

Techniques d'enquête

Modèles bayésiens pour petits domaines sous contraintes d'inégalité avec réconciliation et rétrécissement double

par Balgobin Nandram, Nathan B. Cruze et Andreea L. Erciulescu

Date de diffusion : le 3 janvier 2024



Statistique
Canada

Statistics
Canada

Canada

Comment obtenir d'autres renseignements

Pour toute demande de renseignements au sujet de ce produit ou sur l'ensemble des données et des services de Statistique Canada, visiter notre site Web à www.statcan.gc.ca.

Vous pouvez également communiquer avec nous par :

Courriel à infostats@statcan.gc.ca

Téléphone entre 8 h 30 et 16 h 30 du lundi au vendredi aux numéros suivants :

- | | |
|---|----------------|
| • Service de renseignements statistiques | 1-800-263-1136 |
| • Service national d'appareils de télécommunications pour les malentendants | 1-800-363-7629 |
| • Télécopieur | 1-514-283-9350 |

Normes de service à la clientèle

Statistique Canada s'engage à fournir à ses clients des services rapides, fiables et courtois. À cet égard, notre organisme s'est doté de normes de service à la clientèle que les employés observent. Pour obtenir une copie de ces normes de service, veuillez communiquer avec Statistique Canada au numéro sans frais 1-800-263-1136. Les normes de service sont aussi publiées sur le site www.statcan.gc.ca sous « Contactez-nous » > « [Normes de service à la clientèle](#) ».

Note de reconnaissance

Le succès du système statistique du Canada repose sur un partenariat bien établi entre Statistique Canada et la population du Canada, les entreprises, les administrations et les autres organismes. Sans cette collaboration et cette bonne volonté, il serait impossible de produire des statistiques exactes et actuelles.

Publication autorisée par le ministre responsable de Statistique Canada

© Sa Majesté le Roi du chef du Canada, représenté par le ministre de l'Industrie, 2024

Tous droits réservés. L'utilisation de la présente publication est assujettie aux modalités de l'[entente de licence ouverte](#) de Statistique Canada.

Une [version HTML](#) est aussi disponible.

This publication is also available in English.

Modèles bayésiens pour petits domaines sous contraintes d'inégalité avec réconciliation et rétrécissement double

Balgobin Nandram, Nathan B. Cruze et Andreea L. Erciulescu¹

Résumé

Nous présentons une nouvelle méthodologie pour réconcilier des estimations des totaux des superficies cultivées au niveau du comté à un total prédéfini au niveau de l'État soumis à des contraintes d'inégalité et à des variances aléatoires dans le modèle de Fay-Herriot. Pour la superficie ensemencée du National Agricultural Statistics Service (NASS), un organisme du ministère de l'Agriculture des États-Unis (USDA), il est nécessaire d'intégrer la contrainte selon laquelle les totaux estimés, dérivés de données d'enquête et d'autres données auxiliaires, ne sont pas inférieurs aux totaux administratifs de la superficie ensemencée préenregistrés par d'autres organismes du USDA, à l'exception de NASS. Ces totaux administratifs sont considérés comme fixes et connus, et cette exigence de cohérence supplémentaire ajoute à la complexité de la réconciliation des estimations au niveau du comté. Une analyse entièrement bayésienne du modèle de Fay-Herriot offre un moyen intéressant d'intégrer les contraintes d'inégalité et de réconciliation et de quantifier les incertitudes qui en résultent, mais l'échantillonnage à partir des densités *a posteriori* comprend une intégration difficile; des approximations raisonnables doivent être faites. Tout d'abord, nous décrivons un modèle à rétrécissement unique, qui rétrécit les moyennes lorsque l'on suppose que les variances sont connues. Ensuite, nous élargissons ce modèle pour tenir compte du rétrécissement double par l'emprunt d'information dans les moyennes et les variances. Ce modèle élargi comporte deux sources de variation supplémentaire; toutefois, comme nous rétrécissons à la fois les moyennes et les variances, ce second modèle devrait avoir un meilleur rendement sur le plan de la qualité de l'ajustement (fiabilité) et, possiblement, sur le plan de la précision. Les calculs sont difficiles pour les deux modèles, qui sont appliqués à des ensembles de données simulées dont les propriétés ressemblent à celles des cultures de maïs de l'Illinois.

Mots-clés : Méthode de Devroye; modèle de Fay-Herriot; méthode de grille; modèle hiérarchique bayésien; échantillonneur de Metropolis.

1. Introduction

Pour de nombreux problèmes dans les statistiques officielles, il est nécessaire d'intégrer des contraintes dans l'inférence fondée sur un modèle. Par exemple, dans le cas des estimations pour petits domaines, il peut y avoir des contraintes sur les estimations du modèle, qui doivent être réconciliées à une cible. Il peut s'agir de bornes inférieures (ou supérieures) connues pour les estimations au niveau du comté, qui devraient « s'additionner » pour correspondre à l'estimation au niveau de l'État obtenue précédemment. Un exemple pratique est l'estimation de la superficie ensemencée pour les comtés dans les États, avec une estimation au niveau de l'État obtenue précédemment, lorsqu'il existe des données d'enquête et des données administratives qui peuvent fournir des bornes inférieures aux estimations pour les comtés. Celles-ci doivent être additionnées pour correspondre à l'estimation au niveau de l'État. Même si nous nous concentrons sur une application en agriculture, nous élaborons une méthodologie visant à résoudre le problème selon lequel des estimations pour petits domaines sont nécessaires pour satisfaire à certaines bornes inférieures, et ces estimations sont ensuite réconciliées à une estimation à un niveau plus élevé au moyen de l'approche descendante.

1. Balgobin Nandram, Department of Mathematical Sciences, Worcester Polytechnic Institute, 100 Institute Road, Worcester, MA 01609. Courriels : balnan@wpi.edu et balgobin.nandram@usda.gov; Nathan B. Cruze, NASA Langley Research Center, Mail Stop 290, Hampton, VA 23681. Courriel : nathan.b.cruze@nasa.gov; Andreea L. Erciulescu, Westat, 1600 Research Blvd., Rockville, MD 20850, États-Unis. Courriel : AndreeaErciulescu@westat.com.

Aux États-Unis, les estimations officielles au niveau du comté du rendement des cultures, de la production totale et de la superficie totale publiées par National Agricultural Statistics Service (NASS), un organisme du ministère de l'Agriculture des États-Unis (USDA), sont importantes. Ces estimations officielles permettent de déterminer le montant des paiements à verser aux agriculteurs et aux éleveurs inscrits à plusieurs programmes administrés par d'autres organismes du USDA, notamment Farm Service Agency (FSA) et Risk Management Agency (RMA). En conséquence, NASS s'efforce d'accroître la précision, la fiabilité et la couverture de ses estimations officielles des cultures au niveau du comté. Comme il est décrit dans un rapport intitulé *Improving Crop County Estimates by Integrating Multiple Data Sources* (National Academies of Sciences, Engineering, and Medicine, 2017), l'un des moyens d'y parvenir est d'utiliser des modèles soutenables qui comprennent de multiples sources de variabilité et d'autres données auxiliaires. Le rapport met en évidence un grand nombre de défis auxquels NASS est confronté et le rôle que peut jouer l'inférence fondée sur un modèle dans la publication d'estimations officielles au niveau du comté. Les conclusions du rapport ont été examinées plus en détail dans Cruze, Erciulescu, Nandram, Barboza et Young (2019). Les auteurs ont déterminé que la correspondance des estimations des superficies cultivées avec les totaux administratifs connus des superficies cultivées de la même année est jugée essentielle pour le programme d'estimations des cultures au niveau du comté de NASS.

Les contraintes sur les estimations peuvent se présenter sous la forme de restrictions d'ordre ou de forme (par exemple, Nandram, Sedransk et Smith, 1997; Silvapulle et Sen, 2005; Chen et Nandram, 2022) ou sous la forme de contraintes d'inégalité (Sen et Silvapulle, 2002). Ce second type de restriction est d'un intérêt particulier dans la mesure où il est lié à la cohérence des estimations des cultures totalisées en présence de données administratives disponibles conservées par le USDA. La réconciliation des estimations pour de plus petits domaines géographiques par rapport à celles pour de plus grandes régions géographiques est une forme courante de contrainte d'inégalité présente dans les statistiques officielles. Par exemple, plusieurs études antérieures de NASS y sont parvenues en effectuant un ajustement du ratio (ratisage) après l'analyse des résultats du modèle (par exemple, Erciulescu, Cruze et Nandram, 2018, 2019 et 2020); voir également Steorts, Schmid et Tzavidis (2020) et la bibliographie qui y figure pour obtenir un examen instructif de la réconciliation. Bien que le présent article mette l'accent sur la méthodologie, nous abordons la récente étude de cas et le document d'accompagnement rédigés par NASS (Chen, Nandram et Cruze, 2022) sur le problème de superficie ensemencée sous contrainte selon le modèle à rétrécissement unique. Nous précisons également que, dans le document actuel, nos principales contributions portent sur les contraintes d'inégalité; voir également le rapport de recherche n° RDD-22-02 de la division de la recherche et du développement de NASS (Nandram, Cruze, Erciulescu et Chen, 2022).

Les données non probabilistes ne sont pas exemptes d'erreurs. Tout d'abord, il est entendu que, même si la participation aux programmes de soutien à l'agriculture est populaire aux États-Unis, l'adhésion volontaire aux programmes de FSA et de RMA contribue à une sous-couverture potentielle (un biais vers le bas) dans ces totaux administratifs des superficies. En outre, les taux de participation à ces programmes de soutien peuvent varier d'une année à l'autre, selon la culture de base, selon l'État ou même plus

localement au sein d'un État. On estime toutefois que d'autres erreurs non due à l'échantillonnage sont réduites au minimum par les contrôles de la qualité de FSA et de RMA. Par exemple, les agriculteurs confirment la superficie de leurs terres inscrite aux programmes par des agents de FSA en fonction de limites des terres géolocalisées; les agriculteurs sont passibles de sanctions si leurs rapports sont falsifiés. Compte tenu de ces propriétés, NASS et le USDA considèrent les totaux administratifs disponibles comme des *bornes inférieures informatives*, et la publication de données totalisées cohérentes sur les superficies ensemencées exige : 1) que les totaux des superficies au niveau du comté s'additionnent aux totaux des superficies au niveau de l'État publiés avant la diffusion des estimations des comtés; et 2) que les estimations officielles des superficies ensemencées au niveau du comté respectent la contrainte de borne inférieure dans chaque comté.

En outre, nous prenons en considération les gains possibles résultant d'un rétrécissement double en empruntant simultanément de l'information dans les moyennes et les variances. Les techniques d'estimation des variances échantillonnales fondées sur des modèles fréquentistes et bayésiens ont été prises en considération dans la littérature pour les modèles au niveau du domaine. Par exemple, voir Wang et Fuller (2003); You et Chapman (2006); Gonzalez-Manteiga, Lombardia, Molina, Morales et SantaMaria (2010); Maiti, Ren et Sinha (2014); Dass, Maiti, Ren et Sinha (2012). Récemment, Erciulescu, Cruze et Nandram (2019) ont intégré le rétrécissement double dans les estimations des totaux des superficies récoltées non contraintes.

Soit $\hat{\theta}_i, i = 1, \dots, \ell$, les estimations directes observées de la superficie totale pour ℓ comtés, et $s_i^2, i = 1, \dots, \ell$, les variances observées correspondantes pour les ℓ comtés. Le modèle de Fay-Herriot au niveau du domaine (Fay et Herriot, 1979; voir également Rao et Molina, 2015) est un modèle standard dans l'estimation pour petits domaines de $\hat{\theta}_i$, où

$$\hat{\theta}_i \mid \theta_i \stackrel{\text{ind}}{\sim} \text{Normale}(\theta_i, s_i^2), i = 1, \dots, \ell, \quad (1.1)$$

et, au deuxième degré,

$$\theta_i \mid \boldsymbol{\beta}, \delta^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'_i \boldsymbol{\beta}, \delta^2), i = 1, \dots, \ell, \quad (1.2)$$

où \mathbf{x}_i est un vecteur p de covariables ayant une ordonnée à l'origine et $\boldsymbol{\beta}$ est un vecteur p de coefficients de régression. Dans une analyse entièrement bayésienne de ce modèle, on suppose des distributions *a priori* des paramètres du modèle; *a priori*, nous prenons $\pi(\boldsymbol{\beta}, \delta^2) = \pi(\boldsymbol{\beta}) \pi(\delta^2)$, où $\pi(\delta^2)$ est propre, mais $\pi(\boldsymbol{\beta}) = 1$ est impropre.

Sur le plan procédural, les estimations de NASS de la superficie ensemencée au niveau de l'État (désignons ces cibles de l'État par le scalaire a) sont déterminées et publiées avant la publication des estimations au niveau du comté. Nandram, Erciulescu et Cruze (2019) ont élaboré un modèle de Fay-Herriot entièrement bayésien qui intègre la contrainte de réconciliation $\sum_{i=1}^{\ell} \theta_i = a$ directement dans le modèle. Pour ce faire, le dernier domaine a été supprimé afin de tenir compte de la contrainte de réconciliation. Ils

ont montré de manière empirique que, dans la pratique, le choix du domaine supprimé pour intégrer la contrainte de réconciliation n'a pas beaucoup d'importance. Toutefois, il est plus pratique dans le présent document d'adopter une autre approche, qui ne fait pas appel à la suppression.

Nous souhaitons à présent améliorer ce modèle pour tenir compte des contraintes de réconciliation et d'inégalité sur θ_i . En plus de la contrainte de réconciliation, nous devons ajouter les contraintes d'inégalité propres à chaque comté

$$\theta_i \geq c_i, i=1, \dots, \ell, \quad (1.3)$$

où les c_i sont des quantités fixes connues qui représentent les valeurs administratives fournies par FSA ou RMA. (En pratique, lorsque les deux sources de données sont présentes, la plus grande des deux est utilisée pour établir la borne inférieure, à savoir c_i .) Dans les données de NASS sur les superficies ensemencées, certaines des estimations directes des totaux de la superficie ensemencée peuvent être inférieures de plus d'une ou de deux erreurs-types à leurs c_i correspondantes; par conséquent, il est plus difficile de faire en sorte que les estimations du modèle soient supérieures à c_i . Il convient de noter que $a = \sum_{i=1}^{\ell} \theta_i \geq \sum_{i=1}^{\ell} c_i \equiv c$. En d'autres termes, les processus d'estimation qui génèrent les cibles de l'État respectent aussi les totaux administratifs disponibles au *niveau de l'État*, mais la contrainte de réconciliation peut créer des difficultés supplémentaires lorsque la cible n'est que légèrement supérieure à c , c'est-à-dire lorsque $\frac{c}{a} \rightarrow 1$ par le bas. Nous devons ajouter les contraintes d'inégalité au modèle de Fay-Herriot précisé à l'équation (1.1), à l'équation (1.2) et dans les lois *a priori* pour obtenir la densité *a posteriori* conjointe de $\theta_i, i=1, \dots, \ell$. Afin d'intégrer les contraintes d'inégalité dans le modèle de Fay-Herriot bayésien, nous proposons la simplification suivante. En nous écartant de Nandram, Erciulescu et Cruze (2019), nous intégrons directement les contraintes d'inégalité tout en n'intégrant que partiellement la contrainte de réconciliation dans le modèle de Fay-Herriot bayésien. En d'autres termes, nous intégrons dans le modèle les contraintes, $c_i \leq \theta_i, i=1, \dots, \ell$, ainsi que la restriction selon laquelle $\sum_{i=1}^{\ell} \theta_i < a$. Lorsque cette seconde inégalité est appliquée, un ratissage des estimations du modèle pour qu'elles correspondent au total de l'État a dans une analyse des résultats satisfera toujours à toutes les contraintes d'inégalité de comtés individuels. L'intégration du rétrécissement double dans le modèle à contrainte d'inégalité entraîne d'autres éléments à considérer relativement aux calculs. Par conséquent, nos principales contributions consistent à fournir des estimations pour petits domaines, qui sont soumises à des contraintes d'inégalité et réconciliées à une cible, et nous décrivons un modèle à rétrécissement unique (variances d'échantillon fixes) et deux modèles à rétrécissement double (variances d'échantillon aléatoires).

Dans le présent document, nous examinons une nouvelle méthodologie pour résoudre ces problèmes doubles en modifiant le modèle de Fay-Herriot bayésien décrit dans Nandram, Erciulescu et Cruze (2019) afin de tenir compte des contraintes de réconciliation et d'inégalité dans les modèles bayésiens au niveau du domaine des équations (1.1) et (1.2). En outre, nous élargissons le modèle pour tenir compte du rétrécissement double des moyennes et des variances. À la section 2, nous présentons la méthodologie pour le modèle à rétrécissement unique en présence de totaux avec contrainte d'inégalité. À la section 3, nous décrivons la

méthodologie du modèle à rétrécissement double et de la régression gamma, et le modèle log-linéaire est examiné à l'annexe B. Là encore, les modèles à rétrécissement double intègrent des totaux avec contrainte d'inégalité. L'accent est mis sur les calculs qui facilitent ces approches. Dans la section 4, étant donné que la confidentialité des enquêtes et des données administratives du USDA est source de préoccupation, des ensembles de données simulées dont les propriétés ressemblent à celles de la culture du maïs de l'Illinois sont générés et utilisés pour ajuster et évaluer ces modèles. Nous présentons nos conclusions à la section 5 et constatons que les méthodologies de superficie avec contrainte ont été intégrées avec succès dans les statistiques officielles de NASS à partir de la campagne agricole de 2020.

2. Méthodologie dans le cadre du modèle à rétrécissement unique

Dans la présente section, nous élaborons les méthodologies et les stratégies computationnelles permettant d'intégrer les contraintes d'inégalité et les procédures de réconciliation dans les modèles bayésiens au niveau du domaine des équations (1.1) et (1.2). Cela permet d'obtenir le modèle à rétrécissement unique dans lequel on suppose que les variances d'échantillonnage sont fixes et connues.

Notre stratégie consiste à utiliser la règle de composition (c'est-à-dire la règle de multiplication des probabilités) pour prélever des échantillons de la densité *a posteriori* $\pi(\boldsymbol{\beta}, \delta^2 | \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\sigma}}^2)$, puis des échantillons de $\pi(\boldsymbol{\theta} | \boldsymbol{\beta}, \delta^2, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\sigma}}^2)$. Ces deux problèmes sont difficiles à résoudre. Dans la présente section, nous avons utilisé la loi *a priori* de rétrécissement pour δ^2 (c'est-à-dire $\pi(\delta^2) = 1/(1+\delta^2)^2$, $\delta^2 > 0$) afin d'éviter l'impropriété de la densité *a posteriori*. Soit $\phi = 1/(1+\delta^2)$, alors $\phi \sim \text{Bêta}(1, 1)$ (c'est-à-dire qu'il est uniforme). Il faut noter que si x a une densité suivant une loi demi-Cauchy, la densité de $\sqrt{\delta^2}$ après la transformation $x = \sqrt{\delta^2}$ est alors la loi *a priori* fondée sur Cauchy, $\pi(\delta^2) = \frac{1}{\pi\sqrt{\delta^2}(1+\delta^2)}$, qui se traduit par $\phi \sim \text{Bêta}(0,5; 0,5)$. En outre, les deux densités sont sous la forme de répartition f de Snedecor, où la première densité est une $f(2, 2)$ et la version de Cauchy est une $f(1, 1)$; la $f(2, 2)$ est mathématiquement un peu plus pratique lorsque nous la transformons en $(0, 1)$.

Soit $V = \{\boldsymbol{\theta} : \theta_i \geq c_i, i = 1, \dots, \ell, \sum_{i=1}^{\ell} \theta_i < a\}$. Dans ce cas-ci, cette densité *a posteriori* conditionnelle, $\pi(\boldsymbol{\theta} | \boldsymbol{\beta}, \delta^2, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\sigma}}^2)$, est assujettie à la contrainte d'inégalité et à la contrainte $\sum_{i=1}^{\ell} \theta_i < a$, où a est la cible de réconciliation. Il convient de mentionner que l'inégalité est stricte puisqu'avec l'égalité, l'une des θ_i devient redondante. Cette redondance doit être prise en considération lors de l'ajustement du modèle (un problème beaucoup plus difficile à résoudre), mais avec la contrainte d'inégalité, il n'est pas nécessaire de le faire (un problème beaucoup plus simple à résoudre). En d'autres termes, nous devons tirer $\theta_1, \dots, \theta_{\ell}$ soumises aux contraintes $\theta_i \geq c_i, i = 1, \dots, \ell$ et $\sum_{i=1}^{\ell} \theta_i < a$. Il convient de noter encore une fois que la contrainte de réconciliation n'est que partiellement incluse dans le modèle de Fay-Herriot. Nous utiliserons un échantillonneur de Gibbs pour exécuter cette méthode d'échantillonnage, et la contrainte de réconciliation sera entièrement intégrée dans une analyse des résultats à partir de l'échantillonneur de Gibbs au moyen d'une procédure de ratissage.

La densité *a priori* conjointe est

$$\pi(\boldsymbol{\theta}, \boldsymbol{\beta}, \delta^2) = \pi(\boldsymbol{\beta}, \delta^2) \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}}, \quad \boldsymbol{\theta} \in V, \quad (2.1)$$

où $\phi(\cdot)$ est la densité normale standard. En effet, il s'agit d'une densité *a priori* conjointe très incommode, la constante de normalisation étant une fonction de $(\boldsymbol{\beta}, \delta^2)$. Ensuite, en utilisant le théorème de Bayes, la densité *a posteriori* conjointe est

$$\pi(\boldsymbol{\theta}, \boldsymbol{\beta}, \delta^2 \mid \hat{\boldsymbol{\theta}}, \mathbf{s}^2) \propto \pi(\boldsymbol{\beta}, \delta^2) \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}} \left[\prod_{i=1}^{\ell} \phi\{(\theta_i - \hat{\theta}_i)/s_i\} \right], \quad \boldsymbol{\theta} \in V. \quad (2.2)$$

Il est difficile d'employer les méthodes Monte Carlo par chaîne de Markov pour tirer efficacement des échantillons de $\pi(\boldsymbol{\theta}, \boldsymbol{\beta}, \delta^2 \mid \hat{\boldsymbol{\theta}}, \mathbf{s}^2)$ dans l'équation (2.2).

Nous allons montrer la façon de tirer des échantillons de $\pi(\boldsymbol{\theta}, \boldsymbol{\beta}, \delta^2 \mid \hat{\boldsymbol{\theta}}, \mathbf{s}^2)$ en utilisant l'intégration numérique, l'échantillonneur de Gibbs et l'échantillonneur de Metropolis. (Il est important de noter que, dans la discussion ci-dessous, à l'exception de $\sum_{i=1}^{\ell} \theta_i < a$, l'utilisation des symboles « inférieur ou égal » n'a pas d'importance, car les θ_i sont des variables aléatoires continues.)

Nous montrons d'abord la façon de tirer les θ_i à l'aide de l'échantillonneur de Gibbs. Pour les contraintes, nous avons $c_i \leq \theta_i$, $i=1, \dots, \ell$ et $\sum_{i=1}^{\ell} \theta_i < a$. Cela signifie que $\sum_{i=1}^{\ell} c_i < \sum_{i=1}^{\ell} \theta_i < a$, et que $\max\left(c_i, \sum_{j=1}^{\ell} c_j - \sum_{j=1, j \neq i}^{\ell} \theta_j\right) < \theta_i < a - \sum_{j=1, j \neq i}^{\ell} \theta_j$, $i=1, \dots, \ell$. Par conséquent, le support de la densité *a posteriori* conditionnelle de θ_i , étant donné que $\boldsymbol{\theta}_{(i)} = (\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_{\ell})'$, est

$$\max\left(c_i, \sum_{j=1}^{\ell} c_j - \sum_{j=1, j \neq i}^{\ell} \theta_j\right) < \theta_i < a - \sum_{j=1, j \neq i}^{\ell} \theta_j, \quad i=1, \dots, \ell.$$

Il est facile de démontrer que la densité *a posteriori* conditionnelle est

$$\begin{aligned} \theta_i \mid \boldsymbol{\theta}_{(i)}, \boldsymbol{\beta}, \delta^2, \hat{\boldsymbol{\theta}}, \mathbf{s}^2 &\sim \text{Normale}\left\{\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}'_i \boldsymbol{\beta}, (1 - \lambda_i) \delta^2\right\}, \quad \lambda_i = \delta^2 / (\delta^2 + s_i^2), \\ u_i = \max\left(c_i, \sum_{j=1}^{\ell} c_j - \sum_{j=1, j \neq i}^{\ell} \theta_j\right) &< \theta_i < a - \sum_{j=1, j \neq i}^{\ell} \theta_j = v_i, \quad i=1, \dots, \ell. \end{aligned} \quad (2.3)$$

Maintenant, nous voulons tirer θ_i soumise à la contrainte $u_i \leq \theta_i \leq v_i$. Pour échantillonner $X \sim \text{Normale}(\mu, \sigma^2)$, $a \leq X \leq b$, nous avons le résultat suivant (voir Devroye, 1986),

$$X = \mu + \sigma \Phi^{-1}\left\{(1-U)\Phi\left(\frac{a-\mu}{\sigma}\right) + U\Phi\left(\frac{b-\mu}{\sigma}\right)\right\},$$

où $U \sim \text{Uniforme}(0, 1)$ et où $\Phi(\cdot)$ et $\Phi^{-1}(\cdot)$ sont respectivement la fonction de répartition et la fonction de répartition inverse de la densité normale standard. Nous utilisons l'échantillonneur de Gibbs pour tirer un échantillon $\boldsymbol{\theta}$ dans l'équation (2.3). On obtient cela en tirant $u_i \leq \theta_i \leq v_i$, $i=1, \dots, n$ tour à tour.

L'étape finale consiste à ratisser $\theta_1, \dots, \theta_{\ell}$ selon la cible a pour chaque itération. De cette façon, les itérations finales sont

$$\tilde{\theta}_i = \frac{a}{\sum_{j=1}^{\ell} \theta_j} \theta_i, \quad i=1, \dots, \ell,$$

et l'on peut faire une inférence *a posteriori* sur $\theta_1, \dots, \theta_\ell$ en utilisant ces vecteurs d'itérations ratissés. On comprend maintenant la raison pour laquelle $\sum_{i=1}^{\ell} \theta_i < a$. Il convient de noter encore une fois qu'il s'agit d'une analyse simple et directe des résultats de l'échantillonneur de Gibbs.

Nous montrons ensuite la façon de tirer des échantillons de $\pi(\boldsymbol{\beta}, \delta^2 \mid \hat{\boldsymbol{\theta}}, \mathbf{s}^2)$ à l'aide de l'intégration numérique et de l'échantillonneur de Metropolis. La densité *a posteriori* conjointe de $(\boldsymbol{\beta}, \delta^2)$ est

$$\pi(\boldsymbol{\beta}, \delta^2 \mid \hat{\boldsymbol{\theta}}, \mathbf{s}^2) \propto \pi(\boldsymbol{\beta}, \delta^2) \frac{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta}) / \delta\} \phi\{(\theta_i - \hat{\theta}_i) / s_i\} d\boldsymbol{\theta}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta}) / \delta\} d\boldsymbol{\theta}},$$

ce qui, par développement des carrés, peut être simplifié à

$$\pi(\boldsymbol{\beta}, \delta^2 \mid \hat{\boldsymbol{\theta}}, \mathbf{s}^2) \propto \pi(\boldsymbol{\beta}, \delta^2) \left[\prod_{i=1}^{\ell} \phi\left\{(\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}) / \sqrt{\delta^2 / \lambda_i}\right\} \right] R(\boldsymbol{\beta}, \delta^2), \quad (2.4)$$

avec

$$R(\boldsymbol{\beta}, \delta^2) = \frac{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mu_i) / \tau_i\} d\boldsymbol{\theta}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta}) / \delta\} d\boldsymbol{\theta}},$$

où $\mu_i = \lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}'_i \boldsymbol{\beta}$ et $\tau_i^2 = (1 - \lambda_i) \delta^2$, $i=1, \dots, \ell$. Nous utiliserons l'échantillonneur de Metropolis pour l'ajustement de l'équation (2.4). Deux questions essentielles sont soulevées, à savoir l'élaboration d'une loi instrumentale efficace et le calcul du ratio, $R(\boldsymbol{\beta}, \delta^2)$, des deux intégrales de l'équation (2.4).

Tout d'abord, nous examinons l'élaboration d'une loi instrumentale. Nous disposons d'échantillons de $(\boldsymbol{\beta}, \delta^2)$ du modèle de Fay-Herriot. Nous pouvons maintenant transformer δ^2 en $\beta_{p+1} = \log(\delta^2)$ et l'ajouter comme dernière composante pour obtenir un nouveau vecteur $\boldsymbol{\beta}$ ayant des $p+1$ composantes. Procédons maintenant à l'ajustement d'une densité normale multivariée aux échantillons, $\boldsymbol{\beta} \sim \text{Normale}(\hat{\boldsymbol{\beta}}, \sigma^2 \hat{\Sigma})$, où $\hat{\boldsymbol{\beta}}$ et $\hat{\Sigma}$ sont la moyenne *a posteriori* et la matrice de covariance des échantillons du modèle de Fay-Herriot, et $\eta / \sigma^2 \sim \text{Gamma}(\eta / 2, 1 / 2)$ pour compléter la densité t de la variable aléatoire $(p+1)$ de Student à η degrés de liberté, où η est une constante de réglage.

Ensuite, nous décrivons la façon d'estimer le ratio des intégrales de l'équation (2.4). Soit $\tilde{V} = \{\boldsymbol{\theta} : c_i < \theta_i < \infty, i=1, \dots, \ell\}$; nous avons en fait choisi une borne supérieure pour chaque θ_i . Il convient de noter que $V \subset \tilde{V}$, et peut-être que \tilde{V} n'est pas beaucoup plus grand que V . Soit $I(\boldsymbol{\theta} \in V) = 1$ si $\boldsymbol{\theta} \in V$ et $I(\boldsymbol{\theta} \in V) = 0$ autrement. Alors,

$$R(\boldsymbol{\beta}, \delta^2) = \frac{\int_{\boldsymbol{\theta} \in \tilde{V}} I(\boldsymbol{\theta} \in V) \prod_{i=1}^{\ell} \phi\{(\theta_i - \mu_i) / \tau_i\} d\boldsymbol{\theta}}{\int_{\boldsymbol{\theta} \in \tilde{V}} I(\boldsymbol{\theta} \in V) \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta}) / \delta\} d\boldsymbol{\theta}}.$$

On peut désormais calculer $R(\boldsymbol{\beta}, \delta^2)$ à l'aide de méthodes Monte Carlo. Comme fonction d'importance, nous utilisons les densités *a posteriori* conditionnelles de $\theta_i, i=1, \dots, \ell$, soumises aux contraintes sur \tilde{V} . En d'autres termes,

$$\theta_i | \boldsymbol{\beta}, \delta^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mu_i, \tau_i^2), c_i < \theta_i < \infty, i=1, \dots, \ell. \quad (2.5)$$

Il est maintenant facile de tirer des échantillons $\boldsymbol{\theta}^{(h)}, h=1, \dots, M$ dans l'équation (2.5), où $M \approx 1\,000$ environ; voir Devroye (1986). Alors, un estimateur Monte Carlo de $R(\boldsymbol{\beta}, \delta^2)$ est

$$\widehat{R(\boldsymbol{\beta}, \delta^2)} = \frac{\sum_{h=1}^M I(\boldsymbol{\theta}^{(h)} \in V)}{\sum_{h=1}^M I(\boldsymbol{\theta}^{(h)} \in V) \left[\prod_{i=1}^{\ell} \frac{\phi\{(\theta_i^{(h)} - \mathbf{x}'_i \boldsymbol{\beta})/\delta\}}{\phi\{(\theta_i^{(h)} - \mu_i)/\tau_i\}} \right]}.$$

Il convient de noter que pour chaque h , une fois que l'on a tiré les $\theta_i^{(h)}, i=1, \dots, \ell$ de la loi instrumentale, il suffit de vérifier que $\sum_{i=1}^{\ell} \theta_i^{(h)} < a$. Cependant, cet estimateur Monte Carlo n'existe peut-être pas, ce qui se produit de toute évidence lorsque $\boldsymbol{\theta}^{(h)} \notin V, h=1, \dots, M$ (tout M); dans un tel cas, nous utilisons l'estimateur modifié,

$$\widehat{R_m(\boldsymbol{\beta}, \delta^2)} = \left[\frac{1}{M} \sum_{h=1}^M \prod_{i=1}^{\ell} \frac{\phi\{(\theta_i^{(h)} - \mathbf{x}'_i \boldsymbol{\beta})/\delta\}}{\phi\{(\theta_i^{(h)} - \mu_i)/\tau_i\}} \right]^{-1}.$$

En d'autres termes, nous remplaçons simplement V par \tilde{V} pour former une approximation au cas où l'estimateur Monte Carlo n'existerait pas. Dans les deux cas, nous avons tiré le θ_i comme dans l'équation (2.5), où $\theta_i | \boldsymbol{\beta}, \delta^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mu_i, \tau_i^2), c_i < \theta_i < \infty, i=1, \dots, \ell$. Il est possible que certains des $\boldsymbol{\theta}^{(h)}$ soient dans V , et, dans un tel cas, si le nombre de $\boldsymbol{\theta}^{(h)} \in V$ est au moins égal à $M/2$, nous utilisons l'estimateur antérieur.

Notre procédure nous donne 1 000 échantillons à partir de la densité *a posteriori* de $(\boldsymbol{\beta}, \delta^2)$, en utilisant l'échantillonneur de Metropolis. Les échantillons plus importants de $\theta_1, \dots, \theta_{\ell}$ sont alors obtenus à l'aide de l'échantillonneur de Gibbs. Pour chacune des 1 000 itérations de $(\boldsymbol{\beta}, \delta^2)$ provenant de l'échantillonneur de Metropolis, nous exécutons l'échantillonneur de Gibbs jusqu'à environ 100 itérations et choisissons le dernier ensemble de $\theta_1, \dots, \theta_{\ell}$. Il s'agit de l'échantillonneur dit de Gibbs-within-Metropolis. Ce n'est pas trop coûteux et c'est raisonnablement efficace; nous avons constaté des difficultés semblables dans certains de nos projets (par exemple, Nandram et Choi, 2010; Chen, Nandram et Cruze, 2022).

Dans cette méthode, il n'est pas vraiment nécessaire de contrôler la convergence de l'échantillonneur de Gibbs, car nous n'avons besoin que d'une valeur; toutefois, un « rodage » est nécessaire.

3. Méthodologie dans le cadre des modèles à rétrécissement double

Nous présentons deux modèles à rétrécissement double dans lesquels nous modélisons à la fois les variances et les moyennes d'échantillon. Les contraintes d'inégalité sont également incluses. Dans le cas

présent, l'emprunt d'information se fait aussi bien à partir des moyennes que des variances. Pour la spécification des variances, la première utilise un modèle de régression gamma et la seconde, un modèle log-linéaire. À la section 3, nous modélisons les variances d'échantillon à l'aide de la régression gamma; la section 3.1 décrit la méthode et la section 3.2, le calcul. D'autres calculs sont présentés à l'annexe A. À l'annexe B, nous décrivons le deuxième modèle à rétrécissement double pour les variances d'échantillon en utilisant le modèle log-linéaire. Même un traitement entièrement bayésien du modèle log-linéaire présente des avantages considérables en ce qui concerne les calculs par rapport au modèle de régression gamma.

Nous discutons des raisons de l'existence des deux modèles à rétrécissement double. Les calculs sont difficiles dans les deux modèles. Nous préférons le modèle gamma parce qu'il est plus précis selon les normes de la méthode Monte Carlo par chaîne de Markov. Malheureusement, les calculs nécessitent trop de temps et le système n'est pas opérationnel à NASS. Nous avons cru qu'en passant à un modèle log-linéaire, nous pourrions faire certaines approximations mathématiques, ce qui permettrait à la procédure de rétrécissement double d'être mise en œuvre à NASS ainsi que dans de nombreux autres organismes gouvernementaux. Dans le modèle log-linéaire, nous avons fait deux approximations qui permettent aux calculs d'être très rapides (en secondes) et d'avoir une précision raisonnable. Il est mathématiquement plus difficile de faire des approximations dans le cadre du modèle gamma, mais certains chercheurs peuvent tout de même le préférer.

3.1 Modèle de régression gamma

Pour les domaines ℓ , nous disposons des estimations d'enquête $\hat{\theta}_i$, de leurs erreurs-types s_i et des tailles d'échantillon $n_i \geq 2$ (les tailles d'échantillon doivent être au moins égales à 2). Nous commençons par un modèle pratique qui s'appuie sur notre travail sur le modèle de Fay-Herriot. Nous supposons que

$$\begin{aligned} \hat{\theta}_i \mid \theta_i, \sigma_i^2 &\stackrel{\text{ind}}{\sim} \text{Normale}(\theta_i, \sigma_i^2), \quad i = 1, \dots, \ell, \\ \frac{(n_i - 1) s_i^2}{\sigma_i^2} \mid \sigma_i^2 &\stackrel{\text{ind}}{\sim} \text{Gamma}\left(\frac{n_i - 1}{2}, \frac{1}{2}\right), \quad i = 1, \dots, \ell, \end{aligned}$$

où $X \sim \text{Gamma}(a, b)$ signifie que $f(x) = b^a x^{a-1} e^{-bx} / \Gamma(a)$, $x \geq 0$. Il convient de noter qu'étant donné θ_i et σ_i^2 , nous supposons que $\hat{\theta}_i$ et s_i^2 sont indépendants. Selon la première hypothèse, les θ_i et les σ_i^2 ne sont pas estimables, mais, en regroupant la première et la deuxième hypothèse, les θ_i et les σ_i^2 le sont. Dans ce cas-ci, l'hypothèse du khi carré est raisonnable, mais les degrés de liberté peuvent être légèrement inférieurs à la taille initiale de l'échantillon, car il devrait s'agir de la taille effective de l'échantillon. La taille effective de l'échantillon n'est normalement pas présentée à NASS; en fait, c'est le nombre de rapports pour lesquels les réponses sont positives qui est présenté. Nous avons donc utilisé la taille initiale de l'échantillon; voir Erciulescu, Cruze et Nandram (2019) pour obtenir un modèle similaire sans la contrainte d'inégalité, bien entendu.

A priori, nous supposons que

$$\begin{aligned}\theta_i | \boldsymbol{\beta}, \delta^2 &\stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'\boldsymbol{\beta}, \delta^2), i = 1, \dots, \ell, \\ \sigma_i^{-2} | \alpha, \gamma &\stackrel{\text{ind}}{\sim} \text{Gamma}\left(\frac{\alpha}{2}, \frac{\alpha e^{-x_i\gamma}}{2}\right), i = 1, \dots, \ell.\end{aligned}$$

Ces hypothèses sur θ_i et σ_i^2 assurent un rétrécissement double (rétrécissement à la fois des moyennes et des variances). Nous avons supposé dans ce cas-ci que les deux ensembles de covariables étaient identiques, mais ils peuvent bien sûr être différents.

Il convient de noter que la loi *a priori* pour σ_i^2 est conjuguée, ce qui simplifie un peu les calculs; voir Nandram et Erhardt (2004) pour obtenir des spécifications similaires pour les modèles binomiaux et les modèles de Poisson correspondants. Notre loi *a priori* pour les hyperparamètres est

$$\pi(\boldsymbol{\beta}, \delta^2, \gamma, \alpha) \propto \frac{1}{(1+\delta^2)^2} \frac{1}{(1+\alpha)^2}, \delta^2, \alpha \geq 0.$$

Autrement dit, on suppose des lois *a priori* uniformes pour $\boldsymbol{\beta}$ et γ et des lois *a priori* de rétrécissement (propres) pour δ^2 et α . On suppose également que tous les paramètres sont indépendants. Il convient de noter que δ^2 et α sont non négatifs; nous préférons donc utiliser une loi *a priori* de rétrécissement. À ce stade, il n'y a pratiquement aucun avantage mathématique, informatique ou scientifique à utiliser d'autres lois *a priori* non informatives pour α .

Dans notre modèle, nous incluons la contrainte d'inégalité, $\theta_i > c_i, i = 1, \dots, \ell, \sum_{i=1}^{\ell} \theta_i < a$, où a est la cible. Il est important de noter encore une fois que nous n'incluons que partiellement la contrainte de réconciliation. Le fait qu'il s'agisse de la même région que celle pour le modèle à rétrécissement unique, $V = \{\boldsymbol{\theta}: \theta_i \geq c_i, i = 1, \dots, \ell, \sum_{i=1}^{\ell} \theta_i < a\}$, est pratique. Par conséquent, les densités *a priori* pour les θ_i restent les mêmes,

$$\pi(\boldsymbol{\theta}, \boldsymbol{\beta}, \delta^2) = \pi(\boldsymbol{\beta}, \delta^2) \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\boldsymbol{\beta})/\delta\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}}, \boldsymbol{\theta} \in V,$$

où $\phi(\cdot)$ est la densité normale standard. Il est pratique de définir $\Omega = (\boldsymbol{\beta}, \delta^2, \gamma, \alpha)$. La densité *a priori* conjointe est alors

$$\begin{aligned}\pi(\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \Omega) &= \pi(\Omega) \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\boldsymbol{\beta})/\delta\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}} \\ &\times \prod_{i=1}^{\ell} \left\{ (\alpha e^{-x_i\gamma}/2)^{\alpha/2} (1/\sigma_i^2)^{\alpha/2+1} e^{-(\alpha e^{-x_i\gamma}/2\sigma_i^2)} / \Gamma(\alpha/2) \right\}, \boldsymbol{\theta} \in V.\end{aligned}\tag{3.1}$$

Par indépendance, la densité conjointe de $(\hat{\boldsymbol{\theta}}, \mathbf{s}^2)$ est

$$\begin{aligned}f(\hat{\boldsymbol{\theta}}, \mathbf{S}^2 | \boldsymbol{\theta}, \boldsymbol{\sigma}^2, \Omega) &= \\ \prod_{i=1}^{\ell} \left\{ \frac{1}{\sigma_i} \phi\left\{ \frac{\hat{\theta}_i - \theta_i}{\sigma_i} \right\} \right\} &\times \prod_{i=1}^{\ell} \left\{ \left\{ \left[(n_i - 1)/2\sigma_i^2 \right]^{(n_i-1)/2} (s_i^2)^{(n_i-1)/2-1} e^{-(n_i-1)s_i^2/2\sigma_i^2} \right\} / \Gamma\{(n_i-1)/2\} \right\}.\end{aligned}\tag{3.2}$$

Enfin, en utilisant le théorème de Bayes, la densité *a posteriori* conjointe est proportionnelle au produit des équations (3.1) et (3.2), et l'on peut montrer que

$$\begin{aligned} \pi(\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \Omega | \hat{\boldsymbol{\theta}}, \mathbf{S}^2) &\propto \pi(\boldsymbol{\beta}, \delta^2, \gamma, \alpha) \frac{1}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}} \\ &\times \prod_{i=1}^{\ell} \left\{ (\alpha e^{-x_i \gamma} / 2)^{\alpha/2} (1/\sigma_i^2)^{\alpha/2+1} e^{-\alpha e^{-x_i \gamma} / 2\sigma_i^2} / \Gamma(\alpha/2) \right\} \\ &\times \prod_{i=1}^{\ell} \left\{ \frac{1}{\sqrt{(1-\lambda_i)\delta^2}} \phi\left(\frac{\theta_i - (\lambda_i \hat{\theta}_i + (1-\lambda_i) \mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1-\lambda_i)\delta^2}}\right) \frac{1}{\sqrt{\delta^2/\lambda_i}} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2/\lambda_i}}\right) \right\} \\ &\times \prod_{i=1}^{\ell} \left\{ \left[(n_i - 1) / 2\sigma_i^2 \right]^{(n_i-1)/2} e^{-\{(n_i-1)s_i^2\}/2\sigma_i^2} \right\}, \boldsymbol{\theta} \in V, \end{aligned} \quad (3.3)$$

où $\lambda_i = \delta^2 / (\delta^2 + \sigma_i^2)$, $i = 1, \dots, \ell$.

Il découle maintenant de l'équation (3.3) que les densités *a posteriori* conditionnelles de θ_i sont

$$\theta_i | \boldsymbol{\sigma}^2, \Omega, \hat{\boldsymbol{\theta}}, \mathbf{S}^2 \stackrel{\text{ind}}{\sim} \text{Normale} \left\{ \lambda_i \hat{\theta}_i + (1-\lambda_i) \mathbf{x}'_i \boldsymbol{\beta}, (1-\lambda_i) \delta^2 \right\}, i = 1, \dots, \ell, \boldsymbol{\theta} \in V. \quad (3.4)$$

On peut maintenant éliminer par intégration θ_i de l'équation (3.3) pour obtenir la densité *a posteriori* conditionnelle conjointe de $\boldsymbol{\sigma}^2$,

$$\begin{aligned} \pi(\boldsymbol{\sigma}^2 | \Omega, \hat{\boldsymbol{\theta}}, \mathbf{S}^2) &\propto \int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1-\lambda_i) \mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1-\lambda_i)\delta^2}} \right\} \frac{1}{\sqrt{(1-\lambda_i)\delta^2}} d\boldsymbol{\theta} \\ &\times \prod_{i=1}^{\ell} \left\{ \sqrt{\lambda_i} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2/\lambda_i}}\right) \right\} \prod_{i=1}^{\ell} \left\{ (1/\sigma_i^2)^{(n_i+\alpha-1)/2+1} e^{-\{(n_i-1)s_i^2 + \alpha e^{-x_i \gamma}\}/2\sigma_i^2} \right\}. \end{aligned} \quad (3.5)$$

Il convient de noter que le terme $\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}$ n'est pas une fonction de σ_i^2 et qu'il a été éliminé tout comme d'autres termes de ce type.

On peut maintenant éliminer par intégration σ_i^2 de l'équation (3.3) pour obtenir la densité *a posteriori* conjointe de Ω ,

$$\begin{aligned} \pi(\Omega | \hat{\boldsymbol{\theta}}, \mathbf{S}^2) &\propto \pi(\boldsymbol{\beta}, \delta^2, \gamma, \alpha) \prod_{i=1}^{\ell} \left\{ \frac{\Gamma(\alpha/2)}{(\alpha e^{-x_i \gamma} / 2)^{\alpha/2}} \frac{\Gamma(n_i + \alpha - 1)/2}{\left\{ (n_i - 1) s_i^2 + \alpha e^{-x_i \gamma} / 2 \right\}^{(n_i + \alpha - 1)/2}} \right\} \\ &\times \frac{1}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'_i \boldsymbol{\beta})/\delta\} d\boldsymbol{\theta}} \int_{\boldsymbol{\sigma}^2} \left[\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1-\lambda_i) \mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1-\lambda_i)\delta^2}} \right\} \frac{1}{\sqrt{(1-\lambda_i)\delta^2}} d\boldsymbol{\theta} \right. \\ &\left. \times \prod_{i=1}^{\ell} \left\{ \frac{1}{\sqrt{\delta^2/\lambda_i}} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2/\lambda_i}}\right) \text{IG}_{\sigma_i^2}(a_i, b_i) \right\} \right] d\boldsymbol{\sigma}^2, \end{aligned} \quad (3.6)$$

où $a_i = (n_i + \alpha - 1)/2$ et $b_i = \{(n_i - 1)s_i^2 + \alpha e^{-x_i'}\}/2$, et $IG_c(a, b)$ est la densité gamma inverse, qui est donnée par $f(c) = b^a (\frac{1}{c})^{a+1} e^{-b/c} / \Gamma(a)$, $c > 0$.

3.2 Calcul pour le modèle de régression gamma

Notre stratégie consiste à tirer des échantillons de la densité *a posteriori* conjointe de Ω dans l'équation (3.6). C'est une tâche difficile, mais une fois qu'elle est accomplie, nous pouvons utiliser la règle de multiplication pour tirer des échantillons de σ_i^2 à partir de l'équation (3.5) et ensuite de θ_i à partir de l'équation (3.4). Cette stratégie est utile s'il y a un grand nombre de comtés; l'État du Texas compte 254 comtés. Nous tirons θ_i de la même manière que celle décrite à la section 2. Il est plus difficile de tirer des échantillons de σ_i^2 . Nous décrivons la façon de tirer des échantillons de Ω dans l'équation (3.6). La stratégie de base comporte deux étapes clés.

Tout d'abord, nous ajustons le modèle à rétrécissement double sans les contraintes d'inégalité et la réconciliation. Cela donne un échantillon approximatif de taille $M = 1\,000$ itérations à partir de la densité *a posteriori* de Ω que nous avons obtenue à l'aide d'un échantillonneur de Metropolis. Les détails de cette première étape sont présentés à l'annexe A.

Ensuite, nous convertissons cet échantillon approximatif en un échantillon provenant de la densité *a posteriori* au moyen de la contrainte d'inégalité et de la réconciliation. Nous utilisons les itérations M provenant de la première étape pour établir une densité t de Student multivariée pour $(\beta, \log(\delta^2), \gamma, \log(\alpha))$. À chaque itération obtenue à la première étape, nous exécutons 100 fois un échantillonneur de Metropolis avec la densité t de Student multivariée et choisissons la dernière exécution; voir Nandram et Choi (2010) pour obtenir une procédure semblable. Selon cette manière de type « diviser pour régner », nous réduisons au minimum le risque de blocage de l'échantillonneur de Metropolis. Nous voulons que l'échantillonneur de Metropolis s'éloigne au moins une fois de la valeur de départ; aucun autre contrôle n'est nécessaire; s'il ne s'éloigne pas au moins une fois, nous rejetons l'exécution. Le fait que cette procédure donne un échantillon de M itérations indépendantes de Ω est avantageux. Cependant, cette étape prend du temps; pour les données simulées actuelles, elle a pris environ 16 heures.

Nous décrivons maintenant la façon d'utiliser l'algorithme acceptation-rejet pour tirer des échantillons de σ_i^2 . Nous pouvons réécrire l'équation (3.5) comme suit

$$\begin{aligned} \pi(\sigma^2 \mid \Omega, \hat{\theta}, \mathbf{s}^2) &\propto \int_{\mathbf{0} \in \tilde{V}} \prod_{i=1}^{\ell} \phi \left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}_i' \beta)}{\sqrt{(1 - \lambda_i) \delta^2}} \right\} \frac{1}{\sqrt{(1 - \lambda_i) \delta^2}} d\theta \\ &\times \prod_{i=1}^{\ell} \left\{ \sqrt{\lambda_i} \phi \left(\frac{\hat{\theta}_i - \mathbf{x}_i' \beta}{\sqrt{\delta^2 / \lambda_i}} \right) \right\} \int_{\mathbf{0} \in \tilde{V}} I(\theta \in V) \frac{\prod_{i=1}^{\ell} \phi \left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}_i' \beta)}{\sqrt{(1 - \lambda_i) \delta^2}} \right\} \frac{1}{\sqrt{(1 - \lambda_i) \delta^2}}}{\int_{\mathbf{0} \in \tilde{V}} \prod_{i=1}^{\ell} \phi \left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}_i' \beta)}{\sqrt{(1 - \lambda_i) \delta^2}} \right\} \frac{1}{\sqrt{(1 - \lambda_i) \delta^2}} d\theta} d\theta \quad (3.7) \\ &\times \prod_{i=1}^{\ell} \left\{ (1/\sigma_i^2)^{(n_i + \alpha - 1)/2 + 1} e^{-\{(n_i - 1)s_i^2 + \alpha e^{-x_i'}\}/2\sigma_i^2} \right\}, \end{aligned}$$

où $\tilde{V} \supseteq V$ et \tilde{V} est un plus grand ensemble rectangulaire.

Il est important de noter que les premier et troisième termes de l'équation (3.7) sont des probabilités. Il est également vrai que le deuxième terme de l'équation (3.7) est une probabilité, car

$$\prod_{i=1}^{\ell} \left\{ \sqrt{\lambda_i} \phi \left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2 / \lambda_i}} \right) \right\} \leq \left\{ \frac{1}{\sqrt{2\pi}} \right\}^{\ell}.$$

Par conséquent, nous pouvons utiliser un échantillonneur acceptation-rejet pour tirer σ_i^2 .

Il faut noter que, en raison de l'élaboration, le premier terme de l'équation (3.7) est un produit sur $i = 1, \dots, \ell$. Il en va de même pour le second terme. Donc, si nous faisons abstraction du troisième terme, nous pouvons tirer indépendamment $\sigma_i^2 \stackrel{\text{ind}}{\sim} \text{IG}(a_i, b_i), i = 1, \dots, \ell$ (répartitions non restreintes) selon une probabilité,

$$\int_{\boldsymbol{\theta} \in \tilde{V}} \phi \left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1 - \lambda_i) \delta^2}} \right\} / \sqrt{(1 - \lambda_i) \delta^2} d\boldsymbol{\theta} \times \left\{ \sqrt{\lambda_i} \phi \left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2 / \lambda_i}} \right) \right\}$$

pour compléter l'algorithme acceptation-rejet. Il est possible qu'il y ait plusieurs rejets avant une acceptation, mais cela est rarement le cas. S'il y a 25 rejets, nous tirons simplement les σ_i^2 de leurs répartitions non restreintes, $\sigma_i^2 \stackrel{\text{ind}}{\sim} \text{IG}(a_i, b_i), i = 1, \dots, \ell$.

Il reste alors la question de savoir calculer

$$C = \int_{\boldsymbol{\theta} \in \tilde{V}} I(\boldsymbol{\theta} \in V) \frac{\prod_{i=1}^{\ell} \phi \left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1 - \lambda_i) \delta^2}} \right\} / \sqrt{(1 - \lambda_i) \delta^2}}{\int_{\boldsymbol{\theta} \in \tilde{V}} \prod_{i=1}^{\ell} \phi \left\{ \frac{\theta_i - (\lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1 - \lambda_i) \delta^2}} \right\} / \sqrt{(1 - \lambda_i) \delta^2} d\boldsymbol{\theta}} d\boldsymbol{\theta}.$$

Un estimateur Monte Carlo de C est

$$\hat{C} = \frac{1}{M} \sum_{h=1}^M I(\boldsymbol{\theta}^{(h)} \in V),$$

où

$$\theta_i^{(h)} \stackrel{\text{ind}}{\sim} \text{Normale} \left\{ \lambda_i \hat{\theta}_i + (1 - \lambda_i) \mathbf{x}'_i \boldsymbol{\beta}, (1 - \lambda_i) \delta^2 \right\}, c_i < \theta_i^{(h)} < \infty, h = 1, \dots, M = 1000, i = 1, \dots, \ell.$$

Toutefois, le terme $\frac{1}{M} \sum_{h=1}^M I(\boldsymbol{\theta}^{(h)} \in V)$ est difficile à intégrer dans l'échantillonneur acceptation-rejet. Nous avons surmonté cette difficulté de la manière suivante. Nous avons calculé \hat{C} et avons constaté que plus de 60 % des \hat{C} conduisent à l'acceptation de la totalité des $\sigma_i^2, i = 1, \dots, \ell$. Lorsque les σ_i^2 ne sont pas acceptés, nous tirons des échantillons de leurs répartitions non restreintes, $\sigma_i^2 \stackrel{\text{ind}}{\sim} \text{IG}(a_i, b_i), i = 1, \dots, \ell$.

4. Comparaisons à l'aide d'exemples simulés

Nous comparons nos modèles à l'aide d'exemples simulés. Nous n'effectuons pas d'étude par simulation, où la réplique est importante parce que les modèles sont déjà compliqués. Nous utilisons le coefficient de variation (c.v.) comme mesure de la fiabilité des comparaisons. Nous montrons aussi graphiquement la façon dont les observations dans les données simulées enfreignent les contraintes de borne inférieure et la façon dont nos modèles rectifient ce problème.

Les données d'enquête de NASS et les données administratives du USDA sur les superficies sont assujetties à des mesures de protection de la confidentialité. Par conséquent, nous décrivons un moyen de simuler des données semblables aux données sur la culture de maïs de l'Illinois qui ont été fréquemment utilisées dans les études récentes de NASS sur les estimations des cultures au niveau du comté, et nous les utilisons pour montrer les principales caractéristiques de notre procédure de réconciliation avec contraintes d'inégalité. En pratique, la participation aux programmes de soutien à l'agriculture peut varier selon la culture et l'État. Certaines des estimations d'enquête peuvent déjà satisfaire à la contrainte de borne inférieure, c'est-à-dire certaines $\hat{\theta}_i > c_i$, de sorte que les contraintes de borne inférieure imposées aux estimations du modèle pour ces domaines peuvent être des restrictions peu ou non contraignantes dans ces comtés. Toutefois, dans les États où les taux d'inscription aux programmes de soutien agricole sont élevés, comme dans le cas de la culture du maïs dans l'Illinois, les totaux administratifs peuvent englober une grande partie de la population, de sorte que les estimations directes, sujettes à l'erreur d'échantillonnage, sont inférieures aux totaux administratifs dans de nombreux comtés. Les estimations du modèle pour les comtés doivent être limitées par les bornes inférieures ainsi que par la cible de réconciliation.

À la section 4.1, nous décrivons plusieurs ensembles de données simulées. À la section 4.2, nous présentons les résultats selon le modèle à rétrécissement unique avec les contraintes d'inégalité. À la section 4.3, nous présentons des résultats selon le modèle à rétrécissement double pour le modèle de régression gamma et le modèle log-linéaire, toujours avec les contraintes d'inégalité. Parallèlement, nous avons comparé ces modèles avec les estimations directes (DE), les estimations du modèle de Fay-Herriot bayésien (ME) sans réconciliation ni contraintes d'inégalité et du modèle de Fay-Herriot bayésien avec réconciliation aléatoire (MERB), à la fois au niveau du comté et au niveau du district statistique agricole (dont il est question ci-après).

Il convient de noter que tous les calculs ont été effectués sur une machine équipée du système d'exploitation CentOS (version 6.10), d'un processeur Intel Xeon E5-2690 à 2,90 GHz ayant 16 cœurs logiques et de 128 Go de mémoire vive, et que le logiciel a été compilé au moyen de la version 11.1 d'ifort.

4.1 Description des ensembles de données simulées

Nandram, Erciulescu et Cruze (2019) ont simulé un ensemble de données semblable à celui de Battese, Harter et Fuller (1988); voir également Toto et Nandram (2010) et Nandram; Toto et Choi (2011). Ces données concernent les superficies de maïs et de soja ensemencées pour 37 segments et 12 comtés de l'État de l'Iowa, et il y a deux covariables. (Comme l'Illinois, l'Iowa est un grand État producteur de maïs aux

États-Unis.) Une version bayésienne du modèle pour petits domaines de Battese, Harter et Fuller (1988) est décrite dans Toto et Nandram (2010); voir également Molina, Nandram et Rao (2014).

En procédant à une simulation à partir de ces données, nous pouvons créer un ensemble de données comportant autant de domaines que nous le souhaitons. En particulier, l'Illinois compte $\ell = 102$ comtés regroupés en 9 régions plus petites que l'État appelées « districts statistiques agricoles ». Les données sont traitées pour obtenir les estimations d'enquête et les erreurs-types. Dans nos données simulées, selon la taille réelle des districts statistiques agricoles, nous avons considéré que le premier ensemble de comtés faisait partie du premier district statistique agricole, le deuxième ensemble, du deuxième district statistique agricole, et ainsi de suite, de sorte que les 12 premiers comtés correspondent au premier district statistique agricole, les 11 suivants, au deuxième, et que les districts statistiques agricoles restants comptent respectivement 9, 11, 7, 13, 15, 12 et 12 comtés. Dans le cadre de la simulation des données de superficie, nous avons également ajouté un effet aléatoire pour chaque district statistique agricole. Les tailles d'échantillon dans les comtés sont choisies uniformément dans $(2, 74)$, une fourchette réaliste de tailles d'échantillon dans l'État, comparable aux données réelles sur le maïs de l'Illinois déclarées au cours de la campagne agricole de 2014 (Erculescu, Cruze et Nandram, 2018, 2019). En outre, les coefficients de variation au niveau du comté $c.v._i$ seront simulés uniformément à partir de la fourchette de $(0,08; 0,93)$; ces extrêmes sont comparables aux valeurs rapportées dans Erculescu, Cruze et Nandram (2020) faisant référence à la campagne agricole de 2015. Compte tenu des estimations d'enquête simulées et des coefficients de variation, on obtient les erreurs-types calculées $\hat{\sigma}_i = c.v._i \times \hat{\theta}_i$. Nous disposons donc d'un ensemble de données comprenant les estimations d'enquête, $\hat{\theta}_i$, l'erreur-type d'enquête, $\hat{\sigma}_i$, et les tailles d'échantillon, n_i , pour le i^{e} comté, $i = 1, \dots, \ell$.

Le dernier élément à simuler est constitué des données qui correspondent aux valeurs administratives de superficie, c'est-à-dire les bornes inférieures, c_i . Par souci de simplicité, nous les appelons « valeurs de FSA » tout au long de la simulation. Pour refléter le lien entre les valeurs de FSA et les estimations d'enquête pour l'Illinois, nous supposons que l'équation suivante est valable,

$$c_i = \hat{\theta}_i + U_i \times \hat{\theta}_i, \quad i = 1, \dots, \ell,$$

où $U_i \stackrel{\text{iid}}{\sim} \text{Uniforme}(-s, s)$ et où l'on considère que s est une valeur appropriée (par exemple, $s = 0,10$). Toutefois, le principal problème est de savoir fixer la cible de réconciliation. Dans le problème réel, nous connaissons la cible, mais celle-ci doit être supérieure à la somme des bornes inférieures. Il est donc raisonnable de considérer que la cible est $a = c/d$, où $c \equiv \sum_{i=1}^{\ell} c_i$ et d'établir que $0 < d < 1$. L'exhaustivité des données administratives par rapport au total de l'État peut varier selon l'État et la culture, mais dans l'Illinois, cette valeur est souvent proche de 1.

4.2 Résultats selon le modèle à rétrécissement unique

En appliquant la méthodologie pour un modèle à contrainte d'inégalité ayant des variances fixes élaborée à la section 2, nous spécifions une valeur plausible de $d = 0,99$, indiquant que les données administratives simulées représentent 99 % de la superficie totale ensemencée de maïs au niveau de l'État dans l'Illinois.

Dans ce premier cas, nous limitons la fourchette de coefficients de variation à (0,05; 0,25). La figure 4.1 montre les estimations d'enquête simulées de $\hat{\theta}_i$ par rapport aux valeurs de FSA c_i (partie supérieure) et la moyenne *a posteriori* de θ_i par rapport aux valeurs c_i de FSA selon le modèle de Fay-Herriot bayésien avec contrainte d'inégalité et réconciliation, sans rétrécissement double (appelé le modèle MFSA-NDS). Dans la partie supérieure, nous pouvons voir que de nombreux points se situent au-dessus ou au-dessous de la ligne droite à 45° passant par le point d'origine. (Cela ressemble à une tendance réaliste présentée dans la figure 4.4 de Erciulescu, Cruze et Nandram [2020], semblable à celle appliquée à la récolte de maïs de l'Illinois en 2015.) Lorsque les estimations d'enquête pour de nombreux comtés sont inférieures à leurs valeurs de FSA correspondantes, tous les points du bas sont immédiatement au-dessus de la ligne droite à 45° passant par le point d'origine, ce qui indique que toutes les estimations de MFSA-NDS ne sont pas inférieures à leurs valeurs de FSA correspondantes. En outre, la somme des 102 MFSA est égale au total de l'État, ce qui satisfait à l'exigence de réconciliation par ratissage jusqu'à la cible de l'État.

Dans le tableau 4.1, nous présentons les résultats pour les données simulées de l'Illinois. Nous comparons les résultats avec notre nouveau modèle qui intègre les contraintes d'inégalité (les valeurs de FSA sont les bornes inférieures des estimations du modèle). Plus précisément, nous comparons les estimations à partir de DE, de ME, de MERB et du modèle de Fay-Herriot bayésien à rétrécissement unique avec contrainte d'inégalité et réconciliation (MFSA-NDS).

Les coefficients de variation *a posteriori* minimum, médians et maximum (exprimés en pourcentage, %) sont plus faibles que pour les deux autres modèles (ME et MERB), et encore plus pour les estimations directes (DE). Certes, comme on pouvait s'y attendre, les coefficients de variation des districts statistiques agricoles sont plus faibles que ceux des comtés, à une exception près (5,13 par rapport à 5,31 dans le tableau 4.1, mais il s'agit d'une différence mineure). Il convient de noter que le minimum au niveau du comté et le minimum au niveau du district statistique agricole ne sont pas comparables, car le comté avec le c.v. minimum ne fait pas nécessairement parti du district statistique agricole avec le c.v. minimum. Comme on pouvait s'y attendre, nous constatons que les coefficients de variation sont en ordre décroissant (DE, ME, MERB, MFSA), que la modélisation semble avantageuse et surtout que nous pouvons intégrer les valeurs de FSA dans notre modèle (MFSA) et obtenir des coefficients de variation beaucoup plus faibles.

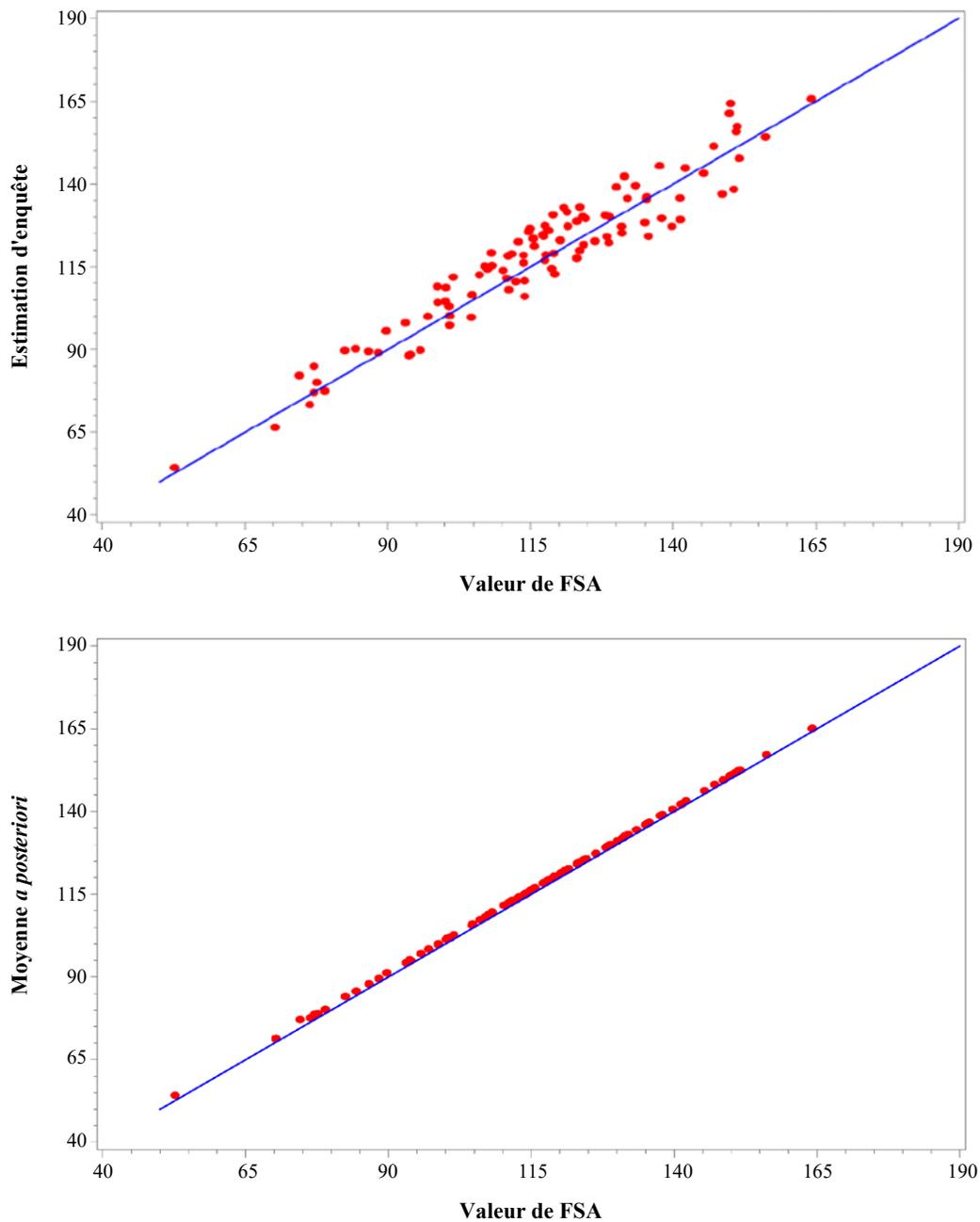
Tableau 4.1
Coefficients de variation (%) pour les données simulées de l'Illinois pour 102 comtés et 9 districts statistiques agricoles, variances fixes.

Niveau	Statistique	DE	ME	MERB	MFSA-NDS
Comté	min.	5,13	4,76	4,79	0,57
	médiane	15,57	10,67	10,58	0,97
	max.	24,93	15,80	15,34	5,22
District statistique agricole	min.	5,31	2,54	2,39	0,24
	médiane	10,60	3,25	3,01	0,30
	max.	14,81	3,92	3,51	0,42

Notes : MFSA est le nouveau modèle de réconciliation avec des valeurs de FSA comme bornes inférieures pour les estimations du modèle, c.v. (0,05 à 0,25) et $d = 0,99$.

c.v. = coefficient de variation; DE = Direct estimates; FSA = Farm Service Agency; ME = Bayesian Fay-Herriot model; MERB = Bayesian Fay-Herriot model with random benchmarking; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking; NDS = Not including double shrinkage.

Figure 4.1 Graphiques des estimations d'enquête (partie supérieure) et des moyennes *a posteriori* (partie inférieure) selon MFSA-NDS pour θ par rapport aux valeurs de FSA pour l'Illinois et les données simulées, sans rétrécissement double, c.v. (0,05 à 0,25) et $d = 0,99$.



Notes : c.v. = coefficient de variation; FSA = Farm Service Agency; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking; NDS = Not including double shrinkage.

4.3 Résultats selon le modèle à rétrécissement double

En ajustant le modèle à rétrécissement double avec les contraintes d'inégalité de la section 3 et en désignant ces estimations comme étant produites par MFSA-DS, nous ajustons le modèle aux données déjà générées de la section 4.2. En d'autres termes, les données pour lesquelles les valeurs simulées $c.v._i \in (0,05; 0,25)$ et $d = 0,99$. Les résumés des coefficients de variation pour MFSA-DS sont présentés dans le tableau 4.2, les quatre premières colonnes étant reprises du tableau 4.1. Nous constatons une légère différence entre le modèle à double rétrécissement et le modèle à rétrécissement unique. Dans les comtés, le c.v. maximum selon le modèle à rétrécissement double est un peu plus petit que celui selon les modèles MFSA-NDS, soit 3,58 % par rapport à 5,22 % pour le cas des variances fixes, mais dans les districts statistiques agricoles (agrégats de comtés), les différences entre les deux approches sont plus faibles.

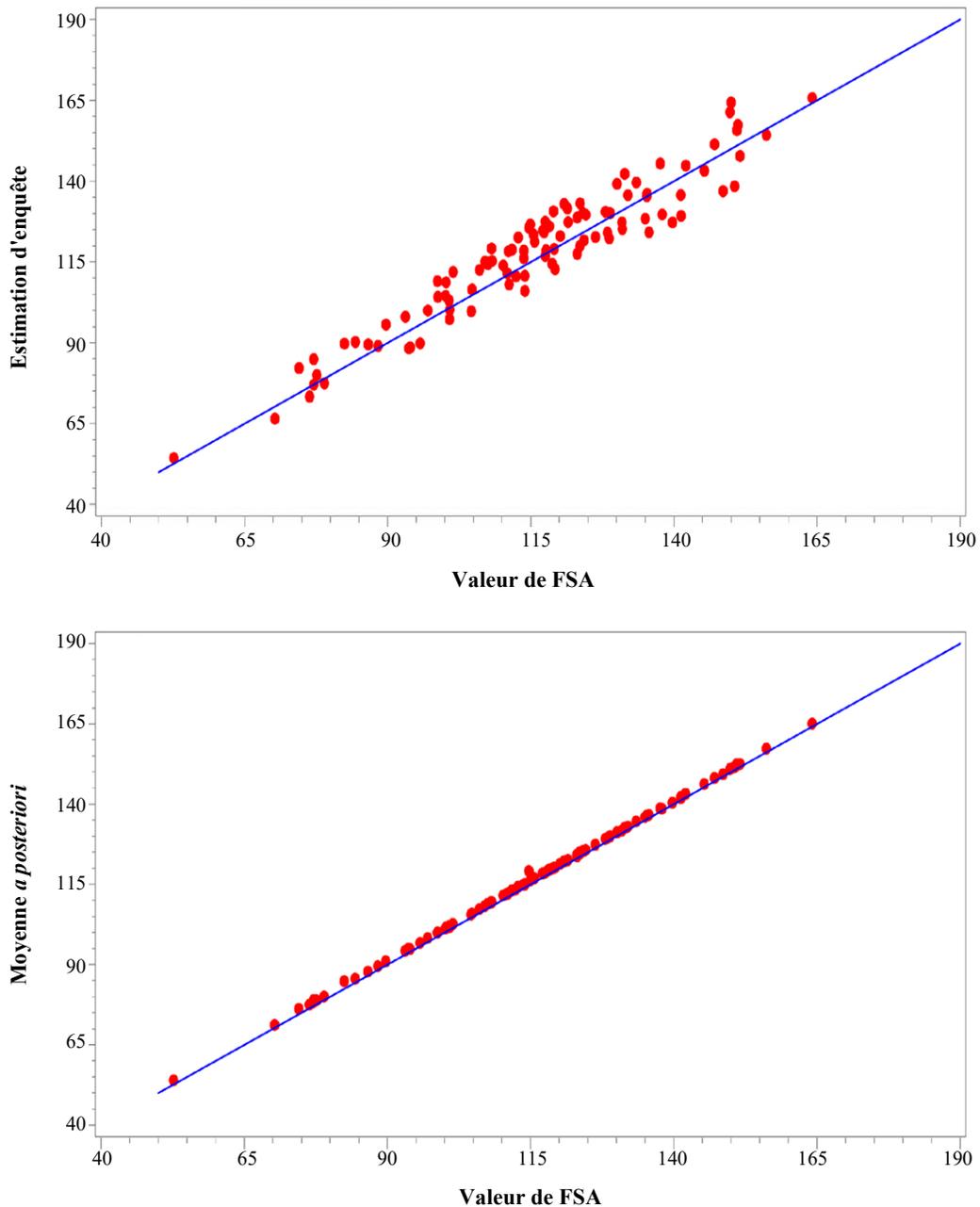
La partie supérieure de la figure 4.2 représente une fois de plus les estimations d'enquête par rapport aux valeurs de FSA (identiques à la partie supérieure de la figure 4.1), et le panneau inférieur représente les moyennes *a posteriori* par rapport aux valeurs de FSA selon le modèle à rétrécissement double avec contraintes de réconciliation et d'inégalité. La partie inférieure de la figure 4.2 n'est que légèrement différente de celle de la figure 4.1, en partie parce que la valeur $d = 0,99$ indique qu'il y a peu de marge entre la cible de l'État et le total des données administratives additionnées pour tous les comtés de l'État.

Tableau 4.2
Coefficients de variation (%) pour les données simulées de l'Illinois pour 102 comtés et 9 districts statistiques agricoles, rétrécissement double, modèle gamma.

Niveau	Statistique	DE	ME	MERB	MFSA-NDS	MFSA-DS
Comté	min.	5,13	4,76	4,79	0,57	0,55
	médiane	15,57	10,67	10,58	0,97	1,01
	max.	24,93	15,80	15,34	5,22	3,58
District statistique agricole	min.	5,31	2,54	2,39	0,24	0,26
	médiane	10,60	3,25	3,01	0,30	0,34
	max.	14,81	3,92	3,51	0,42	0,41

Notes : MFSA est le nouveau modèle de réconciliation comptant des valeurs de FSA comme bornes inférieures pour les estimations du modèle. MFSA-DS fait référence au modèle à rétrécissement double avec réconciliation et contrainte d'inégalité, c.v. (0,05 à 0,25) et $d = 0,99$. c.v. = coefficient de variation; DE = Direct estimates; DS = Double shrinkage; FSA = Farm Service Agency; ME = Bayesian Fay-Herriot model; MERB = Bayesian Fay-Herriot model with random benchmarking; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking; NDS = Not including double shrinkage.

Figure 4.2 Graphiques des estimations d'enquête (partie supérieure) et des moyennes *a posteriori* (partie inférieure) selon MFSA-DS pour θ par rapport aux valeurs de FSA pour l'Illinois et les données simulées, rétrécissement double, régression gamma, c.v. (0,05 à 0,25) et $d = 0,99$.



Notes : c.v. = coefficient de variation; DS = Double shrinkage; FSA = Farm Service Agency; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking.

Pour démontrer le modèle log-linéaire, on a généré un deuxième ensemble de données présentant des caractéristiques légèrement différentes. Notamment, nous spécifions une couverture plus faible des valeurs de FSA ($d = 0,95$) et permettons une fourchette de valeurs plus élevée pour les coefficients de variation de l'enquête, (0,08; 0,93), comparables aux coefficients de variation réels de l'enquête observés au cours de la campagne agricole de 2015. Le tableau 4.3 présente des résumés des coefficients de variation (c.v.). Une fois de plus, nous remarquons une légère différence entre le modèle log-linéaire à rétrécissement double et le modèle à rétrécissement unique au niveau du district statistique agricole. Les différences de coefficients de variation au niveau du comté sont minimales pour la moitié inférieure de tous les comtés, mais le c.v. maximum au niveau des comtés obtenu à partir du modèle à rétrécissement double (23,94 %) est nettement inférieur au c.v. maximum obtenu à partir du modèle à rétrécissement unique (cas de variances fixes) (44,92 %). Bien entendu, comme on pouvait s'y attendre, les coefficients de variation au niveau des districts statistiques agricoles sont inférieurs à ceux au niveau des comtés; il y a une exception pour les estimations directes (8,57 par rapport à 18,90 dans le tableau 4.3). Il convient une fois de plus de noter que le minimum au niveau des comtés et le minimum au niveau des districts statistiques agricoles ne sont pas comparables, car le comté ayant le c.v. minimum ne fait pas nécessairement parti du district statistique agricole ayant le c.v. minimum. Pourtant, les modèles corrigent ce problème.

Dans la figure 4.3, la partie supérieure montre les nouvelles estimations d'enquête simulées par rapport aux valeurs de FSA correspondantes, tandis que la partie inférieure montre les moyennes *a posteriori* du modèle MFSA-DS log-linéaire par rapport aux valeurs de FSA correspondantes. Contrairement à l'ensemble de données avec $d = 0,99$ des sections précédentes, le présent ensemble de données $d = 0,95$ représente une contrainte de borne inférieure plus souple. En conséquence, les estimations de superficie des comtés qui en résultent, et qui s'additionnent aussi pour correspondre au total de l'État, sont toutes visiblement au-dessus de la ligne à 45°. À titre de comparaison, les estimations de MFSA-DS obtenues par régression gamma sont représentées dans la partie inférieure de la figure 4.4. Les deux approches à rétrécissement double donnent des estimations ponctuelles similaires (mais pas identiques) pour la même cible de l'État et les mêmes contraintes administratives de borne inférieure.

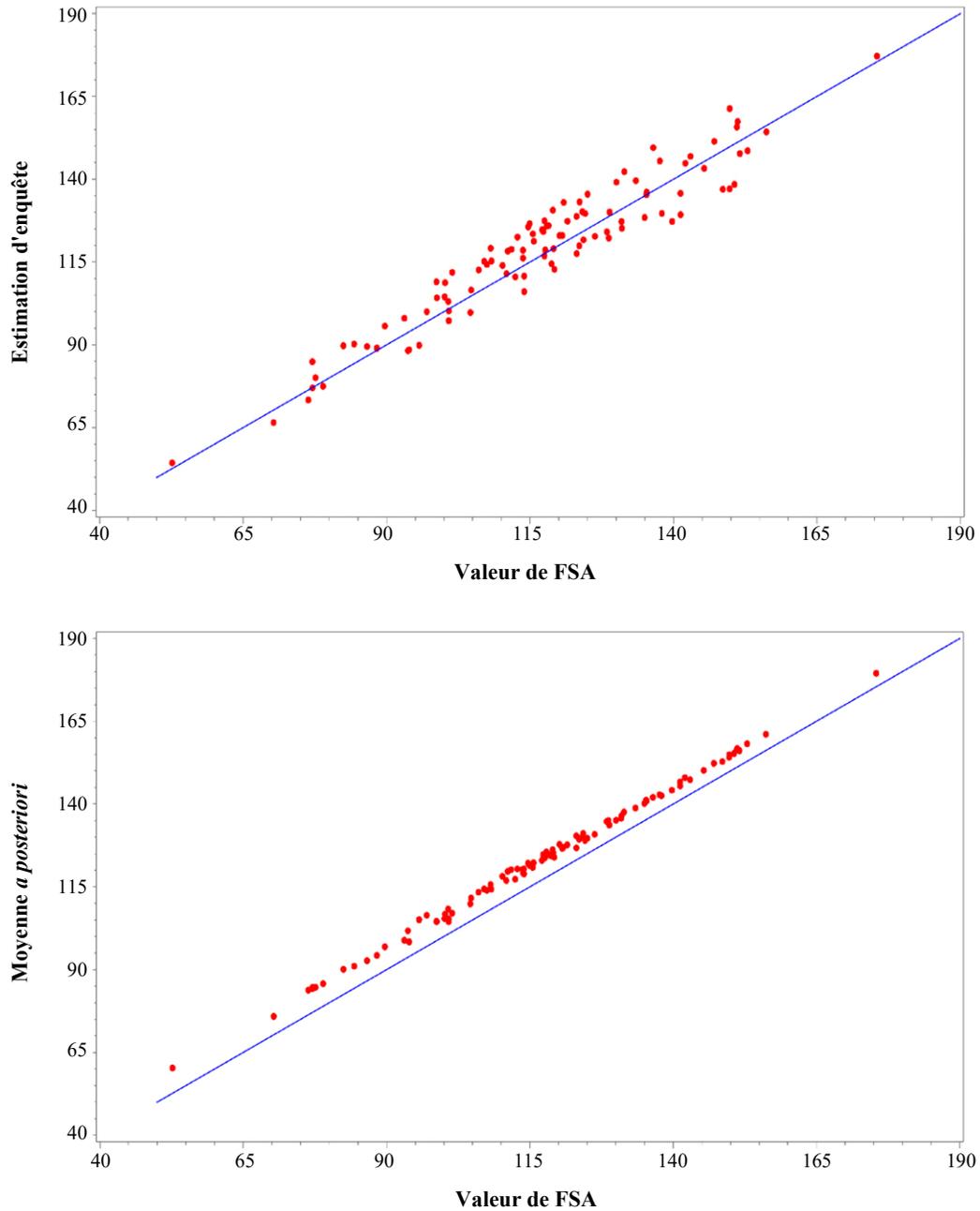
Contrairement à la régression gamma, très vorace en calcul, dont l'exécution a nécessité plus de 16 heures, les résultats du modèle log-linéaire ont été obtenus en quelques minutes, et il est possible d'accélérer davantage le processus grâce au calcul bayésien approximatif décrit à l'annexe B.

Tableau 4.3
Coefficients de variation (%) pour les données simulées de l'Illinois pour 102 comtés et 9 districts statistiques agricoles, rétrécissement double, modèle log-linéaire.

Niveau	Statistique	DE	ME	MERB	MFSA-NDS	MFSA-DS
Comté	min.	8,57	7,73	7,90	2,57	2,54
	médiane	52,90	17,83	17,16	4,56	4,92
	max.	92,70	24,25	24,82	44,92	23,94
District statistique agricole	min.	18,90	5,11	3,78	1,17	1,11
	médiane	37,70	6,15	4,71	1,83	1,43
	max.	52,10	7,19	5,81	2,65	1,63

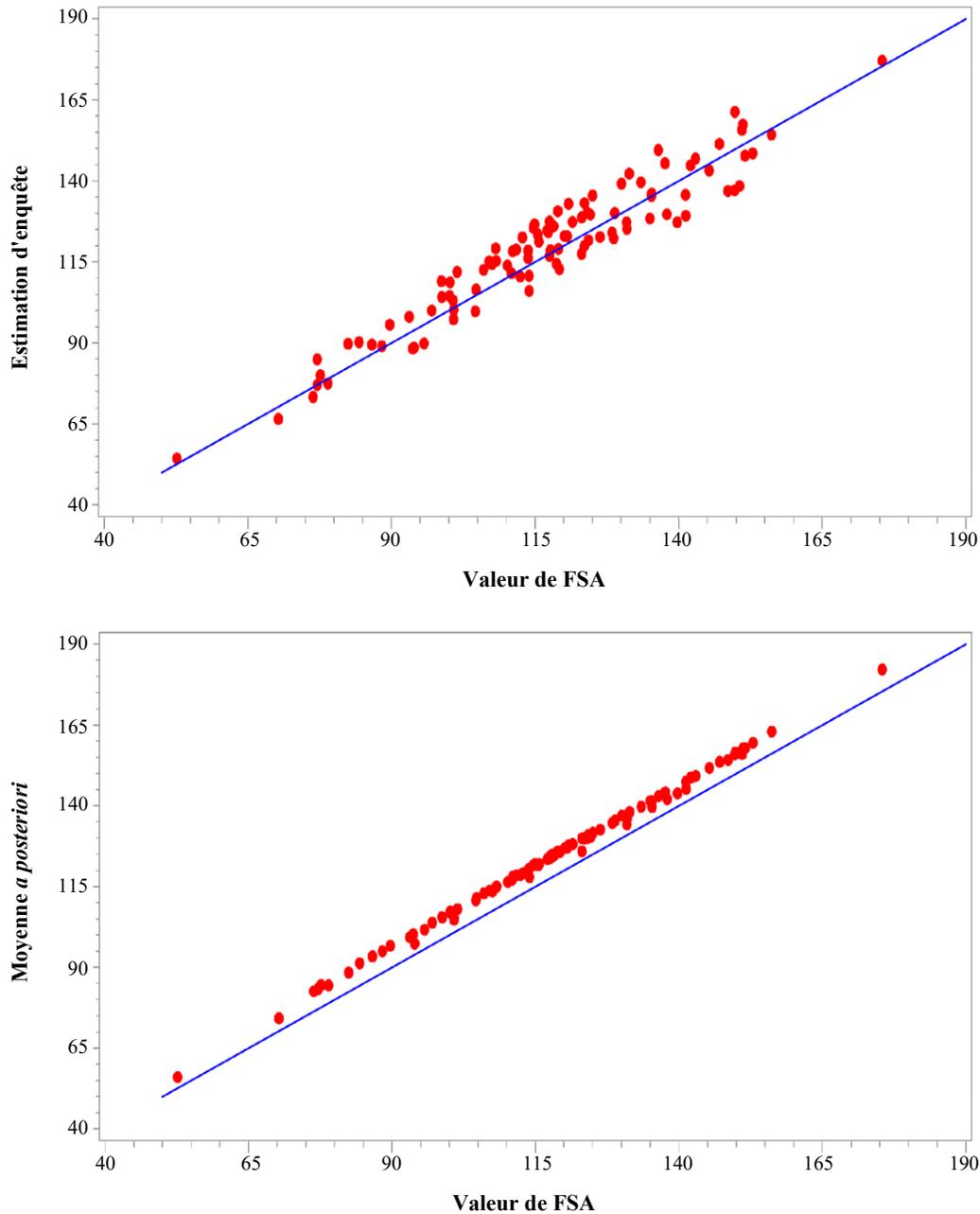
Notes : MFSA est le nouveau modèle de réconciliation comptant des valeurs de FSA comme bornes inférieures pour les estimations du modèle. MFSA-DS fait référence au modèle à rétrécissement double avec réconciliation et contrainte d'inégalité, c.v. (0,08 à 0,93) et $d = 0,95$. c.v. = coefficient de variation; DE = Direct estimates; DS = Double shrinkage; FSA = Farm Service Agency; ME = Bayesian Fay-Herriot model; MERB = Bayesian Fay-Herriot model with random benchmarking; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking; NDS = Not including double shrinkage.

Figure 4.3 Graphiques des estimations d'enquête (partie supérieure) et des moyennes *a posteriori* (partie inférieure) selon MFSA pour θ par rapport aux valeurs de FSA pour l'Illinois et les données simulées, rétrécissement double, modèle log-linéaire, c.v. (0,08 à 0,93) et $d = 0,95$.



Notes : c.v. = coefficient de variation; FSA = Farm Service Agency; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking.

Figure 4.4 Graphiques des estimations d'enquête (partie supérieure) et des moyennes *a posteriori* (partie inférieure) selon MFSA-DS pour θ par rapport aux valeurs de FSA pour l'Illinois et les données simulées, rétrécissement double, régression gamma, c.v. (0,08 à 0,93) et $d = 0,95$.



Notes : c.v. = coefficient de variation; DS = Double shrinkage; FSA = Farm Service Agency; MFSA = Bayesian Fay-Herriot model with inequality constraint and benchmarking.

5. Conclusion

À partir de la campagne agricole de 2020, NASS a converti avec succès son produit de données des estimations par comté en estimations fondées sur un modèle de système de la superficie ensemencée, de la superficie récoltée, de la production totale et du rendement par acre récolté. Il est vrai que nos méthodes peuvent être appliquées directement au rendement; seul un petit ajustement de la réconciliation est nécessaire dans l'analyse des résultats. Les estimations officielles pour 13 produits de base différents cultivés à l'échelle nationale comprennent désormais une réconciliation des estimations des comtés par rapport aux cibles prédéterminées de l'État ainsi que des contraintes de borne inférieure des superficies ensemencées. Motivés par la nécessité que le programme d'estimation des récoltes de NASS permette de produire des tableaux publiés cohérents pour tous les paramètres et selon les données administratives disponibles, nous avons montré la façon d'intégrer les contraintes d'inégalité propres à un domaine *et la réconciliation* dans le modèle de Fay-Herriot. Des modèles à rétrécissement unique et à rétrécissement double sont disponibles. L'exécution des échantillonneurs Metropolis intégraux présente des difficultés de calcul que nous avons surmontées en effectuant des approximations raisonnables supplémentaires dans le modèle à rétrécissement double.

Il est possible d'étendre le modèle hiérarchique bayésien de manière à ce que toutes les contraintes y soient effectivement incluses. En d'autres termes, θ se trouve dans $V = \left\{ \theta : c_i \leq \theta_i, \sum_{i=1}^n \theta_i = a \right\}$, où a est la cible de réconciliation et c_1, \dots, c_n sont les valeurs de FSA. Ainsi, le modèle hiérarchique bayésien (c'est-à-dire la version étendue du modèle de Fay-Herriot bayésien) a $\theta \in V$. Nous avons tenté de le faire pour le modèle le plus simple, le modèle de Fay-Herriot bayésien, mais le problème est extrêmement complexe. Cela nécessite le calcul de probabilités d'orthant (par exemple, Ridgway, 2016; Geweke, 1991; Genz, 1992) à chaque étape d'un échantillonneur Monte Carlo par chaîne de Markov. Aucun problème de ce type n'est mentionné dans Rao et Molina (2015), bien qu'ils aient utilisé la procédure de ratissage pour la réconciliation uniquement, et non pour les contraintes d'inégalité, où $\theta_i > c_i$, ce qui constitue le problème de FSA.

Néanmoins, l'intégration de la contrainte totale dans le modèle hiérarchique bayésien sera bénéfique parce qu'elle permettra de se protéger contre les défaillances du modèle, lesquelles sont si fréquentes dans l'estimation pour petits domaines; il convient d'être prudent à cet égard. Toto et Nandram (2010); Nandram et Sayit (2011); Nandram, Toto et Choi (2011); Nandram, Erciulescu et Cruze (2019) et Janicki et Vesper (2017) ont été en mesure d'intégrer une contrainte beaucoup plus simple (c'est-à-dire $\sum_{i=1}^n \theta_i = a$) dans une analyse bayésienne complète. Toutefois, comme on peut le constater, il est bien plus difficile d'intégrer la contrainte $\theta \in V$. Il s'agit d'un problème que nous aimerions résoudre. Nous pouvons ajouter des effets aléatoires aux moyennes et aux variances pour tenir compte des sous-domaines (comtés au sein des districts statistiques agricoles). Cependant, les calculs sont complexes et des approximations autres que celles fondées sur les méthodes Monte Carlo par chaîne de Markov doivent être envisagées. Nous effectuons actuellement de la recherche à ce sujet.

L'annexe C contient des commentaires sur la généralisation. On peut éviter la contrainte d'inégalité en utilisant une transformation logarithmique, mais cette méthode perd en généralité ou fait une approximation

inutile. Notre solution reste solide tant pour le modèle à rétrécissement unique que pour le modèle à rétrécissement double.

Annexe

A. Ajustement du modèle à rétrécissement double – régression gamma

En abandonnant la contrainte d'inégalité du modèle à rétrécissement double (voir l'équation [3.3]), la densité *a posteriori* conjointe est

$$\begin{aligned} \pi(\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \Omega | \hat{\boldsymbol{\theta}}, \mathbf{s}^2) &\propto \pi(\boldsymbol{\beta}, \delta^2, \gamma, \alpha) \prod_{i=1}^{\ell} \left\{ (\alpha e^{-x_i \gamma} / 2)^{\alpha/2} (1/\sigma_i^2)^{\alpha/2+1} e^{-(\alpha e^{-x_i \gamma} / 2\sigma_i^2)} / \Gamma(\alpha/2) \right\} \\ &\times \prod_{i=1}^{\ell} \left\{ \frac{1}{\sqrt{(1-\lambda_i)\delta^2}} \phi\left(\frac{\theta_i - (\lambda_i \hat{\theta}_i + (1-\lambda_i)\mathbf{x}'_i \boldsymbol{\beta})}{\sqrt{(1-\lambda_i)\delta^2}}\right) \frac{1}{\sqrt{\delta^2/\lambda_i}} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2/\lambda_i}}\right) \right\} \\ &\times \prod_{i=1}^{\ell} \left\{ \left[(n_i - 1) / 2\sigma_i^2 \right]^{(n_i-1)s_i^2/2} e^{-(n_i-1)/2\sigma_i^2} \right\}, \end{aligned} \quad (\text{A.1})$$

où $\lambda_i = \delta^2 / (\delta^2 + \sigma_i^2)$, $i = 1, \dots, \ell$. Sachant $\Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2$, il est clair que (θ_i, σ_i^2) sont indépendants sur $i = 1, \dots, \ell$. C'est la principale différence entre le modèle à rétrécissement double avec et sans contraintes d'inégalité.

Notre stratégie consiste à échantillonner d'abord la densité *a posteriori* $\pi(\Omega | \hat{\boldsymbol{\theta}}, \mathbf{s}^2)$. Ensuite, nous tirons des échantillons de la densité *a posteriori* conditionnelle conjointe de $\pi(\boldsymbol{\sigma}^2 | \Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2)$. Enfin, nous obtenons les échantillons requis à partir de $\pi(\boldsymbol{\theta} | \boldsymbol{\sigma}^2, \Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2)$. Ainsi, après avoir obtenu les tirages de Ω , nous utilisons la règle de multiplication pour obtenir σ_i^2 et θ_i (c'est-à-dire que Ω , σ_i^2 et θ_i sont tirés simultanément).

Il découle de l'équation (A.1) que, sous réserve de $\boldsymbol{\sigma}^2, \Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2$, les valeurs θ_i sont indépendantes et

$$\theta_i | \boldsymbol{\sigma}^2, \Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2 \stackrel{\text{ind}}{\sim} \text{Normale} \left\{ \lambda_i \hat{\theta}_i + (1-\lambda_i) \mathbf{x}'_i \boldsymbol{\beta}, (1-\lambda_i) \delta^2 \right\}, \quad i = 1, \dots, \ell. \quad (\text{A.2})$$

Sachant $\Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2$, les valeurs σ_i^2 sont indépendantes. Par conséquent, en éliminant par intégration θ_i à partir de l'équation (A.1), nous obtenons la densité *a posteriori* conditionnelle de $\boldsymbol{\sigma}^2$, qui est

$$\pi(\boldsymbol{\sigma}^2 | \Omega, \hat{\boldsymbol{\theta}}, \mathbf{s}^2) \propto \prod_{i=1}^{\ell} \left\{ \sqrt{\lambda_i} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}'_i \boldsymbol{\beta}}{\sqrt{\delta^2/\lambda_i}}\right) \left[(1/\sigma_i^2)^{(n_i+\alpha-1)/2+1} e^{-\{(n_i-1)s_i^2 + \alpha e^{-x_i \gamma}\}/2\sigma_i^2} \right] \right\}, \quad (\text{A.3})$$

$i = 1, \dots, \ell$. Il convient de noter que les constantes inutiles sont supprimées (par exemple, les paramètres de conditionnement).

On peut maintenant éliminer par intégration θ_i et σ_i^2 de l'équation (A.1) pour obtenir la densité *a posteriori* conjointe de Ω ,

$$\begin{aligned} \pi(\Omega | \hat{\theta}, \mathbf{s}^2) &\propto \pi(\boldsymbol{\beta}, \delta^2, \boldsymbol{\gamma}, \alpha) \prod_{i=1}^{\ell} \left\{ \frac{\Gamma(\alpha/2)}{(\alpha e^{-x_i \gamma}/2)^{\alpha/2}} \frac{\Gamma((n_i + \alpha - 1)/2)}{\{((n_i - 1)s_i^2 + \alpha e^{-x_i \gamma})/2\}^{(n_i + \alpha - 1)/2}} \right\} \\ &\times \prod_{i=1}^{\ell} \left\{ \int_0^{\infty} \frac{1}{\sqrt{\delta^2/\lambda_i}} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}_i' \boldsymbol{\beta}}{\sqrt{\delta^2/\lambda_i}}\right) \text{IG}_{\sigma_i^2}(a_i, b_i) d\sigma_i^2 \right\}, \end{aligned} \quad (\text{A.4})$$

où $a_i = (n_i + \alpha - 1)/2$ et $b_i = \{(n_i - 1)s_i^2 + \alpha e^{-x_i \gamma}\}/2$. Dans ce cas-ci, $\text{IG}_x(a, b)$ est la densité gamma inverse, qui est donnée par $f(x) = b^a \left(\frac{1}{x}\right)^{a+1} e^{-b/x} / \Gamma(a), x > 0$.

Il est facile d'échantillonner σ_i^2 dans l'équation (A.3) à l'aide de l'échantillonneur acceptation-rejet; il suffit de tirer $\sigma_i^2 | \Omega, \hat{\theta}, \mathbf{S}^2 \sim \text{IG}(a_i, b_i)$ selon la probabilité $\sqrt{\lambda_i} \phi\left(\frac{\hat{\theta}_i - \mathbf{x}_i' \boldsymbol{\beta}}{\delta^2/\lambda_i}\right)$. On voit donc qu'il est facile de tirer les θ_i à partir de l'équation (A.2). Le principal problème est maintenant de savoir échantillonner la densité *a posteriori* conjointe de Ω dans l'équation (A.4). Pour cela, nous utiliserons l'échantillonneur de Metropolis.

Une fois que nous avons obtenu un échantillon à partir de l'équation (A.4), nous le convertissons en un échantillon de l'équation (3.6), ce qui est notre objectif principal. Cela est possible grâce à un autre échantillonneur de Metropolis que nous exécutons d'une nouvelle façon. Nous préférons utiliser des lois instrumentales qui permettent d'avoir des chaînes indépendantes. On obtient cela au moyen de tirages à partir d'une densité t de Student multivariée (à élaborer). Nous n'effectuons pas un long passage parce qu'avec un échantillonneur de Metropolis, la chaîne a tendance à rester bloquée longtemps, ce qui introduit une dépendance de grande portée dans l'échantillon et donne un mélange médiocre et inefficace. Au lieu de cela, nous exécutons plusieurs chaînes, disons $M = 1\,000$ chaînes. Chaque chaîne est exécutée avec un point de départ choisi au hasard à partir d'une densité approximative pour 100 itérations, et la dernière est retenue. Seule une surveillance mineure est nécessaire pour garantir des taux de saut raisonnables. Si la chaîne ne bouge pas depuis le point de départ initial choisi au hasard, elle n'est pas utilisée dans l'échantillon final. En fin de compte, nous obtenons un échantillon aléatoire de M itérations à partir de la densité requise dans l'équation (A.4).

Nous décrivons la façon d'obtenir des échantillons à partir de la densité *a posteriori* de $\Omega = (\boldsymbol{\beta}, \delta^2, \boldsymbol{\gamma}, \alpha)$. Il y a trois étapes. La première étape consiste à obtenir un échantillon de valeurs de départ M , la deuxième consiste à obtenir une loi instrumentale pour l'échantillonneur de Metropolis à chaque valeur de départ et la troisième étape consiste à effectuer un court passage de 100 itérations de chacun des échantillonneurs de Metropolis de la deuxième étape.

Tout d'abord, nous éliminons par intégration θ_i et nous remplaçons σ_i^2 par $s_i^2, i = 1, \dots, \ell$. Étant donné $\hat{\theta}, \mathbf{s}^2, (\boldsymbol{\beta}, \delta^2)$ et $(\boldsymbol{\gamma}, \alpha)$ sont indépendants, de sorte que nous pouvons les échantillonner séparément pour obtenir $M = 1\,000$ points de départ indépendants. Nous avons obtenu ces M points de départ au moyen d'approximations simples.

Par la suite, à chaque point de départ, nous exécutons un échantillonneur de Gibbs pour obtenir $\sigma_i^2, i = 1, \dots, \ell$ et Ω . Pour ce faire, on tire les σ_i^2 de leurs densités *a posteriori* conditionnelles exactes à

l'aide d'un échantillonnage par rejet. Ensuite, étant donné $\sigma^2, s^2, (\boldsymbol{\beta}, \delta^2)$ et $(\boldsymbol{\gamma}, \alpha)$ sont encore indépendants, et les tirages à partir de leurs densités *a posteriori* conjointes respectives sont effectués de la même manière. Il convient de noter que, sachant δ^2 , la distribution de $\boldsymbol{\beta}$ est normale multivariée et $\boldsymbol{\beta}$ peut être éliminée par intégration pour obtenir la densité *a posteriori* conditionnelle de δ^2 qui peut être échantillonnée à l'aide d'une grille. Toutefois, ce n'est pas le cas pour $(\boldsymbol{\gamma}, \alpha)$, car la densité *a posteriori* conditionnelle de $\boldsymbol{\gamma}$ sachant α n'est pas standard (c'est-à-dire qu'elle n'est pas normale multivariée). Ainsi, nous approximations la densité *a posteriori* de $\boldsymbol{\gamma}$ à l'aide d'une densité normale multivariée et, au moyen de cette approximation, l'échantillonnage de $(\boldsymbol{\gamma}, \alpha)$ s'effectue de la même manière que pour $(\boldsymbol{\beta}, \delta^2)$.

Enfin, nous exécutons la deuxième étape 1 000 fois au moyen d'un « rodage » de 100 passages et nous utilisons les $M = 1\,000$ échantillons pour élaborer une densité t de Student multivariée pour $\Omega_a = (\boldsymbol{\beta}, \log(\delta^2), \boldsymbol{\gamma}, \log(\alpha))$, que nous utilisons comme loi instrumentale dans un échantillonneur de Metropolis pour échantillonner la densité *a posteriori* exacte. Cette opération est effectuée 100 fois et la dernière itération est sélectionnée. Chaque point de départ choisi au hasard contribue à l'échantillon de $M = 1\,000$ itérations de Ω_a ou de Ω à partir de la densité *a posteriori* selon le modèle à rétrécissement double sans la contrainte d'inégalité et la réconciliation.

Pour compléter l'ensemble de la procédure, pour chaque Ω_a , nous échantillonons σ_i^2 à partir de leurs densités *a posteriori* conditionnelles en utilisant l'échantillonnage par rejet pour accéder plus efficacement aux densités *a posteriori*. Ensuite, avant tout, les θ_i sont tirés de leurs densités *a posteriori* conditionnelles (normales dans ce cas-ci). L'ensemble de la procédure a duré environ quatre heures, et les taux de saut sont principalement supérieurs à 5 %.

B. Ajustement du modèle à rétrécissement double – modèle log-linéaire

Nous décrivons le modèle log-linéaire à rétrécissement double et montrons la façon de l'ajuster. L'objectif principal est de montrer la présence de gains supplémentaires en matière de vitesse de calcul en utilisant le calcul bayésien approximatif.

Notre modèle est semblable à celui de la section 3, où l'on suppose que $\hat{\theta}_i$ et s_i^2 sont indépendants par paire,

$$\begin{aligned} \hat{\theta}_i \mid \theta_i, \sigma_i^2 &\stackrel{\text{ind}}{\sim} \text{Normale}(\theta_i, \sigma_i^2), i = 1, \dots, \ell, \\ \frac{(n_i - 1) s_i^2}{\sigma_i^2} \mid \sigma_i^2 &\stackrel{\text{ind}}{\sim} \text{Gamma}\left(\frac{n_i - 1}{2}, \frac{1}{2}\right), i = 1, \dots, \ell. \end{aligned}$$

Cependant, *a priori*, nous supposons que

$$\theta_i \mid \boldsymbol{\beta}_1, \delta_1^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'_i \boldsymbol{\beta}_1, \delta_1^2), i = 1, \dots, \ell, \boldsymbol{\theta} \in V,$$

avec le modèle log-linéaire sur σ_i^2 ,

$$\ln(\sigma_i^2) \mid \boldsymbol{\beta}_2, \delta_2^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'\boldsymbol{\beta}_2, \delta_2^2), i = 1, \dots, \ell,$$

où nous supposons également que θ_i et σ_i^2 sont indépendants par paire. Il faut noter que nous avons également la restriction $\boldsymbol{\theta} \in V$. Étant donné que nous utiliserons un échantillonneur de Gibbs approximatif pour ajuster le modèle, nous supposons que $\pi(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \delta_1^2, \delta_2^2) \propto \frac{1}{\delta_1^2} \frac{1}{\delta_2^2}$ (c'est-à-dire que la propriété *a posteriori* n'est pas un problème à condition que la matrice de plan soit de plein rang).

Ensuite, en posant que $D = (\hat{\boldsymbol{\theta}}, \mathbf{s}^2)$, la densité *a posteriori* conjointe de $\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \boldsymbol{\beta}_1, \delta_1^2, \boldsymbol{\beta}_2, \delta_2^2$ est donnée par

$$\begin{aligned} \pi(\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \boldsymbol{\beta}_1, \delta_1^2, \boldsymbol{\beta}_2, \delta_2^2 \mid D) &\propto \prod_{i=1}^{\ell} \left\{ \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-(\hat{\theta}_i - \theta_i)^2 / 2\sigma_i^2} \right\} \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\boldsymbol{\beta}_1) / \delta_1\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\boldsymbol{\beta}_1) / \delta_1\} d\boldsymbol{\theta}} \\ &\times \frac{1}{\delta_1^2} \frac{1}{\delta_2^2} \prod_{i=1}^{\ell} \left\{ \left(\frac{n_i - 1}{\sigma_i^2} \right)^{(n_i - 1)/2} e^{-(n_i - 1)s_i^2 / 2\sigma_i^2} \frac{1}{\sqrt{2\pi\delta_2^2}} e^{-(\ln(\sigma_i^2) - \mathbf{x}'\boldsymbol{\beta}_2)^2 / 2\delta_2^2} \right\}, \boldsymbol{\theta} \in V. \end{aligned}$$

Notre stratégie de calcul consiste à échantillonner la densité *a posteriori* conditionnelle exacte de $\theta_i, i = 1, \dots, \ell$ et de $\sigma_i^2, i = 1, \dots, \ell$. Cependant, nous voulons remplacer les densités *a posteriori* conditionnelles de $\boldsymbol{\beta}_1, \delta_1^2$ et de $\boldsymbol{\beta}_2, \delta_2^2$ par des densités *a posteriori* approximatives. La principale question à ce stade-ci est de savoir effectuer cette dernière tâche.

Nous prenons en considération les deux modèles plus simples pour $\hat{\theta}_i$ et $s_i^2, i = 1, \dots, \ell$. Ce sont

$$\hat{\theta}_i \mid \boldsymbol{\beta}_1, \delta_1^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'\boldsymbol{\beta}_1, \delta_1^2), i = 1, \dots, \ell, \pi(\boldsymbol{\beta}_1, \delta_1^2) \propto 1 / \delta_1^2,$$

et

$$\ln(s_i^2) \mid \boldsymbol{\beta}_2, \delta_2^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'\boldsymbol{\beta}_2, \delta_2^2), i = 1, \dots, \ell, \pi(\boldsymbol{\beta}_2, \delta_2^2) \propto 1 / \delta_2^2.$$

Il est à noter que, dans le modèle complet, nous remplaçons simplement θ_i par $\hat{\theta}_i$ et σ_i^2 par s_i^2 . Dans ce cas-ci, les densités *a posteriori* de $(\boldsymbol{\beta}_1, \delta_1^2)$ et de $(\boldsymbol{\beta}_2, \delta_2^2)$, qui sont indépendantes, ont une forme simple. En posant que X représente la matrice de plan $n \times p$, alors

$$\boldsymbol{\beta}_1 \mid \hat{\boldsymbol{\theta}}, \delta_1^2 \sim \text{Normale}\left\{\hat{\boldsymbol{\beta}}_1, (X'X)^{-1}\delta_1^2\right\}, \delta_1^2 \mid \hat{\boldsymbol{\theta}} \sim \text{IG}\left\{\frac{n-p}{2}, \frac{\sum_{i=1}^n (\hat{\theta}_i - \mathbf{x}'\hat{\boldsymbol{\beta}}_1)^2}{2}\right\},$$

où $\hat{\boldsymbol{\beta}}_1 = (X'X)^{-1}X'\hat{\boldsymbol{\theta}}$. Par conséquent, la densité *a posteriori* de $\boldsymbol{\beta}_1$ est une densité *t* de Student multivariée et, dans ce cas, il est facile de tirer des échantillons de $\boldsymbol{\beta}_1$ et de δ_1^2 . En outre, en posant que $z_i = \ln(s_i^2), i = 1, \dots, \ell$, alors

$$\boldsymbol{\beta}_2 \mid \mathbf{z}, \delta_2^2 \sim \text{Normale}\left\{\hat{\boldsymbol{\beta}}_2, (X'X)^{-1}\delta_2^2\right\}, \delta_2^2 \mid \mathbf{z} \sim \text{IG}\left\{\frac{n-p}{2}, \frac{\sum_{i=1}^n (z_i - \mathbf{x}'\hat{\boldsymbol{\beta}}_2)^2}{2}\right\},$$

où $\hat{\beta}_2 = (X'X)^{-1}X'z$. Dans ce cas-ci encore, la densité *a posteriori* de β_2 est une densité *t* de Student multivariée et il est facile de tirer des échantillons de β_2 et de δ_2^2 . Notre échantillonneur de Gibbs approximatif fonctionne en utilisant ces densités *a posteriori* comme densités *a posteriori* conditionnelles. Nous devons le faire parce que le calcul est difficile et chronophage.

La densité conjointe de $(\theta_i, \sigma_i^2), i = 1, \dots, \ell$ est

$$\begin{aligned} \pi(\boldsymbol{\theta}, \boldsymbol{\sigma}^2 \mid \beta_1, \delta_1^2, \beta_2, \delta_2^2, D) &\propto \prod_{i=1}^{\ell} \left\{ \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-(\hat{\theta}_i - \theta_i)^2 / 2\sigma_i^2} \right\} \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\beta_1) / \delta_1\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\beta_1) / \delta_1\} d\boldsymbol{\theta}} \\ &\times \prod_{i=1}^{\ell} \left\{ \left(\frac{n_i - 1}{\sigma_i^2} \right)^{(n_i - 1)/2} e^{-(n_i - 1)s_i^2 / 2\sigma_i^2} \frac{1}{\sqrt{2\pi\delta_2^2}} e^{-(\ln(\sigma_i^2) - \mathbf{x}'\beta_2)^2 / 2\delta_2^2} \right\}, \boldsymbol{\theta} \in V. \end{aligned}$$

Des difficultés supplémentaires de calcul résident dans cette densité *a posteriori* conditionnelle conjointe. Constatons que, parce que $\boldsymbol{\theta} \in V$, les θ_i ne sont pas indépendants, les σ_i^2 ne sont pas indépendants et les θ_i et les σ_i^2 ne sont pas indépendants par paire. Toutefois, il faut noter que les σ_i^2 sont indépendants dans leur densité *a posteriori* conditionnelle conjointe, mais que les θ_i ne sont pas indépendants dans leur densité *a posteriori* conditionnelle conjointe. Les σ_i^2 sont tirés au moyen de la méthode de grille avec une fourchette $(\frac{1}{10}s_i^2, 10s_i^2)$ assez large, et les θ_i sont tirés au moyen de la méthode de Devroye.

Pour l'échantillonneur de Gibbs, nous avons utilisé 2 500 itérations comme rodage et avons pris chaque troisième itération pour obtenir un échantillon aléatoire de 1 000 itérations. Nous avons constaté que les tests de Geweke pour tous les θ_i et les σ_i^2 ne sont pas significatifs et que les tailles d'échantillon effectives sont toutes proches de la taille réelle de l'échantillon, qui est de 1 000 (la plupart d'entre elles sont de 1 000). Ainsi, nous disposons d'un échantillonneur de Gibbs efficace et, étonnamment, le calcul a pris moins de 20 secondes.

Nous décrivons ensuite une méthode de calcul légèrement différente de celle décrite ci-dessus. Cependant, il nous suffit de préciser la façon de tirer des échantillons à partir des densités *a posteriori* conditionnelles de (β_1, δ_1^2) et de (β_2, δ_2^2) .

La densité *a posteriori* conditionnelle de (β_2, δ_2^2) est simple (il suffit de remplacer s_i^2 par σ_i^2). Ainsi, en posant que $z_i = \ln(\sigma_i^2)$,

$$\beta_2 \mid \mathbf{z}, \delta_2^2 \sim \text{Normale} \left\{ \hat{\beta}_2, (X'X)^{-1} \delta_2^2 \right\}, \delta_2^2 \mid \mathbf{z} \sim \text{IG} \left\{ \frac{n-p}{2}, \frac{\sum_{i=1}^n (z_i - \mathbf{x}'\hat{\beta}_2)^2}{2} \right\}.$$

Il est plus difficile d'échantillonner la densité *a posteriori* conditionnelle de (β_1, δ_1^2) ,

$$\pi(\beta_1, \delta_1^2 \mid \boldsymbol{\theta}, \boldsymbol{\sigma}^2 \beta_2, \delta_2^2, D) \propto \frac{1}{\delta_1^2} \frac{\prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\beta_1) / \delta_1\}}{\int_{\boldsymbol{\theta} \in V} \prod_{i=1}^{\ell} \phi\{(\theta_i - \mathbf{x}'\beta_1) / \delta_1\} d\boldsymbol{\theta}}.$$

Nous avons commencé en utilisant l'échantillonneur de Metropolis. Après avoir utilisé deux lois instrumentales différentes, nous avons constaté une dépendance de grande portée avec des taux de saut faibles, et avons donc abandonné l'échantillonneur de Metropolis. Nous avons décidé d'utiliser des échantillonneurs à grille comme suit. Nous avons ajusté le modèle plus simple, où l'on pose que X représente la matrice de plan $n \times p$,

$$\boldsymbol{\beta}_1 | \hat{\boldsymbol{\theta}}, \delta_1^2 \sim \text{Normale}(\hat{\boldsymbol{\beta}}_1, (XX)^{-1} \delta_1^2), \quad \delta_1^2 | \hat{\boldsymbol{\theta}} \sim \text{IG} \left\{ \frac{n-p}{2}, \frac{\sum_{i=1}^n (\hat{\theta}_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}_1)^2}{2} \right\},$$

avec $\hat{\boldsymbol{\beta}}_1 = (XX)^{-1} X' \hat{\boldsymbol{\theta}}$. Par conséquent, nous pouvons maintenant échantillonner $\boldsymbol{\beta}_1$ et δ_1^2 en utilisant la règle de multiplication. Ensuite, nous trouvons les moyennes *a posteriori* (MP) et les écarts-types *a posteriori* (ETP) de chaque composante de $\boldsymbol{\beta}_1$ et de δ_1^2 ; nous choisissons leurs supports comme étant $\text{MP} \pm 6 * \text{ETP}$ et la borne inférieure de δ_1^2 comme étant $\max(0, \text{MP} - 6 * \text{ETP})$. [Presque tout le support d'une densité unimodale se trouve dans cette fourchette; en fait, nous avons constaté que la procédure n'est pas sensible au choix de 6 à l'inférence sur θ_i .] Nous exécutons maintenant la méthode de la grille dans l'échantillonneur de Gibbs pour tirer $\boldsymbol{\beta}_1$ et δ_1^2 au moyen des supports mentionnés ci-dessus pour $\boldsymbol{\beta}$ et δ_1^2 .

Pour l'échantillonneur de Gibbs, nous avons utilisé 3 500 itérations comme rodage et avons pris chaque quatrième itération pour obtenir un échantillon aléatoire de 1 000 itérations. Nous avons constaté que les tests de Geweke pour tous les θ_i et les σ_i^2 ne sont pas significatifs et que les tailles d'échantillon effectives sont toutes proches de la taille réelle de l'échantillon de 1 000 (la plupart d'entre elles sont de 1 000). Ainsi, nous disposons d'un échantillonneur de Gibbs efficace et, étonnamment, le calcul a pris moins de 40 secondes. C'est le double du temps de l'échantillonneur Gibbs approximatif vu ci-dessus (tout de même rapide).

C. Discussions sur la généralisation

Nous montrons que le problème est plus omniprésent que nous l'avons indiqué dans le présent document. Ensuite, nous discutons des problèmes liés aux solutions standard en utilisant la transformation logarithmique. Rappelons que notre problème consiste à fournir des estimations soumises à la contrainte d'inégalité de la borne inférieure et à la contrainte de réconciliation de l'égalité. Nous discutons principalement de la contrainte d'inégalité.

Le modèle de Fay-Herriot est

$$\hat{\theta}_i | \theta_i \stackrel{\text{ind}}{\sim} \text{Normale}(\theta_i, s_i^2),$$

$$\theta_i | \boldsymbol{\beta}, \delta^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'_i \boldsymbol{\beta}, \delta^2), i = 1, \dots, \ell,$$

selon la loi *a priori* $\pi(\boldsymbol{\beta}, \delta^2)$. Cela est assujéti à la contrainte d'inégalité, $\theta_i \geq c_i, i = 1, \dots, \ell$, et à la contrainte de réconciliation, $\sum_{i=1}^{\ell} \theta_i = a$, où a est la cible. Sachant $\hat{\phi}_i = \hat{\theta}_i - c_i, i = 1, \dots, \ell$ et $c = \sum_{i=1}^{\ell} c_i$. Alors,

$$\hat{\phi}_i \mid \phi_i \stackrel{\text{ind}}{\sim} \text{Normale}(\phi_i, s_i^2), \quad (\text{C.1})$$

$$\phi_i \mid \boldsymbol{\beta}, \delta^2 \stackrel{\text{ind}}{\sim} \text{Normale}(\mathbf{x}'_i \boldsymbol{\beta}, \delta^2), \phi_i \geq 0, i = 1, \dots, \ell, \quad (\text{C.2})$$

avec $\sum_{i=1}^{\ell} \phi_i = a - c$; il faut noter la présence d'un changement dans les coefficients de régression. Par conséquent, nous avons un problème général concernant des contraintes de positivité et une contrainte de réconciliation, et le problème n'est pas propre à l'agriculture. La solution au problème demeure la même que celle que nous avons présentée dans le présent document, mais nous pouvons utiliser la transformation logarithmique pour éviter les contraintes de positivité.

Il y a deux façons de procéder sans les contraintes de positivité :

- a) Transformer le $\hat{\phi}_i$ en remplaçant $\hat{\phi}_i$ par $\log(\hat{\phi}_i)$ dans l'équation (C.1). Il convient de noter que certains $\hat{\phi}_i$ peuvent être négatifs, ce qui entraîne une perte de généralité. Quand ils sont positifs, nous pouvons approximer les moyennes et les variances de la distribution normale dans l'équation (C.1) à l'aide d'une approximation du premier degré en séries de Taylor. Autrement dit, $\log(\hat{\phi}_i) \mid \phi_i \stackrel{\text{ind}}{\sim} \text{Normale}\left(\log(\phi_i), \frac{s_i^2}{\phi_i^2}\right)$. On peut procéder dans l'équation (C.2) par une régression log-normale ou une autre distribution pour les données de taille positive (par exemple, une régression gamma).
- b) Transformer le ϕ_i en remplaçant ϕ_i par e^{ϕ_i} dans l'équation (C.1). Cela crée une situation non conjuguée dans l'équation (C.2), ce qui entraîne des difficultés dans le calcul.

Il faut noter une fois de plus que la réconciliation est effectuée dans une analyse des résultats, comme nous l'avons fait dans le présent document, et qu'il est possible d'utiliser à la fois des modèles à rétrécissement unique et des modèles à rétrécissement double. Lorsque la transformation logarithmique est utilisée, la rétrotransformation vers le ϕ_i original est problématique (par exemple, Manandhar et Nandram, 2021). Toutefois, la méthodologie décrite dans le présent document constitue notre solution de première ligne.

Remerciements et avertissement

M. Cruze et M^{me} Erciulescu ont contribué à la présente recherche dans le cadre de leur mandat à la division de la recherche et du développement de National Agriculture Statistics Service (NASS), un organisme du ministère de l'Agriculture des États-Unis (USDA). Les conclusions et les résultats du présent document sont ceux des auteurs et ne doivent pas être considérés comme une représentation d'une décision ou d'une politique officielle du USDA ou du gouvernement des États-Unis. Cette recherche a été financée par NASS du USDA. Balgobin Nandram a bénéficié d'une subvention de la Simons Foundation (n° 353953, Balgobin Nandram). Les auteurs remercient également M^{me} Linda Young pour ses encouragements et les deux réviseurs pour leurs commentaires instructifs et leur aide.

Bibliographie

- Battese, G.E., Harter, R. et Fuller, W.A. (1988). An error-components model for prediction of county crop areas using survey and satellite data. *Journal of the American Statistical Association*, 83(401), 28-36.
- Chen, L., Nandram, B. et Cruze, N.B. (2022). Hierarchical Bayesian model with inequality constraints for US county estimates. *Journal of Official Statistics*, 38(3), 709-732.
- Chen, X., et Nandram, B. (2022). [Inférence bayésienne pour les données multinomiales issues de petits domaines et intégrant l'incertitude sur la restriction d'ordre](#). *Techniques d'enquête*, 48, 1, 159-192. Article accessible à l'adresse <https://www150.statcan.gc.ca/n1/fr/pub/12-001-x/2022001/article/00004-fra.pdf>.
- Cruze, N.B., Erciulescu, A.L., Nandram, B., Barboza, W.J. et Young, L.J. (2019). Producing official county-level agricultural estimates in the United States: Needs and challenges. *Statistical Science*, 34(2), 301-316.
- Dass, S.C., Maiti, T., Ren, H. et Sinha, S. (2012). [Estimation des intervalles de confiance des paramètres de petit domaine avec rétrécissement des moyennes et des variances](#). *Techniques d'enquête*, 38, 2, 187-203. Article accessible à l'adresse <https://www150.statcan.gc.ca/n1/fr/pub/12-001-x/2012002/article/11756-fra.pdf>.
- Devroye, L. (1986). *Non-Uniform Random Variate Generation*, New York: Springer-Verlag.
- Erciulescu, A.L., Cruze, N.B. et Nandram, B. (2018). Benchmarking a triplet of official estimates. *Environmental and Ecological Statistics*, 25, 523-547.
- Erciulescu, A.L., Cruze, N.B. et Nandram, B. (2019). Model-based county level crop estimates incorporating auxiliary sources of information. *Journal of the Royal Statistical Society, Series A*, 182(1), 283-303.
- Erciulescu, A.L., Cruze, N.B. et Nandram, B. (2020). Statistical challenges in combining survey and auxiliary data to produce official statistics. *Journal of Official Statistics*, 36(1), 63-88.
- Fay, R.E., et Herriot, R.A. (1979). Estimates of income for small places: An application of James-Stein procedures to census data. *Journal of the American Statistical Association*, 74(366), 269-277.
- Genz, A. (1992). Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1, 141-149.

- Geweke, J. (1991). Efficient simulation from the multivariate normal and student-*t* distributions subject to linear constraints and the evaluation of constraint probabilities. *Computing Science and Statistics: Proceedings of the Twenty-Third Symposium on the Interface*, 23, 571-578.
- Gonzalez-Manteiga, W., Lombardia, M.J., Molina, I., Morales, D. et SantaMaria, L. (2010). Small area estimation under Fay-Herriot models with nonparametric estimation of heteroscedasticity. *Statistical Modelling*, 10, 215-239.
- Janicki, R., et Vesper, A. (2017). Benchmarking techniques for reconciling Bayesian small area models at distinct geographical levels. *Statistical Methods and Applications*, 26, 557-581.
- Maiti, T., Ren, H. et Sinha, S. (2014). Prediction error of small area predictors shrinking both means and variances. *Scandinavian Journal of Statistics*, 41, 775-790.
- Manandhar, B., et Nandram, B. (2021). Hierarchical Bayesian models for continuous and positively skewed data from small areas. *Communications in Statistics – Theory and Methods*, 50(4), 944-962.
- Molina, I., Nandram, B. et Rao, J.N.K. (2014). Small area estimation of general parameters with application to poverty indicators: A hierarchical Bayes approach. *The Annals of Applied Statistics*, 8(2), 852-885.
- Nandram, B., et Choi, J.W. (2010). A Bayesian analysis of body mass index data from small domains under nonignorable nonresponse and selection. *Journal of the American Statistical Association*, 105, 120-135.
- Nandram, B., Cruze, N.B., Erciulescu, A.L. et Chen, L. (2022). Bayesian small area models under inequality constraints with benchmarking and double shrinkage. *RDD Research Report, Number RDD-22-02, National Agricultural Statistics Service, USDA*, 1-41.
- Nandram, B., Erciulescu, A.L. et Cruze, N.B. (2019). [Réconciliation bayésienne dans le modèle de Fay-Herriot par suppression aléatoire](#). *Techniques d'enquête*, 45, 2, 389-416. Article accessible à l'adresse <https://www150.statcan.gc.ca/n1/fr/pub/12-001-x/2019002/article/00004-fra.pdf>.
- Nandram, B., et Erhardt, E.B. (2004). Fitting Bayesian two-stage generalized linear models using random samples via the SIR algorithm. *Sankhyā: The Indian Journal of Statistics*, 66(4), 733-755.
- Nandram, B., et Sayit, H. (2011). [Une analyse bayésienne des probabilités de réponse dans les petits domaines sous une contrainte](#). *Techniques d'enquête*, 37, 2, 147-162. Article accessible à l'adresse <https://www150.statcan.gc.ca/n1/fr/pub/12-001-x/2011002/article/11603-fra.pdf>.

- Nandram, B., Sedransk, J. et Smith, S.J. (1997). Order restricted Bayesian estimation of the age composition of a population of Atlantic cod. *Journal of the American Statistical Association*, 92, 33-40.
- Nandram, B., Toto, M.C. et Choi, J.W. (2011). A Bayesian benchmarking of the Scott-Smith model for small areas. *Journal of Statistical Computation and Simulation*, 81, 1593-1608.
- National Academies of Sciences, Engineering, and Medicine (2017). *Improving Crop Estimates by Integrating Multiple Data Sources*, The National Academies Press, Washington, DC. <https://doi.org/10.17226/24892>.
- Rao, J.N.K., et Molina, I. (2015). *Small Area Estimation*, Wiley Series in Survey Methodology.
- Ridgway, J. (2016). Computation of Gaussian orthant probabilities in high dimension. *Statistics and Computing*, 26, 899-916.
- Sen, P.K., et Silvapulle, M.J. (2002). An appraisal of some aspects of statistical inference under inequality constraints. *Journal of Statistical Planning and Inference*, 107, 3-43.
- Silvapulle, M.J., et Sen, P.K. (2005). *Constrained Statistical Inference: Inequality, Order, and Shape Restrictions*, New York: John Wiley & Sons, Inc.
- Steorts, R.C., Schmid, T. et Tzavidis, N. (2020). Smoothing and benchmarking for small area estimation. *Revue Internationale de Statistique*, 88(3), 580-598.
- Toto, M.C.S., et Nandram, B. (2010). A Bayesian predictive inference for small area means incorporating covariates and sampling weights. *Journal of Statistical Planning and Inference*, 140, 2963-2979.
- Wang, J., et Fuller, W.A. (2003). The mean squared error of small area predictors constructed with estimated area variances. *Journal of the American Statistical Association*, 98, 716-723.
- You, Y., et Chapman, B. (2006). [Estimation pour petits domaines au moyen de modèles régionaux et d'estimations des variances d'échantillonnage](https://www150.statcan.gc.ca/n1/fr/pub/12-001-x/2006001/article/9263-fra.pdf). *Techniques d'enquête*, 32, 1, 97-103. Article accessible à l'adresse <https://www150.statcan.gc.ca/n1/fr/pub/12-001-x/2006001/article/9263-fra.pdf>.