

# Limit theorems for adaptive MCMC algorithms

Gersende FORT

LTCI  
CNRS - TELECOM ParisTech

In collaboration with Yves ATCHADE (Univ. Michigan, US), Eric MOULINES  
(TELECOM ParisTech) and Pierre PRIOURET (Univ. Paris 6).

Markov chain Monte Carlo algorithms (MCMC): algorithms to sample from a **target density**  $\pi$

- ▶ in some applications: known up to a (normalizing) constant.
- ▶ complex, so that exact sampling from  $\pi$  is not possible.

**Markov chain Monte Carlo algorithms (MCMC)**: algorithms to sample from a **target density**  $\pi$

- ▶ in some applications: known up to a (normalizing) constant.
- ▶ complex, so that exact sampling from  $\pi$  is not possible.

Define a Markov chain  $\{X_n, n \geq 0\}$  with transition kernel:  $P$

$$\mathbb{E}[f(X_{n+1})|\mathcal{F}_n] = \int f(y) P(X_n, dy)$$

so that

- ▶ for any bounded function  $f$ :  $\lim_n \mathbb{E}_x[f(X_n)] = \pi(f)$ .
- ▶ for any function  $f$  increasing like  $\dots: n^{-1} \sum_{k=1}^n f(X_k) \xrightarrow{a.s.} \pi(f)$ .
- ▶ ...

## I. Adaptive MCMC:

- ▶ why?
- ▶ does the process  $\{X_n, n \geq 0\}$  approximate  $\pi$ ?

# 1.1. Symmetric Random Walk Hastings-Metropolis algorithm

An example of transition kernel  $P$  is described by the algorithm:

- ▶ Choose: a proposal density  $q$
- ▶ Iterate: starting from  $X_n$ 
  - ▶ draw (an increment)  $Y_{n+1} \sim q(\cdot)$
  - ▶ compute the acceptance ratio

$$\alpha(X_n, X_n + Y_{n+1}) := 1 \wedge \frac{\pi(X_n + Y_{n+1})}{\pi(X_n)}$$

- ▶ set

$$X_{n+1} = \begin{cases} Y_{n+1} + X_n & \text{with probability } \alpha(X_n, X_n + Y_{n+1}) \\ X_n & \text{with probability } 1 - \alpha(X_n, X_n + Y_{n+1}) \end{cases}$$

# 1.1. Symmetric Random Walk Hastings-Metropolis algorithm

An example of transition kernel  $P$  is described by the algorithm:

- ▶ Choose: a proposal density  $q$
- ▶ Iterate: starting from  $X_n$ 
  - ▶ draw (an increment)  $Y_{n+1} \sim q(\cdot)$
  - ▶ compute the acceptance ratio

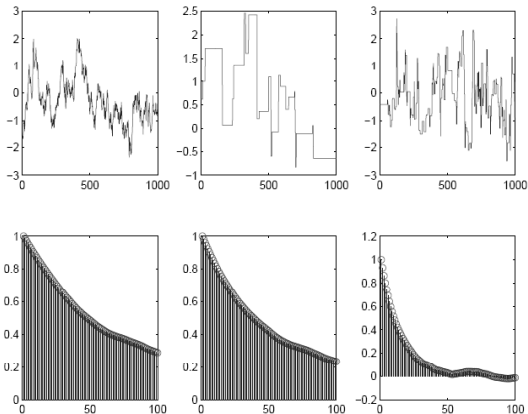
$$\alpha(X_n, X_n + Y_{n+1}) := 1 \wedge \frac{\pi(X_n + Y_{n+1})}{\pi(X_n)}$$

- ▶ set

$$X_{n+1} = \begin{cases} Y_{n+1} + X_n & \text{with probability } \alpha(X_n, X_n + Y_{n+1}) \\ X_n & \text{with probability } 1 - \alpha(X_n, X_n + Y_{n+1}) \end{cases}$$

The efficiency of the algorithm depends upon the proposal  $q$

## 1.2. On the choice of the variance of the proposal distribution



For ex., when  $q$  is Gaussian, how to choose its variance matrix  $\Sigma_q$  ?

- ▶ When  $\pi \sim \mathcal{N}_d(\mu_\pi, \Sigma_\pi)$ , the optimal choice for the variance of  $q$  is

$$\Sigma_q = \frac{(2.38)^2}{d} \Sigma_\pi.$$

*Results obtained by the 'scaling' technique (see also 'fluid limit'). Generalizations exist (other MCMC; relaxing conditions on  $\pi$ )*

ROBERTS-ROSENTHAL (2001); BÉDARD (2007); FORT-MOULINES-PRIOURET (2008).

- ▶ This suggests an **adaptive** procedure: learn  $\Sigma_\pi$  "on the fly" and modify the variance  $\Sigma_q$  **continuously during the run** of the algorithm.



- ▶ When  $\pi \sim \mathcal{N}_d(\mu_\pi, \Sigma_\pi)$ , the optimal choice for the variance of  $q$  is

$$\Sigma_q = \frac{(2.38)^2}{d} \Sigma_\pi.$$

*Results obtained by the 'scaling' technique (see also 'fluid limit'). Generalizations exist (other MCMC; relaxing conditions on  $\pi$ )*

ROBERTS-ROSENTHAL (2001); BÉDARD (2007); FORT-MOULINES-PRIOURET (2008).

- ▶ This suggests an **adaptive** procedure: learn  $\Sigma_\pi$  "on the fly" and modify the variance  $\Sigma_q$  **continuously during the run** of the algorithm.

**Example**: at each iteration, choose  $q$  equal to

$$0.95 \mathcal{N}\left(0, (2.38)^2 d^{-1} \hat{\Sigma}_n\right) + 0.05 \mathcal{N}\left(0, (0.1)^2 d^{-1} \mathbb{I}_d\right)$$

where

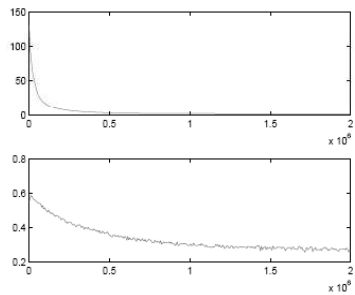
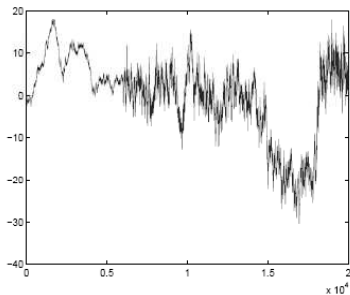
$$\hat{\Sigma}_n = \hat{\Sigma}_{n-1} + \frac{1}{n} \left( \{X_n - \mu_n\} \{X_n - \mu_n\}^T - \hat{\Sigma}_{n-1} \right)$$

$$\mu_n = \mu_{n-1} + \frac{1}{n} (X_n - \mu_{n-1})$$

# Limit theorems for adaptive MCMC algorithms

## └ Motivation

### └ On the choice of the variance of the proposal distribution



## 1.3. Be careful with adaptation !

The previous example illustrates the general framework :

- ▶ Let  $\{P_\theta, \theta \in \Theta\}$  be a family of Markov kernels s.t.  $\pi P_\theta = \pi$  for any  $\theta \in \Theta$ .
- ▶ Define a process  $\{(\theta_n, X_n), n \geq 0\}$  :
  - ▶  $X_{n+1} \sim P_{\theta_n}(X_n, \cdot)$
  - ▶ update  $\theta_{n+1}$  based on  $(\theta_n, X_n, X_{n+1})$  "internal" adaptation

Is it true that the marginal  $\{X_n, n \geq 0\}$  approximates  $\pi$ ?

## 1.3. Be careful with adaptation !

The previous example illustrates the general framework :

- ▶ Let  $\{P_\theta, \theta \in \Theta\}$  be a family of Markov kernels s.t.  $\pi P_\theta = \pi$  for any  $\theta \in \Theta$ .
- ▶ Define a process  $\{(\theta_n, X_n), n \geq 0\}$  :
  - ▶  $X_{n+1} \sim P_{\theta_n}(X_n, \cdot)$
  - ▶ update  $\theta_{n+1}$  based on  $(\theta_n, X_n, X_{n+1})$  "internal" adaptation

Is it true that the marginal  $\{X_n, n \geq 0\}$  approximates  $\pi$  ?

Not always, unfortunately for  $\theta \in ]0,1[$

$$P_\theta = \begin{bmatrix} (1-\theta) & \theta \\ \theta & (1-\theta) \end{bmatrix} \quad \pi = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

Let  $t_1, t_2 \in ]0,1[$ , and set  $\theta_k = t_i$  iff  $X_k = i$ . Then  $\{X_n, n \geq 0\}$  is Markov with invariant probability

$$\tilde{\pi} \propto [t_2 \ t_1]^T \neq \pi$$

## II. Sufficient conditions for convergence of adaptive schemes

$$\{(\theta_n, X_n), n \geq 0\}$$

- ▶ convergence of the marginals  $\{X_n, n \geq 0\}$
- ▶ law of large numbers w.r.t.  $\{X_n, n \geq 0\}$

## 2.1. Convergence of the marginals: Suff Cond

Let

- ▶ a family of Markov kernels  $\{P_\theta, \theta \in \Theta\}$  s.t.  $P_\theta$  has an unique invariant probability measure  $\Pi_\theta$
- ▶ a filtration  $\mathcal{F}_n$  and a process  $\{(X_n, \theta_n), n \geq 0\}$  s.t. for any  $f \geq 0$ ,

$$\mathbb{E}[f(X_{n+1})|\mathcal{F}_n] = \int f(y) P_{\theta_n}(X_n, dy) \quad \mathbb{P} - a.s.$$

Given a target density  $\pi_\star$ , which set of conditions will imply

$$\lim_n \sup_{f, |f|_\infty \leq 1} |\mathbb{E}[f(X_n)] - \pi_\star(f)| = 0 \quad ?$$

Idea :

$$\begin{aligned}\mathbb{E}[f(X_n)] - \pi_*(f) &= \mathbb{E}[\mathbb{E}[f(X_n)|\mathcal{F}_{n-N}]] - \pi_*(f) \\ &= \mathbb{E}\left[\mathbb{E}[f(X_n)|\mathcal{F}_{n-N}] - P_{\theta_{n-N}}^N f(X_{n-N})\right] + \mathbb{E}\left[P_{\theta_{n-N}}^N f(X_{n-N}) - \pi_{\theta_{n-N}}(f)\right] \\ &\quad + \mathbb{E}\left[\pi_{\theta_{n-N}}(f) - \pi_*(f)\right]\end{aligned}$$

Idea :

$$\begin{aligned} \mathbb{E} [f(X_n)] - \pi_*(f) &= \mathbb{E} [\mathbb{E} [f(X_n) | \mathcal{F}_{n-N}] ] - \pi_*(f) \\ &= \mathbb{E} \left[ \mathbb{E} [f(X_n) | \mathcal{F}_{n-N}] - P_{\theta_{n-N}}^N f(X_{n-N}) \right] + \mathbb{E} \left[ P_{\theta_{n-N}}^N f(X_{n-N}) - \pi_{\theta_{n-N}}(f) \right] \\ &\quad + \mathbb{E} [\pi_{\theta_{n-N}}(f) - \pi_*(f)] \end{aligned}$$

i.e. conditions on

- ▶ **(Diminishing Adaptation)** the difference  $\|P_{\theta_n}(x, \cdot) - P_{\theta_{n-1}}(x, \cdot)\|_{\text{TV}}$
- ▶ **(ergodicity of  $P_\theta$  / Containment)** the convergence of  $\|P_\theta^N(x, \cdot) - \pi_\theta\|_{\text{TV}}$  as  $N \rightarrow +\infty$ .
- ▶ **(convergence of the stationary measures)** convergence of  $\pi_{\theta_n}(f) - \pi_*(f)$  as  $n \rightarrow +\infty$ .



Set

$$M_\epsilon(x, \theta) := \inf \{n \geq 1, \|P_\theta^n(x, \cdot) - \pi_\theta\|_{\text{TV}} \leq \epsilon\}.$$

Theorem

Assume

- (i) **D.A. cond**  $\sup_x \|P_{\theta_n}(x, \cdot) - P_{\theta_{n-1}}(x, \cdot)\|_{\text{TV}} \xrightarrow{\mathbb{P}} 0$
- (ii) **C. cond**  $\forall \epsilon > 0, \lim_M \sup_n \mathbb{P}(M_\epsilon(X_n, \theta_n) \geq M) = 0$
- (iii)  $\pi_\theta = \pi_\star$

Then  $\sup_{f, \|f\|_\infty \leq 1} |\mathbb{E}[f(X_n)] - \pi_\star(f)| = 0.$

i.e. conditions on

- ▶ **(Diminishing Adaptation)** the difference  $\|P_{\theta_n}(x, \cdot) - P_{\theta_{n-1}}(x, \cdot)\|_{\text{TV}}$
- ▶ **(ergodicity of  $P_\theta$  / Containment)** the convergence of  $\|P_\theta^N(x, \cdot) - \pi_\theta\|_{\text{TV}}$  as  $N \rightarrow +\infty$ .
- ▶ **(convergence of the stationary measures)** convergence of  $\pi_{\theta_n}(f) - \pi_\star(f)$  as  $n \rightarrow +\infty$ .

Set

$$M_\epsilon(x, \theta) := \inf\{n \geq 1, \|P_\theta^n(x, \cdot) - \pi_\theta\|_{\text{TV}} \leq \epsilon\}.$$

Theorem

Assume

- (i) **D.A. cond**  $\sup_x \|P_{\theta_n}(x, \cdot) - P_{\theta_{n-1}}(x, \cdot)\|_{\text{TV}} \xrightarrow{\mathbb{P}} 0$
- (ii) **C. cond**  $\forall \epsilon > 0, \lim_M \sup_n \mathbb{P}(M_\epsilon(X_n, \theta_n) \geq M) = 0$
- (iii)  $\forall \epsilon > 0, \sup_{f \in \mathcal{F}} \mathbb{P}(|\pi_{\theta_n}(f) - \pi_\star(f)| > \epsilon) \rightarrow 0$

Then  $\sup_{f \in \mathcal{F}} |\mathbb{E}[f(X_n)] - \pi_\star(f)| = 0.$ 

i.e. conditions on

- ▶ **(Diminishing Adaptation)** the difference  $\|P_{\theta_n}(x, \cdot) - P_{\theta_{n-1}}(x, \cdot)\|_{\text{TV}}$
- ▶ **(ergodicity of  $P_\theta$  / Containment)** the convergence of  $\|P_\theta^N(x, \cdot) - \pi_\theta\|_{\text{TV}}$  as  $N \rightarrow +\infty$ .
- ▶ **(convergence of the stationary measures)** convergence of  $\pi_{\theta_n}(f) - \pi_\star(f)$  as  $n \rightarrow +\infty$ .

## 2.2. Convergence of the marginals: in 'practice'

It is sufficient to establish

- ▶ (D.A. cond) problem specific
- ▶ (C. cond) a uniform-in- $\theta$  drift condition (geometric or sub-geometric drift) and a uniform-in- $\theta$  minorization of the transition kernel

(ROBERTS-ROSENTHAL (2007); BAI (2009); ATCHADÉ-FORT (2009))

- ▶ (Cvg  $\pi_{\theta_n}$ )  $\exists \theta_*$  and a set  $A$  s.t.  $\mathbb{P}(A) = 1$  and

$$\forall \omega \in A, \quad \forall x, \forall B \quad P_{\theta_n(\omega)}(x, B) = P_{\theta_*}(x, B).$$

## 3.1. Strong law of large numbers: Suff cond

Let

- ▶ a family of Markov kernels  $\{P_\theta, \theta \in \Theta\}$  s.t.  $P_\theta$  has an unique invariant probability measure  $\pi_\theta$
- ▶ a filtration  $\mathcal{F}_n$  and a process  $\{(X_n, \theta_n), n \geq 0\}$  s.t. for any  $f \geq 0$ ,

$$\mathbb{E}[f(X_{n+1})|\mathcal{F}_n] = \int f(y) P_{\theta_n}(X_n, dy) \quad \mathbb{P}\text{-a.s.}$$

Given a target density  $\pi_*$ , which set of conditions will imply

$$n^{-1} \sum_{k=1}^n f(X_k) \rightarrow \pi_*(f) \quad \mathbb{P} - a.s.$$

for a large class of functions  $f$ ?

Idea :

$$\begin{aligned} n^{-1} \sum_{k=1}^n f(X_k) - \pi_{\star}(f) \\ &= n^{-1} \sum_{k=1}^n \{f(X_k) - \pi_{\theta_{k-1}}(f)\} + n^{-1} \sum_{k=0}^{n-1} \{\pi_{\theta_k}(f) - \pi_{\star}(f)\} \\ &= M_n(f) + R_n(f) + n^{-1} \sum_{k=0}^{n-1} \{\pi_{\theta_k}(f) - \pi_{\star}(f)\} \end{aligned}$$

where  $M_n$  : martingale.

Idea :

$$\begin{aligned}
 n^{-1} \sum_{k=1}^n f(X_k) - \pi_*(f) & \\
 &= n^{-1} \sum_{k=1}^n \{f(X_k) - \pi_{\theta_{k-1}}(f)\} + n^{-1} \sum_{k=0}^{n-1} \{\pi_{\theta_k}(f) - \pi_*(f)\} \\
 &= M_n(f) + R_n(f) + n^{-1} \sum_{k=0}^{n-1} \{\pi_{\theta_k}(f) - \pi_*(f)\}
 \end{aligned}$$

where  $M_n$  : martingale.

i.e. conditions for

- ▶ **a.s. conv for martingales** : from conditions on  $L^p$ -moments for the increments ( $p > 1$ ).
- ▶ **a.s. conv of the residual terms** : from a strengthened diminishing adaptation condition ( $\longleftrightarrow$  conditions on the regularity in  $\theta$  of the Poisson equation)
- ▶ **a.s. conv of the stationary measures** : from the “a.s.” conv of  $P_{\theta_n}(x, B)$  to  $P_{\theta_*}(x, B)$

## 3.2. Strong law of large numbers “in practice”

It is sufficient to establish

- ▶ (strengthened D.A. cond) problem specific
- ▶ (C. cond) a uniform-in- $\theta$  drift condition (geometric or sub-geometric drift) and a uniform-in- $\theta$  minorization of the transition kernel

(ROBERTS-ROSENTHAL (2007); BAI (2009); ATCHADÉ-FORT (2009))

- ▶ (Cvg  $\pi_{\theta_n}$ )  $\exists \theta_*$  and a set  $A$  s.t.  $\mathbb{P}(A) = 1$  and

$$\forall \omega \in A, \quad \forall x, \forall B \quad P_{\theta_n(\omega)}(x, B) = P_{\theta_*}(x, B).$$

## 3.2. Strong law of large numbers “in practice”

It is sufficient to establish

- ▶ (strengthened D.A. cond) problem specific
- ▶ (C. cond) a uniform-in- $\theta$  drift condition (geometric or sub-geometric drift) and a uniform-in- $\theta$  minorization of the transition kernel

(ROBERTS-ROSENTHAL (2007); BAI (2009); ATCHADÉ-FORT (2009))

- ▶ (Cvg  $\pi_{\theta_n}$ )  $\exists \theta_*$  and a set  $A$  s.t.  $\mathbb{P}(A) = 1$  and

$$\forall \omega \in A, \quad \forall x, \forall B \quad P_{\theta_n(\omega)}(x, B) = P_{\theta_*}(x, B).$$

When the drift condition is of the form :

- ▶ (Geom)  $P_\theta V \leq \lambda V + b \mathbb{1}_C$  : strong law of large numbers for functions increasing like  $V^\alpha$  for any  $\alpha \in [0, 1[$ .
- ▶ (Sub-Geom)  $P_\theta V \leq V - c V^{1-\alpha} + b \mathbb{1}_C$  : strong law of large numbers for functions increasing like  $V^\beta$  for any  $\beta \in [0, 1 - \alpha[$ .



## Conclusion

We provide answers to the problem : given

- ▶ a family of Markov kernels  $\{P_\theta, \theta \in \Theta\}$  s.t.  $P_\theta$  has an unique invariant probability distribution  $\pi_\theta$
- ▶ a filtration  $\mathcal{F}_n$  and a process  $\{(X_n, \theta_n), n \geq 0\}$  s.t. for any  $f \geq 0$ ,

$$\mathbb{E}[f(X_{n+1})|\mathcal{F}_n] = \int f(y) P_{\theta_n}(X_n, dy) \quad \mathbb{P}\text{-a.s.}$$

- which set of conditions will imply
  - ▶ convergence of the distribution of  $\{X_n, n \geq 0\}$  to some prob.  $\pi_*$
  - ▶ convergence of the empirical distribution  $n^{-1} \sum_{k=1}^n \delta_{X_k}$
- Appli: convergence of “internal” and “external” adaptive MCMC.
- Details in
  - ▶ **Y. Atchadé, G. Fort** Limit theorems for some adaptive MCMC algorithms with subgeometric kernels, *Accepted in Bernoulli, 2009*
  - ▶ **Y. Atchadé, G. Fort, E. Moulines, P. Priouret** Adaptive MCMC : theory and practice, *submitted*