

Ryan BOUSTANY

Paul COUAIRON

Algorithmes MCMC & MCMC accéléré pour des modèles de régression bayésienne en grande dimension



Cours enseigné par :
Stephan CLÉMENÇON

Année scolaire 2020-2021

Table des matières

Introduction	1
1 Framework et objectifs	2
1.1 Formalisation bayésienne	2
1.2 Objectifs	3
2 Algorithmes	4
3 Résultats	5
3.1 Données simulées	5
3.2 Convergence en cas sur-déterminé	7
3.3 Convergence en grande dimension	8
3.3.1 Évolution de σ	8
3.3.2 Évaluation de la vitesse de l'algorithme	8
4 Conclusion	10

Introduction

Ce projet est inspiré de l'article de recherche de Rui Jin et Aixin Tan : *Fast Markov chain Monte Carlo for high dimensional Bayesian regression models with shrinkage priors*¹, Journal of Computational and Graphical Statistics, 1-36, 2020.

Nous nous intéressons au modèle général de la régression linéaire à p variables et n exemples, qui est représenté comme suit

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon} \quad (1)$$

où $\boldsymbol{\beta}^*$ représente le vecteur des coefficients de la régression de taille $p \times 1$, que nous souhaitons approcher, \mathbf{X} notre matrice de features, de taille $n \times p$, et enfin $\boldsymbol{\varepsilon}$ un bruit gaussien, de taille $n \times 1$. Dans notre cas, nous nous intéressons à un modèle **en grande dimension**, c'est-à-dire que p est très grand, supérieur à n , voire grandement supérieur à n . L'approche standard pour résoudre ce genre de problème est d'utiliser la méthode de régression pénalisée LASSO, qui nous permet d'estimer de nombreuses valeurs de β_j^* à des valeurs nulles. Il est intéressant de quantifier l'incertitude autour de l'estimation de ces coefficients. Une alternative est de proposer des modèles bayésiens, qui nous permettent de quantifier les incertitudes des paramètres et des variables latentes à travers les distributions *à posteriori*. Des méthodes classiques de Monte-Carlo comme celle du Gibbs Sampling nous permettent, grâce à ces distributions, d'approcher le vecteur d'intérêt $\boldsymbol{\beta}^*$. Cependant, en grande dimension, la convergence de l'algorithme de Gibbs est peu satisfaisante. L'idée de R. Jin et A. Tan est de proposer un algorithme de Gibbs alternatif, requérant beaucoup moins de stockage à chaque itération. Dans ce projet, nous nous intéresserons à la comparaison de cet algorithme avec celui de Gibbs classique sur un modèle de **régression groupée**, c'est-à-dire que la matrice de features \mathbf{X} présente des groupes de colonnes avec de fortes corrélations². L'étude de ce type de problème n'est pas anodine : par exemple, en sciences sociales, si nous nous intéressons à plusieurs catégories de variables (liées à l'éducation etc.), un phénomène de corrélation intra-classes apparaîtra. Pour un problème de machine learning, le recodage de certaines variables qualitatives à l'aide d'un One-Hot Encoding fera apparaître de la corrélation dans chaque colonne liée à la variable catégorielle concernée. Dans le cadre de la régression groupée, l'estimateur LASSO est défini par :

$$\hat{\boldsymbol{\beta}}_{group} = \arg \min_{\boldsymbol{\beta}} \left\| \mathbf{Y} - \sum_{k=1}^K \mathbf{X}_{G_k} \boldsymbol{\beta}_{G_k} \right\|^2 + \lambda \sum_{k=1}^K \|\boldsymbol{\beta}_{G_k}\|_2 \quad (2)$$

1. Article disponible sur <https://doi.org/10.1080/10618600.2020.1864383>

2. Les auteurs de l'article appliquent leur modèle pour deux autres cas de grande dimension, que l'on ne traitera pas ici.

où K correspond au nombre de classes (ou groupes) de variables de notre jeu de données, et β_{G_k} et X_{G_k} représentent respectivement les sous-matrices de coefficients et de features pour les classes concernées.

1 Framework et objectifs

1.1 Formalisation bayésienne

En régression, l'approche bayésienne est généralement la même : elle consiste à proposer une distribution pour la vraisemblance du modèle, ainsi que des spécifications pour les distributions *à priori* afin d'en déduire les distributions *à posteriori* pour approcher le vecteur cible β^* . Pour les modèles en grande dimension, la spécificité que l'on souhaite préciser et qui est liée au *shrinkage* (rétrecissement, contraction) fait l'objet d'un paramètre pour lequel on définit une loi *à priori*. Le modèle décrit dans (3) correspond au modèle général utilisé dans le cadre d'estimation bayésienne pour la régression.

$$\begin{aligned} Y|\beta, \sigma^2 &\sim \mathcal{N}_n(X\beta, \sigma^2 I_n) \\ \beta|\eta, \sigma^2 &\sim \mathcal{N}_p(0, \sigma^2 \Sigma_\eta) \\ \eta &\sim p(\eta) \\ \sigma^2 &\sim \text{Inverse-Gamma}(\alpha, \xi) \end{aligned} \tag{3}$$

où Y et X sont définis dans (6), σ^2 correspond à la variance résiduelle eu modèle, et η correspond au paramètre de *shrinkage* dans notre cas. Les paramètres α, ξ concernent quant à eux la distribution à priori que l'on pose pour le paramètre σ . Dans le cas particulier de la **régression groupée**, R. Jin et A. Tan proposent un cadre bayésien s'inscrivant dans le modèle proposé dans (3) :

$$\begin{aligned} Y|\beta, \sigma^2 &\sim \mathcal{N}_n(X\beta, \sigma^2 I_n) \\ \beta_{G_k}|\tau_k^2, \sigma^2 &\stackrel{ind}{\sim} \mathcal{N}_p(0_{m_k}, \sigma^2 \tau_k^2 I_{m_k}) \\ \tau_k^2 &\stackrel{ind}{\sim} \text{Gamma}\left(\frac{m_k + 1}{2}, \frac{\lambda^2}{2}\right) \\ \sigma^2 &\sim \text{Inverse-Gamma}(\alpha, \xi) \end{aligned} \tag{4}$$

où l'indice k correspond au groupe concerné, m_k est la taille du k^e groupe, et τ correspond au paramètre de *shrinkage* (l'équivalent de η dans le modèle 3). Notons que le rôle de la matrice $\sigma_n u$ est joué par $\tau_k^2 I_{m_k}$. Le paramètre τ_k est effectivement propre à chaque groupe ;

il peut représenter la variance intra-groupe par exemple. R. Jin et A. Tan montrent que sous le modèle (4), on peut dériver les lois à posteriori suivantes :

$$\begin{aligned}
\frac{1}{\tau_k^2} | \beta, \sigma^2, Y &\stackrel{ind}{\sim} \text{Inverse-Gaussian} \left(\sqrt{\frac{\lambda^2 \sigma^2}{\|\beta_{G_k}\|^2}}, \lambda^2 \right), k = 1 \dots K \\
\sigma^2 | \beta, \tau^2, Y &\sim \text{Inverse-Gamma} \left(\frac{n + p + 2\alpha}{2}, \frac{\|Y - X\beta\|_2^2 + \beta^T D_\tau^{-1} \beta + 2\xi}{2} \right) \\
\beta | \sigma^2, \tau^2, Y &\sim \mathcal{N}_p \left((X^T X + D_\tau^{-1})^{-1} X^T Y, \sigma^2 (X^T X + D_\tau^{-1})^{-1} \right)
\end{aligned} \tag{5}$$

où $D_\tau = \text{Diag}(\underbrace{\tau_1^2, \dots, \tau_1^2}_{m_1}, \underbrace{\tau_2^2, \dots, \tau_2^2}_{m_1}, \dots, \underbrace{\tau_K^2, \dots, \tau_K^2}_{m_K}) = \tau_k^2 I_{m_k}$.

Maintenant que l'on possède un modèle bayésien précis, avec des distributions à posteriori connues, l'idée est d'appliquer les bonnes propriétés de l'inférence bayésienne dans le but d'approcher le vecteur cible à l'aide de ces lois.

1.2 Objectifs

De la même manière que l'approche fréquentiste, on souhaite approcher un vecteur cible β^* . Dans notre approche, nous ne tentons pas de construire des estimateurs basés sur les observations empiriques de \mathbf{X} et \mathbf{Y} (des réalisations $((X_i, Y_i)_{i=1 \dots n})$), mais de proposer des différentes lois à priori sur les paramètres, un modèle statistique (cette fois-ci grâce aux observations empiriques), d'en déduire les lois à posteriori et enfin de simuler ces dernières de nombreuses fois dans le but d'approcher la vraie distribution de β . L'algorithme classique utilisé pour ce genre de simulation est le Gibbs sampling, qui consiste en des tirages successifs selon les distributions à posteriori dans l'espoir que les paramètres convergent. On rappelle l'algorithme générique de Gibbs, pour deux paramètres β et σ que l'on souhaite estimer, suivant respectivement les distributions à posteriori $\pi(\beta | \mathbf{X}, \mathbf{Y}, \sigma)$ et $\pi(\sigma | \mathbf{X}, \mathbf{Y}, \beta)$:

Algorithm 1: Gibbs sampling

```

input :  $\beta^{[0]} \sim \mathcal{N}_p(0, 1), n$ 
for  $i = 1 \dots n$  do
    | draw  $\sigma^{[k]} \sim \pi(\sigma | \mathbf{X}, \mathbf{Y}, \beta^{[k-1]})$ ;
    | draw  $\beta^{[k]} \sim \pi(\beta | \mathbf{X}, \mathbf{Y}, \sigma^{[k]})$ ;
end
return  $(\beta^{[n]}, \sigma^{2[n]})$ 

```

Sous des conditions concernant les noyaux de transitions des deux distributions, il peut être montré que l'algorithme 1 converge (voir [G.O.Roberts, A.F.M.Smith]). Cependant, la convergence de l'algorithme ne peut être connue en amont ; le nombre d'itérations nécessaires avant d'atteindre un voisinage autour de la solution peut être grand. L'objectif de R. Jin et A. Tan est de proposer une méthode **accélérée** de l'algorithme 1 dans le but d'approcher le vecteur cible dans un contexte de grande dimension.

2 Algorithmes

Le premier algorithme correspond à une application du Gibbs sampling adapté au cadre (5).

Algorithm 2: 3BG

```

input :  $\beta^{[0]} \sim \mathcal{N}_p(0, 1)$ ,  $\sigma^{[0]} \sim \mathcal{N}(0, 1)$ ,  $n$ 
for  $i = 1 \dots n$  do
    draw  $\frac{1}{\tau_k^{2[i]}} | \beta, \sigma^2, Y \stackrel{ind}{\sim}$  Inverse-Gaussian  $\left( \sqrt{\frac{\lambda^2 \sigma^{2[i-1]}}{\|\beta_{G_k}^{[i-1]}\|^2}}, \lambda^2 \right)$ ,  $k = 1 \dots K$ ;
    draw  $\sigma^{2[i]} | \beta, \tau^2, Y \sim$  Inverse-Gamma  $\left( \frac{n+p+2\alpha}{2}, \frac{\|Y - X\beta^{[i]}\|_2^2 + \beta^{T[i]} D_\tau^{-1[i]} \beta^{[i]} + 2\xi}{2} \right)$ ;
    draw
     $\beta^{[i]} | \sigma^2, \tau^2, Y \sim \mathcal{N}_p \left( (X^T X + D_\tau^{-1[i]})^{-1} X^T Y, \sigma^{2[i]} (X^T X + D_\tau^{-1[i]})^{-1} \right)$ 
end
return  $(\beta^{[n]}, \sigma^{2[n]})$ 

```

Le second algorithme est une version modifiée de l'algorithme (2), en modifiant la loi à posteriori pour σ^2 . Les preuves concernant les calculs de cette distribution à posteriori peuvent être trouvées dans le supplementary material fourni par les auteurs de l'article.

Algorithm 3: 2BG

input : $\beta^{[0]} \sim \mathcal{N}_p(0, 1)$, $\sigma^{[0]} \sim \mathcal{N}(0, 1)$, n
for $i = 1 \dots n$ **do**
 draw $\frac{1}{\tau_k^{2[i]}} | \beta, \sigma^2, Y \stackrel{ind}{\sim}$ Inverse-Gaussian $\left(\sqrt{\frac{\lambda^2 \sigma^{2[i-1]}}{\|\beta_{G_K}^{[i-1]}\|^2}}, \lambda^2 \right)$, $k = 1 \dots K$;
 draw $\sigma^{2[i]} | \beta, \tau^2, Y \sim$ Inverse-Gamma $\left(\frac{n}{2} + \alpha, \frac{Y^T (I_n - X(X^T X + D_\tau^{-1})^{-1} X^T) Y}{2} + \xi \right)$
 $\beta^{[i]} | \sigma^2, \tau^2, Y \sim \mathcal{N}_p \left((X^T X + D_\tau^{-1[i]})^{-1} X^T Y, \sigma^{2[i]} (X^T X + D_\tau^{-1[i]})^{-1} \right)$
end
return $(\beta^{[n]}, \sigma^{2[n]})$

Il est important de comprendre que pour 2BG, les tirages de σ^2 et β ne sont pas les mêmes que pour 3BG. En effet, R. Jin et A. Tan proposent une alternative pour la distribution à posteriori de σ^2 , qui ne repose plus sur le vecteur β entièrement ; en ce sens, la densité de transition de la chaîne de Markov liée à (σ^2, β) ne nécessite plus un calcul prenant en compte β pour la distribution conditionnelle, ce qui allège le temps de calcul. De plus, le paramètre de l'inverse-Gaussienne est pris en se concentrant sur le K^e groupe de variables.

3 Résultats

Pour des raisons de puissance de machine, nous avons fixé le nombre d'itérations des samplers à 500. Les auteurs ont pu faire tourner les algorithmes 2BG et 3BG pour **100,000** itérations, car disposant d'une machine de 64 GB RAM, 16 coeurs et processeur de 2.6 GHz. Notre machine contient 8 GB RAM, 4 coeurs et processeur 1,4 GHz, donc une capacité d'exécution drastiquement plus faible.

3.1 Données simulées

Nous avons décidé de tester nos algorithmes sur des données simulées selon un schéma proposé par les auteurs, afin de prendre en compte l'idée de "groupes" dans la matrice de features. On rappelle le modèle général donné dans (6) :

$$Y = X\beta^* + \varepsilon \tag{6}$$

Notre objectif est de travailler en grande dimension et de s'intéresser à l'évolution des performances de nos algorithmes lorsque $\frac{p}{n}$ devient de plus en plus grand. D'autre part, on souhaite dans notre cas prendre en compte la présence de corrélation intra-groupes. Notons n la

taille de l'échantillon d'entrée, et \mathbf{K} un élément de l'ensemble $\mathcal{K} = \{5, 6, 7, 8, 9, 10, 20, 30, 40, 50\}$ (on le fera varier dans nos expériences). Pour la simulation de la matrice de features, R. Jin et A. Tan proposent la démarche suivante :

- Simuler une matrice intermédiaire \mathbf{A} de taille $\mathbf{n} \times \mathbf{K}$ formée de \mathbf{K} colonnes de lois $\mathcal{N}_K(\mathbf{0}_K, \mathbf{I}_K)$.
- Faire apparaître le phénomène de groupe en stackant à la suite de chaque colonne, les colonnes contenant ses valeurs au carré, au cube, puissance 4 et puissance 5.
- On pose alors \mathbf{X} comme étant cette matrice-là, donc de la forme :

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_{11} & \mathbf{X}_{11}^2 & \cdots & \cdots & \mathbf{X}_{11}^5 & \cdots & \cdots & \mathbf{X}_{1K} & \mathbf{X}_{1K}^2 & \cdots & \mathbf{X}_{1K}^5 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{X}_{n1} & \mathbf{X}_{n1}^2 & \cdots & \cdots & \mathbf{X}_{n1}^5 & \cdots & \cdots & \mathbf{X}_{nK} & \mathbf{X}_{nK}^2 & \cdots & \mathbf{X}_{nK}^5 \end{pmatrix}$$

Une fois la matrice \mathbf{X} simulée, nous générons β^* selon un vecteur aléatoire suivant une loi de Student pour les \mathbf{K} premières coordonnées et des $\mathbf{0}$ sur les $4\mathbf{K}$ coordonnées suivant. Puis on génère un réel positif σ^* (généralement pas plus grand que 2). Finalement, on peut générer \mathbf{Y} tel que $\mathbf{Y} = \mathbf{X}\beta^* + \varepsilon$ avec $\varepsilon \sim \mathcal{N}_n(\mathbf{0}_n, \sigma^* \mathbf{I}_n)$.

À noter que l' algorithme **2BG** prend en entrée les éléments suivant :

- \mathbf{X} : de taille $\mathbf{n} \times \mathbf{p}$, $\mathbf{p} = 5\mathbf{K}$.
- \mathbf{Y} : de taille $\mathbf{n} \times 1$.
- **group_size** : nombre d'individus par groupe (dans notre cas 5).
- β : initialisation du vecteur β .
- σ^2 : initialisation du vecteur σ^2 .
- λ : hyperparamètre.
- α : hyperparamètre.
- ξ : hyperparamètre.
- n_{iter} : nombre d'itérations.

3.2 Convergence en cas sur-déterminé

Cas de $p < n$: convergence de σ^2 et β Dans notre exemple présent, nous avons

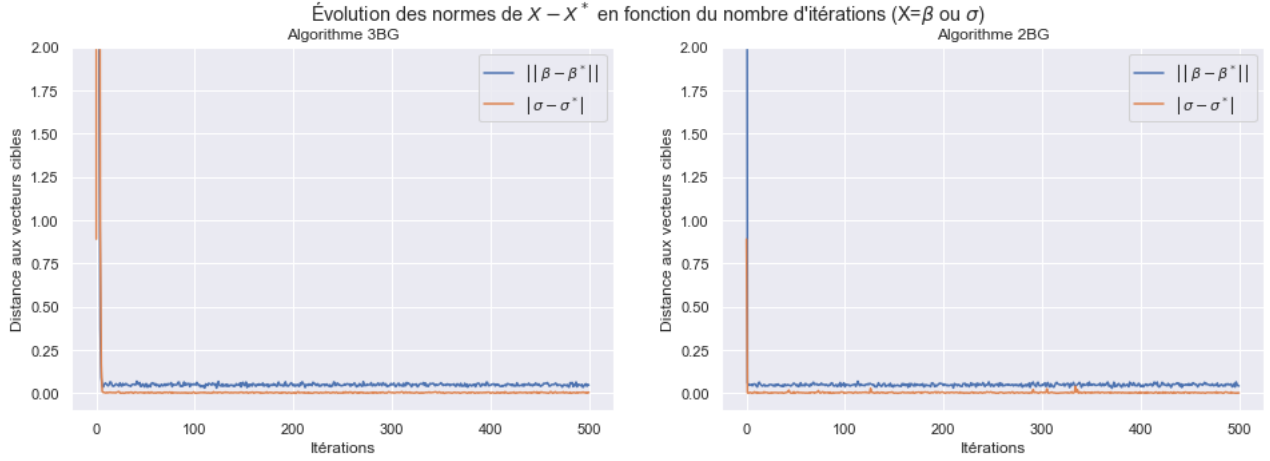


FIGURE 1 – Évolution de $|\sigma - \sigma^*|$ et $\|\beta - \beta^*\|$ en fonction du nombre d'itérations

exécuté les algorithmes **2BG** et **3BG** avec $n = 1000$ et $p = 50$. On note que les deux estimateurs (σ et β) issu des deux algorithmes convergent bien et en peu d'itérations vers β^* et σ^* . Cependant, rien ne nous permet de savoir quel algorithme semble le plus adapté à ce problème.

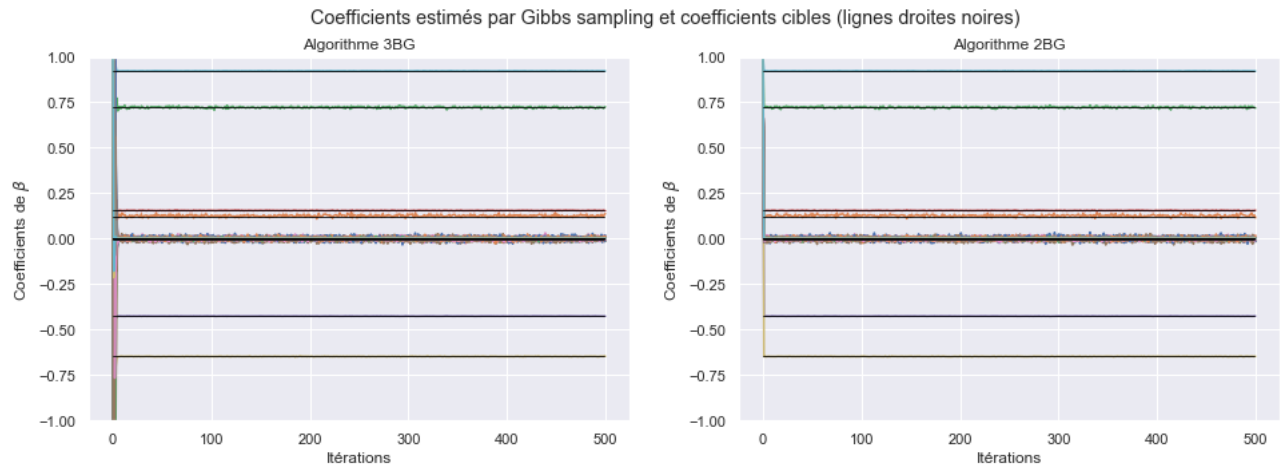


FIGURE 2 – Coefficients estimés par Gibbs sampling et coefficients cibles (lignes droites noires)

3.3 Convergence en grande dimension

L'objectif principal de ce projet est d'étudier les simulations et le comportement des différents algorithmes. On ne regardera pas vraiment quels modèles ou quelles valeurs d'hyperparamètres correspondent le mieux à chaque scénario. Ainsi, nous présentons les résultats empiriques des modèles les plus raisonnables. Nous avons aussi décidé de fixer nous même les hyperparamètres en prenant $\alpha = \xi = 0$ et $\lambda = 1$.

3.3.1 Évolution de σ

Cas de $p > n$: convergence de σ^2

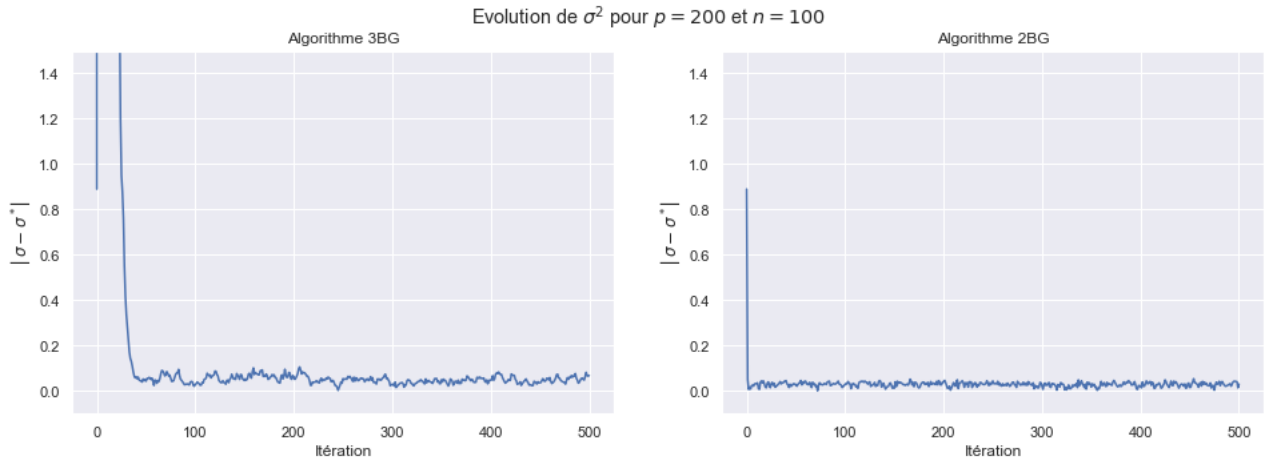


FIGURE 3 – Evolution de σ^2 pour $p = 200$ et $n = 100$

Dans ce cas présent, nous nous sommes concentrés sur σ . Les deux algorithmes nous donnent une approximation de la variance. Au bout de 100^e itération, l'algorithme **2BG** proposé par les auteurs atteint une approximation plus satisfaisante. Cependant, notons que en changeant un hyperparamètre, un Sampler peut performer mieux que l'autre. Les résultats sont globalement comparables. En regardant l'évolution de la chaîne, l'algorithme Gibbs classique semble prendre plus de temps avant de se rapprocher de la valeur cible. 2BG semble plus direct.

3.3.2 Évaluation de la vitesse de l'algorithme

On s'intéresse ici à une autre façon d'interpréter les résultats. On peut aussi se demander si la convergence est sensible au paramètre que nous posons sur l'a priori du paramètre de shrinkage τ . L'idée de R. Tin et A. Tan est alors d'évaluer, en termes de vitesse de convergence, la chaîne de Markov construite par le Gibbs sampler. Nous définissons l'effective sampling size, donnée par :

$$N_{eff} = \frac{N}{1 + 2 \sum_{k=1}^{\infty} \rho_k}$$

où ρ_k est l'autocorrélation de lag k de la chaîne de σ^2 (sous hypothèse de stationnarité).

L'effective sampling size est une estimation du nombre de tirages requis lors de l'algorithme afin d'atteindre le même niveau de précision que si le tirage était aléatoire. Il s'agit d'une mesure importante de l'efficacité d'une méthode de Monte-Carlo comme celles présentées par les algorithmes que les auteurs ont proposé. Cette métrique évalue la qualité du tirage effectué par les algorithmes 2BG et 3BG à chaque itération. Numériquement, plus N_{eff} est élevé, plus ceci indique une meilleure mixité de la chaîne puisqu'il s'agit d'une fonction décroissante des autocorrélations.

Nous allons également prendre le temps de calcul en considération et nous allons examiner l'effective sample size moyen par seconde $\frac{N_{eff}}{T}$. Autrement dit, le temps de calcul de l'effective sample size est également pris en compte afin de normaliser les valeurs de N_{eff} d'un algorithme à l'autre, mais également de donner un sens à la caractérisation dite "rapide" de l'algorithme 2BG proposé par les auteurs.

On va chercher à évaluer, pour chaque valeur de p , l'effective sample size de la chaîne de σ^2 induite par les algorithmes **2BG** et **3BG**. Les autocorrélations seront également évaluées afin de calculer la mixité des chaînes. R. Jin et A. Tan indiquent qu'une faible auto-corrélation implique une bonne mixité de la chaîne induite et donc un meilleur mélange réalisé par l'algorithme.

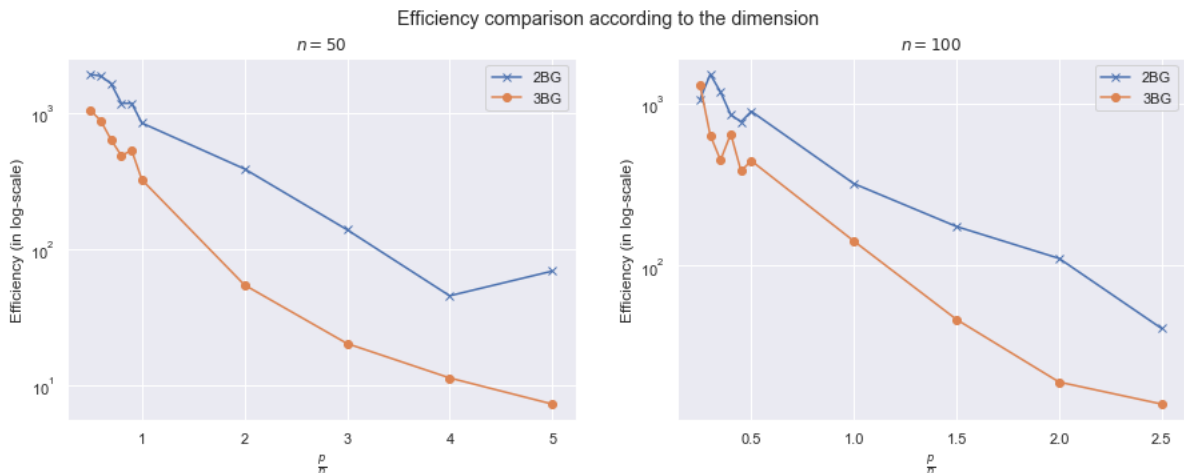


FIGURE 4 – Comparaison de l'efficacité en fonction de la dimension

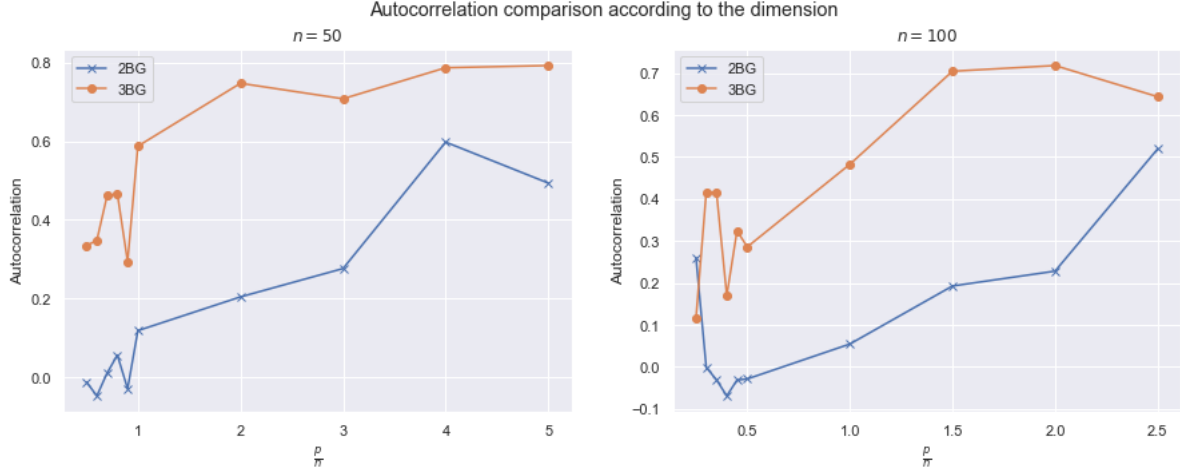


FIGURE 5 – Comparaison d'autocorrélation en fonction de la dimension

La figure 5 montre l'autocorrélation moyenne du décalage de la chaîne σ^2 . Nous pouvons voir que le **2BG** proposé présente des autocorrélations nettement plus faibles que le **3BG** original dans toutes les configurations de (n, p) , ce qui indique un meilleur mélange du **2BG**. En outre, la figure 4 prend en considération le temps de calcul en examinant la taille moyenne effective de l'échantillon par seconde, sur l'échelle logarithmique de chaque échantillonneur. Il est clair que la **2BG** est plus efficace. Par exemple, à $n = p = 50$, la **2BG** produit environ deux fois plus d'échantillons effectifs par seconde que la **3BG**.

4 Conclusion

Dans le modèle bayésien de type Group Lasso, la **2BG** nous semble toujours plus efficace que la **3BG**. L'avantage majeur de la **2BG** est lorsque le nombre de prédicteurs p est supérieur à la taille de l'échantillon de données n . Dans ces problèmes de régression à haute dimension, la **2BG** produit environ **10** à **100** fois plus d'échantillons efficaces par seconde que son concurrent. Le taux de mélange rapide de la **2BG** et le faible coût de calcul par itération de la **2BG**, qui n'est pas supérieur à celui de la **3BG**, contribuent à cet avantage.

Pour aller plus loin, une comparaison avec la méthode fréquentiste du GroupLasso classique, dont l'estimateur est donné par la formule 2, est disponible dans notre Notebook Python joint à ce rapport. Nous verrons notamment que l'efficacité de notre méthode **2BG** (et même **3BG**) dans le cas en grande dimension est largement comparable à l'estimateur fréquentiste, pour le cas de l'estimation de β (pas seulement σ , paramètre sur lequel nous nous sommes intéressés pour le cas en grande dimension dans ce papier).