

PARRELLEL TEMPERING

THEORY AND APPLICATIONS

Mateusz Łącki

Uniwersytet Warszawski

24 April 2013



TODAY'S AGENDA

- 1 WHAT? WHY? WHO?
- 2 METROPOLIS-HASTING ALGORITHM
- 3 METROPOLIS-HASTING LOCALNESS
 - Notation
 - Ideology
 - random walk kernel
 - swap kernel
- 4 THEORY OF PARALLEL TEMPERING
- 5 SWAPPING STRATEGIES
- 6 BIBLIOGRAPHY

WHAT IS PARALLEL TEMPERING? NUTSHELL VIEW.

- Stochastic simulation algorithm
- a.k.a. replica exchange Monte Carlo
 - sampling method
- Extension to Metropolis-Hastings algorithm ...
- ...or rather Metropolis-Hastings-Green algorithm
 - more abstract version: general kernels
 - more freedom: discrete kernels, continuous kernels, different dimensions



WHY PARALLEL TEMPERING?

- Tool to fight *multimodality*
 - "allows good mixing with multimodal target distributions"
 - response to Metropolis-Hastings localness
 - better estimation of integrals $\int_{\mathbb{R}^d} f(x)\pi(x)dx$
- In certain physical models: thermodynamic interpretation
 - Gibbs random-field model

WHO MIGHT BE INTERESTED IN PARALLEL TEMPERING?

- anyone having problems with locality of simulations
- researchers facing model selection problems while exploring the posteriori distribution over some space of models: the g-priors

THE USUAL METROPOLIS-HASTINGS ALGORITHM REVISED

We are given a measure π with density $h : \mathbb{R}^d \mapsto \mathbb{R}_+$.

As. h need not be normalised

$$\int_{\mathbb{R}^d} h(x) dx \in (0, \infty)$$

As. We can evaluate $h(x)$ for all $x \in \mathbb{R}^d$

As. Transitional probability density $q : \underbrace{\mathbb{R}^d}_{\text{currentstate}} \times \underbrace{\mathbb{R}^d}_{\text{proposal}} \mapsto \mathbb{R}_+$

THE USUAL METROPOLIS-HASTINGS ALGORITHM REVISED

As. For all x , $q(x, \circ)$ is a normalised probability density

- \forall_x we can simulate $y \sim q(x, \circ)$
- $\forall_{x,y}$ we can evaluate $q(x, y)$

then the Markov chain $X \equiv \{X^{[k]}\}_{k \geq 0}$ generated by procedure

- 1 $x = X^{[k-1]}$
- 2 $y \sim q(x, \circ)$
- 3 evaluate $R(x, y) = \frac{h(y)q(y, x)}{h(x)q(x, y)}$
- 4 reject y with probability $\alpha(x, y) = 1 \wedge R(x, y)$
 - If rejected $X^{[k]} = x$
 - Otherwise $X^{[k]} = y$

is reversible: its kernel preserves π .

WHAT'S A KERNEL?

A regular version of $\mathbb{E}(\mathbb{I}_A | X = x)$

- measurable function with x for A fixed
- probability distribution with A for any fixed x (stronger than almost everywhere)

in standard MH

$$P(x, A) = \int_A p(x, y) dy + \left(1 - \int_{\mathbb{R}^d} p(x, y) dy\right) \mathcal{I}(x, A)$$

where $p(x, y) = \alpha(x, y)q(x, y)$

and $\mathcal{I}(x, A) = \mathbb{I}_A(x)$ - identity kernel

WHAT'S REVERSIBILITY?

An integral equation

$$\int_A \pi(\mathrm{d}x) P(x, B) = \int_B \pi(\mathrm{d}x) P(x, A)$$

This assures that the chain preserves π

$$\int_{\mathbb{R}^d} \pi(\mathrm{d}x) P(x, B) = \int_B \pi(\mathrm{d}x) = \pi(B)$$

or $\pi P = \pi$

THAT'S FUNCTIONAL ANALYSIS

Reversibility is P self-adjointness:

$$\forall_{f,g \in \mathbb{L}^2(\pi)} \int \pi(dx) P(x, dy) f(x) g(y) = \int \pi(dx) P(x, dy) g(x) f(y)$$

the inner product convention

$$\langle Pf | h \rangle = \langle f | Ph \rangle$$

standard spectral analysis tools applicable

METROPOLIS-HASTINGS-GREEN ALGORITHM

(S, \mathfrak{G}) - measurable space

As. Take a positive measure

$$\rho : \mathfrak{G} \mapsto \mathbb{R}_+$$

not necessarily probabilistic

$$\rho(S) \in (0, \infty)$$

As. $\rho = \rho(S)\pi$

proposal kernel $Q(x, A)$

As. we know how to generate $y \sim Q(x, \circ)$

As. For all $x, y \in S$ the Hastings ratio

$$R(x, y_k) \equiv \frac{\rho(dy) Q(y, dx)}{\rho(dx) Q(x, dy)}$$

known and possible to evaluate for any x and y

CHAIN GENERATION

Given $x = X^{[k-1]}$

- 1 Simulate $y \sim Q(x, \circ)$.
- 2 Calculate $R(x, y)$.
- 3 Accept y with probability $\alpha(x, y) = 1 \wedge R(x, y)$.

GREEN'S RESULT

Then by the Green theorem X 's kernel

$$P(x, A) \equiv \int_A \alpha(x, y) Q(x, dy) + \delta_x(A) \left(1 - \int_{\Omega} \alpha(x, y) Q(x, dy) \right)$$

is reversible with respect to π

Note: kernel is stochastic.

PROBLEMS WITH MH

To assert proposal's adequacy

generate $U \sim \mathcal{U}(0, 1)$

accept y if $U \leq \frac{h(y)q(y,x)}{h(x)q(x,y)} \wedge 1$

if $q(x, y) = q(y, x)$

$$U \leq \frac{h(y)}{h(x)} \wedge 1$$

Problem: if π is multimodal we get stuck in certain region of the state-space

LIANG AND WONG 2001

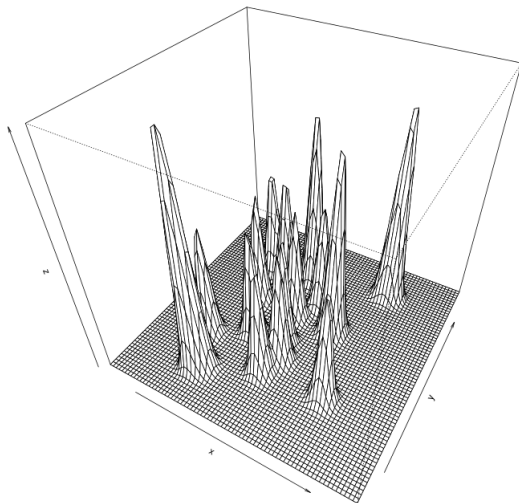
$$f(x) = \sum_{i=1}^{20} \frac{\omega_i}{\sigma_i \sqrt{2\pi}} \exp \left(- \frac{(x - \mu_i)'(x - \mu_i)}{2\sigma_i^2} \right)$$

where $\sigma_1 = \dots = \sigma_{20} = 0.1$, $\omega_1 = \dots = \omega_{20} = 0.05$

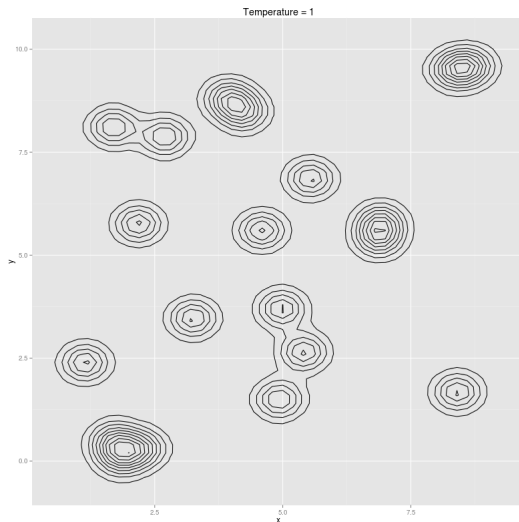
and the means μ_i are given by

1	2	3	4	5	6	7	8	9	10
2.18	8.67	4.24	8.41	3.93	3.25	1.70	4.59	6.91	6.87
5.76	9.59	8.48	1.68	8.82	3.47	0.50	5.60	5.81	5.40
11	12	13	14	15	16	17	18	19	20
5.41	2.70	4.98	1.14	8.33	4.93	1.83	2.26	5.54	1.69
2.65	7.88	3.70	2.39	9.50	1.50	0.09	0.31	6.86	8.11

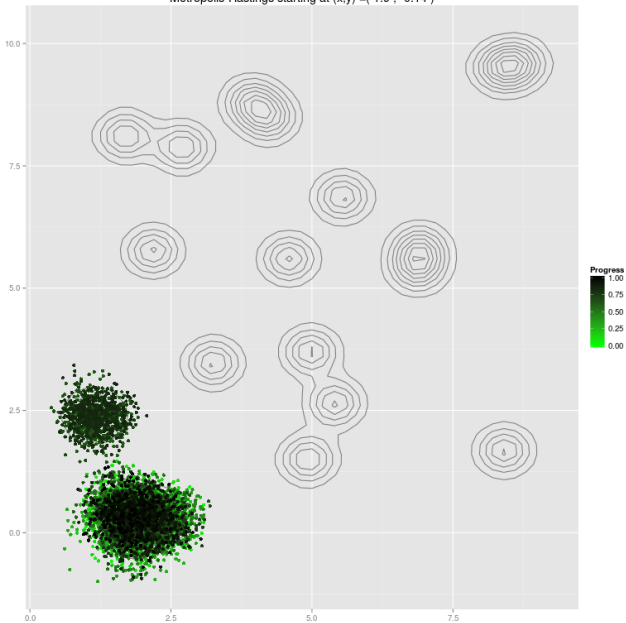
TOY-EXAMPLE VISUALISED

Temperature = 1

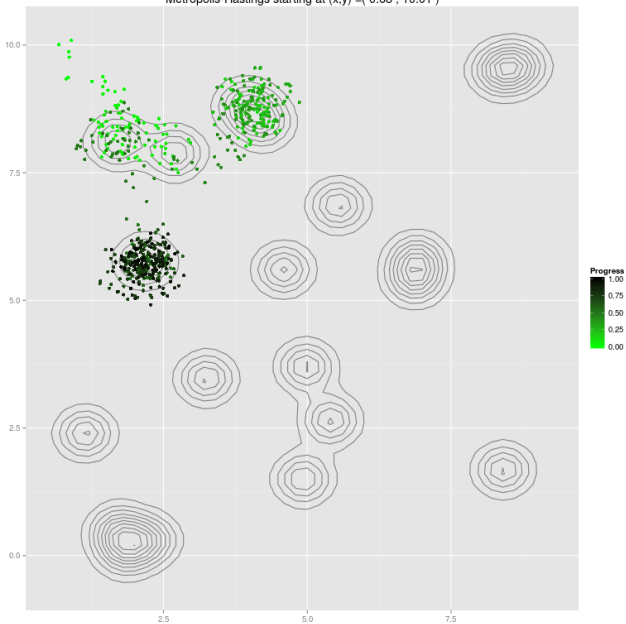
TOY-EXAMPLE VISUALISED AS A CONTOUR PLOT



Metropolis-Hastings starting at $(x,y) = (1.9, -0.14)$



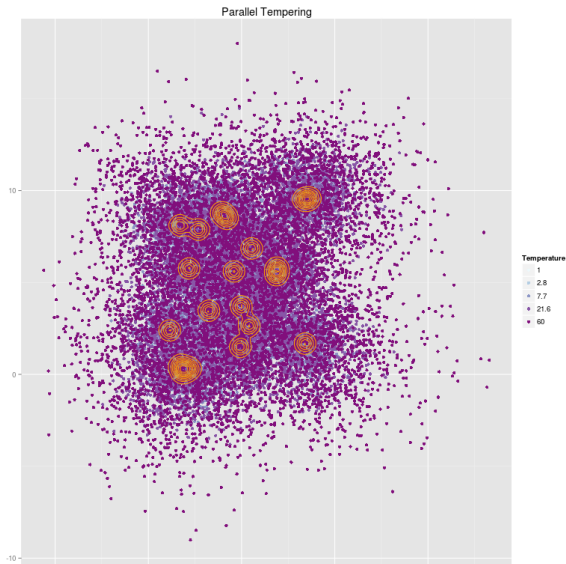
Metropolis-Hastings starting at (x,y)=(0.68 , 10.01)



Metropolis-Hastings starting at $(x,y) = (4.96, 1.33)$



PARALLEL TEMPERING IN ACTION.



GETTING INSIDE PARALLEL TEMPERING

(Ω, \mathfrak{F})	measurable space
Ω	subset of a Polish space
\mathfrak{F}	Borel subsets countably generated
$\mathcal{I} = [0, 1]$	unit interval
$\pi : \mathfrak{F} \mapsto \mathcal{I}$	measure



ASSUMPTIONS AND FURTHER NOTATION

As. π has density w.r. to Lebesgue measure

- π - the density
 - know up to its proportionality factor
 - unnormalised

$$\int_{\Omega} \pi(x) dx \in (0, \infty)$$

SOLUTION'S IDEOLOGY

Space for chains

$$(\Omega^L, \mathfrak{F}^{\otimes L}, \pi_\beta)$$

where $\mathfrak{F}^{\otimes L} \equiv \underbrace{\mathfrak{F} \otimes \cdots \otimes \mathfrak{F}}_{L \text{ times}}$

$$\pi_\beta \propto \pi^{\beta_1} \times \cdots \times \pi^{\beta_L}$$

$\beta = (\beta_1, \dots, \beta_L)$ - inverse temperatures $\beta_i = T_i^{-1}$

and $1 = T_1 < \cdots < T_L < \infty$

no normalisation of π_β coordinates

...ENTERS MARKOV

Markov Chain $X \equiv \{X^{[k]}\}_{k \geq 0}$

- Ω^L state-space for X

$$X^{[k]} = (X_1^{[k]}, \dots, X_L^{[k]})$$

High temperature



Low temperatures coordinates

NON-SWEDISH WAY TO EQUIDISTRIBUTION

As. $\pi > 0$ somewhere between modes

- then $\pi^{\beta_k} > \pi$ for $k \geq 2$

So that if $\pi(y) < \pi(x)$ then

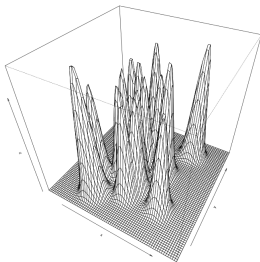
$$\alpha_{\beta_1}(x, y) = 1 \wedge \frac{\pi(y)}{\pi(x)} = \frac{\pi(y)}{\pi(x)} < \left(\frac{\pi(y)}{\pi(x)}\right)^{\beta_k} = 1 \wedge \left(\frac{\pi(y)}{\pi(x)}\right)^{\beta_k} = \alpha_{\beta_k}(x, y)$$

- proposal accepted more often in higher temperatures when

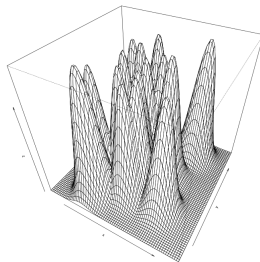
$$\frac{\pi(y)}{\pi(x)} < 1$$

- we enlarge regions from which proposals get accepted

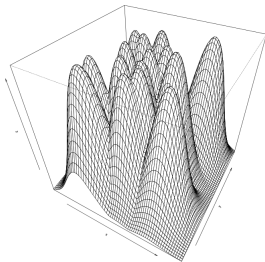
Temperature = 2.8



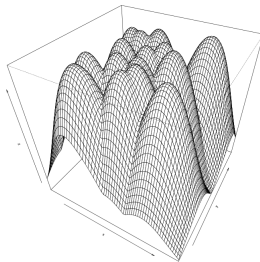
Temperature = 7.7

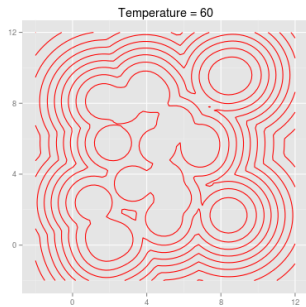
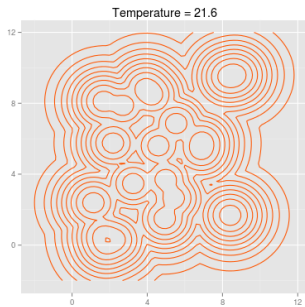
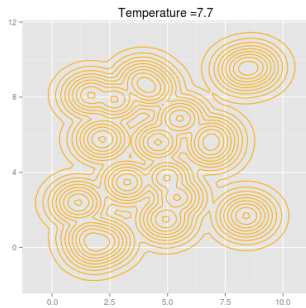
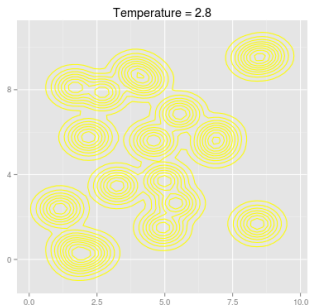


Temperature = 21.6

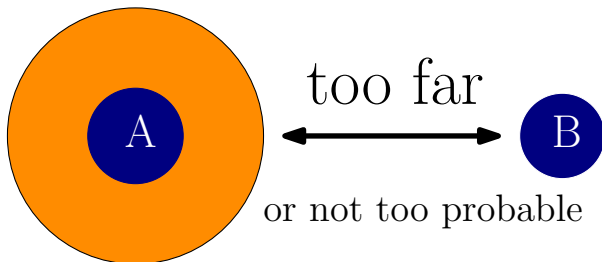


Temperature = 60





POTENTIAL PROBLEMS



POTENTIAL PROBLEMS

- π might be ok
- because of computer's finite arithmetic it is not ok anymore

THEORETICAL DETAILS

algorithm reached n^{th} step - $X^{[n]}$

we act with two kernels

$$X^{[n]} \xrightarrow{\mathcal{S}_\beta} \widetilde{X}^{[n+1]} \xrightarrow{M_{\Sigma,\beta}} X^{[n+1]}.$$

- \mathcal{S}_β - swap kernel
- $M_{\Sigma,\beta}$ - random walk kernel
- Their reversibility is assured by the MHG algorithm reversibility.
so we get π -preservation

$$\mathcal{S}_\beta M_{\Sigma,\beta} \pi = \mathcal{S}_\beta \pi = \pi$$

RANDOM WALK $M_{\Sigma, \beta}$

take $A_i \in \mathfrak{F}$ and $x \in \Omega^L$

then

$$M_{\Sigma, \beta}(x, A_1 \times \cdots \times A_L) = \prod_{l=1}^L M_{\Sigma_l, \beta_l}(x_l, A_l)$$

where $M_{\Sigma_l, \beta_l}(x_l, A_l)$ is equal to

$$\int_A \alpha_{\beta_l}(x_l, y_l) q_{\Sigma_l}(y_l - x_l) dy_l + \delta_x(A) \int [1 - \alpha_{\beta_l}(x_l, y_l)] q_{\Sigma_l}(y_l - x_l) dy_l$$

where α_{β_l} is the acceptance level
and q_{Σ_l} is proposal distribution

OBSERVATIONS AND ASSUMPTIONS

As. q_{Σ_I} - density of $\mathcal{N}(0, \Sigma_I)$

Symmetry $q(x_I, y_I) = q(y_I, x_I)$ implies

$$\alpha_{\beta_I}(x_I, y_I) \equiv 1 \wedge \frac{\pi^{\beta_I}(y_I)}{\pi^{\beta_I}(x_I)}$$

for the chain to preserve π_{β} .

Implementation: independent simulation of M_{Σ_I, β_I} for each $\tilde{X}_I^{[n-1]}$

INTERLACING INDEPENDENT CHAINS WITH \mathcal{S}_β

Swaps

- Less-tempered chains placed in unusual places
- a pair of coordinates $(\tilde{X}_i^{[n-1]}, \tilde{X}_k^{[n-1]})$ drawn at random

As. only one pair per turn

- different strategies possible

Introduce swap operation

$$\mathcal{S}_{ij}x = (x_1, \dots, x_{i-1}, x_j, x_{i+1}, \dots, x_{j-1}, x_i, x_{j+1}, \dots, x_L)$$

PRECISE KERNEL FORM OF \mathcal{S}_β

For any $x \in \Omega^L$ and $A \in \mathfrak{F}^{\otimes L}$

$$\mathcal{S}_\beta(x, A) \equiv$$

$$\sum_{i < j} p_{ij}(x) \alpha_{\text{swap}}(x, S_{ij}x) \mathbb{I}_A(S_{ij}x) + \left(1 - \sum_{i < j} p_{ij}(x) \alpha_{\text{swap}}(x, S_{ij}x)\right) \mathcal{I}(x, A)$$

where

$$\alpha_{\text{swap}}(x, S_{ij}x) = \left[\left(\frac{\pi(x_j)}{\pi(x_i)} \right)^{\beta_i - \beta_j} \frac{p_{ij}(S_{ij}x)}{p_{ij}(x)} \right] \wedge 1$$

is the acceptance level for swaps,

$p_{ij}(x)$ - probability function of a swap given x .

Nota bene: given state x , \mathcal{S}_β has a finite support

$$\mathfrak{S}_x \equiv \{S_{ij}x : i < j\}$$

DIFFERENT POSSIBLE SWAPPING STRATEGIES $p_{ij}(x)$

Strategy I

$$p_{ij}(x) \propto \frac{\pi(x_j)}{\pi(x_i)} \wedge \frac{\pi(x_i)}{\pi(x_j)} = \exp \left(- |\log(\pi(x_j)) - \log(\pi(x_i))| \right)$$

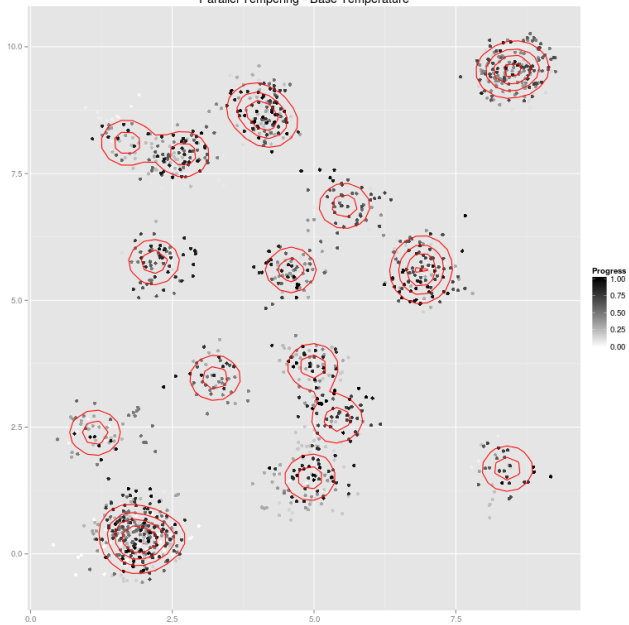
- promotes swaps between coordinates relatively the same, $\pi(x_j) \approx \pi(x_i)$

Strategy II

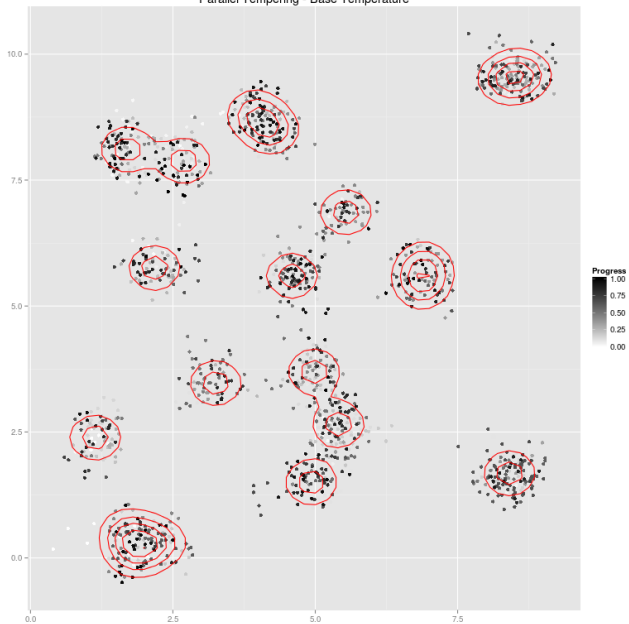
$$p_{ij}(x) \propto \frac{\pi(x_j)}{\pi(x_i)} \wedge 1 = \exp \left(- (\log(\pi(x_j)) - \log(\pi(x_i))) \right) \wedge 1$$

- breaks the symmetry of the previous one

Parallel Tempering - Base Temperature



Parallel Tempering - Base Temperature



DIFFERENT POSSIBLE SWAPPING STRATEGIES $p_{ij}(x)$

Strategy III

$$p_{ij} \propto \left(\frac{\pi(x_j)}{\pi(x_i)} \wedge \frac{\pi(x_i)}{\pi(x_j)} \right)^{\beta_i - \beta_j} = \exp \left(-(\beta_i - \beta_j) |\log(\pi(x_j)) - \log(\pi(x_i))| \right)$$

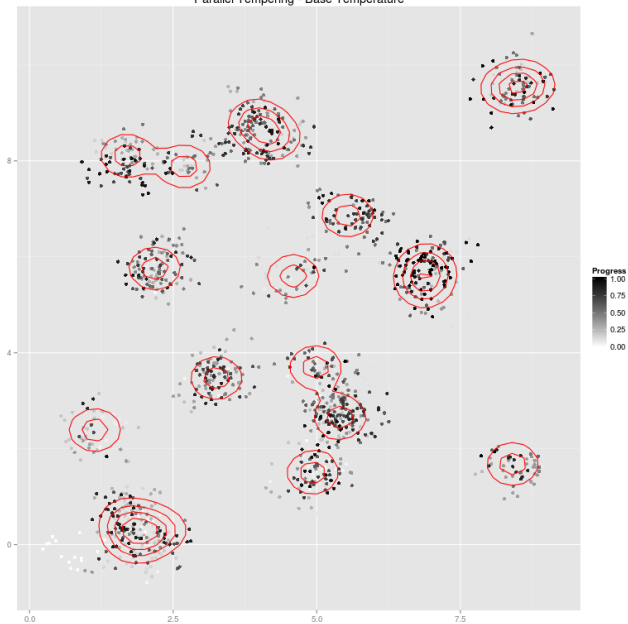
- softens the requirement $\pi(x_j) \approx \pi(x_i)$
- promotes $\beta_i - \beta_j \approx 0$
- promotes swaps between adjacent chains

Strategy IV

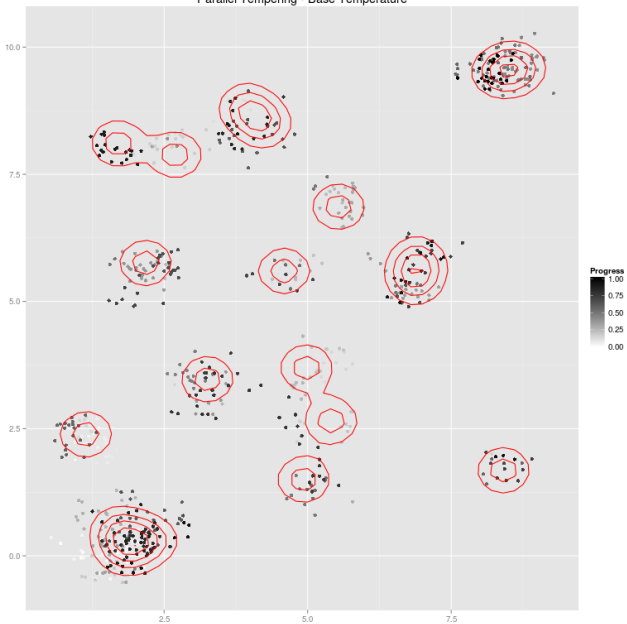
$$p_{ij} \propto \left(\frac{\pi(x_j)}{\pi(x_i)} \wedge \frac{\pi(x_i)}{\pi(x_j)} \right)^{\frac{\beta_i - \beta_j}{1 + \rho(x_i, x_j)}} = \exp \left(-\frac{(\beta_i - \beta_j) |\log(\pi(x_j)) - \log(\pi(x_i))|}{1 + \rho(x_i, x_j)} \right)$$

- added a quasi-metric
- ρ does not require symmetry $\rho(x_i, x_j) = \rho(x_j, x_i)$
- could be of use in the Gibbs random-field model

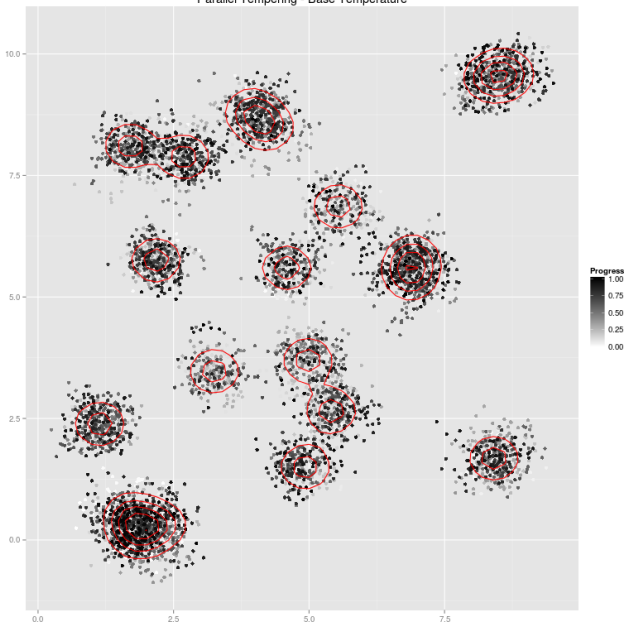
Parallel Tempering - Base Temperature



Parallel Tempering - Base Temperature



Parallel Tempering - Base Temperature



BIBLIOGRAPHY



Błażej Miasojedow, Eric Moulines, Matti Vihola, *Adaptive Parallel Tempering Algorithm*, Arxiv.



Meili Baragatti, Agnès Grimaud, Denys Pommeret *Parallel Tempering with Equi-Energy Moves*.



Charles J. Geyer, *Markov Chain Monte Carlo Lecture Notes*.