

Supplemental Material for “Hierarchical Model Building, Fitting, and Checking: A Behind-the-Scenes Look at a Bayesian Analysis of Arsenic Exposure Pathways”

Full Conditional Distributions

This document provides the full conditional distributions for the models introduced in the paper.

Letting $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\omega})$, the posterior distribution is

$$\begin{aligned} [\boldsymbol{\theta}, \mathbf{X} | \mathbf{Y}] &\propto [\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}] [\mathbf{X} | \boldsymbol{\theta}] [\boldsymbol{\theta}] \\ &= [\mathbf{Y} | \mathbf{X}, \boldsymbol{\omega}] [\mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\tau}] [\boldsymbol{\beta}] [\boldsymbol{\omega}] [\boldsymbol{\tau}], \end{aligned}$$

assuming independence of the four components within $\boldsymbol{\theta}$. The distribution is not available in closed form and so we use a Markov chain Monte Carlo algorithm (hybrid Gibbs sampler with Metropolis steps within Gibbs) to sample from it.

Before providing the full conditional distributions, we first note that for a generic parameter ψ , the quantities a^ψ and b^ψ denote the hyperprior parameters for gamma prior distributions and m^ψ and V^ψ are the hyperprior parameters for normal prior distributions. These hyperparameter parameters were not introduced in the manuscript for the sake of brevity.

1 The local environment to biomarker (LEB) model

1. For each j and i , we draw augmented data, Y_{ji}^{M*} , as follows:

$$Y_{ji}^{M*} \begin{cases} = Y_{ji}^M, & Z_{ji}^M = 0; \\ \sim \text{Truncated } N(X_{ji}^M, 1/\omega_j^M) \text{ on } (-\infty, M_{ji}^M], & Z_{ji}^M = 1; \\ \sim N(X_{ji}^M, 1/\omega_j^M), & Z_{ji}^M = 2. \end{cases} \quad (1)$$

2. For the process variable \mathbf{X}^M , we sample each \mathbf{X}_j^M , $j = 1, \dots, N^M$, conditional on the other process variables $\{\mathbf{X}_k^M : k \neq j\}$ and the other parameters in the model. We have that

$$\begin{aligned} &[\mathbf{X}_j^M | \{\mathbf{X}_k^M : k \neq j\}, \boldsymbol{\mu}^M, \boldsymbol{\beta}^M, \boldsymbol{\tau}^M, \mathbf{Y}^{M*}, \boldsymbol{\omega}^M] \\ &\propto [\mathbf{Y}_j^* | \mathbf{X}_j^M, \boldsymbol{\omega}^M] [\mathbf{X}_j^M | \{\mathbf{X}_k^M : k \neq j\}, \boldsymbol{\beta}_j^M, \tau_j^M] \prod_{l \in \chi_j} [\mathbf{X}_l^M | \{\mathbf{X}_k^M : k \neq l\}, \boldsymbol{\beta}_l^M, \tau_l^M], \quad (2) \end{aligned}$$

where χ_j denotes the columns of the process variable \mathbf{X}^M that depend on \mathbf{X}_j^M in the specification of the conditional mean η_{ji}^M (i.e., those l that belong to S_j^M). This conditional mean η_{ji}^M is defined as:

$$\eta_{ji}^M = \mu_j^M + \sum_{k=1}^{N_j^M} \beta_{jk}^M X_{s_{jk}i}^M, \quad (3)$$

where $\mathbf{s}_j^M = (s_{j1}^M, \dots, s_{jN_j^M}^M)$ is a vector of length N_j^M (the cardinality of S_j^M) containing the columns of the process variable matrix \mathbf{X}^M that affect the mean of the process variable \mathbf{X}_j^M .

The expression (2) above is equal to

$$\prod_{i=1}^{N^I} \left([Y_{ji}^{M*} | X_{ji}^M, \omega_j^M] [X_{ji}^M | \{X_{ki}^M : k \neq j\}, \beta_j^M, \tau_j^M] \prod_{l \in \chi_j} [X_{li}^M | \{X_{ki}^M : k \neq l\}, \beta_l^M, \tau_l^M] \right),$$

from which we can see that each X_{ji}^M , $i = 1, \dots, N^I$, can be sampled independently. Now

$$\begin{aligned} [Y_{ji}^{M*} | X_{ji}^M, \omega_j^M] &\propto \exp\left(\frac{\omega_j^M}{2}(Y_{ji}^{M*} - X_{ji}^M)^2\right), \\ [X_{ji}^M | \{X_{ki}^M : k \neq j\}, \beta_j^M, \tau_j^M] &\propto \exp\left(\frac{\tau_j^M}{2}(X_{ji}^M - \eta_{ji}^M)^2\right), \end{aligned}$$

and for each $l \in \chi_j$,

$$[X_{li}^M | \{X_{ki}^M : k \neq l\}, \beta_l^M] \propto \exp\left(\frac{\tau_l^M}{2}(X_{li}^M - \delta_{li}(j) - \epsilon_{li}(j))^2\right). \quad (4)$$

In (4),

$$\delta_{ji}(n) \equiv \mu_j^M + \sum_{\{k:s_{jk}^M \neq n\}} \beta_{jk}^M X_{s_{jk}i}^M,$$

is the conditional mean of X_{ji}^M minus the regression terms that depend on X_{ni}^M , and

$$\epsilon_{ji}(n) \equiv \sum_{\{k:s_{jk}^M = n\}} \beta_{jk}^M X_{s_{jk}i}^M, \quad (5)$$

is the regression terms that depend on X_{ni}^M . Always, $\delta_{ji}(n) - \epsilon_{ji}(n) = \eta_{ji}^M$.

Rearranging the terms of X_{ji}^M , the distribution of X_{ji}^M conditional on all other parameters is $N(q/p, 1/p)$, where

$$p = \omega_j^M + \tau_j^M + \sum_{l \in \chi_j} \tau_l^M \left(\sum_{\{k:s_{jk}^M = n\}} \beta_{jk}^M \right)^2,$$

and

$$q = \omega_j^M Y_{ji}^{M*} + \tau_j^M \eta_{ji}^M + \sum_{l \in \mathcal{X}_j} \tau_l^M (X_{li}^M - \delta_{li}(j)) \left(\sum_{\{k: s_{jk}^M = n\}} \beta_{jk}^M \right).$$

3. For the regression parameters β^M , we sample each β_j^M conditional on the other parameters in the model. Letting

$$\mathbf{A}_j = \begin{bmatrix} 1 & \left\{ X_{s_{jk}^M 1}^M \right\} \\ 1 & \left\{ X_{s_{jk}^M 2}^M \right\} \\ \vdots & \vdots \\ 1 & \left\{ X_{s_{jk}^M N^I}^M \right\} \end{bmatrix},$$

the full conditional distribution of β_j^M , is

$$\text{N}_{r_j}(\mathbf{R}_j^{-1} \mathbf{q}_j, \mathbf{R}_j^{-1}),$$

where $\mathbf{R}_j = \tau_j^M \mathbf{A}_j^T \mathbf{A}_j + \Sigma_j^{-1}$ and $\mathbf{q}_j = \tau_j^M \mathbf{A}_j^T \mathbf{X}_j^M + \Sigma_j^{-1} \mathbf{m}_j$.

4. For each j , the conditional distribution of the process precision τ_j^M , given the other parameters, is the gamma distribution,

$$\text{Ga} \left(a_j^M + N^I/2, b_j^M + \sum_{i=1}^{N^I} (X_{ji}^M - \eta_{ji}^M)^2/2 \right),$$

where η_{ji}^M is defined in (3) above.

2 The global to local environment (GLE) model

In this section, we provide the full conditional distributions for the parameters in the GLE models. Linking global environmental-media variables to local environmental-media variables will change the full conditional distributions for some of the latent processes and the associated regression coefficients that were given in the previous section. Specifically, now the full conditional distributions of $\mathbf{X}_{m(W)}^M$ depends on \mathbf{X}^W , and $\mathbf{X}_{m(S)}^M$ depends on \mathbf{X}^T and \mathbf{X}^H . The full conditional distributions for other media in the LEB model remain unchanged.

2.1 Global water model

1. With the addition of the global water model, we need to change some of the full conditional distributions in the LEB model for water. For the Water medium we replace equation (1) with:

$$Y_{i,m(W)}^{M*} \begin{cases} = Y_{i,m(W)}^M, & \text{if } Z_{i,m(W)}^M = 0; \\ \sim \text{Truncated } N(X_{i,m(W)}^M, 1/\omega_{m(W)}^M) \text{ on } (-\infty, M_{i,m(W)}^M], & \text{if } Z_{i,m(W)}^M = 1; \\ \sim N(X_{i,m(W)}^M, 1/\omega_{m(W)}^M), & \text{if } Z_{i,m(W)}^M = 2. \end{cases}$$

To sample from the conditional distribution of the water process, $\{X_{i,m(W)}^M : i = 1, \dots, N^I\}$, we let $\mu_{i,m(W)}^M = I(\zeta_{ij}^W = 1)\mu_j^W$. Then $X_{i,m(W)}^M$ is $N(q_i/p_i, 1/p_i)$, where

$$p_i = \omega_{m(W)}^M + \tau_{m(W)}^M \quad \text{and} \quad q_i = \omega^W Y_{i,m(W)}^{W*} + \tau_{m(W)}^M \mu_{i,m(W)}^M,$$

for each individual i .

With the addition of the global water model, the full conditional distribution of the NHEXAS Water model precision, $\tau_{m(W)}^M$, is now

$$\text{Ga} \left(\frac{N^I}{2} + a^{\tau_{m(W)}^M}, \frac{1}{2} \sum_{i=1}^{N^I} \left(X_{i,m(W)}^M - \mu_{i,m(W)}^M \right)^2 + b^{\tau_{m(W)}^M} \right).$$

2. The full conditional distributions of the remaining parameters in the global water model depends on the augmented data Y_{jk}^{W*} , defined as follows:

$$Y_{jk}^{W*} \begin{cases} = Y_{jk}^W, & \text{if } Z_{jk}^W = 0; \\ \sim \text{Truncated } N(X_j^W, 1/\omega^W) \text{ on } (-\infty, M_{jk}^W], & \text{if } Z_{jk}^W = 1; \\ \sim N(X_j^W, 1/\omega^W), & \text{if } Z_{jk}^W = 2. \end{cases}$$

The full conditional distribution of the PWS measurement-error precision, ω^W , is then

$$\text{Ga} \left(\frac{\sum_{j=1}^{N^W} N_j^W}{2} + a^{\omega^W}, \frac{1}{2} \sum_{j=1}^{N^W} \sum_{k=1}^{N_j^W} (Y_{jk}^{W*} - X_j^W)^2 + b^{\omega^W} \right).$$

3. For the global water model, the full conditional distribution of X_{jk}^W for the k th observation ($k = 1, \dots, N_j^W$) of j th PWS ($j = 1, \dots, N^W$) is $N(q_{jk}/p_{jk}, 1/p_{jk})$, where

$$p_{jk} = \omega^W + \tau^W \quad \text{and} \quad q_{jk} = \omega^W Y_{jk}^{W*} + \tau^W \mu_j^W.$$

Next, define $\mathbf{X}_j^W = \{X_{i,m(W)}^M : \zeta_{ij} = 1\}$ and l_j to be the number of elements in \mathbf{X}_j^W . Then, for each PWS $j = 1, \dots, N^W$, the full conditional distribution of the PWS-specific mean μ_j^W is $N(q_j/p_j, 1/p_j)$, with

$$p_j = N_j^W \tau^W + C^W + l_j \tau_{m(W)}^M$$

and

$$q_j = \tau^W \sum_{k=1}^{N_j^W} X_{jk}^W + C^W \alpha^W + \tau_{m(W)}^M \mathbf{1}_{l_k} \mathbf{X}_j^W.$$

The full conditional distribution of the process precision, τ^W , is

$$\text{Ga} \left(\frac{\sum_{j=1}^{N^W} N_j^W}{2} + a\tau^W, \frac{1}{2} \sum_{j=1}^{N^W} \sum_{k=1}^{N_j^W} (X_{jk}^W - \mu_j^W)^2 + b\tau^W \right).$$

4. For the parameters in the global water-LEB linking model, recall that $\zeta_{ij}^W = 1$ implies that individual i is served by PWS j and $\zeta_{ij}^W = 0$ implies that individual i is not served by water-system j . For each individual $i = 1, \dots, N^I$, the full conditional distribution of $\zeta_i^W = (\zeta_{i1}^W, \dots, \zeta_{iN^W}^W)$ is Multinomial($1, (\lambda_{c(i),1}^*, \dots, \lambda_{c(i),N^W}^*)$) where the probability that individual i in country $c(i)$ is served by PWS j , is

$$\lambda_{c(i),j}^* = \frac{\lambda_{c(i),j} n(X_{i,m(W)}^M; \mu_j^W, \tau_{m(W)}^M)}{\sum_{r=1}^{N^W} \lambda_{c(i),r} n(X_{i,m(W)}^M; \mu_r^W, \tau_{m(W)}^M)}.$$

Here, $n(X_{i,m(W)}^M; \mu_j^W, \tau_{m(W)}^M)$ is a normal pdf (with mean μ_j^W and variance $1/\tau_{m(W)}^M$) evaluated at the value of the water process, $X_{i,m(W)}^M$, for each individual i .

5. In the parameters in the global water prior distributions, α^W conditional on all the other parameters is $N(q/p, p)$, where

$$p = N^W C^W + V\alpha^W \quad \text{and} \quad q = \left(C^W \sum_{j=1}^{N^W} \mu_j^W + V\alpha^W m^{\alpha^W} \right).$$

The full conditional distribution of C^W is

$$\text{Ga} \left(\frac{N^W}{2} + aC^W, \frac{1}{2} \sum_{j=1}^{N^W} (\mu_j^W - m^{\alpha^W})^2 + bC^W \right).$$

2.2 Global soil model

1. With the addition of the global topsoil/stream-sediment model, we need to change some of the conditional distributions in the LEB model for soil. For the soil medium, we replace equation (1) with:

$$Y_{i,m(s)}^{M*} \begin{cases} = Y_{i,m(s)}^M, & \text{if } Z_{i,m(s)}^M = 0; \\ \sim \text{Truncated } N(X_{i,m(s)}^M, 1/\omega_{m(s)}^M) \text{ on } (-\infty, M_{i,m(s)}^M], & \text{if } Z_{i,m(s)}^M = 1; \\ \sim N(X_{i,m(s)}^M, 1/\omega_{m(s)}^M), & \text{if } Z_{i,m(s)}^M = 2. \end{cases}$$

2. For parameters appearing in the global topsoil/stream-sediment model, for the topsoil observation at the s th location:

$$Y^{T*}(s) \begin{cases} = Y^T(s), & \text{if } Z^T(s) = 0; \\ \sim N(X^T(s), 1/\omega^T), & \text{if } Z^T(s) = 2 \end{cases}$$

(there is no censored topsoil data).

The full conditional distribution of the topsoil measurement-error precision, ω^T , is

$$\text{Ga}\left(\frac{N^T}{2} + a^{\omega^T}, \frac{1}{2} \sum_{i=1}^{N^T} (Y^{T*}(s_i) - X^T(s_i))^2 + b^{\omega^T}\right).$$

3. For parameters appearing in the stream-sediment data model, for the k th observation in j th HUC8:

$$Y_{jk}^{H*} \begin{cases} = Y_{jk}^H, & \text{if } Z_{jk}^H = 0; \\ \sim \text{Truncated } N(X_j^H, 1/\omega^H) \text{ on } (-\infty, M_{jk}^H], & \text{if } Z_{jk}^H = 1; \\ \sim N(X_j^H, 1/\omega^H), & \text{if } Z_{jk}^H = 2. \end{cases}$$

The full conditional distribution of the stream-sediment measurement-error precision, ω^H , is

$$\text{Ga}\left(\sum_{j=1}^{N^H} \frac{N_j^H}{2} + a^{\omega^H}, \frac{1}{2} \sum_{j=1}^{N^H} \sum_{k=1}^{N_j^H} (Y_{jk}^H - X_j^H)^2 + b^{\omega^H}\right).$$

4. The full conditional distribution of the topsoil process $X^T(s)$ is $N(q_s/p_s, 1/p_s)$, where

$$p_s = \omega^T + \tau^T \quad \text{and} \quad q_s = \tau^T (\beta_0^T + \beta_1^T X_{h(s)}^H) + \omega^T Y^{T*}(s),$$

for $s \in \mathcal{D}$.

The full conditional distribution of the topsoil model precision, τ^T , is

$$\text{Ga}\left(\frac{|\mathcal{D}|}{2} + a^{\tau^T}, \frac{1}{2} \sum_{s \in \mathcal{D}} (X^T(s) - \beta_0^T - \beta_1^T X_{h(s)}^H)^2 + b^{\tau^T}\right).$$

The parameters that appear in the mean function of the topsoil model are $\beta^T = (\beta_0^T, \beta_1^T)'$. Letting $\mathbf{U} = [\mathbf{1}_{|\mathcal{D}|}, \mathbf{G}_D \mathbf{X}^H]$ and $\mathbf{W} = [\mathbf{1}_{N^I}, \zeta^T \mathbf{X}^H]$ we sample β^T from $N_2(P^{-1}Q, P^{-1})$, where

$$P = \tau^T \mathbf{U}' \mathbf{U} + \tau_{m(S)}^M \mathbf{W}' \mathbf{W} + \mathbf{V} \beta^T$$

and

$$Q = \tau^T \mathbf{U}' \mathbf{X}^T + \tau_{m(S)}^M \mathbf{W}' \mathbf{X}_{m(S)}^M + \mathbf{V} \beta^T \mathbf{M} \beta^T.$$

5. The full conditional distribution of the mean of the stream-sediment model, μ^H , is $N(q/p, 1/p)$, where

$$p = \tau^H \mathbf{1}'_{N^H} (\mathbf{I}_{N^H} - \gamma^H \mathbf{A}) \mathbf{1}_{N^H} + V \mu^H,$$

and

$$q = \tau^H \mathbf{1}'_{N^H} (\mathbf{I}_{N^H} - \gamma^H \mathbf{A}) \mathbf{X}^H + V \mu^H m^{\mu^H}.$$

6. Recall that for the unknown parameter ζ_{ij}^T , $\zeta_{ij}^T = 1$ if NHEXAS individual i is located in HUC8 region j , and 0 otherwise. These ζ_{ij}^T 's are updated during our MCMC algorithm using the following scheme:

Sample $\zeta_i^T = (\zeta_{i1}^T, \dots, \zeta_{iN^H}^T) \sim \text{Multinomial}\left(1, (\lambda^*_{c(i),1}, \dots, \lambda^*_{c(i),N^H})'\right)$, for $i = 1, \dots, N^I$, where

$$\lambda^*_{c(i),j} = \frac{\lambda_{c(i),j} n(X_{i,m(S)}^M; \beta_0^T + \beta_1^T X_j^H, \tau_{m(S)}^M)}{\sum_{r=1}^{N^H} \lambda_{c(i),r} n(X_{i,m(S)}^M; \beta_0^T + \beta_1^T X_r^H, \tau_{m(S)}^M)},$$

and $n(X_{i,m(S)}^M; \beta_0^T + \beta_1^T X_j^H, \tau_{m(S)}^M)$ is a normal pdf (with mean $\beta_0^T + \beta_1^T X_j^H$ and variance $1/\tau_{m(S)}^M$) evaluated at $X_{i,m(S)}^M$.

Now define the $|\mathcal{D}| \times N^H$ matrix \mathbf{G}_D such that the (i,j) th entry $G_D(i,j)$ is 1 if topsoil location s_i falls in j th HUC8, and 0 otherwise. Sample \mathbf{X}^H from $N(P^{-1}Q, P^{-1})$, where

$$P = \tau^H (\mathbf{I}_{N^H} - \gamma^H \mathbf{A}) + \tau^T \mathbf{G}'_D \mathbf{G}_D + \tau_{m(S)}^M \zeta^{T'} \zeta^T,$$

and

$$Q = \tau^H (\mathbf{I}_{N^H} - \gamma^H \mathbf{A}) \mu^H \mathbf{1}_{N^H} + \tau^T \mathbf{G}'_D (\mathbf{X}^T - \beta_0^T \mathbf{1}_{|\mathcal{D}|}) + \tau_{m(S)}^M \zeta^{T'} (\mathbf{X}_{m(S)}^M - \beta_0^T \mathbf{1}_{N^I}).$$

7. The full conditional distribution of the stream-sediment precision, τ^H , is

$$\text{Ga}\left(\frac{N^H}{2} + a^{\tau^H}, \frac{1}{2}(\mathbf{X}^H - \mu^H \mathbf{1}_{N^H})'(\mathbf{I}_{N^H} - \gamma^H \mathbf{A})(\mathbf{X}^H - \mu^H \mathbf{1}_{N^H}) + b^{\tau^H}\right).$$

The full conditional distribution of the mean of the stream-sediment model, μ^H , is $N(q/p, 1/p)$, where

$$p = \tau^H \mathbf{1}'_{N^H} (\mathbf{I}_{N^H} - \gamma^H \mathbf{A}) \mathbf{1}_{N^H} + V^{\mu^H},$$

and

$$q = \tau^H \mathbf{1}'_{N^H} (\mathbf{I}_{N^H} - \gamma^H \mathbf{A}) \mathbf{X}^H + V^{\mu^H} m^{\mu^H}.$$

The full conditional distribution of the spatial-dependence parameter in the stream-sediment model, γ^H , is proportional to

$$|\mathbf{I}_{N^H} - \gamma^H \mathbf{A}|^{1/2} \exp\{-0.5\tau^H(\mathbf{X}^H - \mu^H \mathbf{1}_{N^H})'(\mathbf{I}_{N^H} - \gamma^H \mathbf{A})(\mathbf{X}^H - \mu^H \mathbf{1}_{N^H})\} \times \frac{1}{|\lambda_{\max} - \lambda_{\min}|} \mathbf{I}(\gamma^H \in (\lambda_{\min}, \lambda_{\max})),$$

where λ_{\min} , λ_{\max} are the smallest and largest eigenvalues of \mathbf{A} .

8. With the addition of the global-topsoil/stream-sediment model, the full conditional distribution of the LEB Soil log As process $\mathbf{X}^M_{m(S)}$, is modified. Now, this distribution is $N_{N^I}(\mathbf{P}^{-1}\mathbf{Q}, \mathbf{P}^{-1})$, where

$$\mathbf{P} = \omega^M_{m(S)} \mathbf{I}_{N^I} + \tau^M_{m(S)} \mathbf{I}_{N^I},$$

and

$$\mathbf{Q} = \omega^M_{m(S)} \mathbf{Y}^{M*}_{m(S)} + \tau^M_{m(S)} (\beta_0^T \mathbf{1}_{N^I} + \beta_1^T \boldsymbol{\zeta}^{T'} \mathbf{X}^H).$$

In addition, the full conditional distribution of the local soil process precision, $\tau^M_{m(S)}$, is now

$$\text{Ga}\left(\frac{N^I}{2} + a^{\tau^M_{m(S)}}, \frac{1}{2}(\mathbf{X}^M_{m(S)} - \mathbf{W}\boldsymbol{\beta}^T)'(\mathbf{X}^M_{m(S)} - \mathbf{W}\boldsymbol{\beta}^T) + b^{\tau^M_{m(S)}}\right).$$