

Risque d'amplification du biais de l'estimateur par calage généralisé en présence de non-réponse

Éric Lesage

Insee

Mars-Avril 2015

Journées de méthodologie statistique
Insee, Paris

Travail réalisé en collaboration avec David Haziza (Département de Mathématique de l'Université de Montréal) et Xavier D'Haultfoeuille (Crest-Ensa).

Sommaire

- 1 Contexte et rappels
- 2 Le modèle de non-réponse
- 3 Estimation des probabilités de réponse p_i
- 4 Propriétés de l'estimateur par calage généralisé
- 5 Études par simulation

Le traitement de la non-réponse totale par calage est un domaine de recherche déjà ancien mais qui reste très actif :

- Dupont (1996), Deville (1998), Le Guennec et Sautory (2004), Deville (2004) et Sautory (2003).
- Lundström et Särndal (1999), Särndal et Lundström (2005, 2010, 2011).
- Kott (2006), Chang et Kott (2008), Kott et Chang (2008), Kott et Liao (2012).
- Osier (2012), Lesage (2012), Lesage, Haziza et D'haultfoeuille (2015), et Haziza et Lesage (2015)

Introduction II

- Nécessité de modéliser le mécanisme de non-réponse pour choisir la méthode de calage appropriée
 - choix des variables instrumentales, des variables de calage et de la fonction de calage.
- Risques d'amplification du biais et de la variance de l'estimateur par **calage généralisé**.
 - Importance du choix des variables de calage x

Contexte I

- U : population finie de taille N
- y_1, \dots, y_J : J variables d'intérêt collectées par l'enquête
 - $\mathbf{y}_i^\top = (y_{i,1}, \dots, y_{i,J})$ vecteur des caractéristiques associées à l'unité $i \in U$
- Objectif : estimer les totaux dans la population finie

$$t_{yj} = \sum_{i \in U} y_{i,j} \quad j = 1, \dots, J$$

- x_1, \dots, x_P : P variables auxiliaires
 - $\mathbf{x}_i^\top = (x_{i,1}, \dots, x_{i,P})$ vecteur associé à l'unité $i \in U$

Contexte II

- Échantillon s sélectionné selon un **plan de sondage probabiliste** $p(s)$
 - \mathcal{S} : ensemble des échantillons possibles
 - I_i : variable indicatrice d'échantillonnage
 - π_i : probabilité d'inclusion de l'unité i :

$$\pi_i = \mathbb{E}(I_i \mid \mathbf{Y}_i = \mathbf{y}_i, \mathbf{X}_i = \mathbf{x}_i) = \mathbb{E}(I_i \mid \mathbf{X}_i = \mathbf{x}_i)$$

- $d_i = \pi_i^{-1}$: poids d'échantillonnage
- Mécanisme de non-réponse
 - $s_r \subset s$: sous ensemble des répondants
 - R_i : variable indicatrice de réponse

Définition : Estimateur par calage (Deville et Särndal, 1992)

$$\hat{t}_C = \sum_{i \in s} w_i y_i$$

▶ Poids de calage

$$w_i = d_i \times F(\hat{\lambda}^\top \mathbf{x}_i)$$

▶ Fonction de calage

- $F(\cdot)$ une fonction dérivable et monotone,
- $F(0) = 1$ et $F'(0) = 1$.

▶ Équations de calage

$\hat{\lambda}$ est un vecteur de paramètres solution de :

$$\sum_{i \in s} d_i F(\hat{\lambda}^\top \mathbf{x}_i) \mathbf{x}_i = \mathbf{t}_x.$$

• Exemple de fonctions de calage :

- méthode linéaire : $F(\hat{\lambda}^\top \mathbf{x}_i) = 1 + \hat{\lambda}^\top \mathbf{x}_i$,
- méthode exponentielle : $F(\hat{\lambda}^\top \mathbf{x}_i) = \exp(\hat{\lambda}^\top \mathbf{x}_i)$.

Définition : Estimateur par calage généralisé (Deville, 1998)

$$\hat{t}_C = \sum_{i \in s} w_i y_i,$$

- ▶ Poids de calage : $w_i = d_i \times F(\hat{\lambda}^\top \mathbf{z}_i)$,
- ▶ \mathbf{z}_i vecteur des **instruments de calage**, connu pour $i \in s$,
- ▶ $F(\cdot)$ fonction de calage,
- ▶ Équations de calage :

$$\sum_{i \in s} d_i F(\hat{\lambda}^\top \mathbf{z}_i) \mathbf{x}_i = \mathbf{t}_x,$$

\mathbf{x}_i est le vecteur des **variables de calage**

\mathbf{t}_x est le vecteur des marges de calage.

Remarque importante : il n'est pas nécessaire de connaître \mathbf{t}_z .

Les variables explicatives de la non-réponse et de la variable d'intérêt y

- On s'intéresse à **une** variable d'intérêt y_{j0} qu'on note y pour simplifier
- \mathbf{z} : vecteur de Q **variables explicatives de la non-réponse** qui sont également liées à la variable d'intérêt y
 - $\mathbf{z}_i^T = (z_{i,1}, \dots, z_{i,Q})$ vecteur associé à l'unité $i \in U$
 - Nature des variables dans \mathbf{z} : des x_p^o , des x_p^* et des y_j !
 - Il suffit que \mathbf{z}_i soit **connu pour les unités $i \in s_r$** et on n'est pas obligé de connaître les totaux sur la population $t_{z_i,q}$!

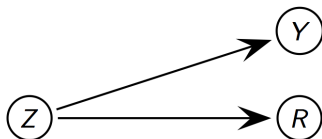


Figure: Relation entre les variables Y , Z and R

Les variables explicatives de la non-réponse et de la variable d'intérêt y

- La variable y n'a pas d'effet propre dans le modèle de non-réponse après avoir tenu compte des variables \mathbf{z} :

$$Y_i \perp R_i \mid \mathbf{Z}_i \quad (1)$$

- On peut dire que y est exclue des variables explicatives de la non-réponse et la relation (1) est appelée **relation d'exclusion** de y pour le modèle de non-réponse
- Modèle de non réponse considéré :

$$\mathbb{E}(R_i \mid \mathbf{Z}_i, Y_i, I_i = 1) = \mathbb{E}(R_i \mid \mathbf{Z}_i, I_i = 1) = h(\boldsymbol{\gamma} ; \mathbf{Z}_i) > 0,$$

où $\boldsymbol{\gamma}$ est un vecteur de même dimension que \mathbf{z} et $h(\cdot)$ une fonction suffisamment régulière.

- $p_i = h(\boldsymbol{\gamma} ; \mathbf{Z}_i)$: probabilité de répondre de l'unité $i \in s_r$

Estimation des probabilités de réponse p_i

- Rappel : dans l'**approche de repondération en deux étapes** pour traiter la non-réponse
 - Le mécanisme de non-réponse est traité comme une ultime phase de sondage
 - $\hat{p}_i = h(\hat{\gamma} ; \mathbf{Z}_i)$: estimateur de p_i , où $\hat{\gamma}$ est un estimateur convergent de γ
 - "*Propensity score adjusted estimator*"

$$\hat{t}_{PSA} = \sum_{i \in s_r} \pi_i^{-1} \hat{p}_i^{-1} y_i$$

- On va montrer
 - on peut estimer les probabilités de réponse p_i à partir d'un système d'équations estimantes qui sont similaires à des équations de calage généralisées
 - l'estimateur par calage généralisé est un estimateur de type PSA.

Estimation de γ

- A partir de notre modèle de non-réponse, on a :

$$\mathbb{E}(I_i R_i \mid \mathbf{Z}_i, Y_i) = \mathbb{E}(I_i \mid \mathbf{Z}_i; Y_i) h(\gamma; \mathbf{Z}_i)$$

- Sans perte de généralité on peut écrire

$$\mathbb{E}(I_i R_i \mid \mathbf{Z}_i, Y_i) = \pi_i h(\gamma; \mathbf{Z}_i)$$

- Deux types d'équations de moments du modèle de non-réponse

- 1 $\mathbb{E}(\pi_i^{-1} h^{-1}(\gamma; \mathbf{z}_i) \mathbf{z}_i I_i R_i) = \mathbb{E}(\mathbf{z}_i)$

- 2 $\mathbb{E}(\{R_i - \pi_i h(\gamma; \mathbf{z}_i)\} \mathbf{z}_i I_i R_i) = \mathbf{0}$

- Deux types d'équations estimantes (contre partie empirique)

- 1 $\sum_{i \in s_r} \pi_i^{-1} h^{-1}(\hat{\gamma}; \mathbf{z}_i) \mathbf{z}_i = \mathbf{t}_z$

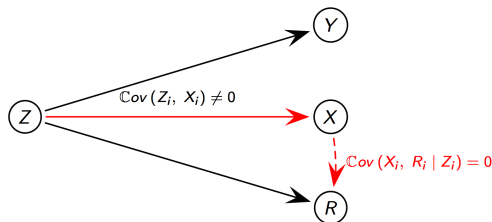
- 2 $\sum_{i \in s} \{R_i - \pi_i h(\hat{\gamma}; \mathbf{z}_i)\} \mathbf{z}_i = \mathbf{0}$

Estimation de γ : problème de covariables inobservées

- Si $z_{i,q}$ n'est pas une variable auxiliaire (x^0 ou x^*) on ne peut résoudre aucune des 2 équations estimantes précédentes.
- La solution est de trouver une variable auxiliaire, $x_{i,q}$, qui serve de proxy pour $z_{i,q}$.
- C'est une solution semblable à celle utilisée pour contourner le problème d'endogénéité en économétrie...
- ... toutefois, notre problème de départ n'est pas la présence d'une variable endogène mais d'une **variable explicative $z_{i,q}$ partiellement observée**.
- On verra que, comme en économétrie, l'utilisation d'une variable instrumentale (notre proxy) fait peser le risque d'avoir un estimateur dont le biais est amplifié par la faible corrélation entre la variable explicative $z_{i,q}$ et sa proxy $x_{i,q}$.

Variable proxy de z

- Considérons simplement que
 - $\mathbf{z}_i^T = (1, z_i)$
- Soit x la variable proxy de z
 - Variable auxiliaire x^* (ou x^0 , on utilise alors $\hat{t}_{x,\pi}$ à la place de t_x)
 - Corrélée à z
 - x n'intervient pas dans le modèle de non-réponse : $X \perp R \mid Z$



Estimation de γ : équations de moment avec variable instrumentale

- Modèle de non-réponse

$$\mathbb{E}(I_i R_i \mid \mathbf{Z}_i, \mathbf{X}_i) = \pi_i h(\gamma ; \mathbf{Z}_i),$$

- $\mathbf{x}_i^\top = (1, x_i)$ et $\mathbf{t}_x^\top = (N, t_x)$

- Nouvelles équations de moments

$$\mathbb{E}(\pi_i^{-1} h^{-1}(\gamma ; \mathbf{z}_i) I_i R_i \mathbf{x}_i) = \mathbb{E}(\mathbf{x}_i)$$

- Équations estimantes (contre partie empirique)

$$\sum_{i \in S_r} \pi_i^{-1} h^{-1}(\hat{\gamma} ; \mathbf{z}_i) \mathbf{x}_i = \mathbf{t}_x$$

Lien avec le calage généralisé

- Le système précédent correspond à un estimateur par calage généralisé où
 - $i \in s_r$, échantillon des répondants
 - x : vecteur des variables de calage
 - t_x : vecteur des marges de calage
 - z : vecteur des **instruments de calage**
 - Le vecteur des paramètres γ correspond au vecteur λ
 - $F(\lambda^\top Z_i) = h^{-1}(\lambda; Z_i)$ où $F(\cdot)$ est une fonction de calage (monotone et continûment dérivable)
- La recherche de l'estimateur $\hat{t}_{PSA} = \sum_{i \in s_r} \pi_i^{-1} \hat{p}_i^{-1} y_i$ nous amène dans ce cas à l'estimateur par calage généralisé $\hat{t}_C = \sum_{i \in s_r} \pi_i^{-1} F(\hat{\lambda}^\top z_i) y_i$
- Voyons maintenant les propriétés de l'estimateur \hat{t}_C ...

Convergence de l'estimateur par calage généralisé

$$\hat{t}_C = \sum_{i \in S_r} \pi_i^{-1} F(\hat{\lambda}^\top \mathbf{z}_i) y_i$$

- Sous les conditions
 - ✓ Modèle de non réponse : $\mathbb{E}(R_i | \mathbf{Z}_i, I_i = 1) = F^{-1}(\gamma^\top \mathbf{Z}_i) > 0$,
 - ✓ Relation d'exclusion sur y : $Y \perp R | \mathbf{Z}$
 - ✓ Relation d'exclusion sur x : $X \perp R | \mathbf{Z}$
 - et quelques conditions techniques
- On a :
 - 1 $\hat{\lambda} \xrightarrow{P} \gamma$
 - 2 $F(\hat{\lambda}^\top \mathbf{z}_i) p_i - 1 \xrightarrow{P} 0$
 - 3 $(\hat{t}_C - t_y) / N \xrightarrow{P} 0$.
- Remarque : ces résultats sont vrais quelque soit le modèle qui lie y et \mathbf{Z} !
- Variance amplifiée par la faible corrélation entre z et x (Osier, 2012)

Non convergence de l'estimateur par calage généralisé

- Sous les conditions

- ✓ Modèle de non réponse : $\mathbb{E}(R_i | \mathbf{Z}_i, I_i = 1) = F^{-1}(\boldsymbol{\gamma}^\top \mathbf{Z}_i) > 0$,

- ✓ Relation d'exclusion sur y : $Y \perp R | \mathbf{Z}$

- ✗ Relation d'exclusion sur x non vérifiée : $X \not\perp R | \mathbf{Z}$
en l'occurrence :

$$\mathbb{E}(R_i | \mathbf{Z}_i, \mathbf{X}_i, I_i = 1) \neq \mathbb{E}(R_i | \mathbf{Z}_i, I_i = 1)$$

- On obtient :

- $F(\hat{\boldsymbol{\lambda}}^\top \mathbf{z}_i)p_i - 1$ ne converge pas vers 0
 - L'estimateur \hat{t}_C n'est pas convergent
 - Le biais est amplifié par la faible corrélation entre z et x

Seconde expression du biais

Proposition :

Le biais approché de \hat{t}_C peut s'écrire :

$$\begin{aligned} \text{Biais}(\hat{t}_C) &\approx - \sum_{k \in \mathcal{U}} (1 - p_i F_i) (y_i - \mathbf{z}_i^\top \boldsymbol{\beta}_{pf}) \\ &+ \frac{\beta_{pf,1}}{\alpha_{pf,1}} \times \sum_{k \in \mathcal{U}} (1 - p_i F_i) (x_i - \mathbf{z}_i^\top \boldsymbol{\alpha}_{pf}), \end{aligned}$$

- ▶ $(y_i - \mathbf{z}_i^\top \boldsymbol{\beta}_{pf})$ résidu de la régression de y sur \mathbf{z}
- ▶ $(x_i - \mathbf{z}_i^\top \boldsymbol{\alpha}_{pf})$ résidu de la régression de x sur \mathbf{z}
- ▶ $\alpha_{pf,1}$ corrélation entre x et \mathbf{z}

Amplification de la variance de l'estimateur \hat{t}_C

Proposition :

La variance liée à la non-réponse approchée de \hat{t}_C est :

$$A\mathbb{V}ar_q(\hat{t}_C) = \sum_{k \in \mathcal{U}} p_i(1 - p_i) F_i^2 \times \left\{ (y_i - \mathbf{z}_i^\top \boldsymbol{\beta}_{pf}) - \frac{\beta_{pf,1}}{\alpha_{pf,1}} (x_i - \mathbf{z}_i^\top \boldsymbol{\alpha}_{pf}) \right\}^2.$$

- Variance due à la non-réponse est faible si :
 - les résidus des modèles de régression $(y_i - \mathbf{z}_i^\top \boldsymbol{\beta}_{pf})$ et $(x_i - \mathbf{z}_i^\top \boldsymbol{\alpha}_{pf})$ sont petits,
 - $\alpha_{pf,1}$ est fort.
- Variance due à la non-réponse instable si : $\alpha_{pf,1}$ faible (i.e., une faible corrélation entre x et z).

Cas particulier des modèles de régression linéaire entre y et le vecteur \mathbf{z} et entre x et \mathbf{z}

- Convergence de l'estimateur par calage généralisé si :

~~✗~~ Modèle de non réponse : $\mathbb{E}(R_i | \mathbf{Z}_i, I_i = 1) = F^{-1}(\gamma^\top \mathbf{Z}_i) > 0,$

✓ Relation d'exclusion sur y : $Y \perp R | \mathbf{Z}$

✓ Relation d'exclusion sur x : $X \perp R | \mathbf{Z}$

✓ Modèle de régression linéaire entre y et \mathbf{z}

$$\mathbb{E}(Y_i | Z_i) = \beta^\top Z_i$$

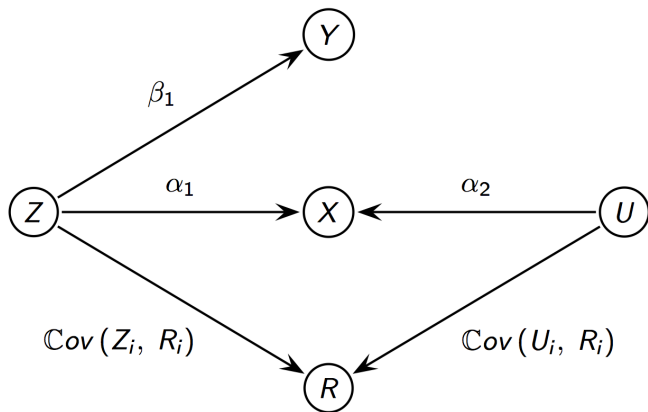
✓ Modèle de régression linéaire entre x et \mathbf{z}

$$\mathbb{E}(Y_i | Z_i) = \alpha^\top Z_i$$

Étude par simulation

On génère une population finie de taille $N = 1\,000$

- instrument de calage z
- variable inobservée u .
- z et u : générées selon une loi uniforme sur $[-\sqrt{3}, \sqrt{3}]$.
- variable d'intérêt y ,
 - $y_i = 10 + 2z_i + \varepsilon_i^y$,
 - où $\varepsilon_i^y \sim \mathcal{N}(0, 1)$, $R^2 \approx 80\%$.
- Jeu de variables proxy $X^{(\alpha_1, \alpha_2)}$
 - $X_i^{(\alpha_1, \alpha_2)} = \alpha_1 z_i + \alpha_2 u_i + \sigma_{(\alpha_1, \alpha_2)} \varepsilon_i^{(\alpha_1, \alpha_2)}$,
 - où $\varepsilon_i^{(\alpha_1, \alpha_2)} \sim \mathcal{N}(0, 1)$,
 - $\alpha_1 \in \{0.2, 0.3, 0.5, 0.7\}$,
 - $\alpha_2 \in \{0, 0.1, 0.3, 0.5\}$.



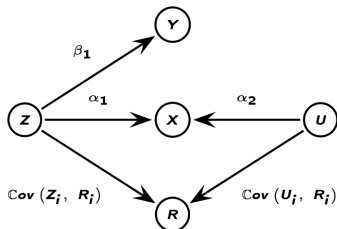
Étude par simulation

- On considère un recensement $n = N = 1\,000$.
- Probabilité de répondre de l'unité i :
 - $\text{logit}(p_i) = 1.5z_i + u_i$
 - Taux de réponse moyen : 50%
- R_i pour $i \in U$ sont générées indépendamment selon des lois de Bernoulli de paramètre p_i .
- Ce processus est répété $K = 10\,000$ fois, menant à $K = 10\,000$ ensembles de répondants.
- On calcule l'estimateur par calage $\hat{t}_C(\alpha_1, \alpha_2)$ avec la méthode linéaire, pour les différentes valeurs de α_1 et α_2 .

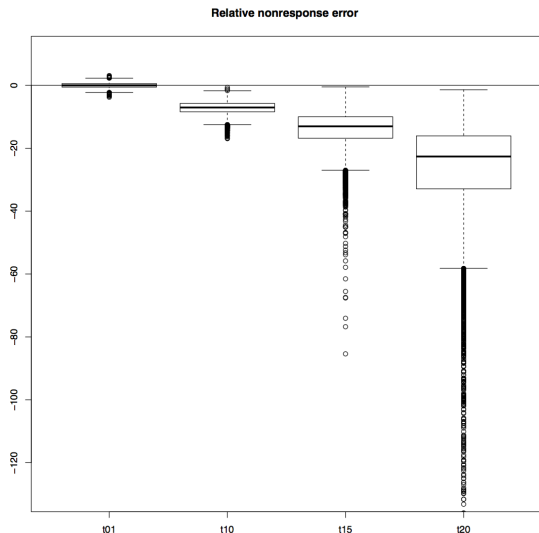
Résultats des simulations

α_1	$\alpha_2 = 0$	$\alpha_2 = 0.1$	$\alpha_2 = 0.3$	$\alpha_2 = 0.5$
0.7	0.02 (0.9)	-0.9 (0.9)	-2.8 (1.0)	-4.9 (1.1)
0.5	-0.1 (1.4)	-1.3 (1.5)	-4.1 (1.7)	-7.2 (2.1)
0.3	-0.2 (2.6)	-2.4 (3.0)	-7.5 (4.1)	-14.0 (5.9)
0.2	-0.6 (4.6)	-4.5 (15.6)	-13.8 (61.9)	-27.4 (65.6)

Table: Biais relatif Monte Carlo et CV Monte Carlo (en %)



Boxplot des erreurs relatives (en %) pour différentes paires (α_1, α_2)



Seconde étude par simulation

- On considère que la variable liée à la non-réponse est la variable d'intérêt y .
- La variable proxy $x^{(\alpha_1, \alpha_2)}$ est généré selon le modèle

$$X_i^{(\alpha_1, \alpha_2)} = \alpha_1 \text{Var}(y)^{-1}(y_i - \beta_0) + \alpha_2 u_i + \sigma_{(\alpha_1, \alpha_2)} \varepsilon_i^{(\alpha_1, \alpha_2)}$$

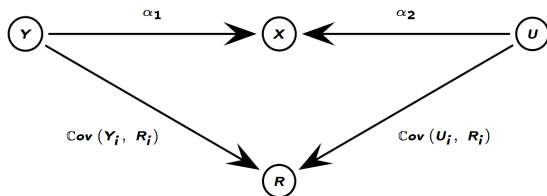
- Probabilité de répondre de l'unité i :

$$\text{logit}(p_i) = -5 + 0.5y_i + u_i.$$

- Taux de réponse moyen : 50%
- On calcule $\hat{t}_C(\alpha_1, \alpha_2)$ avec $z = y$.

Résultats des simulations

$\alpha_1 \alpha_2$	0	0.1	0.3	0.5
0.7	-0.01 (0.8)	-1.0 (0.8)	-3.3 (0.9)	-5.7 (1.0)
0.5	-0.1 (1.4)	-1.5 (1.5)	-4.8 (1.7)	-8.3 (2.0)
0.3	-0.2 (2.6)	-2.7 (2.9)	-8.7 (4.0)	-15.5 (5.5)
0.2	-0.5 (5.0)	-4.7 (7.1)	-15.0 (25.2)	-30.2 (39.1)



Conclusion

- L'utilisation d'un estimateur par **calage généralisé pour traiter la non-réponse** requiert une bonne connaissance du modèle de non-réponse i.e. :
 - 1 des variables “explicatives” de la non-réponse, \mathbf{z} ;
 - 2 de la **forme de la fonction** du modèle linéaire généralisée (qui permet de choisir la fonction de calage).

$$\mathbb{E}(R_i | \mathbf{Z}_i) = F^{-1}(\boldsymbol{\lambda}^\top \mathbf{Z}_i)$$

- et des variables instrumentales \mathbf{x} du modèle de non-réponse qui peuvent servir de proxy pour \mathbf{z} :
 - 1 qui n'interviennent pas dans le modèle de non-réponse (absence de biais)

$$\mathbb{E}(R_i | \mathbf{X}_i, \mathbf{Z}_i) = \mathbb{E}(R_i | \mathbf{Z}_i)$$

- 2 et qui soient fortement corrélées aux variables instrumentales (non-amplification du biais et de la variance).

Merci de votre attention.

bibliographie partielle I



Cassel, C. M., Särndal, C. E., and Wretman, J. H. (1976).

Some results on generalized difference estimation and generalized regression estimation for finite populations.

[Biometrika](#), 63(3), 615-620.



Chang, T. and Kott, P. (2008).

Using calibration weighting to adjust for nonresponse under a plausible model.

[Biometrika](#), 95(3) :555–571.



Deville, J.-C. and Särndal, C.-E. (1992).

Calibration estimators in survey sampling.

[Journal of the American Statistical Association](#), 87(418) :376–382.



Deville, J.-C. (1998).

La correction de la non-réponse par calage ou par échantillonnage équilibré.

In [Actes du colloque de la Société Statistique du Canada, Sherbrooke, Canada](#).



Deville, J.-C. (2004).

La correction de la non-réponse par calage généralisé.

[Actes des journées de méthodologie statistique](#), pages 4–20.



Dupont, F. (1996).

Calage et redressement de la non-réponse totale.

In [Actes des journées de méthodologie statistique, 15 et 16 décembre 1993](#), number 56. INSEE-Méthodes.

bibliographie partielle II



Kott, P. S. (2006).

Using calibration weighting to adjust for nonresponse and coverage errors.
[Survey Methodology](#), 32(2) :133.



Kott, P. S. and Chang, T. (2010).

Using calibration weighting to adjust for nonignorable unit nonresponse.
[Journal of the American Statistical Association](#), 105(491) :1265–1275.



Kott, P. S. and Liao, D. (2012).

Providing double protection for unit nonresponse with a nonlinear calibration-weighting routine.
In [Survey Research Methods](#), volume 6, pages 105–111.



Lesage, E. (2012).

Correction de la non-réponse non ignorable par une approche modèle.
In [Actes des Journées de Méthodologie Statistique de l'INSEE](#).



Le Guennec, J. et Sautory, O. (2004).

Correction de la non-réponse par calage généralisé : une expérimentation.
[Actes des journées de méthodologie statistique, 16 et 17 décembre 2002](#)



Lundström, S. and Särndal, C.-E. (1999).

Calibration as a standard method for treatment of nonresponse.
[Journal of Official Statistics](#), 15(2) :305–327.

bibliographie partielle III



Osier, G.

Dealing with non-ignorable non-response using generalised calibration : A simulation study based on the luxemburgish household budget survey.

[Economie et Statistiques, Working papers du STATEC, \(65\).](#)



Sautory, O. (2003).

Calmar 2 : une nouvelle version du programme calmar de redressement d'échantillon par calage.

In [Recueil : Symposium de Statistique Canada.](#)



Särndal, C.-E. (2011).

Three factors to signal non-response bias with applications to categorical auxiliary variables.

[International Statistical Review, 79\(2\) :233–254.](#)



Särndal, C.-E. and Lundström, S. (2010).

Design for estimation : Identifying auxiliary vectors to reduce nonresponse bias.

[Survey Methodology, 36 :131–144.](#)



Särndal, C.-E. and Lundström, S.(2005).

Estimation in surveys with nonresponse.

Wiley Hoboken, NJ.