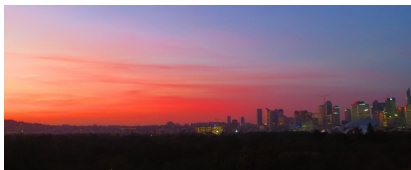


# Rao-Blackwellisation for accelerating Metropolis-Hastings

CHRISTIAN P. ROBERT

Université Paris-Dauphine, Paris & University of Warwick, Coventry

with M. Banterle, G. Casella, R. Douc, C. Grazian, & A. Lee



- 1 Rao-Blackwellisation 101
- 2 Vanilla Rao-Blackwellisation
- 3 Delayed acceptance
  - motivating example
  - Proposed solution
  - Validation of the method
  - Optimizing DA

- 1 Rao-Blackwellisation 101
- 2 Vanilla Rao-Blackwellisation
- 3 Delayed acceptance

Given a density  $f(\cdot)$  to simulate take  $g(\cdot)$  density such that

$$f(x) \leq Mg(x)$$

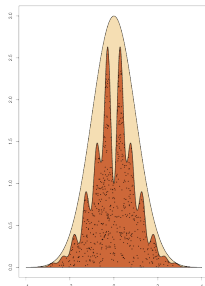
for  $M \geq 1$

To simulate  $X \sim f$ , it is sufficient to generate

$$Y \sim g \quad U|Y = y \sim \mathcal{U}(0, Mg(y))$$

until

$$0 < u < f(y)$$



## [Exercise 3.33, MCSM]

Raw outcome: id sequences  $Y_1, Y_2, \dots, Y_t \sim g$  and

$U_1, U_2, \dots, U_t \sim \mathcal{U}(0, 1)$

Random number of accepted  $Y_i$ 's

$$\mathbb{P}(N = n) = \binom{n-1}{t-1} (1/M)^t (1 - 1/M)^{n-t} ,$$

## [Exercise 3.33, MCSM]

Raw outcome: iid sequences  $Y_1, Y_2, \dots, Y_t \sim g$  and

$U_1, U_2, \dots, U_t \sim \mathcal{U}(0, 1)$

Joint density of  $(N, \mathbf{Y}, \mathbf{U})$

$$\begin{aligned} & \mathbb{P}(N = n, Y_1 \leq y_1, \dots, Y_n \leq y_n, U_1 \leq u_1, \dots, U_n \leq u_n) \\ &= \int_{-\infty}^{y_n} g(t_n)(u_n \wedge w_n) dt_n \int_{-\infty}^{y_1} \dots \int_{-\infty}^{y_{n-1}} g(t_1) \dots g(t_{n-1}) \\ & \quad \times \sum_{(i_1, \dots, i_{t-1})} \prod_{j=1}^{t-1} (w_{i_j} \wedge u_{i_j}) \prod_{j=t}^{n-1} (u_{i_j} - w_{i_j})^+ dt_1 \dots dt_{n-1}, \end{aligned}$$

where  $w_i = f(y_i)/Mg(y_i)$  and sum over all subsets of  $\{1, \dots, n-1\}$  of size  $t-1$

## [Exercise 3.33, MCSM]

Raw outcome: iid sequences  $Y_1, Y_2, \dots, Y_t \sim g$  and

$U_1, U_2, \dots, U_t \sim \mathcal{U}(0, 1)$

Marginal joint density of  $(Y_i, U_i) | N = n, i < n$

$$\begin{aligned} & \mathbb{P}(N = n, Y_1 \leq y, U_1 \leq u_1) \\ &= \binom{n-1}{t-1} \left(\frac{1}{M}\right)^{t-1} \left(1 - \frac{1}{M}\right)^{n-t-1} \\ & \times \left[ \frac{t-1}{n-1} (w_1 \wedge u_1) \left(1 - \frac{1}{M}\right) + \frac{n-t}{n-1} (u_1 - w_1)^+ \left(\frac{1}{M}\right) \right] \int_{-\infty}^y g(t_1) dt_1 \end{aligned}$$

and marginal distribution of  $Y_i$

$$\begin{aligned} m(y) &= t^{-1/n-1} f(y) + n^{-t/n-1} \frac{g(y) - \rho f(y)}{1 - \rho} \\ \mathbb{P}(U_1 \leq w(y) | Y_1 = y, N = n) &= \frac{g(y) w(y) M^{t-1/n-1}}{m(y)} \end{aligned}$$

Accept-reject sample  $(X_1, \dots, X_m)$  associated with  $(U_1, \dots, U_N)$   
and  $(Y_1, \dots, Y_N)$

$N$  is stopping time for acceptance of  $m$  variables among  $Y_j$ 's

Rewrite estimator of  $\mathbb{E}[h]$  as

$$\frac{1}{m} \sum_{i=1}^m h(X_i) = \frac{1}{m} \sum_{j=1}^N h(Y_j) \mathbb{I}_{U_j \leq w_j},$$

with  $w_j = f(Y_j)/Mg(Y_j)$

[Casella and Robert (1996)]

**Rao-Blackwellisation:** smaller variance produced by integrating out the  $U_i$ 's,

$$\frac{1}{m} \sum_{j=1}^N \mathbb{E}[\mathbb{I}_{U_j \leq w_j} | N, Y_1, \dots, Y_N] h(Y_j) = \frac{1}{m} \sum_{i=1}^N \rho_i h(Y_i),$$

where ( $i < n$ )

$$\begin{aligned} \rho_i &= \mathbb{P}(U_i \leq w_i | N = n, Y_1, \dots, Y_n) \\ &= w_i \frac{\sum_{(i_1, \dots, i_{m-2})} \prod_{j=1}^{m-2} w_{i_j} \prod_{j=m-1}^{n-2} (1 - w_{i_j})}{\sum_{(i_1, \dots, i_{m-1})} \prod_{j=1}^{m-1} w_{i_j} \prod_{j=m}^{n-1} (1 - w_{i_j})}, \end{aligned}$$

and  $\rho_n = 1$ .

Numerator sum over all subsets of  $\{1, \dots, i-1, i+1, \dots, n-1\}$  of size  $m-2$ , and denominator sum over all subsets of size  $m-1$

[Casella and Robert (1996)]

## extension to Metropolis–Hastings case

Sample produced by Metropolis–Hastings algorithm

$$x^{(1)}, \dots, x^{(T)}$$

based on two samples,

$$y_1, \dots, y_T \quad \text{and} \quad u_1, \dots, u_T$$

[Casella and Robert (1996)]

## extension to Metropolis–Hastings case

Sample produced by Metropolis–Hastings algorithm

$$x^{(1)}, \dots, x^{(T)}$$

based on two samples,

$$y_1, \dots, y_T \text{ and } u_1, \dots, u_T$$

Ergodic mean rewritten as

$$\delta^{MH} = \frac{1}{T} \sum_{t=1}^T h(x^{(t)}) = \frac{1}{T} \sum_{t=1}^T h(y_t) \sum_{i=t}^T \mathbb{I}_{x^{(i)}=y_t}$$

[Casella and Robert (1996)]

## extension to Metropolis–Hastings case

Sample produced by Metropolis–Hastings algorithm

$$x^{(1)}, \dots, x^{(T)}$$

based on two samples,

$$y_1, \dots, y_T \quad \text{and} \quad u_1, \dots, u_T$$

Conditional expectation

$$\begin{aligned}\delta^{RB} &= \frac{1}{T} \sum_{t=1}^T h(y_t) \mathbb{E} \left[ \sum_{i=t}^T \mathbb{I} X^{(i)} = y_t \middle| y_1, \dots, y_T \right] \\ &= \frac{1}{T} \sum_{t=1}^T h(y_t) \left( \sum_{i=t}^T \mathbb{P}(X^{(i)} = y_t | y_1, \dots, y_T) \right)\end{aligned}$$

with smaller variance

[Casella and Robert (1996)]

Take

$$\rho_{ij} = \frac{f(y_j)/q(y_j|y_i)}{f(y_i)/q(y_i|y_j)} \wedge 1 \quad (j > i),$$

$$\bar{\rho}_{ij} = \rho_{ij}q(y_{j+1}|y_j), \quad \underline{\rho}_{ij} = (1 - \rho_{ij})q(y_{j+1}|y_i) \quad (i < j < T),$$

$$\zeta_{jj} = 1, \quad \zeta_{jt} = \prod_{l=j+1}^t \underline{\rho}_{jl} \quad (i < j < T),$$

$$\tau_0 = 1, \quad \tau_j = \sum_{t=0}^{j-1} \tau_t \zeta_{t(j-1)} \bar{\rho}_{tj}, \quad \tau_T = \sum_{t=0}^{T-1} \tau_t \zeta_{t(T-1)} \rho_{tT} \quad (i < T),$$

$$\omega_T^i = 1, \quad \omega_i^j = \bar{\rho}_{ji} \omega_{i+1}^i + \underline{\rho}_{ji} \omega_{i+1}^j \quad (0 \leq j < i < T).$$

[Casella and Robert (1996)]

## Theorem

The estimator  $\delta^{RB}$  satisfies

$$\delta^{RB} = \frac{\sum_{i=0}^T \varphi_i h(y_i)}{\sum_{i=0}^{T-1} \tau_i \zeta_{i(T-1)}},$$

with  $(i < T)$

$$\varphi_i = \tau_i \left[ \sum_{j=i}^{T-1} \zeta_{ij} \omega_{j+1}^i + \zeta_{i(T-1)} (1 - \rho_{iT}) \right]$$

and  $\varphi_T = \tau_T$ .

[Casella and Robert (1996)]

- 1 Rao-Blackwellisation 101
- 2 Vanilla Rao-Blackwellisation
- 3 Delayed acceptance



# Some properties of the Metropolis–Hastings algorithm

Alternative representation of Metropolis–Hastings estimator  $\delta$  as

$$\delta = \frac{1}{n} \sum_{t=1}^n h(x^{(t)}) = \frac{1}{n} \sum_{i=1}^{M_n} n_i h(z_i),$$

where

- $z_i$ 's are the accepted  $y_j$ 's,
- $M_n$  is the number of accepted  $y_j$ 's till time  $n$ ,
- $n_i$  is the number of times  $z_i$  appears in the sequence  $(x^{(t)})_t$ .

## The "accepted candidates"

Define

$$\tilde{q}(\cdot|z_i) = \frac{\alpha(z_i, \cdot) q(\cdot|z_i)}{p(z_i)} \leq \frac{q(\cdot|z_i)}{p(z_i)}$$

where  $p(z_i) = \int \alpha(z_i, y) q(y|z_i) dy$

To simulate from  $\tilde{q}(\cdot|z_i)$

- 1 Propose a candidate  $y \sim q(\cdot|z_i)$
- 2 Accept with probability

$$\tilde{q}(y|z_i) \bigg/ \left( \frac{q(y|z_i)}{p(z_i)} \right) = \alpha(z_i, y)$$

Otherwise, reject it and starts again.

► this is the transition of the HM algorithm

## The "accepted candidates"

Define

$$\tilde{q}(\cdot|z_i) = \frac{\alpha(z_i, \cdot) q(\cdot|z_i)}{p(z_i)} \leq \frac{q(\cdot|z_i)}{p(z_i)}$$

where  $p(z_i) = \int \alpha(z_i, y) q(y|z_i) dy$

The transition kernel  $\tilde{q}$  admits  $\tilde{\pi}$  as a stationary distribution:

$$\tilde{\pi}(x) \tilde{q}(y|x) = \underbrace{\frac{\pi(x)p(x)}{\int \pi(u)p(u)du}}_{\tilde{\pi}(x)} \underbrace{\frac{\alpha(x, y)q(y|x)}{p(x)}}_{\tilde{q}(y|x)}$$

## The "accepted candidates"

Define

$$\tilde{q}(\cdot|z_i) = \frac{\alpha(z_i, \cdot) q(\cdot|z_i)}{p(z_i)} \leq \frac{q(\cdot|z_i)}{p(z_i)}$$

where  $p(z_i) = \int \alpha(z_i, y) q(y|z_i) dy$

The transition kernel  $\tilde{q}$  admits  $\tilde{\pi}$  as a stationary distribution:

$$\tilde{\pi}(x) \tilde{q}(y|x) = \frac{\pi(x) \alpha(x, y) q(y|x)}{\int \pi(u) p(u) du}$$

## The "accepted candidates"

Define

$$\tilde{q}(\cdot|z_i) = \frac{\alpha(z_i, \cdot) q(\cdot|z_i)}{p(z_i)} \leq \frac{q(\cdot|z_i)}{p(z_i)}$$

where  $p(z_i) = \int \alpha(z_i, y) q(y|z_i) dy$

The transition kernel  $\tilde{q}$  admits  $\tilde{\pi}$  as a stationary distribution:

$$\tilde{\pi}(x) \tilde{q}(y|x) = \frac{\pi(y) \alpha(y, x) q(x|y)}{\int \pi(u) p(u) du}$$

## The "accepted candidates"

Define

$$\tilde{q}(\cdot|z_i) = \frac{\alpha(z_i, \cdot) q(\cdot|z_i)}{p(z_i)} \leq \frac{q(\cdot|z_i)}{p(z_i)}$$

where  $p(z_i) = \int \alpha(z_i, y) q(y|z_i) dy$

The transition kernel  $\tilde{q}$  admits  $\tilde{\pi}$  as a stationary distribution:

$$\tilde{\pi}(x) \tilde{q}(y|x) = \tilde{\pi}(y) \tilde{q}(x|y),$$

Lemma (Douc & X., AoS, 2011)

The sequence  $(z_i, n_i)$  satisfies

- 1  $(z_i, n_i)_i$  is a Markov chain;
- 2  $z_{i+1}$  and  $n_i$  are independent given  $z_i$ ;
- 3  $n_i$  is distributed as a geometric random variable with probability parameter

$$p(z_i) := \int \alpha(z_i, y) q(y|z_i) dy; \quad (1)$$

- 4  $(z_i)_i$  is a Markov chain with transition kernel  $\tilde{Q}(z, dy) = \tilde{q}(y|z)dy$  and stationary distribution  $\tilde{\pi}$  such that

$$\tilde{q}(\cdot|z) \propto \alpha(z, \cdot) q(\cdot|z) \quad \text{and} \quad \tilde{\pi}(\cdot) \propto \pi(\cdot)p(\cdot).$$

Lemma (Douc & X., AoS, 2011)

The sequence  $(z_i, n_i)$  satisfies

- ①  $(z_i, n_i)_i$  is a Markov chain;
- ②  $z_{i+1}$  and  $n_i$  are **independent** given  $z_i$ ;
- ③  $n_i$  is distributed as a geometric random variable with probability parameter

$$p(z_i) := \int \alpha(z_i, y) q(y|z_i) dy; \quad (1)$$

- ④  $(z_i)_i$  is a Markov chain with transition kernel  $\tilde{Q}(z, dy) = \tilde{q}(y|z)dy$  and stationary distribution  $\tilde{\pi}$  such that

$$\tilde{q}(\cdot|z) \propto \alpha(z, \cdot) q(\cdot|z) \quad \text{and} \quad \tilde{\pi}(\cdot) \propto \pi(\cdot)p(\cdot).$$

Lemma (Douc & X., AoS, 2011)

The sequence  $(z_i, n_i)$  satisfies

- ①  $(z_i, n_i)_i$  is a Markov chain;
- ②  $z_{i+1}$  and  $n_i$  are **independent** given  $z_i$ ;
- ③  $n_i$  is distributed as a geometric random variable with probability parameter

$$p(z_i) := \int \alpha(z_i, y) q(y|z_i) dy; \quad (1)$$

- ④  $(z_i)_i$  is a Markov chain with transition kernel  $\tilde{Q}(z, dy) = \tilde{q}(y|z)dy$  and stationary distribution  $\tilde{\pi}$  such that

$$\tilde{q}(\cdot|z) \propto \alpha(z, \cdot) q(\cdot|z) \quad \text{and} \quad \tilde{\pi}(\cdot) \propto \pi(\cdot)p(\cdot).$$

Lemma (Douc & X., AoS, 2011)

The sequence  $(z_i, n_i)$  satisfies

- ①  $(z_i, n_i)_i$  is a Markov chain;
- ②  $z_{i+1}$  and  $n_i$  are independent given  $z_i$ ;
- ③  $n_i$  is distributed as a geometric random variable with probability parameter

$$p(z_i) := \int \alpha(z_i, y) q(y|z_i) dy; \quad (1)$$

- ④  $(z_i)_i$  is a Markov chain with transition kernel  $\tilde{Q}(z, dy) = \tilde{q}(y|z)dy$  and stationary distribution  $\tilde{\pi}$  such that

$$\tilde{q}(\cdot|z) \propto \alpha(z, \cdot) q(\cdot|z) \quad \text{and} \quad \tilde{\pi}(\cdot) \propto \pi(\cdot)p(\cdot).$$

- 1 A natural idea:

$$\delta^* = \frac{1}{n} \sum_{i=1}^{M_n} \frac{h(\mathfrak{z}_i)}{p(\mathfrak{z}_i)},$$

- 1 A natural idea:

$$\delta^* \simeq \frac{\sum_{i=1}^{M_n} \frac{h(\mathfrak{z}_i)}{p(\mathfrak{z}_i)}}{\sum_{i=1}^{M_n} \frac{1}{p(\mathfrak{z}_i)}} = \frac{\sum_{i=1}^{M_n} \frac{\pi(\mathfrak{z}_i)}{\tilde{\pi}(\mathfrak{z}_i)} h(\mathfrak{z}_i)}{\sum_{i=1}^{M_n} \frac{\pi(\mathfrak{z}_i)}{\tilde{\pi}(\mathfrak{z}_i)}}.$$

- 1 A natural idea:

$$\delta^* \simeq \frac{\sum_{i=1}^{M_n} \frac{h(\mathfrak{z}_i)}{p(\mathfrak{z}_i)}}{\sum_{i=1}^{M_n} \frac{1}{p(\mathfrak{z}_i)}} = \frac{\sum_{i=1}^{M_n} \frac{\pi(\mathfrak{z}_i)}{\tilde{\pi}(\mathfrak{z}_i)} h(\mathfrak{z}_i)}{\sum_{i=1}^{M_n} \frac{\pi(\mathfrak{z}_i)}{\tilde{\pi}(\mathfrak{z}_i)}}.$$

- 2 But  $p$  not available in closed form.

- 1 A natural idea:

$$\delta^* \simeq \frac{\sum_{i=1}^{M_n} \frac{h(\mathbf{z}_i)}{p(\mathbf{z}_i)}}{\sum_{i=1}^{M_n} \frac{1}{p(\mathbf{z}_i)}} = \frac{\sum_{i=1}^{M_n} \frac{\pi(\mathbf{z}_i)}{\tilde{\pi}(\mathbf{z}_i)} h(\mathbf{z}_i)}{\sum_{i=1}^{M_n} \frac{\pi(\mathbf{z}_i)}{\tilde{\pi}(\mathbf{z}_i)}}.$$

- 2 But  $p$  not available in closed form.
- 3 The geometric  $\mathbf{n}_i$  is the replacement, an obvious solution that is used in the original Metropolis–Hastings estimate since  $\mathbb{E}[\mathbf{n}_i] = 1/p(\mathbf{z}_i)$ .

The crude estimate of  $1/p(z_i)$ ,

$$n_i = 1 + \sum_{j=1}^{\infty} \prod_{\ell \leq j} \mathbb{I}\{u_{\ell} \geq \alpha(z_i, y_{\ell})\} ,$$

can be improved:

Lemma (Douc & X., AoS, 2011)

If  $(y_j)_j$  is an iid sequence with distribution  $q(y|z_i)$ , the quantity

$$\hat{\xi}_i = 1 + \sum_{j=1}^{\infty} \prod_{\ell \leq j} \{1 - \alpha(z_i, y_{\ell})\}$$

is an unbiased estimator of  $1/p(z_i)$  whose variance, conditional on  $z_i$ , is lower than the conditional variance of  $n_i$ ,  $\{1 - p(z_i)\}/p^2(z_i)$ .

$$\hat{\xi}_i = 1 + \sum_{j=1}^{\infty} \prod_{\ell \leq j} \{1 - \alpha(\mathfrak{z}_i, y_\ell)\}$$

- 1 Infinite sum but finite with at least positive probability:

$$\alpha(x^{(t)}, y_t) = \min \left\{ 1, \frac{\pi(y_t)}{\pi(x^{(t)})} \frac{q(x^{(t)}|y_t)}{q(y_t|x^{(t)})} \right\}$$

For example: take a symmetric random walk as a proposal.

- 2 What if we wish to be sure that the sum is finite?

Finite horizon  $k$  version:

$$\hat{\xi}_i^k = 1 + \sum_{j=1}^{\infty} \prod_{1 \leq \ell \leq k \wedge j} \{1 - \alpha(\mathfrak{z}_i, y_j)\} \prod_{k+1 \leq \ell \leq j} \mathbb{I}\{u_\ell \geq \alpha(\mathfrak{z}_i, y_\ell)\}$$

$$\hat{\xi}_i = 1 + \sum_{j=1}^{\infty} \prod_{\ell \leq j} \{1 - \alpha(\mathfrak{z}_i, y_\ell)\}$$

- ① Infinite sum but finite with at least positive probability:

$$\alpha(x^{(t)}, y_t) = \min \left\{ 1, \frac{\pi(y_t)}{\pi(x^{(t)})} \frac{q(x^{(t)}|y_t)}{q(y_t|x^{(t)})} \right\}$$

For example: take a symmetric random walk as a proposal.

- ② What if we wish to be sure that the sum is finite?

Finite horizon  $k$  version:

$$\hat{\xi}_i^k = 1 + \sum_{j=1}^{\infty} \prod_{1 \leq \ell \leq k \wedge j} \{1 - \alpha(\mathfrak{z}_i, y_j)\} \prod_{k+1 \leq \ell \leq j} \mathbb{I}\{u_\ell \geq \alpha(\mathfrak{z}_i, y_\ell)\}$$

## Proposition (Douc & X., AoS, 2011)

If  $(y_j)_j$  is an iid sequence with distribution  $q(y|z_i)$  and  $(u_j)_j$  is an iid uniform sequence, for any  $k \geq 0$ , the quantity

$$\hat{\xi}_i^k = 1 + \sum_{j=1}^{\infty} \prod_{1 \leq \ell \leq k \wedge j} \{1 - \alpha(z_i, y_\ell)\} \prod_{k+1 \leq \ell \leq j} \mathbb{I}\{u_\ell \geq \alpha(z_i, y_\ell)\}$$

is an unbiased estimator of  $1/p(z_i)$  with an almost sure finite number of terms.

## Proposition (Douc & X., AoS, 2011)

If  $(y_j)_j$  is an iid sequence with distribution  $q(y|z_i)$  and  $(u_j)_j$  is an iid uniform sequence, for any  $k \geq 0$ , the quantity

$$\hat{\xi}_i^k = 1 + \sum_{j=1}^{\infty} \prod_{1 \leq \ell \leq k \wedge j} \{1 - \alpha(z_i, y_\ell)\} \prod_{k+1 \leq \ell \leq j} \mathbb{I}\{u_\ell \geq \alpha(z_i, y_\ell)\}$$

is an unbiased estimator of  $1/p(z_i)$  with an almost sure finite number of terms. Moreover, for  $k \geq 1$ ,

$$\mathbb{V}_{\xi_i^k}^k = \frac{1 - p(z_i)}{p^2(z_i)} - \frac{1 - (1 - 2p(z_i) + r(z_i))^k}{2p(z_i) - r(z_i)} \left( \frac{2 - p(z_i)}{p^2(z_i)} \right) (p(z_i) - r(z_i)),$$

where  $p(z_i) := \int \alpha(z_i, y) q(y|z_i) dy$ . and  $r(z_i) := \int \alpha^2(z_i, y) q(y|z_i) dy$ .

## Proposition (Douc & X., AoS, 2011)

If  $(y_j)_j$  is an iid sequence with distribution  $q(y|z_i)$  and  $(u_j)_j$  is an iid uniform sequence, for any  $k \geq 0$ , the quantity

$$\hat{\xi}_i^k = 1 + \sum_{j=1}^{\infty} \prod_{1 \leq \ell \leq k \wedge j} \{1 - \alpha(z_i, y_\ell)\} \prod_{k+1 \leq \ell \leq j} \mathbb{I}\{u_\ell \geq \alpha(z_i, y_\ell)\}$$

is an unbiased estimator of  $1/p(z_i)$  with an almost sure finite number of terms. Therefore, we have

$$\mathbb{V}_{\hat{\xi}_i z_i} \leq \mathbb{V}_{\hat{\xi}_i^k z_i} \leq \mathbb{V}_{\hat{\xi}_i^0 z_i} = \mathbb{V}_{n_i z_i}.$$

- 1 Rao-Blackwellisation 101
- 2 Vanilla Rao-Blackwellisation
- 3 Delayed acceptance



# Non-informative inference for mixture models

Standard mixture of distributions model

$$\sum_{i=1}^k w_i f(x|\theta_i), \quad \text{with} \quad \sum_{i=1}^k w_i = 1. \quad (1)$$

[Titterington et al., 1985; Frühwirth-Schnatter (2006)]

Jeffreys' prior for mixture not available due to computational reasons : it has not been tested so far

[Jeffreys, 1939]

**Warning:** Jeffreys' prior improper in some settings

[Grazian & Robert, 2015]

# Non-informative inference for mixture models

Grazian & Robert (2015) consider genuine Jeffreys' prior for complete set of parameters in (1), deduced from Fisher's information matrix

Computation of prior density costly, relying on many integrals like

$$\int_{\mathcal{X}} \frac{\partial^2 \log \left[ \sum_{i=1}^k w_i f(x|\theta_i) \right]}{\partial \theta_h \partial \theta_j} \left[ \sum_{i=1}^k w_i f(x|\theta_i) \right] dx$$

Integrals with no analytical expression, hence involving numerical or Monte Carlo (costly) integration

# Non-informative inference for mixture models

When building Metropolis-Hastings proposal over  $(w_i, \theta_i)$ 's, **prior ratio more expensive** than likelihood and proposal ratios

**Suggestion:** split the acceptance rule

$$\alpha(x, y) := 1 \wedge r(x, y), \quad r(x, y) := \frac{\pi(y|\mathcal{D})q(y, x)}{\pi(x|\mathcal{D})q(x, y)}$$

into

$$\tilde{\alpha}(x, y) := \left( 1 \wedge \frac{f(\mathcal{D}|y)q(y, x)}{f(\mathcal{D}|x)q(x, y)} \right) \times \left( 1 \wedge \frac{\pi(y)}{\pi(x)} \right)$$

Simulation from posterior distribution with large sample size  $n$

- Computing time at least of order  $O(n)$
- solutions using likelihood decomposition

$$\prod_{i=1}^n \ell(\theta|x_i)$$

and handling subsets on different processors (CPU), graphical units (GPU), or computers

[Korattikara et al. (2013), Scott et al. (2013)]

- no consensus on method of choice, with instabilities from removing most prior input and uncalibrated approximations

[Neiswanger et al. (2013), Wang and Dunson (2013)]

*“There is no problem an absence of decision cannot solve.” Anonymous*

Given  $\alpha(x, y) := 1 \wedge r(x, y)$ , factorise

$$r(x, y) = \prod_{k=1}^d \rho_k(x, y)$$

under constraint  $\rho_k(x, y) = \rho_k(y, x)^{-1}$

Delayed Acceptance Markov kernel given by

$$\tilde{P}(x, A) := \int_A q(x, y) \tilde{\alpha}(x, y) dy + \left(1 - \int_X q(x, y) \tilde{\alpha}(x, y) dy\right) \mathbf{1}_A(x)$$

where

$$\tilde{\alpha}(x, y) := \prod_{k=1}^d \{1 \wedge \rho_k(x, y)\}.$$

*“There is no problem an absence of decision cannot solve.” Anonymous*

---

**Algorithm 1** Delayed Acceptance

---

To sample from  $\tilde{P}(x, \cdot)$ :

- ① Sample  $y \sim Q(x, \cdot)$ .
  - ② For  $k = 1, \dots, d$ :
    - with probability  $1 \wedge \rho_k(x, y)$  continue
    - otherwise stop and output  $x$
  - ③ Output  $y$
- 

Arrange terms in product so that most computationally intensive ones calculated ‘at the end’ hence least often

*“There is no problem an absence of decision cannot solve.” Anonymous*

---

**Algorithm 1** Delayed Acceptance

---

To sample from  $\tilde{P}(x, \cdot)$ :

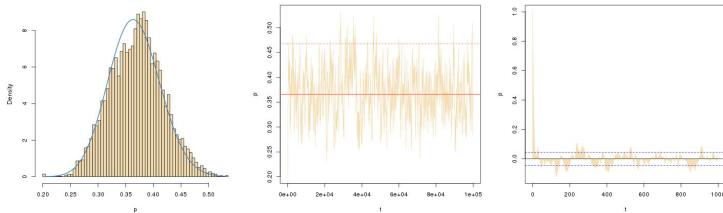
- ① Sample  $y \sim Q(x, \cdot)$ .
  - ② For  $k = 1, \dots, d$ :
    - with probability  $1 \wedge \rho_k(x, y)$  continue
    - otherwise stop and output  $x$
  - ③ Output  $y$
- 

Generalization of Fox and Nicholls (1997) and Christen and Fox (2005), where testing for acceptance with approximation before computing exact likelihood first suggested

More recent occurrences in literature

[Golightly et al. (2014), Shestopaloff and Neal (2013)]

- Delayed Acceptance *efficiently* reduces computing cost only when approximation  $\tilde{\pi}$  is “good enough” or “flat enough”
- Probability of acceptance always smaller than in the original Metropolis–Hastings scheme
- Decomposition of original data in likelihood bits may however lead to deterioration of algorithmic properties without impacting computational efficiency...
- ...e.g., case of a term explosive in  $x = 0$  and computed by itself: leaving  $x = 0$  near impossible



**Figure :** (left) Fit of delayed Metropolis–Hastings algorithm on a Beta-binomial posterior  $p|x \sim Be(x + a, n + b - x)$  when  $N = 100$ ,  $x = 32$ ,  $a = 7.5$  and  $b = .5$ . Binomial  $\mathcal{B}(N, p)$  likelihood replaced with product of 100 Bernoulli terms. Histogram based on  $10^5$  iterations, with overall acceptance rate of 9%; (centre) raw sequence of  $p$ 's in Markov chain; (right) autocorrelogram of the above sequence.

Delayed Acceptance intended for likelihoods or priors, but not a clear solution for “Big Data” problems

- ① all product terms must be computed
- ② all terms previously computed either stored for future comparison or recomputed
- ③ sequential approach limits parallel gains...
- ④ ...unless prefetching scheme added to delays

[Angelino et al. (2014), Strid (2010)]

## Lemma (1)

*For any Markov chain with transition kernel  $\Pi$  of the form*

$$\Pi(x, A) = \int_A q(x, y) a(x, y) dy + \left(1 - \int_X q(x, y) a(x, y) dy\right) \mathbf{1}_A(x),$$

*and satisfying detailed balance, the function  $a(\cdot)$  satisfies (for  $\pi$ -a.e.  $x, y$ )*

$$\frac{a(x, y)}{a(y, x)} = r(x, y).$$

## Lemma (2)

$(\tilde{X}_n)_{n \geq 1}$ , the Markov chain associated with  $\tilde{P}$ , is a  $\pi$ -reversible Markov chain.

## Proof.

From Lemma 1 we just need to check that

$$\begin{aligned} \frac{\tilde{\alpha}(x, y)}{\tilde{\alpha}(y, x)} &= \prod_{k=1}^d \frac{1 \wedge \rho_k(x, y)}{1 \wedge \rho_k(y, x)} \\ &= \prod_{k=1}^d \rho_k(x, y) = r(x, y), \end{aligned}$$

since  $\rho_k(y, x) = \rho_k(x, y)^{-1}$  and  $(1 \wedge a)/(1 \wedge a^{-1}) = a$



The acceptance probability ordering

$$\tilde{\alpha}(x, y) = \prod_{k=1}^d \{1 \wedge \rho_k(x, y)\} \leq 1 \wedge \prod_{k=1}^d \rho_k(x, y) = 1 \wedge r(x, y) = \alpha(x, y),$$

follows from  $(1 \wedge a)(1 \wedge b) \leq (1 \wedge ab)$  for  $a, b \in \mathbb{R}_+$ .

### Remark

*By construction of  $\tilde{P}$ ,*

$$\text{var}(f, P) \leq \text{var}(f, \tilde{P})$$

*for any  $f \in L^2(X, \pi)$ , using Peskun ordering (Peskun, 1973, Tierney, 1998), since  $\tilde{\alpha}(x, y) \leq \alpha(x, y)$  for any  $(x, y) \in X^2$ .*

## Condition (1)

Defining  $A := \{(x, y) \in X^2 : r(x, y) \geq 1\}$ , there exists  $c$  such that

$$\inf_{(x,y) \in A} \min_{k \in \{1, \dots, d\}} \rho_k(x, y) \geq c.$$

Ensures that when

$$\alpha(x, y) = 1$$

then acceptance probability  $\tilde{\alpha}(x, y)$  uniformly lower-bounded by positive constant.

Reversibility implies  $\tilde{\alpha}(x, y)$  uniformly lower-bounded by a constant multiple of  $\alpha(x, y)$  for all  $x, y \in X$ .

## Condition (1)

Defining  $A := \{(x, y) \in X^2 : r(x, y) \geq 1\}$ , there exists  $c$  such that

$$\inf_{(x,y) \in A} \min_{k \in \{1, \dots, d\}} \rho_k(x, y) \geq c.$$

## Proposition (1)

*Under Condition (1), Lemma 34 in Andrieu et al. (2013) implies*

$$\text{Gap}(\tilde{P}) \geq \varrho \text{Gap}(P) \text{ and}$$

$$\text{var}(f, \tilde{P}) \leq (\varrho^{-1} - 1) \text{var}_\pi(f) + \varrho^{-1} \text{var}(f, P)$$

*with  $f \in L_0^2(E, \pi)$ ,  $\varrho = c^{d-1}$ .*

## Proposition (1)

*Under Condition (1), Lemma 34 in Andrieu et al. (2013) implies*

$$\text{Gap}(\tilde{P}) \geq \varrho \text{Gap}(P) \text{ and}$$

$$\text{var}(f, \tilde{P}) \leq (\varrho^{-1} - 1) \text{var}_{\pi}(f) + \varrho^{-1} \text{var}(f, P)$$

*with  $f \in L_0^2(\mathbb{E}, \pi)$ ,  $\varrho = c^{d-1}$ .*

Hence if  $P$  has right spectral gap, then so does  $\tilde{P}$ .

Plus, quantitative bounds on asymptotic variance of MCMC estimates using  $(\tilde{X}_n)_{n \geq 1}$  in relation to those using  $(X_n)_{n \geq 1}$  available

Easiest use of above: modify any candidate factorisation

Given factorisation of  $r$

$$r(x, y) = \prod_{k=1}^d \tilde{\rho}_k(x, y),$$

satisfying the balance condition, define a sequence of functions  $\rho_k$  such that both  $r(x, y) = \prod_{k=1}^d \rho_k(x, y)$  and Condition 1 holds.

Take  $c \in (0, 1]$ , define  $b = c^{\frac{1}{d-1}}$  and set

$$\tilde{\rho}_k(x, y) := \min \left\{ \frac{1}{b}, \max \{b, \rho_k(x, y)\} \right\}, \quad k \in \{1, \dots, d-1\},$$

and

$$\tilde{\rho}_d(x, y) := \frac{r(x, y)}{\prod_{k=1}^{d-1} \tilde{\rho}_k(x, y)}.$$

Then:

### Proposition (2)

*Under this scheme, previous proposition holds with*

$$\varrho = c^2 = b^{2(d-1)}$$

► optimising decomposition For any  $f \in L^2(E, \mu)$  define Dirichlet form associated with a  $\mu$ -reversible Markov kernel  $\Pi : E \times \mathcal{B}(E)$  as

$$\mathcal{E}_\Pi(f) := \frac{1}{2} \int \mu(dx) \Pi(x, dy) [f(x) - f(y)]^2.$$

The (right) spectral gap of a generic  $\mu$ -reversible Markov kernel has the following variational representation

$$\text{Gap}(\Pi) := \inf_{f \in L_0^2(E, \mu)} \frac{\mathcal{E}_\Pi(f)}{\langle f, f \rangle_\mu}.$$

► optimising decomposition

Lemma (Andrieu et al., 2013, Lemma 34)

*Let  $\Pi_1$  and  $\Pi_2$  be  $\mu$ -reversible Markov transition kernels of  $\mu$ -irreducible and aperiodic Markov chains, and assume that there exists  $\varrho > 0$  such that for any  $f \in L_0^2(\mathbb{E}, \mu)$*

$$\mathcal{E}_{\Pi_2}(f) \geq \varrho \mathcal{E}_{\Pi_1}(f) \quad ,$$

*then*

$$\text{Gap}(\Pi_2) \geq \varrho \text{Gap}(\Pi_1)$$

*and*

$$\text{var}(f, \Pi_2) \leq (\varrho^{-1} - 1) \text{var}_{\mu}(f) + \varrho^{-1} \text{var}(f, \Pi_1) \quad f \in L_0^2(\mathbb{E}, \mu).$$

Explorative performances of random-walk MCMC strongly dependent on proposal distribution

Finding optimal scale parameter leads to efficient 'jumps' in state space and smaller...

- ① expected square jump distance (ESJD)
- ② overall acceptance rate ( $\alpha$ )
- ③ asymptotic variance of ergodic average  $\text{var}(f, K)$

[Roberts et al. (1997), Sherlock and Roberts (2009)]

Provides practitioners with 'auto-tune' version of resulting random-walk MCMC algorithm

Quest for optimisation focussing on two main cases:

- 1  $d \rightarrow \infty$ : Roberts et al. (1997) give conditions under which each marginal chain converges toward a Langevin diffusion  
Maximising speed of that diffusion implies minimisation of the ACT and also  $\tau$  free from the functional

Remember  $\text{var}(f, K) = \tau_f \times \text{var}_\pi(f)$  where

$$\tau_f = 1 + 2 \sum_{i=1}^{\infty} \text{Cor}(f(X_0), f(X_i))$$

Quest for optimisation focussing on two main cases:

- 2 finite  $d$ : Sherlock and Roberts (2009) consider unimodal elliptically symmetric targets and show proxy for ACT is Expected Square Jumping Distance (ESJD), defined as

$$\mathbb{E} [\|X' - X\|_{\beta}^2] = \mathbb{E} \left[ \sum_{i=1}^d \beta_i^{-2} (X'_i - X)^2 \right]$$

As  $d \rightarrow \infty$ , ESJD converges to the speed of the diffusion process described in Roberts et al. (1997) [close asymptotia:  $d \gtrsim 5$ ]

*“For a moment, nothing happened. Then, after a second or so, nothing continued to happen.” —  
D. Adams, THGG*

When considering efficiency for Delayed Acceptance, focus on execution time as well

$$\mathbf{Eff} := \text{ESJD} / \text{cost per iteration}$$

similar to Sherlock et al. (2013) for pseudo-Marginal MCMC

Set of assumptions:

**(H1)** Assume [for simplicity's sake] that Delayed Acceptance operates on two factors only, i.e.,

$$r(x, y) = \rho_1(x, y) \times \rho_2(x, y),$$
$$\tilde{\alpha}(x, y) = \prod_{i=1}^2 (1 \wedge \rho_i(x, y))$$

Restriction also considers *ideal* setting where a computationally cheap approximation  $\tilde{f}(\cdot)$  is available and precise enough so that

$$\rho_2(x, y) = r(x, y) / \rho_1(x, y) = \pi(y) / \pi(x) \times \tilde{f}(x) / \tilde{f}(y) = 1$$

**(H2)** Assume that target distribution satisfies (A1) and (A2) in Roberts et al. (1997), which are regularity conditions on  $\pi$  and its first and second derivatives, and that

$$\pi(x) = \prod_{i=1}^n f(x_i)$$

.

**(H3)** Consider only a random walk proposal

$$y = x + \sqrt{\ell^2/d} Z$$

where

$$Z \sim \mathcal{N}(0, I_d)$$

**(H4)** Assume that cost of computing  $\tilde{f}(\cdot)$ ,  $c$  say, proportional to cost of computing  $\pi(\cdot)$ ,  $C$  say, with  $c = \delta C$ .

Normalising by  $C = 1$ , average total cost per iteration of DA chain is

$$\delta + \mathbb{E}[\tilde{\alpha}]$$

and efficiency of proposed method under above conditions is

$$\mathbf{Eff}(\delta, \ell) = \frac{ESJD}{\delta + \mathbb{E}[\tilde{\alpha}]}$$

## Lemma

Under conditions **(H1)–(H4)** on  $\pi(\cdot)$ ,  $q(\cdot, \cdot)$  and on  $\tilde{\alpha}(\cdot, \cdot) = (1 \wedge \rho_1(x, y))$   
As  $d \rightarrow \infty$

$$\text{Eff}(\delta, \ell) \approx \frac{h(\ell)}{\delta + \mathbb{E}[\tilde{\alpha}]} = \frac{2\ell^2\Phi(-\ell\sqrt{I}/2)}{\delta + 2\Phi(-\ell\sqrt{I}/2)}$$

$$a(\ell) \approx \mathbb{E}[\tilde{\alpha}] = 2\Phi(-\ell\sqrt{I}/2)$$

where  $I := \mathbb{E} \left[ \left( \frac{(\pi(x))'}{\pi(x)} \right)^2 \right]$  as in Roberts et al. (1997).

## Proposition (3)

*Under conditions of Lemma 3, optimal average acceptance rate  $\alpha^*(\delta)$  is independent of  $l$ .*

## Proof.

Consider  $\mathbf{Eff}(\delta, \ell)$  in terms of  $(\delta, a(\ell))$ :

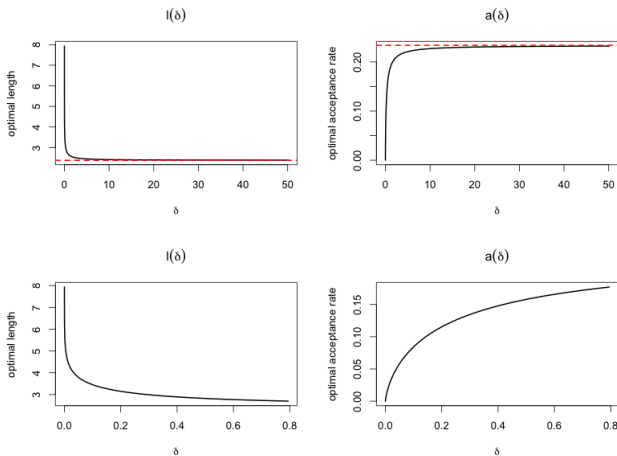
$$a = g(\ell) = 2\Phi\left(-\ell\sqrt{l}/2\right) \quad \ell = g^{-1}(a) = -\Phi^{-1}(a/2) \frac{2}{\sqrt{l}}$$

$$\mathbf{Eff}(\delta, a) = \frac{\frac{4}{l} \left[ \Phi^{-1}\left(\frac{a}{2}\right)^2 a \right]}{\delta + a} = \frac{4}{l} \left\{ \frac{1}{\delta + a} \left[ \Phi^{-1}\left(\frac{a}{2}\right)^2 a \right] \right\}$$



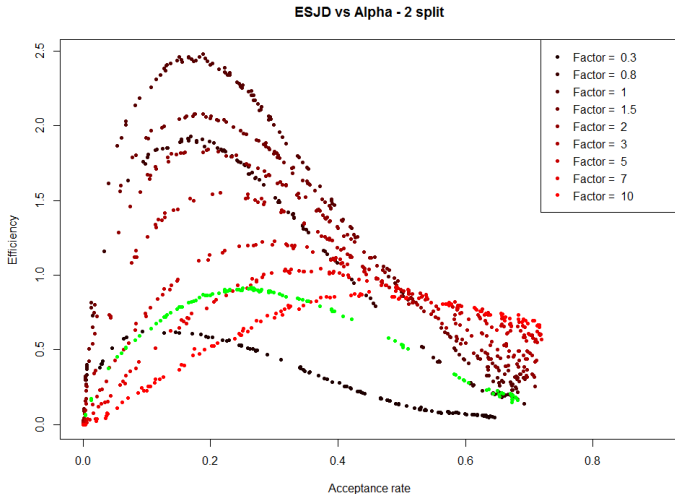
# Delayed Acceptance optimisation

**Figure :** *top panels:*  $\ell^*(\delta)$  and  $\alpha^*(\delta)$  as relative cost varies. For  $\delta \gg 1$  the optimal values converges to values computed for the standard M-H (red, dashed). *bottom panels:* close-up of interesting region for  $0 < \delta < 1$ .



# Robustness wrt (H1)–(H4)

Figure : Efficiency of various DA wrt acceptance rate of the chain; colours represent scaling factor in variance of  $\rho_1$  wrt  $\pi$ ;  $\delta = 0.3$ .



If computing cost comparable for all terms in

$$(x, y) = \prod_{i=1}^K \xi_i(x, y)$$

- rank entries according to the success rates observed on preliminary run
- start with ratios with highest variances
- rank factors by correlation with full Metropolis–Hastings ratio

**Logistic regression:**

- $10^6$  simulated observations with a 100-dimensional parameter space
- optimised Metropolis–Hastings with  $\alpha = 0.234$
- DA optimised via empirical correlation
  - split the data into subsamples of 10 elements
  - include smallest number of subsamples to achieve 0.85 correlation
  - optimise  $\Sigma$  against acceptance rate

algo	ESS (av.)	ESJD (av.)
DA-MH over MH	5.47	56.18

## geometric MALA:

- proposal

$$\theta' = \theta^{(i-1)} + \varepsilon^2 A^T A \nabla_{\theta} \log(\pi(\theta^{(i-1)}|y))/2 + \varepsilon A v$$

with position specific  $A$

[Girolami and Calderhead (2011), Roberts and Stramer (2002)]

- computational bottleneck in computing 3<sup>rd</sup> derivative of  $\pi$  in proposal
- G-MALA variance set to  $\sigma_d^2 = \frac{\ell^2}{d^{1/3}}$
- $10^2$  simulated observations with a 10-dimensional parameter space
- DA optimised via acceptance rate

$$\text{Eff}(\delta, a) = - (2/K)^{2/3} \frac{a\Phi^{-1}(a/2)^{2/3}}{\delta + a(1-\delta)}.$$

algo	accept	ESS/time (av.)	ESJD/time (av.)
MALA	0.661	0.04	0.03
DA-MALA	0.09	0.35	0.31

## Jeffreys for mixtures:

- numerical integration for each term in Fisher information matrix
- split between likelihood (cheap) and prior (expensive) unstable
- saving 5% of the sample for second step
- MH and DA optimised via acceptance rate
- actual averaged gain ( $\frac{ESS_{DA}/ESS_{MH}}{time_{DA}/time_{MH}}$ ) of 9.58

Algorithm	ESS (aver.)	ESJD (aver.)	time (aver.)
MH	1575.963	0.226	513.95
MH + DA	628.767	0.215	42.22