

DOWN-UP METROPOLIS-HASTINGS ALGORITHM FOR MULTIMODALITY

Hyungsuk Tak

Stat310

24 Nov 2015

Joint work with Xiao-Li Meng and David A. van Dyk

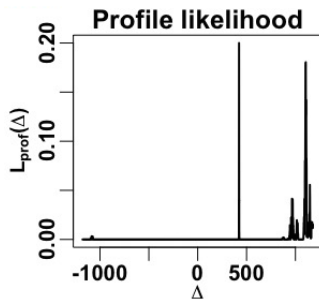
OUTLINE

- ▶ Motivation & idea
- ▶ Down-Up Metropolis-Hastings (DUMH) algorithm
 - ▶ Mathematical specification
 - ▶ Auxiliary variable approach
 - ▶ Algorithmic specification
- ▶ Examples
 - ▶ A mixture of 20 bivariate Gaussian distributions (Kou et al., 2006)
 - ▶ Quasar Q0957+561
- ▶ Discussion

MOTIVATION

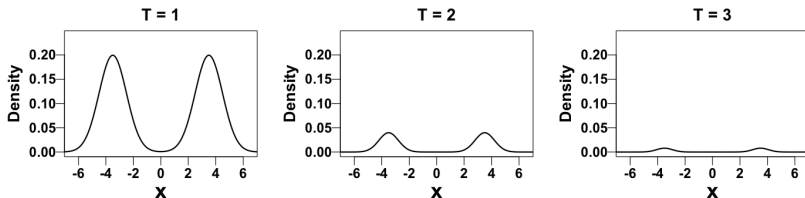
Time delay estimation between gravitationally-lensed quasar images to calculate the current expansion rate of the Universe.

- ▶ Full posterior density function: $\pi(\Delta, \beta, \theta | \text{Data})$.
- ▶ Profile likelihood: $L_{\text{prof}}(\Delta) \propto \pi(\Delta | \text{Data})$.
- ▶ A multimodal posterior distribution of the time delay (Δ in days) for Quasar Q0957+561



MOTIVATION (CONT.)

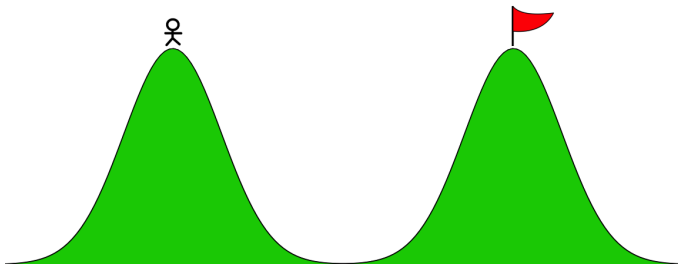
Popular remedy for multimodality is **tempering**, i.e., $\pi(x)^T$.



- ▶ Parallel tempering (Geyer, 1991), simulated tempering (Marinari & Parisi, 1992), tempered transition (Neal, 1996), Equi-energy sampler (Kou et al., 2006), etc.
- ▶ Hard to tune several **temperature-related tuning parameters** off-line
- ▶ **Computationally expensive**
- ▶ Is there a **simple** and **fast** way to expedite jumps between modes?

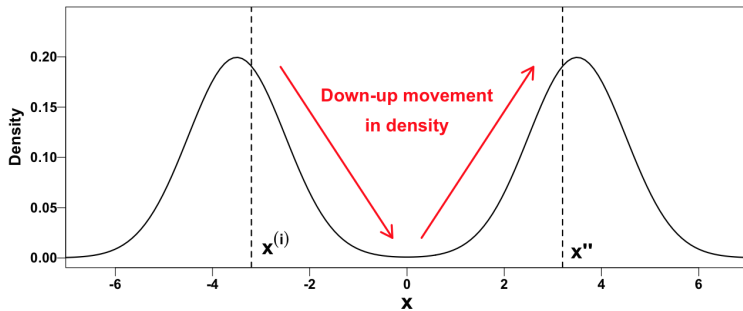
IDEA

How would you go to the top of the other mountain?



IDEA (CONT.)

Let's make a **down-up movement in density** to generate a proposal \mathbf{x}''



A DOWN-UP METROPOLIS-HASTINGS

Two-step Metropolis transitions: $x^{(i)}$: Current state

→ x' : Intermediate proposal

→ x'' : Final proposal

- ▶ How do we make a **downhill** Metropolis transition from $x^{(i)}$ to x' ?

$$\alpha^D(x' | x^{(i)}) = \min \left\{ 1, \frac{\pi(x^{(i)})}{\pi(x')} \right\}$$

- ▶ How do we make an **uphill** Metropolis transition from x' to x'' ?

$$\alpha^U(x'' | x') = \min \left\{ 1, \frac{\pi(x'')}{\pi(x')} \right\}$$

- ▶ For computational stability,

$$\alpha_{\epsilon}^D(x' | x^{(i)}) = \min \left\{ 1, \frac{\pi(x^{(i)}) + \epsilon}{\pi(x') + \epsilon} \right\}, \quad \alpha_{\epsilon}^U(x'' | x') = \min \left\{ 1, \frac{\pi(x'') + \epsilon}{\pi(x') + \epsilon} \right\}$$

A DOWN-UP METROPOLIS-HASTINGS (CONT.)

Proposal density function of DUMH

$$q^{\text{DU}}(x'' | x^{(i)}) = \int q^{\text{D}}(x' | x^{(i)}) q^{\text{U}}(x'' | x') dx',$$

where the downhill kernel density is

$$\begin{aligned} q^{\text{D}}(x' | x^{(i)}) &= N(x' | x^{(i)}, \sigma^2) \alpha_{\epsilon}^{\text{D}}(x' | x^{(i)}) + r^{\text{D}}(x^{(i)}) \delta_{x^{(i)}}(x'), \\ r^{\text{D}}(x^{(i)}) &= 1 - \int N(x' | x^{(i)}, \sigma^2) \alpha_{\epsilon}^{\text{D}}(x' | x^{(i)}) dx', \end{aligned}$$

and similarly the uphill kernel density is

$$\begin{aligned} q^{\text{U}}(x'' | x') &= N(x'' | x', \sigma^2) \alpha_{\epsilon}^{\text{U}}(x'' | x') + r^{\text{U}}(x') \delta_{x'}(x''), \\ r^{\text{U}}(x') &= 1 - \int N(x'' | x', \sigma^2) \alpha_{\epsilon}^{\text{U}}(x'' | x') dx'', \end{aligned}$$

A DOWN-UP METROPOLIS-HASTINGS (CONT.)

Four possibilities of the final proposal:

- ▶ Down-Up: Both accepted
- ▶ Down- : Downhill accepted & uphill rejected (not helpful for jump)
- ▶ -Up: Downhill rejected & uphill accepted (not helpful for jump)
- ▶ : Both rejected ($x^{(i)} = x''$) (waste of iterations)

We only need “Down-up” movement to boost jumps between modes!

A DOWN-UP METROPOLIS-HASTINGS (CONT.)

Exclude the last three possibilities via **forced** Metropolis transitions.

$$q^D(x' | x^{(i)}) = N(x' | x^{(i)}, \sigma^2)\alpha^D(x' | x^{(i)}) + r^D(x^{(i)})\delta_{x^{(i)}}(x'),$$

$$r^D(x^{(i)}) = 1 - \int N(x' | x^{(i)}, \sigma^2)\alpha^D(x' | x^{(i)})dx',$$

$$\rightarrow q^D(x' | x^{(i)}) = \frac{N(x' | x^{(i)}, \sigma^2)\alpha^D(x' | x^{(i)})}{1 - r^D(x^{(i)})}$$

- ▶ How do we generate $x' \sim q^D(x' | x^{(i)})$?

Repeatedly generate $x' \sim q^D(x' | x^{(i)})$ **until it is accepted!**

- ▶ **Intuition:** Flip a coin twice. Restrict the sample space to {HT, HH}. How do we generate {HT} or {HH}? Flip the coin twice until either {HT} or {HH} appears, ignoring {TT} and {TH}.

A DOWN-UP METROPOLIS-HASTINGS (CONT.)

Accept x'' with an MH acceptance probability

$$\begin{aligned}\alpha^{\text{DU}}(x'' | x^{(i)}) &= \min \left\{ 1, \frac{\pi(x'')q^{\text{DU}}(x^{(i)} | x'')}{\pi(x^{(i)})q^{\text{DU}}(x'' | x^{(i)})} \right\} \\ &= \min \left\{ 1, \frac{\pi(x'')(1 - r^{\text{D}}(x^{(i)}))}{\pi(x^{(i)})(1 - r^{\text{D}}(x''))} \right\}.\end{aligned}$$

The second equality holds because

$$\alpha_{\epsilon}^{\text{D}}(x' | x^{(i)}) = \min \left\{ 1, \frac{\pi(x^{(i)}) + \epsilon}{\pi(x') + \epsilon} \right\} = \alpha_{\epsilon}^{\text{U}}(x^{(i)} | x')$$

and thus

$$\begin{aligned}q^{\text{DU}}(x'' | x^{(i)}) &= \int \frac{N(x' | x^{(i)}, \sigma^2)\alpha_{\epsilon}^{\text{D}}(x' | x^{(i)})}{1 - r^{\text{D}}(x^{(i)})} \frac{N(x'' | x', \sigma^2)\alpha_{\epsilon}^{\text{U}}(x'' | x')}{1 - r^{\text{U}}(x')} dx' \\ q^{\text{DU}}(x'' | x^{(i)})\{1 - r^{\text{D}}(x^{(i)})\} &= \int N(x' | x^{(i)}, \sigma^2)\alpha_{\epsilon}^{\text{D}}(x' | x^{(i)}) \frac{N(x'' | x', \sigma^2)\alpha_{\epsilon}^{\text{U}}(x'' | x')}{1 - r^{\text{U}}(x')} dx' \\ &= \int N(x^{(i)} | x', \sigma^2)\alpha_{\epsilon}^{\text{U}}(x^{(i)} | x') \frac{N(x' | x'', \sigma^2)\alpha_{\epsilon}^{\text{D}}(x' | x'')}{1 - r^{\text{U}}(x')} dx' \\ &= \int N(x' | x'', \sigma^2)\alpha_{\epsilon}^{\text{D}}(x' | x'') \frac{N(x^{(i)} | x', \sigma^2)\alpha_{\epsilon}^{\text{U}}(x^{(i)} | x')}{1 - r^{\text{U}}(x')} dx' \\ &= q^{\text{DU}}(x^{(i)} | x'')\{1 - r^{\text{D}}(x'')\}.\end{aligned}$$

A DOWN-UP METROPOLIS-HASTINGS (CONT.)

Useless due to the intractable integrations in the acceptance probability

$$\begin{aligned}\alpha^{\text{DU}}(x'' | x^{(i)}) &= \min \left\{ 1, \frac{\pi(x'')(1 - r^{\text{D}}(x^{(i)}))}{\pi(x^{(i)})(1 - r^{\text{D}}(x''))} \right\} \\ &= \min \left\{ 1, \frac{\pi(x'') \int N(x' | x^{(i)}, \sigma^2) \alpha_{\epsilon}^{\text{D}}(x' | x^{(i)}) dx'}{\pi(x^{(i)}) \int N(x' | x'', \sigma^2) \alpha_{\epsilon}^{\text{U}}(x' | x'') dx'} \right\}.\end{aligned}$$

Is there any way to cancel out **this ratio**?

AUXILIARY VARIABLE APPROACH

If we explore a larger space, then there can be a way to cancel the ratio!
(Møller et al., 2006)

An **auxiliary variable** A such that $\pi^C(A | x)$ is well-defined.

- ▶ Joint target density: $\pi^J(A, x)$
- ▶ Joint proposal density:

$$q^J(A'', x'' | A^{(i)}, x^{(i)}) = q_1(x'' | A^{(i)}, x^{(i)}) q_2(A'' | x'', A^{(i)}, x^{(i)})$$

- ▶ Joint acceptance probability:

$$\begin{aligned} \alpha^J(A'', x'' | A^{(i)}, x^{(i)}) &= \min \left[1, \frac{\pi^J(A'', x'') q^J(A^{(i)}, x^{(i)} | A'', x'')}{\pi^J(A^{(i)}, x^{(i)}) q^J(A'', x'' | A^{(i)}, x^{(i)})} \right] \\ &= \min \left[1, \frac{\pi(x'') \pi^C(A'' | x'') q_1(x^{(i)} | A'', x'') q_2(A^{(i)} | x^{(i)}, A'', x'')}{\pi(x^{(i)}) \pi^C(A^{(i)} | x^{(i)}) q_1(x'' | A^{(i)}, x^{(i)}) q_2(A'' | x'', A^{(i)}, x^{(i)})} \right] \end{aligned}$$

AUXILIARY VARIABLE APPROACH (CONT.)

Let's choose π^C , q_1 , and q_2 to cancel out the intractable ratio.

$$\begin{aligned}q_1(x'' \mid A^{(i)}, x^{(i)}) &= q^{\text{DU}}(x'' \mid x^{(i)}) \\q_2(A'' \mid x'', A^{(i)}, x^{(i)}) &= q^{\text{D}}(A'' \mid x'') \\ \pi^C(A'' \mid x'') &= N(A'' \mid x'', \sigma^2)\end{aligned}$$

Then,

$$\begin{aligned}\alpha^J(A'', x'' \mid A^{(i)}, x^{(i)}) &= \min \left[1, \frac{\pi(x'') N(A'' \mid x'', \sigma^2) q^{\text{DU}}(x^{(i)} \mid x'') q^{\text{D}}(A^{(i)} \mid x^{(i)})}{\pi(x^{(i)}) N(A^{(i)} \mid x^{(i)}, \sigma^2) q^{\text{DU}}(x'' \mid x^{(i)}) q^{\text{D}}(A'' \mid x'')} \right] \\ &= \min \left[1, \frac{\pi(x'') N(A'' \mid x'', \sigma^2) \{1 - r^{\text{D}}(x^{(i)})\} \frac{N(A^{(i)} \mid x^{(i)}, \sigma^2) \alpha_\epsilon^{\text{D}}(A^{(i)} \mid x^{(i)})}{1 - r^{\text{D}}(x^{(i)})}}{\pi(x^{(i)}) N(A^{(i)} \mid x^{(i)}, \sigma^2) \{1 - r^{\text{D}}(x'')\} \frac{N(A'' \mid x'', \sigma^2) \alpha_\epsilon^{\text{D}}(A'' \mid x'')}{1 - r^{\text{D}}(x'')}} \right] \\ &= \min \left[1, \frac{\pi(x'') \alpha_\epsilon^{\text{D}}(A^{(i)} \mid x^{(i)})}{\pi(x^{(i)}) \alpha_\epsilon^{\text{D}}(A'' \mid x'')} \right] = \min \left[1, \frac{\pi(x'') \min\{1, \frac{\pi(x^{(i)}) + \epsilon}{\pi(A^{(i)}) + \epsilon}\}}{\pi(x^{(i)}) \min\{1, \frac{\pi(x'') + \epsilon}{\pi(A'') + \epsilon}\}} \right].\end{aligned}$$

ALGORITHMIC SPECIFICATION

The DUMH algorithm is composed of **eight steps for each iteration**.

Set initial values $A^{(0)}$ and $x^{(0)}$. For $i = 0, 1, \dots, n - 1$,

Step 1: Sample $x' \sim N(x' | x^{(i)}, \sigma^2)$ and $u_1 \sim \text{Unif}(0, 1)$.

Step 2: Repeat Step 1 until $u_1 < \alpha_\epsilon^D(x' | x^{(i)})$ for **forced downhill move**.

Step 3: Sample $x'' \sim N(x'' | x', \sigma^2)$ and $u_2 \sim \text{Unif}(0, 1)$.

Step 4: Repeat Step 3 until $u_2 < \alpha_\epsilon^U(x'' | x')$ for **forced uphill move**.

Step 5: Sample $A'' \sim N(A'' | x'', \sigma^2)$ and $u_3 \sim \text{Unif}(0, 1)$.

Step 6: Repeat Step 5 until $u_3 < \alpha_\epsilon^D(A'' | x'')$ for **forced downhill move**.

Step 7: Sample $u_4 \sim \text{Unif}(0, 1)$.

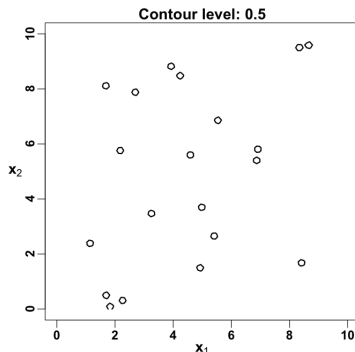
Step 8: Accept (A'', x'') as $(A^{(i+1)}, x^{(i+1)})$ if $u_4 < \alpha^J(A'', x'' | A^{(i)}, x^{(i)})$,
or otherwise set $(A^{(i+1)}, x^{(i+1)})$ to $(A^{(i)}, x^{(i)})$.

EXAMPLE 1

A mixture of 20 bivariate Gaussian distributions (Kou et al., 2006)

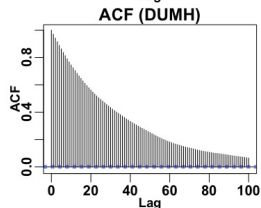
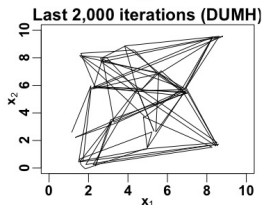
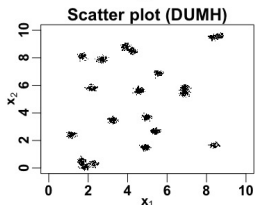
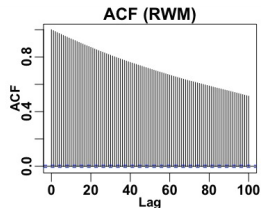
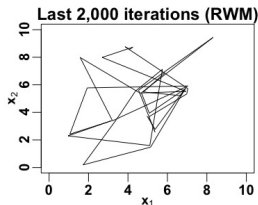
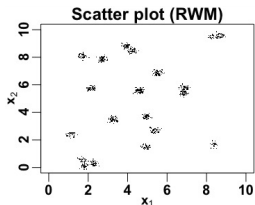
$$\pi(x) \propto \frac{1}{20} \sum_{i=1}^{20} \frac{1}{2\pi\sigma_i^2} \exp\left(-\frac{1}{2\sigma_i^2}(x - \mu_i)^\top(x - \mu_i)\right),$$

where $x = (x_1, x_2)^\top$, $\sigma_i = 0.1$, and the 20 mean vectors, $(\mu_1, \mu_2, \dots, \mu_{20})^\top$, are given in Kou et al. (2006).



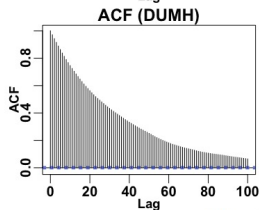
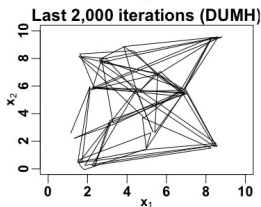
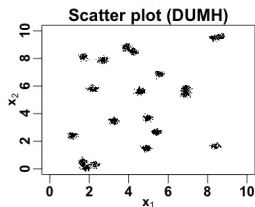
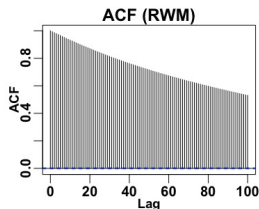
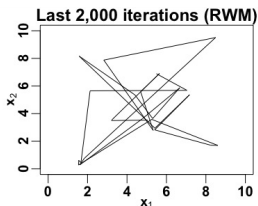
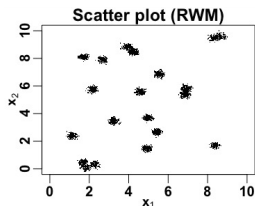
EXAMPLE 1 (CONT.): RWM vs DUMH

50,000 iterations after discarding the first 50,000 iterations.



	σ	Accept. Rate	Sample size	CPU time (sec.)
RWM	4	0.015	50,000	535
DUMH	4	0.045	50,000	2,758

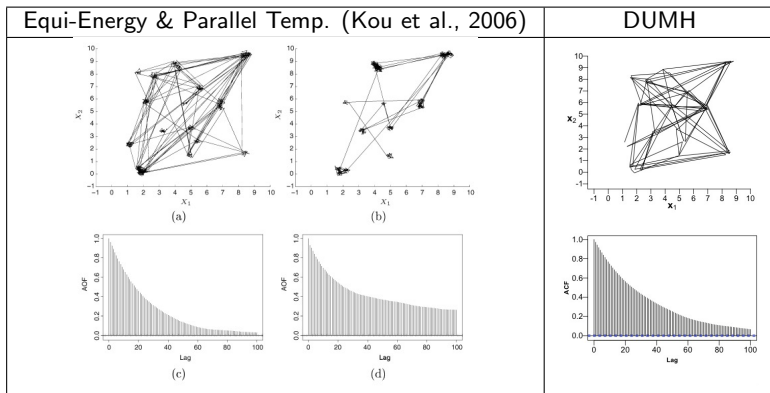
EXAMPLE 1 (CONT.): RWM vs DUMH



	σ	Accept. Rate	Sample size	CPU time (sec.)
RWM	4	0.013	257,757	2,663
DUMH	4	0.046	50,000	2,758

"Do (DU) MH instead of the naive RWM!"

EXAMPLE 1 (CONT.): EE vs PT vs DUMH



	$E(X_1)$	$E(X_2)$	$E(X_1^2)$	$E(X_2^2)$
True value	4.478	4.905	25.605	33.920
DUMH	4.5014 (0.095)	4.8847 (0.141)	25.5858 (0.977)	33.6423 (1.371)
EE	4.5019 (0.107)	4.9439 (0.139)	25.9241 (1.098)	34.4763 (1.373)
PT	4.4185 (0.170)	4.8790 (0.283)	24.9856 (1.713)	33.5966 (2.867)
$\frac{\text{MSE(EE)}}{\text{MSE(DUMH)}}$	1.26	1.03	1.37	1.12
$\frac{\text{MSE(PT)}}{\text{MSE(DUMH)}}$	3.39	3.98	3.47	4.25

Better than the EE and PT in terms of accuracy! (Do MH!)

EXAMPLE 2: QUASAR Q0957+561

MH within Gibbs sampler for $p(\mathbf{X}(\mathbf{t}^\Delta), \Delta, \beta, \theta \mid \text{Data})$.

Step 1: Sample $(\mathbf{X}^{(l)}(\mathbf{t}^{\Delta^{(l)}}), \Delta^{(l)}) \sim p(\mathbf{X}(\mathbf{t}^\Delta), \Delta \mid \beta^{(l-1)}, \theta^{(l-1)})$
 $= p(\mathbf{X}(\mathbf{t}^\Delta) \mid \Delta, \beta^{(l-1)}, \theta^{(l-1)}) \times p(\Delta \mid \beta^{(l-1)}, \theta^{(l-1)})$ by M-H

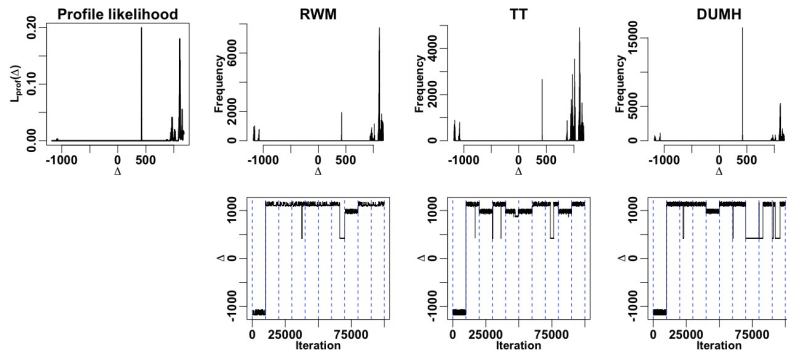
Step 2: Sample $\beta^{(l)} \sim p(\beta \mid \theta^{(l-1)}, \mathbf{X}^{(l)}(\mathbf{t}^{\Delta^{(l)}}), \Delta^{(l)})$

Step 3: Sample $\theta^{(l)} \sim p(\theta \mid \mathbf{X}^{(l)}(\mathbf{t}^{\Delta^{(l)}}), \Delta^{(l)}, \beta^{(l)})$

RWM, tempered transition (Neal, 1996), and DUMH applied to $p(\Delta \mid \beta^{(l-1)}, \theta^{(l-1)})$.

EXAMPLE 2: QUASAR Q0957+561

Each chain (among 10 chains) has 10,000 samples after discarding the first 5,000.



	σ	Accept. Rate	Sample size	CPU time/chain
RWM	600	0.018	100,000	48
TT	300	0.041	100,000	903
DUMH	1100	0.053	100,000	493

DISCUSSION

How do we **fairly** evaluate various MCMC samplers for multimodality?

- ▶ How do we measure **the time for the off-line tuning work** involved in all the temperature-based samplers?
- ▶ How do we **theoretically** evaluate the MCMC samplers?