

Prise en compte des équilibrages dans l'estimation de la variance transversale EEC-TH

version définitive

Karim MOUSSALLAM()*

() Insee, Unité de Méthodologie Statistique-Ménages*

Introduction

L'étude explore comment faire pour mettre en pratique la version 'optimale' de l'approximation de la variance d'un sondage équilibré proposée par Deville et Tillé [DT].

Le tableau 1 des tests de calcul de variance de l'article [DT] suggère assez clairement une supériorité empirique de la technique d'approximation qui identifie les coefficients diagonaux de la forme quadratique approximante avec ceux de l'estimateur sans biais de la variance. En effet, la pondération ainsi optimisée fournit la plus faible sous-estimation 'maximale', en un sens défini dans l'article, pour 6 des 9 plans de sondage simulés, et la plus petite surestimation maximale pour 4 cas. De plus, elle n'est la pire méthode pour aucun des contextes testés, sur ces deux critères. Les résultats obtenus par A.Matei et Y.Tillé pour un tirage de taille fixe à entropie maximale [AMYT] paraissent plus mitigés. Cependant, ils ne sont pas défavorables à cette méthode. Elle évite les surestimations et sous-estimations les plus extrêmes.

Compte tenu de la généralisation de l'équilibrage dans les plans de sondage des ménages réalisés par l'Insee, et de l'absence d'alternative connue à ce mode d'estimation de la variance d'un sondage équilibré, il paraît utile d'étudier les moyens de mettre en œuvre la technique d'optimisation de la pondération, et d'évaluer sa qualité pour l'estimation de variance sur les données réelles d'un sondage particulier.

Le principe de choisir la pondération des résidus qui rapproche la diagonale de la forme quadratique approximante de celle de l'estimateur sans biais de la variance reste éclairant pour améliorer la qualité de l'estimation de variance dans le contexte étudié. Pour l'appliquer, l'algorithme récursif est simple, rapide et produit la pondération la plus minimisante sur les données étudiées, par rapport aux alternatives testées. Au final toutefois, l'amélioration de la qualité de l'estimation de variance par cette optimisation ressort moins clairement que dans les tests de [DT]. L'absence de prise en compte directe de la variance d'atterrissage (cf [GC]) contribue sans doute à limiter l'enjeu de la pondération pour la qualité de l'estimation de variance et complique l'évaluation de sa méthode.

La première partie de l'étude s'intéresse aux propriétés théoriques de l'approximation de variance et à la technique d'optimisation de la pondération. La seconde partie évalue l'efficacité des choix envisagés sur des simulations de l'échantillonnage de première phase de l'enquête Emploi dans la Taxe d'Habitation (appelée EEC-TH dans la suite). Le plan de ce sondage est décrit dans [VL]. Les notations utilisées dans la suite sont explicitées en [Annexe G](#).

1 rappel du principe d'estimation de la variance d'un sondage équilibré

– Pour un sondage sans remise, la variance de l'estimateur d'Horvitz-Thompson du total Y d'une variable d'intérêt y sur une population \mathcal{P} est une forme quadratique en $\frac{y}{\pi}(\mathcal{P})$:

$$Var(\widehat{Y}) = \sum_{i,j \in \mathcal{P}} \frac{y}{\pi}(i) \frac{y}{\pi}(j) (\pi(i,j) - \pi(i)\pi(j)) \quad (\text{e1.1})$$

(où $\pi(i) = p(i \in s) > 0$)

Si la probabilité d'inclusion double ($\pi(i,j) = p(i,j \in s)$) ne s'annule pas sur \mathcal{P}^2 alors l'estimateur sans biais (ESB) de la variance est la forme quadratique en $\frac{y}{\pi}(s)$ ¹ :

$$\widehat{Var}_{HT}(\widehat{Y}) = \widehat{V}_{HT} = \sum_{i,j \in s} \frac{y}{\pi}(i) \frac{y}{\pi}(j) \frac{\pi(i,j) - \pi(i)\pi(j)}{\pi(i,j)} \quad (\text{e1.2})$$

$$= Q^{HT} \left(\frac{y}{\pi}(s) \right)$$

– La forme quadratique (e1.1) est positive, car exactement égale à $Var \left[\sum_s \frac{y}{\pi} \right]$. Par contre, ce n'est pas toujours le cas pour l'estimateur sans biais (e1.2).

→ Si l'échantillonnage est de taille fixe et de probabilités d'inclusion doubles non nulles et si la condition de Yates-Grundy ($\forall i \neq j, \pi(i)\pi(j) \geq \pi(i,j)$) est satisfaite alors :

$$\widehat{V}_{YG} = \frac{1}{2} \sum_{i,j \in s} \left(\frac{y}{\pi}(i) - \frac{y}{\pi}(j) \right)^2 \frac{\pi(i)\pi(j) - \pi(i,j)}{\pi(i,j)} \quad (\text{e1.3})$$

est positif et estime sans biais la variance. En général $\widehat{V}_{YG} \neq \widehat{V}_{HT}$, sauf pour un SAS². La diagonale de la forme quadratique de \widehat{V}_{YG} en fonction de y/π a pour coefficient générique :

$$1 - \pi_i + \sum_{j \in s} \frac{\pi_i \pi_j}{\pi_{i,j}} - |s| \quad (\geq 1 - \pi_i \text{ si taille fixe, vu que } \pi_j \geq \pi_{ij})$$

– Pour un sondage équilibré d'entropie maximale, c'est-à-dire un sondage poissonnien p^* conditionné par l'équation d'équilibrage ($p(s) = p^* [s | \widehat{X} = x(+)]$), la variance asymptotique est approximée en régressant la variable d'intérêt par les variables d'équilibrage avec la pondération $c^* = \frac{\pi^*(1 - \pi^*)}{\pi}$:

$$Var_{p^*} [\widehat{Y} | \widehat{X} = x(+)] \cong Var_{p^*} (\widehat{Y} - \widehat{X}'\beta^*) \text{ avec}$$

$$\beta^* \in \underset{\beta}{\operatorname{argmin}} Var_{p^*} (\widehat{Y} - \widehat{X}'\beta), \text{ et comme } p^* \text{ est poissonnien :}$$

$$Var_{p^*} [\widehat{Y} | \widehat{X} = x(+)] \cong \sum_{\mathcal{P}} \left(\frac{y}{\pi} - \frac{x'}{\pi} \beta^* \right)^2 \pi^*(1 - \pi^*) \text{ où :}$$

$$\beta^* \in \underset{\beta}{\operatorname{argmin}} \sum_{\mathcal{P}} \left(\frac{y}{\pi} - \frac{x'}{\pi} \beta \right)^2 \pi^*(1 - \pi^*)$$

Soit M^* la mesure sur \mathcal{P} pondérée par $\pi^*(1 - \pi^*)$ et \widehat{M}^* la mesure empirique associée au plan de sondage effectif (telle que $E_p(\widehat{M}^*) = M^*$) : $\widehat{M}^*(z) = \sum_s z \frac{\pi^*(1 - \pi^*)}{\pi}$.

1. C'est la seule forme quadratique en $y(s)$ à coefficients fixes /s qui soit sans biais pour toute variable $y(\mathcal{P})$.
2. Pour un sondage aléatoire simple (échantillons de tailles fixes et équiprobables) :

$$\widehat{V}_{HT}(\widehat{Y}) = \frac{N(N-n)}{n} \mathcal{Q}_s^2(y) = \frac{N(N-n)}{n} \frac{1}{n-1} \sum_s (y - \bar{y}_s)^2 = \frac{1}{2} \frac{N(N-n)}{n-1} \frac{1}{n^2} \sum_{i,j \in s} (y_i - y_j)^2$$

→ L'estimation approximative de la variance équilibrée consiste à remplacer la mesure inconnue \widehat{M}^* par une approximation \widetilde{M}^* calculable sur l'échantillon sous la forme : $\widetilde{M}^*(z) = \sum_s cz$. Par la méthode de linéarisation [JCD], l'approximation de variance est formulée ainsi :

$$\widehat{Var}(\widehat{Y}) = \widetilde{M}^* \left[\left(\frac{y}{\pi} - \frac{x'}{\pi} \widehat{\beta} \right)^2 \right] = Q \left(\frac{y}{\pi}(s) \right), \text{ avec :} \quad (\text{e1.4})$$

$$\widehat{\beta} \in \underset{\beta}{\operatorname{argmin}} \widetilde{M}^* \left[\left(\frac{y}{\pi} - \frac{x'}{\pi} \beta \right)^2 \right]$$

⇒ En pratique, la variance est estimée par la somme pondérée des carrés des résidus du ratio de la variable d'intérêt sur la probabilité d'inclusion régressé par le vecteur d'équilibrage. La pondération d'équilibrage est c , et le résidu d'équilibrage est $(y - x\widehat{\beta})/\pi$.

Il faut souligner que la justification méthodologique de ce type d'estimateur de variance est asymptotique. De plus, l'approximation est affectée d'un biais fonction de la proximité entre le plan de sondage effectif et celui équilibré d'entropie maximale. Une approximation moins forte de l'estimation de variance est possible, en estimant par simulation les probabilités d'inclusion double [GC], au lieu d'approcher seulement les probabilités d'inclusion simple.

→ Il est clair que l'enjeu méthodologique essentiel dans le cadre de l'approximation [DT] de la variance est le choix de la pondération $c(s)$ des résidus d'équilibrage. C'est l'objet de l'étude.

2 propriétés de la forme quadratique approximante

2.1 expression comme fonction de la pondération

– L'estimateur de variance défini par (e1.4) est encore une forme quadratique en $\frac{y}{\pi}(s)$.

Avec les notations $u = \frac{x}{\pi}$ et $z = \frac{y}{\pi}$, il peut se formuler ainsi :

$$\widehat{Var}(\widehat{Y}) = \min_{\beta} \left\| \frac{y}{\pi} - \frac{x'}{\pi} \beta \right\|_c^2 = \min_{\beta} \|z - u_s \beta\|_c^2 \quad (\text{e2.5})$$

$$= \left\| \operatorname{proj}_{\operatorname{Im}(u_s)}^{c\perp}(z) \right\|_c^2 \quad (c \text{ est confondue avec la métrique diagonale correspondante, } *c)$$

$$= z^s \operatorname{proj}_{\operatorname{Im}(u_s)}^{c\perp}{}' * c \operatorname{proj}_{\operatorname{Im}(u_s)}^c z_s = z^s * c \operatorname{proj}_{\operatorname{Im}(u_s)}^c z_s = z^s (*c - c \operatorname{proj}_{\operatorname{Im}(u_s)}^c) z_s$$

$$= z^s (*c - cu_s (u^s cu_s)^- u^s c) z_s \text{ d'où la formule explicite :}$$

$$\widehat{Var}(\widehat{Y}) = z^s \Delta(c) z_s \quad (\text{e2.6})$$

$$\text{avec : } \Delta(c) = \operatorname{diag}(c) - (cu)_s [\sum_s cuu']^- (cu)^s \quad (\text{e2.7})$$

– La formule (e2.7) peut également se déduire de l'équation normale :

$$\min_{\beta} \|z - u_s \beta\|_c^2 = \|z - u_s \widehat{\beta}\|_c^2 \text{ où } \widehat{\beta} \text{ est déterminée }^3 \text{ par :}$$

$$\langle z - u_s \widehat{\beta}, u_s \rangle_c = 0 \text{ (équation normale)} \quad (\text{e2.8})$$

$$\text{d'où : } \|z - u_s \widehat{\beta}\|_c^2 = \langle z - u_s \widehat{\beta}, z \rangle_c = \|z\|_c^2 - \langle u_s \widehat{\beta}, z \rangle_c$$

$$= z' * cz - z' * cu_s (\sum cuu')^- u^s * cz = z' * c (\operatorname{id} - u_s (\sum cuu')^- u^s * c) z$$

– Plus explicitement, l'estimateur de la variance équilibrée peut se calculer ainsi :

$$\widehat{V} = \sum_s c \left(\frac{y}{\pi} \right)^2 - \sum_s \frac{yx'}{\pi^2} \left[\sum_s c \frac{xx'}{\pi^2} \right]^- \sum_s c \frac{yx}{\pi^2} \quad (\text{e2.9})$$

3. de manière unique si et seulement si u_s est de plein rang colonne.

– Dans la suite, pour l'indicatrice de l'élément $i \in s$, $\mathbf{1}_i$, le coefficient de la régression est noté $\widehat{\beta}_i$ ou $\widehat{\beta}_c(\mathbf{1}_i)$.

– Vu (e2.7) et (e2.5), les coefficients diagonaux $\delta(c)$ de la forme quadratique approximante se formulent de deux façons, la première utile au calcul et la seconde pour l'étude de leurs propriétés :

$$\delta_i(c) = \Delta_{i,i}(c) = c \left(1 - cu' \left(\sum_c uu'\right)^- u\right)(i) = c_i(1 - u'_i \widehat{\beta}_i) \quad (\text{e2.10})$$

$$= \min_{\beta} \| \mathbf{1}_i - u_s \beta \|_c^2 = \| \mathbf{1}_i - u_s \widehat{\beta}_i \|_c^2 = \sum_s c (\mathbf{1}_i - u'_i \widehat{\beta}_i)^2 \quad (\text{e2.11})$$

$$= \| \mathbf{1}_i \|_c^2 - \| u_s \widehat{\beta}_i \|_c^2 = c_i - \| u_s \widehat{\beta}_i \|_c^2 \quad (\text{e2.12})$$

– Il découle de ces égalités que $c_i(1 - u'_i \widehat{\beta}_i) u'_i \widehat{\beta}_i = \sum_{\neq i} c [u'_i \widehat{\beta}_i]^2$ et donc $0 \leq u'_i \widehat{\beta}_i \leq 1$.

Cet encadrement de la ième coordonnée de la projection de $\mathbf{1}_i$ se construit de manière plus satisfaisante ainsi :

$$u'_i \widehat{\beta}_i = \langle \mathbf{1}_i, \text{proj}_{\text{Im}(u_s)}^c [\mathbf{1}_i] \rangle = \frac{1}{c_i} \langle \mathbf{1}_i, \text{proj}_{\text{Im}(u_s)}^c [\mathbf{1}_i] \rangle_c = \frac{1}{c_i} \| \text{proj}_{\text{Im}(u_s)}^c [\mathbf{1}_i] \|_c^2 \leq \frac{1}{c_i} \| \mathbf{1}_i \|_c^2 \leq 1$$

– L'expression (e2.10) entraîne que $u' \left(\sum_c uu'\right)^- u$ est constant sur les inverses généralisées de $\sum_c uu'$. Il est utile pour la suite de l'étude d'en analyser le mécanisme.

• La démonstration la plus directe utilise le fait que $[\sqrt{cu}]_s \left(\sum_c uu'\right)^- [\sqrt{cu}]_s^s$ est la matrice de la projection orthogonale sur $\text{Im} \{ [\sqrt{cu}]_s \}$ (eC.65), donc ne dépend pas du choix de l'inverse généralisée.

• $\text{Im} \left(\sum_c uu'\right) = \text{Im} \{ [\sqrt{cu}]_s^s [\sqrt{cu}]_s \} = \text{Im} \{ [\sqrt{cu}]_s^s \}$ (par (eA.57)). D'où :

$$\text{Im} \left(\sum_c uu'\right) = \text{Im} [u^s * \sqrt{c}] = \text{Im} (u^s) \quad (\text{vu que } \text{Im} \{ \text{diag} [\sqrt{c}] \} = \mathbb{R}_s \text{ si } c > 0)$$

$$\Rightarrow \forall c > 0, \text{Im} \left(\sum_c uu'\right) = \text{Im} (u^s) = \text{Im} \left(\sum_c uu'\right) \quad (\text{e2.13})$$

$$\rightarrow (\text{e2.13}) \Rightarrow \forall c > 0, \text{rang} \left(\sum_c uu'\right) = \text{rang} (u^s) (= \text{rang} (u_s))$$

$$\Rightarrow \text{si } c > 0 \text{ alors } (\text{e2.13}) \Rightarrow u_i \in \text{Im} \left(\sum_c uu'\right) \text{ car } u_i = u^s \mathbf{1}_i \in \text{Im} (u^s) = \text{Im} \left(\sum_c uu'\right)$$

$$\Rightarrow \text{si } c > 0, \exists a_i / u_i = \sum_c uu' a_i \quad (\text{e2.14})$$

(soit synthétiquement : $\exists A / u_s = A \sum_c uu'$)

$$\Rightarrow u'_i \left(\sum_c uu'\right)^- u_j = a'_i \left(\sum_c uu'\right) \left(\sum_c uu'\right)^- \left(\sum_c uu'\right) a_j = a'_i \left(\sum_c uu'\right) a_j = \sum_c c (a'_i u) (a_j u)$$

donc ne dépend pas du choix de l'inverse généralisée.

→ La forme quadratique approximante (e2.5) est manifestement positive. De plus elle s'annule pour $z_s \in \text{Im} (u_s) \Leftrightarrow \exists \beta / \frac{y}{\pi}(s) = \frac{x}{\pi}(s) \beta \Leftrightarrow y_s \in \text{Im} (x_s)$. C'est une bonne propriété de l'approximation de variance, puisque la variance du total estimé d'une variable qui vérifie cette condition pour tout s est effectivement nulle si l'équilibrage est exact :

$$\sum_s \frac{y}{\pi} = \sum_s \left(\frac{x}{\pi}\right)' \beta = x (+)' \beta = y (+)$$

– L'approximation de variance est une fonction de la pondération c notée $\varphi(c) = \varphi_z(c) = \min_{\beta} \| z - u_s \beta \|_c^2$. Cette fonction peut être définie ainsi sur \mathbb{R}_s^+ . Outre la positivité, elle possède plusieurs propriétés intéressantes. L'étude de celles-ci est également utile pour $\Delta(c)$ et $\delta(c)$.

– Comme $\text{Im} (u_s)$ est un sous-espace vectoriel de dimension finie, le minimum $\varphi(c) = \min_{\beta} \| z - u_s \beta \|_c^2$ est toujours atteint par la projection orthogonale de z_s sur l'image de u_s pour la métrique c . Donc il existe au moins un $\widehat{\beta}_c \in \text{argmin}_{\beta} \| z - u_s \beta \|_c^2$ et avec un tel vecteur $\varphi(c) = \| z - u_s \widehat{\beta}_c \|_c^2 = \| z \|_c^2 - \| u_s \widehat{\beta}_c \|_c^2$. Cette formule de la fonction ne dépend pas du vecteur choisi dans $\text{argmin}_{\beta} \| z - u_s \beta \|_c$ (Annexe C). (Lorsque c'est nécessaire, la notation plus précise $\widehat{\beta}_c(z)$ est utilisée.)

– La fonction φ est croissante au sens où :

$$c \leq \tilde{c} \quad (\Leftrightarrow \forall i, c_i \leq \tilde{c}_i) \Rightarrow \forall \beta \| z - u_s \beta \|_c^2 \leq \| z - u_s \beta \|_{\tilde{c}}^2 \Rightarrow \varphi(c) \leq \varphi(\tilde{c})$$

- Homogénéité de degré 1 :

$$\forall \lambda \geq 0, \varphi(\lambda c) = \min_{\beta} \|z - u_s \beta\|_{\lambda c}^2 = \min_{\beta} \lambda \|z - u_s \beta\|_c^2 = \lambda \min_{\beta} \|z - u_s \beta\|_c^2 = \lambda \varphi(c)$$

Un intérêt pratique de cette propriété est qu'à partir d'une pondération c telle que $\sum_s \delta(c) > 0$, il est toujours possible de construire une pondération \tilde{c} dont la somme $\sum_s \delta(\tilde{c})$ est égale à la trace de l'estimateur sans biais, simplement par un facteur multiplicatif.

- $\varphi(c)$ est concave sur \mathbb{R}_s^+ :

$$\begin{aligned} \forall \alpha \in [0, 1], \forall \beta, \|z - u_s \beta\|_{\alpha c + (1-\alpha)\tilde{c}}^2 &= \alpha \|z - u_s \beta\|_c^2 + (1-\alpha) \|z - u_s \beta\|_{\tilde{c}}^2 \\ &\geq \alpha \varphi(c) + (1-\alpha) \varphi(\tilde{c}) \end{aligned}$$

- Cette propriété équivaut à la convexité en c de $\|u_s \widehat{\beta}_c\|_c^2 (= \sum cu'z (\sum cuu')^{-1} \sum cuz)$, comme somme de deux fonctions convexes : $\|u_s \widehat{\beta}_c\|_c^2 = \|z\|_c^2 - \varphi(c)$ ⁴.

- si $a \geq 0$ et $c \geq 0$ alors $\varphi(a+c) \geq \varphi(a) + \varphi(c)$ car $\varphi(a+c) = \|z - u_s \widehat{\beta}_{a+c}\|_{a+c}^2 = \|z - u_s \widehat{\beta}_{c+a}\|_{c+a}^2 + \|z - u_s \widehat{\beta}_a\|_a^2 \geq \|z - u_s \widehat{\beta}_a\|_a^2 + \|z - u_s \widehat{\beta}_c\|_c^2$

2.2 continuité sur \mathbb{R}_s^+ :

- semi-continuité supérieure : $\varphi(c) = \|z - u_s \widehat{\beta}_c\|_c^2 + \|z - u_s \widehat{\beta}_c\|_{c-\tilde{c}}^2 \geq \varphi(\tilde{c}) + \|z - u_s \widehat{\beta}_c\|_{c-\tilde{c}}^2$

Comme ce dernier terme est une fonction linéaire de $c - \tilde{c}$, $\lim \varphi(c) \leq \varphi(c)$ (\lim désigne la limite supérieure).

- semi-continuité inférieure : $\lim \varphi(c) \geq \varphi(c)$ est vrai pour $c = 0$, car $\varphi(0) = 0$. Supposons $c \neq 0$. Alors pour $\|\tilde{c} - c\|$ assez petit, $\tilde{c} \setminus \{c > 0\} > 0$, et sur ce voisinage de c dans \mathbb{R}_s^+ :

$$\varphi(c) \leq \|z - u_s \widehat{\beta}_c\|_c^2 \leq \|z - u_s \widehat{\beta}_{\tilde{c}}\|_{\tilde{c}}^2 \max_{c>0} \left(\frac{c}{\tilde{c}} \right) = \varphi(\tilde{c}) \max_{c>0} \left(\frac{c}{\tilde{c}} \right)$$

La limite inférieure de la dernière expression égale $\lim \varphi(c)$, car $\lim_{\tilde{c} \rightarrow c} \max_{c>0} \left(\frac{c}{\tilde{c}} \right) = 1 > 0$.

- Une démonstration alternative est éclairante, bien que limitée à un ouvert :

- D_+ l'ensemble des matrices définies positives est un ouvert : si $\min_{\|x\|=1} x'Ax > 0$ alors $\min_{\|x\|=1} x'Bx \geq$

$$\min_{\|x\|=1} x'Ax - \|B - A\| \text{ parce que } x'(B - A)x \geq -\|B - A\| \|x\|^2 \quad 5$$

→ $\varphi(A) = \min_{\beta} \|z - u_s \beta\|_A^2$ est continue sur D_+ :

▷ Comme A est définie positive et que toutes les normes sont équivalentes en dimension finie, $\exists \lambda(A) > 0 / \forall x, \|x\|^2 \leq \lambda(A) \|x\|_A^2$ (avec $\tilde{A} = (A + A') / 2$, $\lambda(A) = \|\tilde{A}^{-1}\|$ est une solution vu que $\|x\|^2 = x' \tilde{A}^{\frac{1}{2}} \tilde{A}^{-1} \tilde{A}^{\frac{1}{2}} x \leq \|\tilde{A}^{-1}\| \|\tilde{A}^{\frac{1}{2}} x\|^2 = \|\tilde{A}^{-1}\| x'Ax$)

▷ soit $\mathcal{O}^{\mathcal{O}}(A)$ un voisinage de A dans D_+ tel que $\|\mathcal{O}^{\mathcal{O}}(A) - A\| < 1 / \lambda(A)$

▷ pour $\tilde{A} \in \mathcal{O}^{\mathcal{O}}(A)$ et pour tout β :

$$\|z - u_s \beta\|_A^2 - \|z - u_s \beta\|_{\tilde{A}}^2 = |(z - u_s \beta)'(A - \tilde{A})(z - u_s \beta)|$$

$$\leq \|\tilde{A} - A\| \|z - u_s \beta\|^2 \leq \|\tilde{A} - A\| \lambda(A) \|z - u_s \beta\|_A^2$$

$$\Rightarrow \|z - u_s \beta\|_A^2 [1 - \|\tilde{A} - A\| \lambda(A)] \leq \|z - u_s \beta\|_{\tilde{A}}^2 \leq \|z - u_s \beta\|_A^2 [1 + \|\tilde{A} - A\| \lambda(A)]$$

$$\Rightarrow \varphi(A) [1 - \|\tilde{A} - A\| \lambda(A)] \leq \varphi(\tilde{A}) \leq \varphi(A) [1 + \|\tilde{A} - A\| \lambda(A)] \text{ (car } 1 - \|\tilde{A} - A\| \lambda(A) > 0)$$

$$\Rightarrow |\varphi(\tilde{A}) - \varphi(A)| \leq \|\tilde{A} - A\| \lambda(A) \varphi(A)$$

4. Incidemment, avec $z = \mathbf{1}_i$, il en résulte que $c_i^2 u_i' (\sum cuu')^{-1} u_i (= \|u_s \widehat{\beta}_c(\mathbf{1}_i)\|_c^2)$ est convexe en c_s .

5. Si $\|\cdot\|$ désigne la norme euclidienne canonique, $|x'(B - A)x| \leq \|x\| \|(B - A)x\| \leq \|x\| \|B - A\| \|x\|$, vu (eE.69).

- majoration explicite de la variation de $\delta_i(c)$, pour $c > 0$:
 $(\min \{\tilde{c}/c\} - 1) \delta_i(c) \leq \delta_i(\tilde{c}) - \delta_i(c) \leq (\max \{\tilde{c}/c\} - 1) \delta_i(c)$ donc :
 $|\delta_i(\tilde{c}) - \delta_i(c)| \leq \max [|\tilde{c}/c - 1|] \delta_i(c)$

2.3 dérivabilité sur \mathbb{R}_s^{+*}

Pour appliquer à $\delta(c)$ un algorithme de minimisation utilisant le gradient, il faut que cette fonction soit dérivable. Le calcul de la dérivée $\dot{\delta}(c)$ est utilisé non seulement pour programmer ces algorithmes mais également pour étudier l'existence et l'unicité de solutions.

- La démonstration est effectuée plus généralement pour $\varphi(c) = \|z - u_s \widehat{\beta}_c\|_{A(c)}$ où $A(c)$ est une fonction dérivable à valeur dans l'ensemble des matrices symétriques définies positives.

hypothèses • $A(c)$ dérivable en c

- A est symétrique définie positive sur $\mathcal{O}(c)$ ⁶

soit $c+h \in \mathcal{O}(c)$

$$\begin{aligned} \varphi(c+h) - \varphi(c) &= \|z - u_s \widehat{\beta}_{c+h}\|_{A(c+h)}^2 - \|z - u_s \widehat{\beta}_c\|_{A(c)}^2 \\ &= \|z - u_s \widehat{\beta}_c\|_{A(c+h)}^2 - \|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{A(c+h)}^2 - \|z - u_s \widehat{\beta}_c\|_{A(c)}^2 \end{aligned} \quad (\mathbf{e2.15})$$

(par Pythagore)

$$\langle u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c), u_s \rangle_{A(c+h)} = \langle z - u_s \widehat{\beta}_c, u_s \rangle_{A(c+h)} \quad (\text{équation normale pour } \widehat{\beta}_{c+h})$$

$$= \langle z - u_s \widehat{\beta}_c, u_s \rangle_{A(c+h) - A(c)} \quad (\text{équation normale pour } \widehat{\beta}_c)$$

$$= \langle A(c+h)^{-1}(A(c+h) - A(c))(z - u_s \widehat{\beta}_c), u_s \rangle_{A(c+h)}$$

$$\Rightarrow \|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{A(c+h)}^2 \leq \|A(c+h)^{-1}(A(c+h) - A(c))(z - u_s \widehat{\beta}_c)\|_{A(c+h)} \|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{A(c+h)}$$

(par l'inégalité de Cauchy-Schwartz)

$$\Rightarrow \|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{A(c+h)} \leq \|A(c+h)^{-1}(A(c+h) - A(c))(z - u_s \widehat{\beta}_c)\|_{A(c+h)}$$

$$= \sqrt{(z - u_s \widehat{\beta}_c)'(A(c+h) - A(c))A(c+h)^{-1}(A(c+h) - A(c))(z - u_s \widehat{\beta}_c)}$$

$$\leq \|z - u_s \widehat{\beta}_c\| \|A(c+h) - A(c)\| \sqrt{\|A(c+h)^{-1}\|}$$

$$= O(h) \quad ^7$$

car le premier terme est constant / h , le second est un $O(h)$ par la dérivabilité de $A(c)$ et le dernier terme converge vers $\sqrt{\|A(c)^{-1}\|} < +\infty$ par les continuités de la racine carrée, de $A(c)$ et de l'inversion sur les matrices inversibles. Comme $\|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{A(c+h)} = O(h) \Rightarrow \|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{A(c+h)}^2 = o(h)$, et la dérivabilité du premier terme de (e2.15) qui s'ensuit entraîne celle de $\varphi(c)$, ainsi que l'expression de la dérivée.

→ La dérivée de φ peut se formuler ainsi :

$$\dot{\varphi}(c) = \|z - u_s \widehat{\beta}_c\|_{dc}^2 \quad \text{où } \widehat{\beta}_c \in \underset{\beta}{\operatorname{argmin}} \|z - u_s \beta\|_c \quad (\mathbf{e2.16})$$

– Plus simplement, la formule de la dérivée peut s'obtenir en traitant $\widehat{\beta}_c$ comme une fonction dérivable de c (c'est effectivement le cas avec le choix $\widehat{\beta}_c = (\sum c u u')^+ \sum c u z$) :

$$\dot{\varphi}(c) = \|z - u_s \widehat{\beta}_c\|_{dc} + 2 \langle z - u_s \widehat{\beta}_c, u_s d\widehat{\beta}_c \rangle_c = \|z - u_s \widehat{\beta}_c\|_{dc} \quad (\text{vu l'équation normale})$$

6. donc $\|\cdot\|_{A(c)}$ est une norme euclidienne

7. $o(h) : \lim_{h \rightarrow 0} \frac{\|o(h)\|}{\|h\|} = 0$ et $O(h) : \lim_{h \rightarrow 0} \frac{\|O(h)\|}{\|h\|} < +\infty$.

(C'est donc l'orthogonalité entre $z - u_s \widehat{\beta}_c$ et u_s qui fait disparaître $d\widehat{\beta}_c$ de l'expression de $\dot{\varphi}(c)$. Il n'est pas rigoureux de considérer que $d\widehat{\beta}_c$ est négligeable en soit.)

- La norme de la dérivée peut s'explicitier ainsi, pour la norme du 'sup' $\|\cdot\|_\infty$:

$$\|\dot{\varphi}(c)\| = \sup_{\|h\|_\infty=1} \left| \sum [z - u' \widehat{\beta}_c]^2 h \right| \quad (\text{e2.17})$$

$$= \sum_s [z - u_s \widehat{\beta}_c]^2 \left(\leq \frac{\varphi(c)}{\min(c)} \text{ et } \geq \varphi(1) \right)$$

- L' **Annexe D** prouve que φ est deux fois dérivable sur $\{c > 0\}$, et même C^∞ .

- Puisque pour tout y , la fonction $\varphi(c) = y' \Delta(c) y$ est dérivable, toutes les composantes de la matrice $\Delta(c)$ sont dérivables en c , et donc Δ est dérivable. Par suite sa diagonale δ est dérivable sur $\{c > 0\}$. Une formule explicite de la dérivée de $\delta(c)$ se déduit de (e2.16) :

$$\dot{\delta}_i(c) = \|\mathbf{1}_i - u_s \widehat{\beta}_c(\mathbf{1}_i)\|_{dc}^2 \quad (\text{e2.18})$$

$$\dot{\delta}_i(c)^j = (\mathbf{1}_i^j - u'_j (\sum c u u')^- c_i u_i)^2 \quad 8$$

$$\dot{\delta}(c) = (\text{id} - (c u)_s (\sum c u u')^- u^s)^{*2} \quad (\text{e2.19})$$

La dérivée se formule aussi : $\dot{\delta}_s(c) = \|\mathbf{1}_s - u \widehat{\beta}_c(\mathbf{1}_s)\|_{dc}^2$ (l'indice intra-norme est occulté ⁹).

→ expression des dérivées partielles de la diagonale (en utilisant (e2.10)) :

$$\frac{\partial \delta_i(c)}{\partial c_i} = (1 - u'_i \widehat{\beta}_c(\mathbf{1}_i))^2 = (1 - u'_i (\sum c u u')^- c_i u_i)^2 = \left[\frac{\delta_i(c)}{c_i} \right]^2 \quad (\text{e2.20})$$

$$i \neq j : \frac{\partial \delta_i(c)}{\partial c_j} = (u'_j \widehat{\beta}_c(\mathbf{1}_i))^2 = (u'_j (\sum c u u')^- c_i u_i)^2 \quad (\text{e2.21})$$

- Pour comparer les deux dérivées ci-dessus, il peut être utile d'utiliser les formules :

$$1 - u'_i \widehat{\beta}_i = \langle \mathbf{1}_i, \text{proj}_{\text{Im}(u_s)}^{c^\perp} [\mathbf{1}_i] \rangle = \frac{1}{c_i} \|\text{proj}_{\text{Im}(u_s)}^{c^\perp} [\mathbf{1}_i]\|_c^2 \quad (\text{e2.22})$$

$$\begin{aligned} \text{si } j \neq i, u'_j \widehat{\beta}_i &= \langle \mathbf{1}_j, \text{proj}_{\text{Im}(u_s)}^{c^\perp} [\mathbf{1}_i] \rangle = \frac{1}{c_j} \langle \mathbf{1}_j, \text{proj}_{\text{Im}(u_s)}^{c^\perp} [\mathbf{1}_i] \rangle_c \\ &= \frac{1}{c_j} \langle \text{proj}_{\text{Im}(u_s)}^{c^\perp} [\mathbf{1}_j], \text{proj}_{\text{Im}(u_s)}^{c^\perp} [\mathbf{1}_i] \rangle_c \end{aligned} \quad (\text{e2.23})$$

→ de (e2.20) et (13.2, ci-après) on déduit :

$$\text{si } \delta_i(1) > 0 \text{ alors } \delta_i(c) \text{ croit strictement en fonction de } c_i > 0 \quad (\text{e2.24})$$

- dérivabilité à droite en $c \geq 0$:

Pour $c \geq 0$ (et $c \neq 0$), la fonction peut encore se formuler ainsi :

$$\varphi(c) = \min_{\beta} \|z - u_s \beta\|_c^2 = \|z - u_s \widehat{\beta}_c\|_c^2$$

et pour $c + h \geq 0$ il reste vrai que :

$$\varphi(c + h) = \|z - u_s \widehat{\beta}_{c+h}\|_{c+h}^2 = \|z - u_s \widehat{\beta}_c\|_{c+h}^2 - \|u_s (\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{c+h}^2$$

d'où $\|u_s (\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{c+h}^2 = -\varphi(c + h) + \varphi(c) + \|z - u_s \widehat{\beta}_c\|_h^2$, qui tend vers 0 lorsque $h \rightarrow 0$ et $c + h \geq 0$, vu la continuité de φ sur \mathbb{R}_s^+ et la linéarité en h du dernier terme. Il en découle que pour $c \geq 0$:

$$\lim_{\substack{h \rightarrow 0 \\ c+h \geq 0}} \|u_s (\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{c+h}^2 = 0.$$

- ou bien $c_i > 0$ et alors $|u'_i (\widehat{\beta}_{c+h} - \widehat{\beta}_c)| \xrightarrow{h \rightarrow 0} 0$ (car $\leq \|u_s (\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{c+h} / (c_i + h_i)$)
- ou bien $c_i = 0$:

$$\rightarrow \text{si } u_i \in \text{Im}(u^{c>0}) \text{ alors comme } u_{c>0} (\widehat{\beta}_{c+h} - \widehat{\beta}_c) \xrightarrow{h \rightarrow 0} 0, u'_i (\widehat{\beta}_{c+h} - \widehat{\beta}_c) \xrightarrow{h \rightarrow 0} 0$$

8. Si $\exists i \neq j / c_i \neq c_j$ et $u'_j (\sum c u u')^- u_i \neq 0$ alors $\dot{\delta}_i(c)^j \neq \dot{\delta}_j(c)^i$ et $\dot{\delta}(c)$ n'est pas une matrice symétrique.

9. une solution pour expliciter le double indiciage ($s \times s$) : $\dot{\delta}_s(c) = \|\mathbf{1}_s - u_{(s)} \widehat{\beta}_c(\mathbf{1}_s)\|_{dc(s)}^2$

→ si $u_i \notin \text{Im}(u^{c>0})$ alors nécessairement $z_i - u_i' \widehat{\beta}_{c+h} = 0$ (si $h_i > 0$), parce qu'il existe γ tel que $\begin{cases} u_i' \gamma = z_i \\ u_{c>0} \gamma = 0 \end{cases}$ (e3.41). Pour la même raison, $\widehat{\beta}_c$ peut toujours être choisi tel que $z_i - u_i' \widehat{\beta}_c = 0$, d'où $u_i'(\widehat{\beta}_{c+h} - \widehat{\beta}_c) = 0$.

⇒ Dans les deux sous-cas, $[u_i'(\widehat{\beta}_{c+h} - \widehat{\beta}_c)]^2 h_i = o(h)$. Ceci prouve la dérivabilité à droite de φ en $c \geq 0$ parce que :

$$\begin{aligned} \|u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c)\|_{c+h}^2 &= \langle u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c), u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c) \rangle_{c+h} = \langle z - u_s \widehat{\beta}_c, u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c) \rangle_{c+h} \\ &= \langle z - u_s \widehat{\beta}_c, u_s(\widehat{\beta}_{c+h} - \widehat{\beta}_c) \rangle_h = \sum h_i (z_i - u_i' \widehat{\beta}_c) u_i'(\widehat{\beta}_{c+h} - \widehat{\beta}_c) = o(h) \end{aligned}$$

L'expression de la dérivée partielle en $c_i = 0$ est :

$$\frac{\partial \varphi(c)}{\partial c_i}(c_i = 0) = \{z_i - u_i' \widehat{\beta}_c(z)\}^2 \text{ si } u_i \in \text{Im}(u^{c>0}) \text{ et } 0 \text{ sinon} \quad (\text{e2.25})$$

3 cible et contraintes de l'optimisation

3.1 critère à minimiser

– Le cadre d'approximation de la variance [DT] est appliqué en choisissant une pondération qui rapproche le plus possible la forme quadratique approximante (e2.5) de l'estimateur sans biais (e1.2). Pour concrétiser cette démarche, l'idée proposée par [DT] est d'exploiter l'information connue sur celui-ci : la diagonale de la forme quadratique $Q^{HT} \left(\frac{y}{\pi}(s) \right)$ égale $(1 - \pi)(s)$. Vu autrement, pour $z = \mathbf{1}_i$, (e2.6) donne $\widehat{V} = \delta_i(c) \mathbf{1}_i$, qui doit estimer $\text{Var}(\widehat{\pi \mathbf{1}_i}) = \text{Var}[\mathbf{1}_i] = \pi_i(1 - \pi_i)$, et comme $\delta_i(c)$ est constante sur les échantillons, $E\{\delta_i(c) \mathbf{1}_i\} = \delta_i(c) \pi_i$.

L'équation de pondération est (e3.26). Comme le montre [DT], il est équivalent¹⁰ de résoudre les deux équations :

$$\delta(c) = 1 - \pi \quad (\text{e3.26})$$

$$\text{et } \frac{\delta(c)}{\sum_s \delta(c)} = \frac{1 - \pi}{\sum_s 1 - \pi}, \text{ soit en posant } w = \frac{c}{\sum_s c} \text{ (pour } c > 0 \text{):}$$

$$\frac{\delta(w)}{\sum_s \delta(w)} = \frac{1 - \pi}{\sum_s 1 - \pi}, w > 0, \sum_s w = 1 \quad (\text{e3.27})$$

En effet, une solution c de (e3.26) donne une solution $w = c / \sum_s c$ de (e3.27). Inversement, si w est solution de (e3.27) alors $c = w \sum_s 1 - \pi / \sum_s \delta(w)$ est solution de (e3.26). Et ces deux applications sont inverses l'une de l'autre sur les deux ensembles de solutions.

– expressions de la diagonale 'réduite' :

$$F_i(w) = \frac{\delta_i(w)}{\sum_s \delta(w)} = \frac{\min_{\beta} \|\mathbf{1}_i - u_s \beta\|_w^2}{\sum_{j \in s} \min_{\beta} \|\mathbf{1}_j - u_s \beta\|_w^2} = \frac{w_i(1 - u_i' \widehat{\beta}_i)}{\sum_{j \in s} w_j(1 - u_j' \widehat{\beta}_j)} = \frac{w_i(1 - w_i u_i' [\sum_s w u u']^- u_i)}{\sum_{j \in s} w_j(1 - w_j u_j' [\sum_s w u u']^- u_j)} \quad (\text{e3.28})$$

La deuxième expression montre que cette fonction est le rapport de deux fonctions concaves. Sa convexité est donc indéterminée.

– Il n'a pas été possible d'étendre le résultat de [DT] sur la borne supérieure de $F(w)$, en raison de la présence dans la formule de la matrice $[\sum_s w u u']^-$, au lieu d'une constante. A fortiori, cette formulation n'a pas permis d'expliciter une solution exacte en fonction de $1 - \pi_s$.

¹⁰. à la condition implicite que $\forall c > 0, \sum_s \delta(c) > 0$

→ La présente étude a privilégié l'équation (e3.26) parce que la fonction δ est plus simple.

Une pondération est qualifiée d'optimale si elle vérifie l'équation (e3.26), ou si elle s'en approche au maximum sous les contraintes imposées à la minimisation et pour une distance à déterminer. Cependant, cette approche naturelle ¹¹ présuppose le choix d'un critère de proximité entre les deux diagonales, si l'égalisation stricte n'est pas réalisable pour le plan de sondage considéré.

Comme exemple pour illustrer la portée de la proximité des diagonales, les deux formes quadratiques de même diagonale $Q = \begin{pmatrix} a & c \\ c & a \end{pmatrix}$ et $\tilde{Q} = \begin{pmatrix} a & -c \\ -c & a \end{pmatrix}$ sont positives si $a \geq 0$ et $c^2 \leq a^2$ (Annexe A). Alors $\sup_{\|u\| \leq 1} |Q - \tilde{Q}|(u) = 2|c|$ peut être arbitrairement grand. Néanmoins, l'écart reste limité relativement à $\sup_{\|u\| \leq 1} |Q|(u) = a + |c| \geq 2|c|$.

– Si seuls sont pris en compte les éléments diagonaux des deux formes quadratiques Q et Q^{HT} , l'écart ainsi tronqué entre les deux estimateurs pour une variable z vaut :

$$\sum_{i \in s} \delta_i(c) z_i^2 - \sum_{i \in s} (1 - \pi_i) z_i^2 = \sum_s a z^2 \quad (\text{e3.29})$$

(en posant $a = \delta(c) - (1 - \pi)$)

→ La pondération 'calée sur la trace' $c_2 = (1 - \pi) \left(\sum_s 1 - \pi \right) / \sum_s \delta(1 - \pi)$ présente trois avantages théoriques :

- Elle minimise l'écart précédent en valeur absolue sur les z tel que z^2 est constante.
- Elle vérifie la propriété attrayante que si $\pi_i = 1$ alors $\delta_i(c_2) = 0 = Q_{i,i}^{HT}$.
- $c_2 \geq 1 - \pi$ car $\delta(1 - \pi) \leq 1 - \pi$ donc le facteur de calage est ≥ 1 .

→ La pondération c_2 exploite une seule information sur la diagonale de l'estimateur sans biais. Comme pour le calage sur marge, exploiter l'information multidimensionnelle connue pourrait rendre plus précise l'estimation de variance.

– Pour minimiser en c le critère $\max_{\sum z^2 \leq 1} \left(\sum a z^2 \right)^2$, l'optimum est :

$$\operatorname{argmin}_c \max_s a^2 = \operatorname{argmin}_c \max_s (\delta(c) - (1 - \pi))^2$$

Cette option présente l'intérêt d'assurer la qualité minimale la meilleure possible (selon le critère 'tronqué' retenu), dans la classe des approximations retenues, pour toutes les variables d'intérêt. Mais elle n'a pas été testée ici, parce qu'elle n'est pas implémentable par les algorithmes programmés pour cette étude.

– Un dernier critère de proximité vise à minimiser la moyenne équipondérée de l'écart (e3.29) au carré sur toutes les variables dont les coordonnées sont comprises entre -1 et 1 :

$$\begin{aligned} \int_{-1}^1 \left(\sum a z^2 \right)^2 dz \frac{1}{2^{|s|}} &= \frac{1}{2^{|s|}} \int_{-1}^1 \left(\sum (a z^2)^2 + \sum_{i \neq j} a_i z_i^2 a_j z_j^2 \right) dz = \sum a^2 \frac{1}{5} + \sum_{i \neq j} a_i a_j \frac{1}{3} \\ &= \frac{1}{45} \left(4 \sum a^2 + 5 \left(\sum a \right)^2 \right) \end{aligned} \quad (\text{e3.30})$$

(L'intégrale ci-dessus est multiple, puisque $z \in [-1, 1]_s$.)

→ Cette troisième approche suggère le critère à minimiser $\sum (\delta(c) - (1 - \pi))^2$, qui est le carré de la norme euclidienne de l'écart entre les deux diagonales ¹². Dans la suite, cette fonction objectif est notée $\tau(c) = \|\delta(c) - (1 - \pi)\|^2$ et $f(c) = \delta(c) - (1 - \pi)$. Avec ce choix, une pondération optimale est une solution du programme $\min_c \|\delta(c) - (1 - \pi)\|^2$.

11. Le rapprochement est discutable si la forme quadratique définie par l'estimateur sans biais n'est pas positive.

12. En fait, le critère a été retenu avant l'analyse ci-dessus. Le critère de proximité alternatif peut se formuler : $\alpha_0 \|\delta(c) - (1 - \pi)\|^2 + \beta_0 \left\{ \sum \delta(c) - \sum (1 - \pi) \right\}^2 = \|A(\delta(c) - (1 - \pi))\|^2$ avec $A = \alpha \operatorname{id} + \beta 1_s 1'_s$.

→ La concavité de $\varphi(c)$ entraîne que pour tout $i \in s$, $1 - \pi_i - \delta_i(c)$ est convexe en c . Comme $x \mapsto x^2$ est convexe et croissante sur \mathbb{R}^+ , pour tout $i \in s$, $\{\delta_i(c) - (1 - \pi_i)\}^2$ est convexe en c sur un voisinage convexe de $1 - \pi$ tel que $\delta(c) \leq 1 - \pi$. Parce qu'une somme de fonctions convexes est convexe, la fonction objectif $\tau(c)$ est convexe en c sur un tel voisinage de $1 - \pi$.

→ dérivée partielle de la fonction objectif (pour $c > 0$) :

$$\begin{aligned} \frac{\partial \tau(c)}{\partial c_i} &= \frac{\partial \|\delta(c) - (1 - \pi)\|^2}{\partial c_i} = \sum_{\#i} 2(\delta(c) - (1 - \pi)) \frac{\partial \delta(c)}{\partial c_i} + 2(\delta_i(c) - (1 - \pi_i)) \frac{\partial \delta_i(c)}{\partial c_i} \\ &= \sum_{\#i} 2(\delta(c) - (1 - \pi)) (u'_i (\sum c u u')^- c u)^2 + 2(\delta_i(c) - (1 - \pi_i)) (1 - u'_i (\sum c u u')^- c_i u_i)^2 \end{aligned} \quad (\text{e3.31})$$

$$= \sum_{\#i} 2(\delta(c) - (1 - \pi)) (u'_i (\sum c u u')^- c u)^2 + 2(\delta_i(c) - (1 - \pi_i)) \left[\frac{\delta_i(c)}{c_i} \right]^2 \quad (\text{e3.32})$$

$$\Rightarrow \frac{\partial \tau(1 - \pi)}{\partial c_i} \leq 0 : C \text{ est cohérent avec } \delta(1 - \pi) \leq 1 - \pi.$$

– Si $c_i = 0, c_s \neq 0$ et $u_i \in \text{Im}(u^{c>0})$ alors $\frac{\partial \delta_i(c)}{\partial c_i} = 1$ (vu (e2.25) et $\widehat{\beta}_c(\mathbb{1}_i) = 0$ dans ce cas).

Mais même ces conditions ne suffisent sans doute pas à assurer que $c_i > 0$ pour une solution minimisante, parce qu'il semble possible que $\sum (\delta(c) - (1 - \pi))_{\#i} [u'_i \widehat{\beta}_{\#i}]^2 > 1 - \pi_i$ pour $c_i = 0$ et $u_i \in \text{Im}(u^{c>0})$.

• La logique de construction de la pondération par rapprochement de l'estimateur sans biais n'est pas la seule envisageable. La pondération de Hajek, ou corrigée de la perte de degrés de liberté, est formulée ainsi :

$$c_h = (1 - \pi) |s| / (|s| - \text{rang}(x_s))$$

Le coefficient correcteur de la perte de degrés de liberté peut s'interpréter par un modèle tel que conditionnellement à $\tilde{u}_s = \sqrt{\frac{\pi^*(1 - \pi^*)}{\pi}} \frac{x}{\pi}(s)$, $\tilde{z}_s = \sqrt{\frac{\pi^*(1 - \pi^*)}{\pi}} \frac{y}{\pi}(s)$ est d'espérance $\tilde{u}_s \beta^*$ et de variance diagonale $\sigma^2 \text{id}_s$. Dans ce cadre : $E \left\{ \sum_s (\tilde{z}_s - \tilde{u}_s \beta^*)^2 \middle| \tilde{u}_s \right\} = \sigma^2 |s|$. Alors avec

$$\text{l'estimateur : } \widehat{\beta}^s = \left\{ \sum_s \frac{\pi^*(1 - \pi^*)}{\pi} \frac{x x'}{\pi^2} \right\}^- \sum_s \frac{\pi^*(1 - \pi^*)}{\pi} \frac{y x}{\pi^2},$$

et $\widehat{\sigma}^2 = \sum_s (\tilde{z}_s - \tilde{u}_s \widehat{\beta}^s)^2 = \|\text{proj}_{\text{Im}(\tilde{u}_s)}^\perp(\tilde{z}_s)\|^2$, on a : $E \left\{ \widehat{\sigma}^2 \middle| \tilde{u}_s \right\} = \sigma^2 (|s| - \text{rang}(x_s))$. Donc dans ce

modèle $\sum_s \frac{\pi^*(1 - \pi^*)}{\pi} \left(\frac{y}{\pi} - \frac{x}{\pi} \beta^* \right)^2$ peut s'estimer par : $\frac{|s|}{|s| - \text{rang}(x_s)} \sum_s \frac{\pi^*(1 - \pi^*)}{\pi} \left(\frac{y}{\pi} - \frac{x}{\pi} \widehat{\beta}^s \right)^2$.

3.2 bornes de la pondération c

– Si c est considéré comme une approximation de $\frac{\pi^*(1 - \pi^*)}{\pi}$ alors il conviendrait d'imposer $0 < c \leq \frac{1}{4\pi}$ ¹³. Pour les algorithmes à gradient, cette borne supérieure a permis de limiter les difficultés informatiques d'inversion de matrice ¹⁴. Par ailleurs, ces contraintes de calcul numérique ont amené à annuler les poids c_i des unités telles que $\delta_i(1) \leq 10^{-6}$.

¹³. Ce majorant peut être multiplié par $|s| / (|s| - \text{rang}(x_s))$ pour une correction pour perte de degrés de liberté.

¹⁴. La fonction *ginv* de SAS/IML échoue lorsque certains coefficients sont trop grands, par exemple pour $(10^{18} \ 10^{18} / 0 \ 10^{18})$. Cette erreur stoppe l'exécution du module IML en cours. L'utilisation des instructions IML *call push/resume/pause* permet de gérer cette erreur. Un calcul alternatif de l'inverse a été programmé pour contourner ce problème, ainsi que pour améliorer la qualité de l'inverse, pour une matrice symétrique. Mais il n'a pas été retenu car plus lent que *ginv* et surtout dans certains cas de valeurs propres plus instables d'une itération d'un calcul d'une même inverse à l'autre.

- L'équation (e2.11) fournit un encadrement de $\delta(c)$:

$$0 \leq \delta_i(c) = \min_{\beta} \|\mathbf{1}_i - u_s \beta\|_c^2 \leq \|\mathbf{1}_i\|_c^2 = c_i \text{ d'où :}$$

$$0 \leq \delta(c) \leq c \quad (\text{e3.33})$$

- $c_i = 0 \Rightarrow \delta_i(c) = 0$. Donc une solution exacte de l'équation de pondération vérifie nécessairement $c_i > 0$, pour i tel que $\pi_i < 1$.

- $\delta_i(c) = 1 - \pi_i \Rightarrow c_i \geq 1 - \pi_i$ donc si c est une solution exacte alors $c \geq 1 - \pi$. Ceci suggère de borner la minimisation à $\{c \geq 1 - \pi\}$.

- Le cas où $\delta_i(c) = c_i$ est exclu si $c u(i) \neq 0$:

$$\delta_i(c) = c_i \Leftrightarrow 0 \in \operatorname{argmin}_{\beta} \|\mathbf{1}_i - u_s \beta\|_c^2 \Leftrightarrow \langle \mathbf{1}_i - u'0, u_s \rangle_c = 0 \Leftrightarrow c_i u_i = 0$$

$$\rightarrow \text{si } c_i u_i \neq 0 \text{ alors } \delta_i(c) < c_i \quad (\text{e3.34})$$

- Cette propriété implique que si $c > 0$ est une solution de $\delta(c) = 1 - \pi$ et si le vecteur de calage ne s'annule pas sur s alors $\forall i \in s / \pi_i < 1, c_i > 1 - \pi_i$. En effet, comme $\delta_i(c) \leq c_i$ une solution vérifie nécessairement $c_i \geq 1 - \pi_i$ et de plus $c_i = 1 - \pi_i \Rightarrow \delta_i(c) = c_i \Rightarrow c_i u_i = 0$.

- La croissance de la fonction δ et son homogénéité de degré 1 implique l'encadrement :

$$\min(c) \delta(1) \leq \delta(c) \leq \max(c) \delta(1) \quad (\text{e3.35})$$

- Soit c une solution exacte, donc telle que $\delta(c) = 1 - \pi$. Si $\delta(1) > 0$, cet encadrement implique que :

$$\min(c) \leq \min \left[\frac{1 - \pi}{\delta(1)} \right] \text{ et } \max(c) \geq \max \left[\frac{1 - \pi}{\delta(1)} \right]$$

- Cette dernière inégalité prouve que la pondération maximale d'une solution (e3.26) est d'autant plus grande que le minimum de $\delta(1)$ sur $\{\pi < 1\}$ est petit.

- Si $u_i \neq 0$ alors la croissance du i -ième coefficient diagonal en fonction du i ème poids est bornée, à autres poids fixés et finis :

$$\lim_{c_i \rightarrow +\infty} \delta_i(c) = \min_{u'_i \beta = 1} \sum_{\neq i} c(u_s \beta)^2 = \min_{u'_i \beta = 1} \|u_{\neq i} \beta\|_{c_{\neq i}}^2 \quad 15 \quad (\text{e3.36})$$

- ↳ $\delta_i(c)$ est fonction croissante de c_i , et ce coefficient est majoré par $\min_{u'_i \beta = 1} \sum_{\neq i} c(u_s \beta)^2$.

- ↳ Pour tout $\eta > 0$, $\delta_i(c) \geq \eta \wedge \min_{|1 - u'_i \beta| \leq \sqrt{\frac{\eta}{c_i}}} \sum_{\neq i} c(u_s \beta)^2$ (minimum des deux termes).

En effet, pour tout β , ou bien $(1 - u'_i \beta)^2 \geq \frac{\eta}{c_i}$ et alors $c_i (1 - u'_i \beta)^2 + \sum_{\neq i} c(u_s \beta)^2 \geq \eta$, ou bien

$(1 - u'_i \beta)^2 < \frac{\eta}{c_i}$ et alors $\sum_{\neq i} c(u_s \beta)^2 \geq \min_{|1 - u'_i \beta| \leq \sqrt{\frac{\eta}{c_i}}} \sum_{\neq i} c(u_s \beta)^2$. La suite suppose le choix $\eta < c_i$.

$$\hookrightarrow \min_{|1 - u'_i \beta| \leq \sqrt{\frac{\eta}{c_i}}} \sum_{\neq i} c(u_s \beta)^2 = \min_{\begin{cases} u'_i \beta = 1 \\ |1 - \alpha| \leq \sqrt{\frac{\eta}{c_i}} \end{cases}} \sum_{\neq i} c(u_s \beta)^2 \alpha^2 = \min_{u'_i \beta = 1} \sum_{\neq i} c(u_s \beta)^2 \left[1 - \sqrt{\frac{\eta}{c_i}} \right]^2$$

- pour $j \neq i$:

$$\lim_{c_i \rightarrow +\infty} \delta_j(c) = c_j \left[1 - c_j u'_j \operatorname{proj}_{u_i}^+ \left[\sum_{\neq i} c \operatorname{proj}_{u_i}^+ u u' \operatorname{proj}_{u_i}^+ \right]^+ \operatorname{proj}_{u_i}^+ u_j \right] \quad (\text{e3.37})$$

- ↳ On peut écrire $u_s = \lambda_s u'_i + v_s$ avec $v_s u_i = 0$ (donc $v_i = 0$ et $\lambda_i = 1$).

- ↳ Comme $\beta = \alpha u_i / \|u_i\|^2 + h$ avec $h \perp u_i$:

$$\begin{aligned} \delta_j(c) &= \min_{\beta} \|\mathbf{1}_j - u_s \beta\|_c^2 = \min_{\alpha \in \mathbb{R}, h \perp u_i} \|\mathbf{1}_j - \alpha \lambda_s - v_s h\|_c^2 \\ &= \|\mathbf{1}_j - \hat{\alpha} \lambda_s - v_s \hat{h}\|_c^2 = \|\mathbf{1}_j - \hat{\alpha} \lambda_s - v_s \hat{h}\|_{c_{\neq i}}^2 + c_i \hat{\alpha}^2 \end{aligned}$$

15. $\sum_{\neq i}$ signifie $\sum_{j \in s, j \neq i}$.

↳ Donc $c_i \widehat{\alpha}^2 \leq \delta_j(c) \leq c_j < +\infty$ et par suite $\lim_{c_i \rightarrow \infty} \widehat{\alpha} = 0$.

↳ Soit $\tilde{h} \in \operatorname{argmin}_{h \perp u_i} \|\mathbb{1}_j - v_s h\|_{c_{\neq i}}$

↳ $\sqrt{\delta_j(c)} \leq \|\mathbb{1}_j - v_s \tilde{h}\|_{c_{\neq i}}$ (vu en prenant $\beta = \tilde{h}$)

↳ $\sqrt{\delta_j(c)} \geq \|\mathbb{1}_j - \lambda_s \widehat{\alpha} - v_s \tilde{h}\|_{c_{\neq i}} \geq \|\mathbb{1}_j - v_s \tilde{h}\|_{c_{\neq i}} - \|\lambda_s \widehat{\alpha}\|_{c_{\neq i}} \geq \|\mathbb{1}_j - v_s \tilde{h}\|_{c_{\neq i}} - \|\lambda_s\|_{c_{\neq i}} |\widehat{\alpha}|$

↳ Comme le dernier terme ci-dessus tend vers 0, $\lim_{c_i \rightarrow \infty} \sqrt{\delta_j(c)} = \|\mathbb{1}_j - v_s \tilde{h}\|_{c_{\neq i}}$ \square

$\Rightarrow \lim_{c_i \rightarrow \infty} \tau = \left[\min_{u'_i \beta = 1} \sum_{\neq i} c(u_s \beta)^2 - (1 - \pi_i) \right]^2 + \sum_{\neq i} \left[c \left(1 - cv' \left[\sum_{\neq i} cvv' \right]^+ v \right) - (1 - \pi) \right]^2$: la

fonction objectif tend vers une limite finie lorsqu'un seul coefficient de pondération tend vers l'infini.

Plus généralement :

lemme I3.1 limite à l'infini de δ

si $\delta(1) > 0$ ¹⁶, $1 - \pi_s > 0$, u_s ne s'annule pas, $c^{\mathbb{N}}$ est une séquence de pondérations convergente vers $c^\infty \in]0, +\infty]_s$ et $\delta_s(c^{\mathbb{N}})$ est bornée dans \mathbb{R}_s alors, avec $s_\infty = \{i \in s / c_i^\infty = +\infty\}$, $\Delta s = s \setminus s_\infty$ son complémentaire et $s_\infty^i = s_\infty \setminus \{i\}$:

▷ $\Delta s \neq \emptyset$ (ie au moins un poids fini)

$$\triangleright \forall i \in s_\infty, \lim \delta_i(c^{\mathbb{N}}) = \min_{\substack{u'_i \beta = 1 \\ u_{s_\infty^i} \beta = 0}} \sum_{\Delta s} c^\infty(u' \beta)^2 = \min_{\substack{u'_i \beta = 1 \\ u_{s_\infty^i} \beta = 0}} \|u_{\Delta s} \beta\|_{c_\infty}^2 \quad (\mathbf{e3.38})$$

$$\triangleright \forall j \in \Delta s, \lim \delta_j(c^{\mathbb{N}}) = \min_{\beta} \|\mathbb{1}_j - v_{\Delta s} \beta\|_{c_{\Delta s}^\infty}^2, \text{ où } v^{\Delta s} = \operatorname{proj}_{\operatorname{Im}(u^{s_\infty})}^\perp(u^{\Delta s}) \quad (\mathbf{e3.39})$$

↳ Si $\Delta s = \emptyset$ alors $\delta_s(1 - \pi) > 0$ implique que $\lim_{\lambda \rightarrow \infty} \lambda \delta(1 - \pi) = (+\infty)_s$, et donc $\lim \delta_s(c^{\mathbb{N}}) =$

$(+\infty)_s$, ce qui contredit l'hypothèse que $\delta_s(c^{\mathbb{N}})$ est bornée.

↳ Pour $i \in s_\infty$ $\delta_i(c^n) \leq \min_{\substack{u'_i \beta = 1 \\ u_{s_\infty^i} \beta = 0}} \sum_{\Delta s} c^n(u' \beta)^2$. Comme $c_{\Delta s}^{\mathbb{N}} \rightarrow c_{\Delta s}^\infty$ et $0 < c_{\Delta s}^\infty < +\infty$ par

hypothèse, pour tout $0 < \epsilon < 1$ et n assez grand $c_{\Delta s}^\infty(1 - \epsilon) \leq c_{\Delta s}^n \leq c_{\Delta s}^\infty(1 + \epsilon)$, d'où on déduit que : $\overline{\lim} \delta_i(c^{\mathbb{N}}) \leq \min_{\substack{u'_i \beta = 1 \\ u_{s_\infty^i} \beta = 0}} \sum_{\Delta s} c^\infty(u' \beta)^2 (< +\infty)$.

↳ On peut décomposer : $\widehat{\beta}^n = \widehat{\beta}_{c^n}(\mathbb{1}_i) = u^{s_\infty} \widehat{\gamma}^n + \Delta \widehat{\beta}^n$ avec $\Delta \widehat{\beta}^n \perp \operatorname{Im}(u^{s_\infty})$.

$$\overline{\lim} c_i^{\mathbb{N}} (1 - u'_i u^{s_\infty} \widehat{\gamma}^{\mathbb{N}})^2 + \sum_{s_\infty^i} c^{\mathbb{N}} (u' u^{s_\infty} \widehat{\gamma}^{\mathbb{N}})^2 \leq \overline{\lim} \delta_i(c^{\mathbb{N}}) < +\infty \Rightarrow \lim u_{s_\infty} u^{s_\infty} \widehat{\gamma}^{\mathbb{N}} = \begin{pmatrix} 1 \\ 0_{s_\infty^i} \end{pmatrix}$$

↳ minoration :

$$\delta_i(c^n) = \min_{\Delta \beta \in \operatorname{Im}(u^{s_\infty})^\perp} \sum_s c^n \{ \mathbb{1}_i - u' (u^{s_\infty} \widehat{\gamma}^n + \Delta \beta) \}^2 \geq \min_{\Delta \beta \in \operatorname{Im}(u^{s_\infty})^\perp} \sum_{\Delta s} c^n [u' (u^{s_\infty} \widehat{\gamma}^n + \Delta \beta)]^2$$

Ce minorant est manifestement une fonction continue en $(c_{\Delta s}^n, u^{s_\infty} \widehat{\gamma}^n)$.

↳ $\lim u_{s_\infty} u^{s_\infty} \widehat{\gamma}^{\mathbb{N}} = \begin{pmatrix} 1 \\ 0_{s_\infty^i} \end{pmatrix} \Rightarrow \begin{pmatrix} 1 \\ 0_{s_\infty^i} \end{pmatrix} \in \operatorname{Im}(u_{s_\infty})$ (car ce sous-espace vectoriel est fermé),

$\Rightarrow \exists \gamma^\infty / u_{s_\infty} u^{s_\infty} \gamma^\infty = \begin{pmatrix} 1 \\ 0_{s_\infty^i} \end{pmatrix}$ Comme u_{s_∞} est linéaire et injective sur $\operatorname{Im}(u^{s_\infty})$ (de dimension finie), $\lim u_{s_\infty} u^{s_\infty} \widehat{\gamma}^{\mathbb{N}} = u_{s_\infty} u^{s_\infty} \gamma^\infty \Rightarrow \lim u^{s_\infty} \widehat{\gamma}^{\mathbb{N}} = u^{s_\infty} \gamma^\infty$.

16. La démonstration fonctionne également pour $\tilde{s} \subset s$ tel que $\delta_{s \setminus \tilde{s}}(1) = 0$, $\delta_{\tilde{s}}(1) > 0$ et $1 - \pi_{\tilde{s}} > 0$.

$$\hookrightarrow \text{d'où } \underline{\lim} \delta_i(c^{\mathbb{N}}) \geq \min_{\Delta\beta \in \text{Im}(u^{s_\infty})^\perp} \sum_{\Delta s} c^\infty [u'(u^{s_\infty} \gamma^\infty + \Delta\beta)]^2 \geq \min_{u_s \beta = \begin{pmatrix} 1 \\ 0_{s_\infty} \end{pmatrix}} \sum_{\Delta s} c^\infty (u' \beta)^2$$

$\Rightarrow \delta_i(c^{\mathbb{N}})$ est convergente vers cette limite.

\hookrightarrow La propriété (e3.39) se démontre comme (e3.37), en écrivant $u_s = \Lambda u_{s_\infty} + v_s / v_s u^{s_\infty} = 0$ et $\beta = \alpha + \gamma$ avec $\alpha \in \text{Im}(u^{s_\infty})$, $\gamma \in \text{Im}(u^{s_\infty})^\perp$. \square

\rightarrow En général, ces éléments théoriques n'assurent pas que la minimisation ne fasse pas tendre des coefficients de la pondération vers l'infini (voir une condition suffisante pour la pondération 'réursive').

3.3 existence de solutions

– Le lemme 13.2 exclut l'existence d'une solution exacte de (e3.26) si l'indicatrice d'une des unités échantillonnées est dans l'image des vecteurs d'équilibrage rapportés à la probabilité d'inclusion. Ceci correspond à une unité atypique pour le vecteur d'équilibrage.

lemme 13.2 condition de constance de δ_i

$$\delta_i(1) = 0 \tag{e3.40}$$

$$\Leftrightarrow \mathbb{1}_i \in \text{Im}(u_s) \Leftrightarrow \exists c > 0 / \delta_i(c) = 0 \Leftrightarrow \forall c > 0, \delta_i(c) = 0$$

\hookrightarrow Ce lemme découle immédiatement du fait que pour $c > 0$ $\delta_i(c)$ est le carré d'une distance entre $\mathbb{1}_i$ et $\text{Im}(u_s)$, selon (e2.11).

– Le cas $\delta_i(1) = 0$ est exclu lorsque $u_s = 1_s$ et $|s| > 1$, vu qu'alors $\delta_i(1) = 1 - 1/\sum_s 1 = 1 - 1/|s| > 0$. Il n'apparaît donc pas dans les contextes d'une seule contrainte d'équilibrage pour un échantillonnage de taille fixe.

– La condition $\exists i \in s / \mathbb{1}_i \in \text{Im}(u_s)$ ne semble pas très forte, comme $\mathbb{1}_s$ est une base de \mathbb{R}_s dont $\text{Im}(u_s)$ est un sous-espace vectoriel. Elle peut être précisée ainsi :

$$\exists \beta / u_s \beta = \mathbb{1}_i \Leftrightarrow \begin{cases} u_{\neq i} \beta = 0 \\ u'_i \beta = 1 \end{cases} \Leftrightarrow \begin{cases} \beta' u^{\neq i} = 0 \\ u'_i \beta \neq 0 \end{cases} \Leftrightarrow \begin{cases} \beta \in \text{Im}(u^{\neq i})^\perp \\ u'_i \beta \neq 0 \end{cases} \Leftrightarrow u_i \notin \text{Im}(u^{\neq i}) \tag{e3.41}$$

En effet $u_i \notin \text{Im}(u^{\neq i}) \Leftrightarrow u_i \notin \text{Im}(u^{\neq i})^{\perp\perp} \Leftrightarrow \exists \beta \in \text{Im}(u^{\neq i})^\perp / u'_i \beta \neq 0$. Donc :

$$\mathbb{1}_i \in \text{Im}(u_s) \Leftrightarrow u_i \notin \text{Im}(u^{\neq i}) \tag{e3.42}$$

Cette condition peut s'exprimer ainsi : x_i n'est pas une combinaison linéaire des vecteurs d'équilibrage du reste de l'échantillon. En effet : $\text{Im}(u^{\neq i}) = \text{Im}(x^{\neq i}) \text{diag} \left[\left(\frac{1}{\pi} \right)^{\neq i} \right] = \text{Im}(x^{\neq i})$.

\rightarrow Il y a donc un lien via ces unités atypiques entre la difficulté d'optimiser la pondération et celle de l'équilibrage : si $\begin{cases} x_i \notin \text{Im}(x^{\neq i}) \\ \pi_i < 1 \end{cases}$ alors il n'est pas possible d'équilibrer exactement

le plan de sondage, puisque pour un échantillon $s \neq i$, $\sum_s \frac{x}{\pi} \neq \sum_{\mathcal{P}} x$.

– Les unités concernées par (e3.42) n'interviennent pas dans les coefficients diagonaux des autres unités :

$$\mathbb{1}_i \in \text{Im}(u_s) \Rightarrow \forall j \neq i, \langle \mathbb{1}_j - u_s \widehat{\beta}_j, \mathbb{1}_i \rangle = 0 \Rightarrow u'_i \widehat{\beta}_c(\mathbb{1}_{\neq i}) = 0 \text{ (si } c > 0) \tag{e3.43}$$

$$u_i \notin \text{Im}(u^{\neq i}) \Rightarrow \forall j \neq i, u'_j \left(\sum_c c u u' \right)^- u_j = u'_j \left[\sum_{\neq i} c u u' \right]^+ u_j \tag{e3.44}$$

La deuxième propriété découle de ce que pour tout β , il existe γ tel que $u_{\neq i} \gamma = 0$ et $u'_i \gamma = -u'_i \beta$, vu (e3.41) ¹⁷. Il s'ensuit que $\forall j \neq i, \forall \beta, \|\mathbb{1}_j - u_s \beta\|_c^2 \geq \|\mathbb{1}_j - u_s \beta - u_s \gamma\|_c^2 = \|\mathbb{1}_j - u_s \beta\|_{c_{\neq i}}^2$.

17. deux autres preuves : $\mathbb{1}_i \in \text{Im}(u_s) \Rightarrow \|\text{proj}_{\text{Im}(u_s)}^{\perp c} [\mathbb{1}_i]\|_c^2 = c_i (1 - u'_i \widehat{\beta}_i)^2 + \sum_{\neq i} c [u' \widehat{\beta}_i]^2 = 0$ d'une part et $\Rightarrow \forall j \neq i, \langle \mathbb{1}_i, \text{proj}_{\text{Im}(u_s)}^{\perp c} [\mathbb{1}_j] \rangle_c = 0$ d'autre part.

– La condition $\exists i \in s / \delta_i(1) = 0$ (**e3.40**) est vérifiée pour une strate du tirage de première phase de l'enquête Emploi de taille d'échantillon égale à la dimension d'équilibrage, cas de la Corse ¹⁸, ou lorsque des variables d'équilibrage sont très concentrées. C'est le cas par exemple pour l'indicatrice de commune rurale en Ile-de-France. Ces deux régions sont les seules où le nombre d'unités $|\{\delta(1) \leq 10^{-6}\}|$ dépasse 3 pour au moins une des 10 000 simulations. Pour la Corse, toutes les unités sont dans ce cas dans 80% des simulations.

Un troisième cas de coefficients diagonaux quasiment nuls est observé pour le plan de sondage étudié, pour des régions où deux unités $i \neq j$ représentent la totalité d'une variable d'équilibrage sur l'échantillon. En effet, l'équilibrage du complémentaire de l'échantillon implique qu'une sous-matrice de u_s est du type :

$$u_{i,j}^{k,l} = \begin{pmatrix} \frac{\alpha}{\pi_j} & \frac{\alpha}{1 - \pi_i} \\ \frac{\beta}{\pi_j} & \frac{\beta}{1 - \pi_j} \end{pmatrix} \text{ et } u_{\{i,j\}^c}^{k,l} = 0$$

qui est de déterminant non nul si $\pi_i \neq \pi_j$ et $\alpha\beta \neq 0$. Alors $\mathbf{1}_i$ et $\mathbf{1}_j \in \text{Im}(u_s^{k,l}) \subset \text{Im}(u_s)$.

– Les coordonnées i telles que $\delta_i(1) = 0$ ne sont pas prises en compte dans l'approximation de variance ¹⁹. C'est donc un facteur de biais de celle-ci, si une unité concernée est de probabilité d'inclusion $0 < \pi_i < 1$.

– Un deuxième problème soulevé par la condition (**e3.40**) est d'inciter l'algorithme de minimisation à la hausse des pondérations d'autres unités, lorsque l'égalité à 0 est approchée. En effet pour $i \neq j$ la dérivée partielle croisée est positive (**e2.21**).

Ce phénomène de compensation est accentué lorsque la pondération est plafonnée. Il explique sans doute que les coefficients diagonaux ne soient pas tous inférieurs à 1 pour les minimisations testées ²⁰. Par contre, la pondération simple $c_1 = 1 - \pi$ vérifie cette propriété souhaitable, parce que $\delta(1 - \pi) \leq 1 - \pi$ (vu (**e3.33**)).

– Comme vu précédemment, pour une unité i telle que $\delta_i(1) = 0$:

- $\forall c > 0, \begin{cases} u'_i \widehat{\beta}_c(\mathbf{1}_i) = 1 \\ u'_{\neq i} \widehat{\beta}_c(\mathbf{1}_i) = 0 \end{cases}$
- pour $j \neq i, u'_i \widehat{\beta}_c(\mathbf{1}_j) = 0$

\Rightarrow (**e3.31**) entraîne alors que $\forall c > 0, \frac{\partial \tau(c)}{\partial c_i} = 0$. Donc pour une telle unité, si l'égalité

du coefficient à 0 est exacte, le poids peut rester constant dans une minoration.

– Les cas où $\delta_i(1)$ est très proche de 0 mais pas exactement nul induisent en pratique une erreur bloquante dans l'inversion de la matrice $\sum c u u'$ au cours de la minimisation, parce qu'ils mènent rapidement à des poids très élevés, au delà d'un million ²¹. C'est pourquoi ils ont été systématiquement écartés des minimisations par une méthode de gradient, ce qui correspond à forcer leur poids c à 0. (Ces poids ont été mis à la valeur $1 - \pi$ pour effectuer l'estimation de variance, afin d'éviter que la PROC REG ne s'interrompt pour cause de poids nuls. Ceci ne doit pas avoir d'incidence sur l'estimation de variance.) Théoriquement, la mise à l'écart de

18. Pour la Corse, $|s| = \text{rang}(x_s)$ dans la plupart des simulations de tirage. L'observation de la constance de $\delta(c)$ qui en découle, et donc de l'absence d'effet des algorithmes d'optimisation de la pondération, a mené à la découverte du lemme de constance du coefficient diagonal.

19. Pour ces cas, comme $\Delta \geq 0, \forall j \in s, |\Delta_{i,j}(c)| \leq \sqrt{\delta_i(c) \delta_j(c)} = 0$. Donc l'estimateur de variance ne comprend pas de terme en y_i . Par suite, il est impossible que l'estimateur soit sans biais $\forall y(\mathcal{P})$.

20. S'il n'y avait pas de compensation et si $\delta_i(c) > 1$ alors baisser c_i rapprocherait de $1 - \pi_i$, vu que $\frac{\partial \delta_i(c)}{\partial c_i} = (1 - u'_i (\sum c u u')^{-1} c_i u_i)^2 = [\delta_i(c) / c_i]^2 > 0$, sans modifier $\delta_{\neq i}(c)$, et il serait possible ainsi d'atteindre $1 - \pi_i$ puisque δ est continue et que pour $c_i = 0, \delta_i(c) = 0$ (vu (**e2.10**)).

$\tilde{s} = \{i \in s / \delta_i(1) \cong 0\} = s_{\delta(1)=0}$ devrait être suffisante, parce que, pour $j \in s \setminus \tilde{s}$:

$$\begin{cases} \forall i \in \tilde{s}, u_i \notin \text{Im}(u^{*i}) \\ u_j \in \text{Im}(u^{*j}) \end{cases} \Rightarrow u_j \in \text{Im}(u^{*\tilde{s},j})$$

En effet, si $u_j = u_{\tilde{s}}\lambda_{\tilde{s}} + u^{s \setminus \tilde{s},j}\beta$ et $\exists i \in \tilde{s} / \lambda_i \neq 0$ alors $u_i \in \text{Im}(u^{*i})$.

– Du fait sans doute de l’incertitude numérique, des cas (3/2000 en moyenne) sont observés où $\delta_i(1) > 10^{-6}$ mais $\delta_i(1 - \pi) = 0$. Plus surprenant, il se trouve également que $\delta_i(1 - \pi) > 10^{-6}$ mais $\delta_i(c) = 0$ pour c grand.

– Même lorsque $\delta(1)$ n’a pas de coordonnée nulle, l’existence de $c \geq 0$ tel que $\delta(c) = 1 - \pi$ peut dépendre du vecteur $(1 - \pi)_s$.

• L’ **Annexe E** montre que :

▷ Le rang de $\dot{\delta}(c)$ est constant et égal à $\dim \{\text{Vect}[h^2, h \in \text{Im}(x_s)^+]\}$.

▷ Le rang de $\dot{\delta}(c)$ est $\geq |s| - \text{rang}(x_s)$.

▷ Pour que $\dot{\delta}(c)$ soit inversible, il suffit qu’il existe une sous-matrice de la matrice d’équilibrage x_s de rang (x_s) colonnes, composée d’une première sous-matrice carrée de dimension et de rang égaux à $\text{rang}(x_s)$ et la deuxième sous-matrice telle que la suppression d’une ligne quelconque donne une matrice de rang égal à $\text{rang}(x_s)$. Cette condition suffisante peut s’interpréter approximativement comme celle d’une taille d’échantillon suffisamment grande par rapport au nombre de contraintes d’équilibrage.

– Vu l’étude de la dérivée ‘à droite’ de $\tau(c) = \|\delta(c) - (1 - \pi)\|$, il ne semble pas évident qu’une solution minimisante donne toujours des poids non nuls aux unités i telles que $\delta_i(1) > 0$.

→ Pour une solution minimisante telle que $c > 0$, la dérivée de la fonction objectif est nécessairement nulle. Le i -ème coefficient de cette dérivée peut s’exprimer ainsi :

$$2 \sum_j [\delta_j(c) - (1 - \pi_j)] \left\{ \mathbb{1}_i^j - u_i' \left(\sum c u u' \right)^- c_j u_j \right\}^2 \quad (\text{e3.45})$$

→ Cet expression suggère que si le minimum n’est pas une solution exacte et si $\delta(1) > 0$ alors le signe de $\delta(c) - (1 - \pi)$ n’est pas constant sur les coordonnées du minimum.

– Plus rigoureusement, si $\delta(c) \geq 1 - \pi$ et s’il existe i tel que $\delta_i(c) > 1 - \pi_i$ alors : $\frac{\partial \tau}{\partial c_i} \geq$

$2(\delta_i(c) - (1 - \pi_i)) \left(1 - u_i' \left(\sum c u u' \right)^- c_i u_i \right)^2 = 2(\delta_i(c) - (1 - \pi_i)) \left[\frac{\delta_i(c)}{c_i} \right]^2 > 0$. Il est alors possible de diminuer strictement la fonction objectif en réduisant c_i .

– Au final, il n’a pas été possible de déterminer des conditions suffisantes d’existence d’une solution $0 < c < +\infty$ qui minimise $\|\delta(c) - (1 - \pi)\|$. La difficulté est de prouver qu’il est possible de limiter la minimisation à un compact, autrement dit que la procédure de minimisation ne fait pas tendre des coordonnées de la pondération vers l’infini ou vers zéro ²².

4 comment optimiser la pondération des résidus d’équilibrage

4.1 algorithmes de minimisation de $\|\delta(c) - (1 - \pi)\|$

4.1.1 un algorithme récursif

– Une procédure de minimisation de mise en œuvre simple recherche un point fixe de la fonction $\xi(c) = c - \delta(c) + 1 - \pi$ ²³. Son calcul ne nécessite donc pas la dérivée de δ . Cette

21. La probabilité d’inclusion du sondage étudié est comprise entre 1.1% et 4.2%, avec une médiane de 1.6%. S’il existe $i \in s / \delta_i(1) = 10^{-6}$ et $\pi = 1\%$ alors le poids maximal d’une solution exacte dépasse 990 000.

22. Plus précisément, il suffirait que le rapport du poids maximal sur le poids minimal soit borné (si $(\bar{c} / \underline{c}) \leq (1 / \epsilon) < +\infty : \lambda(c) = \langle \delta(c), 1 - \pi \rangle / \|\delta(c)\|^2$ est définie pour $c \in [\epsilon, 1]_s$, si $\delta(1) > 0$, et $\|\delta[\lambda(c)c] - (1 - \pi)\|$ est une fonction continue sur ce compact.)

méthode est identique à celle décrite par [AMYT] page 553, à l'initialisation et au traitement en cas de non convergence près.

→ $\forall c \geq 0, \xi(c) \geq 1 - \pi$, vu que $\delta(c) \leq c$. Cette propriété est satisfaisante, parce qu'elle est vérifiée par une solution exacte (éventuelle), et utile pour étudier $\xi^{\mathbb{N}}$.

→ Si $\pi_i = 1$ alors l'initialisation de c à $1 - \pi$ donne $[\xi^{\mathbb{N}}(1 - \pi)]_i = 0$.

→ si $\delta_i(1) = 0$ alors $\xi_i(c) = c_i + 1 - \pi_i$ donc $\xi_i^n(1 - \pi) = (n + 1)(1 - \pi_i)$

▷ Il peut être préférable d'écarter ces unités de la procédure récursive, si leur pondération croissante à l'infini pose un problème numérique.

– $\xi_i(c) > c_i \Leftrightarrow \delta_i(c) < 1 - \pi_i$ donc à chaque itération, la procédure :

▷ augmente (strictement) le poids des $i / \delta_i(c) < 1 - \pi_i$

▷ diminue le poids des $i / \delta_i(c) > 1 - \pi_i$

▷ s'arrête si la nouvelle pondération ne fait pas strictement baisser la fonction objectif

lemme 14.3 condition suffisante de majoration des poids récursifs

(pour une infinité d'itérations)

$$\text{si } u_i \neq 0 \text{ et } \min_{u'_i \beta = 1} \|u_{\neq i} \beta\|_{1-\pi_{\neq i}}^2 > 1 - \pi_i \text{ alors } \begin{cases} \overline{\lim} \xi_i^{\mathbb{N}}(1 - \pi) < +\infty \\ \underline{\lim} \delta_i^{\mathbb{N}}(1 - \pi) < +\infty \end{cases} \quad (\mathbf{e4.46})$$

Posons $\delta_i^{\mathbb{N}} = \delta_i[\xi^{\mathbb{N}}(1 - \pi)]$ et $\xi_i^{\mathbb{N}} = \xi_i^{\mathbb{N}}(1 - \pi)$.

↳ Supposons que : $\overline{\lim} \xi_i^{\mathbb{N}} = +\infty$ (e4.47)

– soit $\sigma_1(\mathbb{N}) \subset \mathbb{N} / \lim \xi_i^{\sigma_1(\mathbb{N})} = +\infty$

→ $\sigma_2(n) = \sup \{k \leq \sigma_1(n) / \delta_i^k \leq 1 - \pi_i\}$ est bien définie car $\delta_i^0 = \delta_i(1 - \pi) \leq 1 - \pi_i$.

→ si $\sigma_2(n) < \sigma_1(n)$ alors ξ_i^k décroît pour k allant de $\sigma_2(n) + 1$ à $\sigma_1(n)$, parce que $\delta_i^k > 1 - \pi_i \Rightarrow \xi_i^{k+1} < \xi_i^k$, donc $\xi_i^{\sigma_2(n)+1} \geq \xi_i^{\sigma_1(n)}$ et $\xi_i^{\sigma_2(n)+1} = \xi_i^{\sigma_2(n)} + 1 - \pi_i - \delta_i[\xi^{\sigma_2(n)}] \leq \xi_i^{\sigma_2(n)} + 1 - \pi_i$

→ donc dans tous les cas $\xi_i^{\sigma_2(n)} \geq \xi_i^{\sigma_1(n)} - (1 - \pi_i)$ et donc $\lim \xi_i^{\sigma_2(\mathbb{N})} = +\infty$

→ $\delta_i[\xi^{\sigma_2(\mathbb{N})}] \geq \min_{\beta} \|\mathbb{1}_i - u_s \beta\|_{\xi_i^{\sigma_2(n)}, 1-\pi_{\neq i}}^2$ donc $\underline{\lim} \delta_i[\xi^{\sigma_2(\mathbb{N})}] \geq \min_{u'_i \beta = 1} \|u_{\neq i} \beta\|_{1-\pi_{\neq i}}^2$, en ap-

pliquant (e3.36). Par la condition (e4.46), ce minorant est strictement supérieur à $1 - \pi_i$, alors que $\delta_i[\xi^{\sigma_2(\mathbb{N})}] \leq 1 - \pi_i$, d'où une contradiction. Donc $\overline{\lim} \xi_i^{\mathbb{N}} = +\infty$ est impossible sous la condition (e4.46).

↳ $\overline{\lim} \delta_i(\xi^{\mathbb{N}}) \leq \overline{\lim} \xi_i^{\mathbb{N}} < +\infty$ □

→ La condition suffisante de majoration des poids (e4.46) peut être explicitée ainsi :

$$\forall i \in s, u_i \neq 0 \text{ et } (1 - \pi_i) u'_i \{u_{\neq i} \star (1 - \pi)_{\neq i} u_{\neq i}\}^- u_i < 1 \quad (\mathbf{e4.48})$$

Cette formulation suggère l'interprétation qu'aucune unité ne doit être prépondérante ou trop singulière selon le vecteur de calage.

→ La condition (e4.46) constitue une sorte de généralisation de celle en haut de la page 580 de [DT], qui correspond au cas où $u_s = 1_s$.

→ Pour une unité i qui vérifie (e4.46), $\exists! \bar{c}_i / \delta_i[\bar{c}_i, 1 - \pi_{\neq i}] = 1 - \pi_i$ et $\xi_i^{\mathbb{N}} \leq \bar{c}_i + 1 - \pi_i$

↳ La fonction $f(c_i) = \delta(c_i, 1 - \pi_{\neq i})$ est continue sur \mathbb{R}^{+*} et strictement croissante (parce que $\delta_i(1) > 0$), $f(1 - \pi_i) < 1 - \pi_i$ (si $(1 - \pi_i)u_i \neq 0$) et $\lim_{+\infty} f > 1 - \pi_i$ par hypothèse.

↳ La propriété $\xi_i^n \leq \bar{c}_i + 1 - \pi_i$ est vraie pour $n = 0$. Supposons la vraie pour un $n \in \mathbb{N}$. Ou bien $\xi_i^n > \bar{c}_i$ et alors $\delta_i(\xi^n) \geq \delta_i[\bar{c}_i, 1 - \pi_{\neq i}] = 1 - \pi_i$ donc $\xi_i^{n+1} = \xi_i^n - \delta_i(\xi^n) + 1 - \pi_i \leq \xi_i^n \leq \bar{c}_i + 1 - \pi_i$. Ou bien $\xi_i^n \leq \bar{c}_i$ et alors $\xi_i^{n+1} = \xi_i^n - \delta_i(\xi^n) + 1 - \pi_i \leq \xi_i^n + 1 - \pi_i \leq \bar{c}_i + 1 - \pi_i$.

◁ Cependant, la condition (e4.46) s'avère trop forte pour le sondage étudié. En effet, sur le premier échantillon simulé, elle n'est vérifiée que pour l'Ile-de-France, même après avoir écarté les unités de $\delta_i(1) \leq 10^{-6}$. Même lorsque les variables d'équilibrage du complémentaire sont omises, la condition de majoration des poids récursifs n'est vérifiée que pour 8 régions.

23. Alternativement, si l'objectif est de minimiser $\|A[\delta(c) - (1 - \pi)]\|$, avec par exemple $A = \alpha \text{id} + \beta 1_s 1'_s$, la fonction contractante de récursion pourrait être du type : $\xi(c) = c - \lambda A(\delta(c) - (1 - \pi))$, avec $\lambda > 0$.

- Si $\xi^{\mathbb{N}}$ converge vers $c^\infty \in [1 - \pi, +\infty[$ alors c^∞ est une solution exacte.

↳ Comme ξ est continue sur \mathbb{R}_s^+ , $\xi(\lim \xi^{\mathbb{N}}) = \lim \xi^{\mathbb{N}}$ donc $\xi(c^\infty) = c^\infty$ ce qui équivaut à $\delta(c^\infty) = 1 - \pi$.

- Une condition de baisse de la fonction objectif par récurrence sur ξ doit assurer :

$$\|\delta[\xi(c)] - (1 - \pi)\| = \|\delta[c - \delta(c) + 1 - \pi] - (1 - \pi)\| \leq \|\delta(c) - (1 - \pi)\| \quad (\mathbf{e4.49})$$

L'intérêt de cette propriété est que la fonction objectif baisserait en suivant la trajectoire $\xi^{\mathbb{N}}(c)$, tant qu'elle est vérifiée.

- Or :

$$\triangleright c - \xi(c) = \delta(c) - (1 - \pi)$$

$$\triangleright \delta[\xi(c)] - (1 - \pi) = \delta[\xi(c)] - \xi(c) + \xi(c)(1 - \pi) = (\delta - \text{id})(\xi(c)) - (\delta - \text{id})(c)$$

Donc la condition de baisse par récurrence sur ξ peut s'écrire comme une contraction de $\text{id} - \delta$:

$$\begin{aligned} \|\delta[\xi(c)] - (1 - \pi)\| &\leq \|\delta(c) - (1 - \pi)\| \\ &\Leftrightarrow \|(\text{id} - \delta)(\xi(c)) - (\text{id} - \delta)(c)\| \leq \|\xi(c) - c\| \end{aligned} \quad (\mathbf{e4.50})$$

- Une formulation alternative de la condition de baisse :

$$\begin{aligned} \|\delta[\xi(c)] - (1 - \pi)\|^2 &\leq \|\delta(c) - (1 - \pi)\|^2 \\ &\Leftrightarrow \|\delta[\xi(c)] - \delta(c)\|^2 \leq -2 \langle \delta[\xi(c)] - \delta(c), \delta(c) - (1 - \pi) \rangle \end{aligned} \quad (\mathbf{e4.51})$$

→ Il y a baisse de la fonction objectif entre c et $\xi(c)$ si le raisonnement suivant est valide :

si $\bullet \|\delta(c) - (1 - \pi)\|$ est assez petit pour approximer $\delta[\xi(c)] - \delta(c)$ par $\dot{\delta}(c)[1 - \pi - \delta(c)]$

- Les termes non diagonaux de cette dérivée sont négligeables pour $1 - \pi_s - \delta_s(c)$ alors :

$$\begin{aligned} &\|\delta[\xi(c)] - \delta(c)\|^2 + 2 \langle \delta[\xi(c)] - \delta(c), \delta(c) - (1 - \pi) \rangle \\ &\cong \sum_{i \in \mathcal{S}} \left\{ \frac{\partial \delta_i(c)}{\partial c_i} (1 - \pi_i - \delta_i(c)) \right\}^2 - 2 \sum_{i \in \mathcal{S}} \frac{\partial \delta_i(c)}{\partial c_i} [1 - \pi_i - \delta_i(c)]^2 \leq 0 \\ &\text{car } \frac{\partial \delta_i(c)}{\partial c_i} = \left[\frac{\delta_i(c)}{c_i} \right]^2 \leq 1 \end{aligned}$$

– La condition de baisse (**e4.50**) est vérifiée si $\|\text{id} - \dot{\delta}\| \leq 1$ ²⁴ sur $[c, c - \delta(c) + 1 - \pi]$, c'est-à-dire que cette dérivée soit contractante sur cet intervalle :

$$\|(\text{id} - \delta)(\xi(c)) - (\text{id} - \delta)(c)\| \leq \sup_{\tilde{c} \in [c, \xi(c)]} \|(\text{id} - \dot{\delta})(\tilde{c})(\xi(c) - c)\| \quad (\mathbf{e4.52})$$

$$\leq \sup_{[c, \xi(c)]} \|\text{id} - \dot{\delta}\| \|\xi(c) - c\| \quad (\mathbf{e4.53})$$

- L'Annexe E montre que $\text{id} - \dot{\delta}$ est contractante en c constante, mais pas généralement.

- Néanmoins, l'algorithme récursif possède les propriétés avantageuses suivantes :

- ▷ baisse de la fonction objectif par rapport à la pondération $1 - \pi$ (Annexe F)
- ▷ La pondération produite est minorée par $1 - \pi$, et donc ne s'annule pas.
- ▷ Sous la condition (**e4.46**) l'algorithme ne fait pas tendre de poids vers l'infini.

- Cependant :

◁ Il n'est pas démontré que l'algorithme réduise la fonction objectif indéfiniment.

◁ L'optimalité de la pondération produite n'est pas assurée, même lorsqu'il y a une solution exacte²⁵.

◁ La condition suffisante de majoration des poids s'avère trop forte pour le sondage étudié.

→ Au final, la justification de cet algorithme repose essentiellement sur son efficacité constatée empiriquement à diminuer la fonction objectif, dans le contexte d'un plan de sondage où il ne semble pas exister de solution exacte à l'équation de pondération.

24. La norme d'une application linéaire est définie en Annexe E.

25. Une tentative pour appliquer un raisonnement analogue à celui de [CDL] p.468 n'a pas abouti.

La simplicité de l'algorithme récursif, sans calcul de la dérivée, a pour contrepartie notamment une croissance lente des poids maximaux. Par suite, cette version bute le plus souvent sur le nombre maximal d'itérations. Une première solution à cet inconvénient est d'utiliser une fonction de récursion multiplicative. L'autre solution est le recours à un algorithme 'à gradient'.

4.1.2 algorithme récursif multiplicatif

Pour accélérer ou au moins prolonger la croissance de la pondération des unités de petits $\delta_i(1)$, la fonction de récursion 'multiplicative' $\tilde{\xi}(c) = (1 - \pi)/(\delta(c)/c)$ est envisageable²⁶. Effectivement, elle s'avère fournir un algorithme quatre fois plus rapide que la fonction 'additive' pour une pondération plus minimisante pour la majorité des régions.

Cependant, son paramétrage est délicat. Si un incrément maximal de la pondération par itération n'est pas imposé ou trop élevé, l'algorithme s'arrête très rapidement par interruption de la baisse de la fonction objectif. En ce cas, le performance de la minimisation est très variable selon la région, avec des baisses de la distance amples par rapport à l'algorithme précédent (>10%), mais également des augmentations (modérées). Ces résultats laissent supposer qu'il reste une marge d'amélioration sensible de la minimisation.

D'autre part, la version multiplicative de l'algorithme récursif accroît les coefficients diagonaux excessifs ($\delta(1) > 1$) si l'incrément toléré est supérieur à 10. Une interprétation est que les variations entre itérations sont trop fortes pour un ajustement fin. L'ajout d'une séquence additive ne change pas sensiblement le résultat. Si le blocage sur un point final problématique est dû au fait que l'itération suivant n'envisage qu'une agrégation de modifications dont le bilan est défavorable, une autre piste serait de n'effectuer qu'une partie des variations de poids $\xi(c) - c$. Mais les résultats obtenus dans cette direction n'ont pas été probants.

L'arrêt de l'algorithme récursif est déclenché par l'interruption de la baisse de la fonction objectif plutôt que par le nombre maximal d'itérations. De plus il n'est pas en général causé par l'impossibilité de calculer la fonction. La dérivée de la fonction objectif, calculée pour cette pondération sur le premier échantillon simulé, n'est quasiment nulle que pour l'Ile-de-France. C'est également le cas pour la version multiplicative. Ce résultat suggère que l'algorithme récursif n'aboutit pas à une solution minimisante.

4.1.3 algorithmes utilisant le gradient

Trois types de méthode de minimisation utilisant le gradient, donc uniquement une dérivée première, ont été testés. A la k -ième étape du processus un algorithme de minimisation fournit un incrément $\tilde{\delta}$ tel que $\|f(c_k + \tilde{\delta})\| < \|f(c_k)\|$.

• descente du gradient (GDA) : Comme $\|f(c + \tilde{\delta})\|^2 = \|f(c)\|^2 + 2\langle f(c), \dot{f}(c)\tilde{\delta} \rangle + o(\tilde{\delta})$, à la condition que $\dot{f}(c)'f(c)$ ne soit pas nul, il est assuré qu'il existe $\lambda > 0$ tel que :

$\|f\{c - \lambda \dot{f}(c)'f(c)\}\|^2 < \|f(c)\|^2$, et de plus ce vecteur donne la direction de baisse maximale

en c , au sens où $\frac{\dot{f}(c)'f(c)}{\|\dot{f}(c)'f(c)\|} \in \operatorname{argmax}_{\|u\|=1} \frac{\partial \|f(c + \lambda u)\|^2}{\partial \lambda} (\lambda = 0)$.

Au voisinage d'une solution, $\dot{f}(c)'f(c)$ se rapproche de 0. Il s'avère en pratique que les baisses de la fonction objectif obtenues par l'algorithme GDA sont fortes pour les premières itérations, mais s'affaiblissent rapidement avant l'arrivée à une solution aussi bonne que d'autres algorithmes de temps d'exécution comparable.

26. Cette approche correspond à la transformation logarithmique effectuée dans [CDL] avant la récursion.

- Gauss-Newton (GNA) : Si f était deux fois dérivable en c alors le développement limité de la fonction objectif serait :

$$\|f(c + \tilde{\delta})\|^2 = \|f(c)\|^2 + 2 \langle f(c), \dot{f}(c) \tilde{\delta} \rangle + \langle \dot{f}(c) \tilde{\delta}, \dot{f}(c) \tilde{\delta} \rangle + \langle f(c), \ddot{f}(c) \tilde{\delta}^2 \rangle + o(\|\tilde{\delta}\|^2)$$

La méthode GNA choisit l'incrément qui minimise la fonction convexe suivante :

$$\tilde{\delta}^* \in \underset{\tilde{\delta}}{\operatorname{argmin}} 2 \langle f(c), \dot{f}(c) \tilde{\delta} \rangle + \langle \dot{f}(c) \tilde{\delta}, \dot{f}(c) \tilde{\delta} \rangle$$

$$\Rightarrow \dot{f}(c)' \dot{f}(c) \tilde{\delta}^* = -\dot{f}(c)' f(c) \Rightarrow \tilde{\delta}^* = -\{\dot{f}(c)' \dot{f}(c)\}^{-1} \dot{f}(c)' f(c)$$

C'est bien une direction de baisse au sens large, parce que :

$$\dot{f}(c) \tilde{\delta}^* = -\operatorname{proj}_{\operatorname{Im}[\dot{f}(c)]} [f(c)] \Rightarrow \langle f(c), \dot{f}(c) \tilde{\delta} \rangle = -\|\operatorname{proj}_{\operatorname{Im}[\dot{f}(c)]} [f(c)]\|^2 \leq 0$$

Bien que cet incrément soit encore nul lorsque le gradient $\dot{f}(c)' f(c)$ s'annule, la méthode tend à gonfler l'incrément dans certaines directions où la dérivée est petite, ce qui peut accélérer la baisse de la fonction objectif. Un peu plus précisément, en posant $A = \dot{f}(c)$, soit v une base orthonormée de vecteurs propres de AA' , donc une base de $\operatorname{Im}(A)$ (eA.57). Pour tout x , $\|A'x\|^2 = \sum \lambda(v) \langle x, v \rangle^2$ alors que $\|\operatorname{proj}_{\operatorname{Im}(A)}(x)\|^2 = \sum \langle x, v \rangle^2$. Donc si $\operatorname{proj}_{\operatorname{Im}(A)}(x)$ est dans une somme d'espaces propres de $AA' = \dot{f}(c)' \dot{f}(c)$ de valeurs propres inférieures à 1 alors $\|\operatorname{proj}_{\operatorname{Im}(A)}(x)\| \geq \|A'x\|$. Or $f(c)$ est susceptible d'être dans ce cas si $\|A'f(c)\|$ est petit par rapport à $\|f(c)\|$, donc là où la performance de GDA est faible. A noter que si $A'f(c)$ n'est pas nul alors $x = \operatorname{proj}_{\operatorname{Im}(A)} [f(c)]$ non plus, parce que $\operatorname{Ker} \{\operatorname{proj}_{\operatorname{Im}(A)}\} = \operatorname{Im}(A)^\perp = \operatorname{Ker}(A')$ (eA.56).

- Levenberg-Marquardt (LMA) : Ce troisième algorithme vise à combiner l'avantage de la quasi-assurance d'une baisse de la norme par la méthode du gradient avec la plus grande rapidité de la méthode de Newton au voisinage d'une solution. La formule de l'incrément diffère de celle de l'algorithme GNA par le produit d'un facteur d'atténuation (variable) α avec la diagonale de la matrice à inverser, dans l'équation déterminant l'incrément $\tilde{\delta}$:

$$\{\dot{f}(c)' \dot{f}(c) + \alpha \operatorname{diag} \{\dot{f}(c)' \dot{f}(c)\}\} \tilde{\delta} = -\dot{f}(c)' f(c)$$

→ Parmi les algorithmes à gradient étudiés, il a semblé plus efficace de combiner une première étape de minimisation par descente du gradient d'au plus 10 itérations et tant que le taux de décroissance de la fonction objectif excède 0.1%, avec une deuxième étape de minimisation par GNA. En particulier, ce type de programme a paru plus rapide qu'un algorithme de Levenberg et Marquardt, dans la version testée pour celui-ci.

4.1.4 nécessité d'imposer des contraintes à la minimisation

Sans contraintes, les algorithmes à gradient sont rapidement bloqués par un problème numérique d'inversion de matrice. De ce fait, il s'est avéré nécessaire d'imposer un plafonnement à la pondération. Une alternative intéressante est de contrôler que $\delta < 1$. L'intérêt est que cette contrainte s'introduit naturellement compte tenu de l'objectif de rapprochement de $1 - \pi_s$, donc moins arbitraire. Mais sans algorithme qui intègre efficacement ce type de contrainte, l'étude ne l'a testée (par la version n_3 ci-après) que par rejet des incréments qui ne satisfont pas la contrainte. Ceci mène rapidement à l'arrêt de l'algorithme. La non-optimalité de cette méthode est indiquée par une dérivée ($\|\dot{f}(c)' f(c)\|$) plus éloignée de 0 que les autres méthodes à gradient.

4.1.5 changement de variable pour plafonner la pondération

Pour mettre en œuvre une contrainte de plafonnement de la pondération, un changement de variable est apparu nettement préférable à une restriction de l'incrément.

- Si la contrainte est $c \in]\underline{c}, \bar{c}[$, avec la restriction que $-\infty < \underline{c} < \bar{c} < +\infty$, une fonction bijective (et C^1) de \mathbb{R} vers $] \underline{c}, \bar{c}[$ est : $\varphi(t) = \underline{c} + (\bar{c} - \underline{c}) \frac{1}{1 + e^{-2t}}$ ²⁷. La dérivée du changement de variable se calcule ainsi :

$$\begin{aligned} \dot{\varphi}(t) &= (\bar{c} - \underline{c}) \left[-\frac{1}{(1 + e^{-2t})^2} \right] (-2e^{-2t}) = (\bar{c} - \underline{c}) \left[\frac{1}{1 + e^{-2t}} - \frac{1}{(1 + e^{-2t})^2} \right] 2 \\ &= 2(c - \underline{c}) - 2 \frac{(c - \underline{c})^2}{\bar{c} - \underline{c}} \end{aligned}$$

L'inverse du changement de variable vaut : $\varphi^{-1}(c) = -\frac{1}{2} \log \left(\frac{\bar{c} - c}{c - \underline{c}} \right)$.

note 1: Pour $\underline{c} = 0, \bar{c} = \frac{1}{4\pi}$ et π prenant les valeurs 0.01, 0.02, 0.03 et 0.04, le taux d'écart entre $\varphi^{-1} \varphi \varphi^{-1} \varphi(t)$ et t est inférieur à 0.05% pour t allant de -50 à 16. Au delà de cette borne, l'imprécision numérique devient plus sensible. Ces résultats suggèrent que le changement de variable ne pose pas de problèmes informatiques, à condition d'éviter les extrémités.

- sans borne supérieure pour la pondération c :

$$\varphi(t) = \underline{c} + e^t \quad \dot{\varphi}(t) = e^t = c - \underline{c} \quad \varphi^{-1}(c) = \log(c - \underline{c})$$

• Dans la suite, les résultats chiffrés sont issus de 10 000 simulations du tirage de première phase de l'enquête Emploi. Ce sondage de grappes de logements est stratifié par région et équilibré sur 32 variables, dont la moitié pour équilibrer l'échantillon complémentaire.

27. Elle peut se construire à partir de la tangente hyperbolique, vu que $\frac{1}{1 + e^{-2t}} = \frac{\text{th}(t) + 1}{2}$. Cette fonction présente l'intérêt d'être bijective de \mathbb{R} dans $]0, 1[$.

4.1.6 Comparaison de l'efficacité de la minimisation selon l'algorithme

La combinaison de premières étapes GDA avec des étapes finales par GNA paraît clairement préférable à l'algorithme de pure descente du gradient pour l'objectif opérationnel à minimiser (Tableau 1), tout en étant plus rapide :

- dérivée plus proche de 0
- sensibilité nettement moindre à l'initialisation (à une région près)
- minimisation plus poussée (distance et dérivée plus petites, pour toutes les régions)

Tableau 1 – Comparaison des performances pour la minimisation de deux algorithmes 'à gradient', sur la première simulation d'échantillonnage, par région

	région	distance finale			norme de la dérivée		taux impact initialis.	
		GDA	GDA +GNA	taux écart	GDA	GDA +GNA	GDA	GDA +GNA
11	Île-de-France	0,984	0,984	0,0	3,02E-04	1,54E-08	0,000	0,000
21	Champagne-Ardenne	1,877	1,870	-0,4	1,65E-02	2,12E-05	-0,189	0,000
22	Picardie	1,217	1,199	-1,5	2,19E-02	8,92E-07	-0,146	0,000
23	Haute-Normandie	1,782	1,770	-0,6	3,03E-02	4,82E-05	-0,183	0,000
24	Centre	1,580	1,568	-0,8	3,29E-02	4,75E-06	-0,096	0,000
25	Basse-Normandie	1,531	1,518	-0,9	3,53E-02	6,37E-06	-0,073	0,000
26	Bourgogne	2,158	2,151	-0,3	2,57E-02	1,19E-06	-0,146	-0,412
31	Nord-Pas-de-Calais	1,021	1,013	-0,8	1,43E-02	4,72E-07	-0,124	0,000
41	Lorraine	1,682	1,669	-0,8	4,18E-02	4,25E-06	-0,125	0,000
42	Alsace	1,703	1,689	-0,8	2,24E-02	2,73E-06	-0,090	0,000
43	Franche-Comté	1,858	1,841	-0,9	2,88E-02	1,46E-05	-0,204	0,000
52	Pays de la Loire	1,178	1,169	-0,7	2,64E-02	2,14E-07	-0,110	0,000
53	Bretagne	0,886	0,876	-1,2	1,21E-02	1,21E-05	-0,111	0,000
54	Poitou-Charentes	2,025	2,014	-0,5	4,04E-02	8,49E-06	-0,085	0,000
72	Aquitaine	1,798	1,787	-0,6	1,90E-02	1,18E-05	-0,081	0,000
73	Midi-Pyrénées	1,444	1,440	-0,3	1,62E-02	3,04E-06	-0,016	0,000
74	Limousin	2,156	2,151	-0,3	1,26E-02	7,74E-06	-0,090	0,000
82	Rhône-Alpes	0,877	0,872	-0,6	1,02E-02	4,21E-08	-0,185	0,000
83	Auvergne	2,303	2,296	-0,3	2,69E-02	4,15E-06	-0,113	0,000
91	Languedoc-Roussillon	0,986	0,971	-1,5	2,93E-02	1,40E-06	-0,128	0,000
93	Provence-Alpes-Côte d'Azur	0,431	0,414	-3,9	1,04E-02	3,84E-06	-0,195	0,000

Note :

– Un maximum de 200 itérations est programmé pour l'algorithme de descente du gradient et 100 pour la combinaison d'algorithmes GDA+GNA. La durée d'exécution est de 1'45" pour le premier programme contre moins de 1'05" pour le second. La lenteur relative des algorithmes à gradient testés ici pourrait être affectée par la partie du programme qui effectue une recherche itérative du premier incrément qui fait baisser la fonction objectif. Une autre possibilité est que le plafonnement de la pondération imposée à ces méthodes induit un mouvement vers une frontière du domaine contraint.

– Les deux premières colonnes fournissent la valeur finale de la fonction objectif ($\|\delta(c) - (1 - \pi)\|$).

– La norme de la dérivée est $\|\dot{f}(c)' f(c)\|$. Elle est calculée sans les coordonnées telles que $\delta_i(c) \leq 10^{-6}$, qui sont exclues des itérations de la minimisation. Mais celles-ci sont prises en compte dans la distance finale.

– Le taux d'impact de l'initialisation est calculé entre les distances finales obtenues avec les initialisations respectivement à $1 - \pi$ et $\min\left[\frac{1 - \pi}{\delta(1)}, \frac{1}{4\pi}\right]$.

L'algorithme récursif a été appliqué avec une contrainte de décroissance minimale de la fonction objectif (0.001%) et une limite de 500 itérations. Ainsi paramétré, cet algorithme n'a pas nécessité de plafonner directement la pondération ²⁸ ou d'annuler des poids pour éviter les problèmes numériques. C'est avantageux par rapport aux autres méthodes de minimisation testées. Il bénéficie également de la simplicité de sa programmation (pas de calcul du gradient $\dot{\delta}(c)$). La condition de poursuite de l'algorithme tant que la fonction objectif baisse est nécessaire aussi pour les autres algorithmes. Elle présente l'avantage de ne pas être arbitraire, contrairement à l'imposition de bornes à c ou $\delta(c)$. Cette technique s'avère empiriquement très efficace pour optimiser la pondération, au sens où elle mène plus rapidement que les autres algorithmes testés à des baisses plus amples de $\|\delta(c) - (1 - \pi)\|$. Sur les 10 000 simulations, les algorithmes récursifs minimisent davantage la fonction objectif que les algorithmes à gradient testés pour toutes les régions (Tableau 2). En contrepartie, la pondération maximale est beaucoup plus élevée.

Tableau 2 – Médiane calculée sur les 10 000 simulations, pour la fonction objectif de la minimisation et pour le poids maximal, par région

région	distance $\ \delta(c) - (1 - \pi)\ $							pondération max (c)							$\frac{1 - \pi^*}{\delta(1)}$	
	v1	v2	n1	n2	n3	ren	rem	v1	v2	n1	n2	n3	ren	ren *		rem
11	1,89	1,78	1,14	1,16	1,26	1,13	1,11	1,0	1,0	12,5	5,0	2,7	225	21	71	2,1
21	3,81	3,05	1,91	2,00	3,02	1,83	1,76	1,0	1,7	13,4	5,5	1,9	434	433	65 182	23,9
22	4,17	3,21	2,07	2,21	3,38	2,00	1,87	1,0	2,2	20,3	5,9	2,4	392	392	32 047	20,1
23	3,91	3,06	1,82	1,94	3,13	1,77	1,68	1,0	1,8	19,4	5,5	2,0	414	414	47 052	15,4
24	3,50	2,66	1,40	1,50	2,57	1,37	1,32	1,0	1,5	19,0	5,2	1,6	400	400	71 487	10,3
25	4,20	3,38	2,32	2,44	3,55	2,21	2,08	1,0	2,2	18,7	6,0	2,5	423	405	43 611	39,7
26	4,11	3,22	2,09	2,22	3,37	2,02	1,91	1,0	2,1	19,3	5,8	2,3	414	413	40 157	27,6
31	2,99	2,41	1,14	1,22	2,01	1,12	1,11	1,0	1,2	17,3	5,1	1,5	374	374	95 871	7,3
41	3,62	2,73	1,49	1,60	2,70	1,45	1,38	1,0	1,6	18,8	5,3	1,7	399	399	64 378	10,7
42	3,76	2,98	1,75	1,86	2,97	1,71	1,64	1,0	1,6	18,3	5,4	1,8	438	435	64 090	20,2
43	4,20	3,39	2,22	2,34	3,55	2,12	2,01	1,0	2,2	15,5	6,0	2,5	412	411	42 035	27,3
52	3,04	2,34	1,20	1,27	2,01	1,18	1,17	1,0	1,3	18,1	5,1	1,5	399	398	89 271	9,8
53	3,18	2,38	1,11	1,20	2,14	1,09	1,06	1,0	1,3	18,3	5,1	1,5	359	359	87 802	6,7
54	4,01	3,21	2,12	2,24	3,31	2,05	1,92	1,0	2,0	20,4	5,7	2,1	475	373	46 809	13,8
72	3,18	2,48	1,31	1,38	2,20	1,29	1,27	1,0	1,3	18,4	5,1	1,5	430	416	97 736	14,1
73	3,40	2,67	1,63	1,70	2,51	1,60	1,51	1,0	1,4	18,8	5,2	1,6	459	400	81 885	15,1
74	3,93	3,10	2,12	2,18	3,17	2,00	1,92	1,0	1,9	11,1	5,7	2,2	486	401	72 921	15,0
82	2,44	1,93	0,59	0,66	1,29	0,56	0,56	1,0	1,1	14,9	5,0	1,5	238	238	78 327	3,5
83	4,38	3,73	2,71	2,85	3,94	2,54	2,36	1,0	2,7	22,9	6,5	3,0	298	274	22 037	34,4
91	3,32	2,58	1,42	1,50	2,40	1,40	1,36	1,0	1,4	18,2	5,2	1,5	409	409	79 751	11,6
93	2,54	2,12	0,89	0,95	1,59	0,87	0,87	1,0	1,1	17,1	5,0	1,4	331	331	101 251	5,2

Note :

- Les codes des méthodes sont explicités en note de bas du [Tableau 9](#), sauf rem : récursif de fonction de récursion multiplicative. Les unités i telles que $\delta_i(c) \leq 10^{-6}$ sont prises en compte dans le calcul de distance.
- La durée d'exécution moyenne du programme récursif est de 37" pour la version additive et de 8" pour la multiplicative (comparé à 1" pour GDA+GNA).
- L'ajout d'un contrôle sur la qualité de l'inversion de matrice n'a pas d'effet sensible sur cette minimisation.
- * : hors unités de $\delta_i(1) \leq 10^{-6}$

28. De fait, le nombre maximal d'itérations \bar{n} borne la pondération récursive à $\bar{n}(1 - \pi)$. Avec 10 000 itérations sans conditions d'arrêt, la pondération maximale dépasse 7 000.

- La méthode récursive donne la distance minimale dans la grande majorité des simulations croisées avec les régions (Tableau 3), et jamais la pire.

Tableau 3 – Rang des méthodes selon la distance obtenue, par simulation et région, hors Corse

	1	2	3	4	5	6	7
v1					50	347	209 708
v2				5	88 095	122 005	
n1	2 422	989	204 020	2 645	29		
n2	85	101	2 686	207 191	42		
n3			37	30	121 888	87 753	397
ren	25 581	181 718	2 593	213			
rem	182 017	27 297	769	21	1		

La baisse de la distance $\|\delta(c) - (1 - \pi)\|$ entre les pondérations v_2 et ren n'est pas seulement observée en médiane des simulations, mais sur toutes les simulations. En fait, elle est d'au moins 15% pour toutes les simulations et toutes les régions (hors Corse).

Avec un autre critère de qualité de la minimisation, l'algorithme récursif paraît très largement préférable selon le cumul des coefficients diagonaux supérieurs à 1 (Tableau 4).

Tableau 4 – Cumul des coefficients diagonaux supérieurs à 1 sur 10 000 simulations

	v2	n1	n2	n3	ren	rem
$\sum (\delta - 1)^+$	1 856 162	630 236	720 287	1 631 884	41 338	69 444
$\max(\delta)$	3,05	2,22	2,50	3,12	1,70	2,85

C'est remarquable comme la pondération ren semble plus élevée que celle de v_2 (vu le Tableau 2), et que δ est croissante. En fait, si la pondération maximale de la méthode récursive est très supérieure à celle de v_2 et la pondération moyenne également (+18), pour la majorité des unités la pondération 'récursive' est plus basse : $c_{ren} < c_2 - 10^{-6}$ pour 16 millions de poids individuels sur les 24 millions simulés. Alors que le calage sur la trace augmente uniformément la pondération de départ, l'optimisation fournit une pondération plus modulée.

- Sur 500 itérations de l'algorithme récursif sur le premier échantillon simulé, sans condition d'arrêt, δ_i^N dépasse $1 - \pi_i$ d'au plus 6%, quelque soit la région.

- De plus, la distance $\|\delta(\xi^N) - (1 - \pi)\|$ baisse d'une itération à l'autre pour tous les individus et presque toutes les 500 itérations.

- L'évolution des poids d'une itération à l'autre n'apparaît pas erratique (Annexe F).

- Plus précisément, sur les données étudiées, l'algorithme récursif parvient en au plus 500 itérations à rapprocher $\delta_i(c)$ de $1 - \pi_i$ à un taux d'écart de moins de 0.4% pour les trois quarts des unités de l'ensemble des simulations (Tableau 5). La majorité des unités ont un coefficient ν supérieur à 0.2. Pour cette tranche, 95% des unités ont un taux d'écart de moins de 2.1%. Le fait que cette catégorie comporte également des unités de $\delta_i(c)$ quasiment nul, contrairement à leur coefficient ν , s'explique sans doute par une incertitude numérique.

Tableau 5 – Distribution sur 10 000 simulations du taux d'écart entre $\delta_i(c)$ et $1 - \pi_i$ selon le coefficient $\nu_i = \min_{u_i, \beta=1} \|u_i \beta\|_{1-\pi_i}^2 - (1 - \pi_i)$

ν	n	q1	median	q3	p90	p95	max
<-0.2	3 164 930	0,7	25,1	60,6	86,0	100,0	100
0.2-<0	724 824	0,0	0,0	0,0	2,2	5,3	100
0-<0.2	677 321	0,0	0,0	0,0	2,0	4,9	100
>=0.2	19 047 551	0,0	0,0	0,0	0,9	2,1	100
ensemble	23 614 626	0,0	0,0	0,4	6,0	42,7	100

⇒ Sur les données étudiées, l’algorithme récursif apparaît nettement préférable aux autres pondérations testées pour l’objectif de minimiser la distance choisie entre les deux diagonales, du moins si l’incertitude numérique sur cet indicateur est ignorée (voir ci-dessous). Il semble capable d’égaliser pratiquement les deux diagonales pour les unités dont le coefficient ν n’est pas trop petit, en un nombre limité d’itérations.

4.2 les coefficients diagonaux contrôlent les résidus

Les poids très élevés paraissent susceptibles d’accroître la sensibilité de l’estimation de variance à des points extrêmes, dans le sens d’une surestimation de la variance ²⁹. Mais le contrôle des coefficients diagonaux permettent de borner l’impact de la valeur individuelle d’une variable sur l’estimateur de variance. Donc les unités de poids très élevé mais de faible coefficient diagonal doivent peser peu dans la variance estimée.

C’est effectivement démontrable, sous certaines conditions :

$$\begin{aligned} \widehat{e} &= (\text{id} - \text{proj}_{\text{Im}(u_s)}^c)z \Rightarrow c\widehat{e} = c(\text{id} - \text{proj}_{\text{Im}(u_s)}^c)z = \Delta(c)z \\ \text{donc : } |c_i \widehat{e}_i| &= \left| \sum_j \Delta_{i,j}(c) z_j \right| & (\text{e4.54}) \\ &\leq \sum_j |z(j)| \sqrt{\delta_i(c) \delta_j(c)} \quad (\text{car } \Delta \geq 0 \text{ (eA.58)}) \end{aligned}$$

$$\text{soit : } |c_i \widehat{e}_i| \leq \sqrt{\delta_i(c)} \sum_s |z| \sqrt{\delta(c)} \quad (\text{e4.55})$$

Cette majoration du résidu en fonction des coefficients diagonaux implique que même lorsque le poids d’une unité c_i est très élevé, sa contribution à l’estimation de variance ($c_i \widehat{e}_i^2$) est faible si $\delta_i(c)$ est petit et que le deuxième terme de ce majorant n’est pas trop grand.

4.3 calage final

Pour les pondérations produites par minimisation, la non-convergence des coefficients diagonaux pour les unités de faible coefficient ν peut être contrecarrée par un calage final sur la trace de l’estimateur sans biais, c’est-à-dire $\sum_s 1 - \pi$.

Pour la version retenue de l’algorithme récursif, le critère de proximité des diagonales a été privilégié. C’est pourquoi ce calage final n’a été effectué dans la version récursive ren que s’il baissait la fonction objectif. Ce choix est discutable : pour cette procédure, l’arrêt après un nombre fini d’itérations a sans doute principalement pour effet de sous-pondérer certaines unités (selon leur coefficient diagonal). En moyenne sur les simulations, l’écart (négatif) entre la somme des coefficients diagonaux et celle de $1 - \pi$ dépasse 10% pour 7 régions et 20% pour une région (Auvergne).

⇒ Un calage final de la pondération d’équilibrage peut être prudent pour prémunir d’une sous-estimation de la variance.

29. L’idée est qu’augmenter le poids relatif d’une unité atypique dans une moyenne pondérée de carrés de résidus accroît celle-ci, bien que l’unité tire davantage l’estimation du coefficient β :

$$\frac{\partial \frac{M[\omega(y - x'\widehat{\beta})^2]}{M(\omega)}}{\partial \omega_i} = \frac{1}{M(\omega)} \left\{ (y_i - x_i'\widehat{\beta})^2 - \frac{M(\omega(y - x'\widehat{\beta})^2)}{M(\omega)} \right\}$$

4.4 effet du nombre de variables d'équilibrage sur la minimisation

– La réduction de la norme par l'algorithme de minimisation est sensible même lorsque les variables d'équilibrage du complémentaire de l'EEC-TH sont prises en compte dans la fonction objectif (Tableau 6), mis à part la Corse. Cependant il n'y a plus de région pour lesquelles une solution quasiment exacte est trouvée, alors que sans les variables complémentaires la norme est presque annulée pour 7 régions.

Tableau 6 – Réduction de l'objectif $\|f(c)\|$ par rapport à $\|f(1 - \pi_1)\|$ selon la prise en compte des variables d'équilibrage de l'échantillon complémentaire de l'EEC-TH

région	sans les var.complémentaires			avec les var.complémentaires			s	rang (x_s)	min [$\delta(1)$]	
	distance initiale	distance finale	taux di- min %	distance initiale	distance finale	taux di- min %			sans	avec
11	1,1670	0,9844	-15,6	1,8645	0,9846	-47,2	535	25	0,0E+00	0,0E+00
21	1,9304	0,0226	-98,8	3,8652	1,8703	-51,6	63	28	0,41	0,03
22	2,1368	0,2749	-87,1	4,0142	1,1915	-70,3	51	28	0,34	0,10
23	1,9943	0,0912	-95,4	3,9615	1,7997	-54,6	58	28	0,43	0,04
24	1,7369	0,0000	-100,0	3,5548	1,5702	-55,8	76	28	0,57	0,14
25	2,0147	0,1293	-93,6	4,0035	1,5197	-62,0	53	28	0,40	0,01
26	2,0110	0,0843	-95,8	4,0316	2,1936	-45,6	53	28	0,41	0,07
31	1,3730	0,0000	-100,0	3,0112	1,0151	-66,3	138	28	0,70	0,14
41	1,7816	0,0000	-100,0	3,6465	1,6839	-53,8	73	28	0,48	0,10
42	1,8961	0,0467	-97,5	3,8392	1,7591	-54,2	65	28	0,45	0,07
43	2,1276	0,0195	-99,1	4,2186	1,8638	-55,8	48	28	0,42	0,07
52	1,4864	0,2660	-82,1	3,1284	1,1710	-62,6	117	28	0,40	0,10
53	1,5527	0,0000	-100,0	3,2522	0,8944	-72,5	101	28	0,56	0,16
54	2,0397	0,7308	-64,2	3,9911	2,0343	-49,0	58	28	0,18	4,4E-11
72	1,5825	0,4929	-68,9	3,3243	1,7952	-46,0	111	28	0,30	0,03
73	1,6794	0,4019	-76,1	3,3315	1,4453	-56,6	90	28	0,33	0,01
74	2,0240	0,7901	-61,0	3,8823	2,1557	-44,5	59	28	0,11	0,0E+00
82	1,1507	0,0000	-100,0	2,5801	0,8654	-66,5	207	28	0,68	0,12
83	2,3102	0,6565	-71,6	4,4218	2,3521	-46,8	42	28	0,22	9,7E-10
91	1,5786	0,0000	-100,0	3,3111	0,9824	-70,3	97	28	0,63	0,15
93	1,1515	0,0000	-100,0	2,6027	0,3912	-85,0	226	28	0,65	0,31
94	3,0975	3,0975	0,0	3,0975	3,0975	0,0	10	10	0,0E+00	0,0E+00

Note :

– La distance est la norme de la fonction à annuler $\|f(c)\| = \sqrt{\sum_{i \in s} [\delta_c(i) - (1 - \pi_1(i))]^2}$, pour la première simulation du tirage de l'échantillon de première phase de l'EEC-TH.

– $\text{rang}(\sum uu') = \text{rang}\left(\sum_s \frac{xx'}{\pi^2}\right) = \text{rang}\left[\begin{pmatrix} x \\ \pi \end{pmatrix}_s\right] = \text{rang}(x_s)$, où x est le vecteur d'équilibrage. Il a été calculé via la fonction `homogen` de SAS/IML, qui fournit une base du noyau ($\text{Ker}(\sum uu')$). Le nombre total de variables d'équilibrages est de 32.

– Les deux dernières colonnes fournissent le minimum régional du coefficient diagonal pour la pondération par 1. (Les cinq coefficients très proches de 0 ne sont pas stables d'une itération du calcul par `ginv` à l'autre, tout en restant quasiment nuls.)

– L'algorithme utilisé pour ce tableau est récursif. Les unités de $\delta_i(1)$ très petit n'ont pas été exclues.

4.5 incertitude numérique

– Le [Tableau 2](#) suggère une relation entre la distance minimisée et le niveau maximal de la pondération, avec un ordre similaire $ren/n_1/n_2/n_3/v_2$. C'est particulièrement clair pour les méthodes n_1 et n_2 : l'abaissement du plafond de la pondération de $\frac{1}{4\pi}$ à 5 toutes choses égales par ailleurs induit une croissance allant de 2% à 12% de la distance obtenue en médiane. Or la qualité de l'inversion de la matrice $\sum c_{uu'}$ semble décroître en fonction de $\max(c)$, ou sans doute plutôt de la dispersion du poids d'équilibrage. Par exemple, pour la Picardie, la pondération maximale augmente en fonction de l'itération de l'algorithme GDA+GNA et la qualité (absolue) de l'inverse tend à se dégrader corrélativement ([Tableau 7](#)).

– Il a été envisagé d'introduire des contraintes qui stoppent l'algorithme lorsque la qualité de l'inverse tombe en dessous d'un seuil. Les tests effectués n'ont pas mené à retenir cette option pour les méthodes à gradient, au profit d'un plafonnement de la pondération, dont le choix est peut-être moins arbitraire.

→ Il reste cependant que l'incertitude numérique induite par les poids élevés relativise les performances de minimisation et affecte la comparaison des différentes méthodes. En outre, la relation entre l'importance de la diminution de la distance entre les diagonales et la qualité de l'estimation de variance est inconnue. Il est donc indispensable de compléter la comparaison des pondérations d'équilibrage sur des simulations de l'estimation de variance.

Tableau 7 – Dégradation de la qualité de l'inverse de $A = \sum c_{uu'}$ en fonction de l'itération de l'algorithme de minimisation, pour la Picardie (reg=22), sur la première simulation

itération	$\ f(c)\ $	$\ \dot{f}(c)' f(c)\ $	$\max(c)$	$\ A - AA^{-1}A\ $	taux diminution distance
	0,4601		10,1	1,7	
1	0,4010	0,3094	10,1	6,2	-12,8
2	0,3794	0,1743	10,1	4,0	-5,4
3	0,3775	0,1408	10,1	12,6	-0,5
4	0,3701	0,1404	10,1	7,6	-1,9
5	0,3572	0,1499	10,1	11,0	-3,5
6	0,3537	0,0454	10,1	3,4	-1,0
7	0,3507	0,0457	10,1	12,1	-0,8
8	0,3481	0,0469	10,1	6,0	-0,8
9	0,3460	0,0497	10,1	10,3	-0,6
10	0,3439	0,0529	10,2	7,1	-0,6
11	0,3395	0,0581	15,3	11,3	-1,3
12	0,3315	0,0574	17,0	137,1	-2,3
13	0,3312	0,0947	17,1	20,6	-0,1
14	0,3258	0,0956	17,1	80,1	-1,6
15	0,3218	0,1090	17,1	28,5	-1,2
16	0,3215	0,1162	17,1	53,6	-0,1
17	0,3148	0,1162	17,1	37,8	-2,1
18	0,3148	0,1308	17,1	332,4	0,0
19	0,2952	0,1308	17,1	114,7	-6,2
20	0,2817	0,1293	17,1	918,4	-4,6
21	0,2778	0,1139	17,1	445,9	-1,4
22	0,2458	0,1077	17,1	369,5	-11,5
23	0,2442	0,0212	17,1	419,6	-0,7
24	0,2441	0,0020	17,1	368,4	0,0
25	0,2441	0,0009	17,1	602,4	0,0
26	0,2441	0,0009	17,1	384,1	0,0
27	0,2441	0,0000	17,1	450,5	0,0

5 enjeu de la pondération des résidus d'équilibrage pour la variance de première phase de l'enquête emploi

Le résultat le plus sensible sur la qualité simulée de l'estimation de la variance de première phase de l'enquête Emploi est que les pondérations des résidus les plus simples, par $1 - \pi$ ou par équipondération sous-estiment plus fréquemment et plus largement la variance (Tableau 9). En effet, pour trois des quatre variables testées, ces méthodes sous-estiment l'écart quadratique moyen (calculé sur les 10 000 simulations) dans plus de la moitié des simulations. La sous-estimation par la pondération $1 - \pi$ (version v_1) est conforme aux attentes théoriques, vu que $\delta(1 - \pi) < 1 - \pi$, si $(1 - \pi)x \neq 0$ ³⁰. Il semble donc prudent de choisir une pondération plus élaborée.

L'équipondération avec calage sur la trace de l'ESB, v_a , est indistinguable de la pondération $1 - \pi$ calée pareillement (v_2). Cette faible sensibilité de l'approximation de variance ($\|z - u_s \widehat{\beta}_c(z)\|_c^2$) aux petites fluctuations de la pondération peut être rapprochée de la formule de sa dérivée en fonction de la pondération : $\|z - u_s \widehat{\beta}_{c+\Delta c}(z)\|_{c+\Delta c}^2 \cong \|z - u_s \widehat{\beta}_c(z)\|_c^2 + \|z - u_s \widehat{\beta}_c(z)\|_{\Delta c}^2$ et $\|z - u_s \widehat{\beta}_c(z)\|_{\Delta c}^2 \ll \|z - u_s \widehat{\beta}_c(z)\|_c^2$ si $|\Delta c| \ll c$.

Le classement des méthodes selon la qualité de l'estimation de variance dépend de la variable testée. Autrement dit, aucune des méthodes testées ne permet de réduire sensiblement les sous-estimations médianes sans accroître les surestimations médianes.

– Les deux dernières variables, dont l'estimation de variance est moins précise, ont un écart quadratique moyen qui se stabilise plus lentement (Tableau 8). L'ampleur de la surestimation de la dernière variance paraît reliée à la baisse sensible de l'écart quadratique moyen mesuré respectivement sur 2 000 et 10 000 simulations.

Tableau 8 – Stabilisation de la racine carrée de l'écart quadratique moyen (reqm)

	stabilisation reqm simulation	reqm 10000/1000	reqm 10000/2000
foy rev	3 691	0,8	0,9
nbpi2	1 315	-1,7	0,0
mar	7 348	2,9	0,3
isf	5 820	-3,8	-3,0

Note : La première colonne donne le nombre de simulations nécessaires pour stabiliser le reqm à moins de 1% de sa limite calculée sur les 10 000 simulations.

– L'estimation de variance de la première variable est très largement sous- ou sur-estimée dans un faible nombre de simulations, pour toutes les pondérations testées. L'ampleur de l'écart entre les indicateurs des écarts-types et des variances estimées signifie que le coefficient de variation entre les simulations de l'estimation de l'écart-type est élevé (vu que $Var_{sim}(\widehat{\sigma}) = \frac{\sum \widehat{\sigma}^2}{|sim|} - \left(\frac{\sum \widehat{\sigma}}{|sim|}\right)^2$). Cette variabilité de l'estimateur de la première variance est également observée sur l'intervalle inter-quartile, nettement plus ample que pour les autres variables. Enfin c'est la seule variable pour laquelle l'écart entre la moyenne et la médiane des estimations de variance est notable.

30. Un autre argument, sans référence à l'estimateur sans biais, est que si $\pi^* = \pi$ alors l'estimateur $\sum_s (1 - \pi) \left(\frac{y}{\pi} - \frac{x}{\pi} \widehat{\beta}\right)^2$ est toujours inférieur à l'estimateur 'théorique' $\sum_s (1 - \pi) \left(\frac{y}{\pi} - \frac{x}{\pi} \beta^*\right)^2$, donc susceptible de sous-estimer systématiquement la variance.

→ notes de lecture du Tableau 9 :

– Le taux d'erreur est mesuré par rapport à l'erreur moyenne (rmse) calculée sur les 10 000 simulations du tirage de première phase $\left(\sigma_+^2(\widehat{Y}) = \sum_{sim} (\widehat{Y} - Y)^2 / |sim| \right)$.

– La colonne 'moy' donne la moyenne des taux d'erreurs de l'écart-type estimé : $\sum_{sim} \frac{\widehat{\sigma}(\widehat{Y}) - \sigma_+(\widehat{Y})}{\sigma_+(\widehat{Y})} / |sim|$, égal au taux d'erreur de la moyenne des estimations d'écart-type.

– Dans la colonne 'moy var' figure la racine carrée de la moyenne des erreurs sur la variance estimée $\sqrt{\sum_{sim} \widehat{Var}(\widehat{Y}) / |sim| / \sigma_+^2(\widehat{Y})} - 1$, en %.

– La dernière colonne fournit le rang selon la médiane du taux d'erreur.

– Pour toutes les méthodes, le poids d'une unité de $c \leq 0$ est forcé à $c = 1 - \pi$.

– $v_1 : c = 1 - \pi$

– $v_0 : c = 1$

– $v_2 : c = (1 - \pi) \sum_s (1 - \pi) / \sum_s \delta(1 - \pi)$

– $v_a : c = \sum_s (1 - \pi) / \sum_s \delta(1)$ (équipondération avec calage sur la trace de l'ESB)

– $v_b : c = \sum_s \delta(1) (1 - \pi) / \sum_s \delta(1)^2$ (valeur ν qui minimise $\|\nu \delta(1) - (1 - \pi)\|$)

– $v_h : c = (1 - \pi) |s| / |s| - \text{rang}(x_s)$

– $n_1 : c < \frac{1}{4\pi}, \min_c \|\delta(c) - (1 - \pi)\|_{\delta(1)}$ ³¹, calage final sur la trace de l'ESB

– $n_2 : c < 5, \min_c \|\delta(c) - (1 - \pi)\|_{\delta(1)}$, calage sur la trace de l'ESB

– $n_3 : \text{pas de plafonnement de } c \text{ mais contrainte } \delta(c) < 1, \min_c \|\delta(c) - (1 - \pi)\|_{\delta(1)}$,

calage final sur la trace de l'ESB

– *ren* : algorithme récursif d'au plus 500 itérations³². Un calage final sur la trace de l'ESB est effectué s'il diminue strictement la fonction objectif³³.

– *rel* : Cet algorithme récursif diffère du précédent uniquement par les variables d'équilibrage complémentaire, qui ne sont pris en compte que pour les régions de taille d'échantillon supérieure à 100, pour le calcul à la fois de c et des résidus d'équilibrage.

– *cen* : calage systématique de *ren* sur la trace de l'ESB

– L'indexation de la table contenant les 10 000 échantillons simulés joue crucialement sur la performance du programme. A défaut, la rapidité de lecture d'un échantillon se détériore avec l'indice de sa position. Par exemple, pour $i = 4532$, l'indexation réduit la durée de lecture d'un échantillon simulé de 2 minutes et 48 secondes à 0.10 seconde.

– Les variables étudiées sont issues de la base de sondage, donc de totaux connus :

- *foy_rev*=revenu du foyer fiscal
- *nb_pi2*=nombre de logements principaux TH de deux pièces
- *mar*=(définition incertaine) nombre de foyers de couples mariés
- *isf*=nombre de foyers redevables de l'ISF

31. La pondération de la norme par $\delta(1)$ vise à limiter l'augmentation de poids c_i qui accroît $\delta_j(c)$, où $j \neq i$, lorsque c_j atteint le taquet \bar{c} .

32. Les autres critères d'arrêt sont : $(\sum cuu')^-$ pas calculable par *ginv*, $\|\delta(c) - (1 - \pi)\| < 10^{-8}$ et $\|\Delta\delta(c)\| < 10^{-8}$ (la variation du vecteur des coefficients diagonaux entre les deux dernières itérations).

33. Le calage final sur la trace est effectué ainsi, pour une strate $h : \tilde{c} = c \sum_{s_h} (1 - \pi) / \sum_{s_h} \delta(c)$. Il tend à augmenter l'estimation de variance.

Tableau 9 – Distribution sur 10 000 simulations du taux d'écart à l'erreur quadratique moyenne

variable		moy	moyvar	min	p5	p10	q1	médian	q3	p90	p95	max	rang
foy_rev	v1	4,7	8,5	-40,5	-28,5	-24,1	-15,1	-2,0	16,2	43,6	65,2	155,9	1
foy_rev	v0	5,5	9,3	-40,0	-27,9	-23,5	-14,5	-1,2	17,1	44,8	66,5	157,9	2
foy_rev	v2	10,1	13,9	-37,0	-24,4	-19,9	-10,4	3,3	22,3	50,0	72,0	163,7	11
foy_rev	va	10,1	13,9	-37,0	-24,4	-19,9	-10,4	3,3	22,3	50,1	72,0	163,7	10
foy_rev	vb	9,0	12,8	-37,7	-25,3	-20,7	-11,3	2,2	21,2	48,8	70,8	162,5	4
foy_rev	vh	11,4	15,3	-36,1	-23,3	-18,8	-9,2	4,6	23,9	51,8	73,7	165,5	12
foy_rev	n1	9,3	13,0	-37,3	-24,7	-20,1	-10,7	2,9	21,4	48,7	70,0	160,1	8
foy_rev	n2	9,3	13,0	-37,3	-24,7	-20,1	-10,7	2,8	21,4	48,7	70,0	160,1	7
foy_rev	n3	9,4	13,1	-37,1	-24,7	-20,2	-10,7	2,8	21,6	48,8	70,5	160,5	6
foy_rev	ren	8,8	12,5	-37,7	-25,1	-20,6	-11,2	2,3	20,8	48,0	69,5	159,4	5
foy_rev	cen	9,3	13,0	-37,3	-24,7	-20,1	-10,6	2,9	21,4	48,7	70,0	160,1	9
foy_rev	rel	7,8	11,5	-38,3	-26,0	-21,4	-12,3	1,2	19,8	47,1	68,9	159,1	3
nbpi2	v1	-7,8	-7,8	-16,6	-11,4	-10,6	-9,3	-7,8	-6,3	-4,9	-4,1	0,7	1
nbpi2	v0	-7,1	-7,0	-15,9	-10,7	-9,8	-8,6	-7,1	-5,6	-4,2	-3,4	1,5	2
nbpi2	v2	3,0	3,1	-6,5	-0,9	-0,1	1,4	3,0	4,6	6,1	7,1	13,5	6
nbpi2	va	3,0	3,1	-6,5	-0,9	-0,1	1,4	3,0	4,6	6,1	7,1	13,5	7
nbpi2	vb	0,5	0,5	-8,7	-3,4	-2,5	-1,1	0,5	2,1	3,6	4,5	10,3	4
nbpi2	vh	5,8	5,8	-4,3	1,7	2,6	4,1	5,8	7,5	9,0	10,0	17,4	12
nbpi2	n1	4,3	4,3	-5,7	0,2	1,0	2,5	4,3	6,0	7,6	8,6	14,1	10
nbpi2	n2	4,2	4,2	-5,7	0,1	1,0	2,5	4,2	5,9	7,5	8,4	14,0	9
nbpi2	n3	3,5	3,5	-6,0	-0,5	0,4	1,8	3,5	5,1	6,7	7,6	14,3	8
nbpi2	ren	2,9	3,0	-6,6	-1,1	-0,3	1,2	2,9	4,6	6,2	7,2	12,3	5
nbpi2	cen	4,3	4,3	-5,7	0,2	1,1	2,6	4,3	6,0	7,6	8,6	14,1	11
nbpi2	rel	-0,4	-0,3	-9,9	-4,3	-3,5	-2,0	-0,4	1,3	2,8	3,7	9,0	3
mar	v1	-18,1	-18,1	-24,8	-20,9	-20,3	-19,3	-18,1	-16,9	-15,9	-15,3	-9,7	1
mar	v0	-17,4	-17,4	-24,2	-20,3	-19,6	-18,6	-17,4	-16,3	-15,2	-14,6	-9,0	2
mar	v2	-6,0	-6,0	-13,4	-9,2	-8,6	-7,4	-6,0	-4,6	-3,4	-2,7	3,5	5
mar	va	-6,0	-6,0	-13,4	-9,2	-8,6	-7,4	-6,0	-4,6	-3,4	-2,7	3,5	6
mar	vb	-8,9	-8,9	-16,1	-12,0	-11,4	-10,2	-8,9	-7,6	-6,4	-5,7	0,3	4
mar	vh	-2,8	-2,8	-10,2	-6,2	-5,5	-4,3	-2,8	-1,4	-0,1	0,6	6,9	12
mar	n1	-3,9	-3,9	-12,5	-7,5	-6,8	-5,5	-4,0	-2,4	-1,0	-0,2	5,9	10
mar	n2	-4,0	-4,0	-12,5	-7,6	-6,9	-5,6	-4,1	-2,6	-1,2	-0,3	6,1	9
mar	n3	-5,3	-5,3	-12,9	-8,6	-7,9	-6,7	-5,3	-3,9	-2,6	-1,9	4,4	8
mar	ren	-5,6	-5,6	-13,8	-9,1	-8,4	-7,1	-5,6	-4,1	-2,7	-1,9	4,2	7
mar	cen	-3,9	-3,9	-12,4	-7,5	-6,8	-5,5	-4,0	-2,4	-1,0	-0,1	5,9	11
mar	rel	-9,8	-9,8	-17,7	-13,1	-12,4	-11,2	-9,9	-8,5	-7,1	-6,4	-0,4	3
isf	v1	13,6	13,7	-9,7	3,5	5,6	9,3	13,3	17,6	21,8	24,3	39,9	1
isf	v0	14,4	14,6	-9,0	4,3	6,4	10,1	14,2	18,6	22,7	25,2	41,0	2
isf	v2	20,2	20,4	-3,4	10,0	12,2	15,9	20,0	24,4	28,5	31,0	47,0	8
isf	va	20,2	20,4	-3,4	10,0	12,2	15,9	20,0	24,4	28,5	31,0	47,0	7
isf	vb	18,9	19,0	-4,7	8,7	10,8	14,5	18,7	23,0	27,2	29,7	45,7	4
isf	vh	21,9	22,1	-1,7	11,7	13,8	17,6	21,7	26,1	30,2	32,7	48,8	12
isf	n1	20,7	20,9	-3,5	10,3	12,5	16,3	20,4	25,0	29,1	31,7	48,6	10
isf	n2	20,7	20,8	-3,6	10,2	12,4	16,3	20,4	25,0	29,0	31,6	48,6	9
isf	n3	20,1	20,3	-3,8	9,8	12,0	15,8	19,9	24,4	28,4	31,1	48,0	6
isf	ren	20,0	20,2	-4,2	9,6	11,8	15,6	19,7	24,3	28,4	31,0	48,0	5
isf	cen	20,7	20,9	-3,5	10,3	12,5	16,3	20,4	25,0	29,1	31,7	48,6	11
isf	rel	18,7	18,8	-5,6	8,2	10,4	14,2	18,4	23,0	27,1	29,7	46,4	3

– La variance estimée diffère selon que les variables d'intérêt sont régressées séparément par la procédure REG de SAS ou simultanément. L'estimation est toujours plus élevée avec la régression 'simultanée' des 4 variables. Ceci affecte sensiblement les pondérations à valeurs élevées. Par exemple, le taux d'erreur maximal sur la variance du revenu par la méthode *ren* atteint 1353%. Plutôt qu'un signal de non robustesse, cette anomalie est causée par les observations de valeur manquante sur l'une des variables dépendantes qui sont exclues par la PROC REG du calcul du coefficient ($\widehat{\beta}$), mais pas de celui des résidus. Seule la variable *isf* a des valeurs manquantes dans la base de sondage ³⁴, d'où l'absence de différence entre les deux calculs de l'approximation de variance pour cette variable.

→ Du fait de poids maximaux plus élevés, la pondération optimisée est beaucoup plus sensible à cette incohérence entre les champs du calcul du coefficient et du résidu de la régression. Il faut donc veiller à l'éviter.

La méthode d'optimisation de la pondération par récursion (*ren*) donne une estimation de variance préférable à la meilleure non-optimisée (v_2) dans plus de la moitié des simulations pour les quatre variables (Tableau 10). Mais le bilan est plus mitigé pour les autres pondérations optimisées. Donc il n'y a pas de relation simple entre la hiérarchie des méthodes selon la minimisation de la distance entre diagonales obtenue (Tableau 2) et celle de la qualité de l'estimation de variance.

Tableau 10 – Nombre de simulations où l'estimation de variance est préférable à celle de v_2

variable	v_h	n1	n2	n3	ren	cen	rel
foy_rev	4 399	5 625	5 595	5 549	5 335	5 617	5 179
nbpi2	1 038	1 152	1 137	1 049	5 424	1 151	4 382
mar	9 986	9 984	9 985	9 986	7 193	9 984	0
isf	3	3 100	3 211	6 944	6 975	3 071	9 400

Note : Pour une simulation, une méthode d'estimation de la variance m_1 est jugée préférable à m_2 si m_1 surestime l'écart quadratique moyen (calculé sur l'ensemble des simulations) et m_2 soit le surestime encore davantage soit le sous-estime strictement, ou alors m_1 sous-estime le reqm et m_2 encore plus.

Au total, dans le cas de l'échantillon de première phase d'EEC-TH, la minimisation de la distance entre les diagonales de l'approximation de variance et de l'estimateur sans biais fournit des pondérations de qualité de l'estimation de variance acceptable par rapport au simple calage sur la trace. L'algorithme récursif, qui 'pousse' cette minimisation plus loin, semble fournir une estimation meilleure, en médiane des simulations. De plus, il s'avère possible de moduler cet algorithme. Par exemple, pour le critère de proximité $\|\delta(c) - (1 - \pi)\|^2 + (\sum \delta(c) - \sum 1 - \pi)^2 / |s|$, une variante de l'algorithme récursif a fourni une distribution des précisions sur les 10 000 simulations très proche de n_2 ³⁵, mais avec une pondération maximale encore plus élevée que *ren*. Un calage final systématique de la pondération récursive donne des estimations presque indistinguables de celles de n_1 . C'est vrai également pour la version multiplicative.

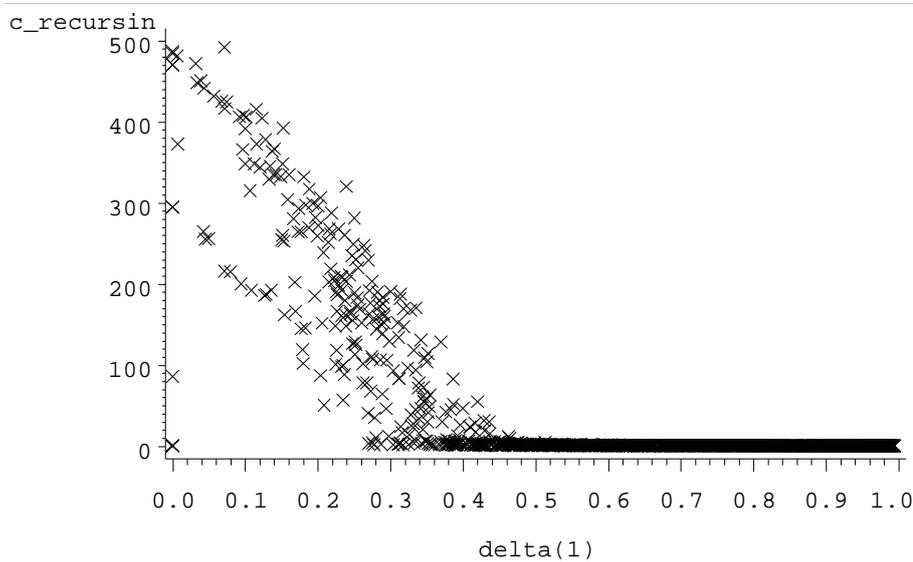
L'incidence de l'optimisation de la pondération est limitée, selon la comparaison des versions v_h et *ren*. L'écart entre les médianes représente 2 à 3% de la racine carrée de l'écart quadratique moyen, c'est-à-dire de l'objectif à estimer pour le calcul de variance. Mais cette différence n'est pas négligeable en comparaison au taux d'erreur de cette estimation pour 3 des 4 variables testées.

34. Il n'y a qu'une dizaine d'unités concernées dans la base de sondage.

35. On se serait attendu plutôt à des résultats intermédiaires entre les versions v_2 et *ren*.

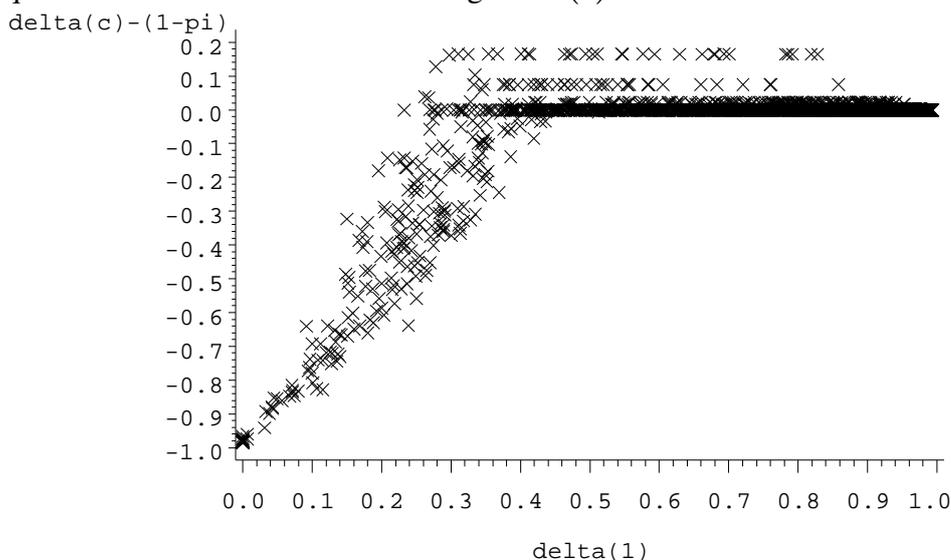
Le rôle du coefficient $\delta(1)$ dans les pondérations élevées produites par les méthodes de minimisation est illustré de manière claire sur la première simulation de la méthode récursive (Graphique 1). Les pondérations de plus de 100 correspondent à des coefficients $\delta(1) \leq 0.4$.

Graphique 1 – Pondération c 'minimisante' de la méthode récursive en fonction du coefficient $\delta(1)$, pour la première simulation d'échantillonnage



Même la méthode récursive ne parvient pas à pousser la minimisation suffisamment pour les plus petits coefficients $\delta(1)$, c'est-à-dire les unités d'indicatrice proche de $\text{Im}(u_s)$ (Graphique 2). De plus, l'augmentation de poids de ces unités a une incidence à la hausse des coefficients diagonaux des autres unités, si la compensation par le poids de celles-ci n'est pas suffisante ou est bloquée par le rejet de l'itération suivante. Ce mécanisme explique sans doute les coefficients diagonaux excédant $1 - \pi$. Ce phénomène traduit en partie l'absence de solution exacte, si l'algorithme fonctionne correctement ³⁶, ou le nombre insuffisant d'itérations.

Graphique 2 – Écart entre le coefficient diagonal $\delta(c)$ et $1 - \pi$ en fonction du coefficient $\delta(1)$



36. Des coefficients diagonaux excessifs ou insuffisants sont vraisemblablement inévitables dès qu'il n'y a pas de solution exacte à l'équation de pondération mais que la dérivée est nulle pour une solution minimisante, vu (e3.45).

– La faiblesse des tailles d'échantillon de certaines strates régionales, ainsi que la dégradation de la qualité numérique de la diagonale calculée lorsque la pondération augmente, pourrait limiter l'avantage de la pondération optimisée (*ren*) dans le contexte étudié. En effet, la taille d'échantillon comparée à la dimension de l'espace d'équilibrage joue notamment sur la fréquence des unités i telles que $\mathbb{1}_i \in \text{Im}(u_s)$, donc $\delta_i(1) = 0$.

– La variable testée de variance la plus surestimée par la méthode des résidus, *isf*, a le coefficient de variation le plus élevé (Tableau 11), et inversement pour la plus sous-estimée (*mar*). Si le plan de sondage est assimilé à un SAS, on s'attendrait à ce que l'approximation poissonnienne surestime davantage les variables de ratio $\mathcal{S}^2(y) / \frac{\sum y^2}{|\mathcal{P}|}$ plus bas, mais l'inverse est observé.

Tableau 11 – Distribution des variables testées dans la base de sondage

statistique	foy rev	nbpi2	mar	isf
CV	74,7	69,8	57,9	186,5
moyenne	5 198 633	25	78,1	2,4
$100\mathcal{S}^2(y) / \frac{\sum y^2}{ \mathcal{P} }$	35,8	32,7	25,1	77,7
R^2	0,77	0,84	0,98	0,52
reqm/total vrai, en %	1,1	0,9	0,3	2,4

– Le lien de la variable d'intérêt avec les variables de l'équilibrage pourrait expliquer la sous-estimation de l'écart quadratique moyen de la variable *mar*. Le R^2 très élevé de cette variable pourrait correspondre à un poids relatif de la variance d'atterrissage plus important, selon G. Chauvet [GC]. Cependant, cette sous-estimation porte sur un reqm très petit relativement à la moyenne de la variable.

– L'écart quadratique moyen augmente sensiblement entre 1 000 et 10 000 simulations pour la variable *mar*, alors qu'il diminue amplement pour *isf* (Tableau 8). Cette instabilité est aussi observée en comparant l'écart quadratique moyen calculé sur les deux moitiés des simulations (Tableau 12). Ces résultats suggèrent qu'une autre piste d'explication des différences observées entre les variables pour la qualité d'estimation de la variance est que l'écart quadratique lui-même serait sur- ou sous-évalué sur les 10 000 simulations.

Tableau 12 – Taux d'écart entre les reqm des 5 000 premières et dernières simulations

reg	total	11	21	22	23	24	25	26	31	41	42	43	52	53	54	72	73	74	82	83	91	93	94
foy rev	2,0	2,0	-5,5	-2,9	-0,6	3,9	-0,3	10,0	3,4	-2,3	1,8	1,9	5,3	-9,8	8,2	6,9	-1,9	-2,8	-7,8	0,9	-4,0	2,8	0,6
nbpi2	-0,7	2,8	5,0	4,4	-6,6	2,3	0,4	-1,4	2,1	-3,2	0,0	-4,5	-1,7	-5,1	1,3	-2,3	-0,1	2,8	-0,2	2,3	-3,6	1,5	-0,8
mar	3,4	3,6	-2,4	-1,4	-4,4	3,2	5,9	3,0	1,5	0,4	1,2	-2,1	4,1	-2,9	-2,7	-2,2	3,8	-0,4	3,4	0,7	5,3	-2,9	2,0
isf	-4,6	-2,7	0,7	-1,1	-10,6	-3,8	-3,0	-1,3	2,2	4,2	-2,5	-3,6	-0,9	0,7	2,4	-2,4	3,1	-5,0	1,3	-2,7	6,2	2,2	7,9

5.1 efficacité pour les quatre plus grandes régions

La comparaison des algorithmes pour l’Ile-de-France, Rhône-Alpes, Provence-Alpes-Côte d’Azur et le Nord-Pas-de-Calais paraît utile pour évaluer le rôle de la faiblesse de la taille d’échantillon régional dans l’incidence limitée de l’optimisation des poids des résidus d’équilibrage sur la qualité de l’approximation de variance. En effet, l’échantillon francilien est de plus de 500 pour la première phase de l’enquête Emploi. Parallèlement, les algorithmes de minimisation annulent presque la distance $\delta(c) - (1 - \pi)$ pour cette région, hors unités de très faible $\delta(1)$ le cas échéant. Pour les trois régions, le poids des coefficients diagonaux supérieurs à 1 est quasiment nulle pour l’algorithme récursif. Donc ces régions pourraient permettre d’évaluer le rôle de l’absence de solution exacte à l’équation de pondération dans les résultats limités de l’optimisation de la pondération.

En fait, les deux pondérations retenues ici ne sont pas clairement hiérarchisées par la qualité de l’estimation de variance pour ces trois régions (Tableau 13). (Les résultats obtenus pour la réunion des trois plus grandes régions sont proches de ceux de l’Ile-de-France.)

Tableau 13 – Distribution de l’erreur d’estimation de la variance pour les quatre régions

région	variable		moy	moyvar	min	p5	p10	q1	médian	q3	p90	p95	max
IdF	foy_rev	v2	13,2	18,5	-41,7	-26,2	-21,5	-10,8	4,5	26,5	63,4	89,6	202,0
	foy_rev	ren	12,1	17,4	-42,2	-26,9	-22,1	-11,6	3,7	25,4	61,2	87,1	197,8
	nbpi2	v2	8,2	8,3	-6,5	1,4	2,9	5,4	8,1	10,9	13,5	15,2	27,6
	nbpi2	ren	8,4	8,5	-6,6	1,6	3,1	5,6	8,3	11,2	13,8	15,5	28,7
	mar	v2	4,0	4,1	-11,6	-2,1	-0,9	1,4	4,0	6,5	9,0	10,5	18,4
	mar	ren	4,9	5,0	-11,3	-1,6	-0,3	2,1	4,8	7,5	10,2	11,8	19,6
	isf	v2	36,1	36,4	0,7	21,1	24,3	29,7	36,0	42,4	48,3	51,9	72,8
	isf	ren	35,4	35,8	-0,3	20,3	23,5	29,0	35,2	41,8	47,7	51,4	73,4
RA	foy_rev	v2	-9,9	2,8	-60,2	-48,0	-45,3	-39,2	-26,1	-2,3	42,0	84,3	322,0
	foy_rev	ren	-10,5	1,3	-60,5	-47,9	-45,3	-39,1	-26,1	-1,9	38,5	80,2	304,8
	nbpi2	v2	3,9	4,2	-19,4	-8,4	-6,0	-1,9	3,2	8,7	14,5	18,5	44,0
	nbpi2	ren	4,9	5,3	-19,5	-7,8	-5,4	-1,1	4,2	10,0	16,3	20,3	48,7
	mar	v2	1,2	1,4	-20,8	-9,2	-7,1	-3,4	0,9	5,4	9,7	12,5	31,8
	mar	ren	3,4	3,7	-20,1	-8,0	-5,6	-1,7	2,9	8,0	13,1	16,6	44,8
	isf	v2	2,9	4,3	-41,9	-22,7	-18,2	-9,2	1,5	13,9	25,7	32,9	69,9
	isf	ren	4,0	5,4	-42,2	-22,4	-17,6	-8,6	2,3	15,4	28,0	35,6	77,5
PACA	foy_rev	v2	-3,8	4,7	-47,8	-30,7	-27,2	-21,0	-12,7	-0,8	15,4	34,6	315,8
	foy_rev	ren	-3,9	4,2	-47,5	-30,8	-27,4	-21,0	-12,6	-0,7	15,4	35,0	306,4
	nbpi2	v2	6,9	7,1	-16,3	-3,7	-1,6	2,3	6,7	11,3	15,6	18,4	37,2
	nbpi2	ren	7,6	7,8	-16,7	-3,5	-1,2	2,7	7,5	12,2	16,6	19,7	38,3
	mar	v2	-3,8	-3,6	-27,2	-15,1	-12,8	-9,0	-4,3	0,9	5,8	8,7	25,8
	mar	ren	-1,7	-1,3	-26,8	-14,0	-11,5	-7,4	-2,3	3,3	9,2	13,1	35,4
	isf	v2	1,7	3,4	-43,0	-22,5	-18,6	-11,0	-1,6	9,9	27,9	40,3	99,7
	isf	ren	2,8	4,6	-42,6	-22,5	-18,3	-10,6	-0,7	11,8	30,8	42,3	101,5
NPDC	foy_rev	v2	-6,6	-5,0	-43,9	-29,4	-26,1	-19,1	-9,4	2,6	16,0	25,1	83,7
	foy_rev	ren	-5,7	-4,1	-42,9	-29,3	-25,8	-18,6	-8,6	4,4	17,3	26,8	97,4
	nbpi2	v2	-1,5	-0,9	-32,0	-17,3	-14,4	-8,9	-2,2	5,0	12,6	17,2	37,6
	nbpi2	ren	-0,4	0,2	-31,8	-17,0	-13,9	-8,2	-1,3	6,3	14,5	19,5	42,3
	mar	v2	-8,3	-8,1	-34,1	-19,4	-16,9	-12,9	-8,5	-3,8	0,4	3,1	20,7
	mar	ren	-7,3	-7,0	-33,2	-19,1	-16,5	-12,2	-7,5	-2,5	2,1	5,1	22,0
	isf	v2	-12,6	-8,5	-58,8	-46,1	-41,8	-32,8	-18,5	1,0	27,6	45,3	112,1
	isf	ren	-9,0	-4,3	-59,2	-45,3	-40,8	-31,3	-15,8	6,1	34,4	51,0	132,2

Conclusion

– Sur le plan théorique, l'étude a mis en évidence une condition nécessaire d'existence d'une solution exacte à l'équation de pondération, qui égalise les deux diagonales. Elle reflète que les unités atypiques entravent le rapprochement des diagonales, comme l'équilibrage. Parce que la condition nécessaire n'est pas toujours vérifiée lorsque la dimension d'équilibrage est supérieure à 1, il a semblé judicieux d'élargir l'objectif à la minimisation d'une distance entre les deux diagonales. Dans cette optique, l'algorithme récursif assure une baisse très sensible de la distance inter-diagonales par rapport à la pondération par $1 - \pi$, même calée sur la trace. Le deuxième résultat méthodologique est une condition suffisante pour qu'une pondération produite par cette procédure soit finie.

– En revanche, l'existence d'une solution minimisant la distance entre les diagonales des deux estimateurs de variance n'a pas pu être démontrée. La convergence de l'algorithme vers une solution exacte n'est pas assurée non plus. De plus, les propriétés de la pondération minimisante (éventuelle) restent en partie inconnues. Il n'y a pas de garantie théorique qu'une telle pondération soit non nulle et finie et a fortiori qu'elle soit préférable pour la qualité de l'estimation de variance. Une autre piste de recherche méthodologique est d'améliorer le critère de proximité, par exemple en prenant en compte la trace. En effet, le choix de distance joue un rôle crucial dans le compromis entre les 'excès' et les 'défauts' des coefficients diagonaux, en l'absence de solution exacte à l'équation de pondération.

– Les résultats des différentes versions testées suggèrent qu'il est encore possible d'améliorer sensiblement l'algorithme de minimisation de la distance entre diagonales.

– Empiriquement, selon l'objectif technique de rapprochement des diagonales, l'algorithme récursif ressort clairement comme le meilleur choix parmi ceux évalués pour optimiser la pondération. Cependant, l'ampleur et la régularité de la minimisation par cet algorithme restent inexplicables. Donc il n'est pas sûr que la méthode soit aussi performante sur d'autres données.

– Pour l'objectif final d'améliorer la qualité de l'estimation de variance, les résultats obtenus confortent globalement le choix d'optimiser la pondération selon le principe de minimisation de la distance des diagonales. Cette méthode réduit l'erreur d'estimation de la variance de manière relativement sensible. Les poids maximaux élevés ne paraissent pas affecter la robustesse de la méthode. Mais l'enjeu par rapport à la pondération simple 'calée sur la trace' est limité relativement à la variance, sur les données étudiées. La supériorité de la pondération optimisée ressort donc moins clairement que dans les tests de [DT], peut-être du fait du contexte d'absence de solution exacte. L'écart entre la pondération optimisée et celle de Hajek est de 1 à 3% de l'écart-type à estimer. Le contrôle des coefficients diagonaux peut importer pour l'estimation de variance d'un sondage en deux phases.

La conclusion opérationnelle de l'étude est qu'il faut éviter de pondérer l'estimateur de variance par 1 ou par $1 - \pi$, parce que cela mène en général à une sous-estimation sensible de la variance. Le choix parmi les autres pondérations est moins tranché. La pondération optimisée par la procédure récursive constitue une alternative crédible aux pondérations plus 'simples' et peut être préférée au calage de $1 - \pi$ sur la trace pour le plan de sondage étudié. Cependant, la justification théorique et les résultats obtenus dans la présente étude peuvent être jugés insuffisants pour préférer ce choix à la pondération de Hajek préconisée par [AMYT]. Celle-ci limite davantage les sous-estimations, au prix de surestimations accentuées par ailleurs.

Références

- [DT] Variance approximation under balanced sampling, JC Deville, Y Tillé, Journal of statistical planning and inference 2005
- [AMYT] Evaluation of variance approximations and estimators in maximum entropy sampling with unequal probability and fixed sample size, A.Matei, Y.Tillé, Journal of Official Statistics vol.21 N° 4, 2005
- [JCD] Variance estimation for complex statistics and estimators : linearisation and residual techniques, JC Deville, Survey Methodology, vol.25 N° 2, 1999
- [VL] La construction du nouvel échantillon de l'enquête emploi en continu à partir des fichiers de la taxe d'habitation, V.Loonis, Actes des JMS 2009
- [CDL] Weighted finite population sampling to maximize entropy, XH Chen, AP Dempster, JS Liu, Biometrika, 1994
- [GC] On variance estimation for the French master sample, Journal of Official Statistics, vol.27 N° 4, 2011

Annexe A matrices symétriques

- Pour une matrice A d'un espace vectoriel $E(= \mathbb{R}_s)$ dans un espace F :
 - $\text{rang}(A) + \dim[\text{Ker}(A)] = \dim(E)$: Il existe un sous-espace $\Delta E / E = \text{Ker}(A) \oplus \Delta E$ et alors A est bijective de ΔE sur $\text{Im}(A)$.

$$\text{Im}(A) = \text{Ker}(A')^\perp \quad (\text{eA.56})$$
 - car $\text{Im}(A)^\perp = \{\langle AE, \cdot \rangle = 0\} = \{\langle E, A' \cdot \rangle = 0\} = \{A' \cdot = 0\} = \text{Ker}(A')$ (Il en découle que $\text{rang}(A') = \dim[F - \text{Ker}(A')] = \dim[\text{Ker}(A')^\perp] = \text{rang}(A)$.)
 - $$\text{Im}(A) = \text{Im}(AA') \quad (\text{eA.57})$$
 - $\hookrightarrow \text{Im}(AA') \subset \text{Im}(A)$
 - $\hookrightarrow \text{Ker}(AA') = \text{Ker}(A')$ car $\forall h, h'AA'h = \|A'h\|^2$
 - $\hookrightarrow \Rightarrow \text{rang}(AA') = \text{rang}(A') = \text{rang}(A)$
 - Si A est une matrice symétrique ($A' = A$) alors :
 - $\text{Im}(A) = \text{Ker}(A)^\perp$
 - A est bijective sur $\text{Im}(A)$ (car surjective et $\text{Im}(A) \cap \text{Ker}(A) = 0 \Rightarrow A$ injective sur $\text{Im}(A)$)
 - Si A est une matrice symétrique alors elle est diagonalisable dans une base orthonormée. Donc il existe une matrice orthogonale U , c'est-à-dire que $UU' = U'U = \text{id}$, et une matrice D diagonale telles que $A = UDU'$.
 - Par définition, une matrice carrée est positive si $\forall x, x'Ax \geq 0$ et définie positive si $\forall x \neq 0, x'Ax > 0$. Pour une matrice définie positive A , $\|x\|_A = \sqrt{x'Ax}$ est une norme ³⁷.
 - Pour deux matrices symétriques A et B , $A \leq B \stackrel{\text{déf}}{\Leftrightarrow} \forall x, x'Ax \leq x'Bx$
 - Si tous les mineurs principaux d'une matrice M (non nulle) symétrique sont positifs, alors le déterminant s'écrit $|M - \lambda \text{id}| = \sum_k (-\lambda)^k a_k$ avec $a_k \geq 0$, et $a \neq 0$. Donc ce polynôme en λ n'admet pas de racines strictement négatives. (Par ailleurs, le critère de Sylvester, plus fort, énonce que si les n mineurs principaux du type $|M_{\leq k}|$ sont > 0 alors M est définie positive.)

37. A est (définie) positive si la matrice symétrique $A + A'$ est (définie) positive, car $x'Ax = x'A'x = x'(A + A')x/2$ et donc $\|x\|_A = \|x\|_{(A+A')/2}$. Néanmoins, l'extension de la définition peut être utile parce que l'ensemble des matrices symétriques n'est pas un ouvert.

- Si Δ est une matrice symétrique et positive alors pour tous les indices i, j :

$$\forall x, (x\mathbb{1}_i + \mathbb{1}_j)' \Delta (x\mathbb{1}_i + \mathbb{1}_j) \geq 0$$

$$\text{donc } \forall x, x^2 \Delta_{i,i} + 2\Delta_{i,j}x + \Delta_{j,j} \geq 0$$

$$\text{donc le discriminant est négatif : } \Delta_{i,j}^2 - \Delta_{i,i}\Delta_{j,j} \leq 0 \quad (\text{eA.58})$$

→ Cette condition n'est pas suffisante si $|s| > 2$: soit $a_s > 0$ et $\Delta_{i,j} = (2\mathbb{1}_{i=j} - 1)\sqrt{a_i}\sqrt{a_j}$

$$\sum_i \sqrt{a_i}^2 \frac{1}{\sqrt{a_i}} \frac{1}{\sqrt{a_i}} - \sum_{i \neq j} \sqrt{a_i}\sqrt{a_j} \frac{1}{\sqrt{a_i}} \frac{1}{\sqrt{a_j}} = |s| - |s|(|s| - 1) = |s|(2 - |s|) < 0$$

Annexe B inverses généralisées

- A^- désigne une inverse généralisée de la matrice carrée A , définie par la propriété : $AA^-A = A$. Cette définition est telle que si $b \in \text{Im}(A)$ alors $a = A^-b$ est une solution de l'équation $Aa = b$.

→ Si $B[\text{Im}(A)] \subset \text{Ker}(A)$ alors $A^- + \widetilde{B}$ est une inverse généralisée.

- Deux inverses généralisées de A , A^- et $\widetilde{A^-}$, diffèrent par une matrice B telle que $B[\text{Im}(A)] \subset \text{Ker}(A)$, vu que $A(A^- - \widetilde{A^-})A = 0$.

- La qualité d'une inverse généralisée est mesurée dans la présente étude par $\|A - AA^-A\|$, avec la norme euclidienne, c'est-à-dire la racine carrée de la somme des carrés des écarts entre les coefficients des deux matrices.

- Pour une matrice A , A^+ est l'inverse généralisée de Moore-Penrose (MP). Elle est définie par les quatre propriétés :

$$AA^+A = A \quad (\text{mp1})$$

$$A^+AA^+ = A^+ \quad (\text{mp2})$$

$$(A^+A)' = A^+A \quad (\text{mp3})$$

$$(AA^+)' = AA^+ \quad (\text{mp4})$$

$$\begin{aligned} \text{Im}(A) \stackrel{(\text{mp1})}{\subset} \text{Im}(AA^+) \stackrel{(\text{mp4})}{\supseteq} \text{Im}[(A^+)'A'] \subset \text{Im}[(A^+)''] \stackrel{(\text{mp2})}{\subset} \text{Im}[(AA^+)'] \stackrel{(\text{mp4})}{\supseteq} \text{Im}(AA^+) \subset \text{Im}(A) \\ \Rightarrow \text{Im}(A) = \text{Im}[(A^+)'] = \text{Ker}(A^+)^{\perp} \quad (\text{eB.59}) \end{aligned}$$

• La première égalité de (eB.59) entraîne que $\text{rang}[A^+] = \text{rang}(A)$.

→ Comme A et A^+ jouent des rôles symétriques dans les propriétés de Moore-Penrose, on peut permuter dans (eB.59), d'où l'image de l'inverse MP :

$$\text{Im}(A^+) = \text{Im}[A'] = \text{Ker}(A)^{\perp} \quad (\text{eB.60})$$

⇒ Si A est une matrice symétrique alors $\text{Im}(A^+) = \text{Im}(A)$ et $\text{Ker}(A^+) = \text{Ker}(A)$.

• unicité : Si A^+ et $\widetilde{A^+}$ vérifient les quatre conditions de Moore-Penrose, alors elles ont même noyau $(\text{Im}(A^+))^{\perp}$, vu (eB.59). Par conséquent, il existe une matrice inversible B telle que $\widetilde{A^+} = BA^+$.

(mp1) $\Leftrightarrow ABA^+A = A$. D'autre part, (mp3) $\Rightarrow (BA^+A)' = BA^+A$ soit : $A^+AB' = BA^+A \Rightarrow AA^+AB' = ABA^+A \Rightarrow AB' = ABA^+A = A$, d'après la première égalité. Donc $AB' = A \Rightarrow A^+AB' = A^+A \stackrel{(\text{mp3})}{\Rightarrow} BA^+A = A^+A$ et en multipliant à droite par A^+ , il vient finalement $A^+ = BA^+ = \widetilde{A^+}$

⇒ Les conditions de Moore-Penrose définissent une inverse unique (si elle existe). Il en découle :

- $(A^+)' = (A^+)^{\perp}$

- si $\alpha \neq 0$ alors $(\alpha A)^+ = \frac{1}{\alpha}A^+$

- si A et B sont deux matrices symétriques et si $AB = 0 = BA$ alors :

$$(A + B)^+ = A^+ + B^+ \quad (\text{eB.61})$$

$$\left\{ \begin{array}{l} \text{Im}(A^+) = \text{Im}(A) \subset \text{Ker}(B) = \text{Ker}(B^+) \\ \text{Im}(B^+) = \text{Im}(B) \subset \text{Ker}(A) = \text{Ker}(A^+) \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} A^+B = 0 \\ B^+A = 0 \\ BA^+ = 0 \\ AB^+ = 0 \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} (A^+ + B^+)(A + B) = A^+A + B^+B \\ (A + B)(A^+ + B^+) = AA^+ + BB^+ \end{array} \right.$$

Les six dernières identités entraînent que $A^+ + B^+$ vérifient les propriétés de Moore et Penrose pour $A + B$.

→ Si A est une matrice symétrique alors l'inverse de Moore-Penrose est donnée par $A^+ = UD^+U'$, où D^+ est la matrice diagonale constituée de l'inverse des coefficients non nuls de D (une diagonalisation de A) et de ses coefficients nuls. L'inverse généralisée peut également dans ce cas se construire à partir du fait qu'une matrice symétrique A est bijective sur $\text{Im}(A)$ et nulle sur $\text{Im}(A)^\perp$.

– si $A = \alpha uu'$ avec $\alpha \neq 0$ et $\|u\| \neq 0$ alors $A^+ = \frac{1}{\alpha} \frac{uu'}{\|u\|^4}$

– décroissance de l'inversion de Moore et Penrose :

Si A et B sont deux matrices symétriques de même image :

$$\text{si } A \leq B \text{ alors } B^+ \leq A^+ \quad (\text{eB.62})$$

↳ Un même changement de base orthonormée permet de se ramener à deux matrices symétriques et inversibles :

$$A = U \begin{pmatrix} \tilde{A} & 0 \\ 0 & 0 \end{pmatrix} U' \text{ et } B = U \begin{pmatrix} \tilde{B} & 0 \\ 0 & 0 \end{pmatrix} U'$$

$$\tilde{A} \leq \tilde{B} \Rightarrow \tilde{B}^{-\frac{1}{2}} \tilde{A} \tilde{B}^{-\frac{1}{2}} \leq \text{id} \text{ (car } A \leq B \Rightarrow \forall C, CAC' \leq CBC')$$

$$\Rightarrow \tilde{B}^{\frac{1}{2}} \tilde{A}^{-1} \tilde{B}^{\frac{1}{2}} \geq \text{id} \text{ (en utilisant la propriété précédente)}$$

$$\Rightarrow \tilde{A}^{-1} \geq \tilde{B}^{-1} \Rightarrow A^+ \geq B^+$$

Annexe C projection orthogonale

Un projecteur M -orthogonal, p , est défini par les deux propriétés :

$$\rightarrow (\text{id} - p)'Mp = 0 \quad (\text{eC.63})$$

$$\rightarrow p^2 = p \text{ (idempotence)} \quad (\text{eC.64})$$

$$\Rightarrow \text{Ker}(p) = \{p = 0\} = \{\|p\|_M = 0\} = \{p'Mp = 0\} = \{x / x'p'Mp = 0\} \stackrel{(\text{eC.63})}{=} \{x / x'Mp = 0\} = \text{Im}(p)^{\perp M}$$

– $\text{proj}_{\text{Im}(u_s)}(z) = u_s \beta$ tel que $\langle z - u_s \beta, u_s \rangle = 0$ soit $\langle u'_s z - u'_s u_s \beta, \rangle = 0$ et donc $u'_s u_s \beta = u'_s z$.
 $\beta = (u'_s u_s)^{-1} u'_s z$ est une solution, d'où finalement :

$$\text{proj}_{\text{Im}(u_s)} = u_s (u'_s u_s)^{-1} u'_s \quad (\text{eC.65})$$

– AA^+ est le projecteur orthogonal sur $\text{Im}(A)$ car $\text{Im}(AA^+) = \text{Im}(A)$, $AA^+AA^+ = AA^+$ et $(AA^+)' = AA^+$.

– u_s est constante sur l'ensemble $\text{argmin}_\beta \|z - u_s \beta\|_c$, de valeur $\text{proj}_{\text{Im}(u_s)}^c(z)$:

si $\tilde{\beta}, \tilde{\beta} \in \text{argmin}_\beta \|z - u_s \beta\|_c$ alors :

$$\|z - u_s \tilde{\beta}\|_c^2 = \|z - u_s \tilde{\beta}\|_c^2 = \|z - u_s \tilde{\beta}\|_c^2 + \|u_s \tilde{\beta} - u_s \tilde{\beta}\|_c^2 \Rightarrow \|u_s \tilde{\beta} - u_s \tilde{\beta}\|_c^2 = 0$$

Annexe D dérivabilité infinie de $\delta(c)$

• $\text{proj}_{\text{Im}(u_s)}^c(z)$ est C^∞ sur $\{c > 0\}$:

Soit \tilde{u}_s une sous-matrice de u_s de plein rang colonne et surjective sur $\text{Im}(u_s)$. Alors $\text{proj}_{\text{Im}(u_s)}^c(z) = \tilde{u}_s \tilde{\beta}$ avec $\langle z - \tilde{u}_s \tilde{\beta}, \tilde{u}_s \rangle_c = 0$ et $\varphi(c) = \left\| z - \tilde{u}_s \tilde{\beta} \right\|_c^2$

– Une démonstration de la dérivabilité infinie ($\delta \in C^\infty(\mathbb{R}_+^*)$) par le théorème d'inversion locale utilise le fait que la fonction $\Phi(c, \tilde{\beta}) = \langle z - \tilde{u}_s \tilde{\beta}, \tilde{u}_s \rangle_c = \langle z - \tilde{u}_s \tilde{\beta}, c \tilde{u}_s \rangle$ est C^∞ (comme produit scalaire de deux fonctions linéaires), et comme \tilde{u}_s est de plein rang colonne et $c > 0$, $\frac{\partial \Phi}{\partial \tilde{\beta}}(c, \tilde{\beta}) = -\langle \tilde{u}_s, \tilde{u}_s \rangle_c$

est inversible pour tout $c > 0$. Donc $\Psi(c, \tilde{\beta}) = (c, \Phi(c, \tilde{\beta}))$ est un difféomorphisme local, et $\tilde{\beta}(c) = \text{proj}_{\text{Vect}(\tilde{\beta})} \Psi^{-1}(c, 0)$ est C^∞ .

– Une démonstration plus spécifique :

$\widehat{\beta} = [\widehat{u}^s * c\widehat{u}_s]^{-1} \widehat{u}^s * cz$ est le produit matriciel de deux fonctions C^∞ , parce que la fonction $A \mapsto A^{-1}$ est C^∞ sur l'ensemble des matrices inversibles. (Si A est inversible et si $\|H\| < 1/\|A^{-1}\|$ alors $(A+H)^{-1} = A^{-1}(\text{id}+A^{-1}H)^{-1} = A^{-1} \sum (-A^{-1}H)^{\mathbb{N}}$).

→ Il s'ensuit que $\text{proj}_{\text{Im}(u_s)}^c(z) = \widehat{u}_s \widehat{\beta} = u_s \widehat{\beta}_c$ est C^∞ sur $\{c > 0\}$ et par suite $\varphi(c) = \langle [z - u_s \widehat{\beta}_c]^2, c \rangle = \langle z - u_s \widehat{\beta}_c, cz \rangle$ aussi (comme produit scalaire de deux fonctions C^∞).

• $(\sum cuu')^+$ est dérivable (et C^∞) en c sur $\mathbb{R}_+^*(s)$:

↳ Comme $\|z\|_c^2 - \varphi_z(c) = \|u_s \widehat{\beta}_c(z)\|_c^2$ est dérivable (car différence de deux dérivables) :

$$\begin{aligned} \|u_s \widehat{\beta}_c(z)\|_c^2 &= \sum_j c_j [u_j' (\sum cuu')^- \sum czu]^2 \\ &= (\sum czu') (\sum cuu')^- (\sum cuu') (\sum cuu')^- (\sum czu) \\ &= (\sum czu') (\sum cuu')^- (\sum czu) \\ &= \sum_{i,j} c_i c_j z_i z_j u_i' (\sum cuu')^- u_j \end{aligned}$$

est dérivable en $c > 0$ pour tout z .

⇔ $\forall i, j \in s, u_i' (\sum cuu')^- u_j$ est dérivable, donc $\sum uu' (\sum cuu')^- \sum uu'$ est dérivable.

↳ Il existe une (unique) matrice Φ inversible telle que :

$$\begin{cases} \sum uu' \Phi = \text{id} = \Phi \sum uu' \text{ sur } \text{Im}(\sum uu') \\ \Phi = \text{id} \text{ sur } \text{Ker}(\sum uu') \end{cases} \quad (\text{eD.66})$$

Ceci découle de ce que comme $\sum uu'$ est symétrique, $E = \text{Im}(\sum uu') \oplus \text{Ker}(\sum uu')$ et $\sum uu'$ est bijective sur $\text{Im}(\sum uu')$.

↳ La matrice Φ est symétrique, du fait de l'unicité de la solution de (eD.66).

↳ Comme Φ ne dépend pas de c , il en découle que $\Gamma(c) = \Phi' \sum uu' (\sum cuu')^- \sum uu' \Phi$ est dérivable en c .

▷ sur $\text{Ker}(\sum uu')$, $\Gamma(c) = 0$ car $\sum uu' \Phi [\text{Ker}(\sum uu')] = \sum uu' \text{Ker}(\sum uu') = \{0\}$ donc $\Gamma(c) = (\sum cuu')^+$ sur ce sous-espace vectoriel (car $\text{Ker}(\sum uu') = \text{Ker}((\sum uu')^+)$).

▷ Pour $x \in \text{Im}(\sum uu')$:

$$\begin{aligned} \Gamma(c)x &= \Phi' \sum uu' (\sum cuu')^+ \sum uu' \Phi x \\ &= \Phi \sum uu' (\sum cuu')^+ x \text{ (car } \sum uu' \Phi = \text{id sur } \text{Im}(\sum uu')) \\ &= (\sum cuu')^+ x \end{aligned}$$

(car $\text{Im}[(\sum cuu')^+] = \text{Im}(\sum cuu') = \text{Im}(\sum uu')$ et $\Phi \sum uu' = \text{id sur } \text{Im}(\sum uu')$)

⇒ $\Gamma(c) = (\sum cuu')^+$, et cette dernière fonction de c est donc également dérivable.

• une formulation de la dérivée seconde de φ :

$$\begin{aligned} \ddot{\varphi}(c) &= \frac{d}{dc} \|z - u_s \widehat{\beta}_c\|_{dc} = 2 \langle z - u_s \widehat{\beta}_c, -d(u_s \widehat{\beta}_c) \rangle_{dc} \\ \langle u_s (\widehat{\beta}_c - \widehat{\beta}_{c+h}), u_s \rangle_{c+h} &= \langle z - u_s \widehat{\beta}_{c+h}, u_s \rangle_{c+h} - \langle z - u_s \widehat{\beta}_c, u_s \rangle_{c+h} \\ &= - \langle z - u_s \widehat{\beta}_c, u_s \rangle_{c+h} \text{ (car } \widehat{\beta}_{c+h} \text{ vérifie l'équation normale pour } \langle \rangle_{c+h}) \\ &= - \langle z - u_s \widehat{\beta}_c, u_s \rangle_h \text{ (car } \widehat{\beta}_c \text{ vérifie l'équation normale pour } \langle \rangle_c) \text{ d'où :} \\ &\quad \langle u_s (\widehat{\beta}_{c+h} - \widehat{\beta}_c), u_s \rangle_{c+h} = \langle z - u_s \widehat{\beta}_c, u_s \rangle_h \end{aligned} \quad (\text{eD.67})$$

(eD.67) permet de formuler la dérivée du projecteur ainsi :

$$\langle d(u_s \widehat{\beta}_c), u_s \rangle_c = \langle z - u_s \widehat{\beta}_c, u_s \rangle_{dc} \quad (\text{eD.68})$$

(Cette équation détermine effectivement $d(u_s \widehat{\beta}_c) : d(u_s \widehat{\beta}_c)h = u_s (u^s * cu_s)^- u^s * h [z - u \widehat{\beta}_c]_s$).

Il découle de (eD.68) :

$$\begin{aligned} \langle z - u_s \widehat{\beta}_c, u_s \widehat{\beta}_{c+h} \rangle_{dc} &= \langle d(u_s \widehat{\beta}_c), u_s \widehat{\beta}_{c+h} \rangle_c \\ \Rightarrow \langle z - u_s \widehat{\beta}_c, d(u_s \widehat{\beta}_c) \rangle_{dc} &= \langle d(u_s \widehat{\beta}_c), d(u_s \widehat{\beta}_c) \rangle_c \text{ d'où :} \\ \ddot{\varphi}(c) &= -2 \langle z - u_s \widehat{\beta}_c, d(u_s \widehat{\beta}_c) \rangle_{dc} = -2 \langle d(u_s \widehat{\beta}_c), d(u_s \widehat{\beta}_c) \rangle_c \\ \ddot{\varphi}(c) hh &= -2 [z - u \widehat{\beta}_c]^s h u_s (u^s * c u_s)^- u^s h [z - u \widehat{\beta}_c]_s \\ &= -2 \left\| \text{Proj}_{\text{Im}(u_s)}^c \left\{ \frac{h}{c} [z - u \widehat{\beta}_c]_s \right\} \right\|_c^2 \end{aligned}$$

Le signe de cette dernière formule permet de retrouver la concavité de φ .

Annexe E rang de $\dot{\delta}(c)$ et contraction de $\text{id} - \dot{\delta}(c)$

E.1 norme d'une fonction linéaire

Si $f \in \mathcal{L}(E, F)$ ³⁸, E de dimension finie de base e_s et F normé alors :

$$\|f\| \stackrel{\text{déf}}{=} \sup_{h \neq 0} \|fh\| / \|h\| \quad (\text{eE.69})$$

→ $\|f\| < +\infty$:

↳ première approche :

$\sup_{h \neq 0} \|fh\| / \|h\| = \sup_{\|h\|=1} \|fh\| < \infty$ car f est continue et $\{\|h\|=1\}$ compact parce que $\dim(E) < \infty$

↳ deuxième approche (sans admettre la continuité de f) :

$$\begin{aligned} \sup_{\|x\|=1} \|f(x)\| &\leq \sup_{\|x\|=1} \sum |x_i| \|f(e_i)\| \\ &\leq \nu(x) \sum \|f(e_i)\| \quad (\text{où } \nu(x) = \sup(x_i, i \in s)) \\ &\leq \eta \|x\| \text{ avec } \eta < +\infty \text{ (en admettant l'équivalence des normes)} \\ &\Rightarrow \sup_{\|x\|=1} \|f(x)\| < +\infty \end{aligned}$$

$$\Rightarrow f \text{ continue : } \|f(y) - f(x)\| = \|f(y-x)\| \leq \|f\| \|y-x\|$$

$$- \forall h, |h'Ah| \leq \|h\| \|Ah\| \leq \|A\| \|h\|^2$$

• si $f \in \mathcal{L}(\mathbb{R}_s)$ muni de la distance euclidienne canonique alors $\|f\| = \sup \sqrt{\Lambda(f'f)}$ ³⁹. En effet, soit v une base orthonormée de vecteurs propres de la matrice symétrique positive $f'f$. Pour $\|x\| = 1$, $x'f'fx = \sum \lambda x^2 \leq \sup[\lambda] \sum x^2 = \sup \Lambda(f'f) \Rightarrow \|f'f\| = \sup \Lambda(f'f)$.

$$\rightarrow \|f\|^2 = \sup_{\|h\|=1} (\|fh\|)^2 = \sup_{\|h\|=1} h'f'fh = \|f'f\| \text{ donc :}$$

$$\|f\|^2 = \|f'f\| = \sup \Lambda(f'f) \quad (\text{eE.70})$$

• Les espaces propres de $f'f'$ de valeurs propres non nulles sont isomorphes à ceux de $f'f$: Soit $\lambda \in \Lambda(f'f) \setminus 0$

→ f est injective sur $E_\lambda(f'f)$, vu que $E_\lambda(f'f) \cap \text{Ker}(f) = 0$ pour $\lambda \neq 0$

→ $f[E_\lambda(f'f)] \subset E_\lambda(f'f)$: pour $h \in E_\lambda(f'f) \setminus 0$, $f(h) \neq 0$ et $f'f'f(h) = \lambda f(h)$ donc $f(h) \in E_\lambda(f'f)$

→ par symétrie, $f'E_\lambda(f'f) \subset E_\lambda(f'f)$

→ $f[E_\lambda(f'f)] = E_\lambda(f'f)$ car f' est injective sur $E_\lambda(f'f)$ (par symétrie de l'injectivité de f sur $E_\lambda(f'f)$), donc $\dim[E_\lambda(f'f)] \leq \dim[E_\lambda(f'f)]$ et symétriquement $\dim[E_\lambda(f'f)] \leq \dim[E_\lambda(f'f)]$. Par suite $\dim[E_\lambda(f'f)] = \dim[E_\lambda(f'f)]$.

38. $\mathcal{L}(E, F)$ est l'ensemble des applications linéaires entre les espaces vectoriels E et F .

39. Λ désigne l'ensemble des valeurs propres et VecP celui des vecteurs propres.

- pour la norme du sup :

$$\begin{aligned}\|\text{id} - \dot{\delta}(c)\|_\infty &= \sup_i \sup_{\|h\|_\infty=1} |h_i - \|\mathbf{1}_i - u_s \widehat{\beta}_i(c)\|_1|^2 \\ &= \sup_i 1 - \{1 - u'_i \widehat{\beta}_i(c)\}^2 + \sum_{j \neq i} \{u'_j \widehat{\beta}_i(c)\}^2 \\ &= 1 + \|\mathbf{1}_i - u_s \widehat{\beta}_i(c)\|_1^2 - 2 \left[\frac{\delta_i(c)}{c_i} \right]^2\end{aligned}$$

Donc $\text{id} - \dot{\delta}(c)$ est contractante pour cette norme si $\forall i \in s, \left[\frac{\delta_i(c)}{c_i} \right]^2 > \frac{1}{2} \|\mathbf{1}_i - u_s \widehat{\beta}_i(c)\|_1^2 = \frac{1}{2} \delta_i(1)$

\Rightarrow Pour c constante, $\text{id} - \dot{\delta}(c)$ est contractante ssi $\forall i \in s, \delta_i(1) > \sqrt{2} / 2$.

E.3.1 expression de $\dot{\delta}(c)$ en fonction d'un projecteur

Dans la suite, $q = q_c$ désigne le projecteur orthogonal (pour le produit scalaire) canonique sur $\text{Im}[\sqrt{c}u_s]^\perp$, q^{*2} est la matrice composée des carrés des coefficients de q , les matrices diagonales sont confondues avec les vecteurs correspondants, et la commutativité du produit de matrices diagonales est utilisée.

- expression de la dérivée en fonction de la projection orthogonale canonique :

$$\begin{aligned}\text{proj}_{\text{Im}(u_s)}^{c\perp} &= \text{id} - u_s \left(\sum c u u' \right)^- u^s * c = [\sqrt{c}]^{-1} q * \sqrt{c} \\ \dot{\delta}(c)_{i,j} &= \left[\mathbf{1}_i^j - u'_j \widehat{\beta}_c(\mathbf{1}_i) \right]^2 = \left[\text{proj}_{\text{Im}(u_s)}^{c\perp} [\mathbf{1}_i]_j \right]^2 \\ &= \left\{ [\sqrt{c_i}]^{-1} q_{i,j} \sqrt{c_j} \right\}^2 = (c_i)^{-1} q_{i,j}^2 c_j \\ \Rightarrow \dot{\delta}(c) &= c^{-1} q^{*2} * c\end{aligned}\tag{eE.72}$$

- soit v_K une base orthonormée (BON) de $\text{Im}(q)$ alors $q = \sum_K v v'$ et :

$$\begin{aligned}h' c^{-1} q^{*2} * c h &= \sum_{i,j} h_i \frac{1}{c_i} \sum_{k,l} v_k^i v_k^j v_l^i v_l^j c_j h_j \\ &= \sum_{k,l} \left\langle h, \frac{v_k v_l}{c} \right\rangle \langle h, c v_k v_l \rangle\end{aligned}$$

\rightarrow Cette approche permet de prouver (avec $c = 1$) que la matrice des coefficients au carré d'une somme de vecteurs colonnes multipliés par leurs transposés, c'est-à-dire d'une matrice symétrique positive quelconque, est positive, puisqu'il apparait ainsi une somme de carrés ⁴⁰.

\rightarrow En outre, elle prouve qu'il n'est pas toujours vrai que $c^{-1} q^{*2} * c \geq 0$. En effet, si $\text{Im}(q) = \mathbb{R}v$ et si $(\frac{v^2}{c}, cv^2)$ ne sont pas colinéaires, alors il existe h tel que $\left\langle h, \frac{v^2}{c} \right\rangle \langle h, cv^2 \rangle < 0$. Or $(\frac{v^2}{c}, cv^2)$ ne

sont pas colinéaires s'il existe $i \neq j$ tel que $v_i^2 v_j^2 \left(\frac{c_j}{c_i} - \frac{c_i}{c_j} \right) \neq 0$, c'est-à-dire deux coordonnées de v non nulles pour lesquelles la pondération n'est pas identique. C'est possible même si c est très proche d'une constante, sauf si v a une seule coordonnée non nulle.

40. Plus généralement, si M et \widetilde{M} sont des matrices symétriques positives de même dimension alors $M * \widetilde{M} \geq 0$:

$$h' M * \widetilde{M} h = \sum_{i,j,k,l} h_i \sqrt{M_{i,k}} \sqrt{M_{j,k}} \sqrt{\widetilde{M}_{i,l}} \sqrt{\widetilde{M}_{j,l}} h_j = \sum_{k,l} \left\langle h, \sqrt{M}^k \sqrt{\widetilde{M}}^l \right\rangle^2$$

- majoration de q^{*2} :

$$\begin{aligned} \sum_{i,j} (q_{i,j})^2 h_i h_j &= \sum_{i,j,k,l} h_i h_j v_k^i v_k^j v_l^i v_l^j = \sum_{k,l} \langle h, v_k v_l \rangle^2 \\ &= \sum_{k,l} \langle h v_k, v_l \rangle^2 = \sum_{k \in K} (h v_k)' \sum_{l \in K} v_l v_l' (h v_k) = \sum_{k \in K} (h v_k)' q (h v_k) = \sum_{k \in K} (h v_k)' q^2 (h v_k) \\ &= \sum_K \|q(h v)\|^2 \leq \sum_K \|h v\|^2 = \sum_{k \in K} \sum_i (h_i v_k^i)^2 \\ &\leq \sum_i h_i^2 \sum_{k \in K} (v_k^i)^2 = \sum_i h_i^2 q_i \quad \text{d'où un majorant :} \end{aligned}$$

$$h' q^{*2} h \leq h' \text{diag}(q) h$$

⇒ Les deux résultats précédents mènent à un encadrement de la matrice des carrés des coefficients d'un projecteur orthogonal canonique :

$$0 \leq q^{*2} \leq \text{diag}(q) \leq \text{id} \quad (\text{eE.73})$$

E.3.2 rang de $\hat{\delta}(c)$

- Le rang de $\hat{\delta}(c)$ est constant sur $\{c > 0\}$:

$$\text{rang}[\hat{\delta}(c)] = \text{rang}[(q_1)^{*2}] = \dim \{ \text{Vect}[h^2, h \in \text{Im}(x_s)^{\perp}] \} \quad (\text{eE.74})$$

↳ Comme la matrice diagonale c est bijective, $\text{rang}(c^{-1} q^{*2} * c) = \text{rang}(q^{*2})$

$$q^{*2} h = 0 \Leftrightarrow h' q^{*2} h = 0 \quad (\text{car } q^{*2} \text{ est symétrique et positive vu (eE.73)})$$

$$h' q^{*2} h = \sum_{k,l \in K} \langle h, v_k v_l \rangle^2 = 0 \quad (\text{où } v_K \text{ est toujours une BON de } \text{Im}(q))$$

$$\Leftrightarrow \forall k, l \in K, \langle h, v_k v_l \rangle = 0 \Leftrightarrow \forall x, y \in \text{Im}(q), \langle h, xy \rangle = 0 \quad (\text{car } v_K \text{ base de } \text{Im}(q))$$

$$\Rightarrow \text{Ker}(q^{*2}) = \{xy, x, y \in \text{Im}(q)\}^{\perp}$$

$$\Rightarrow \text{Im}(q^{*2}) = \text{Vect}[xy, x, y \in \text{Im}(q)] \quad (\text{car } q^{*2} \text{ est symétrique donc } \text{Im}(q^{*2}) = \text{Ker}(q^{*2})^{\perp})$$

$$\Rightarrow \text{Im}(q^{*2}) = \text{Vect}[x^2, x \in \text{Im}(q)] \quad \text{car } xy = ((x+y)^2 - x^2 - y^2) / 2$$

$$\Rightarrow \text{rang}(q^{*2}) = \dim \{ \text{Vect}[h^2, h \in \text{Im}(q)] \} \quad \text{où } \text{Im}(q) = \text{Im}[\sqrt{c} u_s]^{\perp}$$

↳ $(q_c)^{*2}$ et $(q_1)^{*2}$ ont même rang :

$$h \in \text{Im}[\sqrt{c} u_s]^{\perp} \Leftrightarrow u_s' * \sqrt{c} h = 0$$

$$\Leftrightarrow \sqrt{c} h \in \text{Im}(u_s)^{\perp}$$

$$\{h^2, h \in \text{Im}[\sqrt{c} u_s]^{\perp}\} = \{[\sqrt{c} h]^2, h \in \text{Im}(u_s)^{\perp}\} = \{c h^2, h \in \text{Im}(u_s)^{\perp}\}$$

$$= c \{h^2, h \in \text{Im}(u_s)^{\perp}\}$$

$$\Rightarrow \text{Vect}\{h^2, h \in \text{Im}[\sqrt{c} u_s]^{\perp}\} = c \text{Vect}\{h^2, h \in \text{Im}(u_s)^{\perp}\}$$

$$\Rightarrow \text{Im}[(q_c)^{*2}] = c \text{Im}[(q_1)^{*2}] \quad \text{qui sont de mêmes dimensions car } c > 0$$

↳ $\dim \{ \text{Vect}[h^2, h \in \text{Im}(u_s)^{\perp}] \} = \dim \{ \text{Vect}[h^2, h \in \text{Im}(x_s)^{\perp}] \}$:

$$\text{Im}(u_s)^{\perp} = \left\{ h / h' \frac{x}{\pi}(s) = 0 \right\} = \left\{ \pi_s \frac{h}{\pi} / \frac{h'}{\pi} x_s = 0 \right\} = \pi_s \text{Im}(x_s)^{\perp}$$

$$\Rightarrow [\text{Im}(u_s)^{\perp}]^{*2} = (\pi_s)^{*2} [\text{Im}(x_s)^{\perp}]^{*2}$$

$$\Rightarrow \text{Vect} \{ [\text{Im}(u_s)^{\perp}]^{*2} \} = (\pi_s)^{*2} \text{Vect} \{ [\text{Im}(x_s)^{\perp}]^{*2} \}$$

• détermination du rang de $\hat{\delta}(c) = \dim \{ \text{Vect}[h^2, h \in \text{Im}(u_s)^{\perp}] \}$:

– Il existe deux sous-ensembles K d'indices colonnes et $s_K \subset s$ de cardinaux égaux au rang de u_s et tels que $u_{s_K}^K$ soit inversible et $\exists B / u_s = u_s^K B$. Dans la suite, Δ_s désigne le complémentaire de s_K dans s .

– Il est possible de paramétrer les éléments de $\text{Im}(u_s)^\perp$ par leurs coordonnées sur le sous-ensemble des indices Δs : soit $x \in \text{Im}(u_s)^\perp$

$$\begin{aligned} x'u_s &= 0 \Leftrightarrow x'u_s^K = 0 \text{ car } u_s = u_s^K B \\ \Leftrightarrow x'_{s_K} u_{s_K}^K + x'_{\Delta s} u_{\Delta s}^K &= 0 \\ \Leftrightarrow x_{s_K} &= -[u_{s_K}^{s_K}]^{-1} u_{s_K}^{\Delta s} x_{\Delta s} \\ \Rightarrow x &= \begin{pmatrix} \text{id}_{\Delta s} \\ A \end{pmatrix} x_{\Delta s} \text{ où } A = -[u_{s_K}^{s_K}]^{-1} u_{s_K}^{\Delta s} \end{aligned}$$

– détermination de l'orthogonal de $\text{Vect}[x^2, x \in \text{Im}(u_s)^\perp]$:

$$\begin{aligned} h'x^2 &= h' \left[\begin{pmatrix} x_{\Delta s} \\ Ax_{\Delta s} \end{pmatrix} \right]^2 \\ &= \sum_{j \in \Delta s} h_j (x_j)^2 + \sum_{i \in s_K} h_i \left[\sum_{j \in \Delta s} A_{i,j} x_j \right]^2 \\ &= \sum_{j \in \Delta s} (x_j)^2 \left(h_j + \sum_{i \in s_K} h_i (A_{i,j})^2 \right) + \sum_{i \in s_K} h_i \sum_{j \neq \bar{j}} A_{i,j} A_{i,\bar{j}} x_j x_{\bar{j}} \end{aligned}$$

La condition pour que ceci soit nul pour tout $x_{\Delta s}$ est :

$$\begin{cases} \forall j \in \Delta s, h_j = - \sum_{i \in s_K} h_i (A_{i,j})^2 \\ \forall j_1 \neq j_2 \in \Delta s, \sum_{i \in s_K} h_i A_{i,j_1} A_{i,j_2} = 0 \end{cases}$$

$$\Rightarrow \dim \{ \{x^2, x \in \text{Im}(u_s)^\perp\}^\perp \} \leq \text{rang}(u_s) (= |s_K|)$$

$$\Rightarrow \text{rang}[\dot{\delta}(c)] \geq |s| - \text{rang}(u_s)$$

⇒ Le rang de $\dot{\delta}(c)$ est déterminé par celui de la famille de (produits de) vecteurs $(A^{j_1} A^{j_2})_{j_1 < j_2 \in \Delta s}$.

• Si $\forall j_1 \in \Delta s, u_{\Delta s \neq j_1}^K$ est de rang = $|K|$ alors $\text{rang}[\dot{\delta}(c)] = |s|$ et $\dot{\delta}(c)$ est inversible. En effet :

→ $\text{rang}(A^{\Delta s \neq j_1}) = \text{rang}[u_{s_K}^{\Delta s \neq j_1}]$, car la première matrice est le produit de la seconde par une matrice inversible.

$$\rightarrow \text{si } \forall j_1 \in \Delta s, \text{rang}(A^{\Delta s \neq j_1}) = |K| = |s_K| \text{ alors } \sum_{i \in s_K} h_i A_{i,j_1} A_i^{\Delta s \neq j_1} = 0 \Rightarrow \sum_{i \in s_K} h_i A_{i,j_1} = 0$$

$(\forall j_1 \in \Delta s) \Rightarrow \sum_{i \in s_K} h_i A_i = 0 \Rightarrow h_{s_K} = 0$ (En particulier, pour une seule contrainte d'équilibrage, δ est un

difféomorphisme local si toutes les coordonnées de u_s sont non nulles pour tout échantillon.)

E.3.3 conditions pour que $\text{id} - \dot{\delta}(c)$ soit contractante

• si c est constante alors $\text{id} - \dot{\delta}(c) = \text{id} - \dot{\delta}(1)$ est contractante pour la norme euclidienne :

– $\dot{\delta}(c) = c^{-1} q^{*2} * c = q^{*2}$: En effet, comme c est une constante, elle peut être factorisée.

– $q_c = q_1$ est évident car $\text{Im}[\sqrt{c} u_s] = \text{Im}(u_s)$ si c constante > 0

– $\text{id} - \dot{\delta}(c)$ contractante $\Leftrightarrow \forall h, 2h' \dot{\delta}(c) h \geq \|\dot{\delta}(c) h\|^2 \Leftrightarrow 2\dot{\delta}(c) \geq \dot{\delta}(c)' \dot{\delta}(c)$ (par définition) donc :

$$\text{id} - \dot{\delta}(c) \text{ contractante} \Leftrightarrow 2c^{-1} q^{*2} * c \geq c q^{*2} * c^{-2} q^{*2} * c$$

→ Si c est constante, la contraction revient donc à : $2q^{*2} \geq q^{*2} q^{*2}$. Cette propriété est vérifiée, parce que (vu (eE.73)) $q^{*2} \leq \text{id} \Rightarrow q^{*2} q^{*2} \leq q^{*2}$. Cette implication découle de ce que q^{*2} est symétrique et positive, donc admet une racine carrée.

– Plus précisément, $\text{id} - \dot{\delta}(1)$ est strictement contractante sur $\text{Ker}[\dot{\delta}(1)]^\perp$: Comme $\dot{\delta}(1) = q^{*2}$, qui est symétrique et positive, si $q^{*2} h \neq 0$ alors $h' q^{*2} h \neq 0$ et donc $h' q^{*2} q^{*2} h \leq h' q^{*2} h < 2h' q^{*2} h$.

• si c_0 est constante alors il existe un voisinage $\mathcal{U}^0(c_0)$ tel que $\text{id} - \dot{\delta}(c)$ soit contractante 'sur la majeure partie de E ' (en un sens limité apparaissant ci-dessous) :

– Comme $\text{id} - \dot{\delta}(1)$ est strictement contractante sur $\text{Ker} [\dot{\delta}(1)]^\perp$ et que $\{h \in \text{Ker} [\dot{\delta}(1)]^\perp, \|h\| = 1\}$ est compact :

$$\sup_{\left\{ \begin{array}{l} h \in \text{Ker} [\dot{\delta}(1)]^\perp \\ \|h\| = 1 \end{array} \right.} \|(\text{id} - \dot{\delta}(1))h\| = \eta < 1$$

– Par la continuité de la fonction $\text{id} - \dot{\delta}(c)$ sur $\{c > 0\}$, pour $\eta < \nu < 1$ il existe $\mathcal{O}^\eta(c_0)$ tel que :

$$\sup_{\left\{ \begin{array}{l} h \in \text{Ker} [\dot{\delta}(1)]^\perp \\ \|h\| = 1 \end{array} \right.} \|(\text{id} - \dot{\delta}(c))h\| \leq \nu < 1$$

– $\text{Ker} [\dot{\delta}(1)]^\perp \cap \text{Ker} [\dot{\delta}(c)] = 0$, puisque pour $h \neq 0$ dans $\text{Ker} [\dot{\delta}(1)]^\perp$, $\|(\text{id} - \dot{\delta}(c))h\| \leq \nu \|h\|$, donc $\dot{\delta}(c)h = 0$ impossible.

– Comme $\dim \{\text{Ker} [\dot{\delta}(c)]\} = \dim \{\text{Ker} [\dot{\delta}(1)]\}$, $E = \text{Ker} [\dot{\delta}(1)]^\perp \oplus \text{Ker} [\dot{\delta}(c)]$.

• En général, la condition $\dot{\delta}(c) \geq 0$ (nécessaire pour une contraction) a peu de chances d'être vérifiée (sur tout l'espace E) lorsque c n'est pas constante et $\dot{\delta}(c)$ n'est pas inversible ($\Rightarrow \text{Ker} (q^{*2}) \neq 0$) :

$$h' \dot{\delta}(c) h = h' c^{-1} q^{*2} * ch$$

$$= (k_1 + k_2)' c^{-2} q^{*2} k_1 \text{ où } ch = k_1 + k_2 \text{ et } k_1 \in \text{Im} (q^{*2}), k_2 \in \text{Im} (q^{*2})^\perp$$

$$\forall h, h' \dot{\delta}(c) h \geq 0 \Leftrightarrow \forall k_1 \in \text{Im} (q^{*2}), k_2 \in \text{Im} (q^{*2})^\perp, k_1' c^{-2} q^{*2} k_1 + k_2' c^{-2} q^{*2} k_1 \geq 0$$

$$\Rightarrow \forall k_1 \in \text{Im} (q^{*2}), k_2 \in \text{Im} (q^{*2})^\perp, k_2' c^{-2} q^{*2} k_1 = 0 \text{ (sinon ce terme pourrait être dilaté vers } -\infty \text{ via } k_2)$$

$$\Rightarrow \forall k_1 \in \text{Im} (q^{*2}), k_2 \in \text{Im} (q^{*2})^\perp, k_2' c^{-2} k_1 = 0 \text{ (parce que } q^{*2} \text{Im} (q^{*2}) = \text{Im} (q^{*2}))$$

Cette dernière condition de préservation de l'orthogonalité est très forte. Elle est vérifiée ssi $\text{Im} (q^{*2})$ est stable par c^{-2} . En effet, si f est une application linéaire et p un projecteur orthogonal canonique, f conserve l'orthogonalité entre $\text{Im} (p)$ et $\text{Im} (p)^\perp$ ssi $\forall x, y, [f(\text{id} - p)x]' fpy = 0 \Leftrightarrow (\text{id} - p)f' f p = 0$. Si $u \in \text{Im} (p)$ et $f' f u \notin \text{Im} (p)$ alors $(\text{id} - p)f' f u \neq 0$ donc $(\text{id} - p)f' f p u \neq 0$. Donc $[f \text{Im} (p)]^\perp \perp [f \text{Im} (\text{id} - p)] \Leftrightarrow (\text{id} - p)f' f p = 0 \Leftrightarrow f' f \text{Im} (p) \subset \text{Im} (p)$.

→ Pour expliquer l'efficacité de la minimisation par l'algorithme récursif, il est sans doute nécessaire de recourir à une propriété moins forte que la contraction 'pure'.

• En conclusion de cette étude théorique de la contraction de $\text{id} - \dot{\delta}(c)$, celle-ci est contractante lorsque c est constante et 'généralement contractante' pour c assez proche d'une constante. Il n'a pas été possible de mettre en évidence plus clairement la raison de l'efficacité de la minimisation observée pour l'algorithme récursif.

→ Empiriquement, pour valider et compléter l'étude de la contraction de $\text{id} - \dot{\delta}$, un calcul a été effectué de la plus grande valeur propre de :

$$[\text{id} - \dot{\delta}(c)]' (\text{id} - \dot{\delta}(c)) = \left[\text{id} - (\text{id} - (cu)_s (\sum c u u')^{-1} u^s)^{*2} \right]' \left[\text{id} - (\text{id} - (cu)_s (\sum c u u')^{-1} u^s)^{*2} \right]$$

→ Sur les 10 000 échantillons simulés, les pondérations $c = 1$ et $c = 1 - \pi$ donnent pour cette matrice des valeurs propres maximales de l'ordre de 1^{41} , à $4 \cdot 10^{-6}$ près.

→ En revanche pour la pondération $c = \frac{1 - \pi}{\delta(1)}$ (pour les unités telles que $\delta_i(1) > 10^{-6}$, et $c_i = 1 - \pi_i$ sinon), la valeur propre maximale est plus élevée, sans toutefois dépasser à 1.05 (à l'exception de 315 des 10 000 simulations).

→ Sur le premier échantillon, lorsqu'un poids de 2 est attribué à la moitié de chaque échantillon régional, et 1 pour le reste, la valeur propre maximale dépasse sensiblement 1 pour 4 régions, dont Pays de la Loire et PACA.

⇒ Sur les données du sondage étudié, la fonction ξ n'est plus contractante lorsque la pondération est trop dispersée, en fonction de la taille de l'échantillon et de la spécificité régionale. Ce résultat est conforme aux attentes théoriques. Cependant, la condition de contraction dans toutes les directions n'est pas nécessaire pour l'efficacité relative de l'algorithme récursif.

41. La valeur 1 est atteinte dès que $\text{Ker} [\dot{\delta}(c)] \neq 0$.

Annexe F compléments sur l'algorithme récursif

• Tant que $\delta(c) \leq 1 - \pi$ et $\delta[\xi(c)] \leq 1 - \pi + 1 - \pi - \delta(c)$ et qu'il existe $i \in s$ tel que ces deux inégalités soient strictes et $\delta_i(1) > 0$, la condition de poursuite de l'algorithme est vérifiée. En effet, ces conditions impliquent que $\xi(c) \geq c$ et $\delta(c) \leq \delta[\xi(c)] \leq 1 - \pi + 1 - \pi - \delta(c)$ (inégalités vectorielles) donc $-(1 - \pi - \delta(c)) \leq \delta[\xi(c)] - (1 - \pi) \leq 1 - \pi - \delta(c)$ et ces deux inégalités sont strictes pour i ⁴², ce qui implique que $\|1 - \pi - \delta[\xi(c)]\| < \|1 - \pi - \delta(c)\|$ (pour la norme euclidienne).

⇒ Si $\exists i \in s / (1 - \pi_i)u_i \neq 0$ et $\delta_i(1) > 0$ alors la première étape de l'algorithme récursif fonctionne toujours, c'est-à-dire que :

$$\|\delta[\xi(1 - \pi)] - (1 - \pi)\| < \|\delta(1 - \pi) - (1 - \pi)\| \quad (\mathbf{eF.75})$$

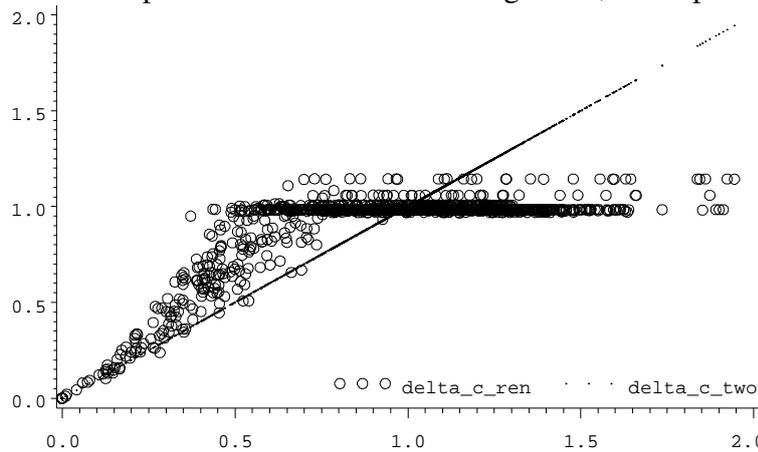
↳ Comme $\forall c > 0, \delta(c) \leq c, \delta(1 - \pi) \leq 1 - \pi$, et cette inégalité est stricte pour la coordonnée i car $(1 - \pi_i)u_i \neq 0$.

↳ de même $\delta[\xi(1 - \pi)] \leq \xi(1 - \pi) = 1 - \pi + 1 - \pi - \delta(1 - \pi)$. L'inégalité est stricte pour i car $\xi_i(1 - \pi)u_i \neq 0$ vu que $u_i \neq 0$ et $\xi_i(1 - \pi) \neq 0$ car $1 - \pi_i > 0$ d'une part et $\delta_i(1 - \pi) < 1 - \pi_i$ d'autre part.

⇒ La pondération récursive est préférable à la pondération par $1 - \pi$ selon la distance inter-diagonales, à la seule condition qu'il existe $i \in s / (1 - \pi_i)u_i \neq 0$ et $\delta_i(1) > 0$.

• Le **Graphique 3** montre que la pondération récursive comparée au calage sur la trace de l'ESB donne plus de poids aux unités 'atypiques', celles de plus faibles coefficients diagonaux. D'autre part le maximum du coefficient diagonal sur l'ensemble des simulations est de 1.70 pour la pondération récursive et la médiane de 0.985, contre 3.05 et 1.014.

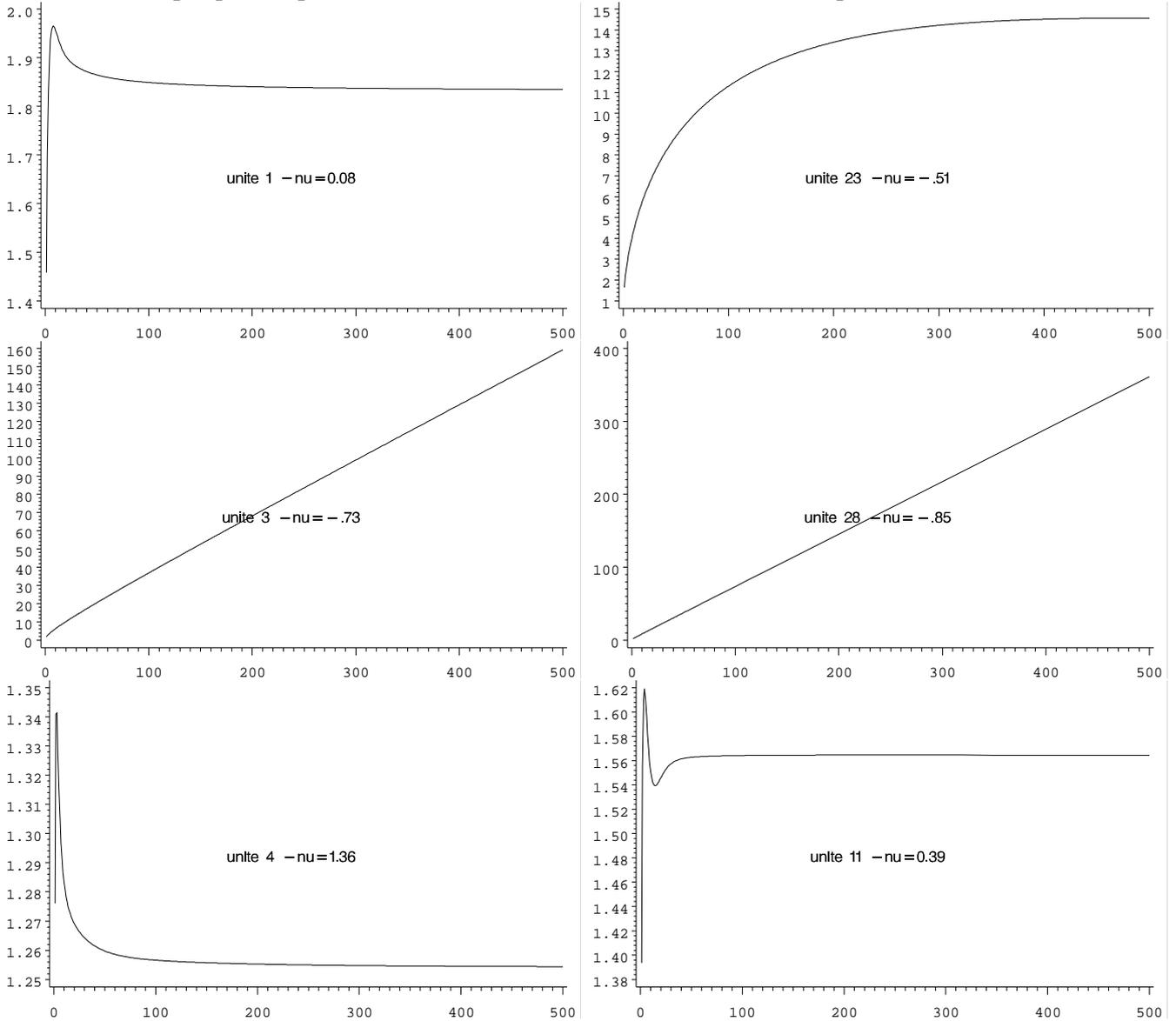
Graphique 3 – Comparaison des coefficients diagonaux, sur le premier échantillon



• Sur l'exemple des 58 unités d'une région (Poitou-Charentes), l'évolution des $\xi_i(\mathbb{N})$ paraît très régulière. Trois profils peuvent être distingués (**Graphique 4**) : stabilisation rapide, croissance puis stabilisation progressive et enfin croissance régulière.

42. en utilisant :
$$\begin{cases} \xi(c) \geq c \\ \xi_i(c) < c_i \Rightarrow \delta_i(c) < \delta_i[\xi(c)], \text{ vu } (\mathbf{e2.24}) \\ \delta_i(1) > 0 \end{cases}$$

Graphique 4 – pondération récursive en fonction de l'itération, pour une unité



Note : La pondération est produite par l'algorithme récursif sur le premier échantillon simulé, pour la région Poitou-Charentes. Les unités représentées sont représentatives des différents profils d'évolution de la pondération selon l'itération. Ils se distinguent par le paramètre $\nu_i = \min_{u_i^* \beta = 1} \|u^{*i} \beta\|_{1-\pi_i}^2 - (1 - \pi_i)$, ainsi que par $\delta_i (1)$.

Annexe G notations

- Le complémentaire d'un sous-ensemble $A \subset \Omega$ est noté $A^c = \Omega \setminus A$.
- L'ensemble des applications de l'ensemble A dans B est noté B^A ou B_A selon la représentation de ses éléments, respectivement en ligne (donc indicée en colonne) ou en vecteur-colonne
- Pour une matrice A , la transposée est notée A' , A^- est une inverse généralisée et A^+ est l'inverse généralisée de Moore et Penrose (??).
- $\text{Im}(A)$ est l'image d'une matrice A ($\text{Im}(A) = AE$, où E est l'espace de départ), $\text{Ker}(A)$ est son noyau ($\text{Ker}(A) = \{A = 0\}$), $|A|$ son déterminant et $\text{rang}(A) = \dim[\text{Im}(A)]$
- x_s ou $x(s)$ désigne une matrice indicée en ligne par l'ensemble s . Celle-ci est confondue avec la famille de vecteurs lignes $((x'_i)_{i \in s})$ et avec l'application linéaire $(\beta \mapsto (x'_i \beta)_{i \in s})$. Si nécessaire, l'ensemble des indices des colonnes est précisé sous la forme $x_s(K)$.

- x^s désigne une matrice indicée en colonne par l'ensemble s (donc $x^s = (x_s)'$) (Une entorse à la logique de la notation des indices est faite dans le cas d'un indice unique : u_i ou $u(i)$ désignera un vecteur colonne et u'_i un vecteur ligne, pour éviter la confusion possible avec la puissance i .)
- Les indices sont omis lorsque le contexte n'est pas ambigu.
- Pour un ensemble B :
 - $\text{Vect}(B)$ est le plus petit sous-espace vectoriel qui contient B
 - $B^\perp = \{x \in E / \langle x, B \rangle = 0\} = \{x \in E / \forall y \in B, \langle x, y \rangle = 0\}$ est le sous-espace vectoriel orthogonal à B
- s est l'échantillon tiré dans l'univers \mathcal{P} . Alternativement, en fonction du contexte, cette notation désigne l'indicatrice de présence dans l'échantillon ($s \in \{0, 1\}^{\mathcal{P}}$).
- $|s|$ est le cardinal de l'ensemble s , soit ici la taille de l'échantillon
- $\mathbb{1}_i$ est le vecteur de i -ième coordonnée égale à 1 et dont les autres coordonnées sont nulles. ($\mathbb{1}_s$ est donc la matrice formée de ces indicatrices, à ne pas confondre avec le vecteur de coordonnées égales à 1, noté $\mathbb{1}_s$)
- L'espace vectoriel d'intérêt pour l'échantillon ($\ni y(s)$) est noté $E = \mathbb{R}_s$.
- Pour deux vecteurs x et y du même espace vectoriel \mathbb{R}_s :
 - $x \leq y$ si $\forall i, x_i \leq y_i$
 - $x < y$ si $\forall i, x_i < y_i$ (en particulier, $c > 0$ signifie : $\forall i \in s, c_i > 0$)
- Le produit de deux vecteurs est le vecteur constitué des produits de leurs coordonnées : $(xy)_i = x_i y_i$ et en particulier $x^2 = ((x_i)^2)_{i \in s}$. Parallèlement, la matrice de diagonale x est confondue avec le vecteur x .
- Plus généralement deux types de produits termes à termes de matrices sont à distinguer :
 - produit de deux matrices A et B de mêmes dimensions : $(A * B)_{i,j} = A_{i,j} B_{i,j}$ (Cette notation vise à éviter la confusion avec le produit matriciel usuel.)
 - produit d'un vecteur par une matrice :
 - en lignes : $(xA)_{i,j} = x_i A_{i,j}$
 - en colonnes : $(A * x)_{i,j} = A_{i,j} x_j$
- Il faut faire attention lorsqu'une expression combine des produits matriciels usuels avec des produits termes à termes. Par exemple, il n'est pas vrai en général que $(AB) * C = A(B * C)$.
- A^{*2} désigne la matrice dont les coordonnées sont les carrés de celles de la matrice A .
- Pour une variable réelle, $y : s \mapsto \mathbb{R}$, soit $y \in \mathbb{R}_s$, ou vectorielle : $\sum y = Y = y(+)$ est le total sur l'univers, $\sum_s y = \sum_{i \in s} y_i$ est le total sur l'échantillon, y_d est le total sur une unité de d -ième phase
- $\|x\|_M$ est la norme euclidienne d'un vecteur x pour la métrique M (ie une matrice symétrique, définie positive), soit $\sqrt{x' M x}$. $\langle \cdot, \cdot \rangle_M$ est son produit scalaire (La mention de la métrique est omise si c'est l'identité.) En particulier, pour un vecteur $c (> 0)$, $\|x\|_c$ est la norme pour la métrique diagonale définie par c , ce qui correspond simplement à la somme pondérée : $\|x\|_c^2 = \sum c x^2$. Cette dernière formule permet d'étendre la définition de $\|x\|_c^2$ au cas où $c \in \mathbb{R}_s$.
- proj_F^M est la projection orthogonale pour la métrique M sur le sous-espace-vectoriel F . (La mention de $M (= \text{id})$ est encore omise pour la métrique euclidienne canonique.) $\text{proj}_F^{\perp M} = \text{proj}_{F^\perp}^M$ est la projection sur l'orthogonal de F , soit : $\text{proj}_F^{\perp M} = \text{id} - \text{proj}_F^M$.
- m est la mesure de comptage réduite sur l'univers \mathcal{P} : $m(y) = \frac{Y}{N}$, où $N = |\mathcal{P}|$