

# Supplementary Material to Bayesian two-part spatial models for semicontinuous data with application to emergency department expenditures

BRIAN NEELON\*

*Department of Public Health Sciences, Medical University of South Carolina, 135 Cannon Street  
Suite 303 MSC 835, Charleston, SC 29425 USA*

neelon@musc.edu

LI ZHU

*Department of Biostatistics, University of Pittsburgh, 130 De Soto Street, Pittsburgh,  
Pennsylvania 15261, USA*

SARA E. BENJAMIN NEELON

*Department of Community and Family Medicine, Duke University School of Medicine, 2200 W.  
Main Street, Durham, North Carolina 27705, USA*

\*To whom correspondence should be addressed.

**Table S1:** *Empirical Kullback-Leibler divergence estimates comparing true and model-estimated densities for four selected block groups in the illustrative example.*

Block Group	LN	LSN	3-df LST	16-df LST	DPLN- $\omega_5$
Average	1.518	0.100	0.076	0.093	0.056
3	1.695	0.115	0.077	0.106	0.065
35	2.172	0.215	0.132	0.223	0.085
148	1.685	0.075	0.059	0.073	0.043

**Table S2:** *Summary statistics for DSR data ( $N=44722$ ).*

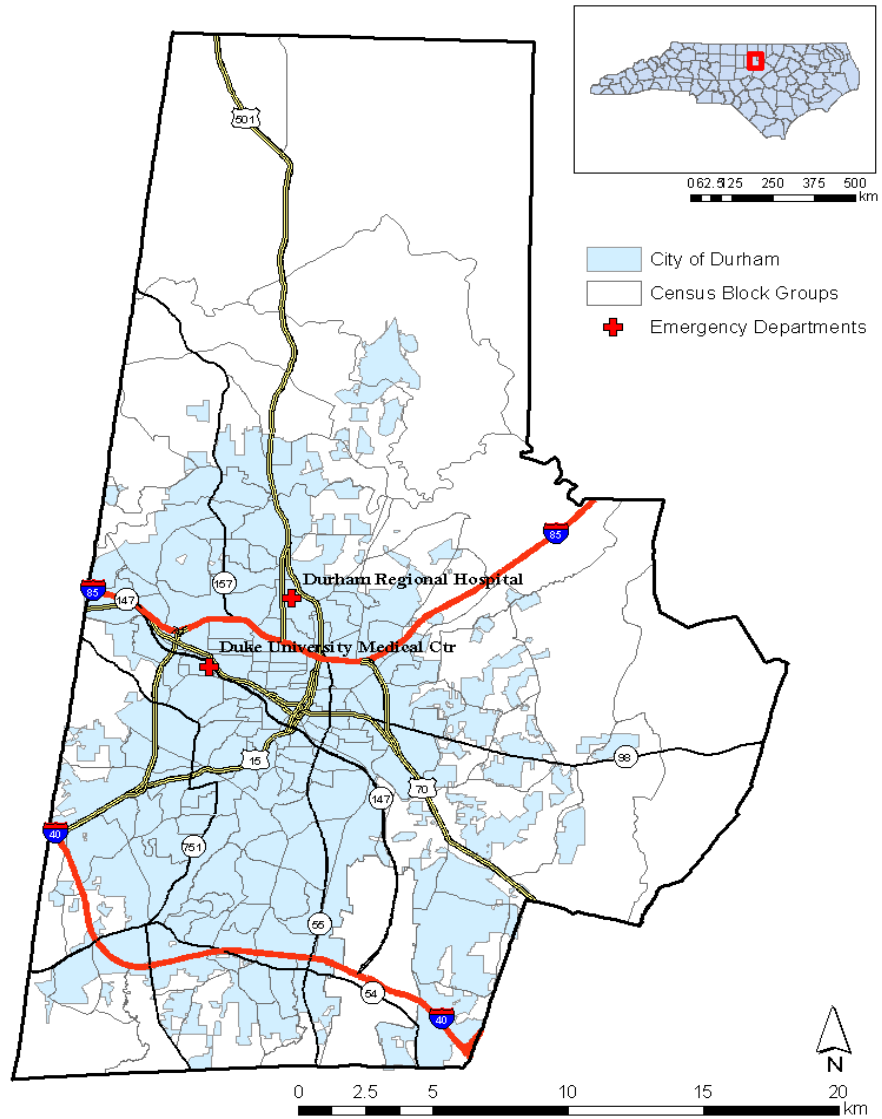
Variable	$n$	%
ED users	11083	25
Male gender	16425	36
Race		
Non-Hispanic white	22986	51
Non-Hispanic black	19432	44
Hispanic	2304	5
Insurance		
Private	24044	54
Medicare	8669	19
Medicaid	3909	9
Uninsured	8100	18
	Median	Range
Age (years)	47	(18, 105)
Block group sample size	267	(1, 803)
Block group median household income	45640	(5934, 157300)

**Table S3:** *DPLN-estimated  $\Pr(Y > 0)$  and  $E(Y|Y > 0)$  for the four block groups highlighted in Figure 3. 95% credible intervals are given in parentheses.*

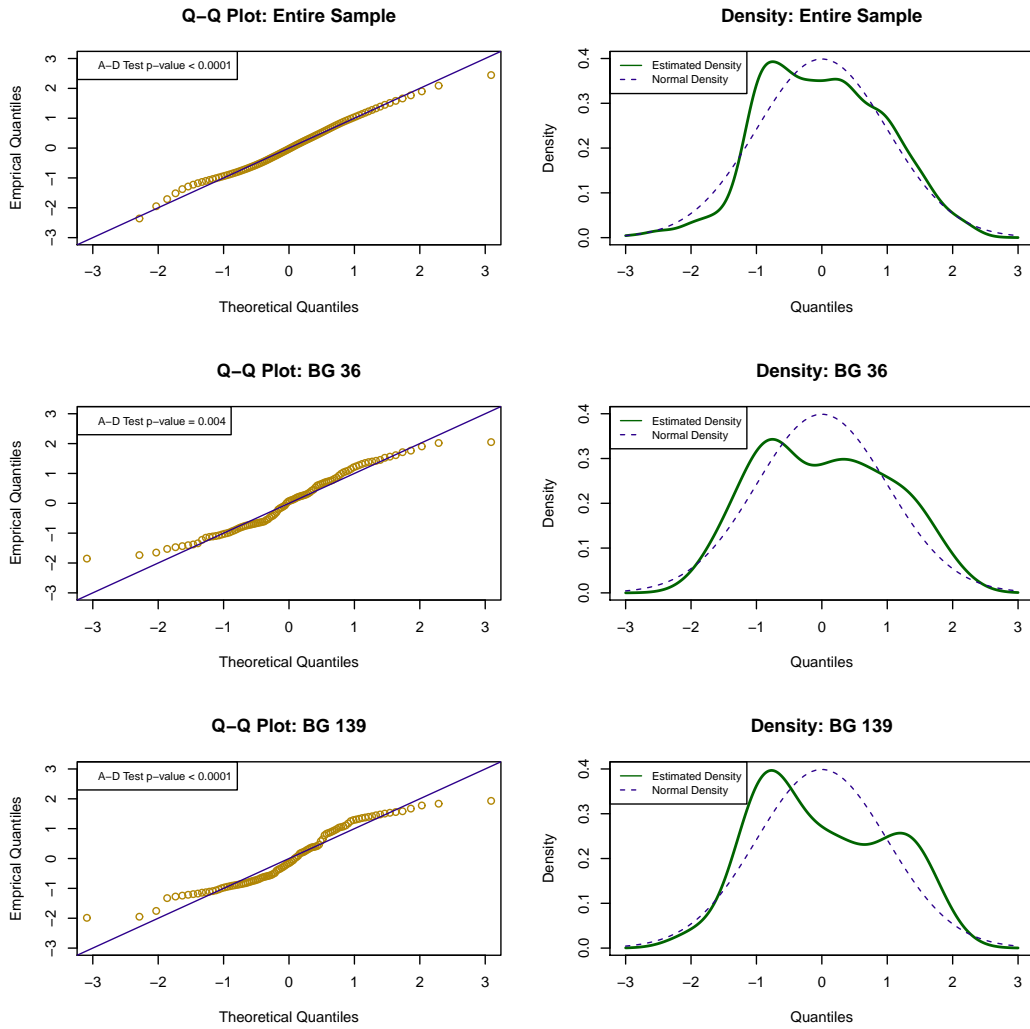
Block Group	$\Pr(Y > 0)$	$E(Y Y > 0)$
84	0.082 (0.061, 0.104)	3182 (2720, 3809)
86	0.081 (0.057, 0.108)	3713 (3082, 4315)
122	0.037 (0.025, 0.054)	3175 (2482, 3935)
24	0.028 (0.017, 0.046)	3910 (3052, 4905)

**Table S4:** *DPLN- and LST-estimated median and 90th quantiles ( $Q_{90}$ ) for block groups 24 and 122. 95% credible intervals are given in parentheses.*

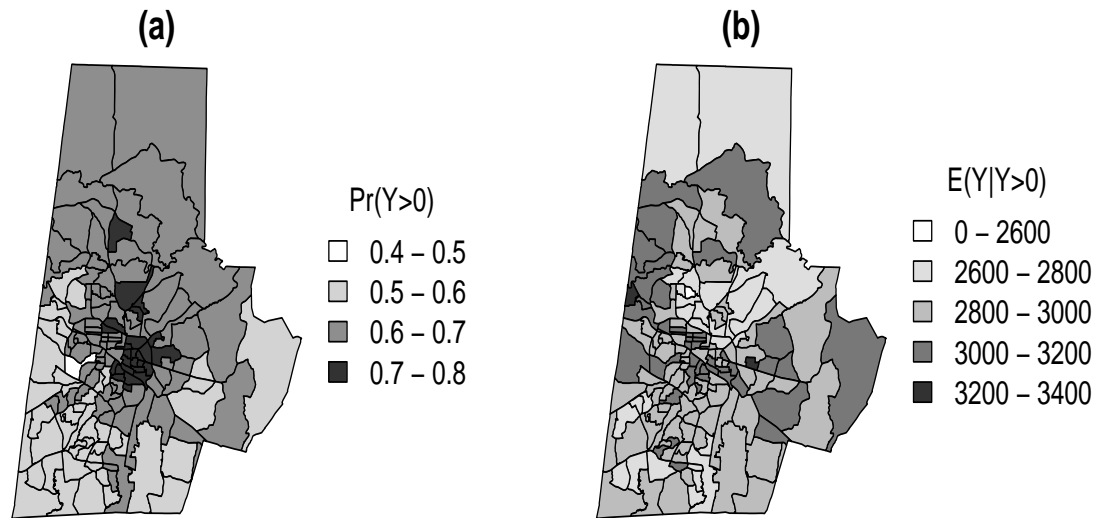
Model	Block Group	Median	$Q_{90}$
DPLN ( $\omega = 1$ )	24	2760 (2160, 3720)	8300 (6320, 10962)
	122	2200 (1760, 2821)	6662 (5040, 8642)
LST ( $\nu = 16$ )	24	3557 (2927, 4272)	10393 (8555, 12514)
	122	2819 (2324, 3436)	8238 (6786, 10067)



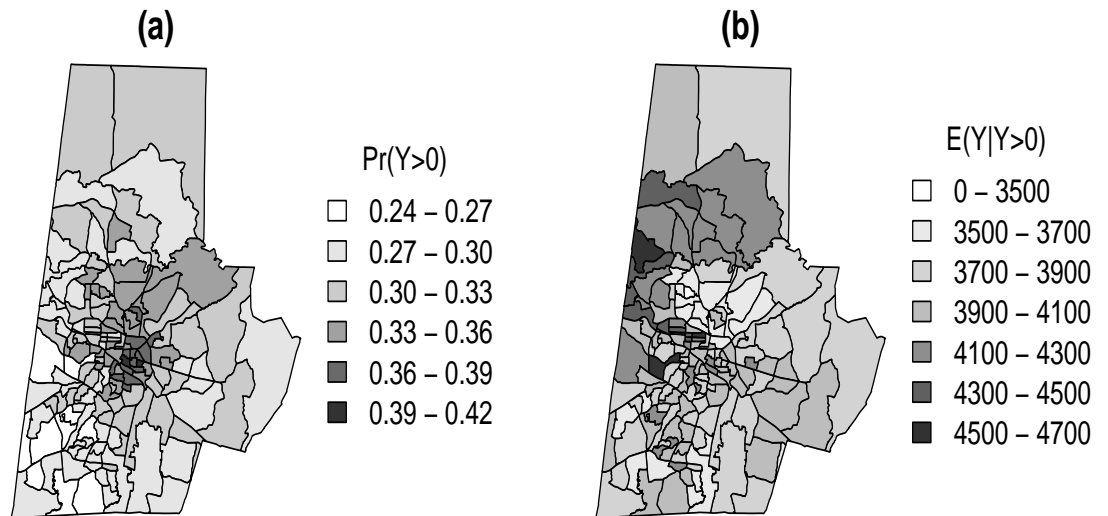
**Figure S1:** Map of Durham County block groups, city limits, and emergency department locations.



**Figure S2:** Residual diagnostics on log scale from an LN spatial model fitted to the nonzero expenditures in the ED application. Left panels: Q-Q plots of standardized residuals for the entire sample and for two select block groups. Right panels: Estimated residual densities with standard normal densities provided as reference. P-values are from Anderson-Darling (A-D) tests of normality. BG: Block Group.



**Figure S3:** DPLN-estimated (a)  $\Pr(Y > 0)$  and (b)  $E(Y|Y > 0)$  for patient group 2 from the ED expenditure analysis.



**Figure S4:** DPLN-estimated (a)  $\Pr(Y > 0)$  and (b)  $E(Y|Y > 0)$  for patient group 3 from the ED expenditure analysis.

## MCMC ALGORITHMS

*MCMC Updates for Two-Part Lognormal Model*

1. Update  $U_{ij}$ . Assuming a probit link for the binary part in equation (2.2), the full conditional for the latent variable  $U_{ij}$  ( $i = 1, \dots, n; j = 1, \dots, n_i$ ) is

$$U_{ij} | - \sim \text{TN}_{[0, +\infty)}(\mathbf{x}'_{ij} \boldsymbol{\gamma} + s_1(v_{ij}) + b_{1i}, 1), \quad \text{if } y_{ij} > 0$$

$$U_{ij} | - \sim \text{TN}_{(-\infty, 0]}(\mathbf{x}'_{ij} \boldsymbol{\gamma} + s_1(v_{ij}) + b_{1i}, 1), \quad \text{if } y_{ij} = 0.$$

2. Update  $\boldsymbol{\gamma}$  and the spline coefficients associated with  $s_1(v_{ij})$ . For brevity, let  $\boldsymbol{\gamma}$  denote a  $p \times 1$  vector of fixed effect *and* spline coefficients, where  $p$  denotes the total number of coefficients. Assuming a joint  $N_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  prior, the joint full conditional for  $\boldsymbol{\gamma}$  is given by  $N(\boldsymbol{\mu}_\boldsymbol{\gamma}, \boldsymbol{\Sigma}_\boldsymbol{\gamma})$ , where

$$\boldsymbol{\Sigma}_\boldsymbol{\gamma} = (\boldsymbol{\Sigma}_0^{-1} + \mathbf{X}'\mathbf{X})^{-1} \quad \text{and}$$

$$\boldsymbol{\mu}_\boldsymbol{\gamma} = \boldsymbol{\Sigma}_\boldsymbol{\gamma} \{ \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'(\mathbf{u} - \mathbf{D}\mathbf{b}_1) \}.$$

Here,  $\mathbf{u}$  is an  $N \times 1$  vector of sampled values,  $u_{ij}$ , from the previous step,  $\mathbf{X}$  is an  $N \times p$  design matrix combining fixed effects covariates and spline basis functions, and  $\mathbf{D}$  is a  $N \times n$  random effects design matrix for the  $n \times 1$  vector  $\mathbf{b}_1$ .

3. Update  $\mathbf{b}_1$ . Given the conditional prior specification in equation (3.3), the full conditional for  $\mathbf{b}_1$  is  $N(\boldsymbol{\mu}_{\mathbf{b}_1}, \boldsymbol{\Sigma}_{\mathbf{b}_1})$ , where

$$\boldsymbol{\Sigma}_{\mathbf{b}_1} = \left\{ \frac{\mathbf{Q}}{(1 - \rho^2)\Lambda_{11}} + \mathbf{D}'\mathbf{D} \right\}^{-1} \quad \text{and}$$

$$\boldsymbol{\mu}_{\mathbf{b}_1} = \boldsymbol{\Sigma}_{\mathbf{b}_1} \left\{ \frac{\rho\mathbf{Q}\mathbf{b}_2}{(1 - \rho^2)\sqrt{\Lambda_{11}\Lambda_{22}}} + \mathbf{D}'(\mathbf{u} - \mathbf{X}\boldsymbol{\gamma}) \right\}.$$

Note that this update relies on data ( $\mathbf{u}$  and  $\mathbf{X}$ ) from all  $N$  observations.

4. Update  $\boldsymbol{\beta}$  and the spline coefficients associated with  $s_2(v_{ij})$ . As with  $\boldsymbol{\gamma}$  above, let  $\boldsymbol{\beta}$  denote a  $p \times 1$  vector of fixed effect and spline coefficients. Assuming an  $N_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  prior, the full

conditional for  $\beta$  is  $N(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta)$ , where

$$\boldsymbol{\Sigma}_\beta = \left( \boldsymbol{\Sigma}_0^{-1} + \frac{\mathbf{X}'_1 \mathbf{X}_1}{\sigma^2} \right)^{-1} \quad \text{and}$$

$$\boldsymbol{\mu}_\beta = \boldsymbol{\Sigma}_\beta \left\{ \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \frac{\mathbf{X}'_1 (\ln \mathbf{y}_1 - \mathbf{D}_1 \mathbf{b}_2)}{\sigma^2} \right\}.$$

Here,  $\mathbf{y}_1$  is an  $N_1 \times 1$  vector of nonzero expenditures, with  $N_1$  denoting the number of nonzero observations;  $\mathbf{X}_1$  is an  $N_1 \times p$  design matrix containing fixed effects covariates and spline bases for the nonzero observations; and  $\mathbf{D}_1$  is the corresponding  $N_1 \times n$  random effects design matrix for the nonzero observations.

5. Update  $\sigma^2$ . Assuming an  $IG(g, s)$  prior for  $\sigma^2$ , the full conditional is  $IG(g^*, s^*)$ , where  $g^* = g + \frac{N_1}{2}$ , and  $s^* = s + \frac{1}{2} \{ \ln \mathbf{y}_1 - \mathbf{X}_1 \boldsymbol{\beta} - \mathbf{D}_1 \mathbf{b}_2 \}' \{ \ln \mathbf{y}_1 - \mathbf{X}_1 \boldsymbol{\beta} - \mathbf{D}_1 \mathbf{b}_2 \}$ .
6. Update  $\mathbf{b}_2$ . Given the conditional prior in equation (3.4), update  $\mathbf{b}_2$  from its  $N(\boldsymbol{\mu}_{\mathbf{b}_2}, \boldsymbol{\Sigma}_{\mathbf{b}_2})$  full conditional, where

$$\boldsymbol{\Sigma}_{\mathbf{b}_2} = \left\{ \frac{\mathbf{Q}}{(1 - \rho^2) \Lambda_{22}} + \frac{\mathbf{D}'_1 \mathbf{D}_1}{\sigma^2} \right\}^{-1} \quad \text{and}$$

$$\boldsymbol{\mu}_{\mathbf{b}_2} = \boldsymbol{\Sigma}_{\mathbf{b}_2} \left\{ \frac{\rho \mathbf{Q} \mathbf{b}_2}{(1 - \rho^2) \sqrt{\Lambda_{11} \Lambda_{22}}} + \frac{\mathbf{D}'_1 (\ln \mathbf{y}_1 - \mathbf{X}_1 \boldsymbol{\beta})}{\sigma^2} \right\}.$$

Note that this update relies only on data ( $\mathbf{y}_1$  and  $\mathbf{X}_1$ ) for the  $N_1$  positive observations.

7. Update  $\boldsymbol{\Lambda}$ . Assuming an  $IW(\kappa_0, \mathbf{S}_0)$  prior, the full conditional for  $\boldsymbol{\Lambda}$  is  $IW(\kappa^*, \mathbf{S}^*)$ , where  $\kappa^* = n + \kappa_0 - 1$ , and  $\mathbf{S}^* = \mathbf{S}_0 + \mathbf{b}^{*'} \mathbf{Q} \mathbf{b}^*$ . Here,  $\mathbf{b}^*$  denotes an  $n \times 2$  matrix with first column equal to  $\mathbf{b}_1$  and second column equal to  $\mathbf{b}_2$ . Use  $\boldsymbol{\Lambda}$  to extract  $\Lambda_{11}$ ,  $\Lambda_{22}$  and  $\rho$ .
8. Apply sum-to-zero constraints to  $\mathbf{b}_1$  and  $\mathbf{b}_2$ .

#### *MCMC Updates for Two-Part Log Skew-Normal Model*

For the two-part LSN model, the updates of the binary part are similar to those given for the two-part LN model. The continuous part is updated as follows:



1. Update the latent variable  $Z_{ij}$  in equation (2.4). If we assume a truncated normal  $\text{TN}_{[0,\infty)}(0, 1)$  prior for all  $Z_{ij}$  such that  $y_{ij} > 0$ , the full conditional for  $Z_{ij}$  is  $N(\mu_{z_{ij}}, \sigma_{z_{ij}}^2)$ , where

$$\sigma_{z_{ij}}^2 = \frac{\xi^2}{\psi^2 + \xi^2} \quad \text{and}$$

$$\mu_{z_{ij}} = \sigma_{z_{ij}}^2 \psi \{ \ln y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta} - b_{2i} \} / \xi^2.$$

Concatenate the sampled values,  $z_{ij}$ , to form an  $N_1 \times 1$  vector  $\mathbf{z}_1$ .

2. Update  $\psi$ ,  $\boldsymbol{\beta}$  and the spline coefficients associated with  $s_2(v_{ij})$ . Combine the  $N_1 \times p$  matrix  $\mathbf{X}_1$  defined earlier with the vector  $\mathbf{z}_1$  columnwise to form an  $N_1 \times (p + 1)$  matrix  $\mathbf{X}_1^*$  that now includes fixed effect covariates, spline basis functions, and the latent vector  $\mathbf{z}_1$ . Likewise, combine the vector  $\boldsymbol{\beta}$  defined earlier and the parameter  $\psi$  defined in equation (2.4) to form a  $(p + 1) \times 1$  vector  $\boldsymbol{\beta}^*$  that includes fixed and spline coefficients as well as  $\psi$ . Assigning a joint prior  $N_{p+1}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ , the full conditional for  $\boldsymbol{\beta}^*$  is  $N(\boldsymbol{\mu}_{\boldsymbol{\beta}^*}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}^*})$ , where

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}^*} = \left( \boldsymbol{\Sigma}_0^{-1} + \frac{\mathbf{X}_1^{*'} \mathbf{X}_1^*}{\sigma^2} \right)^{-1} \quad \text{and}$$

$$\boldsymbol{\mu}_{\boldsymbol{\beta}^*} = \boldsymbol{\Sigma}_{\boldsymbol{\beta}^*} \left\{ \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \frac{\mathbf{X}_1^{*'} (\ln \mathbf{y}_1 - \mathbf{D}_1 \mathbf{b}_2)}{\sigma^2} \right\}.$$

3. Update  $\xi^2$ ,  $\mathbf{b}_2$  and  $\boldsymbol{\Lambda}$ . Updates for  $\xi^2$ ,  $\mathbf{b}_2$  and  $\boldsymbol{\Lambda}$  are identical to updating  $\sigma^2$ ,  $\mathbf{b}_2$  and  $\boldsymbol{\Lambda}$  in the LN model if we replace  $\mathbf{X}_1$  and  $\boldsymbol{\beta}$  with  $\mathbf{X}_1^*$  and  $\boldsymbol{\beta}^*$ , respectively. As before, apply sum-to-zero constraints to  $\mathbf{b}_1$  and  $\mathbf{b}_2$ .

#### MCMC Updates for Two-Part Log Skew-t Model

The updates of the binary part of the two-part LST are also similar to those given in two-part LN model. The updates for the continuous part are similar to those for the LSN, but with the following modifications:

1. Update  $Z_{ij}$  in equation (2.6). Given a truncated normal  $\text{TN}_{[0,\infty)}(0, 1/w_{ij})$  prior for all  $Z_{ij}$

such that  $y_{ij} > 0$ , the full conditional of  $Z_{ij}$  is  $N(\mu_{z_{ij}}, \sigma_{z_{ij}}^2)$ , where

$$\sigma_{z_{ij}}^2 = \frac{\xi^2}{(\psi^2 + \xi^2)w_{ij}} \quad \text{and}$$

$$\mu_{z_{ij}} = \sigma_{z_{ij}}^2 \psi w_{ij} \{ \ln y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta} - b_{2i} \} / \xi^2.$$

2. Update  $W_{ij}$ . Assuming a  $\text{Ga}(\frac{\nu}{2}, \frac{\nu}{2})$  prior, for all  $y_{ij} > 0$ , the full conditional for  $W_{ij}$  is:

$$\text{Ga} \left( \frac{\nu}{2} + 1, \quad \frac{\nu}{2} + \frac{\{ \ln y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta} - b_{2i} \}^2}{2\xi^2} + \frac{z_{ij}^2}{2} \right).$$

Concatenate the sampled values,  $w_{ij}$ , to form an  $N_1 \times 1$  vector  $\mathbf{w}_1$ .

3. Update other parameters. Define  $\widetilde{\mathbf{X}}_1 = \mathbf{X}_1^* \sqrt{\mathbf{w}_1}$ ,  $\widetilde{\mathbf{D}}_1 = \mathbf{D}_1 \sqrt{\mathbf{w}_1}$ , and  $\widetilde{\mathbf{y}}_1 = \ln(\mathbf{y}_1) \sqrt{\mathbf{w}_1}$ , where  $\mathbf{X}_1^*$  and  $\mathbf{D}_1$  are defined earlier. Updates for  $\boldsymbol{\beta}^*$ ,  $\xi^2$ ,  $\mathbf{b}_2$  and  $\boldsymbol{\Lambda}$  are identical to updating  $\boldsymbol{\beta}^*$ ,  $\xi^2$ ,  $\mathbf{b}_2$  and  $\boldsymbol{\Lambda}$  in the LSN model if we replace  $\mathbf{X}_1^*$ ,  $\mathbf{D}_1$  and  $\mathbf{y}_1$  with  $\widetilde{\mathbf{X}}_1$ ,  $\widetilde{\mathbf{D}}_1$  and  $\widetilde{\mathbf{y}}_1$ , respectively. Apply sum-to-zero constraints as before.

#### *MCMC Updates for DP Mixture of Two-Part Lognormal Models*

For the DPLN, we follow the stick-breaking construction outlined in Section 2.4 and iterate through following steps:

1. Update the component index  $C_{ij}$ . For all  $i$  and  $j$ , sample  $C_{ij}$  from its multinomial full conditional:

$$\text{Pr}(C_{ij} = k | -) = \frac{\{v_k \prod_{l < k} (1 - v_l)\} f(y_{ij} | C_{ij} = k)}{\sum_{r=1}^K \{v_r \prod_{s < r} (1 - v_s)\} f(y_{ij} | C_{ij} = r)},$$

where  $f(y_{ij} | C_{ij} = k)$  is defined in equation (2.10) and  $K$  is the total number of mixture components after truncation.

2. Update the stick breaking weights  $v_k$ . Assuming a  $\text{Be}(1, \omega)$  prior, the full conditional for  $v_k$  ( $k = 1, \dots, K$ ) is:

$$v_k | - \sim \text{Be} \left( 1 + \sum_{i=1}^n \sum_{j=1}^{n_i} \mathbb{1}(C_{ij} = k), \quad \omega + \sum_{i=1}^n \sum_{j=1}^{n_i} \mathbb{1}(C_{ij} > k) \right).$$

3. Conditional on  $C_{ij} = k$ , update the latent variable  $U_{ij}$  from its truncated normal full conditional, which is analogous the update for the LN model but with component- $k$  likelihood and parameters.
4. For  $k = 1, \dots, K$ , update the component-specific parameters  $\gamma_k, \beta_k, \mathbf{b}_{1k}, \mathbf{b}_{2k}, \mathbf{\Lambda}_k$  and  $\sigma_k^2$  following full conditionals similar to those for the LN model, but now with component-specific likelihood and parameters. For all  $k$ , apply sum-to-zero constraints to  $\mathbf{b}_{1k}$  and  $\mathbf{b}_{2k}$ .