

Drzewa losowe i algorytmy Monte Carlo (Random trees and Monte Carlo algorithms)

Wojciech Niemirow¹

Kollokwium Wydziałowe, MIMUW, październik 2019

¹Mostly based on joint work with Tomasz Cąkała and Błażej Miasojedow (MIMUW).

Introduction and Simplified Example

Importance Sampling

MCMC: Metropolis-Hastings Algorithm

Lemma

Poisson Tree MCMC

Hidden Markov Model

Poisson Tree Particle Filter

Extended proposal & extended target

What is new?

Parallel computations

Continuous time models

IS, SMC, MCMC, PMCMC

Algorithms for sampling from (complicated) probability distributions:

Importance (Weighted) Sampling, IS

Markov Chain Monte Carlo, MCMC

Sequential Monte Carlo, SMC

Particle MCMC



```
graph LR; IS[Importance (Weighted) Sampling, IS] --> PMCMC[Particle MCMC]; MCMC[Markov Chain Monte Carlo, MCMC] --> PMCMC; SMC[Sequential Monte Carlo, SMC] --> PMCMC;
```

Introduction: Simplified Example

$p(\cdot)$ – probability density on space \mathcal{X} . Let

$$\pi(x) = \frac{p(x)w(x)}{z}.$$

Introduction: Simplified Example

$p(\cdot)$ – probability density on space \mathcal{X} . Let

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Assume that:

- ▶ $p(\cdot)$ is easy to sample from.
- ▶ $w(\cdot) \geq 0$ is a weight function.
- ▶ $\pi(\cdot)$ is of interest (difficult *target* distribution).
- ▶ $z = \int p(\mathbf{x})w(\mathbf{x})d\mathbf{x}$ is a normalizing constant (intractable).

Introduction: Simplified Example

$p(\cdot)$ – probability density on space \mathcal{X} . Let

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Assume that:

- ▶ $p(\cdot)$ is easy to sample from.
- ▶ $w(\cdot) \geq 0$ is a weight function.
- ▶ $\pi(\cdot)$ is of interest (difficult *target* distribution).
- ▶ $z = \int p(\mathbf{x})w(\mathbf{x})d\mathbf{x}$ is a normalizing constant (intractable).

We can sample $\mathbf{X} \sim p(\cdot)$ and compute $w(\mathbf{x})$ for a given $\mathbf{x} \in \mathcal{X}$.

We want to sample from $\pi(\cdot)$ and to compute z .

Simplified Example: statistical motivation

Bayes formula:

$$p(x|y) = \frac{p(x)p(y|x)}{p(y)}.$$

Simplified Example: statistical motivation

Bayes formula:

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x})p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y})}.$$

Bayesian statistics: \mathbf{X} is hidden, we observe $\mathbf{Y} = \mathbf{y}$

- ▶ $p(\mathbf{x})$ is the *prior* distribution.
- ▶ $w(\mathbf{x}) = p(\mathbf{y}|\mathbf{x})$ is the *likelihood* function.
- ▶ $\pi(\mathbf{x}) = p(\mathbf{x}|\mathbf{y})$ is the *posterior* distribution.

(\mathbf{y} is fixed and thus omitted.)

Simplified Example: Importance Sampling

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Simplified Example: Importance Sampling

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Sampling scheme:

- ▶ $N \sim \text{Poiss}(\lambda)$, i.e. $\mathbf{P}(N = n) = e^{-\lambda} \frac{\lambda^n}{n!}$, $n = 0, 1, 2, \dots$
- ▶ If $N = 0$ then $\hat{\mathbf{Z}} := \mathbf{0}$ end.
- ▶ If $N > 0$ then sample $\mathbf{X}_1, \dots, \mathbf{X}_N \sim_{\text{iid}} p(\cdot)$,

Simplified Example: Importance Sampling

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Sampling scheme:

- ▶ $N \sim \text{Pois}(\lambda)$, i.e. $\mathbf{P}(N = n) = e^{-\lambda} \frac{\lambda^n}{n!}$, $n = 0, 1, 2, \dots$
- ▶ If $N = 0$ then $\hat{\mathbf{Z}} := \mathbf{0}$ end.
- ▶ If $N > 0$ then sample $\mathbf{X}_1, \dots, \mathbf{X}_N \sim_{\text{iid}} p(\cdot)$,

$$\hat{\mathbf{Z}} := \frac{1}{\lambda} \sum_{j=1}^N w(\mathbf{X}_j), \quad \hat{\pi}(\cdot) = \frac{1}{\lambda} \sum_{j=1}^N \delta_{\mathbf{X}_j}(\cdot) w(\mathbf{X}_j).$$

Simplified Example: Importance Sampling

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

- ▶ $\mathbf{E}w(\mathbf{X}_j) = \int p(\mathbf{x})w(\mathbf{x})d\mathbf{x} = z,$
- ▶ $\mathbf{E}\left(\sum_{j=1}^N w(\mathbf{X}_j) \mid N\right) = Nz$ and $\mathbf{E}N = \lambda,$
- ▶ $\mathbf{E}\sum_{j=1}^N w(\mathbf{X}_j) = \lambda z.$

We thus have $\mathbf{E}\hat{\mathbf{Z}} = z.$

Simplified Example: Importance Sampling

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

- ▶ $\mathbf{E}w(\mathbf{X}_j) = \int p(\mathbf{x})w(\mathbf{x})d\mathbf{x} = z,$
- ▶ $\mathbf{E}\left(\sum_{j=1}^N w(\mathbf{X}_j) \mid N\right) = Nz$ and $\mathbf{E}N = \lambda,$
- ▶ $\mathbf{E}\sum_{j=1}^N w(\mathbf{X}_j) = \lambda z.$

We thus have $\mathbf{E}\hat{\mathbf{Z}} = z.$

Similarly, $\mathbf{E}\hat{\pi}(\mathbf{B}) = \pi(\mathbf{B})$ for every $\mathbf{B} \subseteq \mathcal{X}.$

Simplified Example: Importance Sampling

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

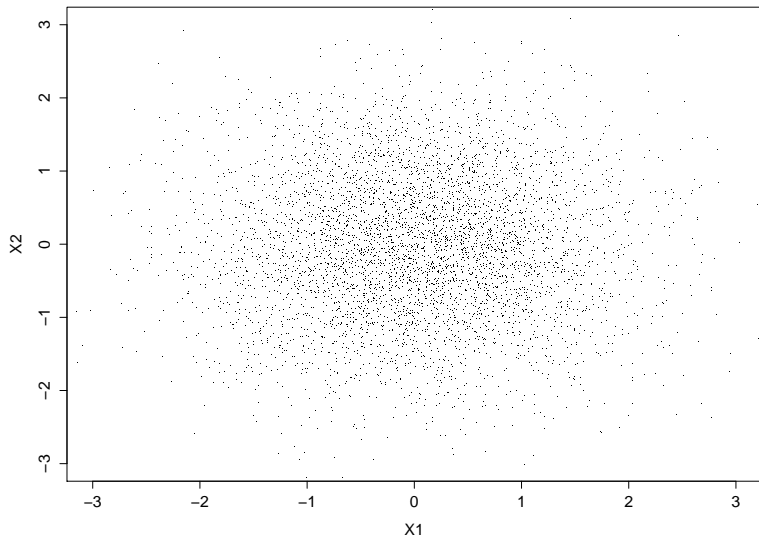
- ▶ $\mathbf{E}w(\mathbf{X}_j) = \int p(\mathbf{x})w(\mathbf{x})d\mathbf{x} = z,$
- ▶ $\mathbf{E}\left(\sum_{j=1}^N w(\mathbf{X}_j) \mid N\right) = Nz$ and $\mathbf{E}N = \lambda,$
- ▶ $\mathbf{E}\sum_{j=1}^N w(\mathbf{X}_j) = \lambda z.$

We thus have $\mathbf{E}\hat{\mathbf{Z}} = z.$

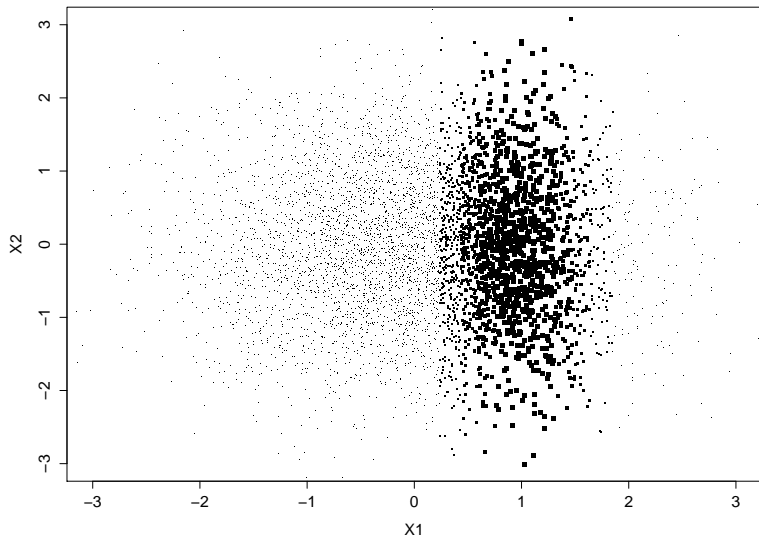
Similarly, $\mathbf{E}\hat{\pi}(\mathbf{B}) = \pi(\mathbf{B})$ for every $\mathbf{B} \subseteq \mathcal{X}.$

If $\lambda \rightarrow \infty$ then $N \rightarrow \infty$ and $\hat{\pi}(\mathbf{B}) \rightarrow \pi(\mathbf{B})$ a.s.

Sample from $p(\cdot)$



Weighted sample from $p(\cdot)$ approximates $\pi(\cdot)$



MCMC: Metropolis-Hastings Algorithm

Target distribution: $\phi(\xi)$ on space Ξ .

We generate Markov chain $\xi^{(0)}, \xi^{(1)}, \dots, \xi^{(m)}, \dots \rightarrow \phi(\cdot)$:

MCMC: Metropolis-Hastings Algorithm

Target distribution: $\phi(\xi)$ on space Ξ .

We generate Markov chain $\xi^{(0)}, \xi^{(1)}, \dots, \xi^{(m)}, \dots \rightarrow \phi(\cdot)$:

- ▶ Proposal distribution (transition density): $\psi(\xi, \xi')$.
- ▶ Acceptance probability:

$$\alpha(\xi, \xi') = \frac{\phi(\xi')\psi(\xi', \xi)}{\phi(\xi)\psi(\xi, \xi')} \wedge 1.$$

In one step we sample

$$\xi^{(m)} \sim T(\xi^{(m-1)}, \cdot),$$

where, for $\xi' \neq \xi$,

$$T(\xi, \xi') = \psi(\xi, \xi')\alpha(\xi, \xi').$$

and $T(\xi, \{\xi\}) = 1 - \int_{\xi' \neq \xi} T(\xi, \xi') d\xi'$.

MCMC: Metropolis-Hastings Algorithm

One step of MHA from $\xi^{(m-1)}$ to $\xi^{(m)}$:

Sample $\xi' \sim \psi(\xi^{(m-1)}, \cdot)$ { proposal }

Sample $U \sim \mathbf{Unif}(0, 1)$

if $U \leq \alpha(\xi^{(m-1)}, \xi')$ then

$\xi^{(m)} := \xi'$ { move accepted with probability α }

else

$\xi^{(m)} := \xi^{(m-1)}$ { move rejected with probability $1 - \alpha$ }

end if

MCMC: Metropolis-Hastings Algorithm

Theorem

Transition kernel of MHA is ϕ -reversible, i.e.

$$\phi(\xi)T(\xi, \xi') = \phi(\xi')T(\xi', \xi)$$

Consequently, ϕ is the stationary (equilibrium) distribution of the chain $\xi^{(m)}$.

Metropolis et al. (1953), Hastings (1970).

Independent Metropolis-Hastings (IMHA)

Proposal distribution $\psi(\xi, \xi') = \psi(\xi')$ is independent of ξ .

Acceptance probability:

$$\alpha(\xi, \xi') = \frac{\phi(\xi')\psi(\xi)}{\phi(\xi)\psi(\xi')} \wedge 1.$$

Simplified Example

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Simplified Example

Target distribution:

$$\pi(\mathbf{x}) = \frac{p(\mathbf{x})w(\mathbf{x})}{z}.$$

Sampling scheme continued. We draw $(N, \mathbf{X}_1, \dots, \mathbf{X}_N, \mathbf{S})$:

- ▶ $N \sim \text{Pois}(\lambda)$, i.e. $\mathbf{P}(N = n) = e^{-\lambda} \frac{\lambda^n}{n!}$, $n = 0, 1, 2, \dots$
- ▶ If $N = 0$ then $\hat{\mathbf{Z}} := 0$ end.
- ▶ If $N > 0$ then sample $\mathbf{X}_1, \dots, \mathbf{X}_N \sim_{\text{iid}} p(\cdot)$,

$$\hat{\mathbf{Z}} := \frac{1}{\lambda} \sum_{j=1}^N w(\mathbf{X}_j).$$

Choose $\mathbf{S} \in \{1, \dots, N\}$ at random:

$$\mathbf{P}(\mathbf{S} = \mathbf{s} | N, \mathbf{X}_1, \dots, \mathbf{X}_N) = \frac{w(\mathbf{X}_s)}{\sum_{j=1}^N w(\mathbf{X}_j)}.$$

Simplified Example

ψ – joint probability distribution of all random variables:

$$\psi(\mathbf{n}, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}) = e^{-\lambda} \frac{\lambda^n}{n!} \prod_{j=1}^n p(\mathbf{x}_j) \frac{w(\mathbf{x}_s)}{\sum_{j=1}^n w(\mathbf{x}_j)}.$$

for $\mathbf{n} > \mathbf{0}$ and $\psi(\mathbf{0}) = e^{-\lambda}$.

Simplified Example

ψ – joint probability distribution of all random variables:

$$\psi(\mathbf{n}, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}) = e^{-\lambda} \frac{\lambda^n}{n!} \prod_{j=1}^n p(\mathbf{x}_j) \frac{w(\mathbf{x}_s)}{\sum_{j=1}^n w(\mathbf{x}_j)}.$$

for $\mathbf{n} > \mathbf{0}$ and $\psi(\mathbf{0}) = e^{-\lambda}$.

Lemma

$$\psi(\mathbf{n}, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}) = \frac{\mathbf{z}}{\hat{\mathbf{z}}} \phi(\mathbf{n}, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}),$$

where ϕ is a probability distribution such that $\phi(\mathbf{0}) = \mathbf{0}$ and

$$\phi(\mathbf{x}_s) = \pi(\mathbf{x}_s).$$

Marginal of ϕ is the target π !

Simplified Example

Proof of Lemma:

$$\begin{aligned}\psi(n, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{s}) &= e^{-\lambda} \frac{\lambda^n}{n!} \prod_{j=1}^n p(\mathbf{x}_j) \frac{w(\mathbf{x}_s)}{\sum_{j=1}^n w(\mathbf{x}_j)} \\ &= \frac{1}{z} p(\mathbf{x}_s) w(\mathbf{x}_s) \frac{z}{\hat{z}} \frac{1}{n} e^{-\lambda} \frac{\lambda^{n-1}}{(n-1)!} \prod_{j \neq s} p(\mathbf{x}_j) \\ &= \underbrace{\frac{\pi(\mathbf{x}_s)}{\hat{z}}}_{\text{target}} \cdot \frac{1}{n} \psi(n-1, \mathbf{x}_{-s}).\end{aligned}$$

□

Simplified Example: IMHA

Independent Metropolis-Hastings chain on the **extended space** – space of configurations $\xi = (\mathbf{n}, \mathbf{X}, \mathbf{S})$:

- ▶ If the current state is $\xi^{(m-1)} = (\mathbf{N}, \mathbf{X}, \mathbf{S})$ then
- ▶ Draw a proposal: $(\mathbf{N}', \mathbf{X}', \mathbf{S}') \sim \psi$ (sampling scheme as described). Compute

$$\alpha := \frac{\phi(\mathbf{N}', \mathbf{X}', \mathbf{S}')\psi(\mathbf{N}, \mathbf{X}, \mathbf{S})}{\phi(\mathbf{N}, \mathbf{X}, \mathbf{S})\psi(\mathbf{N}', \mathbf{X}', \mathbf{S}')} \wedge \mathbf{1} = \frac{\hat{\mathbf{Z}}'}{\hat{\mathbf{Z}}} \wedge \mathbf{1}.$$

- ▶ with probability α accept: $\xi^{(m)} := (\mathbf{N}', \mathbf{X}', \mathbf{S}')$;
- ▶ with probability $\mathbf{1} - \alpha$ reject: $\xi^{(m)} := (\mathbf{N}, \mathbf{X}, \mathbf{S})$.

Simplified Example: IMHA

Independent Metropolis-Hastings chain on the **extended space** – space of configurations $\xi = (\mathbf{n}, \mathbf{x}, \mathbf{s})$:

- ▶ If the current state is $\xi^{(m-1)} = (\mathbf{N}, \mathbf{X}, \mathbf{S})$ then
- ▶ Draw a proposal: $(\mathbf{N}', \mathbf{X}', \mathbf{S}') \sim \psi$ (sampling scheme as described). Compute

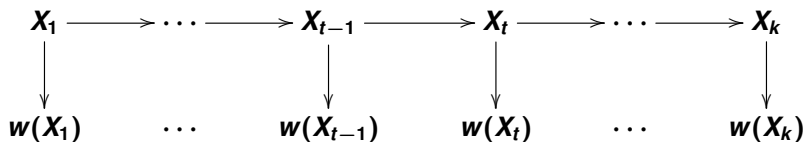
$$\alpha := \frac{\phi(\mathbf{N}', \mathbf{X}', \mathbf{S}')\psi(\mathbf{N}, \mathbf{X}, \mathbf{S})}{\phi(\mathbf{N}, \mathbf{X}, \mathbf{S})\psi(\mathbf{N}', \mathbf{X}', \mathbf{S}')} \wedge \mathbf{1} = \frac{\hat{\mathbf{Z}}'}{\hat{\mathbf{Z}}} \wedge \mathbf{1}.$$

- ▶ with probability α accept: $\xi^{(m)} := (\mathbf{N}', \mathbf{X}', \mathbf{S}')$;
- ▶ with probability $\mathbf{1} - \alpha$ reject: $\xi^{(m)} := (\mathbf{N}, \mathbf{X}, \mathbf{S})$.

The chain preserves $\phi(\mathbf{n}, \mathbf{x}, \mathbf{s})$ thus marginally $\pi(\mathbf{x}_s)$. It converges to the target distribution.

Extended Proposal: ψ , Extended Target: ϕ .

Hidden Markov Model



- ▶ $\mathbf{X} = \mathbf{X}_{1:k} = (\mathbf{X}_1, \dots, \mathbf{X}_k)$ is a hidden Markov chain with transition kernel $p(\mathbf{x}_{t-1}, \mathbf{x}_t) = p(\mathbf{x}_t | \mathbf{x}_{t-1})$.
- ▶ Likelihood weights : $w(\mathbf{x}_t) = p(\mathbf{y}_t | \mathbf{x}_t)$, where \mathbf{y}_t is observed.

Target (posterior) distribution:

$$\pi(\mathbf{x}_{1:k}) = \frac{1}{Z} \prod_{t=1}^k p(\mathbf{x}_{t-1}, \mathbf{x}_t) w(\mathbf{x}_t).$$

Algorithm PTPF (Poisson Tree Particle Filter)

Directed tree with “marked” nodes:

$$(\mathcal{V}, \mathcal{E}, \mathbf{X} = \{\mathbf{X}_v, v \in \mathcal{V}\}, \mathbf{S}).$$

Notations:

- ▶ $\mathbf{0}$ – fictitious root,
- ▶ $\mathbf{ch}(v)$ – children of node $v \in \mathcal{V}$,
- ▶ $\mathbf{an}(v)$ – ancestors of $v \in \mathcal{V}$ (with v , without $\mathbf{0}$),
- ▶ \mathcal{V}_t – t th generation ($t = 1, \dots, k$),
- ▶ \mathbf{S} – selected final node $\in \mathcal{V}_k$.

Algorithm PTPF (Poisson Tree Particle Filter)

Sampling scheme: For $t = 0, 1, \dots, k$, for every node $v \in \mathcal{V}_t$

- ▶ Choose Λ_v depending on history,
- ▶ $N_v \sim \text{Poiss}(\Lambda_v w(X_v))$,
- ▶ Create set $\text{ch}(v)$ of cardinality N_v ,
- ▶ For every $u \in \text{ch}(v)$ sample $X_u \sim p(X_v, \cdot)$ propagate,
- ▶ Compute $W_u := w(X_u)$ weigh

until $t = k$ or $\sum_{v \in \mathcal{V}_t} N_t = 0$;

Algorithm PTPF (Poisson Tree Particle Filter)

Sampling scheme: For $t = 0, 1, \dots, k$, for every node $v \in \mathcal{V}_t$

- ▶ Choose Λ_v depending on history,
- ▶ $N_v \sim \text{Poiss}(\Lambda_v w(X_v))$,
- ▶ Create set $\text{ch}(v)$ of cardinality N_v ,
- ▶ For every $u \in \text{ch}(v)$ sample $X_u \sim p(X_v, \cdot)$ propagate,
- ▶ Compute $W_u := w(X_u)$ weigh

until $t = k$ or $\sum_{v \in \mathcal{V}_t} N_t = 0$;

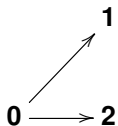
$\hat{Z} := \sum_{u \in \mathcal{V}_k} w(X_u) / C_u$, where $C_u = \lambda_0 \prod_{i \in \text{an}(u) - u} \Lambda_v$,

Select $S \in \mathcal{V}_k$: $\mathbf{P}(S = s) \propto w(X_s) / C_s$.

Example: *extended proposal*

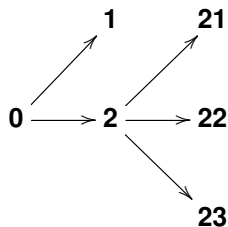
0

Example: *extended proposal*



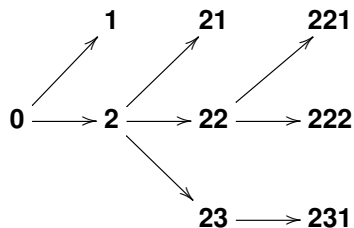
$$\psi = \exp[-\lambda_0] (\lambda_0)^2 p(x_0, x_1) p(x_0, x_2)$$

Example: *extended proposal*



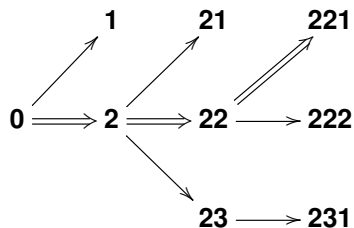
$$\begin{aligned}\psi &= \exp[-\lambda_0] (\lambda_0)^2 p(x_0, x_1) p(x_0, x_2) \\ &\times \exp[-\lambda_1 w_1] \\ &\times \exp[-\lambda_2 w_2] (\lambda_2 w_2)^3 p(x_2, x_{21}) p(x_2, x_{22}) p(x_2, x_{23})\end{aligned}$$

Example: *extended proposal*



$$\begin{aligned} \psi &= \mathbf{exp}[-\lambda_0] (\lambda_0)^2 \mathbf{p}(\mathbf{x}_0, \mathbf{x}_1) \mathbf{p}(\mathbf{x}_0, \mathbf{x}_2) \\ &\times \mathbf{exp}[-\lambda_1 w_1] \\ &\times \mathbf{exp}[-\lambda_2 w_2] (\lambda_2 w_2)^3 \mathbf{p}(\mathbf{x}_2, \mathbf{x}_{21}) \mathbf{p}(\mathbf{x}_2, \mathbf{x}_{22}) \mathbf{p}(\mathbf{x}_2, \mathbf{x}_{23}) \\ &\times \mathbf{exp}[-\lambda_{21} w_{21}] \\ &\times \mathbf{exp}[-\lambda_{22} w_{22}] (\lambda_{22} w_{22})^2 \mathbf{p}(\mathbf{x}_{22}, \mathbf{x}_{221}) \mathbf{p}(\mathbf{x}_{22}, \mathbf{x}_{222}) \\ &\times \mathbf{exp}[-\lambda_{23} w_{23}] (\lambda_{23} w_{23})^1 \mathbf{p}(\mathbf{x}_{23}, \mathbf{x}_{231}) \end{aligned}$$

Example: *extended proposal*



$$\begin{aligned} \psi &= \exp[-\lambda_0] (\lambda_0)^2 p(x_0, x_1) p(x_0, x_2) \\ &\times \exp[-\lambda_1 w_1] \\ &\times \exp[-\lambda_2 w_2] (\lambda_2 w_2)^3 p(x_2, x_{21}) p(x_2, x_{22}) p(x_2, x_{23}) \\ &\times \exp[-\lambda_{21} w_{21}] \\ &\times \exp[-\lambda_{22} w_{22}] (\lambda_{22} w_{22})^2 p(x_{22}, x_{221}) p(x_{22}, x_{222}) \\ &\times \exp[-\lambda_{23} w_{23}] (\lambda_{23} w_{23})^1 p(x_{23}, x_{231}) \\ &\times \frac{1}{\hat{z}} \frac{w_{221}}{\lambda_0 \lambda_2 \lambda_{22}}; \quad \hat{z} = \frac{w_{221}}{\lambda_0 \lambda_2 \lambda_{22}} + \frac{w_{222}}{\lambda_0 \lambda_2 \lambda_{22}} + \frac{w_{231}}{\lambda_0 \lambda_2 \lambda_{23}} \end{aligned}$$

Extended proposal and extended target

Extended proposal ψ is the joint probability distribution of all the variables produced by PTPF.

Lemma

If $\mathcal{V}_k \neq \emptyset$ then

$$\psi(\mathcal{V}, \mathcal{E}, \mathbf{x}, \mathbf{s}) = \phi(\mathcal{V}, \mathcal{E}, \mathbf{x}, \mathbf{s}) \frac{\mathbf{z}}{\hat{\mathbf{z}}},$$

where ϕ (**extended target**) is such that the marginal of $\mathbf{x}_{\text{an}(\mathbf{s})}$ is the target: $\phi(\mathbf{x}_{\text{an}(\mathbf{s})}) = \pi(\mathbf{x}_{\text{an}(\mathbf{s})})$ and $\phi(\mathbf{0}) = \mathbf{0}$.

Poisson Tree IMHA

Independent Metropolis-Hastings chain on the **extended space** – space of **trees** $\xi = (\mathcal{V}, \mathcal{E}, \mathbf{x} = \{\mathbf{x}_v, v \in \mathcal{V}\}, \mathbf{s})$:

- ▶ If the current state is $\xi^{(m-1)}$ then
- ▶ Draw a proposal: run PTPF to obtain $\xi' \sim \psi$. Compute $\hat{\mathbf{Z}}'$.

$$\alpha := \frac{\phi(\xi')\psi(\xi)}{\phi(\xi)\psi(\xi')} \wedge \mathbf{1} = \frac{\hat{\mathbf{Z}}'}{\hat{\mathbf{Z}}} \wedge \mathbf{1}.$$

- ▶ with probability α accept: $\xi^{(m)} := \xi'$;
- ▶ with probability $\mathbf{1} - \alpha$ reject: $\xi^{(m)} := \xi^{(m-1)}$.

Poisson Tree IMHA

Independent Metropolis-Hastings chain on the **extended space** – space of **trees** $\xi = (\mathcal{V}, \mathcal{E}, \mathbf{x} = \{\mathbf{x}_v, v \in \mathcal{V}\}, \mathbf{s})$:

- ▶ If the current state is $\xi^{(m-1)}$ then
- ▶ Draw a proposal: run PTPF to obtain $\xi' \sim \psi$. Compute $\hat{\mathbf{Z}}'$.

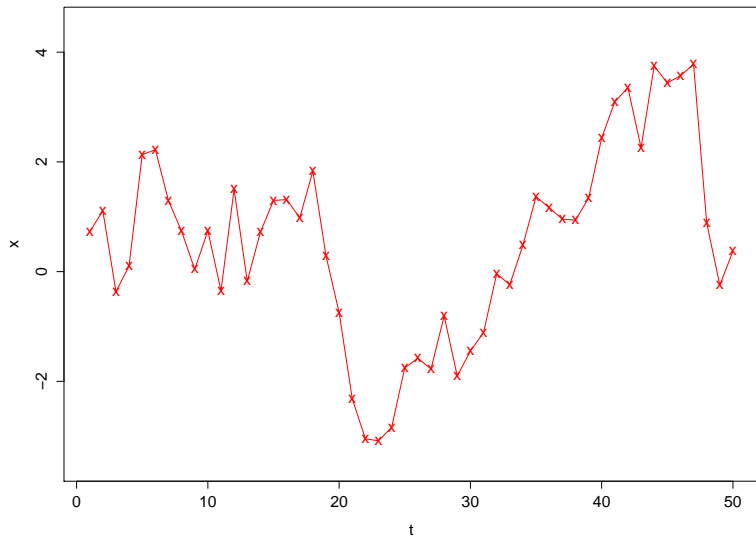
$$\alpha := \frac{\phi(\xi')\psi(\xi)}{\phi(\xi)\psi(\xi')} \wedge \mathbf{1} = \frac{\hat{\mathbf{Z}}'}{\hat{\mathbf{Z}}} \wedge \mathbf{1}.$$

- ▶ with probability α accept: $\xi^{(m)} := \xi'$;
- ▶ with probability $\mathbf{1} - \alpha$ reject: $\xi^{(m)} := \xi^{(m-1)}$.

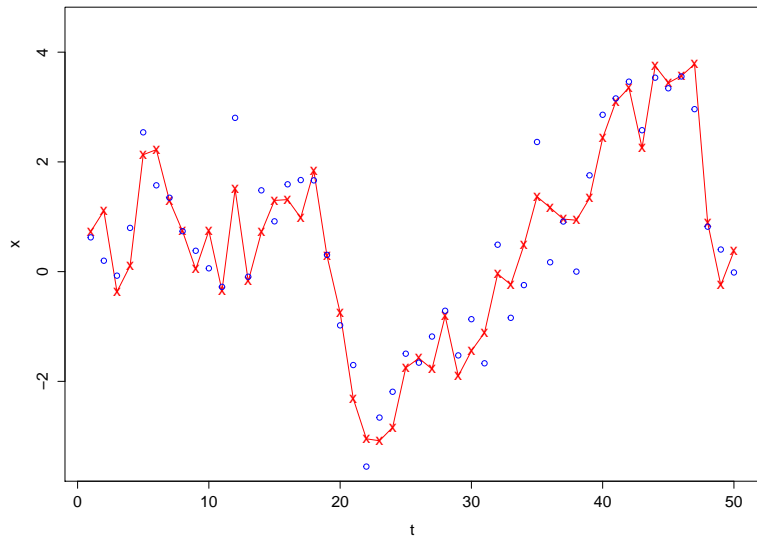
The chain preserves $\phi(\xi)$ thus marginally $\pi(\mathbf{x}_{\text{an}(\mathbf{s})})$. It converges to the target distribution.

Extended Proposal: ψ , Extended Target: ϕ .

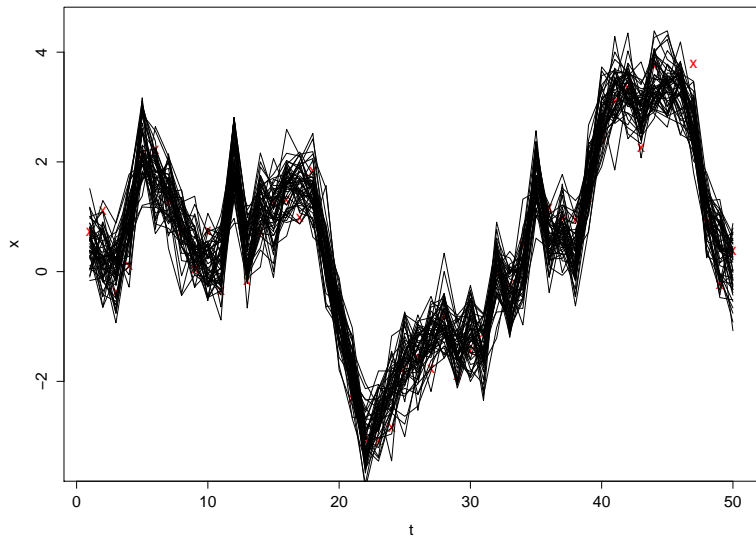
Hidden process



Hidden process and observations



Trajectories sampled from the posterior via PTMC



What is new in PTPF?

Parallelization of computations:

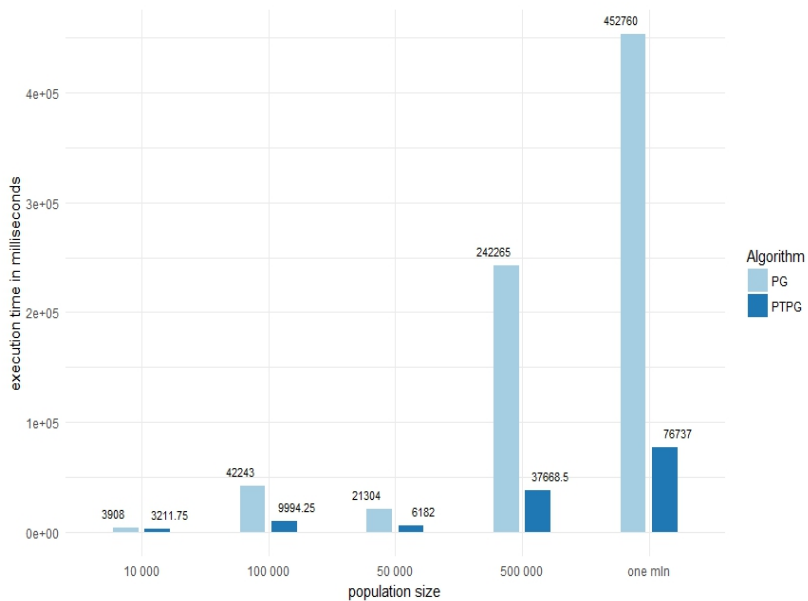
- ▶ If $\Lambda_{\mathbf{v}}$ depends only on $\mathbf{an}(\mathbf{v}) - \mathbf{v}$ then the branches of the tree evolve **completely independently**.

What is new in PTPF?

Parallelization of computations:

- ▶ If $\Lambda_{\mathbf{v}}$ depends only on $\mathbf{an}(\mathbf{v}) - \mathbf{v}$ then the branches of the tree evolve **completely independently**.
- ▶ If the branches evolve **partly independently**, we can control their number.

Parallel computations are more efficient



What is new in PTPF?

Continuous time models:

- ▶ PTMC can be directly applied to hidden piecewise deterministic (quasi-) Markov processes.

What is new in PTPF?

Continuous time models:

- ▶ PTMC can be directly applied to hidden piecewise deterministic (quasi-) Markov processes.
- ▶ If the branches evolve **partly independently**, we can control their number.

A tree with space-time nodes

$[t_{\min}]$



0

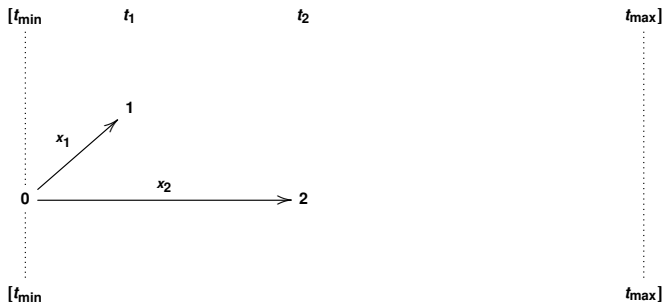
$[t_{\min}]$

t_{\max}

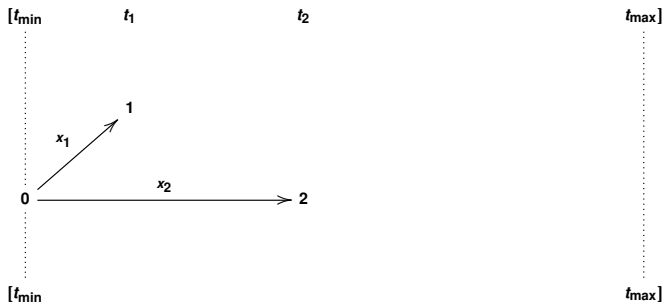


t_{\max}

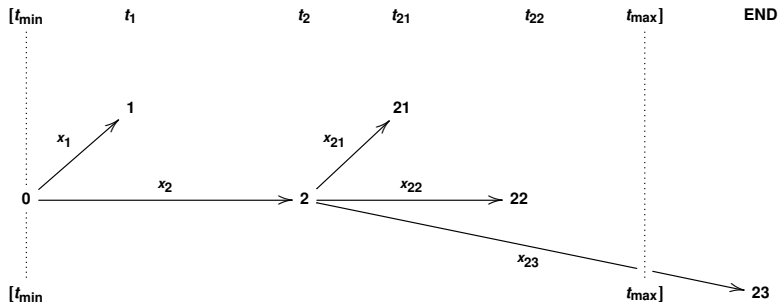
A tree with *space-time* nodes



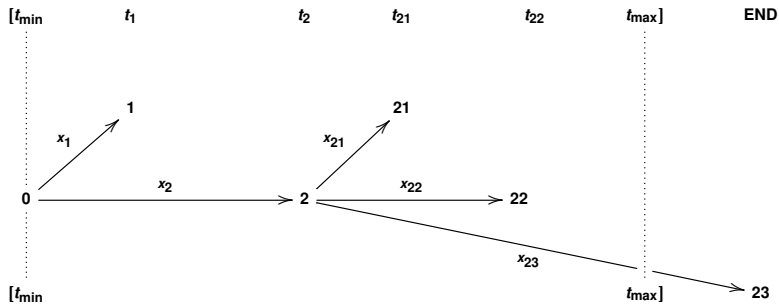
A tree with *space-time* nodes



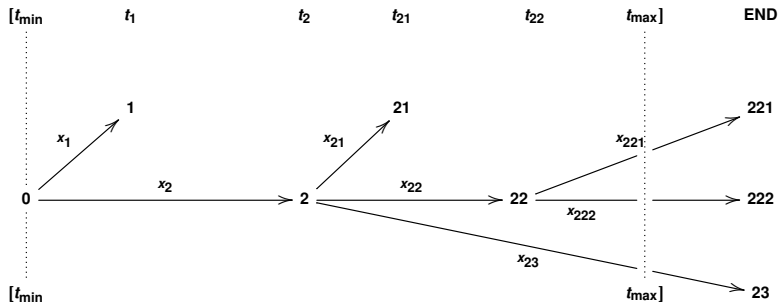
A tree with *space-time* nodes



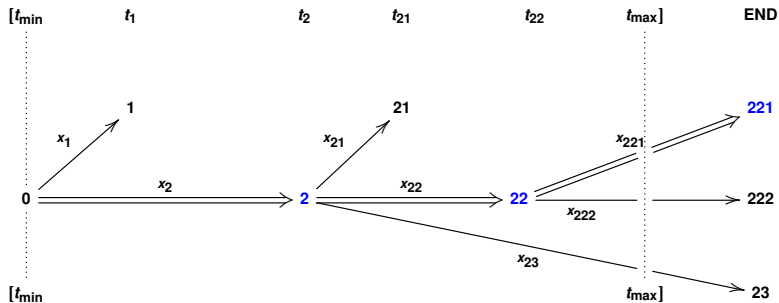
A tree with *space-time* nodes



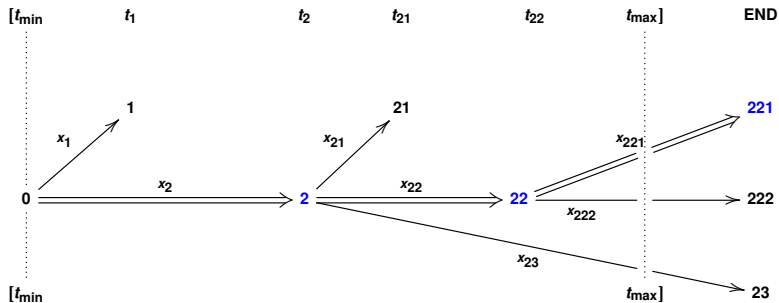
A tree with *space-time* nodes







A tree with *space-time* nodes



A tree with *space-time* nodes



References

-  C. Andrieu, A. Doucet, R. Holenstein: Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society B*, 2010
-  F. Lindsten, M.I. Jordan, T.B. Schön: Particle Gibbs with Ancestor Sampling, *Journal of Machine Learning Research*, 2014
-  F. Lindsten, R. Douc, E. Moulines: Uniform ergodicity of the Particle Gibbs sampler, *Scandinavian Journal of Statistics*, 2015
-  T. Çakała, B. Miasojedow, W. Niemiro: Particle MCMC algorithms with Poisson resampling: parallelization and continuous time models. [arXiv:1707.01660v2](https://arxiv.org/abs/1707.01660v2))