat{\mathbf{c}}_r) \hat{R}_o (x_k - (\hat{\mathbf{X}}_r / \hat{\mathbf{c}}_r) \mathbf{c}_k) \right)^T,$$

et

$$\text{cov}_o(\mathbf{z}_i^T) = d_i(s) \begin{pmatrix} 0 & 0 & 0 \\ 0 & o_i(1 - \hat{\mathbf{x}}_{io}) & o_i(1 - \hat{\mathbf{x}}_{ir}) \\ 0 & o_i(1 - \hat{\mathbf{x}}_{ir}) & r_i(1 - \hat{\mathbf{x}}_{ir}) \end{pmatrix}.$$

Conclusion

Nous avons présenté une nouvelle méthode d'estimation de la variance en cas de réponses manquantes. Nous donnons un estimateur valide de la variance pour diverses méthodes d'ajustement de la pondération utilisées fréquemment pour tenir compte de la non-réponse complète, ainsi que pour l'imputation basée sur des fonctions lisses des valeurs observées, en particulier l'imputation par quotient, qui est utilisée fréquemment pour tenir compte de la non-réponse partielle. L'extension à l'imputation par le plus proche voisin et aux enquêtes par panel est à l'étude.

Bibliographie

- [1] Binder, D., « Linearization methods for single phase and two-phase samples: a cookbook approach », *Survey Methodology*, 22, pp 17-22, 1996.
- [2] Demnati, A. and Rao, J. N. K., Linearization variance estimators for survey data, Methodology Branch Working Paper, SSMD-2001-010E. Statistics Canada, 2001.
- [3] Fay, R. E., « A design-based perspective on missing data variance », in *Proceeding of the 1991 Annual Research Conference, US Bureau of the census*, pp 429-440, 1991.
- [4] Rao, J. N. K., « On variance estimation with imputed survey data (with discussion) », *Journal of the American Statistical Association*, 91, pp 499-520, 1996.
- [5] Royall, R. M., and Cumberland, W. G., « An empirical study of the ratio estimator and estimators of its variance », *Journal of the American Statistical Association*, pp 76, 66-77, 1981.
- [6] Särndal, C.-E., Swensson, B., and Wretman, J.H., «The Weighted residual technique for estimating the variance of the general regression estimator of the finite population total », *Biometrika*, 76, pp 527-537, 1989.
- [7] Shao, J. and Steel, P., « Variance estimation for survey data with composite imputation and nonnegligible sampling fractions », *Journal of the American Statistical Association*, 94, pp 254-265, 1999.
- [8] Valliant, R., « Postsratification and conditional variance estimation », *Journal of the American Statistical Association*, 88, pp 89-96, 1993.
- [9] Yung, W. and Rao, J. N. K., « Jackknife linearization variance estimators under stratified multi-Stage sampling », *Survey Methodology*, 22, pp 23-31, 1996.