

Algorithme de Métropolis Hasting Pseudo Marginal : Théorie et implémentation

Enzo Cabaret

I Introduction aux chaines de Markov

L'algorithme de Métropolis Hasting fait partie des méthodes dites "Markov Chains Monte Carlo" (MCMC), ces méthodes permettent de simuler des lois de probabilité dont on ne connaîtrait pas l'intégralité des paramètres. Ces méthodes ont deux aspects, Monte Carlo : il s'agit d'estimer une quantité à partir de tirage aléatoires, Markov Chains: on va chercher à simuler les réalisations d'une chaîne de Markov dont la loi limite est notre loi cible. Dans cette première section on donne les éléments de base pour comprendre les aspects des chaînes de markovs utiles pour appréhender les MCMC.

On se place dans un espace mesurable (X, \mathcal{X})

Definition 1

On dit que $P : X \times \mathcal{X} \rightarrow \mathbb{R}^+$ est un kernel de Markov si $\forall (x, A) \in X, \mathcal{X}$ on a:

- $X \ni y \mapsto P(y, A)$ est $\mathcal{B}(\mathbb{R}^+)$ mesurable
- $\mathcal{X} \ni B \mapsto P(x, B)$ est une mesure de probabilité sur (X, \mathcal{X})

Definition 2

Soit $(X_k)_{k \in \mathbb{N}}$ une séquence de variables aléatoires sur un même espace $(\Omega, \mathcal{G}, \mathbb{P})$ et prenant des valeurs dans X on dit que $(X_k)_{k \in \mathbb{N}}$ est une chaîne de markov si et seulement si on a:

- $\forall (k, A) \in \mathbb{N} \times \mathcal{X}$ on a $\mathbb{P}(X_{k+1} \in A | X_{0:k}) = P(X_k, A)$ \mathbb{P} -ps
- $\mathbb{P}(X_0 \in A) = \nu(A)$ avec $\nu \in M_1(\mathcal{X})$

Notations additionnelles

On introduit ici des notations en lien avec les kernels de markovs pour simplifier l'expression des calculs par la suite.

Soit $\mu \in M_1(\mathcal{X})$, P, Q des kernel de markov et h une fonction bornée sur X on note :

- $\mu P: \mathcal{X} \ni A \mapsto \mu P(A) = \int \mu(dx) P(x, A)$
- PQ le kernel de markov: $(x, A) \mapsto \int_X P(x, dy) Q(y, A)$

Lemma 1

$(X_k)_{k \in \mathbb{N}}$ une chaîne de markov sur $(\Omega, \mathcal{G}, \mathbb{P})$ et qui prend des valeurs dans X avec un kernel P et une distribution initiale ν alors $\forall n \in \mathbb{N}$ la loi de $X_{0:n}$ est $\nu(dx_0) \prod_{i=1}^{n-1} P(x_i, dx_{i+1})$. De même X_n suit la loi νP^n

Proof

Facile. Pour la deuxième assertion il suffit d'intégrer le résultat de la première. ■

Probabilités invariantes

Definition 3

On appelle probabilité invariante pour le kernel de Markov P toute mesure π vérifiant : $\pi = \pi P$

Definition 4

Soit π une mesure de probabilité et P un kernel de markov on dit que P est π -reversible si

$$\pi(dx)P(x, dy) = \pi(dy)P(y, dx) \quad (1)$$

Proposition 2

Soit P un kernel de Markov et π une probabilité. Si P est π -reversible alors π est invariante pour P .

II Algorithme d'Acceptation Rejet et Métropolis Hastings

Rappels des simulations de variables aléatoires

On note ici X une variable aléatoire (v.a.) réelle et F sa fonction de répartition : $F(x) := \mathbb{P}(X \leq x)$

On rappelle que F est croissante, continue à droite et limité à gauche.

Definition 5

On appelle inverse généralisée de F ou fonction quantile la fonction définie par :

$$\forall u \in]0, 1[\quad F^{-1}(u) = \inf x \in \mathbb{R} \mid F(x) \geq u \quad (2)$$

Theorem 3

Soient X une v.a. réelle et U une v.a. de loi uniforme sur $]0, 1[$. Alors $F^{-1}(U)$ et X ont la même loi.

Proof

(si on a le temps de l'écrire) ■

Algorithme d'acceptation rejet

La méthode d'acceptation rejet s'explique simplement. On suppose que l'on sait générer un vecteur aléatoire M_1 de \mathbb{R}^d suivant une loi μ . On génère une suite iid M_1, \dots, M_n, \dots en s'arrêtant au premier d'entre eux qui vérifie une certaine condition \mathcal{H}_0 . On note T l'indice correspondant. le vecteur aléatoire M_T suit une nouvelle loi ν . Si l'on parvient facilement à simuler M et à vérifier la condition \mathcal{H}_0 on possède donc une méthode pour simuler une variable suivant la loi ν .

Proposition 4

Soit $(M_n)_{n \in \mathbb{N}}$ une suite de vecteur aléatoires $\Omega \rightarrow \mathbb{R}^d$ indépendant de même loi μ . Soit $B \subset \mathbb{R}^d$ un borélien tel que $\mu(B) > 0$. $\forall \omega \in \Omega$ on pose :

$$T(\omega) \triangleq \inf \{i \in \mathbb{N}^*; M_i(\omega) \in B\} \quad (3)$$

Avec la convention $\inf \emptyset = +\infty$. On définit $M_T : \Omega \rightarrow \mathbb{R}^d$ par

$$M_T(\omega) \triangleq \begin{cases} M_T(\omega) & \text{si } T(\omega) < +\infty \\ 0 & \text{sinon} \end{cases} \quad (4)$$

Dans ces conditions :

- a) T est une va suivant la loi géométrique de paramètre $p = \mu(B)$
- b) M_T est un vecteur aléatoire de loi ν donnée par

$$\forall A \in \text{Bor}(\mathbb{R}^d), \quad \nu(A) = \frac{\mu(A \cap B)}{\mu(B)} \quad (5)$$

Proposition 5

Soit f une densité de probabilité sur \mathbb{R}^d et G son hypographe. Soit $M = (Z, Y)$ un vecteur aléatoire de $\mathbb{R}^d \times \mathbb{R}$ de loi uniforme sur G . Alors Z a pour densité f par rapport à λ_d

Proposition 6

Soit $X : \Omega \rightarrow \mathbb{R}^d$ un vecteur aléatoire dont la loi a pour densité g par rapport à λ_d . Posons

$$M \triangleq (X, cg(X)U) \quad (6)$$

où U est une variable aléatoire uniforme sur $[0, 1]$ indépendante de X et $c > 0$ une constante. on note H l'hypographe de cg . alors M suit la loi uniforme sur H .

Dans le cas d'une variable à densité l'algorithme d'acceptation rejet se décrit de la manière suivante. On suppose que l'on sait simuler un vecteur de densité g et on souhaite simuler un autre vecteur de densité f . Supposons que l'on ait c une constante positive telle que $f < cg$. On a $\lambda_{d+1}(G) = \int_{\mathbb{R}^d} f d\lambda_d = 1$ et $\lambda_{d+1}(H) = \int_{\mathbb{R}^d} cg d\lambda_d = c$. On en déduit que $\lambda_{d+1}(H \setminus G) = \int_{\mathbb{R}^d} cg - f d\lambda_d$, ce qui nous permet de considérer que $\lambda_{d+1}(H \setminus G) > 0$ car le cas échéant $f = cg = g$ car f, g sont des densités de probabilité. Par hypothèse on sait simuler $(X_i)_{i \geq 1}, (U_i)_{i \geq 1}$, la première car c'est un vecteur de densité g la deuxième car elle est de loi uniforme sur $[0, 1]$ en considérant $M_i \triangleq (X_i, g(X_i)U_i)$ qui est uniforme sur H par la proposition 6. Posons:

$$T(\omega) \triangleq \inf\{i \in \mathbb{N}^*; M_i(\omega) \in G\} \quad (7)$$

et

$$M_T(\omega) \triangleq \begin{cases} M_T(\omega) & \text{si } T(\omega) < +\infty \\ 0 & \text{sinon} \end{cases} \quad (8)$$

Les propositions 4 et 5 nous permettent alors de conclure que $Z = X_T$ est un vecteur aléatoire de densité f

Algorithme de Métropolis Hasting

Dans les méthodes MCMC nous allons chercher à simuler les réalisations d'une chaîne de Markov de kernel P acceptant pour probabilité invariante la loi que l'on cherche à simuler π . Toute la problématique revient à trouver le kernel P adéquat. L'approche de Métropolis Hasting nous permet de construire ce kernel.

Comme dans la méthode d'acceptation rejet on suppose que l'on dispose d'une densité nous permettant de générer des candidats pour l'état suivant de notre chaîne de Markov y à l'état courant x que l'on note $q(x, y)$. Si q vérifie la condition de réversibilité (insérer ref) notre objectif est déjà atteint. Cependant il n'y a pas de raison que ce soit le cas. Supposons $x, y \in X$ tels que :

$$\pi(x)q(x, y) > \pi(y)q(y, x) \quad (9)$$

Dans ce cas la notre chaîne de Markov bouge de x vers y trop souvent on introduit donc une probabilité $\alpha(x, y) < 1$ pour que y soit pris comme état suivant moins fréquemment. Cela revient à choisir comme kernel de Markov le kernel de densité.

$$p_{MH}(x, y) = q(x, y)\alpha(x, y) \text{ avec } x \neq y \quad (10)$$

On souhaite donc introduire $\alpha(x, y)$ et $\alpha(y, x)$ tels que l'équation d'équilibre (insérer ref) soit vérifiée cad :

$$\pi(x)q(x, y)\alpha(x, y) = \pi(y)q(y, x)\alpha(y, x) \quad (11)$$

ce qui nous donne l'identité suivante pour définir α :

$$\alpha(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}\alpha(y, x) \quad (12)$$

Dans l'algorithme de métropolis hasting on cherche à ce que $\alpha(x, y)$ soit maximale pour une exploration des nouveaux états plus rapides ce qui revient donc à :

$$\alpha(x, y) = \min\left(\frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}, 1\right) \quad (13)$$

Pour finir la définition du Kernel de Métropolis Hasting on doit considérer le cas ou notre chaîne de Markov reste dans le même état car actuellement on se plaçait dans le cas $x \neq y$. En utilisant notre raisonnement précédent la probabilité que la chaîne reste dans le même état est :

$$r(x) = 1 - \int_X q(x, y)\alpha(x, y)dy \quad (14)$$

Finalement on définit le kernel de Métropolis Hasting par :

$$P_{MH}(x, y) = q(x, y)\alpha(x, y)dy + \left[1 - \int_{\mathcal{X}} q(x, y)\alpha(x, y)dy\right] \delta_x(dy) \quad (15)$$

Ce qui nous amène à l'implémentation de l'algorithme de Métropolis Hasting :

- Initialiser x_0 arbitrairement.
- Pour $i = 0, \dots, n - 1$ faire :
 - Tirer y à avec $q(x, \cdot)$ et u avec $\mathcal{U}(0, 1)$
 - Si $u < \alpha(x_i, y)$ alors $x_{i+1} = y$
 - Sinon $x_{i+1} = x_i$
- Fin Pour
- Retourner x_0, \dots, x_n

Nous illustrons l'algorithme par un exemple simple :

On se place dans le cas d'une distribution normale centrée réduite perturbée, le problème ici c'est que l'on ne peut pas exprimer la constante de normalisation toutefois l'algorithme de métropolis hasting nous permet quand meme de générer des tirages.

- fonction cible : $\pi(x) \propto \sin^2(x) \times \sin^2(2x) \times \exp(-x^2/2)$
- kernel : $q : x \mapsto \mathcal{U}(x - \mu, x + \mu)$ avec $\mu > 0$

In [2]:

```
import matplotlib.pyplot as plt
import numpy as np
import numpy.random as npr
import scipy.stats as stats
from scipy.stats import expon, geom, norm
from math import pi
from pylab import *
% matplotlib inline
```

UsageError: Line magic function `%` not found.

In [3]:

```
import numpy as np
accpeted = []
def pi(x):
    return (np.sin(x)**2)*(np.sin(2*x)**2)*np.exp(-(x**2)/2)

def metropolis_hastings(x_0,n,mu=1):
    X = np.empty(n)
    X[0] = x_0
    for i in range(n-1):
        y = np.random.uniform(X[i]-mu,X[i]+mu)
        u = np.random.uniform(0,1)
        if u < min(pi(y)/pi(X[i]),1):
            X[i+1] = y
        else:
            X[i+1] = X[i]
    return X
```

In [4]:

```
import matplotlib.pyplot as plt

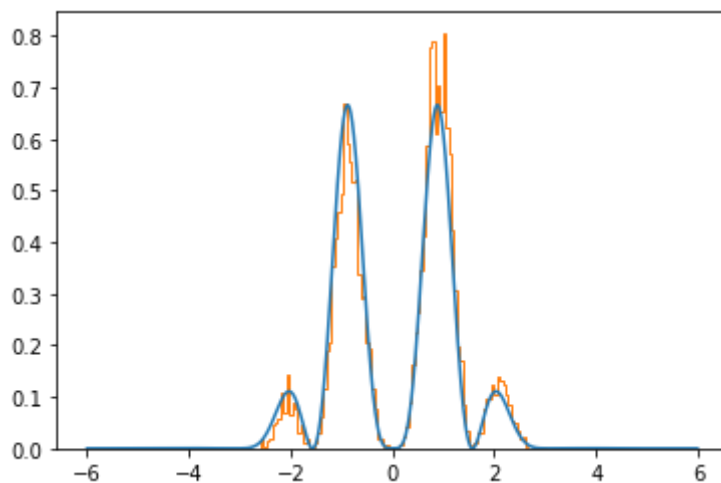
#paramètres
x_0 = 3
n = 10**4

x_MH = metropolis_hastings(x_0,n)

#graphe de la fonction cible :
x = np.arange(-6,6,0.01)
y = pi(x)
y = [i/(sum(y)*0.01) for i in y]
plt.plot(x,y)

#histogramme des échantillons obtenus par MH
n,bins,patches = plt.hist(x_MH,100,density = True,histtype = 'step')

plt.show()
```



II. Metropolis Hastings Méthode pseudo-marginales

Généralisation de l'algorithme MH

Soit $\pi \in M_1(X)$ et Q un noyau de Markov sur $X \times \mathcal{X}$. Dans la section précédente on a regardé l'algorithme Metropolis-Hastings quand π et $Q(x)$ ont toutes les deux des densités par rapport à λ . Dans cette partie on ne fait pas cette supposition et l'expression de α^{MH} , n'est plus disponible. La supposition que nous feront est la suivante :

$$\mu_0(dxdy) = \pi(dx)Q(x, dy) \quad (16)$$

$$\mu_1(dxdy) = \pi(dy)Q(x, dy) \quad (17)$$

Proposition 7

Il existe une fonction $(x, y) \Rightarrow r(x, y)$ tel que $r(x, y) > 0, \mu_0$ a.s et $\forall h \in F_+(X^2)$, $\int h(x, y)\mu_1(dx, dy) = \int h(x, y)r(x, y)\mu_0(dx, dy)$

On montre grâce à cette propriété que ces deux mesures sont équivalentes. Dans ce cas, on remplace dans la version du noyau MH où α^{MH} est donnée dans le **Lemme 5**, par $\alpha(x, y) = r(x, y) \wedge 1$ qui est π -reversible.

Lemma 8

On assumant la propriété 1, on pose $\alpha(x, y) = r(x, y)$, le noyau de Markov est alors : $P_{\langle \pi, Q \rangle}^{MH} = Q(x, dy)\alpha(x, y) + \bar{\alpha}(x)\delta_x(dy)$ ou $\bar{\alpha}(x) = 1 - \int_X Q(x, dy)\alpha(x, y)$ est π -reversible

Idée de la preuve : Détailler la condition d'équilibre (1).

Techniques pseudo marginales de Monte Carlo

Proof

On assume dans cette partie que π et Q sont dominées par une mesure commune λ et on note $\pi(dx) = \pi(x)\lambda(dx)$ et $Q(x, y) = q(x, y)\lambda(dx)$.

On a normalement besoin (à une constante multiplicative près) d'une expression explicite de la distribution $\pi(x) \forall x \in X$ comme distribution cible de l'algorithme de Metropolis-Hastings. Néanmoins il peut arriver qu'on ne puisse pas calculer la distribution $\pi(x)$ même pas à une constante multiplicative près. Dans ce cas là, **on assume alors que l'on peut avoir un estimateur non-biaisé de la distribution $\pi(x)$** .

Pour obtenir cet estimateur non-biaisé, on tire $W \sim R(x, dw)$ où R est un noyau de Markov de $X \longrightarrow \mathbb{R}^+$ tel que la condition d'estimateur non-biaisé soit respectée. ■

On obtient alors l'algorithme MH Pseudo marginal suivant :

Input : n

Output : X_0, \dots, X_n

À $t = 0$, tirage de X_0 selon une distribution arbitraire $W_0 \sim R(X_0, \cdot)$

for $t \leftarrow 0$ **to** $n - 1$ **do**

• Tirage de $\tilde{X}_{t+1} \sim Q(X_t, \cdot)$ et après $\tilde{W}_{t+1} \sim \tilde{R}(X_{t+1}, \cdot)$

• On pose $(X_{t+1}, W_{t+1}) = \begin{cases} (\tilde{X}_{t+1}, \tilde{W}_{t+1}) \text{ avec prob. } \frac{\tilde{W}_{t+1} q(\tilde{X}_{t+1}, X_t)}{W_t q(X_t, \tilde{X}_{t+1})} \wedge 1 \\ (X_t, W_t) \text{ avec prob. } 1 - \frac{\tilde{W}_{t+1} q(\tilde{X}_{t+1}, X_t)}{W_t q(X_t, \tilde{X}_{t+1})} \wedge 1 \end{cases}$

Cette version de l'algorithme de Metropolis-Hastings est très proche de la version du chapitre précédent mais en remplaçant $\pi(x)$ par un estimateur non-biaisé. On justifie l'utilisation des techniques Pseudo-marginales de Monte-Carlo en montrant que c'est en fait un algorithme MH généralisé en considérant la chaîne de Markov $(\bar{X}_k)_{k \in \mathbb{N}} = (X_k, W_k)_{k \in \mathbb{N}}$ sur un nouvel espace $\bar{X} = X \times \mathbb{R}_*^+$, avec une nouvelle distribution cible $\Pi(d\bar{x}) = \Pi(dx dw) = w R(x, dw) \lambda(dx)$. On pose $\mu_0(d\bar{x}, d\bar{x}') = w R(x, dw) \lambda(dx) Q(x, dx') R(x', dw')$ et $\mu_1(d\bar{x}, d\bar{x}') = w' R(x', dw') \lambda(dx) Q(x', dx) R(x, dw)$ pour satisfaire la propriété 1. On écrit ensuite $Q(x, dy) = q(x, y) \lambda(dy)$, on obtient $\forall h \in F_X(\bar{X}^2)$:

$$\int_{\bar{X}^2} h(\bar{x}, \bar{x}') \mu_1(d\bar{x}', d\bar{x}) = \int_{\bar{X}^2} h(\bar{x}, \bar{x}') r(\bar{x}, \bar{x}') \mu_0(dx, dy) \quad (18)$$

Comme $r > 0$, on applique le lemme 1 avec $\alpha(\bar{x}, \bar{x}') = r(\bar{x}, \bar{x}') \wedge 1$ on a que $P_{\langle \pi, Q \rangle}^{MH}(\bar{x}, d\bar{x}')$ est Π -reversible. En appliquant l'algorithme 2 on a que $P_{\langle \pi, Q \rangle}^{MH}(\bar{x}, d\bar{x}')$ est un noyau de Markov. **En conclusion :**
 $(\bar{X}_k)_{k \in \mathbb{N}} = (X_k, W_k)_{k \in \mathbb{N}}$ produite par l'algorithme MH Pseudo marginal où la distribution cible est Π admet π comme distribution marginale pour X_k , $(X_k)_{k \in \mathbb{N}}$ n'est plus une chaîne de Markov mais $(\bar{X}_k)_{k \in \mathbb{N}}$ oui.

Méthode pseudo marginale avec un mélange de k gaussiennes

Mélange de gaussienne avec un estimateur sans biais

In [19]:

```
#gaussian mixture of 5 gaussian
mu = [-3,0,3]
sig = [0.5,0.5,0.5]

def pi(x):
    result = 0
    for m,s in zip(mu,sig):
        result+= np.exp(-(x-m)**2)/(2*s**2)
    return result

#estimateur sans biais de pi(x)
def unbiased_estimator(pi,x):
    return np.random.exponential()*pi(x)

#pseudo marginal metropolis hasting
def metropolis_hastings(x_0,n,pi_hat,mu=2.5):
    X = np.empty(n)
    X[0] = x_0
    for i in range(n-1):
        y = np.random.uniform(X[i]-mu,X[i]+mu)
        u = np.random.uniform(0,1)
        if u < min(pi_hat(y)/pi_hat(X[i]),1):
            X[i+1] = y
        else:
            X[i+1] = X[i]
    return X
```

In [20]:

```
pi_hat = lambda x: unbiased_estimator(pi,x)

x_MH = metropolis_hastings(0,10**4,pi_hat)

#graphe de la fonction cible :
x = np.arange(-6,6,0.01)
y = pi(x)
y = [i/(sum(y)*0.01) for i in y]
plt.plot(x,y)

#hisogramme des échantillons obtenus par MH
n,bins,patches = plt.hist(x_MH,100,density = True,histtype = 'step')

plt.show()
```

