

By similar identifications of terms we have

$$M_{21} = \tilde{K}_{21} L_{11}^{-\top}, \quad (5b)$$

$$M_{31} = L_{31}, \quad (5c)$$

$$M_{22}^2 = \tilde{K}_{22} - M_{21} M_{21}^\top, \quad (5d)$$

$$M_{32} = \frac{1}{M_{22}} \left(\tilde{K}_{32} - M_{31} M_{21}^\top \right), \quad (5e)$$

$$M_{33} M_{33}^\top = K_{33} - M_{31} M_{31}^\top - M_{32} M_{32}^\top, \quad (5f)$$

Using back-substitution it is straightforward to find the factors M_{21} and M_{31} by exploiting the fact that the L_{11} matrix is lower triangular. The term M_{22} is scalar and found simply via a square root operation and the term M_{32} is given by (5e). However, computing M_{33} by a direct Cholesky decomposition is too expensive and again it comes down to exploiting the structure inherent in the problem. Let us start by inserting (4f) into (5f), resulting in

$$M_{33} M_{33}^\top = L_{33} L_{33}^\top + L_{31} L_{31}^\top + L_{32} L_{32}^\top - M_{31} M_{31}^\top - M_{32} M_{32}^\top \quad (6a)$$

$$= L_{33} L_{33}^\top + L_{32} L_{32}^\top - M_{32} M_{32}^\top, \quad (6b)$$

where the last equality follows (5c). Since $L_{32} L_{32}^\top$ and $M_{32} M_{32}^\top$ are both rank one, we can now compute M_{33} by one rank-one update and one rank-one downdate of L_{33} . See [1] for an overview of rank-one update/downdate methods. This concludes our work in finding the decomposition $\tilde{K} = M M^\top$.

2 The CPF-AS Algorithm

The basic idea underlying PMCMC is to use SMC to construct a Markov kernel leaving the exact joint smoothing distribution invariant. Hence, we seek a family of ergodic Markov kernels on \mathcal{X}^{T+1} ,

$$\{M_\theta : \theta \in \Theta\}, \quad (7)$$

such that, for each θ , $M_\theta(\mathbf{x}_{0:T} | \tilde{\mathbf{x}}_{0:T})$ leaves $p(\mathbf{x}_{0:T} | \theta, \mathbf{y}_{0:T})$ invariant. In PGAS, these kernels are constructed using a procedure referred to as a conditional particle filter with ancestor sampling (CPF-AS). This procedure is particularly suitable for non-Markovian latent variable models [2], as it relies only on a forward recursion.

CPF-AS is similar to a standard SMC sampler, but with the important difference that one particle at each time step is specified *a priori*. Let these particles be denoted $\tilde{\mathbf{x}}_{0:T} = \{\tilde{\mathbf{x}}_0, \dots, \tilde{\mathbf{x}}_T\}$. More precisely, we condition on the event that $\tilde{\mathbf{x}}_t$ is contained in the collection of particles $\{\mathbf{x}_t^i\}_{i=1}^N$, generated at time t . To accomplish this, we sample according to $\mathbf{x}_t^i \sim p(\mathbf{x}_t | \theta, \mathbf{x}_{0:t-1}^{\mathbf{a}^i})$ only for $i = 1, \dots, N-1$. The N th particle is then set deterministically: $\mathbf{x}_t^N = \tilde{\mathbf{x}}_t$. The CPF-AS is given in Algorithm 1.

The conditioning on a pre-specified collection of particles implies an invariance property of the CPF-AS, which is key to its applicability in an MCMC sampler.

Proposition 1. *Let the support of the target density be a subset of the support of the proposal density. Then, for any θ and any $N \geq 2$, the procedure*

- (i) Run Algorithm 1 conditionally on $\tilde{\mathbf{x}}_{0:T}$;
- (ii) Sample $\tilde{\mathbf{x}}'_{0:T}$ with $\mathbb{P}(\tilde{\mathbf{x}}'_{0:T} = \mathbf{x}_{0:T}^i) = \mathbf{w}_T^i$;

defines an irreducible and aperiodic Markov kernel M_θ^N on \mathcal{X}^T , with invariant distribution $p(\mathbf{x}_{0:T} | \theta, \mathbf{y}_{0:T})$.

Proof. The invariance property follows by the construction of the CPF-AS in [2], and the fact that the law of $\tilde{\mathbf{x}}'_{0:T}$ is independent of permutations of the particle indices. This allows us to always place the conditioned particles at the N th position. Irreducibility and aperiodicity follows from [3, Theorem 5]. \square

Algorithm 1 CPF-AS, conditioned on $\tilde{\mathbf{x}}_{0:T}$

1. Initialize:

- (a) Draw $\mathbf{x}_0^i \sim p(\mathbf{x}_0 \mid \theta, \mathbf{y}_0)$ for $i = 1, \dots, N - 1$.
- (b) Set $\mathbf{x}_0^N = \tilde{\mathbf{x}}_0$.
- (c) For $i = 1, \dots, N$, set $\mathbf{w}_0^i \propto p(\mathbf{y}_0 \mid \theta, \mathbf{x}_0^i)$, where the weights are normalized to sum to 1.

2. For $t = 1, \dots, T$ do:

- (a) Draw \mathbf{a}_t^i with $P(\mathbf{a}_t^i = j) = \mathbf{w}_{t-1}^j$ for $i = 1, \dots, N - 1$.
 - (b) Draw $\mathbf{x}_t^i \sim p(\mathbf{x}_t \mid \theta, \mathbf{x}_{0:t-1}^i)$ for $i = 1, \dots, N - 1$.
 - (c) Draw \mathbf{a}_t^N with $\mathbb{P}(\mathbf{a}_t^N = j) \propto \mathbf{w}_{t-1}^j p(\tilde{\mathbf{x}}_{t:T} \mid \theta, \mathbf{x}_{1:t-1}^j)$.
 - (d) Set $\mathbf{x}_t^N = \tilde{\mathbf{x}}_t$.
 - (e) For $i = 1, \dots, N$, set $\mathbf{w}_t^i \propto p(\mathbf{y}_t \mid \theta, \mathbf{x}_t^i)$, where the weights are normalized to sum to 1.
-

Consequently, if $\tilde{\mathbf{x}}_{0:T} \sim p(\mathbf{x}_{0:T} \mid \theta, \mathbf{y}_{0:T})$ and we sample $\tilde{\mathbf{x}}'_{0:T}$ according to the procedure given in Proposition 1, then, for any number of particles N , it holds that $\tilde{\mathbf{x}}'_{0:T} \sim p(\mathbf{x}_{0:T} \mid \theta, \mathbf{y}_{0:T})$. For $N = 1$ we get, by construction, $\tilde{\mathbf{x}}'_{0:T} = \tilde{\mathbf{x}}_{0:T}$, i.e. the trajectories are perfectly correlated (this is why we need $N \geq 2$ to get an irreducible kernel). On the other hand, as $N \rightarrow \infty$, the conditioning will have a negligible effect on the CPF-AS. Hence, $\tilde{\mathbf{x}}'_{0:T}$ will be effectively independent of $\tilde{\mathbf{x}}_{0:T}$ and (with an infinite number of particles) distributed according to the exact smoothing distribution. The number of particles N will thus affect the mixing of the Markov kernel M_θ^N . The invariance property of the kernel holds for any N , but the larger we take N , the smaller the correlation will be between $\tilde{\mathbf{x}}'_{0:T}$ and $\tilde{\mathbf{x}}_{0:T}$. However, it has been experienced in practice that the correlation drops off very quickly as N increases [2, 4], and for many models a moderate N is enough to obtain a rapidly mixing kernel.

3 An Additional plot for the Nonlinear System Benchmark

In Figure 1 we show the values of the hyper-parameters that are learnt during the experiment.

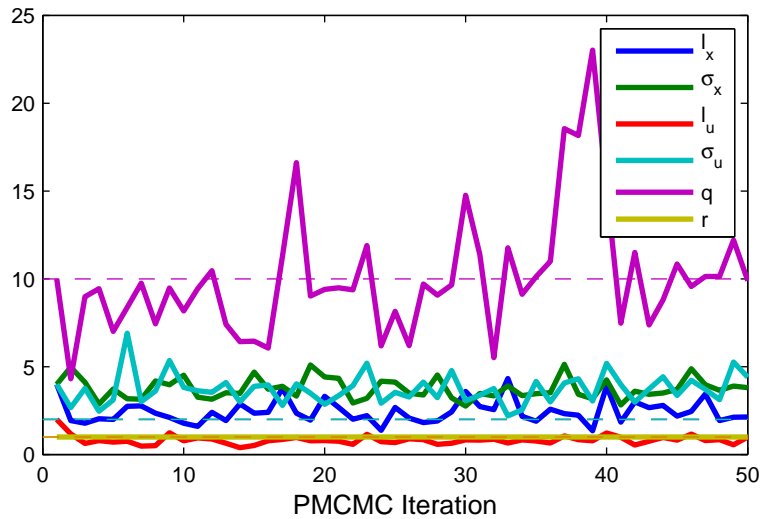


Figure 1: Hyper-parameter samples.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

References

- [1] P. Gill, G. Golub, W. Murray, and M. Saunders, “Methods for modifying matrix factorizations,” *Mathematics of Computation*, vol. 126, no. 28, pp. 505–535, 1974.
- [2] F. Lindsten, M. Jordan, and T. B. Schön, “Ancestor sampling for particle Gibbs,” in *Advances in Neural Information Processing Systems 25*, P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds., 2012, pp. 2600–2608.
- [3] C. Andrieu, A. Doucet, and R. Holenstein, “Particle Markov chain Monte Carlo methods,” *Journal of the Royal Statistical Society: Series B*, vol. 72, no. 3, pp. 269–342, 2010.
- [4] F. Lindsten and T. B. Schön, “On the use of backward simulation in the particle Gibbs sampler,” in *Proceedings of the 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Kyoto, Japan, Mar. 2012.