

# Supplementary Materials

## 1 Proof of Proposition 1

**Proof:** We first consider the fixed-directional jump. In the following, we calculate the transition probability  $\mathbb{P}((\boldsymbol{\theta}_j, \mathcal{M}_j) \mid (\boldsymbol{\theta}_i, \mathcal{M}_i))$ . Suppose the transition from  $(\boldsymbol{\theta}_i, \mathcal{M}_i)$  to  $(\boldsymbol{\theta}_j, \mathcal{M}_j)$  is achieved by some jumping distance  $r$ , then we have

$$\boldsymbol{\theta}_j = \mathbf{u} + r(\widehat{\boldsymbol{\theta}}_j - \widehat{\mathbf{u}}) \quad \text{and} \quad \mathbf{v} = \boldsymbol{\theta}_i + r(\widehat{\mathbf{v}} - \widehat{\boldsymbol{\theta}}_i).$$

It follows that

$$\begin{aligned} \mathbb{P}((\boldsymbol{\theta}_j, \mathcal{M}_j) \mid (\boldsymbol{\theta}_i, \mathcal{M}_i)) &= \int q_i(\mathbf{u}) \times \frac{p_j(\mathbf{v}, \boldsymbol{\theta}_j, \mathcal{M}_j \mid \mathbf{y})}{\sum_{k=1}^m p_j(\mathbf{v}^{(k)}, \boldsymbol{\theta}_j^{(k)}, \mathcal{M}_j \mid \mathbf{y})} \times \prod_{k=1}^{m-1} p(r^{(k)}) p(r) dr^{(1)} \dots dr^{(m-1)} \\ &\quad \times \min \left\{ 1, \frac{\sum_{k=1}^m p_j(\mathbf{v}^{(k)}, \boldsymbol{\theta}_j^{(k)}, \mathcal{M}_j \mid \mathbf{y})}{\sum_{k=1}^m p_i(\boldsymbol{\theta}_i^{(k)}, \mathbf{u}^{(k)}, \mathcal{M}_i \mid \mathbf{y})} \right\} \\ &= \int \min \left\{ \left[ \sum_{k=1}^m p_j(\mathbf{v}^{(k)}, \boldsymbol{\theta}_j^{(k)}, \mathcal{M}_j \mid \mathbf{y}) \right]^{-1}, \left[ \sum_{k=1}^m p_i(\boldsymbol{\theta}_i^{(k)}, \mathbf{u}^{(k)}, \mathcal{M}_i \mid \mathbf{y}) \right]^{-1} \right\} \\ &\quad \times \prod_{k=1}^{m-1} p(r^{(k)}) dr^{(1)} \dots dr^{(m-1)} \times p(\boldsymbol{\theta}_j, \mathcal{M}_j \mid \mathbf{y}) q_i(\mathbf{u}) q_j(\mathbf{v}) p(r), \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\theta}_j^{(k)} &= \mathbf{u} + r^{(k)}(\widehat{\boldsymbol{\theta}}_j - \widehat{\mathbf{u}}), \quad \mathbf{v}^{(k)} = \boldsymbol{\theta}_i + r^{(k)}(\widehat{\mathbf{v}} - \widehat{\boldsymbol{\theta}}_i), \\ \boldsymbol{\theta}_i^{(k)} &= \mathbf{v} - r^{(k)}(\widehat{\mathbf{v}} - \widehat{\boldsymbol{\theta}}_i), \quad \mathbf{u}^{(k)} = \boldsymbol{\theta}_j - r^{(k)}(\widehat{\boldsymbol{\theta}}_j - \widehat{\mathbf{u}}). \end{aligned}$$

Note that when we jump back from  $(\mathbf{v}, \boldsymbol{\theta}_j)$  to  $(\boldsymbol{\theta}_i, \mathbf{u})$ , we flip the sign of the jumping direction, that is, changing the jumping direction to  $(\widehat{\mathbf{u}} - \widehat{\boldsymbol{\theta}}_j, \widehat{\boldsymbol{\theta}}_i - \widehat{\mathbf{v}})$ , and keep the same jumping distance  $r$ . This is the reason why we do not require  $p(r)$  to be symmetric and centered at 0. Since the Jacobian between  $(\mathbf{v}, \boldsymbol{\theta}_j)$  and  $(\boldsymbol{\theta}_i, \mathbf{u})$  is simply 1, and the multiple integral is symmetric with respect to the index  $i$  and  $j$ , the transition kernel satisfies the detailed balance condition, thus leaves  $p(\boldsymbol{\theta}_k, \mathcal{M}_k \mid \mathbf{y})$  invariant.

For the adaptive-directional jump, we first standardize the jumping direction, that is, setting  $\mathbf{e} = (\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u}) / \|(\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u})\|$ , so that the jumping distance is independent to the current

state  $(\boldsymbol{\theta}_i, \mathbf{u})$ . Besides, because the jumping direction always points to the mode of the augmented posterior distribution  $p(\mathbf{v}, \boldsymbol{\theta}_j, M_j \mid \mathbf{y})$ , when we jump back from  $(\mathbf{v}, \boldsymbol{\theta}_j)$  to  $(\boldsymbol{\theta}_i, \mathbf{u})$ , we should flip the sign of the jumping distance as we won't flip the sign of the jumping direction. Consequently,  $p(r)$  is required to be symmetric and centered at 0.

The proof of the reversibility of the transition kernel equipped with the adaptive-directional jump follows similarly (thus is omitted) as the case of the fixed-directional jump, but requires an additional calculation of the Jacobian, which is detailed as below. Suppose the transition from  $(\boldsymbol{\theta}_i, \mathcal{M}_i)$  to  $(\boldsymbol{\theta}_j, \mathcal{M}_j)$  is achieved by some jumping distance  $r$ , then we have

$$\boldsymbol{\theta}_j = \mathbf{u} + r \frac{\widehat{\boldsymbol{\theta}}_j - \mathbf{u}}{\|(\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u})\|} \quad \text{and} \quad \mathbf{v} = \boldsymbol{\theta}_i + r \frac{\widehat{\mathbf{v}} - \boldsymbol{\theta}_i}{\|(\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u})\|}.$$

For notational convenience, we define  $\mathbf{x} \in \mathbb{R}^{d_i+d_j}$  as follows. For  $k \in [d_j]$  and  $l \in [d_i]$ , let

$$x_k = \frac{(\widehat{\theta}_{jk} - u_k)}{\|(\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u})\|^{3/2}}, \quad x_{d_j+l} = \frac{(\widehat{v}_l - \theta_{il})}{\|(\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u})\|^{3/2}}.$$

Then we have

$$\left| \frac{\partial(\boldsymbol{\theta}_j, \mathbf{v})}{\partial(\mathbf{u}, \boldsymbol{\theta}_i)} \right| = \det((1 - r\|\mathbf{x}\|^2)I_{d_j} + r\mathbf{x}\mathbf{x}^\top) = \left[ 1 - \frac{r}{\|(\widehat{\mathbf{v}} - \boldsymbol{\theta}_i, \widehat{\boldsymbol{\theta}}_j - \mathbf{u})\|} \right]^{d_i+d_j-1}.$$

## 2 More details of the Log-Gaussian Cox process

We provide more details of the Log-Gaussian Cox process discussed in Section 5.1. We first detail the adaptive SMC algorithm in Algorithm 1. We then examine the effect of the threshold  $c$  used in the flat histogram criterion. Figure 1 shows that the WL mixture method is robust to the choice of  $c$  in the region  $[0.1, 0.3]$ .

## 3 Variational approximation for the Bayesian Lasso

Algorithm 2 details the coordinate ascent variational inference (CAVI) algorithm used in constructing the surrogate distribution for the Bayesian Lasso example (see Section 5.2). We use the mean-field variational family, in which we assume

$$q(\beta_j) \sim N(m_j, s_j^2), \quad q(\eta_j) \sim N(\phi_j, \zeta_j^2) \quad \text{for } j \in [p], \quad \text{and} \quad q(\xi) \sim N(u, v^2).$$

The definition of  $\boldsymbol{\beta}, \boldsymbol{\eta}$  and  $\xi$  can be found in Section 5.2.

## References

Zhou, Y., A. M. Johansen, and J. A. D. Aston (2016). Toward automatic model comparison: an adaptive sequential Monte Carlo approach. *Journal of Computational and Graphical Statistics* 25(3), 701–726.

---

**Algorithm 1** An adaptive sequential Monte Carlo (SMC) sampler (Zhou et al., 2016).

---

**Input:** proposal distribution  $q(\boldsymbol{\theta})$ , Markov kernels  $\{K_t\}$ .

1. Initialization.

- (a) Sample  $\boldsymbol{\theta}_0^{(i)}$  from  $q(\boldsymbol{\theta})$  for  $i \in [n]$  independently.
- (b) Set  $w_0^{(i)} = 1/n$  for  $i \in [n]$ . Set  $\lambda_0 = 0$  and  $t = 0$ .

2. While  $\lambda_t < 1$ , iterate between the following steps.

- (a) Set  $t \leftarrow t + 1$ .
- (b) For some pre-specified  $\kappa \in (0, 1)$ , using binary search to find  $\lambda_t \in (\lambda_{t-1}, 1]$  such that

$$n^{-1}\text{CESS}_t(\lambda_t) = \frac{\left(\sum_{i=1}^n w_{t-1}^{(i)} (\gamma/q)(\boldsymbol{\theta}_{t-1}^{(i)})^{\lambda_t - \lambda_{t-1}}\right)^2}{\sum_{i=1}^n w_{t-1}^{(i)} (\gamma/q)(\boldsymbol{\theta}_{t-1}^{(i)})^{2(\lambda_t - \lambda_{t-1})}} = \kappa.$$

If  $n^{-1}\text{CESS}_t(1) > \kappa$ , set  $\lambda_t = 1$  and  $T = t$ .

- (c) Compute the unnormalized weights  $w_t^{(i)} = (\gamma_t/\gamma_{t-1})(\boldsymbol{\theta}_{t-1}^{(i)})$  for  $i \in [n]$ . The geometric sequence of the auxiliary distributions  $\gamma_t(\boldsymbol{\theta})$  is defined as

$$\gamma_0(\boldsymbol{\theta}) = q(\boldsymbol{\theta}), \quad \gamma_T(\boldsymbol{\theta}) = \gamma(\boldsymbol{\theta}), \quad \text{and} \quad \gamma_t(\boldsymbol{\theta}) = q(\boldsymbol{\theta})^{\lambda_t} \gamma(\boldsymbol{\theta})^{1-\lambda_t}.$$

- (d) Compute  $r_t = \sum_{i=1}^n w_{t-1}^{(i)} w_t^{(i)}$ , which is an estimate of  $Z_t/Z_{t-1}$ .  $Z_t$  is the normalizing constant of  $\gamma_t$ .
- (e) Normalize the weights  $\{w_t^{(i)}\}_{i \in [n]}$  to sum 1.
- (f) Compute the (normalized) effective sample size

$$\text{ESS}_t = \frac{1}{n \sum_{i=1}^n (w_t^{(i)})^2}.$$

If  $\text{ESS}_t \leq 0.5$ , resample particles using systematic resampling, and set  $w_t^{(i)} = 1/n$  for  $i \in [n]$ .

- (g) Move particles  $\boldsymbol{\theta}_t^{(i)}$  according to the Markov kernel  $K_t$  for  $i \in [n]$ .

**Output:** normalizing constant estimate  $\hat{Z}_\gamma = \prod_{t=1}^T r_t$ .

---

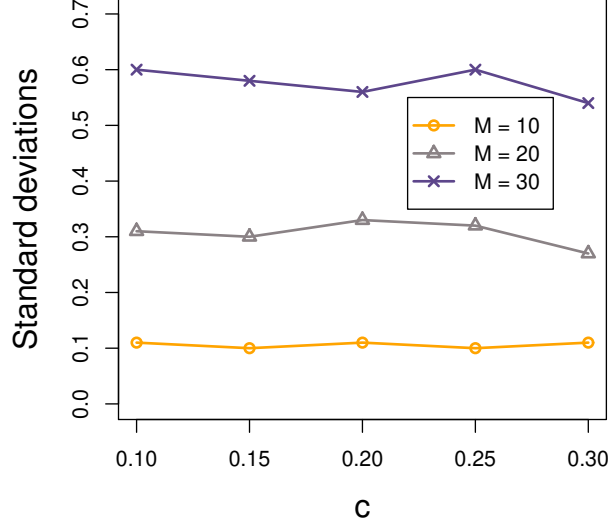


Figure 1: Demonstration of the effect of the threshold  $c$  used in the flat histogram criterion for the Log-Gaussian Cox process. The  $y$ -axis denotes the empirical standard deviation of the log normalizing constant estimates over 10 independent runs. The unit square is discretized into an  $M \times M$  regular grid (see Section 5.1).

---

**Algorithm 2** The CAVI updates for the Bayesian Lasso example.

---

1. Given  $q(\beta_i)$  for  $i \neq j$ ,  $q(\boldsymbol{\eta})$  and  $q(\xi)$ , update  $q(\beta_j)$ .

$$(m_j, s_j) = \arg \max_{(m, s > 0)} \left\{ e^{-u+v^2/2} \left[ A_{-j}m - \frac{1}{2} \left( (X^\top X)_{jj} + e^{-\phi_j + \zeta_j^2/2} \right) (m^2 + s^2) \right] + \frac{1}{2} \log s^2 \right\},$$

where  $A_{-j} = \sum_{k=1}^n X_{kj}y_k - \sum_{k \neq j} (X^\top X)_{kj}m_k$ .

2. Given  $q(\eta_i)$  for  $i \neq j$ ,  $q(\boldsymbol{\beta})$  and  $q(\xi)$ , update  $q(\eta_j)$ .

$$(\phi_j, \zeta_j) = \arg \max_{(\phi, \zeta > 0)} \left\{ \phi - \lambda^2 e^{\phi + \zeta^2/2} - (m_j^2 + s_j^2) e^{-u+v^2/2} e^{-\phi + \zeta^2/2} + \log \zeta^2 \right\}.$$

3. Given  $q(\boldsymbol{\beta})$ ,  $q(\boldsymbol{\eta})$ , update  $q(\xi)$ .

$$(u, v) = \arg \max_{(u, v > 0)} \left\{ -(n+p)u - B e^{-u+v^2/2} + \log v^2 \right\},$$

where  $B = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top X\mathbf{m} + \text{tr}[(\mathbf{m}\mathbf{m}^\top + \text{diag}(s_k^2))X^\top X] + \sum_{k=1}^p (m_k^2 + s_k^2) e^{-\phi_k + \zeta_k^2/2}$ .

---