

Un algorithme de Hastings–Metropolis avec apprentissage séquentiel

Didier CHAUVEAU, Pierre VANDEKERKHOVE

Analyse et mathématiques appliquées, Université de Marne-la-Vallée, 5, boulevard Descartes, Champs-sur-Marne, 77454 Marne-la-Vallée cedex 2, France
Courriel : {chauveau,vandek}@math.univ-mlv.fr

(Reçu le 6 janvier 1999, accepté après révision le 15 avril 1999)

Résumé. Nous proposons un nouvel algorithme de type Hastings–Metropolis, dans lequel la loi instrumentale apprend progressivement la densité cible, accélérant ainsi la convergence. L'algorithme utilise comme loi instrumentale au pas n un histogramme de la densité de l'algorithme à cet instant, construit à partir d'un nombre croissant de chaînes parallèles i.i.d. L'algorithme non homogène obtenu converge en n avec une vitesse géométrique plus forte qu'un algorithme de Hastings–Metropolis classique utilisant une loi instrumentale arbitraire. © Académie des Sciences/Elsevier, Paris

A Hasting–Metropolis algorithm with learning proposal

Abstract. We present a non-homogeneous Hastings–Metropolis algorithm for which the proposal density approximates the target density, as the number of iterations increases. The proposal at the n -th step is a non-parametric estimate of the density of the algorithm, and uses an increasing number of i.i.d. copies of the Markov chain. The resulting algorithm converges (in n) geometrically faster than a Hastings–Metropolis algorithm with an arbitrary proposal. © Académie des Sciences/Elsevier, Paris

1. Introduction

L'algorithme de Hastings–Metropolis (voir [2]) permet de générer une chaîne de Markov $(x^{(n)})$ homogène de loi stationnaire de densité f , à partir d'une densité $q(y|x)$ appelée *loi instrumentale*, et de la connaissance analytique de f à une constante près. Les propriétés d'ergodicité et de convergence ont été largement étudiées dans la littérature (voir [5] pour un résumé), et Holden [3] a proposé récemment une mesure qualitative de la vitesse de convergence en terme de proximité de q à f , sous des hypothèses classiques assurant la convergence géométrique sur un compact.

Nous supposons disposer ici d'une loi instrumentale q_0 assurant cette convergence géométrique avec une vitesse déterminée. L'idée naturelle que nous présentons consiste à remplacer, après quelques itérations, q_0 par des lois instrumentales q_n approchant progressivement f , en exploitant la connaissance de f dont on dispose au travers de la densité p^n de la loi de l'algorithme à l'itération n .

Note présentée par Paul DEHEUVELS.

La convergence de p^n vers f , assurée pour q_0 , suggère de prendre pour q_n un estimateur non paramétrique de p^n construit à partir de $N(n)$ copies i.i.d. de $(x^{(n)})$. L'algorithme inhomogène résultant converge avec une vitesse géométrique meilleure que celle donnée par q_0 .

2. Un algorithme de Hastings–Metropolis inhomogène

Soient $\Omega \subset \mathbb{R}^s$ et f la densité cible, strictement positive sur Ω , minorée par une constante $\alpha > 0$. Cette hypothèse sur f (déjà présente dans [3]) est restrictive, mais l'extension à des situations plus générales est à l'étude, et l'algorithme présenté ici est utilisable empiriquement dans les applications. Nous donnons dans cette section le comportement théorique d'une version de l'algorithme de Hastings–Metropolis utilisant une suite q_n de lois instrumentales (que nous précisons à la section 3), résultant en un algorithme inhomogène de mesure invariante de densité f (l'utilisation des résultats de Holden [3] dans un contexte inhomogène différent a également été proposé par Vandekerkhove [4]). Une itération $x^{(n)} \rightarrow x^{(n+1)}$ de l'algorithme de loi instrumentale q_n , pour une chaîne initialisée suivant $x^{(0)} \sim p^0$, est :

1. tirer $y \sim q_n(\cdot)$
2. calculer $\alpha_n(y, x^{(n)}) = \min \left\{ 1, \frac{f(y) q_n(x^{(n)})}{f(x^{(n)}) q_n(y)} \right\}$
3. prendre $x^{(n+1)} = \begin{cases} y & \text{avec probabilité } \alpha_n(y, x^{(n)}), \\ x^{(n)} & \text{avec probabilité } 1 - \alpha_n(y, x^{(n)}). \end{cases}$

Nous utiliserons des densités q_n strictement positives sur Ω . Posons :

$$Q_n(x, y) = \min \left\{ q_n(x); \frac{q_n(y) f(x)}{f(y)} \right\}, R^n(x) = \left(\frac{p^n(x)}{f(x)} - 1 \right) \text{ et } R_M^n = \sup_{x \in \Omega} \{ |R^n(x)| \}.$$

LEMME 1. – Soit q_n une densité sur Ω . Alors, pour tout $y \in \Omega$, $\int_{\Omega} Q_n(x, y) dx \leq 1$, et

$$R^{n+1}(y) = R^n(y) \left[1 - \int_{\Omega} Q_n(x, y) dx \right] + \int_{\Omega} R^n(x) Q_n(x, y) dx \quad (1)$$

Démonstration. – Elle est semblable à celle donnée dans [3] pour le cas homogène, avec ici une loi instrumentale q_n et une quantité Q_n dépendante du pas de l'algorithme. \square

PROPOSITION 1. – Si, au pas n , $q_n(x) \geq a_n f(x)$ pour tout $x \in \Omega$, avec $a_n \in]0, 1[$. Alors,

$$\left| \frac{p^{n+1}(y)}{f(y)} - 1 \right| \leq (1 - a_n) \sup_{x \in \Omega} \left\{ \left| \frac{p^n(x)}{f(x)} - 1 \right| \right\}. \quad (2)$$

Démonstration. – Analogue à celle de Holden [3], en utilisant (1) et en vérifiant que la condition $q_n(x) \geq a_n f(x)$ implique $Q_n(x, y) \geq a_n f(x)$. \square

3. L'histogramme comme densité de la variable instrumentale

Considérons une infinité de copies i.i.d. de l'algorithme de Hastings–Metropolis inhomogène vu à la section précédente et défini pour une suite de lois instrumentales de densités q_n construites de la manière suivante : soit $T = (0, t_1, \dots, t_i, t_{i+1}, \dots)$ une partition arbitraire de \mathbb{N} . On définit alors la suite q_n relativement à T par :

$$q_n(x) = \left(H_{N(t_i)}(x) \mathbb{1}_{C(t_i)} + \tilde{H}_{N(t_i)}(x) \mathbb{1}_{\bar{C}(t_i)} \right) \mathbb{1}_{n \in [t_i, t_{i+1}[}, \quad (3)$$

où $H_{N(0)}$ est une loi q_0 arbitraire satisfaisant $q_0(x) \geq a_0 f(x)$, $x \in \Omega$, avec $a_0 \in]0, 1[$; $H_{N(t_i)}$ est l'histogramme de la densité de l'algorithme de Hastings–Metropolis inhomogène à l'instant t_i basé sur $N(t_i)$ copies i.i.d. de l'algorithme, et $\tilde{H}_{N(t_i)}$ une modification de cet histogramme si une certaine condition $C(t_i)$ n'est pas vérifiée par l'échantillon. Notons que les chaînes utilisées pour construire $H_{N(t_i)}$ au pas t_i doivent ensuite être éliminées, afin de préserver l'indépendance entre les chaînes restantes.

Rappelons la définition de l'histogramme dans notre cadre. Nous identifierons sans perdre en généralité notre compact Ω à $[0, 1]^s$. Soit $\pi_{N(n),q}$ une partition rectangulaire de $[0, 1]^s$ définie par :

$$\pi_{N(n),q} = [(q_1 - 1)h_{N(n)}, q_1 h_{N(n)}] \times \cdots \times [(q_s - 1)h_{N(n)}, q_s h_{N(n)}], \quad (4)$$

où $q = (q_1, \dots, q_s) \in \mathbb{Z}^s$ et $h_{N(n)}$ est un réel positif dépendant de $N(n)$ désignant la largeur des classes de la partition et destiné à tendre vers 0. L'histogramme s'écrit alors

$$H_{N(n)}(x) = \sum_q \frac{\mu_{N(n)}(\pi_{N(n),q})}{N(n)h_{N(n)}^s} \mathbb{1}_{\pi_{N(n),q}}(x), \quad (5)$$

où $\mu_N(B)$ représente le nombre de points de l'échantillon qui appartiennent au borélien B . Il peut arriver que dans certaines classes le résultat de ce comptage soit nul. Cette éventualité est marquée par la condition $\bar{C}(n)$ (et $C(n)$ pour son contraire) dans l'expression (3). Faute de pouvoir prendre exactement l'histogramme pour q_n lorsque $\bar{C}(n)$ est vérifiée, car alors la condition de minoration de la proposition 1 est violée, nous proposons de modifier légèrement l'histogramme de manière à le rendre positif partout, et notons cet histogramme modifié $\tilde{H}_{N(n)}$. Nous avons alors la proposition suivante, dont la démonstration est directe à partir de (2) :

PROPOSITION 2. – Pour la suite de densités (q_n) définies en (3) il existe une suite de nombres (a_n) , $a_n \in]0, 1[$, telle que $q_n(x) \geq a_n f(x)$, $x \in \Omega$, et $|R^{n+1}(y)| \leq \prod_{k=1}^n (1 - a_k) R_M^0$.

La proposition ci-dessous donne une inégalité exponentielle à distance finie pour l'histogramme $H_{N(n)} = H_N$ basé sur N réalisations i.i.d. de p^n :

PROPOSITION 3. – Soient f une densité strictement positive C -lipschitzienne sur Ω , H_N l'histogramme donné par (5), et $\varepsilon > 0$. Posons $\delta_{N,n} = 2\alpha(1 - a_0)^n R_M^0 + \sqrt{s} h_N C$. Alors les conditions

$$h_N \rightarrow 0, \quad N h_N^{3s} \geq (20/(\varepsilon - \delta_{N,n})^2) \quad \text{pour } N > N_0, \quad n > n_0, \quad (6)$$

où n_0 et N_0 sont tels que $(\varepsilon - \delta_{N_0, n_0}) > 0$ et $(\varepsilon - \delta_{N_0, n_0}) h_{N_0}^s \leq 1$, entraînent :

$$\mathbb{P} \left[\sup_{x \in \Omega} |H_N(x) - p^n(x)| > \varepsilon \right] \leq 3 \exp \left[-N h_N^{2s} (\varepsilon - \delta_{N,n})^2 / 25 \right], \quad \text{pour } n > n_0 \text{ et } N > N_0. \quad (7)$$

Démonstration. – On identifie comme précédemment Ω à $[0, 1]^s$. On a :

$$\sup_{x \in \Omega} |H_N(x) - p^n(x)| \leq \sup_{x \in \Omega} |H_N(x) - \mathbb{E}[H_N(x)]| + \sup_{x \in \Omega} |\mathbb{E}[H_N(x)] - p^n(x)|.$$

Le premier terme est contrôlé à partir d'une inégalité exponentielle pour la loi multinomiale ([1], p. 174), qui exige que h_N ne tende pas trop vite vers 0 (condition (6)), et que N et n soient assez grands. Le second terme est contrôlé grâce à l'hypothèse f lipschitzienne et à la convergence géométrique de p^n vers f avec vitesse $(1 - a_0)^n$ qui permet d'avoir, sur chaque classe $\pi_{N,q}$:

$$\sup_{x \in \pi_{N,q}} |\mathbb{E}[H_N(x)] - p^n(x)| \leq \sup_{x,y \in \pi_{N,q}} |p^n(x) - p^n(y)| \leq \delta_{N,n}. \quad \square$$

4. Optimalité de l'algorithme

Nous montrons que l'algorithme avec apprentissage séquentiel converge plus rapidement vers f , en n , que tout algorithme de Hastings–Metropolis homogène usuel utilisant la loi instrumentale arbitraire q_0 (cet algorithme classique converge avec la vitesse $(1 - a_0)^n$ vers f , voir [3]). Au pas n de l'algorithme séquentiel, q_n est soit $H_{N(n)}$, soit $\tilde{H}_{N(n)}$ (on suppose ici que $T = (0, 1, 2, \dots)$ pour alléger les notations). Posons alors $A_n = \{\mathbb{1}_{C(n)} = 0\}$, $B_n = \{a_n < a_0 \mid A_n^c\}$ et $C_n = \{a_n < a_0\}$, l'événement B_n signifiant que q_0 est un meilleur candidat au sens de (2) que l'histogramme $H_{N(n)}$ choisi au pas n , et C_n signifiant que q_0 est un meilleur candidat que q_n . La proposition ci-dessous exprime le fait que l'on utilisera infiniment souvent une loi instrumentale meilleure que q_0 (plus proche de f et donc de constante de minoration $a_n > a_0$), obtenant ainsi un algorithme plus rapide.

PROPOSITION 4. – Si le nombre $N(n)$ de copies i.i.d. de l'algorithme séquentiel, et la finesse $h_{N(n)}$ de la partition utilisés pour construire $H_{N(n)}$ vérifient les conditions de la proposition 3, et

$$N(n)h_{N(n)}^{2s} \geq c \log(n), \quad (8)$$

où c est une constante calculable, alors $\mathbb{P}(\overline{\lim} A_n) = \mathbb{P}(\overline{\lim} B_n) = \mathbb{P}(\overline{\lim} C_n) = 0$, et l'on a :

$$\mathbb{P} \left[\overline{\lim}_{n \rightarrow \infty} \left\{ \prod_{k=1}^n (1 - a_k) > (1 - a_0)^n \right\} \right] = 0. \quad (9)$$

Démonstration. – On se ramène aux probabilités de déviation pour l'histogramme :

$$\mathbb{P}(A_n) \leq \mathbb{P} \left[\sup_{x \in \Omega} |H_{N(n)}(x) - p^n(x)| > \alpha \left(1 - \prod_{k=1}^{n-1} (1 - a_k) R_M^0 \right) \right], \quad (10)$$

et (7) permet de conclure grâce à un argument de Borel–Cantelli car $\sum_{n \geq 0} \mathbb{P}(A_n) < +\infty$ pour une vitesse adaptée $N(n)h_{N(n)}^{2s} \geq c_{N,n}(\alpha) \log(n)$, où $c_{N,n}(\alpha) \approx 25\beta/\alpha^2$, $\beta > 1$. Une expression analogue à (10) est obtenue pour B_n , enfin $\mathbb{P}(C_n) \leq \mathbb{P}(A_n) + (1 - \mathbb{P}(A_n))\mathbb{P}(B_n)$. On en déduit alors (9). \square

5. Mise en œuvre

Une implémentation de cet algorithme approchant la situation théorique est possible en considérant un nombre a priori fixé de chaînes i.i.d., divisé en k paquets. À l'instant t_i , $i = 1, \dots, k$, on utilise le i -ème paquet de $N(t_i)$ chaînes pour construire $H_{N(t_i)}$. L'algorithme inhomogène utilise donc des lois instrumentales qui apprennent f de mieux en mieux, et il devient homogène de loi instrumentale $H_{N(t_k)}$ après t_k (le k -ième paquet pouvant être réduit à une seule chaîne d'un algorithme plus rapide que celui associé à q_0). Cet algorithme semble particulièrement adapté aux situations dans lesquelles la loi f est multimodale. Dans ce cas en effet, les lois instrumentales proposées favoriseront rapidement les sauts entre modes déjà « découverts », accélérant ainsi l'exploration du domaine d'intérêt.

Références bibliographiques

- [1] Bosq D., Lecoutre J.-P., Théorie de l'estimation fonctionnelle, Economica, Paris, 1987.
- [2] Hastings W.K., Monte Carlo sampling methods using Markov Chains and their applications, Biometrika 57 (1970) 97–109.
- [3] Holden L., Geometric Convergence of the Metropolis–Hastings Simulation Algorithm, Preprint, Norwegian Computing Center, Oslo, Norway, 1996.
- [4] Vandekerkhove P., A Sequential Metropolis–Hastings Algorithm, Preprint, Università di Pavia, Italy, 1998.
- [5] Robert C.P., Méthodes de Monte-Carlo par chaînes de Markov, Economica, Paris, 1996.