# Exploration Vs. Exploitation in Adaptive Monte Carlo Sampling

Scott C. Schmidler

Department of Statistical Science Duke University

AdapSkIII Conference January 4, 2010

- Mixing times and finite time convergence of adaptive MCMC.
- Combining adaptive strategies.
- Some cautionary notes about MIS kernels.

Available from "www.stat.duke.edu/scs":

- SS & Woodard (2010). Lower Bounds on the Convergence Rates of Adaptive MCMC Methods. (Submitted, under revision)
- Wang & SS (2010). Exploration vs Exploitation: Hybrid Strategies for Adaptive MCMC. (in preparation) see poster
- Ji & SS (2010). Adaptive Markov Chain Monte Carlo for Bayesian Variable Selection. *J. Comp. Graph. Stat.*, (to appear).
- SS & Wiehe (2010). Reservoir Exchange and Adaptive Monte Carlo. (submitted)

• Many ways of "adapting" an MCMC algorithm based on sample path exist; many can be shown to satisfy LLNs.

4 3 b

- Many ways of "adapting" an MCMC algorithm based on sample path exist; many can be shown to satisfy LLNs.
- Purpose of adaptation is to improve *rate* of convergence.

- Many ways of "adapting" an MCMC algorithm based on sample path exist; many can be shown to satisfy LLNs.
- Purpose of adaptation is to improve *rate* of convergence.
- Convergence of MC estimators involves *both* bias *and* variance.

- Many ways of "adapting" an MCMC algorithm based on sample path exist; many can be shown to satisfy LLNs.
- Purpose of adaptation is to improve *rate* of convergence.
- Convergence of MC estimators involves *both* bias *and* variance.
- Different adaptation strategies can be understood as improving one or the other.

- Many ways of "adapting" an MCMC algorithm based on sample path exist; many can be shown to satisfy LLNs.
- Purpose of adaptation is to improve *rate* of convergence.
- Convergence of MC estimators involves *both* bias *and* variance.
- Different adaptation strategies can be understood as improving one or the other.
- Can obtain improved algorithms by combining strategies of different types.

伺 ト イ ヨ ト イ ヨ ト

# High dimensional, multimodal target distributions

- Molecular simulation
- Bayesian variable selection/model selection
- Mixture models
- Non-linear physics-based models



Two approaches developed by various authors

Adaptive random-walk proposals

$$q_{n+1}(x,\cdot) = (1-\alpha)N(x,\hat{\Sigma}_n) + \alpha N(x,\Sigma_0)$$

e.g. Haario et al, Roberts & Rosenthal

Adaptive independence proposals (AMIS)

$$q_{n+1}(x,\cdot) = g(\cdot;\hat{\theta}_n) \quad \hat{\theta}_n = \theta(X_1,\ldots,X_n)$$

e.g. Andrieu & Moulines, Ji & Schmidler, etc.

# Adaptive Metropolized independence sampler (AMIS) [Ji and Schmidler, 2009]

Finite mixture proposal distribution:

$$q(x) = \lambda N(x; \tilde{\mu}, \tilde{\Sigma}) + (1 - \lambda) \sum_{m=1}^{M} w_m N(x; \mu_m, \Sigma_m)$$

(see also Andrieu & Moulines 2005, others)

Point-mass mixture proposal for variable selection:

$$q(x) = (1 - \lambda) \Big[ w_0 \delta(x) + \sum_{m=1}^{M} w_m N(\mu_m, \Sigma_m) \Big] + \lambda N(x; \tilde{\mu}, \tilde{\Sigma})$$

Adapt parameters  $\psi = \{w_m, \mu_m, \Sigma_m\}_{m=1}^M$  to approximate  $\pi(x)$ .

Adaptive strategy: Minimize  $\mathcal{D}[\pi(x) \parallel q(x; \psi)] = \mathbb{E}_{\pi}\left[\log \frac{\pi(x)}{q(x;\psi)}\right]$ 

 $\psi^*$  obtained as a root of derivative:

$$h(\psi) = -\int \frac{\pi(x)}{q(x;\psi)} \frac{\partial}{\partial \psi} q(x;\psi) = 0$$

Approximate  $h(\psi)$  by Monte Carlo integration:

$$h(\psi) \approx \frac{1}{K} \sum_{k=1}^{K} f(X^{(k)}, \psi) \quad \text{for} \quad f(x, \psi) = \frac{\partial}{\partial \psi} [\log \frac{\pi(x)}{q(x; \psi)}]$$

where  $X^{(k)} \sim \pi(x)$ .

 $\hat{h}(X^{(1:K)};\psi)$ : estimate of  $h(\psi)$  based on sample path  $X^{(1:k)}$ 

Stochastic Approximation algorithm [Robbins and Monro, 1951].

$$\psi_{n+1} = \psi_n + r_{n+1}(h(\psi_n) + \xi_{n+1})$$
  
=  $\psi_n + r_{n+1} \hat{h}(X_n^{(1:K)}; \psi_n)$ 

 $\{r_n\}$  decreasing step-sizes satisfying  $\sum_n r_n = \infty$  and  $\sum_n r_n^2 < \infty$ 

Resulting chain is non-Markovian, but can be shown to satisfy a WLLN using results of [Roberts and Rosenthal, 2007]

Bayesian logistic regression model,

$$y_i \mid x_i, \beta \sim \mathsf{Bernoulli}\left(g^{-1}(x_i\beta)\right) \qquad \beta \sim \pi_0(\beta)$$

 $y_i \in \{0,1\}; g(u)$  logistic link

Simulated data set:

- 200 observations
- r = 10 covariates

•  $\beta_{1:10} = [-.01, -1.5, .15, .5, -.15, -.2, -.6, .25, 1.5, -.05]$ 

A 3 A

## Bayesian logistic regression



Figure: Autocorrelation of  $\beta_{1:10}$  under data-augmentation Gibbs sampler [Holmes and Held, 2006] (blue), and adaptive MCMC algorithm (red).

Constructs *I* processes  $X^{(i)}$  with tempered target densities  $\pi^{(i)} \propto \pi^{\beta_i}$  for inverse temperatures  $1 = \beta_1 > \ldots > \beta_l \ge 0$ . (Also truncation  $\pi^{(i)} \propto \pi^{\beta_i} \wedge c_i$ )

For each *i*, bin sample history  $(X_{0:n}^{(i)})$  according to energy.

Process  $X^{(i)}$  occasionally proposes to move to a state previously visited by  $X^{(i+1)}$  lying in same energy bin.

These "equi-energy" moves can be non-local in the state space, potentially enabling transitions between distinct modes of  $\pi$ .

(1日) (日) (日) 日

### MRAM Processes

Let  $X^{(1)}, \ldots, X^{(l)}$  discrete time stochastic processes on  $\mathcal{X}$ . So  $X^{(i)} = X_0^{(i)}, X_1^{(i)}, \ldots$ 

Generated by time-inhomogeneous sequences of transition kernels:

$$K_{i,n} = \alpha T_i + (1 - \alpha)R_{i,n}$$

with  $\alpha \in [0, 1]$ ,  $T_i$  an ergodic time-homogeneous Markov  $\pi^{(i)}$ -reversible transition kernel, and  $R_{i,n}$  is a *resampling* kernel with proposal:

$$Q_{i,n}(X_{n-1}^{(i)}, y) = \sum_{i'=1}^{l} \sum_{j=0}^{n-1} w_{i'j} \delta(y - X_j^{(i')})$$

(Proposes new state from the set of previous samples  $X_{0:n-1}^{(1:I)}$ .)

ヨッ イヨッ イヨッ

Mulitchain resampling adaptive Metropolis (MRAM):

- Equi-Energy Sampler
- Importance-Resampling from the Past (Atchadé)
- Gelfand-Sahu

Let  $\{T_{\theta}\}_{\theta \in \Theta^{(i)}}$  be a set of ergodic,  $\pi^{(i)}$ -reversible Markov kernels.  $T_{\theta_{i,n}}$  time-inhomogeneous *but*  $\pi^{(i)}$ -*invariant* transition kernels  $\theta_{i,n} = g_i(X_{0:n-1}^{(1:l)})$ 

Examples: Haario algorithm and similar variants; multi-chain algorithm of Rosenthal *et al*.

 $X_1, \ldots, X_n$  no longer a Markov chain.

Under what conditions does  $\hat{f}_n = \frac{1}{n} \sum_{i=1}^n f(X_i)$  converge?

- Haario et al 2001: WLLN, using "mixingales"
- Andrieu & Robert (2001): SA interpretation of Haario algorithm
- Andrieu & Moulines (2005), Atchade& Rosenthal (2005): generalizations to other algorithms (and a CLT)
- Roberts & Rosenthal (2007): Simplified conditions, coupling

伺い イラト イラト

Nearly all theory to date deals with *ergodicity* (LLN). A few give conditions for CLTs (e.g Andrieu & Moulines (2005)).

This was needed and a major breakthrough IMO. But all *asymptotic* theory.

Adaptation is only interesting if it improves rates!

### Statistical Efficiency: $var(\hat{f})$

Under reasonably weak conditions<sup>\*</sup>, for any function f with  $var_{\pi}(f) \leq \infty$ , we obtain a CLT:

$$\begin{split} \sqrt{n}(\bar{f}_n - \mu_f) &\to N(0, \sigma_{\bar{f}_n}^2) \end{split}$$
for  $\sigma_{\bar{f}_n}^2 = \sigma_f^2(1 + 2\sum_{j=1}^n (1 - \frac{j}{n})\rho_j) \text{ and} \cr \rho_j &= \frac{1}{\sigma_f^2} E\left((f(X^{(n)}) - \mu_f)(f(X^{(n+j)}) - \mu_f)\right) \end{split}$ 

lag-j autocorrelation.

Finite sample efficiency:

Convergence as well as autocorrelation

$$\mathsf{MSE}(\hat{ heta}) = \mathsf{Bias}^2(\hat{ heta}) + \mathsf{Var}(\hat{ heta})$$

*Finite sample* efficiency: Convergence as well as autocorrelation

$$\mathsf{MSE}(\hat{ heta}) = \mathsf{Bias}^2(\hat{ heta}) + \mathsf{Var}(\hat{ heta})$$

For multimodal targets, bias can dominate in MCMC.

*Finite sample* efficiency: Convergence as well as autocorrelation

$$\mathsf{MSE}(\hat{ heta}) = \mathsf{Bias}^2(\hat{ heta}) + \mathsf{Var}(\hat{ heta})$$

For multimodal targets, bias can dominate in MCMC. For good adaptive MCMC algorithms, bias *will* dominate.

- Ergodicity: SLLN under usual conditions ( $\phi$ -irred, aper,  $\pi$ -invariant)
- <u>Geometric</u>:  $\exists \lambda \in [0,1)$  and  $M(x) < \infty$   $(\pi a.e. x \in \mathcal{X})$  s.t.

$$\|\mu K^n - \pi\| \le M(x)\lambda^n$$

Requires minorization, drift conditions. Implies CLT.

- <u>Uniform</u>:  $M(x) \equiv M$
- Rapid mixing:  $\lambda$  grows at most polynomially in d(Note G.E. requires only  $\lambda^* > 0$ ; e.g.holds for any  $|\mathcal{X}| < \infty$ )
- Quantitative: e.g. Rosenthal 1995

・吊 ・ ・ ラ ・ ・ ラ ト ・ ラ

Let  $(\mathcal{X}^{(d)}, \mathcal{F}^{(d)}, \lambda^{(d)})$  a sequence of measure spaces, and  $\pi^{(d)}$  densities wrt  $\lambda^{(d)}$  for  $d \in \mathbb{N}$  the *problem size*.

#### Mixing time

$$\tau_{\epsilon} = \sup_{\pi_0} \min\{n : \|\pi_{n'} - \pi\|_{\mathsf{TV}} < \epsilon \quad \forall n' \ge n\}.$$

where

$$\|\pi_n - \pi\|_{\mathsf{TV}} = \sup_{A \subset \mathcal{X}} |\pi_n(A) - \pi(A)|$$

*P* is *rapidly mixing* if  $\tau_{\epsilon}$  is bounded above by a polynomial in *d*.

*P* is *torpidly mixing* if  $\tau_{\epsilon}$  is bounded below by an exponential in *d*.

4 同 1 4 三 1 4 三 1 4 二

For  ${\mathcal X}$  finite or compact, we have

A sequence of transition kernels  $P^{(d)}$  is rapidly mixing if **Gap** $(P^{(d)})$  decreases at most polynomially in *d*.

 $P^{(d)}$  is torpidly (or slowly) mixing if **Gap**( $P^{(d)}$ ) decreases exponentially in *d*.

Compare to geometric ergodicity, which requires only  $\operatorname{Gap}(P^{(d)}) > 0$ . (true for any  $|\mathcal{X}| < \infty$ .)

Convergence bounds for Markov chains:

- Spectral bounds (reversibilize if needed); or operator norm
  - E.g. conductance and canonical paths
- Coupling (minorization/drift)

Adaptive algorithms aren't Markov chains!

Produce non-Markovian, time-inhomogeneous, irreversible stochastic processes.

How to obtain bounds?

We obtain lower bounds on mixing times via the *hitting time* for subsets  $A \subset \mathcal{X}$ :

$$H_A = \min_i H_A^{(i)}$$
  $H_A^{(i)} = \min\{n : X_n^{(i)} \in A\}$ 

and involving the familiar *conductance* of a  $\pi$ -reversible Markov kernel T:

$$\Phi_{\mathcal{T}} = \inf_{\substack{A \subset \mathcal{X}:\\ 0 < \pi(A) < 1}} \Phi_{\mathcal{T}}(A) \qquad \Phi_{\mathcal{T}}(A) = \frac{\int_{A} \pi(dv) \mathcal{T}(v, A^{c})}{\pi(A) \pi(A^{c})}$$

 $\Phi_T(A)$  captures the probability of moving between A and  $A^c$  $\Phi_T$  quantifies the worst "bottleneck".

$$\Pr(H_A \le n) \le \pi(A) - \epsilon \quad \Rightarrow \quad \|\pi_n - \pi\|_{\mathsf{TV}} \ge \epsilon$$
$$\Rightarrow \quad \tau_\epsilon > n$$

 $\Rightarrow$   $\;$  We can lower bound mixing times via bounds on hitting

times.

- ₹ 🖬 🕨

- ₹ 🖬 🕨

#### Theorem (SW09)

For any  $\epsilon > 0$  and any  $A \subset \mathcal{X}$  such that  $0 < \pi^{(i)}(A) < 1$  for all *i*, the mixing time  $\tau_{\epsilon}^*$  of the MRAM satisfies:

$$au_{\epsilon}^* \geq (\pi(A) - \epsilon) \left[ cI \max_i \gamma(A, i) \Phi_{T_i}(A) \right]^{-1}$$

Here  $\gamma(A, i) = \min\{1, \pi^{(i)}(A)/\pi(A)\}$  is the *persistence* defined by Woodard,Schmidler,Huber (2007).

#### Note appearance of the conductance:

#### Corollary

For any  $0 < \epsilon < 1/4$ , the mixing time  $\tau_{\epsilon}^*$  of an adaptive sampler based on T, with I = 1, satisfies:

$$\tau_{\epsilon}^* \geq \frac{1}{4\Phi_{\mathcal{T}}}.$$

#### Corollary

Slow mixing of the Markov chain with transition kernel T implies slow mixing of any MRAM process based on T that has I = 1.

伺 ト く ヨ ト く ヨ ト

Note also the similarity of the bound for MRAM processes:

$$au_{\epsilon}^* \geq (\pi(A) - \epsilon) \left[ c I \max_i \gamma(A, i) \Phi_{T_i}(A) \right]^{-1}.$$

to the bound obtained by Woodard, Schmidler, Huber (2007) for non-adaptive swapping processes:

$$\tau_{\epsilon}^* \geq 2^{-8} \ln(2\epsilon)^{-1} \left[ \max_i \gamma(A, i) \Phi_{T_i}(A) \right]^{-1/2}$$

#### Mixtures of normals

$$\pi(z) = \frac{1}{2} N_M(z; -1_M, \sigma_1^2 \mathsf{I}_M) + \frac{1}{2} N_M(z; 1_M, \sigma_2^2 \mathsf{I}_M)$$

Theorem (WSH07a): Tempering is rapidly mixing for  $\sigma_1 = \sigma_2$ . Theorem (WSH07b): Tempering is torpidly mixing for  $\sigma_1 \neq \sigma_2$ .

Theorem (SW09): EES is torpidly mixing for  $\sigma_1 \neq \sigma_2$ .

伺 ト イヨト イヨト

### Similar (hitting time) argument gives:

### Theorem (SW09)

Haario and multi-chain samplers are torpidly mixing on mixture-of-normals problem.

-∢ ≣ ▶

*Finite sample* efficiency: Convergence as well as autocorrelation

$$\mathsf{MSE}(\hat{ heta}) = \mathsf{Bias}^2(\hat{ heta}) + \mathsf{Var}(\hat{ heta})$$

 $\Rightarrow$  MRAM and IAMC sampling can only improve autocorrelation piece!

- Mixing times and finite time convergence of adaptive MCMC.
- Combining adaptive strategies.
- Some cautionary notes about MIS kernels.

Suggests considering alternative "adaptation" strategies.

- I: IAMC (adaptive random walks, AMIS)
- II MRAM (equi-energy)

Suggests considering alternative "adaptation" strategies.

- I: IAMC (adaptive random walks, AMIS)
- II MRAM (equi-energy)
- (III) Modifying the stationary distribution
  - Wang-Landau, and generalizations (Atchade & Liu, Liang)
  - Multi-canonical
  - Metadynamics (Parisi et al)

Have received much interest in physics literature; recently adopted for statistical problems. (Liang, Atchade & Liu).

## Generalized Wang-Landau (Atchade & Liu, 2009)

Partition state space  $\mathcal{X} = \mathcal{X}_0 \cup \ldots \cup \mathcal{X}_k$  according to predefined energy levels  $-\infty \leq e_0 < e_1 < \cdots < e_k \leq \infty$ .

Goal: Sample from  $\tilde{\pi}(x) = \sum_{i=1}^{k} \frac{\pi(x)}{\pi(X_i)} \mathbf{1}_{X_i}(x)$  uniform energy

**Algorithm:** Adaptively estimate  $\hat{\pi}_n(i) \approx \pi(\mathcal{X}_i)$  by SA:  $\{\gamma_n\}$  a sequence of decreasing positive numbers. Initialize  $\phi_0(i) > 0$  for  $i = 1, \dots, k$ , and  $\hat{\pi}_0(i) = \frac{\phi_0(i)}{\sum_{i} \phi_0(i)}$ (i) Sample  $X_{n+1} \sim \sum_{i=1}^k \frac{\pi(x)}{\hat{\pi}_n(i)} \mathbf{1}_{\mathcal{X}_i}(x)$  by MH. (ii) Set  $\phi_{n+1}(i) = \phi_n(i) \left( 1 + \gamma_{a_n} \mathbf{1}_{\{X_{n+1} \in \mathcal{X}_i\}} \right); \ \hat{\pi}_{n+1}(i) = \frac{\phi_{n+1}(i)}{\sum_i \phi_{n+1}(j)}.$ (iii) If  $\max_{i} \left| v_{\kappa,n+1}(i) - \frac{1}{k} \right| \leq \frac{c}{k}$  where  $v_{\kappa,n}(i) = \frac{1}{n-\kappa} \sum_{i=\kappa+1}^{n} \mathbf{1}_{\{X_j \in \mathcal{X}_i\}}$ then set  $\kappa = n + 1$  and  $a_{n+1} = a_n + 1$ , otherwise  $a_{n+1} = a_n$ .

- 4 周 ト 4 月 ト 4 月 ト - 月

These ways of adapting address fundamentally different problems:

<u>I & II</u>: Improve mixing of chain among regions of target distribution *already visited* 

- Improves autocorrelation of chain
- In general cannot help in exploring previously unseen regions

Call these *Exploitation* methods.

III: Tries to push chain away from points "like" those already seen.

- Can help in finding new regions; improve mixing time.
- May suffer from high autocorrelation.

Call these *Exploration* methods.

伺下 イヨト イヨト

Note:

• Not rigorous statement for *all* IAMC methods on *all* targets; depends on form of kernel. But clear that method's power is essentially limited by these choices.

- ∢ ≣ ▶

- A 🗐 🕨

Note:

- Not rigorous statement for *all* IAMC methods on *all* targets; depends on form of kernel. But clear that method's power is essentially limited by these choices.
- Some authors (Craiu *et al*, Heaton & Schmidler) use multiple parallel chains to aid exploration in IAMC. Can help in practice but ultimately limited by ability to initialize well.

Note:

- Not rigorous statement for *all* IAMC methods on *all* targets; depends on form of kernel. But clear that method's power is essentially limited by these choices.
- Some authors (Craiu *et al*, Heaton & Schmidler) use multiple parallel chains to aid exploration in IAMC. Can help in practice but ultimately limited by ability to initialize well.
- No method will work for *all* problems some are provably hard (see e.g Schmidler & Woodard, in prep). Can hope for improved behavior on practical problems.

伺 ト イ ヨ ト イ ヨ ト

Can we combine types to achieve best of both? Yes but requires some care.

One approach: Mixture kernels

$$\mathcal{K}_{ ext{adapt}} = lpha \mathcal{K}_{ ext{exploit}} + (1 - lpha) \mathcal{K}_{ ext{explore}}$$

Suffers problems in multimodal examples (Wiehe & Schmidler, 2010).

Alternative approach:

Run exploration chain independently in parallel, but use samples to augment AMIS approximation.

# Exploration/Exploitation Algorithm (Wang & SS, 2010)

- Run two chains in parallel:  $X^{WL}$  and  $X^{AMIS+}$
- Solution Section 2 Every  $N_c$  iterations, update the proposal distribution for  $X^{\text{AMIS}+}$ .
- At iteration n = m \* N<sub>c</sub>, let E<sub>n</sub> be the energy ring of X<sup>AMIS+</sup><sub>n-1</sub>.
  Form KDE f by adding the samples {X<sup>WL</sup><sub>1</sub>,...,X<sup>WL</sup><sub>n</sub>} to those in E<sub>n</sub>.
- Propose  $X_n^{\text{AMIS}+}$  from  $\hat{f}_c$ .
- At other iterations, run the two chains independently.

But . . .

Problem 1: Performance of the WL algorithm depends heavily on a good choice of the energy rings  $E_0, \ldots, E_k$ : number, spacing, max.

Recommended heuristics:

Estimate highest energy, lowest, form geometric progression.



Figure: Normal mixture with modes at (-5,-5) and (5,5)



< Ξ

Example

Figure: Example 2, modes at (-5,-5) and (5,5)



Scott C. Schmidler

Exploration Vs. Exploitation in Adaptive Monte Carlo Sampling



(b) d = 4, fixed energy levels

Conductance argument yields provably slow mixing.

## Energy level adaptation scheme

Performance of the WL algorithm depends heavily on a good choice of the energy rings  $E_0, \ldots, E_k$ .

We introduce an adaptive scheme to make updating energy levels fully automatic:

Initialize by a geometric progression:

$$e_0 = \inf_x E(x) = 0, \ e_1 = 1, \ e_2 = r_e, \dots, E_{k-1} = r_e^{k-2}, E_k = \infty.$$

- Every  $n_{\text{split}}$  iterations: if any  $|\log(\phi_i) \log(\phi_{i+1})| > E$ , divide the *i*-th energy ring by adding a new  $e_{i+1}^* = e_i \times \sqrt{\frac{e_{i+1}}{e_i}}$ , again using geometric progression. Set  $\log(\phi_{i+1}^*) = 0$ .
- Also update the second largest  $e_i$ ;

$$E_{k-1}^* = \frac{E_{k-1}^2}{E_k}$$

Set  $\log(\phi_k^*) = 0$ .

**Algorithm:** Adaptively estimate  $\hat{\pi}_n(i) \approx \pi(\mathcal{X}_i)$  by SA:  $\{\gamma_n\}$  a sequence of decreasing positive numbers. Initialize  $\phi_0(i) > 0$  for  $i = 1, \dots, k$ , and  $\hat{\pi}_0(i) = \frac{\phi_0(i)}{\sum_i \phi_0(i)}$ (i) Sample  $X_{n+1} \sim \sum_{i=1}^k \frac{\pi(x)}{\hat{\pi}_n(i)} \mathbf{1}_{\mathcal{X}_i}(x)$  by MH. (ii) Set  $\phi_{n+1}(i) = \phi_n(i) (1 + \gamma_{a_n} \mathbf{1}_{\{X_{n+1} \in \mathcal{X}_i\}})$  and  $\hat{\pi}_{n+1}(i) = \frac{\phi_{n+1}(i)}{\sum \phi_{n+1}(i)}.$ (iii) If  $\max_{i} \left| v_{\kappa,n+1}(i) - \frac{1}{k} \right| \leq \frac{c}{k}$  where  $v_{\kappa,n}(i) = \frac{1}{n-\kappa} \sum_{i=\kappa+1}^{n} \mathbf{1}_{\{X_j \in \mathcal{X}_i\}}$ then set  $\kappa = n + 1$  and  $a_{n+1} = a_n + 1$ , otherwise  $a_{n+1} = a_n$ . (iv)\* For every  $n_{\text{split}}$  iterations, adaptively update  $E = \{E_i\}$ .

伺い イラト イラト



(c) d = 4, update internal energy levels

◆□ > ◆□ > ◆臣 > ◆臣 > ○臣 ○ のへで

Problem 2: Wang-Landau for MC integration: converges in limit, but can be slow. WL inefficient for MC integration.

Reweighting complicated due to WL process.

We use importance sampling:

- **(**) Use all samples to find a kernel density estimate  $\hat{f}$ .
- **2** Importance resampling: compute importance weights  $w_i = \frac{h(x_i)}{\hat{f}(x_i)}$  and resample  $x_1, \ldots, x_m$ .
- **③** Form kernel density estimate  $\hat{\pi}$  from  $x_1, \ldots, x_m$ .

### Importance sampling vs ergodic averaging



(d)

э

- Run two chains in parallel:  $X^{AE-WL}$  and  $X^{AMIS+}$
- Solution Section 2 Every  $N_c$  iterations, update the proposal distribution for  $X^{\text{AMIS}+}$ .
- At iteration n = m \* N<sub>c</sub>, let E<sub>n</sub> be the energy ring of X<sub>n-1</sub><sup>AMIS+</sup>.
  Form KDE f̂ by adding the samples {X<sub>1</sub><sup>AE-WL</sup>,...,X<sub>n</sub><sup>AE-WL</sup>} to those in E<sub>n</sub>.
- Propose  $X_n^{\text{AMIS}+}$  from  $\hat{f}_c$ .
- At other iterations, run the two chains independently.

同 ト イヨ ト イヨ ト 三 ヨ

### Example: Trimodal distribution in d = 2

$$\pi(x) = \frac{1}{3}N(x; [-3, -3]^T, I) + \frac{1}{3}N(x; [7, 7]^T, I) + \frac{1}{3}N(x; [5, -5]^T)$$





(e) AMIS

< 17 >

Э

2



0

10

20

Lag

(f) AMIS + WL

20

Lag

30

40

Э

50

æ





(h)

<ロ> <同> <同> < 同> < 同>

2

Bimodal distribution in d = 3:

$$\pi(x) = \frac{1}{2}N(x; [-7, -7, -7]^T, I) + \frac{1}{2}N(x; [7, 7, 7]^T, I)$$

(4回) (4 回) (4 回)

-2



(i) AMIS

<ロ> <同> <同> < 同> < 同>

2



(j) WL

イロン イロン イヨン イヨン

2



(k) AMIS + WL

Scott C. Schmidler Exploration Vs. Exploitation in Adaptive Monte Carlo Sampling

æ

Э



(I) KL divergence from target

2

3

- Many ways of "adapting" an MCMC algorithm based on sample path exist; many can be shown to satisfy LLNs.
- Purpose of adaptation is to improve *rate* of convergence.
- Convergence of MC estimators involves *both* bias *and* variance.
- Different adaptation strategies can be understood as improving one or the other.
- Can obtain improved algorithms by combining strategies of different types.

伺 ト イ ヨ ト イ ヨ ト

- Dawn Woodard (Cornell)
- Jianyu Wang (Duke Statistics)
- Chunlin Ji (Kuang-Chi Inst)
- Kevin Wiehe (Duke postdoc)



Atchadé, Y. F. (2009).

Resampling from the past to improve on MCMC algorithms. *Far East Journal of Theoretical Probability*, 27:81–99.



#### Holmes, C. and Held, L. (2006).

Bayesian auxiliary variable models for binary and multinomial regression. *Bayesian Analysis*, 1:145–168.



Ji, C. and Schmidler, S. C. (2009). Adaptive Markov chain Monte Carlo for Bayesian variable selection. *J. Comp. Graph. Stat.*, (to appear).



Kou, S. C., Zhou, Q., and Wong, W. H. (2006). Equi-energy sampler with applications in statistical inference and statistical

Equi-energy sampler with applications in statistical inference and statis mechanics.

Ann. Statist., 34(4):1581–1619.



Robbins, H. and Monro, S. (1951). A stochastic approximation method. *Ann. Math. Stat.*, 22:400–407.



Roberts, G. O. and Rosenthal, J. S. (2007). Coupling and ergodicity of adaptive MCMC. *J. Appl. Prob.*, 44:458–475.

・ 同 ト ・ ヨ ト ・ ヨ ト