

When Markov Chains control Monte Carlo sampling

Gersende Fort

CNRS

Institut de Mathématiques de Toulouse, France



Processus Markoviens, semi-markoviens et leurs applications - Montpellier, 2023.

Joint works with

- Yves Atchade, Univ. Michigan, USA
- Benjamin Jourdain, ENPC, France
- Estelle Kuhn, INRAE, France
- Tony Lelièvre, ENPC, France
- Eric Moulines, Ecole Polytechnique, France
- Pierre Priouret, Univ. Paris 6, France
- Amandine Schreck, Telecom ParisTech, France
- Gabriel Stoltz, ENPC, France
- Pierre Vandekerkhove, Univ. Marne-la-Vallée, France

Convergence of adaptive and interacting Markov chain Monte Carlo algorithms. Fort, Moulines, Priouret (2011, Ann Stat).

A Central limit theorem for adaptive and interacting Markov Chains. Fort, Moulines, Priouret, Vandekerkhove (2014, Bernoulli).

Self-Healing Umbrella Sampling: Convergence and Efficiency. Fort, Jourdain, Lelièvre and Stoltz (2017, Stat & Comput)

Adaptive Equi-Energy sampler: Convergence and Illustration. Schreck, Fort, Moulines (2013, TOMACS).

A control, why ?

To improve Monte Carlo methods targetting: $d\pi = \pi d\mu$

- The "naive" MC sampler depends on design parameters in \mathbb{R}^p or in infinite dimension θ
- Theoretical studies characterize an optimal choice of these parameters θ_* by

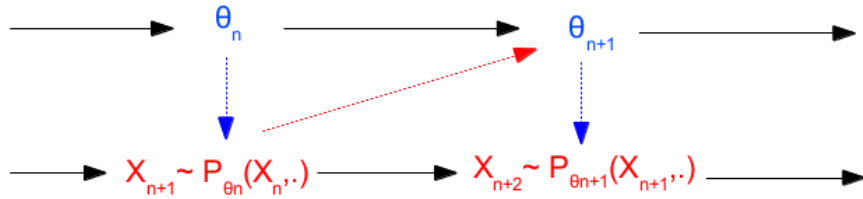
$$\theta_* \in \Theta \text{ s.t. } \int H(\theta, x) d\pi(x) = 0$$

or

$$\theta_* \in \operatorname{argmin}_{\theta \in \Theta} \int C(\theta, x) d\pi(x) = 0.$$

- Strategies:
 - Strategy 1: a preliminary "machinery" for the approximation of θ_* ; **then** run the MC sampler with $\theta \leftarrow \theta_*$
 - Strategy 2: learn θ and sample **concomitantly**

In this talk, Monte Carlo sampling !



$$\mathbb{E} \left[f(X_{n+1}) | \mathcal{F}_n \right] = P_{\theta_n} f(X_n)$$

- from the Monte Carlo point of view:
which conditions on the updating scheme for convergence of the sampler ?
Case: Markov chain Monte Carlo sampler
- from the optimization point of view:
which conditions on the Monte Carlo approximation for convergence of the stochastic optimization ?
Case: Stochastic Approximation methods with Markovian inputs

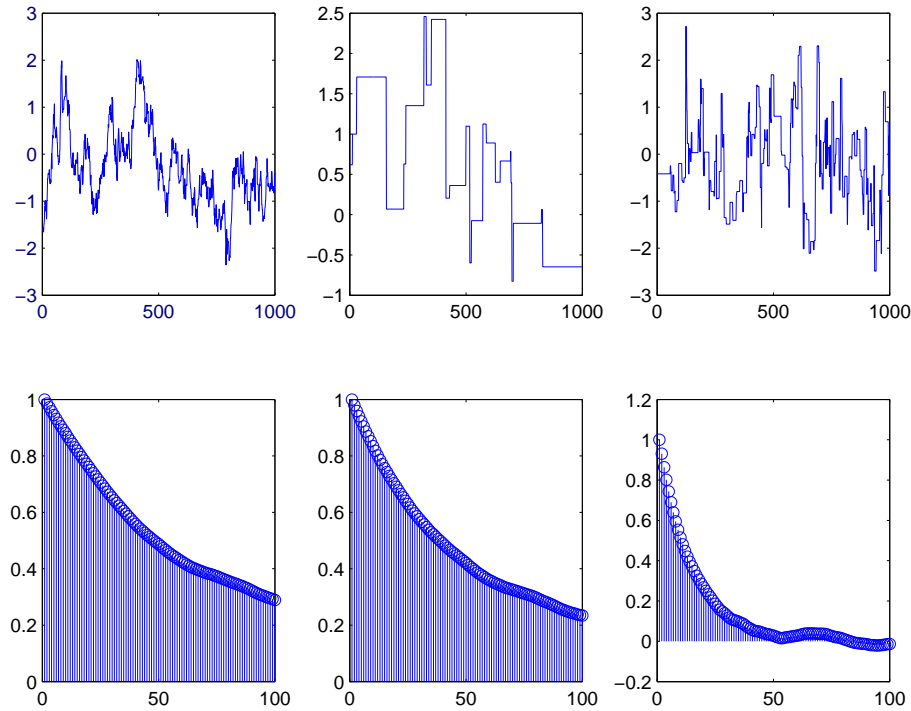
Outline

- Part I. Motivating examples: adaptive and interacting Markov chain Monte Carlo samplers.
- Part II. Ergodicity and limit theorems.

Part I: Motivating examples

1st Ex. Adaptive Hastings-Metropolis (1/3)

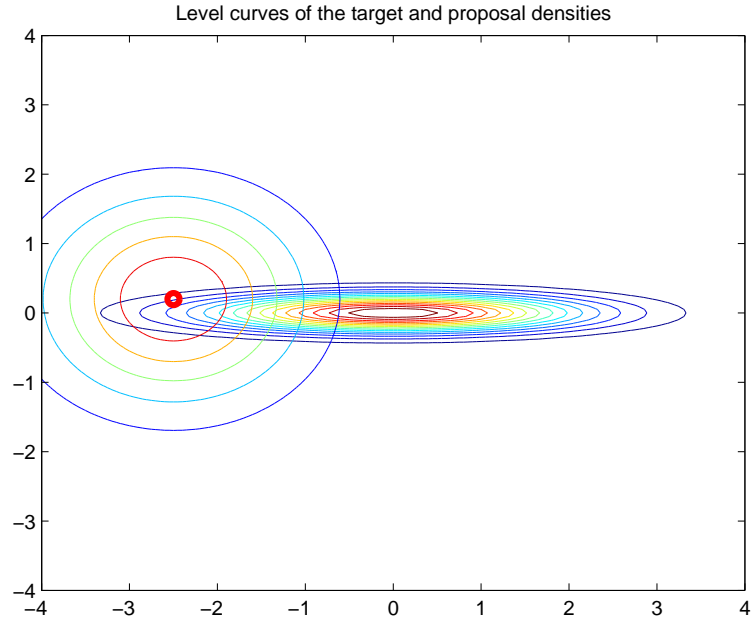
Symmetric Random Walk : proposal $X_{t+1/2} \sim X_t + \mathcal{N}(0, \Gamma)$



(d=1) Different values of Γ : [top] a path of the Markov chain. [bottom] auto-correlation function

1st Ex. Adaptive Hastings-Metropolis (2/3)

Pioneering works: Gelman, Roberts and Gilks (1996)



($d = 2$) The level curves of the target density, and of the proposal distribution $\mathcal{N}(X_t, Id)$.

The optimal rule

$$\Gamma = \frac{2.38^2}{d} \Sigma_{\pi} \quad \text{But } \Sigma_{\pi} \text{ is unknown}$$

1st Ex. Adaptive Hastings-Metropolis (3/3)

Pioneering works: Haario, Saksman and Tamminen (1999).

Iterate

- Sample X_{t+1} by a HM kernel with proposal $\mathcal{N}(X_t, \hat{\Gamma}_t)$: $X_{t+1} \sim P_{\hat{\Gamma}_t}(X_t, \cdot)$.
- Update the unknown parameter

$$\begin{aligned}\hat{\Gamma}_{t+1} &= \frac{1}{t+1} \sum_{\ell=1}^{t+1} (X_\ell - \hat{\mu}_\ell)(X_\ell - \hat{\mu}_\ell)^\top \\ &= \hat{\Gamma}_t + \frac{1}{t+1} \left\{ (X_{t+1} - \hat{\mu}_{t+1})(X_{t+1} - \hat{\mu}_{t+1})^\top - \hat{\Gamma}_t \right\}\end{aligned}$$

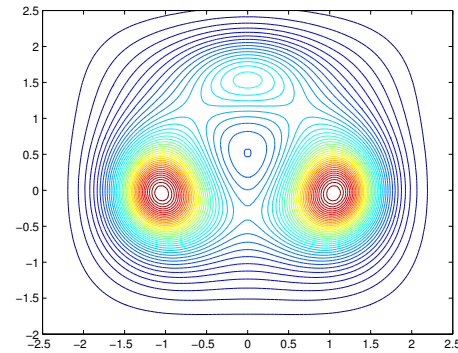
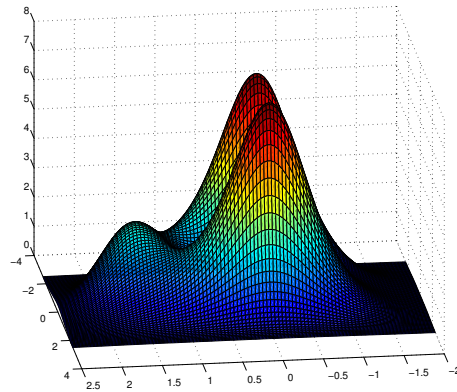
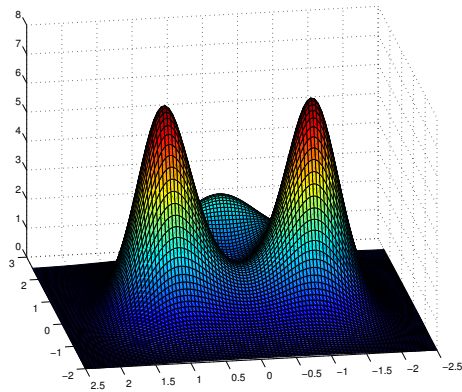
by one step of a *Stochastic Approximation* algorithm designed to solve

$$\mathbb{E}_\pi \left[(X - \mathbb{E}_\pi[X])(X - \mathbb{E}_\pi[X])^\top \right] - \Gamma = 0.$$

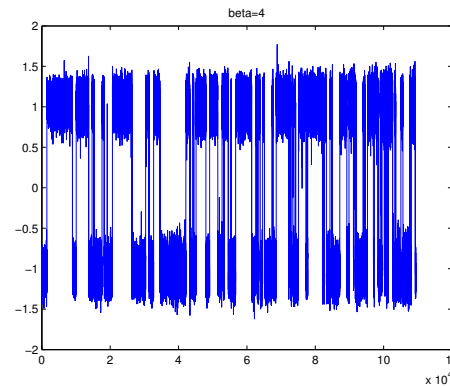
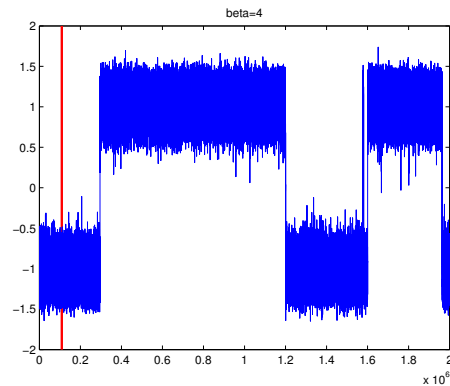
2nd Ex. Adaptive Importance Sampling (1/6)

The problem

- A highly multimodal target density $d\pi$ on $\mathcal{X} \subseteq \mathbb{R}^d$.



- Two samplers with different behaviors (plot: the x -path of a chain in \mathbb{R}^2)



2nd Ex. Adaptive IS by Wang Landau approaches (2/6)

The strategy for choosing the proposal mechanism

- A family of proposal mechanisms obtained by biasing locally the target:
 - given a partition $\mathcal{X}_1, \dots, \mathcal{X}_I$ of \mathcal{X} ,
 - for any weight vector $\theta = (\theta(1), \dots, \theta(I))$

$$d\pi_\theta(x) = \frac{1}{\sum_{i=1}^I \frac{\theta_\star(i)}{\theta(i)}} \sum_{i=1}^I \mathbf{1}_{\mathcal{X}_i}(x) \frac{d\pi(x)}{\theta(i)}, \quad \text{with } \theta_\star(i) := \int_{\mathcal{X}_i} d\pi(u).$$

- Optimal proposal: $d\pi_{\theta_\star}$ <proof>
- Unfortunately, θ_\star unavailable.

2nd Ex. Adaptive IS by Wang Landau approaches (3/6)

If π_{θ_\star} were available

- The algorithm would be:

- Sample X_1, \dots, X_n, \dots i.i.d. with distribution $d\pi_{\theta_\star}$ (or a MCMC with target $d\pi_{\theta_\star}$)

- Compute the importance ratio

$$\frac{d\pi}{d\pi_{\theta_\star}}(X_k) = I \sum_{i=1}^I \mathbf{1}_{\mathcal{X}_i}(X_k) \theta_\star(i)$$

- When approximating an expectation, set

$$\int \phi d\pi \approx \frac{I}{T} \sum_{t=1}^T \left(\sum_{i=1}^I \mathbf{1}_{\mathcal{X}_i}(X_t) \theta_\star(i) \right) \phi(X_t).$$

2nd Ex. Adaptive IS by Wang Landau approaches (4/6)

θ_* and therefore $d\pi_{\theta_*}$ are unknown, so ?

- $\theta_* \in \mathbb{R}^I$ collects $\int \mathcal{X}_i d\pi$ for all $i \in \{1, \dots, I\}$,
- θ_* the unique root of $\theta \mapsto \int_{\mathcal{X}} H(\theta, x) d\pi_{\theta}(x) \in \mathbb{R}^I$ where for all $i \in \{1, \dots, I\}$

$$H_i(\theta, x) := \theta(i) \mathbf{1}_{\mathcal{X}(i)}(x) - \theta(i) \sum_{j=1}^I \mathbf{1}_{\mathcal{X}(j)}(x) \theta(j).$$

thus suggesting the use of a Stochastic Approximation procedure: $\theta_* \approx \lim_t \theta_t$

$$\theta_{t+1} = \theta_t + \gamma_{t+1} H(\theta_t, X_{t+1}) \quad X_{t+1} \sim d\pi_{\theta_t}$$

- This update scheme is a normalized counter of the number of visits to \mathcal{X}_i

2nd Ex. Adaptive IS by Wang Landau approaches (5/6)

The algorithm: Wang-Landau based procedures

- Initialisation: a weight vector θ_0

Repeat for $t = 1, \dots, T$

- sample a point $X_{t+1} \sim d\pi_{\theta_t}$
- update the estimate of θ_*

$$\theta_{t+1} = \theta_t + \gamma_{t+1} H(\theta_t, X_{t+1}) \quad .$$

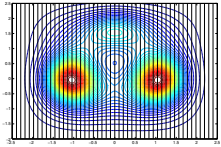
where $X_{t+1} \sim P_{\theta_t}(X_t, \cdot)$ and P_{θ} inv. wrt $d\pi_{\theta}$.

- Expected:

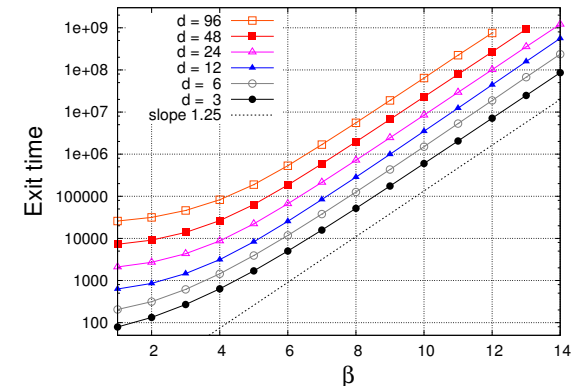
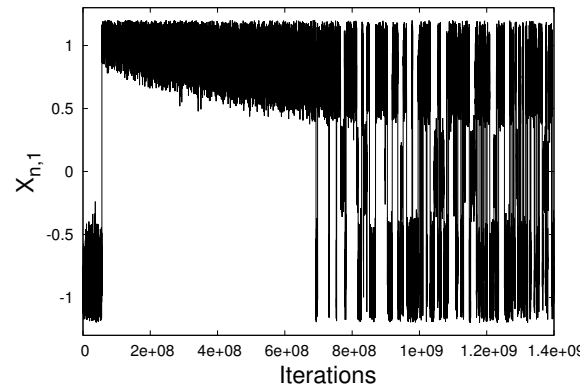
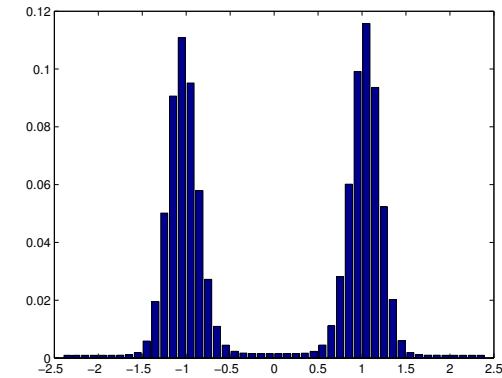
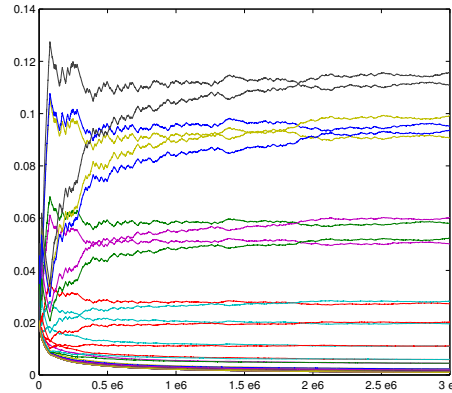
- the convergence of θ_t to θ_* : SA scheme, fed with adaptive (controlled) MCMC sampler,
- the convergence of the distribution of X_t to $d\pi_{\theta_*}$

2nd Ex. Adaptive IS by Wang Landau approaches (6/6)

Does it work ? Plot: convergence of θ_t and first exit times from one mode



► see F, Kuhn, Jourdain, Lelièvre, Stoltz (2014); F, Jourdain, Lelièvre, Stoltz (2015,2017,2018) for studies of these Wang-Landau bases algorithms; including self-tuned SA update rules (γ_t is random).



Conclusion of the 2nd example

- Iterative sampler
- Each iteration combines : (i) a sampling step $X_{t+1} \sim P_{\theta_t}(X_t, \cdot)$; and (ii) a step from an optimization algo. to update the knowledge of some optimal parameter.
- The points $\{X_1, \dots, X_t, \dots\}$ can be seen as the output of a controlled Markov chain

$$\mathbb{E} \left[f(X_{t+1}) | \mathcal{F}_t \right] = P_{\theta_t}(X_t, \cdot) \quad \mathcal{F}_t := \sigma(X_{0:t}, \theta_0)$$

where P_{θ} has $d\pi_{\theta}$ as its unique invariant distribution.

- The convergence of the parameter θ_t is the convergence of a SA scheme with "controlled Markovian" dynamics

$$\theta_{t+1} = \theta_t + \gamma_{t+1} H(\theta_t, X_{t+1})$$

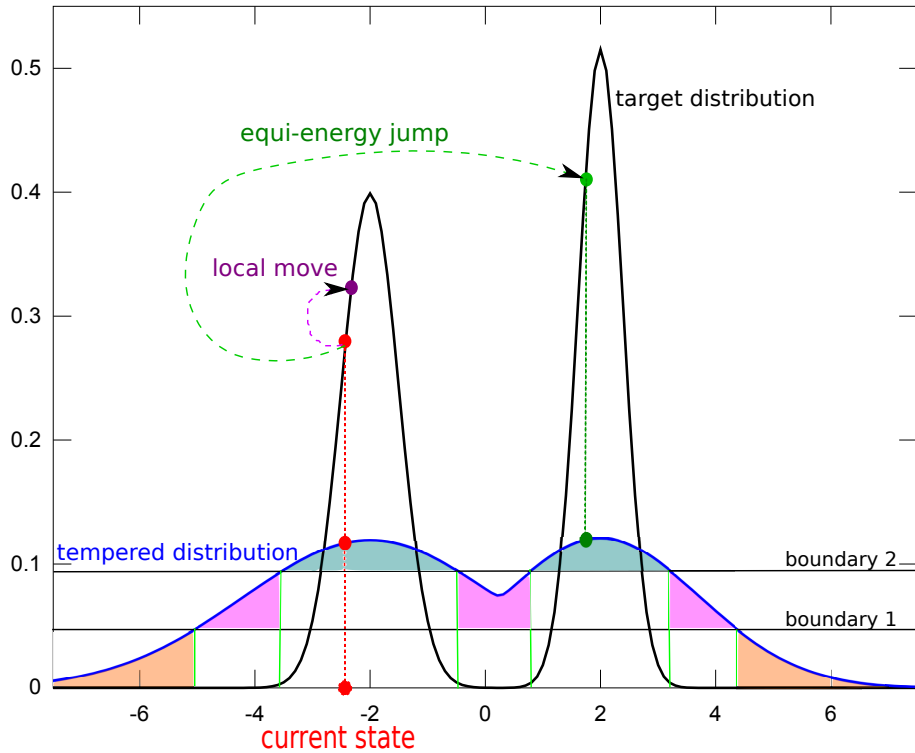
3rd Ex. the Adaptive Equi-Energy sampler (1/4)

extend the EE sampler by Kou-

Zhou-Wong, 2006

- We discussed the case when $\theta \in \mathbb{R}^p$. But there are more general situations: θ may be a distribution case of "interacting" MCMC. (Del Moral-Doucet, 2010; F.-Moulines-Priouret, 2012; Schreck-F.-Moulines, 2013; F.-Moulines-Priouret-Vandekerkhove, 2016)
- Both **interacting** and **tempering** and **adaptive** algorithm.
 - Interacting: run K chains in parallel, s.t. chain $\#k$ is built by using the points of chain $\#(k-1)$. Except the chain $\#1$.
 - Tempering: given $\beta_1 < \dots < \beta_K = 1$, chain $\#k$ is designed to target $d\pi^{\beta_k}$.
 - Adaptive: the mechanism of interaction is learnt on the fly.

3rd Ex. the Adaptive Equi-Energy sampler (2/4)



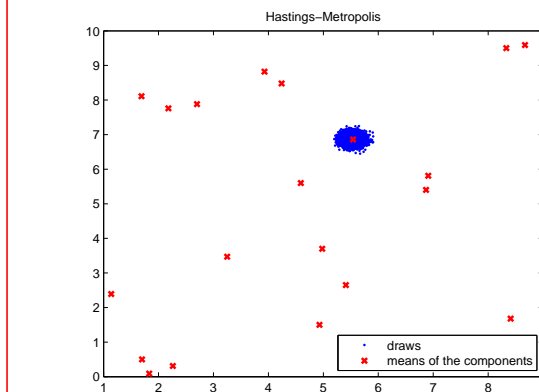
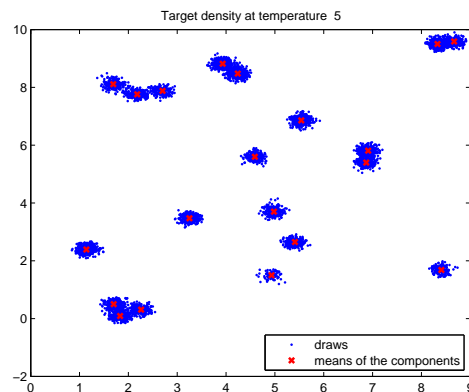
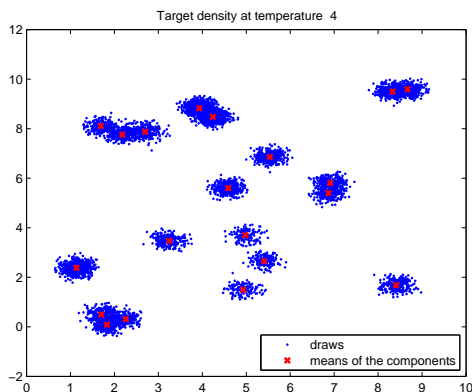
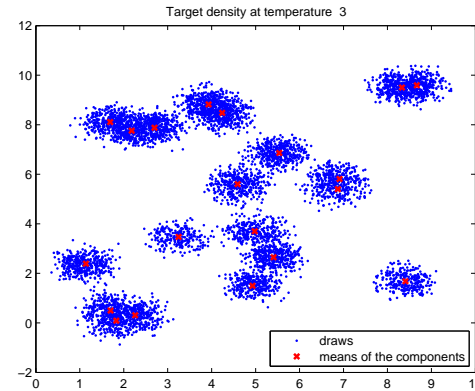
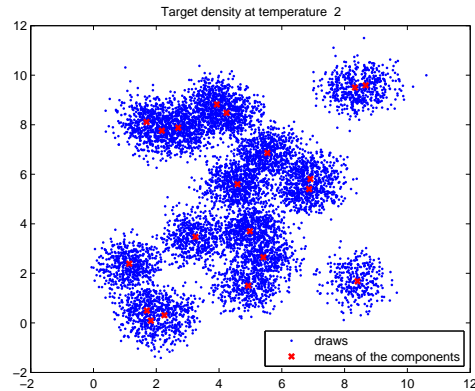
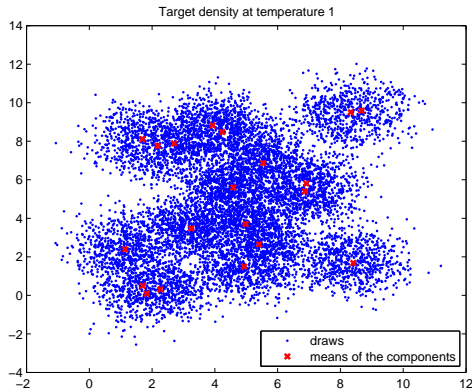
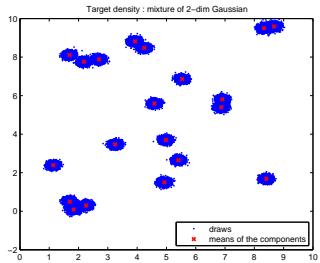
The equi-energy jump: (i) adaptive definition of the equi-energy rings as an estimation of the quantiles of $-\log \pi^{\beta_k}(Z)$ with $Z \sim \pi^{\beta_k}$; (ii) acceptance-rejection ratio;

From chain $X^{(k)}$ to $X^{(k+1)}$:

$$P_{\theta_{k,t}}(X_t^{(k+1)}, \cdot) = (1 - \epsilon) \underbrace{Q(X_t^{(k+1)}, \cdot)}_{\text{MCMC with target } \pi^{\beta_{k+1}}} + \epsilon \underbrace{\tilde{Q}_{\theta_{k,t}}(X_t^{(k+1)}, \cdot)}_{\substack{\text{kernel depending on} \\ \text{the empirical distribution } \theta_{k,t} \\ \text{of the auxiliary process } X_{1:t}^{(k)}}}$$

3rd Ex. the Adaptive Equi-Energy sampler (3/4)

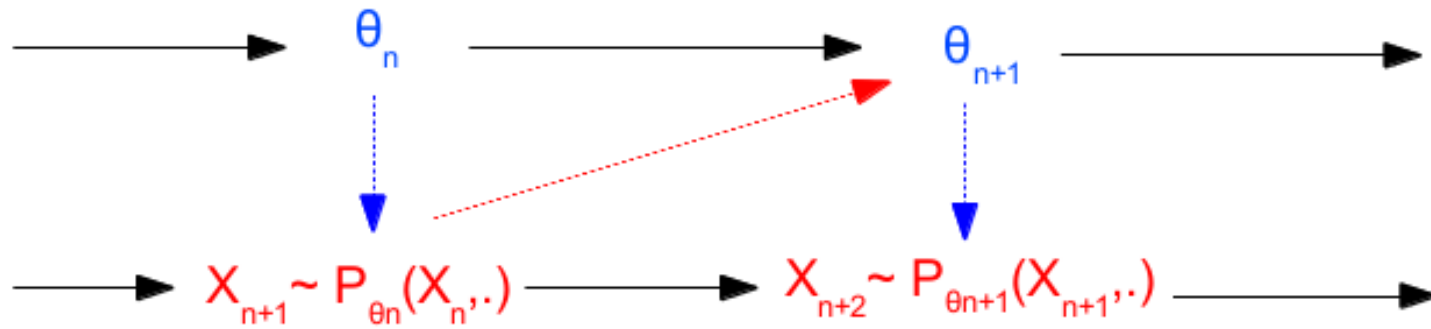
- target density : $\pi = \sum_{i=1}^{20} \mathcal{N}_2(\mu_i, \Sigma_i)$
- 5 parallel processes with target distribution π^{β_k}
($\beta_5 = 1$)



3rd Ex. the Adaptive Equi-Energy sampler (4/4)

- In this example, θ is homogeneous to an empirical distribution (random probability measure).
- For this adaptive sampler Schreck-F.-Moulines (2013): convergence in distribution, law of large numbers.
- For the non-adaptive sampler: convergence analysis in Kou-Zhou-Wong, 2006; Atchadé, 2010; Andrieu-Jasra-Doucet-Del Moral, 2011; F.-Moulines-Priouret, 2012; F.-Moulines-Priouret-Vandekerkhove, 2014
- General results when θ is not necessarily in \mathbb{R}^p : convergence in distribution, law of large numbers, CLT in F.-Moulines-Priouret, 2012; F.-Moulines-Priouret-Vandekerkhove, 2014

Conclusion of this first part (1/3): is a theory required ?



Conclusion of this first part (2/3): Yes !

YES ! convergence can be lost by the adaption mechanism

Even in a simple case when

$$\forall \theta \in \Theta, \quad P_\theta \text{ invariant wrt } d\pi,$$

one can define a simple adaption mechanism

$$X_{t+1} | \text{past}_{1:t} \sim P_{\theta_t}(X_t, \cdot) \quad \theta_t \in \sigma(X_{1:t})$$

such that

$$\lim_t \mathbb{E} [f(X_t)] \neq \int f \, d\pi.$$

Conclusion of the first part (3/3): look !

Fix $t_0, t_1 \in (0, 1)$ s.t. $t_0 + t_1 = 1$.

A $\{0, 1\}$ -valued chain $\{X_t\}_t$ defined by $X_{t+1} \sim P_{X_t}(X_t, \cdot)$ where the transition matrices are

$$P_0 = \begin{bmatrix} t_0 & (1 - t_0) \\ (1 - t_0) & t_0 \end{bmatrix} \quad P_1 = \begin{bmatrix} t_1 & (1 - t_1) \\ (1 - t_1) & t_1 \end{bmatrix}$$

Then

- P_0 and P_1 are invariant w.r.t $[1/2, 1/2]$
- **But** $\{X_t\}$ is a Markov chain invariant w.r.t. $[t_1, t_0]$.

$\{X_t\}$ does not have the same invariant distribution as the one of P_0 and P_1 .

Part II: Convergence of Adaptive/Controlled Markov chains

Convergence results

- The framework:

- a filtration $\{\mathcal{F}_t, t \geq 0\}$ on $(\Omega, \mathcal{A}, \mathbb{P})$
- a \mathcal{F}_t -adapted $\mathcal{X} \times \Theta$ -valued process $\{(X_t, \theta_t), t \geq 0\}$ defined on (Ω, \mathcal{A})
- a family of transition kernels $\{P_\theta, \theta \in \Theta\}$ on a general state space $(\mathcal{X}, \mathcal{X})$
- a conditional distribution satisfying

$$\mathbb{E} \left[f(X_{t+1}) | \mathcal{F}_t \right] = \int P_{\theta_t}(X_t, dx) f(x) \quad f \text{ bounded continuous}$$

and a convergence (in some sense) of the kernels $\{P_{\theta_t}, t \geq 0\}$

BEWARE: the chain $\{X_t\}_t$ is NOT a Markov chain

- Questions:

- convergence in distribution of X_t ?
- limit theorems (SLLN, CLT)

- Hereafter:

- focus on the convergence in distribution; then few words on CLT.
- focus first on $\theta \in \Theta \subset \mathbb{R}^p$; then few words on a more general situation.

Assumptions (1/3) Invariant distribution

$\forall \theta \in \Theta, \exists \pi_\theta$ s.t. the kernel P_θ invariant wrt π_θ

Assumptions (2/3) (Generalized) Containment condition

- Uniform-in- θ ergodicity condition

$$\sup_{\theta \in \Theta} \|P_\theta^r(x; \cdot) - \pi_\theta\|_{\text{TV}} \leq C\rho^r \quad \rho \in (0, 1).$$

In practice: a drift and a minorization condition \rightarrow explicit control of ergodicity

$$P_\theta V \leq \lambda_\theta V + b_\theta, \quad P_\theta(x, \cdot) \geq \delta_\theta \nu_\theta(\cdot) \text{ for } x \in \{V \leq 2b_\theta(1 - \lambda_\theta)^{-1} - 1\}$$

- [Weakened Cond.] for any $\epsilon > 0$, there exists a non-decreasing sequence r_ϵ s.t. $\lim_t r_\epsilon(t)/t = 0$ and

$$\limsup_t \mathbb{E} \left[\|P_{\theta_{t-r_\epsilon(t)}}^{r_\epsilon(t)}(X_{t-r_\epsilon(t)}; \cdot) - \pi_{\theta_{t-r_\epsilon(t)}}\|_{\text{TV}} \right] \leq \epsilon$$

- Controlled rate of growth-in- θ

here, $r_\epsilon(t) = t^\bullet$

$$\|P_\theta^r(x; \cdot) - \pi_\theta\|_{\text{TV}} \leq C_\theta \rho_\theta^r$$

$$t^{-\tau} \|\theta_t\| < \infty \text{ a.s.}$$

$$\limsup_t t^{-\tilde{\tau}} (C_{\theta_t} \vee (1 - \rho_{\theta_t})^{-1}) < \infty \text{ a.s.}$$

Assumptions (3/3) (Generalized) Diminishing adaptation condition

- Uniform-in- θ ergodic condition,

$$\lim_t \mathbb{E} [D(\theta_t, \theta_{t-1})] = 0$$

where $D(\theta, \theta') = \sup_x \|P_\theta(x, \cdot) - P_{\theta'}(x, \cdot)\|_{TV}$.

- [Weakened cond.] For any $\epsilon > 0$,

$$\lim_t \mathbb{E} \left[\sum_{j=1}^{r_\epsilon(t)-1} D(\theta_{t-r_\epsilon(t)+j}, \theta_{t-r_\epsilon(t)}) \right] = 0$$

In practice

- Prove a Lipschitz property $D(\theta, \theta') \leq C \|\theta - \theta'\|$
- Use the definition of θ_t as a function of $(X_\ell)_{\ell \leq t}$ and possibly other "external" sampled points
- Require controls of the form $\mathbb{E}[W(X_\ell)]$, solved e.g. by drift inequalities

$$\mathbb{E}[W(X_\ell) | \mathcal{F}_{\ell-1}] = P_{\theta_{\ell-1}} W(X_{\ell-1}) \leq \lambda_{\theta_{\ell-1}} W(X_{\ell-1}) + b_{\theta_{\ell-1}}$$

Convergence in distribution (1/3)

When $\pi_\theta = \pi$ for any θ

Under these conditions, for any bounded function f ,

$$\lim_t \mathbb{E} [f(X_t)] = \int f(x) \, d\pi(x)$$

Sketch of proof :

$$\mathbb{E} [f(X_t)] - \int f(x) \, d\pi(x) = \mathbb{E} [f(X_t)] - \mathbb{E} [P_{\theta_{t-r}}^r f(X_{t-r})] + \mathbb{E} [P_{\theta_{t-r}}^r f(X_{t-r})] - \int f(x) \, d\pi(x)$$

$$\left| \mathbb{E} [f(X_t)] - \mathbb{E} [P_{\theta_{t-r}}^r f(X_{t-r})] \right| \leq \|f\|_\infty \mathbb{E} \left[\sup_\theta \|P_\theta^r(X_{t-r}, \cdot) - \pi\|_{\text{TV}} \right]$$

$$\left| \mathbb{E} [f(X_t)] - \mathbb{E} [P_{\theta_{t-r}}^r f(X_{t-r})] \right| \leq \sum_{j=1}^{r-1} \mathbb{E} [D(\theta_{t-r+j}, \theta_{t-r})]$$

Convergence in distribution (2/3)

When each kernel P_θ has its own invariant distribution π_θ , with an explicit expression

- Under these three conditions, and
 - there exists a constant α s.t. $\lim_t \int f d\pi_{\theta_t} = \alpha$ a.s.

then

$$\lim_t \mathbb{E} [f(X_t)] = \alpha.$$

- Corollary: if $\{\pi_{\theta_t}\}_t$ converges weakly to π a.s., then $\alpha = \int f d\pi$ for any bounded continuous function f .

Convergence in distribution (3/3)

When π_θ exists but its expression is unknown

It is the most technical case: how to prove the convergence of $\int f \, d\pi_{\theta_t}$ when only properties on the kernels P_{θ_t} are available ?

We write

$$\begin{aligned} \int f \, d\pi_{\theta_t} - \int f \, d\pi_{\theta_\star} &= \left(\int f \, d\pi_{\theta_t} - \int f(y) P_{\theta_t}^k(x, dy) \right) \\ &\quad + \left(\int P_{\theta_t}^k(x, dy) f(y) - \int P_{\theta_\star}^k(x, dy) f(y) \right) + \left(\int P_{\theta_\star}^k(x, dy) f(y) - \int f \, d\pi_{\theta_\star} \right) \end{aligned}$$

and control the blue terms by a condition on the ergodicity of the transition kernels. For the red one,

$$\begin{aligned} P_{\theta_t}^k f(x) - P_{\theta_\star}^k f(x) &= \int \left(P_{\theta_t}(x, dy) - P_{\theta_\star}(x, dy) \right) P_{\theta_\star}^{k-1} f(y) \\ &\quad + \int P_{\theta_t}(x, dy) \left(P_{\theta_t}^{k-1} f(y) - P_{\theta_\star}^{k-1} f(y) \right) \end{aligned}$$

Convergence in distribution (3/3) (to follow)

Starting from :

$$\forall x \in \mathcal{X}, A \in \mathcal{X}, \quad \exists \Omega_{x,A}, \quad \mathbb{P}(\Omega_{x,A}) = 1 \quad \forall \omega \in \Omega_{x,A} \quad \lim_t P_{\theta_t(\omega)}(x, A) = P_{\theta_*}(x, A),$$

the steps are:

$$\forall x \in \mathcal{X}, \quad \exists \Omega_x, \quad \mathbb{P}(\Omega_x) = 1 \quad \forall \omega \in \Omega_x \quad \lim_t P_{\theta_t(\omega)}(x, \cdot) \xrightarrow{w} P_{\theta_*}(x, \cdot)$$

↪ Tool: separable metric space \mathcal{X} (ex. Polish)

$$\exists \Omega', \quad \mathbb{P}(\Omega') = 1 \quad \forall \omega \in \Omega', x \in \mathcal{X} \quad \lim_t P_{\theta_t(\omega)}(x, \cdot) \xrightarrow{w} P_{\theta_*}(x, \cdot),$$

↪ Tool: Polish space \mathcal{X} + equicontinuity of $\{P_\theta f - P_{\theta_*} f, \theta \in \Theta\}$

$$\exists \Omega_*, \quad \mathbb{P}(\Omega_*) = 1 \quad \forall \omega \in \Omega_* \quad \lim_t P_{\theta_t(\omega)}^k(x, \cdot) \xrightarrow{w} P_{\theta_*}^k(x, \cdot),$$

↪ Tool: Feller properties of the kernels $\{P_\theta, \theta \in \Theta\}$.

See F.-Moulines-Priouret, 2012)

In the literature

(Roberts-Rosenthal,2007; Atchadé-F.-Moulines-Priouret, 2011; F.-Moulines-Priouret,2012; F.-Moulines-Priouret-Vandekerkhove, 2012)

- Extensions of the sufficient conditions for "convergence in distribution" to the case
 - when NO uniform-in- θ ergodic behavior of the transition kernels $\{P_\theta\}_\theta$ i.e. neither the state space \mathcal{X} nor the parameter space Θ have to be finite / countable / compact
 - without requiring convergence of the sequence $\{\theta_t\}_t$ as a preliminary step for the proof (when $\pi_\theta = \pi$)
 - without assuming the stability of the sequence $\{\theta_t\}_t$ as a preliminary step for the proof
 - each kernel may have its own invariant distribution, explicitly known or not.
- Based on strenghtened "containment" and "diminishing adaptation" conditions,
 - strong Law of Large Numbers for $\{f(X_t)\}_t$ and $\{f(\theta_t, X_t)\}_t$
 - Central Limit Theorem for $\{f(X_t)\}_t$ (see below)

Strong Law of Large Numbers

Under additional assumptions strengthening the conditions between

- the diminishing adaptation condition

$$D_V(\theta_t, \theta_{t-1}) = \sup_x \frac{\|P_{\theta_t}(x, \cdot) - P_{\theta_{t-1}}(x, \cdot)\|_V}{V(x)}$$

- the rate of convergence of the kernels P_θ to stationarity

- the stability (control of growth in t) of the sequence $\{\theta_t\}_t$

• for any measurable function f such that $\sup_x |f|/V < \infty$

$$\lim_T \frac{1}{T} \sum_{t=1}^T f(X_t) = \lim_t \int f(x) d\pi_{\theta_t}(x) \text{ a.s.}$$

when the RHS exists a.s.

• Extensions: SLLN for $(x, \theta) \mapsto f(x, \theta)$.

Central Limit Theorem (1/2)

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \left(f(X_t) - \int f d\pi_{\theta_*} \right) = \frac{1}{\sqrt{T}} \sum_{t=1}^T \left(f(X_t) - \int f d\pi_{\theta_{t-1}} \right) + \frac{1}{\sqrt{T}} \sum_{t=1}^T \left(\int f d\pi_{\theta_{t-1}} - \int f d\pi_{\theta_*} \right)$$

Under the assumptions

- each kernel P_θ is geometrically ergodic (drift, minorization)
- Trade off: diminishing adaptation, moment conditions, stability of $\{\theta_t\}_t$
- "containment": rate of ergodicity, moment conditions, stability of $\{\theta_t\}_t$

- CLT for the first part, with limiting variance given by

$$\sigma^2(f) = \lim_T \frac{1}{T} \sum_{t=1}^T F(\theta_t, X_t)$$

where $\langle \text{comment on the Poisson equation} \rangle$

$$F(\theta, x) = P_\theta(\Lambda_\theta f)^2 - (P_\theta \Lambda_\theta f)^2, \quad \Lambda_\theta f = (I - P_\theta)^{-1} f$$

Central Limit Theorem (2/2)

For the second part:

- Restricted to algorithms satisfying <comment>

$$\mathbb{E} [f(X_{0:t}) | \theta_{0:t-1}] = \int f(x_{0:t}) \, d\nu(x_0) \prod_{j=1}^t P_{\theta_{j-1}}(x_{j-1}, dx_j)$$

- Upon noting the linearization

$$\pi_{\theta}(f) - \pi_{\theta_{\star}}(f) = \pi_{\theta_{\star}}(P_{\theta} - P_{\theta_{\star}}) \Lambda_{\theta_{\star}} f + \pi_{\theta}(P_{\theta} - P_{\theta_{\star}}) \Lambda_{\theta_{\star}} (P_{\theta} - P_{\theta_{\star}}) \Lambda_{\theta_{\star}} f$$

- Assuming: a CLT with variance $\gamma^2(f)$ for the first part, and a cvg in Prob to zero for the second part

• A global CLT with additive variance

$$\frac{1}{\sqrt{T}} \sum_{t=1}^T \left(f(X_t) - \int f d\pi_{\theta_{\star}} \right) \xrightarrow{d} \mathcal{N} \left(0, \sigma^2(f) + \gamma^2(f) \right)$$

As a conclusion of this part II

- A family of ergodic kernels $\{P_\theta\}_{\theta \in \Theta}$; to adapt the parameters θ_t , a strategy based on the past of the algorithm.
- The easiest situation:
 - uniform-in- θ ergodicity conditions (*i.e. roughly: may be true if the sequence $\{\theta_t\}_t$ remains in a compact set ...*)
- Far more flexible but also more technical:
 - an ergodic behavior depending on θ
 - and the rate of growth of $t \mapsto |\theta_t|$ is controlled
- In both cases,
 - the updating rule $\theta_t \longrightarrow \theta_{t+1}$ is s.t. the adaption is diminishing along iterations.