

# When Monte Carlo and Optimization met in a Markovian dance

Gersende Fort

CNRS

Institut de Mathématiques de Toulouse, France



ICTS "Advances in Applied Probability", Bengaluru, August 2019.

**Intertwined, why ?**

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

- The "naive" MC sampler depends on design parameters  $\theta$ , in  $\mathbb{R}^p$  or in infinite dimension
- Theoretical studies characterize an optimal choice of these parameters  $\theta_\star$  by

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

or

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

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

# To make optimization methods tractable

- Intractable objective function

$$\theta \text{ s.t. } h(\theta) = 0 \quad \text{when } h \text{ is not explicit } h(\theta) = \int_X H(\theta, x) \, d\pi_\theta(x)$$

or

$$\operatorname{argmin}_{\theta \in \Theta} f(\theta) \quad f(\theta) := \int_X C(\theta, x) \, d\pi_\theta(x)$$

- Intractable auxiliary quantities

Ex-1 Gradient-based methods

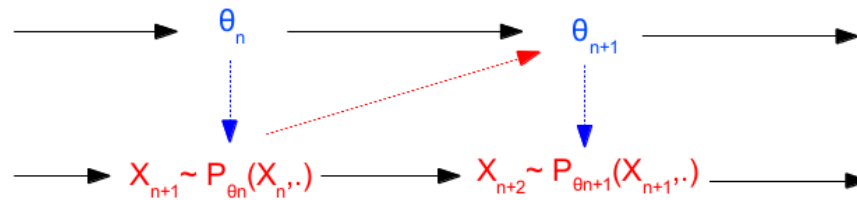
$$\nabla f(\theta) = \int_X H(\theta, x) \, d\pi_\theta(x)$$

Ex-2 Majorize-Minimization methods

$$\text{at iteration } t, \quad f(\theta) \leq F_t(\theta) = \int_X H_t(\theta, x) \, d\pi_{t,\theta}(x)$$

- Strategies: Use Monte Carlo to approximate the unknown quantities.

## In this talk, Markov !



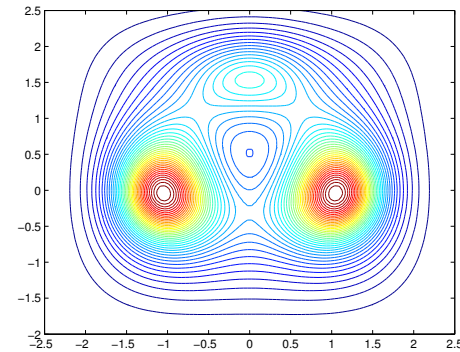
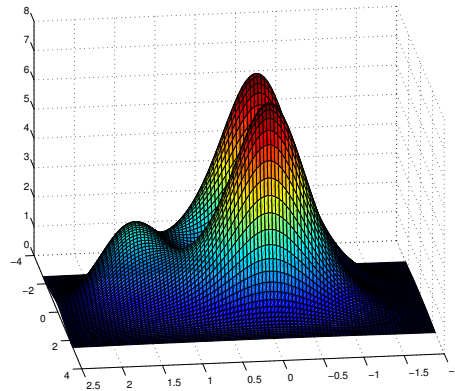
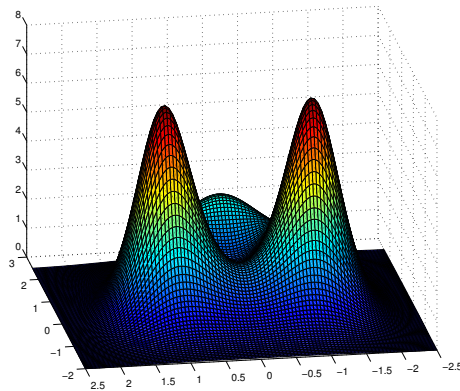
- from the Monte Carlo point of view:  
which conditions on the updating scheme for convergence of the sampler ?  
Case: Markov chain Monte Carlo sampler
- from the optimization point of view:  
which conditions on the Monte Carlo approximation for convergence of the stochastic optimization ?  
Case: Stochastic Approximation methods with Markovian inputs
- (Talk) Application to a Computational Machine Learning problem: penalized Maximum Likelihood through Stochastic Proximal-Gradient based methods

# **Part I: Motivating examples**

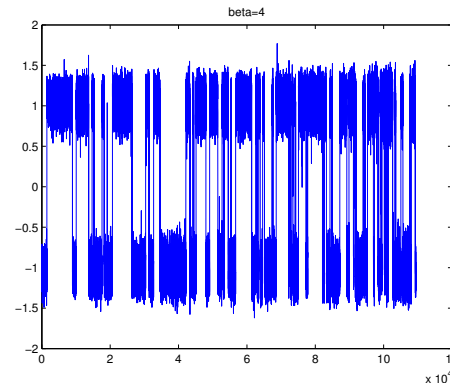
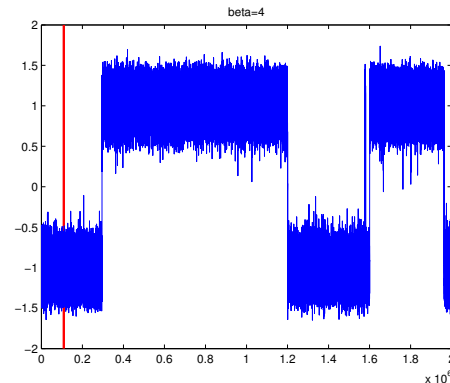
# 1st Ex. Adaptive Importance sampling by Wang-Landau approaches (1/6)

## The problem

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



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



## 1st Ex. (2/6)

### The strategy for choosing the proposal mechanism

- A family of proposal mechanisms obtained by biasing locally the target:
  - given a partition  $X_1, \dots, X_I$  of  $X$ ,
  - for any weight vector  $\theta = (\theta(1), \dots, \theta(I))$

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

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



## 1st Ex. (3/6)

If  $\pi_{\theta_\star}$  were available

- The algorithm would be:
  - Sample  $X_1, \dots, X_t, \dots$  i.i.d. with distribution  $d\pi_{\theta_\star}$  (or a MCMC with target  $d\pi_{\theta_\star}$ )
  - Compute the importance ratio

$$\frac{d\pi}{d\pi_{\theta_\star}}(X_t) = I \sum_{i=1}^I 1_{X_i}(X_t) \theta_\star(i)$$

- When approximating an expectation, set

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

## 1st Ex. (4/6)

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

- $\theta_\star \in \mathbb{R}^I$  collects  $\int_{X_i} d\pi$  for all  $i \in \{1, \dots, I\}$ ,
- $\theta_\star$  the unique root of  $\theta \mapsto \int_X H(\theta, x) d\pi_\theta(x) \in \mathbb{R}^I$  where for all  $i \in \{1, \dots, I\}$

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

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

$$\theta_{t+1} = \theta_t + \gamma_{t+1} H(\theta_t, X_{t+1}) \quad X_{t+1} \sim d\pi_{\theta_t} \quad \text{or } X_{t+1} \sim \text{one-step MCMC}$$

- This update scheme is a normalized counter of the number of visits to  $X_i$   
<proof>

## 1st Ex. (5/6)

### The algorithm: Wang-Landau based procedures

- Initialisation: a weight vector  $\theta_0$

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

- sample a point  $X_{t+1} \sim P_{\theta_t}(X_t, \cdot)$  and  $P_{\theta}$  inv. wrt  $d\pi_{\theta}$ .
- update the estimate of  $\theta_{\star}$

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

- Expected:

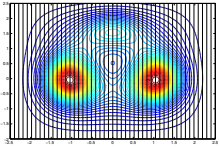
- the convergence of  $\theta_t$  to  $\theta_{\star}$ : SA scheme, fed with adaptive (controlled) MCMC sampler,
- the convergence of the distribution of  $X_t$  to  $d\pi_{\theta_{\star}}$

thus allowing

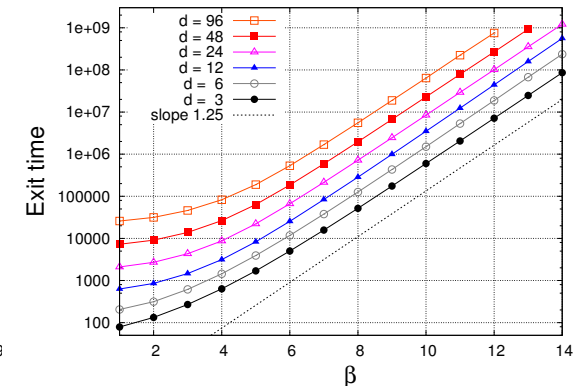
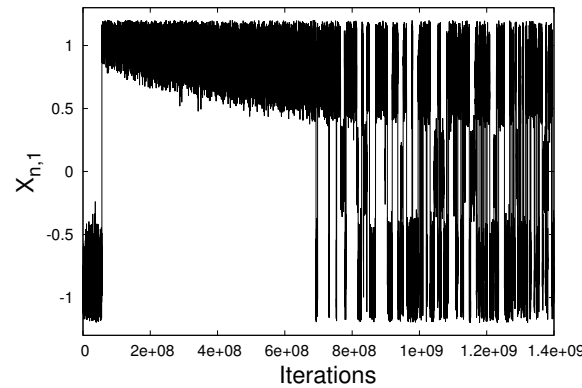
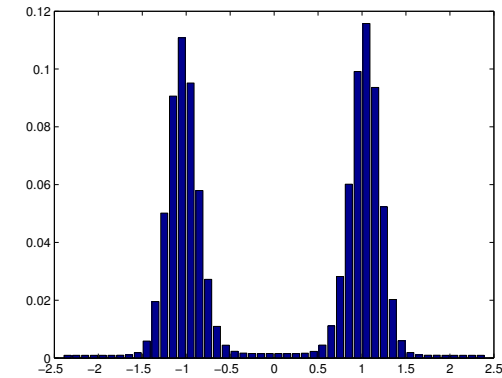
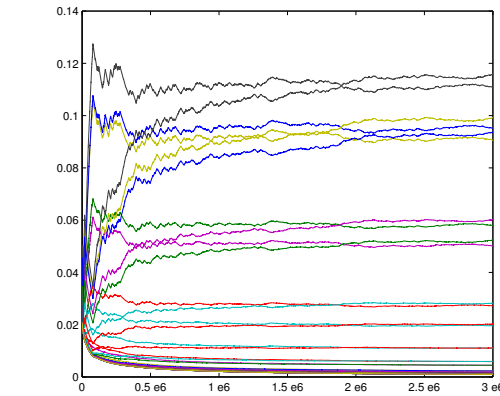
$$\int \phi d\pi \approx \frac{I}{T} \sum_{t=1}^T \left( \sum_{i=1}^I 1_{X_i}(X_t) \theta_t(i) \right) \phi(X_t).$$

## 1st Ex. (6/6)

Does it work ? Plot: convergence of  $\theta_t$  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 ( $\gamma_t$  is random).



## Conclusion of the 1st example

- Iterative sampler
- Each iteration combines : (i) a sampling step  $X_{t+1} \sim P_{\theta_t}(X_t, \cdot)$ ; and (ii) an optimization step 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  $P_{\theta}$  has  $d\pi_{\theta}$  as its unique invariant distribution.

- The convergence of the parameter  $\theta_t$  is the convergence of a SA scheme with "controlled Markovian" dynamics

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

## 2nd Example: penalized ML in latent variable models (1/6)

- An example from Pharmacokinetic:
  - $N$  patients.
  - At time 0: dose  $D$  of a drug.
  - For patient  $\#i$ , observations  $Y_{i1}, \dots, Y_{iJ_i}$  giving the evolution of the concentration at times  $t_{i1}, \dots, t_{iJ_i}$ .

- The model:

$$Y_{ij} = \mathcal{F}(t_{ij}, X_i) + \epsilon_{ij} \quad \epsilon_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

where  $X_i \in \mathbb{R}^L$  is modeled as

$$X_i = Z_i \beta + d_i \in \mathbb{R}^L \quad d_i \stackrel{i.i.d.}{\sim} \mathcal{N}_L(0, \Omega) \text{ and independent of } \epsilon_{\bullet}$$

and  $Z_i$  known matrix s.t. each row of  $X_i$  has in intercept (fixed effect) and covariates.

- Statistical analysis: (i) estimation of  $\theta = (\beta, \sigma^2, \Omega)$ , under sparsity constraints on  $\beta$ ; (ii) selection of the covariates based on  $\hat{\beta}$ .

## 2nd Ex. (2/6)

### Penalized Maximum Likelihood

- The likelihood of  $Y := \{Y_{ij}, 1 \leq i \leq N, 1 \leq j \leq J_i\}$  is not explicit:
  - The distribution of  $Y_{i,j}$  given  $X_i$  is simple; the distribution of  $X_i$  is simple.
  - The joint distribution has an explicit expression - It is an example of latent variable model:

$$\log L(Y; \theta) = \log \int p(Y, x_{1:N}; \theta) \, d\nu(x_{1:N})$$

- Sparsity constraints on the parameter  $\theta$ : through a penalty term  $g(\theta)$
- The penalized ML is of the form

$$\operatorname{argmin}_{\Theta} (-\log L(Y; \theta) + g(\theta))$$

with an intractable objective function.

## 2nd Ex. (3/6)

What about first-order methods for solving the optimization ?

- On the likelihood term:

- Usually regular enough so that the Gradient exists and <proof>

$$\begin{aligned}\nabla_{\theta} \log L(Y; \theta) &= \int \frac{\partial_{\theta} p(Y, x; \theta)}{p(Y, x; \theta)} \frac{p(Y, x; \theta) d\mu(x)}{\int p(Y, z; \theta) d\mu(z)} \\ &= \int \partial_{\theta} (\log p(Y, x; \theta)) \underbrace{d\pi_{\theta}(x)}_{\substack{\text{the a posteriori distribution of } x \text{ given } Y \\ \text{the dep upon } Y \text{ is omitted}}}\end{aligned}$$

- the a posteriori distribution is known up to a normalizing constant.

- On the penalty term

- May be non smooth, but: convex and lower semi-continuous

- Hence a Proximal operator (implicit gradient) is associated - <See the talk>.



## 2nd Ex. (4/6)

### What about EM-like methods for solving the optimization ?

- Expectation-Maximization Dempster-Laird-Rubin, 1977 introduced to solve

$$\operatorname{argmin}_{\theta \in \Theta} \left( -\log \int_{\mathcal{X}} p(x; \theta) d\mu(x) + g(\theta) \right)$$

where the first part is intractable; by iterating two steps

- Expectation step

$$Q(\theta, \theta_t) := \int \log p(x; \theta) \frac{p(x; \theta_t) d\mu(x)}{\int p(z; \theta_t) d\mu(z)} = \int \log p(x; \theta) d\pi_{\theta_t}(x)$$

- Minimization step

$$\theta_{t+1} := \operatorname{argmin}_{\theta} (-Q(\theta, \theta_t) + g(\theta)).$$

- $\theta \mapsto Q(\theta, \theta_t)$  is an integral which is intractable;  $d\pi_{\theta}$  is known up to a normalizing constant.

## 2nd Ex. (5/6)

- Both in EM-like approaches and in gradient-based approaches,
  - faced with intractable auxiliary quantities of the form

$$\int_{\mathcal{X}} H(\theta, x) \, d\pi_{\theta_t}(x) \tag{1}$$

at iteration  $t$  of the optimization algorithm.

- intractable integral;  $d\pi_{\theta}$  is often known up to a normalizing constant.

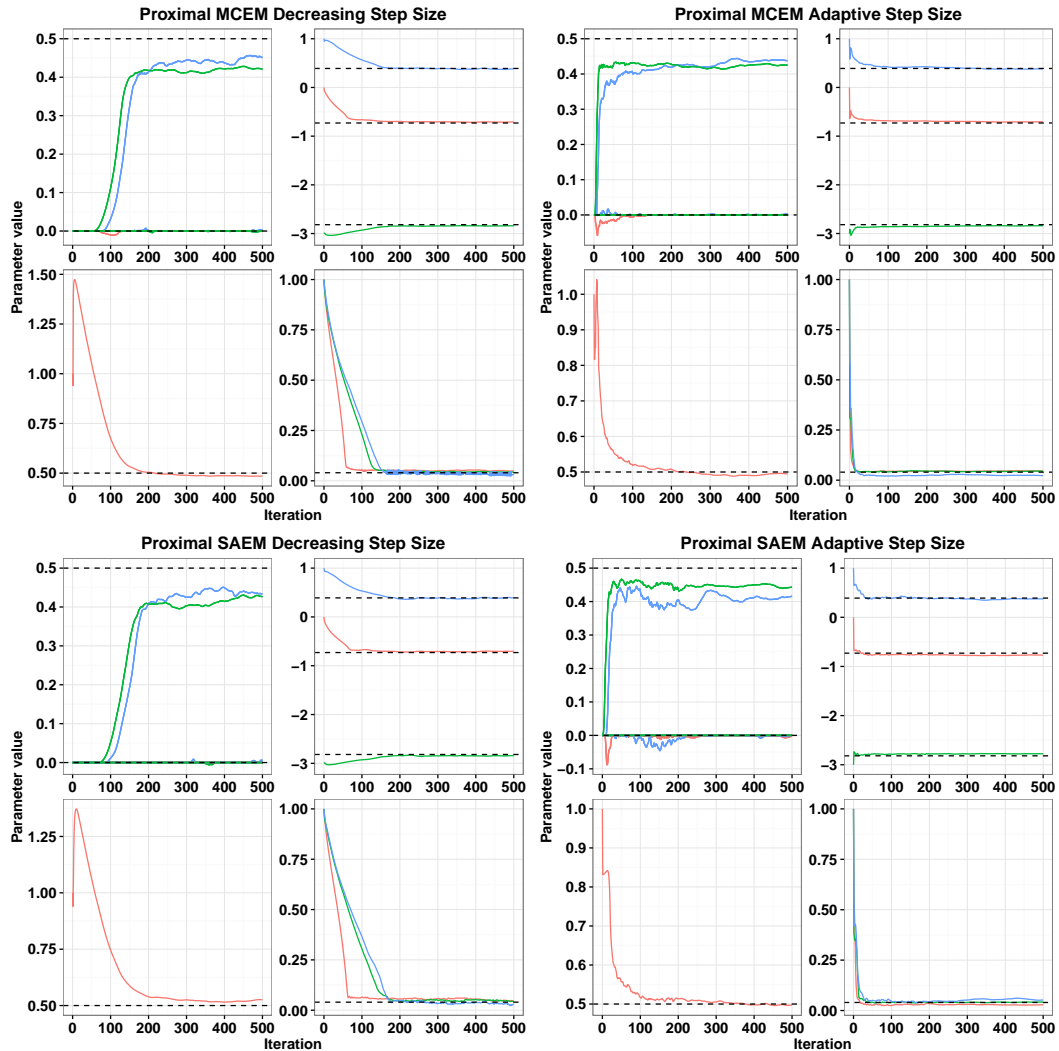
- What kind of approximation of the integral (1) at iteration  $t$  ?

- Quadrature techniques: poor behavior w.r.t. the dimension of  $\mathcal{X}$
- I.i.d. samples from  $\pi_{\theta_t}$  to define a Monte Carlo approximation: not possible, in general.
- use  $m$  samples from a MCMC sampler  $\{X_{j,t+1}, j \leq m\}$  with unique inv. dist.  $d\pi_{\theta_t}$ .

## 2nd Ex. (6/6)

Does it work ?

see F-Moulines (2003)  
for EM-like approaches;  
see Atchadé-F.-Moulines  
(2017) and  
F.-Ollier-Samson (2018)  
for gradient-based  
approaches;  
see F.-Ollier-Samson  
(2018) for the parallel  
between EM-like  
and Gradient-based  
techniques



## Conclusion of the 2nd example

- Iterative optimization technique
- Each iteration combines : (i) an update of the parameter; (ii) a sampling step  $X_{j+1,t+1} \sim P_{\theta_t}(X_{j,t+1}, \cdot)$  to approximate auxiliary quantities.
- The convergence of  $\{\theta_t\}_t$  is the convergence of a stochastically perturbed iterative optimization algorithm. At each iteration: an exact quantity  $\int H(\theta, x) \, d\pi_{\theta_t}(x)$  is approximated by a Monte Carlo sum

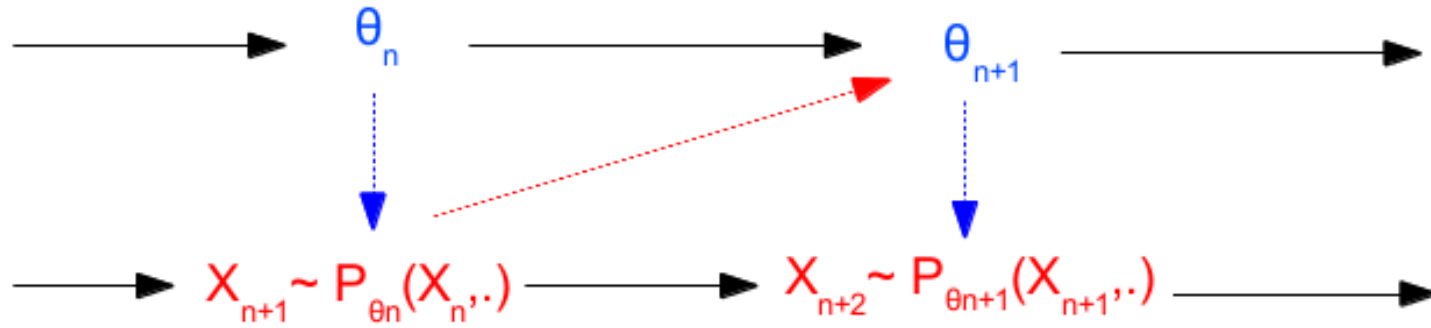
$$\int H(\theta, x) \, d\pi_{\theta_t}(x) \approx \frac{1}{m_{t+1}} \sum_{j=1}^{m_{t+1}} H(\theta, X_{j,t+1})$$

- The points  $\{X_{j,t+1}\}_j$  satisfy

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

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

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



**Conclusion of this first part (2/3): is a theory required when sampling ?**

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

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

---

<proof> 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.  $[t_1, t_0]$

## Conclusion of this first part (3/3): is a theory required when optimizing ?

YES ! Unfortunately ,

- a biased approximation <proof>

$$\mathbb{E} \left[ \frac{1}{m_{t+1}} \sum_{j=1}^{m_{t+1}} H(\theta, X_{j,t+1}) \middle| \mathcal{F}_t \right] = ? \neq \int_{\mathcal{X}} H(\theta, x) \mathrm{d}\pi_{\theta_t}(x)$$

- For a reduced computational cost: a bias which we would like NOT vanishing i.e.  $m_t = m(= 1)$ .

Ex. Stochastic Approximation with controlled Markovian dynamics

$$\begin{aligned} \theta_{t+1} &= \theta_t + \gamma_{t+1} H(\theta_t, X_{t+1}) & X_{t+1} &\sim P_{\theta_t}(X_t, \cdot) \\ &= \theta_t + \gamma_{t+1} \underbrace{\int H(\theta_t, x) \mathrm{d}\pi_{\theta_t}(x)}_{h(\theta_t)} + \gamma_{t+1} \underbrace{\left( H(\theta_t, X_{t+1}) - h(\theta_t) \right)}_{\text{non centered}} \end{aligned}$$