

Summary

Parallel tempering (PT) is a class of MCMC algorithms that constructs a path of distributions annealing between a reference, π_0 , and intractable target, π_1 . States along the path are swapped to improve mixing in the target.

Problem: Past work on PT has only used linear paths with a fixed **reference** that is often different than the target. PT swapping can be computationally expensive.

Contribution: We extend to annealing paths with a variational reference and optimize the choice of reference. This improves PT swapping significantly.

Parallel tempering

Annealing path: π_{β} is a path of distributions between π_0 and π_1 . A linear path with a fixed reference is typically used: $\pi_{\beta}(x) \propto \pi_{0}^{1-\beta}(x) \cdot \pi_{1}^{\beta}(x)$

Run N+1 chains targeting π_{β_n} . Alternate between local exploration and communication moves.

Local exploration: Update each chain according to an MCMC algorithm.



Reference

Communication: Swap states between chains n and n+1 with probability α_n



Objective: Maximize the **restart rate** τ (rate at which samples from the reference chain travel to the target chain)

$$\tau(\pi_0, \pi_1) = \left(2 + 2\sum_{n=0}^{N-1} \frac{r_n}{1 - r_n}\right)^{-1}, \quad r_n =$$

Parallel tempering with a variational reference Nikola Surjanovic¹, Saifuddin Syed², Alexandre Bouchard-Côté¹, Trevor Campbell¹

 $1 - \mathbb{E}[\alpha_n]$

Suboptimality of a fixed reference

When the reference and target do not overlap much, PT will often reject communication swaps between chains.

Proposition (problem with standard PT): For an exchangeable* Bayesian model with m conditionally i.i.d. data points, PT with a fixed reference satisfies restart rate $\xrightarrow{a.s.} 0, \quad m \to \infty$

*I.e., satisfies usual technical assumptions for a Bernstein-von Mises result, etc.

Prior is far away from posterior in data limit (low restart rate)



Annealing paths with a variational reference

We introduce a variational reference family, $\{q_{\phi} : \phi \in \Phi\}$. The linear annealing path with the modified reference is

 $\pi_{\phi,\beta}(x) \propto q_d^{\perp}$

We consider exponential family reference distributions

and choose a reference distribution close to the target.

Prio

Proposition: For an exchangeable Bayesian model with m conditionally i.i.d. data points, PT with a normal variational reference that minimizes the forward KL satisfies restart rate \xrightarrow{P}

Variational reference tuning

We minimize the forward (inclusive) KL divergence,

 $\mathrm{KL}(\pi_1||q_\phi) = \mathbb{E}_{\pi_1}[\log \pi_1(X)] - \mathbb{E}_{\pi_1}[\log q_\phi(X)]$

and use a gradient-free procedure to tune the reference.

²Department of Statistics, University of Oxford

Posterior

Concentrates to a point mass

$$q_{\phi}^{1-\beta}(x)\cdot\pi_{1}^{\beta}(x)$$

 $q_{\phi}(x) = c(\phi)h(x)\exp(\phi^{\top}\eta(x))$



Posterior

Variational reference is close to posterior (high restart rate)

$$1/2, m \to \infty$$



Algorithm:

Note: To stabilize the tuning of the variational reference, we introduce PT with two references—one fixed and one variational. (Details provided in paper.)

Theorem (outline): The variational parameter estimates converge to the forward KL minimizer almost surely.

The restart rate at the forward KL minimum is bounded in terms of the flexibility of the variational family. (More details in paper.)

PT with a variational reference empirically outperforms PT with a fixed reference. Below: variational PT with a normal reference (mean-field approximation and full covariance) versus NRPT [Syed et al. 2021].



Variational reference tuning (...)

Use samples to tune the reference with moment matching

Choose ϕ so that $\mathbb{E}_{q_{\phi}}[\eta(X)] = T^{-1} \sum_{t=1}^{T} \eta(X_t)$

1) Run the non-reversible PT algorithm (**NRPT**) [Syed et al. 2021] 2) Use obtained samples X_1, X_2, \ldots, X_T from the target chain to update the annealing schedule using the procedure in [Syed et al. 2021] 3) Update ϕ so that $\mathbb{E}_{q_{\phi}}[\eta(X)] = T^{-1} \sum_{t=1}^{T} \eta(X_t)$ 4) Repeat 1-3 until the computational budget is depleted

Experiments

Email: nikola.surjanovic@stat.ubc.ca