

Parameter elimination in particle Gibbs sampling

Anna Wigren (UU), Riccardo Sven Risuleo (UU), Lawrence Murray (Uber AI), Fredrik Lindsten (LiU)

The paper in 30 seconds

The performance of particle Gibbs samplers is limited by the underlying Gibbs sampler that is approximated. We can improve on this by eliminating the parameters from the state update using marginalization.

Marginalization is possible for models with conjugacy between the parameter prior and the complete data likelihood. It yields a non-Markovian state update, but can still be done in linear time.

Marginalizing by hand is cumbersome, but can be automated using probabilistic programming.

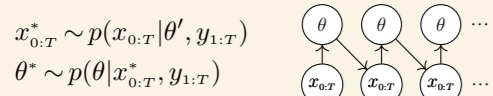
Take with you

- Particle Gibbs is limited by Gibbs
- We can improve using marginalization
- Automatic marginalization is possible through probabilistic programming

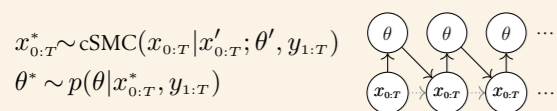
Gibbs sampling

We aim to infer the joint posterior $p(x_{0:T}, \theta | y_{1:T})$ of a state-space model.

Not analytically tractable in general, instead we use **Gibbs sampling**.



Sampling the high-dimensional state trajectory is challenging. **Particle Gibbs** uses conditional sequential Monte Carlo (cSMC) to do this.



Conditional sequential Monte Carlo: Represents the target $p(x_{0:T} | \theta^i, y_{1:T})$ as N particles $\{x_{0:T}^i, w_T^i\}_{i=1}^N$ generated from iterating, for $t = 1, \dots, T$:

- Resample a_t^i using w_{t-1}^i , set $a_t^N = N$
- Draw $x_t^i \sim q(x_t | x_{0:t-1}^i)$, set $x_t^N = x_t^i$
- Update $w_t^i \propto \frac{p(x_t^i | x_{0:t-1}^i, \theta)}{q(x_t^i | x_{0:t-1}^i)}$

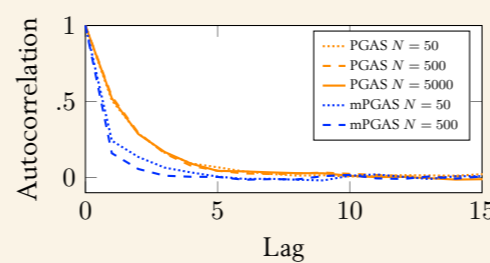


Motivation - why do we want to eliminate the parameters?

In Gibbs sampling there is a dependence between $x_{0:T}$ and θ which leads to correlated samples.

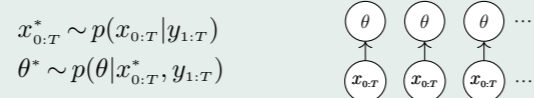
Particle Gibbs is an **exact approximation** of the Gibbs sampler and can never do better than the sampler it approximates.

By **marginalizing** out the parameters from the state update we can improve on this!

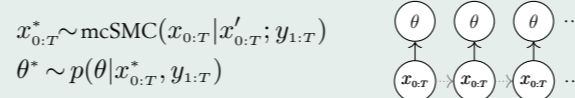


When can we marginalize?

Marginalized Gibbs sampling amounts to



Marginalized cSMC in the state update yields



Problem: State trajectory update is **non-Markovian**, in particular $w_t^i \propto \frac{p(x_t^i, y_t | x_{0:t-1}^i, y_{1:t-1})}{q(x_t^i | x_{0:t-1}^i)}$

Solution: Possible if we have **conjugacy** between parameter prior and complete data likelihood!

Restricted exponential family:

Complete data likelihood:

$$p(x_t, y_t | x_{t-1}, \theta) = h_t \exp(\theta^T s_t - A^T(\theta) r_t)$$

Prior:

$$\pi(\theta | \chi_0, \nu_0) = g(\chi_0, \nu_0) \exp(\theta^T \chi_0 - A^T(\theta) \nu_0)$$

$$p(x_t, y_t | x_{0:t-1}, y_{1:t-1}) = \frac{g(\chi_{t-1}, \nu_{t-1})}{g(\chi_t, \nu_t)} h_t$$

$$\chi_t = \chi_0 + \sum_{k=1}^t s_k = \chi_{t-1} + s_t$$

$$\nu_t = \nu_0 + \sum_{k=1}^t r_k = \nu_{t-1} + r_t$$

For the restricted exponential family marginalized cSMC amounts to iterating

- Update hyperparameters χ_t^i, ν_t^i
- Resample a_t^i using w_{t-1}^i , set $a_t^N = N$
- Draw $x_t^i \sim q(x_t | x_{0:t-1}^i)$, set $x_t^N = x_t^i$
- Update $w_t^i \propto \frac{g(\chi_{t-1}^i, \nu_{t-1}^i) h_t^i}{g(\chi_t^i, \nu_t^i) q(x_t^i | x_{0:t-1}^i)}$

We can run this in **linear time!**

Extension: Particle Gibbs with ancestor sampling (PGAS)

Particle Gibbs is known to mix slowly. To improve we can use marginalized **ancestor sampling**, where the ancestor for the reference trajectory is sampled according to $\tilde{w}_{t-1}^i = \bar{w}_{t-1}^i \frac{p([x_{0:t-1}^i, x_{t:T}^i] | y_{1:T})}{p(x_{0:t-1}^i | y_{1:t-1})}$

For the restricted exponential family we get $\tilde{w}_{t-1}^i \propto \bar{w}_{t-1}^i \frac{g(\chi_{t-1}^i, \nu_{t-1}^i)}{g(\chi_T^i, \nu_T^i)} h_t^i$ with $\chi_T^i = \chi_{t-1}^i + s_t(x_t^i, x_{t-1}^i, y_t) + s_{t+1:T}^i$ and $\nu_T^i = \nu_{t-1}^i + r_t(x_{t-1}^i) + r_{t+1:T}^i$

Additional steps in marginalized cSMC with ancestor sampling are

- Update statistics $s_{t+1:T}^i, r_{t+1:T}^i$
- Update hyperparameters χ_T^i, ν_T^i
- Sample a_t^N from \tilde{w}_{t-1}^i

We can run this in **linear time!**

Contact:

anna.wigren@it.uu.se
lawrence.murray@uber.com

riccardo.risuleo@it.uu.se
fredrik.lindsten@liu.se

Links:

Webpage Birch: birch-lang.org/



Probabilistic programming and Birch

Marginalization by hand is time-consuming, we do it automatically using probabilistic programming!

What is probabilistic programming?

- Tool for encoding probabilistic models and performing inference in them.
- Separates the model and the inference algorithm.
- Written in a probabilistic programming language.

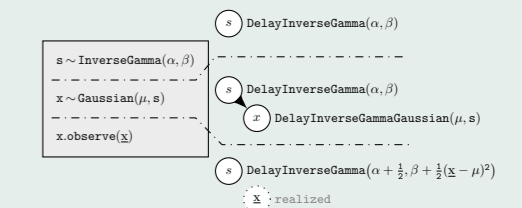
```

beta:Real <- 0.9;
sigma2:Real <- 1.0;
T:Integer <- 10;
x:Random<Real>[T];
y:Random<Real>[T];

x[1] ~ Gaussian(0.0, sigma2);
y[1] ~ Gaussian(x[1], sigma2);
for (t:Integer in 2..T) {
  x[t] ~ Gaussian(beta*x[t-1], sigma2);
  y[t] ~ Gaussian(x[t], sigma2);
}
    
```

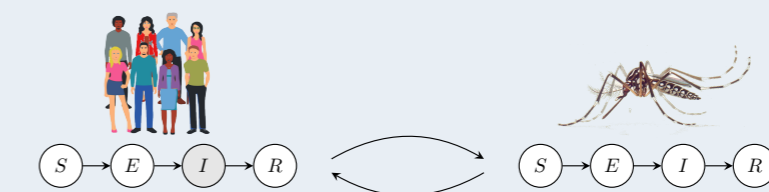
Birch:

- Object-oriented probabilistic programming language
- Uses inference methods based on SMC
- Employs **delayed sampling** to find and use conjugacy relations, enables **automatic marginalization!**

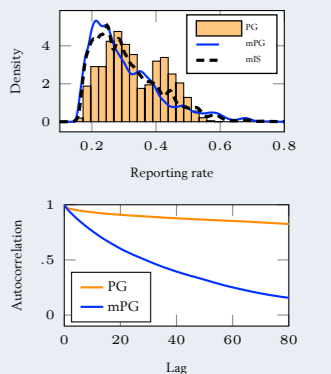


Example: Dengue fever dataset

The dataset contains the number of newly infected people during an outbreak of dengue fever on the island Yap in Micronesia.



Transitions between compartments are binomial, with beta priors on the parameters \rightarrow **conjugacy!**



Extension: Partial marginalization

For models where only some parameters are conjugate, split them into θ_m (conjugate) and θ_u (not conjugate).

Marginalize out θ_m and sample θ_u using any MCMC method.

We use marginalized PMMH in Birch.

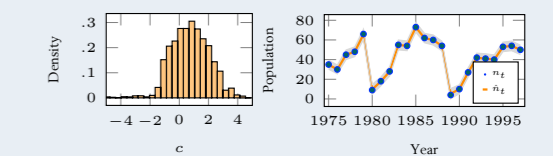
Models where some parameters can be marginalized out are common in practice.

Example: Song sparrow dataset

The dataset contains the number of observed song sparrows on an island during 20 years.

$$\log n_{t+1} = \log n_t + [1 - (n_t)^c] b + \sigma_w v_t$$

$$y_t = n_t + \sigma_w w_t$$



Extension: Diffuse priors and blocking

A diffuse prior on the parameters leads to a diffuse prior on the states when marginalizing.

This can lead to poor mixing and numerical problems for the first time steps.

We propose to separate the state trajectory into two overlapping blocks and do Gibbs updates of each block separately.

- Sample $x_{0:B+L}$ using standard cSMC
- Sample $x_{B+1:T}$ using marginalized cSMC
- Sample $\theta \sim p(\theta | x_{0:T}, y_{1:T})$

