

# Fast Simulation of Hyperplane-Truncated Multivariate Normal Distributions

Yulai Cong<sup>\*</sup>, Bo Chen<sup>†§</sup>, and Mingyuan Zhou<sup>‡§</sup>

**Abstract.** We introduce a fast and easy-to-implement simulation algorithm for a multivariate normal distribution truncated on the intersection of a set of hyperplanes, and further generalize it to efficiently simulate random variables from a multivariate normal distribution whose covariance (precision) matrix can be decomposed as a positive-definite matrix minus (plus) a low-rank symmetric matrix. Example results illustrate the correctness and efficiency of the proposed simulation algorithms.

**Keywords:** Cholesky decomposition, conditional distribution, equality constraints, high-dimensional regression, structured covariance/precision matrix

## 1 Introduction

We investigate the problem of simulation from a multivariate normal (MVN) distribution whose samples are restricted to the intersection of a set of hyperplanes, which is shown to be inherently related to the simulation of a conditional distribution of a MVN distribution. A naive approach, which linearly transforms a random variable drawn from the conditional distribution of a related MVN distribution, requires a large number of intermediate variables that are often computationally expensive to instantiate. To address this issue, we propose a fast and exact simulation algorithm that directly projects a MVN random variable onto the intersection of a set of hyperplanes. We further show that sampling from a MVN distribution, whose covariance (precision) matrix can be decomposed as the sum (difference) of a positive-definite matrix, whose inversion is known or easy to compute, and a low-rank symmetric matrix, may also be made significantly fast by exploiting this newly proposed stimulation algorithm for hyperplane-truncated MVN distributions, avoiding the need of Cholesky decomposition that has a computational complexity of  $O(k^3)$  (Golub and Van Loan 2012), where  $k$  is the dimension of the MVN random variable.

Related to the problems under study, the simulation of MVN random variables subject to certain constraints (Gelfand et al. 1992) has been investigated in many other different settings, such as multinomial probit and logit models (Albert and Chib 1993;

---

<sup>\*</sup>National Laboratory of Radar Signal Processing and Collaborative Innovation Center of Information Sensing & Understanding, Xidian University, Xi'an, Shaanxi 710071, China, yulai.cong@163.com

<sup>†</sup>National Laboratory of Radar Signal Processing and Collaborative Innovation Center of Information Sensing & Understanding, Xidian University, Xi'an, Shaanxi 710071, China, bchen@mail.xidian.edu.cn

<sup>‡</sup>McCombs School of Business, The University of Texas at Austin, Austin, TX 78712, USA, mingyuan.zhou@mcombs.utexas.edu

<sup>§</sup>Corresponding authors.

McCulloch et al. 2000; Imai and van Dyk 2005; Train 2009; Holmes and Held 2006; Johndrow et al. 2013), Bayesian isotonic regression (Neelon and Dunson 2004), Bayesian bridge (Polson et al. 2014), blind source separation (Schmidt 2009), and unmixing of hyperspectral data (Altmann et al. 2014; Dobigeon et al. 2009a). A typical example arising in these different settings is to sample a random vector  $\mathbf{x} \in \mathbb{R}^k$  from a MVN distribution subject to  $k$  inequality constraints as

$$\mathbf{x} \sim \mathcal{N}_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \mathcal{S} = \{\mathbf{x} : \mathbf{l} \leq \mathbf{G}\mathbf{x} \leq \mathbf{u}\}, \quad (1)$$

where  $\mathcal{N}_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  represents a MVN distribution truncated on the sample space  $\mathcal{S}$ ,  $\boldsymbol{\mu} \in \mathbb{R}^k$  is the mean,  $\boldsymbol{\Sigma} \in \mathbb{R}^{k \times k}$  is the covariance matrix,  $\mathbf{G} \in \mathbb{R}^{k_2 \times k}$  is a full-rank matrix,  $\mathbf{l} \in \mathbb{R}^{k_2}$ ,  $\mathbf{u} \in \mathbb{R}^{k_2}$ , and  $\mathbf{l} < \mathbf{u}$ . If the elements of  $\mathbf{l}$  and  $\mathbf{u}$  are permitted to be  $-\infty$  and  $+\infty$ , respectively, then both single sided and fewer than  $k$  inequality constraints are allowed. Equivalently, as in Geweke (1991, 1996), one may let  $\mathbf{x} = \boldsymbol{\mu} + \mathbf{G}^{-1}\mathbf{z}$  and use Gibbs sampling (Geman and Geman 1984; Gelfand and Smith 1990) to sample the  $k$  elements of  $\mathbf{z}$  one at a time conditioning on all the others from a univariate truncated normal distribution, for which efficient algorithms exist (Robert 1995; Damien and Walker 2001; Chopin 2011). To deal with the case that the number of linear constraints imposed on  $\mathbf{x}$  exceed its dimension  $k$  and to obtain better mixing, one may consider the Gibbs sampling algorithm for truncated MVN distributions proposed in Rodriguez-Yam et al. (2004). In addition to Gibbs sampling, to sample truncated MVN random variables, one may also consider Hamiltonian Monte Carlo (Pakman and Paninski 2014; Lan et al. 2014) and a minimax tilting method proposed in Botev (2016).

## 2 Hyperplane-truncated and conditional MVNs

For the problem under study, we express a  $k$ -dimensional MVN distribution truncated on the intersection of  $k_2 < k$  hyperplanes as

$$\mathbf{x} \sim \mathcal{N}_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \mathcal{S} = \{\mathbf{x} : \mathbf{G}\mathbf{x} = \mathbf{r}\}, \quad (2)$$

where

$$\mathbf{G} \in \mathbb{R}^{k_2 \times k}, \quad \mathbf{r} \in \mathbb{R}^{k_2},$$

and  $\text{Rank}(\mathbf{G}) = k_2$ . The probability density function can be expressed as

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{G}, \mathbf{r}) = \frac{1}{Z} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \delta(\mathbf{G}\mathbf{x} = \mathbf{r}), \quad (3)$$

where  $Z$  is a constant ensuring  $\int p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{G}, \mathbf{r}) d\mathbf{x} = 1$ , and  $\delta(x) = 1$  if the condition  $x$  is satisfied and  $\delta(x) = 0$  otherwise. Let us partition  $\mathbf{G}$ ,  $\mathbf{x}$ ,  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Sigma}$ , and  $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$  as

$$\mathbf{G} = (\mathbf{G}_1, \mathbf{G}_2), \quad \mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\Lambda} = \begin{bmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{bmatrix},$$

whose sizes are

$$(k_2 \times k_1, k_2 \times k_2), \quad \begin{bmatrix} k_1 \times 1 \\ k_2 \times 1 \end{bmatrix}, \quad \begin{bmatrix} k_1 \times 1 \\ k_2 \times 1 \end{bmatrix}, \quad \begin{bmatrix} k_1 \times k_1 & k_1 \times k_2 \\ k_2 \times k_1 & k_2 \times k_2 \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} k_1 \times k_1 & k_1 \times k_2 \\ k_2 \times k_1 & k_2 \times k_2 \end{bmatrix},$$

respectively, where  $k = k_1 + k_2$ ,  $\boldsymbol{\Sigma}_{21} = \boldsymbol{\Sigma}_{12}^T$ , and  $\boldsymbol{\Lambda}_{21} = \boldsymbol{\Lambda}_{12}^T$ .

A special case that frequently arises in real applications is when  $\mathbf{G}_1 = \mathbf{0}_{k_2 \times k_1}$  and  $\mathbf{G}_2 = \mathbf{I}_{k_2}$ , which means  $(\mathbf{0}_{k_2 \times k_1}, \mathbf{I}_{k_2})\mathbf{x} = \mathbf{x}_2 = \mathbf{r}$  and the need is to simulate  $\mathbf{x}_1$  given  $\mathbf{x}_2 = \mathbf{r}$ . For a MVN random variable  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , it is well known, *e.g.*, in Tong (2012), that the conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2 = \mathbf{r}$ , *i.e.*, the distribution of  $\mathbf{x}$  restricted to  $\mathcal{S} = \{\mathbf{x} : (\mathbf{0}_{k_2 \times k_1}, \mathbf{I}_{k_2})\mathbf{x} = \mathbf{r}\}$ , can be expressed as

$$\mathbf{x}_1 | \mathbf{x}_2 = \mathbf{r} \sim \mathcal{N}[\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{r} - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}]. \quad (4)$$

Alternatively, applying the Woodbury matrix identity to relate the entries of the covariance matrix  $\boldsymbol{\Sigma}$  to those of the precision matrix  $\boldsymbol{\Lambda}$ , one may obtain the following equivalent expression as

$$\mathbf{x}_1 | \mathbf{x}_2 = \mathbf{r} \sim \mathcal{N}[\boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12}(\mathbf{r} - \boldsymbol{\mu}_2), \boldsymbol{\Lambda}_{11}^{-1}]. \quad (5)$$

In a general setting where  $\mathbf{G} \neq (\mathbf{0}_{k_2 \times k_1}, \mathbf{I}_{k_2})$ , let us define a full rank linear transformation matrix  $\mathbf{H} \in \mathbb{R}^{k \times k}$ , with  $(\mathbf{H}_1, \mathbf{H}_2)$  as the  $(k \times k_1, k \times k_2)$  partition of  $\mathbf{H}$ , where the columns of  $\mathbf{H}_1 \in \mathbb{R}^{k \times k_1}$  span the null space of the  $k_2$  rows of  $\mathbf{G}$ , making  $\mathbf{GH} = (\mathbf{GH}_1, \mathbf{GH}_2) = (\mathbf{0}_{k_2 \times k_1}, \mathbf{GH}_2)$ , where  $\mathbf{GH}_2$  is a  $k_2 \times k_2$  full rank matrix. For example, a linear transformation matrix  $\mathbf{H}$  that makes  $\mathbf{GH} = (\mathbf{0}_{k_2 \times k_1}, \mathbf{I}_{k_2})$  can be constructed using the command  $\mathbf{H} = \text{inv}([\text{null}(\mathbf{G})'; \mathbf{G}])$  in Matlab and  $\mathbf{H} \leftarrow \text{solve}(\text{rbind}(\text{t}(\text{Null}(\text{t}(\mathbf{G}))), \mathbf{G}))$  in R. With  $\mathbf{H}$  and  $\mathbf{H}^{-1}$ , one may re-express the constraints as  $\mathcal{S} = \{\mathbf{x} : (\mathbf{0}_{k_2 \times k_1}, \mathbf{GH}_2)(\mathbf{H}^{-1}\mathbf{x}) = \mathbf{r}\}$ . Denote  $\mathbf{z} = \mathbf{H}^{-1}\mathbf{x}$ , then we can generate  $\mathbf{x}$  by letting  $\mathbf{x} = \mathbf{H}\mathbf{z}$ , where

$$\mathbf{z} \sim \mathcal{N}_{\mathcal{D}}[\mathbf{H}^{-1}\boldsymbol{\mu}, \mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^{-1})^T], \quad \mathcal{D} = \{\mathbf{z} : \mathbf{GH}_2\mathbf{z}_2 = \mathbf{r}\} = \{\mathbf{z} : \mathbf{z}_2 = (\mathbf{GH}_2)^{-1}\mathbf{r}\}. \quad (6)$$

More specifically, denoting  $\boldsymbol{\Lambda} = [\mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^{-1})^T]^{-1} = \mathbf{H}^T\boldsymbol{\Sigma}^{-1}\mathbf{H}$  as the precision matrix for  $\mathbf{z}$ , we have

$$\begin{bmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{bmatrix} = \mathbf{H}^T\boldsymbol{\Sigma}^{-1}\mathbf{H} = \begin{bmatrix} \mathbf{H}_1^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_1 & \mathbf{H}_1^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_2 \\ \mathbf{H}_2^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_1 & \mathbf{H}_2^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_2 \end{bmatrix}, \quad (7)$$

and hence  $\mathbf{x}$  truncated on  $\mathcal{S}$  can be naively generated using the following algorithm, whose computational complexity is described in Table 1 of the Appendix.

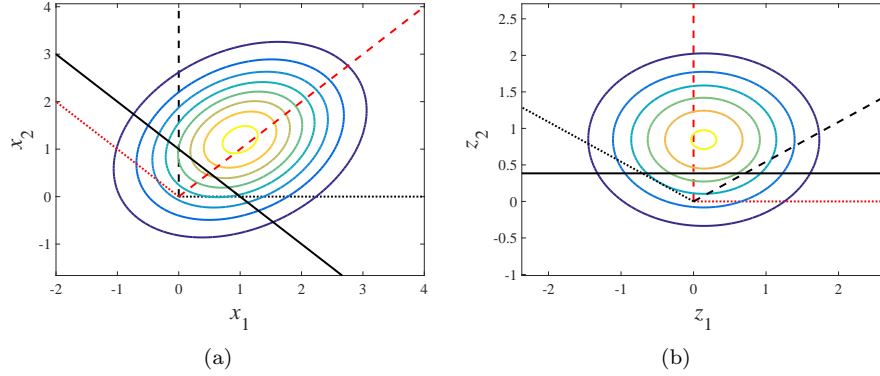


Figure 1: Illustration of (a)  $p(\mathbf{x})$  in (3), where  $\boldsymbol{\mu} = (1, 1.2)^T$ ,  $\boldsymbol{\Sigma} = [(1, 0.3)^T, (0.3, 1)^T]$ ,  $\mathbf{G} = (1, 1)$ , and  $\mathbf{r} = 1$ , and (b)  $p(\mathbf{z})$  in (6), where  $\mathbf{H}_1 = (-0.7071, 0.7071)^T$ ,  $\mathbf{H}_2 = (1.3, 1.3)^T$ , and  $\mathbf{H}^{-1} = [(-0.7071, 0.3846)^T, (0.7071, 0.3846)^T]$ . The coordinate systems of  $\mathbf{x}$  and  $\mathbf{z}$  are shown in black and red, respectively, and the first and second axes of a coordinate system are shown as dotted and dashed lines, respectively.

---

**Algorithm 1** Simulation of the hyperplane truncated MVN distribution  $\mathbf{x} \sim \mathcal{N}_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\mathcal{S} = \{\mathbf{x} : \mathbf{G}\mathbf{x} = \mathbf{r}\}$ , by transforming a random variable drawn from the conditional distribution of another MVN distribution.

---

- Find  $\mathbf{H} = (\mathbf{H}_1, \mathbf{H}_2)$  that satisfies  $\mathbf{G}\mathbf{H} = (\mathbf{G}\mathbf{H}_1, \mathbf{G}\mathbf{H}_2) = (\mathbf{0}_{k_2 \times k_1}, \mathbf{G}\mathbf{H}_2)$ , where  $\mathbf{G}\mathbf{H}_2$  is a full rank matrix;
- Let  $\mathbf{z}_2 = (\mathbf{G}\mathbf{H}_2)^{-1}\mathbf{r}$ ,  $\boldsymbol{\Lambda}_{11} = \mathbf{H}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{H}_1$ , and  $\boldsymbol{\Lambda}_{12} = \mathbf{H}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{H}_2$ ;
- Sample  $\mathbf{z}_1 \mid \mathbf{z}_2 = (\mathbf{G}\mathbf{H}_2)^{-1}\mathbf{r} \sim \mathcal{N}(\boldsymbol{\mu}_{z_1}, \boldsymbol{\Lambda}_{11}^{-1})$ , where

$$\boldsymbol{\mu}_{z_1} = (\mathbf{I}_{k_1}, \mathbf{0}_{k_1 \times k_2}) \mathbf{H}^{-1} \boldsymbol{\mu} - \boldsymbol{\Lambda}_{11}^{-1} \boldsymbol{\Lambda}_{12} [(\mathbf{G}\mathbf{H}_2)^{-1}\mathbf{r} - (\mathbf{0}_{k_2 \times k_1}, \mathbf{I}_{k_2}) \mathbf{H}^{-1} \boldsymbol{\mu}];$$

- Return  $\mathbf{x} = \mathbf{H}\mathbf{z} = \mathbf{H}_1\mathbf{z}_1 + \mathbf{H}_2(\mathbf{G}\mathbf{H}_2)^{-1}\mathbf{r}$ .
- 

For illustration, we consider a simple 2-dimensional example with  $\boldsymbol{\mu} = (1, 1.2)^T$ ,  $\boldsymbol{\Sigma} = [(1, 0.3)^T, (0.3, 1)^T]$ ,  $\mathbf{G} = (1, 1)$ , and  $\mathbf{r} = 1$ . If we choose  $\mathbf{H}_1 = (-0.7071, 0.7071)^T$  and  $\mathbf{H}_2 = (1.3, 1.3)^T$ , then we have  $\mathbf{z}_2 = (\mathbf{G}\mathbf{H}_2)^{-1}\mathbf{r} = (2.6)^{-1} = 0.3846$ ,  $\boldsymbol{\Lambda}_{11} = 1.4285$ , and  $\boldsymbol{\Lambda}_{12} = 0$ ; as shown in Figure 1, we may generate  $\mathbf{x}$  using

$$\mathbf{x} = (-0.7071, 0.7071)^T z_1 + (1.3, 1.3)^T z_2$$

where  $z_1 \sim \mathcal{N}(0.1414, 0.7)$  and  $z_2 = 0.3846$ .

For high dimensional problems, however, Algorithm 1 in general requires a large number of intermediate variables that could be computationally expensive to compute.

In the following discussion, we will show how to completely avoid instantiating these intermediate variables.

### 3 Fast and exact simulation of MVN distributions

Instead of using Algorithm 1, we first provide a theorem to show how to efficiently and exactly simulate from a hyperplane-truncated MVN distribution. In the Appendix, we provide two different proofs. The first proof facilitates the derivations by employing an existing algorithm of Hoffman and Ribak (1991) and Doucet (2010), which describes how to simulate from the conditional distribution of a MVN distribution shown in (4) without computing  $\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$  and its Cholesky decomposition. Note it is straightforward to verify that the algorithm in Hoffman and Ribak (1991) and Doucet (2010), as shown in the Appendix, can be considered as a special case of the proposed algorithm with  $\mathbf{G} = [\mathbf{0}, \mathbf{I}]$ .

---

**Algorithm 2** Simulation of the hyperplane truncated MVN distribution  $\mathbf{x} \sim \mathcal{N}_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\mathcal{S} = \{\mathbf{x} : \mathbf{G}\mathbf{x} = \mathbf{r}\}$ , by transforming a random variable drawn from  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

---

- Sample  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ;
  - Return  $\mathbf{x} = \mathbf{y} + \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}(\mathbf{r} - \mathbf{G}\mathbf{y})$ , which can be realized using
    - Solve  $\boldsymbol{\alpha}$  such that  $(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)\boldsymbol{\alpha} = \mathbf{r} - \mathbf{G}\mathbf{y}$ ;
    - Return  $\mathbf{x} = \mathbf{y} + \boldsymbol{\Sigma}\mathbf{G}^T\boldsymbol{\alpha}$ .
- 

**Theorem 1.** *Suppose  $\mathbf{x}$  is simulated with Algorithm 2, then it is distributed as  $\mathbf{x} \sim \mathcal{N}_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\mathcal{S} = \{\mathbf{x} : \mathbf{G}\mathbf{x} = \mathbf{r}\}$ , where  $\mathbf{G} \in \mathbb{R}^{k_2 \times k}$ ,  $\mathbf{r} \in \mathbb{R}^{k_2}$ , and  $\text{Rank}(\mathbf{G}) = k_2 < k$ .*

The above algorithm and theorem, whose computational complexity is described in Table 2 of the Appendix, show that one may draw  $\mathbf{y}$  from the unconstrained MVN as  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and directly map it to a vector  $\mathbf{x}$  on the intersection of hyperplanes using  $\mathbf{x} = \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} + [\mathbf{I} - \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{G}]\mathbf{y}$ . For illustration, with the same  $\boldsymbol{\mu}$ ,  $\boldsymbol{\Sigma}$ ,  $\mathbf{G}$ , and  $\mathbf{r}$  as those in Figure 1, we show in Figure 2 a simple two dimensional example, where the unrestricted Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is represented with a set of ellipses, and the constrained sample space  $\mathcal{S}$  is represented as a straight line in the two-dimensional setting. With  $\boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} = (0.5, 0.5)^T$ ,  $[\mathbf{I} - \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{G}] = [(0.5, -0.5)^T, (-0.5, 0.5)^T]$ , one may directly map a sample  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  to a vector on the constrained space. For example, if  $\mathbf{y} = (1, 2)^T$ , then it would be mapped to  $\mathbf{x} = (0, 1)^T$  on the straight line.

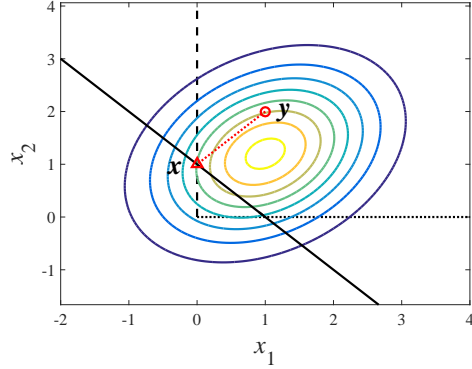


Figure 2: A two dimensional demonstration of Algorithm 2 that maps a random sample from  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  to a sample in the constrained space using  $\mathbf{x} = \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} + [\mathbf{I} - \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{G}]\mathbf{y}$ . For example, if  $\boldsymbol{\mu} = (1, 1.2)^T$ ,  $\boldsymbol{\Sigma} = [(1, 0.3)^T, (0.3, 1)^T]$ ,  $\mathbf{G} = (1, 1)$ , and  $\mathbf{r} = 1$ , then  $\mathbf{y} = (1, 2)^T$  would be mapped to  $\mathbf{x} = (0, 1)^T$  on a straight line using Algorithm 2.

### 3.1 Fast simulation of MVN distributions with structured covariance or precision matrices

For fast simulation of MVN distributions with structured covariance or precision matrices, our idea is to relate them to higher-dimensional hyperplane-truncated MVN distributions, with block-diagonal covariance matrices, that can be efficiently simulated with Algorithm 2. We first introduce an efficient algorithm for the simulation of a MVN distribution, whose covariance matrix is a positive-definite matrix subtracted by a low-rank symmetric matrix. Such kind of covariance matrices commonly arise in the conditional distributions of MVN distributions, as shown in (4). We then further extend this algorithm to the simulation of a MVN distribution whose precision (inverse covariance) matrix is the sum of a positive-definite matrix and a low-rank symmetric matrix. Such kind of precision matrices commonly arise in the conditional posterior distributions of the regression coefficients in both linear regression and generalized linear models.

**Theorem 2.** *The probability density function (PDF) of the MVN distribution*

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}), \quad (8)$$

is the same as the PDF of the marginal distribution of  $\mathbf{x}_1 = (x_1, \dots, x_{k_1})^T$  in  $\mathbf{x} = (\mathbf{x}_1^T, x_{k_1+1}, \dots, x_k)^T$ , whose PDF is expressed as

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\mu}, \tilde{\boldsymbol{\Sigma}}, \mathbf{G}, \mathbf{r}) &= \mathcal{N}_{\{\mathbf{x}: \mathbf{G}\mathbf{x}=\mathbf{r}\}}(\boldsymbol{\mu}, \tilde{\boldsymbol{\Sigma}}) \\ &= \frac{1}{Z} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \tilde{\boldsymbol{\Sigma}}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right] \delta(\mathbf{G}\mathbf{x} = \mathbf{r}), \end{aligned} \quad (9)$$

where  $Z$  is a normalization constant,  $\mathbf{G}_1 = \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}$  is a matrix of size  $k_2 \times k_1$ ,  $\mathbf{G}_2$  is a user-specified full rank invertible matrix of size  $k_2 \times k_2$ ,  $\mathbf{r} \in \mathbb{R}^{k_2}$  is a user-specified vector, and

$$\mathbf{G} = (\mathbf{G}_1, \mathbf{G}_2) \in \mathbb{R}^{k_2 \times k}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \in \mathbb{R}^k, \quad \tilde{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0} & \tilde{\boldsymbol{\Sigma}}_{22} \end{bmatrix} \in \mathbb{R}^{k \times k}, \quad (10)$$

where

$$\boldsymbol{\mu}_2 = \mathbf{G}_2^{-1}(\mathbf{r} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1), \quad (11)$$

$$\tilde{\boldsymbol{\Sigma}}_{22} = \mathbf{G}_2^{-1}(\boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})(\mathbf{G}_2^{-1})^T. \quad (12)$$

The above theorem shows how the simulation of a MVN distribution, whose covariance matrix is a positive-definite matrix minus a symmetric matrix, can be realized by the simulation of a higher-dimensional hyperplane-truncated MVN distribution. By construction, it makes the covariance matrix  $\tilde{\boldsymbol{\Sigma}}$  of the truncated-MVN be block diagonal, but still preserves the flexibility to customize the full-rank matrix  $\mathbf{G}_2$  and the vector  $\mathbf{r}$ . While there are infinitely many choices for both  $\mathbf{G}_2$  and  $\mathbf{r}$ , in the following discussion, we remove that flexibility by specifying  $\mathbf{G}_2 = \mathbf{I}_{k_2}$ , leading to  $\mathbf{G} = (\mathbf{G}_1, \mathbf{G}_2) = (\boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}, \mathbf{I}_{k_2})$ , and  $\mathbf{r} = \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1$ . This specific setting of  $\mathbf{G}_2$  and  $\mathbf{r}$  leads to the following Corollary that is a special case of Theorem 2. Note that while we choose this specific setting in the paper, depending on the problems under study, other settings may lead to even more efficient simulation algorithms.

**Corollary 3.** *The PDF of the MVN distribution*

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}) \quad (13)$$

is the same as the PDF of the marginal distribution of  $\mathbf{x}_1$  in  $\mathbf{x} = (\mathbf{x}_1^T, x_{k_1+1}, \dots, x_k)^T$ , whose PDF is expressed as

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}_{\{\mathbf{x}: \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\mathbf{x}_1 + \mathbf{x}_2 = \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1\}}(\boldsymbol{\mu}, \tilde{\boldsymbol{\Sigma}}) \\ &= \frac{1}{Z} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \tilde{\boldsymbol{\Sigma}}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \delta(\boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\mathbf{x}_1 + \mathbf{x}_2 = \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1), \end{aligned} \quad (14)$$

where  $\mathbf{x}_2 = (x_{k_1+1}, \dots, x_k)^T$ ,  $Z$  is a normalization constant, and

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^k, \quad \tilde{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \end{bmatrix} \in \mathbb{R}^{k \times k}. \quad (15)$$

Further applying Theorem 1 to Corollary 3, as described in detail in the Appendix, a MVN random variable  $\mathbf{x}$  with a structured covariance matrix can be generated as in Algorithm 3, where there is no need to compute  $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$  and its Cholesky decomposition. Suppose the covariance matrix  $\boldsymbol{\Sigma}_{11}$  admits some special structure that makes it easy to invert and computationally efficient to simulate from  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{11})$ , then Algorithm 3 could lead to a significant saving in computation if  $k_2 \ll k_1$ . On the other

hand, when  $k_2 \gg k_1$  and  $\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$  admits no special structures, Algorithm 3 may not bring any computational advantage and hence one may resort to the naive Cholesky decomposition based procedure. Detailed computational complexity analyses for both methods are provided in Tables 3 and 4 of the Appendix, respectively.

---

**Algorithm 3** Simulation of the MVN distribution

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}).$$

- 
- Sample  $\mathbf{y}_1 \sim \mathcal{N}(\mathbf{0}, \Sigma_{11})$  and  $\mathbf{y}_2 \sim \mathcal{N}(\mathbf{0}, \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})$  ;
  - Return  $\mathbf{x}_1 = \boldsymbol{\mu}_1 + \mathbf{y}_1 - \Sigma_{12}\Sigma_{22}^{-1}(\Sigma_{21}\Sigma_{11}^{-1}\mathbf{y}_1 + \mathbf{y}_2)$ , which can be realized using
    - Solve  $\boldsymbol{\alpha}$  such that  $\Sigma_{22}\boldsymbol{\alpha} = \Sigma_{21}\Sigma_{11}^{-1}\mathbf{y}_1 + \mathbf{y}_2$ ;
    - Return  $\mathbf{x}_1 = \boldsymbol{\mu}_1 + \mathbf{y}_1 - \Sigma_{12}\boldsymbol{\alpha}$ .
- 

**Corollary 4.** *A random variable simulated with Algorithm 3 is distributed as  $\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$ .*

The efficient simulation algorithm for a MVN distribution with a structured covariance matrix can also be further extended to a MVN distribution with a structured precision matrix, as described below, where  $\boldsymbol{\beta} \in \mathbb{R}^p$ ,  $\boldsymbol{\mu}_\beta \in \mathbb{R}^p$ ,  $\Phi \in \mathbb{R}^{n \times p}$ , and both  $\mathbf{A} \in \mathbb{R}^{p \times p}$  and  $\Omega \in \mathbb{R}^{n \times n}$  are positive-definite matrices. Computational complexity analyses for both the naive Cholesky decomposition based implementation and Algorithm 4 are provided in Table 5 and 6 of the Appendix, respectively. Similar to Algorithm 3, Algorithm 4 may bring a significant saving in computation when  $p \gg n$  and  $\mathbf{A}$  admits some special structure that makes it easy to invert and computationally efficient to simulate  $\mathbf{y}_1$ .

---

**Algorithm 4** Simulation of the MVN distribution

$$\boldsymbol{\beta} \sim \mathcal{N}\left[\boldsymbol{\mu}_\beta, (\mathbf{A} + \Phi^T\Omega\Phi)^{-1}\right].$$

- 
- Sample  $\mathbf{y}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$  and  $\mathbf{y}_2 \sim \mathcal{N}(\mathbf{0}, \Omega^{-1})$  ;
  - Return  $\boldsymbol{\beta} = \boldsymbol{\mu}_\beta + \mathbf{y}_1 - \mathbf{A}^{-1}\Phi^T(\Omega^{-1} + \Phi\mathbf{A}^{-1}\Phi^T)^{-1}(\Phi\mathbf{y}_1 + \mathbf{y}_2)$ , which can be realized using
    - Solve  $\boldsymbol{\alpha}$  such that  $(\Omega^{-1} + \Phi\mathbf{A}^{-1}\Phi^T)\boldsymbol{\alpha} = \Phi\mathbf{y}_1 + \mathbf{y}_2$ .
    - Return  $\boldsymbol{\beta} = \boldsymbol{\mu}_\beta + \mathbf{y}_1 - \mathbf{A}^{-1}\Phi^T\boldsymbol{\alpha}$ .
-



**Corollary 5.** *The random variable obtained with Algorithm 4 is distributed as  $\beta \sim \mathcal{N}(\boldsymbol{\mu}_\beta, \boldsymbol{\Sigma}_\beta)$ , where  $\boldsymbol{\Sigma}_\beta = (\mathbf{A} + \boldsymbol{\Phi}^T \boldsymbol{\Omega} \boldsymbol{\Phi})^{-1}$ .*

## 4 Illustrations

Below we provide several examples to illustrate Theorem 1, which shows how to efficiently simulate from a hyperplane-truncated MVN distribution, and Corollary 4 (Corollary 5), which shows how to efficiently simulate from a MVN distribution with a structured covariance (precision) matrix. We run all our experiments on a 2.9 GHz computer.

### 4.1 Simulation of hyperplane-truncated MVNs

We first compare Algorithms 1 and 2, whose generated random samples follow the same distribution, as suggested by Theorem 1, to highlight the advantages of Algorithm 2 over Algorithm 1. We then employ Algorithm 2 for a real application whose data dimension is high and sample size is large.

#### Comparison of Algorithms 1 and 2

We compare Algorithms 1 and 2 in a wide variety of settings by varying the data dimension  $k$ , varying the number of hyperplane constraints  $k_2$ , and choosing either a diagonal covariance matrix  $\boldsymbol{\Sigma}$  or a non-diagonal one. We generate random diagonal covariance matrices using the MATLAB command `diag(0.05 + rand(k, 1))` and random non-diagonal ones using `U.' * diag(0.05 + rand(k, 1)) * U`, where `rand(k, 1)` is a vector of  $k$  uniform random numbers and  $U$  consists of a set of  $k$  orthogonal basis vectors. The elements of  $\boldsymbol{\mu}$ ,  $\mathbf{r}$ , and  $\mathbf{G}$  are all sampled from  $\mathcal{N}(0, 1)$ , with the singular value decomposition applied to  $\mathbf{G}$  to check whether  $\text{Rank}(\mathbf{G}) = k_2$ .

First, to verify Theorem 1, we conduct an experiment with  $k = 5000$  data dimension,  $k_2 = 20$  hyperplanes, and a diagonal  $\boldsymbol{\Sigma}$ . Contour plots of two randomly selected dimensions of the 10,000 random samples simulated with Algorithms 1 and 2 are shown in the top and bottom rows of Figure 3, respectively. The clear matches between the contour plots of these two different algorithms suggest the correctness of Theorem 1.

To demonstrate the efficiency of Algorithm 2, we first carry out a series of experiments with the number of hyperplane constraints fixed at  $k_2 = 20$  and the data dimension increased from  $k = 50$  to  $k = 5000$ . The computation time of simulating 10,000 samples averaged over five random trials is shown in Figure 4(a) for non-diagonal  $\boldsymbol{\Sigma}$ 's and in Figure 4(d) for diagonal ones. It is clear that, when the data dimension  $k$  is high, Algorithm 2 has a clear advantage over Algorithm 1 by avoiding computing unnecessary intermediate variables, which is especially evident when  $\boldsymbol{\Sigma}$  is diagonal. We then carry out a series of experiments where we vary not only  $k$ , but also  $k_2$  from  $0.1k$  to  $0.9k$  for each  $k$ . As shown in Figure 4, it is evident that Algorithm 2 dominates Algorithm 1 in all scenarios, which can be explained by the fact that Algorithm 2 needs to compute much fewer intermediate variables. Also observed is that a larger  $k_2$  leads to slower

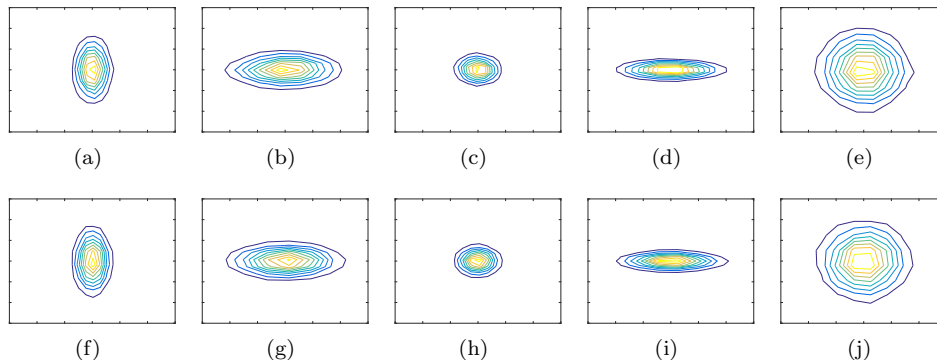


Figure 3: Comparison of the contour plots of two randomly selected dimensions of the 10,000  $k = 5000$  dimensional random samples simulated with Algorithm 1 (top row) and Algorithm 2 (bottom row). Each of the five columns corresponds to a random trial.

simulation for both algorithms, but to a much lesser extent for Algorithm 2. Moreover, the curvatures of those curves indicate that Algorithm 2 is more practical in a high dimensional setting. Note that since Algorithm 2 can naturally exploit the structure of the covariance matrix  $\Sigma$  for fast simulation, it is clearly more capable of benefiting from having a diagonal or block-diagonal  $\Sigma$ , demonstrated by comparing Figures 4(b) and 4(c) with Figures 4(e) and 4(f). All these observations agree with our computational complexity analyses for Algorithms 1 and 2, as shown in Table 1 and 2 of the Appendix, respectively.

### A practical application of Algorithm 2

In what follows, we extend Algorithm 2 to facilitate simulation from a MVN distribution truncated on a probability simplex  $\mathbb{S}^k = \{\mathbf{x} : \mathbf{x} \in \mathbb{R}^k, \mathbf{1}^T \mathbf{x} = 1, x_i \geq 0, i = 1, \dots, k\}$ . This problem frequently arises when unknown parameters can be interpreted as fractions or probabilities, for instance, in topic models (Blei et al. 2003), admixture models (Pritchard et al. 2000; Dobigeon et al. 2009b; Bazot et al. 2013), and discrete directed graphical models (Heckerman 1998). With Algorithm 2, one may remove the equality constraint to greatly simplify the problem.

More specifically, we focus on a big data setting in which the globally shared simplex-constrained model parameters could be linked to some latent counts via the multinomial likelihood. When there are tens of thousands or millions of observations in the dataset, scalable Bayesian inference for the simplex-constrained globally shared model parameters is highly desired, for example, for inferring the topics' distributions over words in latent Dirichlet allocation (Blei et al. 2003; Hoffman et al. 2010) and Poisson factor analysis (Zhou et al. 2012, 2016).

Let us denote the  $\kappa$ th model parameter vector constrained on a  $V$ -dimensional simplex by  $\phi_\kappa \in \mathbb{S}^V$ , which could be linked to the latent counts  $n_{vj\kappa} \in \mathbb{Z}$  of the  $j$ th

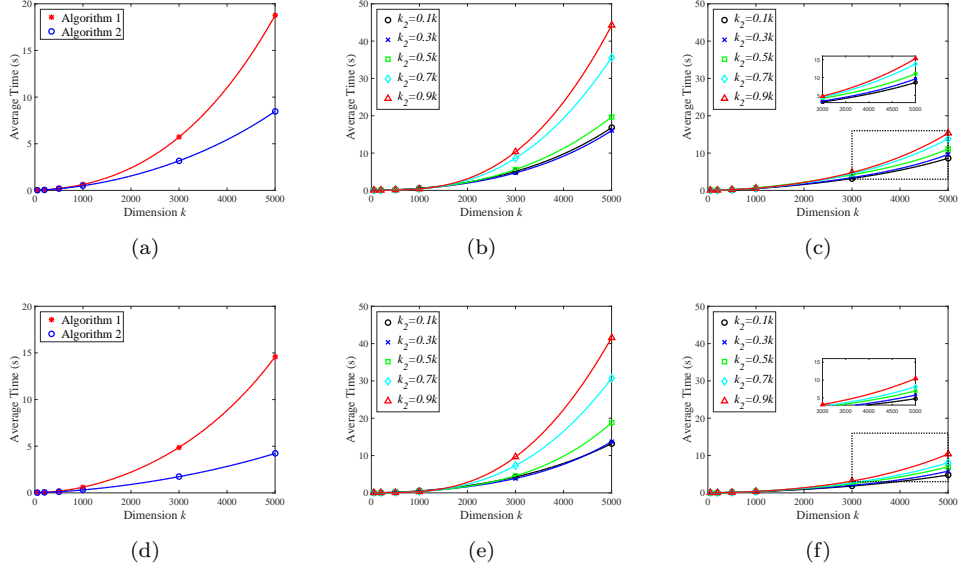


Figure 4: Average time of simulating 10,000 hyperplane-truncated MVN samples over five random trials in different dimensions with non-diagonal covariance matrixes (top row) and diagonal ones (bottom row). (a)(d) Comparison with fixed  $k_2 = 20$ . (b)(e) Algorithm 1 with varying  $k_2$ . (c)(f) Algorithm 2 with varying  $k_2$ .

document under a multinomial likelihood as  $(n_{1j\kappa}, \dots, n_{Vj\kappa}) \sim \text{Mult}(n_{\cdot j\kappa}, \phi_\kappa)$ , where  $\mathbb{Z} = \{0, 1, 2, \dots\}$ ,  $v \in \{1, \dots, V\}$ ,  $\kappa \in \{1, \dots, K\}$ , and  $j \in \{1, \dots, N\}$ . In topic modeling, one may consider  $K$  as the total number of latent topics and  $n_{vj\kappa}$  as the number of words at the  $v$ th vocabulary term in the  $j$ th document that are associated with the  $\kappa$ th latent topic. Note that the dimension  $V$  in real applications is often large, such as tens of thousands in topic modeling. Given the observed counts  $n_{vj}$  for the whole dataset, in a batch-learning setting, one typically iteratively updates the latent counts  $n_{vj\kappa}$  conditioning on  $\phi_\kappa$ , and updates  $\phi_\kappa$  conditioning on  $n_{vj\kappa}$ .

However, this batch-learning inference procedure would become inefficient and even impractical when the dataset size  $N$  grows to a level that makes it too time consuming to finish even a single iteration of updating all local variables  $n_{vj\kappa}$ . To address this issue, we consider constructing a mini-batch based Bayesian inference procedure that could make substantial progress in posterior simulation while the batch-learning one may still be waiting to finish a single iteration.

Without loss of generality, in the following discussion, we drop the latent factor/topic index  $\kappa$  to simplify the notation, focusing on the update of a single simplex-constrained global parameter vector. More specifically, we let the latent local count vector  $\mathbf{n}_j = (n_{1j}, \dots, n_{Vj})^T$  be linked to the simplex-constrained global parameter vector  $\phi \in \mathbb{S}^V$  via the multinomial likelihood as  $\mathbf{n}_j \sim \text{Mult}(n_{\cdot j}, \phi)$ , and impose a Dirichlet distribution

prior on  $\phi$  as  $\phi \sim \text{Dir}(\eta \mathbf{1}_V)$ .

Instead of waiting for all  $\mathbf{n}_j$  to be updated before performing a single update of  $\phi$ , we develop a mini-batch based Bayesian inference algorithm under a general framework for constructing stochastic gradient Markov chain Monte Carlo (SG-MCMC) (Ma et al. 2015), allowing  $\phi$  to be updated every time a mini-batch of  $\mathbf{n}_j$  are processed. For the sake of completeness, we concisely describe the derivation for a SG-MCMC algorithm, as outlined below, for simplex-constrained globally shared model parameters. We refer the readers to Cong et al. (2017) for more details on the derivation and its application to scalable inference for topic modeling.

Using the reduced-mean parameterization of the simplex constrained vector  $\phi$ , namely  $\varphi = (\phi_1, \dots, \phi_{V-1})^T$ , where  $\varphi \in \mathbb{R}_+^{V-1}$  is constrained with  $\varphi \leq 1$ , we develop a SG-MCMC algorithm that updates  $\varphi$  for the  $t$ th mini-batch as

$$\varphi_{t+1} = \left[ \varphi_t + \frac{\varepsilon_t}{M} [(\rho \bar{\mathbf{n}}_{\cdot\cdot} + \eta) - (\rho n_{\cdot\cdot} + \eta V) \varphi_t] + \mathcal{N} \left( \mathbf{0}, \frac{2\varepsilon_t}{M} [\text{diag}(\varphi_t) - \varphi_t \varphi_t^T] \right) \right]_{\Delta}, \quad (16)$$

where  $\varepsilon_t$  are annealed step sizes,  $\rho$  is the ratio of the dataset size  $N$  to the mini-batch size,  $\mathbf{n}_{\cdot\cdot} = (\mathbf{n}_{1\cdot}, \dots, \mathbf{n}_{V\cdot})^T = \sum_{j \in \mathcal{I}_t} \mathbf{n}_j$ ,  $\bar{\mathbf{n}}_{\cdot\cdot} = (\bar{\mathbf{n}}_{1\cdot}, \dots, \bar{\mathbf{n}}_{(V-1)\cdot})^T$ ,  $[\cdot]_{\Delta}$  denotes the constraint that  $\varphi \in \mathbb{R}_+^{V-1}$  and  $\varphi \leq 1$ , and  $M := \mathbb{E} \left[ \sum_{j=1}^N n_{\cdot j} \right]$  is approximated along the updating using  $M = (1 - \varepsilon_t) M + \varepsilon_t \rho \mathbb{E} [n_{\cdot\cdot}]$ . Alternatively, we have an equivalent update equation for  $\phi$  as

$$\phi_{t+1} = \left[ \phi_t + \frac{\varepsilon_t}{M} [(\rho \mathbf{n}_{\cdot\cdot} + \eta) - (\rho n_{\cdot\cdot} + \eta V) \phi_t] + \mathcal{N} \left( \mathbf{0}, \frac{2\varepsilon_t}{M} \text{diag}(\phi_t) \right) \right]_{\angle}, \quad (17)$$

where  $[\cdot]_{\angle}$  represents the constraint that  $\phi \in \mathbb{R}_+^V$  and  $\mathbf{1}^T \phi = 1$ .

It is clear that (16) corresponds to simulation of a  $V - 1$  dimensional truncated MVN distribution with  $V$  inequality constraints. Since the number of constraints is larger than the dimension, previously proposed iterative simulation methods such as the one in Botev (2016) are often inappropriate. Note that, by omitting the non-negative constraints, the update in (17) corresponds to simulation of a hyperplane-truncated MVN simulation with a diagonal covariance matrix, which can be efficiently sampled as described in the following example.

**Example 1:** *Simulation of a hyperplane-truncated MVN distribution as*

$$\mathbf{x} \sim \mathcal{N}_{\mathcal{S}}[\boldsymbol{\mu}, a \text{diag}(\boldsymbol{\phi})], \quad \mathcal{S} = \{\mathbf{x} : \mathbf{1}^T \mathbf{x} = 1\},$$

where  $\mathbf{x} \in \mathbb{R}^k$ ,  $\boldsymbol{\mu} \in \mathbb{R}^k$ ,  $\mathbf{1}^T \mathbf{x} = \sum_{i=1}^k x_i$ ,  $\boldsymbol{\phi} \in \mathbb{R}^k$ ,  $a > 0$ ,  $\phi_i > 0$  for  $i \in \{1, \dots, k\}$ , and  $\mathbf{1}^T \boldsymbol{\phi} = \sum_{i=1}^k \phi_i = 1$ , can be realized as follows.

- Sample  $\mathbf{y} \sim \mathcal{N}[\boldsymbol{\mu}, a \text{diag}(\boldsymbol{\phi})]$ ;
- Return  $\mathbf{x} = \mathbf{y} + (1 - \mathbf{1}^T \mathbf{y}) \boldsymbol{\phi}$ .

The sampling steps in Example 1 directly follow Algorithm 2 and Theorem 1 with the distribution parameters specified as  $\Sigma = \text{adiag}(\phi)$ ,  $\mathbf{G} = \mathbf{1}^T$ , and  $\mathbf{r} = 1$ . Accordingly, we present the following fast sampling procedure for (16).

**Example 2:** *Simulation from (16) can be approximately but rapidly realized as*

- Sample  $\mathbf{y} \sim \mathcal{N}[\phi_t + \frac{\varepsilon_t}{M} [(\rho \mathbf{n}_{\cdot} + \eta) - (\rho \mathbf{n}_{\cdot} + \eta V) \phi_t], \frac{2\varepsilon_t}{M} \text{diag}(\phi_t)]$ ;
- Calculate  $\mathbf{z} = \mathbf{y} + (1 - \mathbf{1}^T \mathbf{y}) \phi_t$ ;
- If  $\mathbf{z} \in \mathbb{S}$ , return  $\varphi_{t+1} = (z_1, \dots, z_{V-1})^T$ ; else calculate  $\mathbf{d} = \max(\varepsilon, \mathbf{z})$  with a small constant  $\varepsilon \geq 0$ , let  $\mathbf{e} = \mathbf{d} / \sum_{i=1}^V d_i$ , and return  $\varphi_{t+1} = (e_1, \dots, e_{V-1})^T$ .

To verify Example 2, we conduct an experiment using multinomial-distributed data vectors of  $V = 2000$  dimensions, which are generated as follows: considering that the simplex-constrained vector  $\phi$  is usually sparse in a high-dimensional application, we sample a  $V = 2000$  dimensional vector  $\mathbf{f}$  whose elements are uniformly distributed between 0 and 1, randomly select 40 dimensions and reset their values to be 100, and set  $\phi = \mathbf{f} / \sum_{i=1}^V f_i$ ; we simulate  $N = 10,000$  samples, each  $\mathbf{n}_j$  of which is generated from the multinomial distribution  $\text{Mult}(n_{\cdot j}, \phi)$ , where the number of trials is random and generated as  $n_{\cdot j} \sim \text{Pois}(50)$ . We set  $\varepsilon_t = t^{-0.99}$  and use mini-batches, each of which consists of 10 data samples, to stochastically update global parameters via SG-MCMC.

For comparison, we choose the same SG-MCMC inference procedure but consider simulating (16), as performed every time a mini-batch of data samples are provided, either as in Example 2 or with the Gibbs sampler of Rodriguez-Yam et al. (2004). Simulating (16) with the Gibbs sampler of Rodriguez-Yam et al. (2004) is realized by updating all the  $V$  dimensions, one dimension at a time, in each Gibbs sampling iteration. We set the total number of Gibbs sampling iterations for (16) in each mini-batch based update as 1, 5, or 10. Note that in practice, the  $\mathbf{n}_j$  belonging to the current mini-batch are often latent and are updated conditioning on the data samples in the mini-batch and  $\phi$ . For simplicity, all  $\mathbf{n}_j$  here are simulated once and then fixed.

Using  $\phi_{post}^* = (\sum_{j=1}^N \mathbf{n}_j + \eta) / (\sum_{j=1}^N n_{\cdot j} + \eta V)$ , the posterior mean of  $\phi$  in a batch-learning setting, as the reference, we show in Figure 5 how the residual errors for the estimated  $\phi^*$ , defined as  $\|\phi^* - \phi\|_2$ , change both as a function of the number of processed mini-batches and as a function of computation time under various settings of the mini-batch based SG-MCMC algorithm. The curves shown in Figure 5 suggest that for each mini-batch, to simulate (16) with the Gibbs sampler of Rodriguez-Yam et al. (2004), it is necessary to have more than one Gibbs sampling iteration to achieve satisfactory results. It is clear from Figure 5(a) that the Gibbs sampler with 5 or 10 iterations for each mini-batch, even though each mini-batch has only 10 data samples, provides residual errors that quickly approach that of the batch posterior mean with a tiny gap, indicating the effectiveness of the SG-MCMC updating in (16). While simulating (16) with Gibbs sampling could in theory lead to unbiased samples if the number of Gibbs sampling iterations is large enough, it is much more efficient to simulate (16) with the procedure described in Example 2, which provides a performance that is undis-

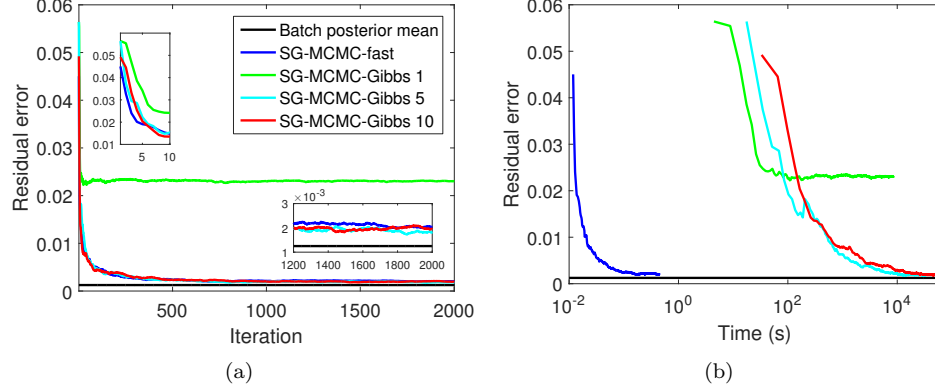


Figure 5: Comparisons of the residual errors of the simplex-constrained parameter vector, estimated under various settings of the stochastic-gradient MCMC (SG-MCMC) algorithm, as a function of (a) the number of processed mini-batches and (b) time. The curves labeled as “Batch posterior mean”, “SG-MCMC-fast”, and “SG-MCMC-Gibbs” correspond to the batch posterior mean, SG-MCMC with (16) simulated as in Example 2, and SG-MCMC with (16) simulated with the Gibbs sampler of Rodriguez-Yam et al. (2004), respectively. The digit following “SG-MCMC-Gibbs” represents the number of Gibbs sampling iterations to simulate (16) for each mini-batch.

tinguishable from those of the Gibbs sampler with as many as 5 or 10 iterations for each mini-batch, but at the expense of a tiny fraction of a single Gibbs sampling iteration.

## 4.2 Simulation of MVNs with structured covariance matrices

To illustrate Corollary 4, we mimic the truncated MVN simulation in (16) and present the following simulation example with a structured covariance matrix.

**Example 3:** Simulation of a MVN distribution as

$$\mathbf{x}_1 \sim \mathcal{N}[\boldsymbol{\mu}_1, a \text{diag}(\boldsymbol{\phi}_1) - a \boldsymbol{\phi}_1 \boldsymbol{\phi}_1^T],$$

where  $\mathbf{x}_1 \in \mathbb{R}^{k-1}$ ,  $\boldsymbol{\mu}_1 \in \mathbb{R}^{k-1}$ ,  $a > 0$ ,  $\boldsymbol{\phi}_1 = (\phi_1, \dots, \phi_{k-1})^T$ ,  $\phi_i > 0$  for  $i \in \{1, \dots, k-1\}$ , and  $\sum_{i=1}^{k-1} \phi_i < 1$ , can be realized as follows.

- Sample  $\mathbf{y}_1 \sim \mathcal{N}[\mathbf{0}, a \text{diag}(\boldsymbol{\phi}_1)]$  and  $\mathbf{y}_2 \sim \mathcal{N}(0, a^{-1} \phi_k)$ , where  $\phi_k = 1 - \sum_{i=1}^{k-1} \phi_i$ ;
- Return  $\mathbf{x}_1 = \boldsymbol{\mu}_1 + \mathbf{y}_1 - (\mathbf{1}^T \mathbf{y}_1 + a \mathbf{y}_2) \boldsymbol{\phi}_1$ .

Denoting  $\mathbf{x} = (\mathbf{x}_1^T, x_k)^T$ ,  $\boldsymbol{\phi} = (\boldsymbol{\phi}_1^T, \phi_k)^T$ ,  $\boldsymbol{\mu} = (\boldsymbol{\mu}_1^T, \mu_k)^T$ , and  $\mu_k = 1 - \mathbf{1}^T \boldsymbol{\mu}_1$ , the above sampling steps can also be equivalently expressed as follows.

- Sample  $\mathbf{y} \sim \mathcal{N}[\boldsymbol{\mu}, a \text{diag}(\boldsymbol{\phi})]$ ;

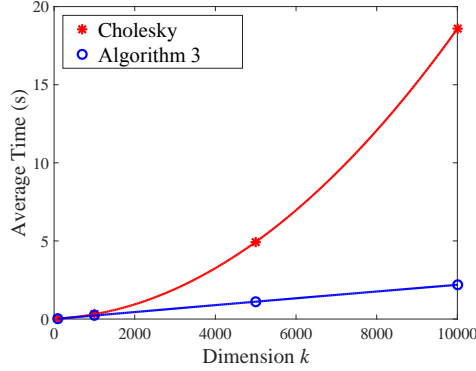


Figure 6: Comparison of the naive Cholesky decomposition based implementation and Algorithm 3 in terms of the average time of generating 10,000  $k$ -dimensional random samples from  $\mathbf{x}_1 \sim \mathcal{N}[\boldsymbol{\mu}_1, a \text{diag}(\boldsymbol{\phi}_1) - a \boldsymbol{\phi}_1 \boldsymbol{\phi}_1^T]$ . The distribution parameters are randomly generated and computation time averaged over 100 random trials is displayed.

- Return  $\mathbf{x}_1 = \mathbf{y}_1 + (1 - \mathbf{1}^T \mathbf{y}) \boldsymbol{\phi}_1$ .

Directly following Algorithm 3 and Corollary 4, the first sampling approach for the above example can be derived by specifying the distribution parameters as  $\boldsymbol{\Sigma}_{11} = a \text{diag}(\boldsymbol{\phi}_1)$ ,  $\boldsymbol{\Sigma}_{12} = \boldsymbol{\phi}_1$ ,  $\boldsymbol{\Sigma}_{21} = \boldsymbol{\phi}_1^T$ , and  $\boldsymbol{\Sigma}_{22} = a^{-1}$ , while the second approach can be derived by specifying  $\boldsymbol{\Sigma}_{11} = a \text{diag}(\boldsymbol{\phi}_1)$ ,  $\boldsymbol{\Sigma}_{12} = a \boldsymbol{\phi}_1$ ,  $\boldsymbol{\Sigma}_{21} = a \boldsymbol{\phi}_1^T$ , and  $\boldsymbol{\Sigma}_{22} = a$ .

To illustrate the efficiency of the proposed algorithms in Example 3, we simulate from the MVN distribution  $\mathbf{x}_1 \sim \mathcal{N}[\boldsymbol{\mu}_1, a \text{diag}(\boldsymbol{\phi}_1) - a \boldsymbol{\phi}_1 \boldsymbol{\phi}_1^T]$  using both a naive implementation via Cholesky decomposition of the covariance matrix and the fast simulation algorithm for a hyperplane-truncated MVN random variable described in Example 3. We set the dimension from  $k = 10^2$  up to  $k = 10^4$  and set  $\boldsymbol{\mu} = (1/k, \dots, 1/k)$  and  $a = 0.5$ . For each  $k$  and each simulation algorithm, we perform 100 independent random trials, in each of which  $\boldsymbol{\phi}$  is sampled from the Dirichlet distribution  $\text{Dir}(1, \dots, 1)$  and 10,000 independent random samples are simulated using that same  $\boldsymbol{\phi}$ .

As shown in Figure 6, for the proposed Algorithm 3, the average time of simulating 10,000 random samples increases linearly in the dimension  $k$ . By contrast, for the naive Cholesky decomposition based simulation algorithm, whose computational complexity is  $O(k^3)$  (Golub and Van Loan 2012), the average simulation time increases at a significantly faster rate as the dimension  $k$  increases.

For explicit verification, with the 10,000 simulated  $k = 10^4$  dimensional random samples in a random trial, we randomly choose two dimensions and display their joint distribution using a contour plot. As in Figure 7, shown in the first row are the contour plots of five different random trials for the naive Cholesky implementation, whereas shown in the second row are the corresponding ones for the proposed Algorithm 3. As expected, the contour lines of the two figures in the same column closely match each other.

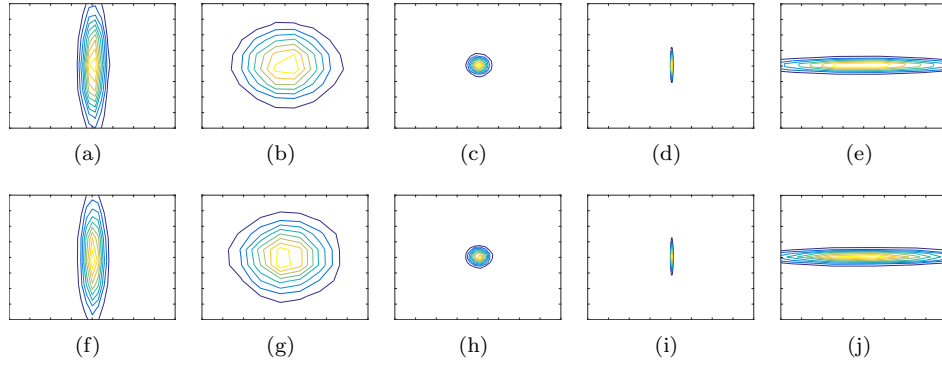


Figure 7: Comparison of the contour plots of two randomly selected dimensions of the 10,000  $k = 10^4$  dimensional random samples simulated with the naive Cholesky implementation (top row) and Algorithm 3 (bottom row). Each of the five columns corresponds to a random trial.

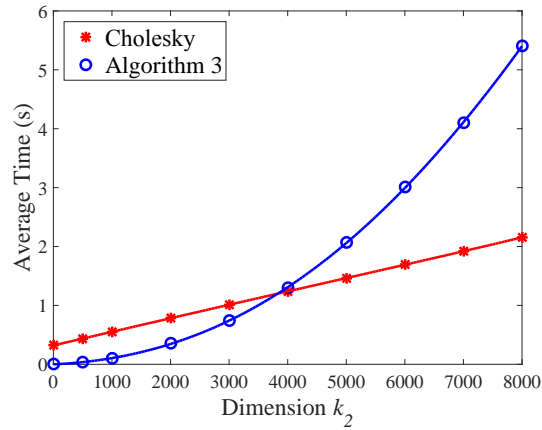


Figure 8: Comparison of the naive Cholesky decomposition based implementation and Algorithm 3 in terms of the average time of generating one  $k_1 = 4000$  dimensional sample from  $\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})$ , with diagonal  $\boldsymbol{\Sigma}_{11}$  and  $\boldsymbol{\Sigma}_{22}$ . The distribution parameters are randomly generated and computation time averaged over 50 random trials is displayed.

To further examine when to apply Algorithm 3 instead of the naive Cholesky decomposition based implementation in a general setting, we present the computational complexity analyses in Tables 3 and 4 of the Appendix for the naive approach and Algorithm 3, respectively. In addition, we mimic the settings in Section 4.1 to conduct a set of experiments with randomly generated  $\boldsymbol{\Sigma}_{12}$ , diagonal  $\boldsymbol{\Sigma}_{11}$ , and diagonal  $\boldsymbol{\Sigma}_{22}$ . We fix  $k_1 = 4000$  and vary  $k_2$  from 1 to 8000. The computation time for one sample averaged over 50 random trials is presented in Figure 8. It is clear from Tables 3 and 4 and Figure 8 that, as a general guideline, one may choose Algorithm 3 when  $k_2$  is



smaller than  $k_1$  and  $\Sigma_{11}$  admits some special structure that makes it easy to invert and computationally efficient to simulate from  $\mathcal{N}(\mathbf{0}, \Sigma_{11})$ .

### 4.3 Simulation of MVNs with structured precision matrices

To examine when to apply Algorithm 4 instead of the naive Cholesky decomposition based procedure, we first consider a series of random simulations in which the sample size  $n$  is fixed while the data dimension  $p$  is varying. We then show that Algorithm 4 can be applied for high-dimensional regression whose  $p$  is often much larger than  $n$ .

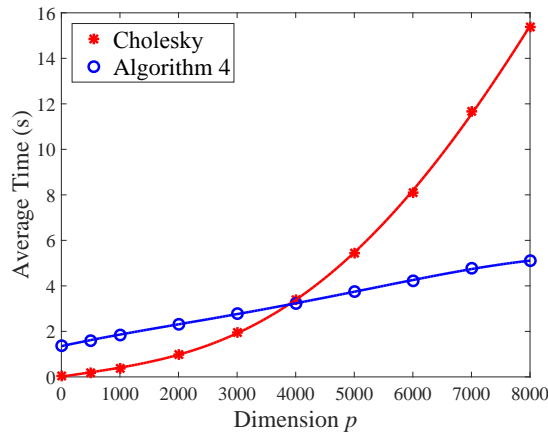


Figure 9: Comparison of the naive Cholesky decomposition based implementation and Algorithm 4 in terms of the average time of generating one  $p$  dimensional sample from  $\beta \sim \mathcal{N}[\mu_\beta, (\mathbf{A} + \Phi^T \Omega \Phi)^{-1}]$ , with diagonal  $\mathbf{A}$  and  $\Omega$ . The distribution parameters are randomly generated and computation time averaged over 50 random trials is displayed.

We fix  $n = 4000$ , vary  $p$  from 1 to 8000, and mimic the settings in Section 4.1 to randomly generate  $\Phi$ , diagonal  $\mathbf{A}$ , and diagonal  $\Omega$ . As a function of dimensions  $p$ , the computation time for one sample averaged over 50 random trials is shown in Figure 9. It is evident that, identical to the complexity analysis in Tables 5 and 6, Algorithm 4 has a linear complexity with respect to  $p$  under these settings, which will bring significant acceleration in a high-dimensional setting with  $p \gg n$ . If the sample size  $n$  is large enough that  $n > p$ , then one may directly apply the naive Cholesky decomposition based implementation.

Algorithm 4 could be slightly modified to be applied to high-dimensional regression, where the main objective is to efficiently sample from the conditional posterior of  $\beta \in \mathbb{R}^{p \times 1}$  in the linear regression model as

$$t \sim \mathcal{N}(\Phi \beta, \Omega^{-1}), \quad \beta \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \quad (18)$$

where  $\Phi \in \mathbb{R}^{n \times p}$ ,  $\Omega \in \mathbb{R}^{n \times n}$ , and different constructions on  $\mathbf{A} \in \mathbb{R}^{p \times p}$  lead to a wide variety of regression models (Caron and Doucet 2008; Carvalho et al. 2010; Polson et al.

2014). The conditional posterior of  $\beta$  is directly derived and shown in the following example, where its simulation algorithm is summarized by further generalizing Corollary 5.

**Example 4:** *Simulation of the MVN distribution*

$$\beta \sim \mathcal{N} \left[ (\mathbf{A} + \Phi^T \Omega \Phi)^{-1} \Phi^T \Omega t, (\mathbf{A} + \Phi^T \Omega \Phi)^{-1} \right]$$

can be realized as follows.

- Sample  $\mathbf{y}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$  and  $\mathbf{y}_2 \sim \mathcal{N}(\mathbf{0}, \Omega^{-1})$  ;
- Return  $\beta = \mathbf{y}_1 + \mathbf{A}^{-1} \Phi^T (\Omega^{-1} + \Phi \mathbf{A}^{-1} \Phi^T)^{-1} (t - \Phi \mathbf{y}_1 - \mathbf{y}_2)$ , which can be realized using
  - Solve  $\alpha$  such that  $(\Omega^{-1} + \Phi \mathbf{A}^{-1} \Phi^T) \alpha = t - \Phi \mathbf{y}_1 - \mathbf{y}_2$ ;
  - Return  $\beta = \mathbf{y}_1 + \mathbf{A}^{-1} \Phi^T \alpha$ .

Note that if  $\Omega = \mathbf{I}_n$ , then the simulation algorithm in Example 4 reduces to the one in Proposition 2.1 of Bhattacharya et al. (2016), which is shown there to be significantly more efficient than that of Rue (2001) for high-dimensional regression if  $p \gg n$ .

## 5 Conclusions

A fast and exact simulation algorithm is developed for a multivariate normal (MVN) distribution whose sample space is constrained on the intersection of a set of hyperplanes, which is shown to be inherently related to the conditional distribution of a unconstrained MVN distribution. The proposed simulation algorithm is further generalized to efficiently simulate from a MVN distribution, whose covariance (precision) matrix can be decomposed as the sum (difference) of a positive-definite matrix and a low-rank symmetric matrix, using a higher dimensional hyperplane-truncated MVN distribution whose covariance matrix is block-diagonal.

## References

- Albert, J. H. and Chib, S. (1993). “Bayesian analysis of binary and polychotomous response data.” *J. Amer. Statist. Assoc.*, 88(422): 669–679.
- Altmann, Y., McLaughlin, S., and Dobigeon, N. (2014). “Sampling from a multivariate Gaussian distribution truncated on a simplex: a review.” In *Statistical Signal Processing (SSP), 2014 IEEE Workshop on*, 113–116. IEEE.
- Bazot, C., Dobigeon, N., Tourneret, J.-Y., Zaas, A. K., Ginsburg, G. S., and Hero III, A. O. (2013). “Unsupervised Bayesian linear unmixing of gene expression microarrays.” *BMC Bioinformatics*, 14(1): 1.

- Bhattacharya, A., Chakraborty, A., and Mallick, B. K. (2016). “Fast sampling with Gaussian scale mixture priors in high-dimensional regression.” *Biometrika*, 103(4): 985.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). “Latent Dirichlet allocation.” *JMLR*, 3: 993–1022.
- Botev, Z. (2016). “The normal law under linear restrictions: simulation and estimation via minimax tilting.” *J. Roy. Statist. Soc.: Series B*.
- Caron, F. and Doucet, A. (2008). “Sparse Bayesian nonparametric regression.” In *ICML*, 88–95. ACM.
- Carvalho, C. M., Polson, N. G., and Scott, J. G. (2010). “The horseshoe estimator for sparse signals.” *Biometrika*, 97(2): 465–480.
- Chopin, N. (2011). “Fast simulation of truncated Gaussian distributions.” *Statistics and Computing*, 21(2): 275–288.
- Cong, Y., Chen, B., and Zhou, M. (2017). “Topic-layer-adaptive stochastic gradient Riemannian (TLASGR) MCMC for deep latent Dirichlet allocation.” *Preprint*.
- Damien, P. and Walker, S. G. (2001). “Sampling truncated normal, beta, and gamma densities.” *Journal of Computational and Graphical Statistics*, 10(2): 206–215.
- Dobigeon, N., Moussaoui, S., Coulon, M., Tourneret, J.-Y., and Hero, A. O. (2009a). “Joint Bayesian endmember extraction and linear unmixing for hyperspectral imagery.” *IEEE Transactions on Signal Processing*, 57(11): 4355–4368.
- Dobigeon, N., Moussaoui, S., Tourneret, J.-Y., and Carteret, C. (2009b). “Bayesian separation of spectral sources under non-negativity and full additivity constraints.” *Signal Processing*, 89(12): 2657–2669.
- Doucet, A. (2010). “A note on efficient conditional simulation of Gaussian distributions.” *Departments of Computer Science and Statistics, University of British Columbia*.
- Gelfand, A. E., Smith, A. F., and Lee, T.-M. (1992). “Bayesian analysis of constrained parameter and truncated data problems using Gibbs sampling.” *J. Amer. Statist. Assoc.*, 87(418): 523–532.
- Gelfand, A. E. and Smith, A. F. M. (1990). “Sampling-based approaches to calculating marginal densities.” *J. Amer. Statist. Assoc.*, 85(410): 398–409.
- Geman, S. and Geman, D. (1984). “Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images.” *IEEE Trans. Pattern Anal. Mach. Intell.*, 721–741.
- Geweke, J. (1991). “Efficient simulation from the multivariate normal and student-t distributions subject to linear constraints and the evaluation of constraint probabilities.” In *Computing science and statistics: Proceedings of the 23rd symposium on the interface*, 571–578.

- Geweke, J. F. (1996). “Bayesian inference for linear models subject to linear inequality constraints.” In *Modelling and Prediction Honoring Seymour Geisser*, 248–263. Springer.
- Girolami, M. and Calderhead, B. (2011). “Riemann manifold Langevin and Hamiltonian Monte Carlo methods.” *J. Roy. Statist. Soc.: B*, 73(2): 123–214.
- Golub, G. H. and Van Loan, C. F. (2012). *Matrix computations*, volume 3. JHU Press.
- Heckerman, D. (1998). “A tutorial on learning with Bayesian networks.” In *Learning in graphical models*, 301–354. Springer.
- Hoffman, M., Blei, D., and Bach, F. (2010). “Online learning for latent Dirichlet allocation.” In *NIPS*.
- Hoffman, Y. and Ribak, E. (1991). “Constrained realizations of Gaussian fields-A simple algorithm.” *The Astrophysical Journal*, 380: L5–L8.
- Holmes, C. C. and Held, L. (2006). “Bayesian auxiliary variable models for binary and multinomial regression.” *Bayesian Analysis*, 1(1): 145–168.
- Imai, K. and van Dyk, D. A. (2005). “A Bayesian analysis of the multinomial probit model using marginal data augmentation.” *Journal of Econometrics*, 124(2): 311–334.
- Johndrow, J., Dunson, D., and Lum, K. (2013). “Diagonal orthant multinomial probit models.” In *AISTATS*, 29–38.
- Lan, S., Zhou, B., and Shahbaba, B. (2014). “Spherical Hamiltonian Monte Carlo for constrained target distributions.” In *ICML*, 629–637.
- Ma, Y., Chen, T., and Fox, E. (2015). “A complete recipe for stochastic gradient MCMC.” In *NIPS*, 2899–2907.
- McCulloch, R. E., Polson, N. G., and Rossi, P. E. (2000). “A Bayesian analysis of the multinomial probit model with fully identified parameters.” *Journal of Econometrics*, 99(1): 173–193.
- Neelon, B. and Dunson, D. B. (2004). “Bayesian isotonic regression and trend analysis.” *Biometrics*, 60(2): 398–406.
- Pakman, A. and Paninski, L. (2014). “Exact Hamiltonian Monte Carlo for truncated multivariate Gaussians.” *Journal of Computational and Graphical Statistics*, 23(2): 518–542.
- Patterson, S. and Teh, Y. W. (2013). “Stochastic gradient Riemannian Langevin dynamics on the probability simplex.” In *NIPS*, 3102–3110.
- Polson, N. G., Scott, J. G., and Windle, J. (2014). “The Bayesian bridge.” *J. Roy. Statist. Soc.: Series B*, 76(4): 713–733.

- Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). “Inference of population structure using multilocus genotype data.” *Genetics*, 155(2): 945–959.
- Robert, C. P. (1995). “Simulation of truncated normal variables.” *Statistics and Computing*, 5(2): 121–125.
- Rodriguez-Yam, G., Davis, R. A., and Scharf, L. L. (2004). “Efficient Gibbs sampling of truncated multivariate normal with application to constrained linear regression.” *Technical report*.
- Rue, H. (2001). “Fast sampling of Gaussian Markov random fields.” *J. Roy. Statist. Soc.: Series B*, 63(2): 325–338.
- Schmidt, M. (2009). “Linearly constrained Bayesian matrix factorization for blind source separation.” In *NIPS*, 1624–1632.
- Tong, Y. L. (2012). *The multivariate normal distribution*. Springer Science & Business Media.
- Train, K. E. (2009). *Discrete choice methods with simulation*. Cambridge University Press.
- Zhou, M., Cong, Y., and Chen, B. (2016). “Augmentable Gamma Belief Networks.” *J. Mach. Learn. Res.*, 17(163): 1–44.
- Zhou, M., Hannah, L., Dunson, D. B., and Carin, L. (2012). “Beta-negative binomial process and Poisson factor analysis.” In *AISTATS*, 1462–1471.

### Acknowledgments

The authors would like to thank the editor-in-chief, editor, associate editor, and two anonymous referees for their comments and suggestions, which have helped us improve the paper substantially. M. Zhou thanks Yingbo Li and Xiaojing Wang for helpful discussions. B. Chen thanks the support of the Thousand Young Talent Program of China, NSFC (61372132), NCET-13-0945, and NDPR-9140A07010115DZ01015.

## Appendix

---

**Algorithm 5** (Hoffman and Ribak 1991; Doucet 2010) Simulation of the conditional distribution of  $\mathbf{x}_1$  given  $\mathbf{x}_2 = \mathbf{r}$  as  $\mathbf{x}_1 | \mathbf{x}_2 = \mathbf{r} \sim \mathcal{N}[\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{r} - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}]$ , where the joint distribution of  $\mathbf{x} = (\mathbf{x}_1^T, \mathbf{x}_2^T)$  follows  $\mathbf{x} \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

---

- Sample  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and denote  $\mathbf{y}_1 = (y_1, \dots, y_{k_1})^T$  and  $\mathbf{y}_2 = (y_{k_1+1}, \dots, y_k)^T$ ;
  - Return  $\mathbf{x}_1 = \mathbf{y}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{r} - \mathbf{y}_2)$ .
- 

### Proofs

*Proof of Theorem 1.* Let us denote  $\boldsymbol{\Lambda} = \mathbf{H}^T\boldsymbol{\Sigma}^{-1}\mathbf{H}$  as a precision matrix that can be partitioned as in (7). For Algorithm 1, instead of directly simulating  $\mathbf{z}_1$  given  $\mathbf{z}_2$  using the conditional distribution of the MVN, we apply Algorithm 5 (Hoffman and Ribak 1991; Doucet 2010) to modify its sampling steps as follows.

- Sample  $\tilde{\mathbf{z}} \sim \mathcal{N}[\mathbf{H}^{-1}\boldsymbol{\mu}, \mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^{-1})^T]$ , and denote  $\tilde{\mathbf{z}}_1 = (z_1, \dots, z_{k_1})^T$  and  $\tilde{\mathbf{z}}_2 = (z_{k_1+1}, \dots, z_k)$ ;
- Let  $\mathbf{z} = (\mathbf{z}_1^T, \mathbf{z}_2^T)^T$ , where  $\mathbf{z}_2 = (\mathbf{GH}_2)^{-1}\mathbf{r}$  and  $\mathbf{z}_1 = \tilde{\mathbf{z}}_1 - \boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12}(\mathbf{z}_2 - \tilde{\mathbf{z}}_2)$ , and return

$$\begin{aligned} \mathbf{x} &= \mathbf{Hz} = \mathbf{H}_1\mathbf{z}_1 + \mathbf{H}_2(\mathbf{GH}_2)^{-1}\mathbf{r} \\ &= \mathbf{H}_1\tilde{\mathbf{z}}_1 + \mathbf{H}_1\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12}\tilde{\mathbf{z}}_2 + (\mathbf{H}_2 - \mathbf{H}_1\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12})(\mathbf{GH}_2)^{-1}\mathbf{r} \\ &= (\mathbf{H}_1, \mathbf{H}_1\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12})\tilde{\mathbf{z}} + (\mathbf{H}_2 - \mathbf{H}_1\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12})(\mathbf{GH}_2)^{-1}\mathbf{r}. \end{aligned} \quad (19)$$

Therefore, we can equivalently generate  $\mathbf{x}$  as follows.

- Sample  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .
- Return  $\mathbf{x} = (\mathbf{H}_1, \mathbf{H}_1\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12})\mathbf{H}^{-1}\mathbf{y} + (\mathbf{H}_2 - \mathbf{H}_1\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12})(\mathbf{GH}_2)^{-1}\mathbf{r}$ .

The computation can be significantly simplified if  $\boldsymbol{\Lambda}_{12} = \mathbf{0}$ , which means

$$\boldsymbol{\Lambda}_{12} = \mathbf{H}_1^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_2 = \mathbf{0}.$$

Since  $\mathbf{H}_1^T\mathbf{G}^T = \mathbf{0}$  by definition, to make  $\boldsymbol{\Lambda}_{12} = \mathbf{0}$ , if and only if we have  $\mathbf{H}_2$  as

$$\mathbf{H}_2 = \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M},$$

where  $\mathbf{M} \in \mathbb{R}^{k_2 \times k_2}$  is an arbitrary full rank matrix, under which we have

- Sample  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .
- Return  $\mathbf{x} = (\mathbf{H}_1, \mathbf{0}_{k \times k_2})\mathbf{H}^{-1}\mathbf{y} + \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r}$ , or return
 
$$\mathbf{x} = \mathbf{y} - (\mathbf{0}_{k \times k_1}, \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M})\mathbf{H}^{-1}\mathbf{y} + \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r}. \quad (20)$$

Let us denote  $(\mathbf{0}_{k \times k_1}, \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M})\mathbf{H}^{-1} = \mathbf{C}$ . We have

$$(\mathbf{0}_{k \times k_1}, \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M}) = \mathbf{C}\mathbf{H} = (\mathbf{C}\mathbf{H}_1, \mathbf{C}\boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M})$$

and hence  $\mathbf{C}\mathbf{H}_1 = \mathbf{0}$  and  $\mathbf{C}\boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M} = \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M}$ . Since  $\mathbf{G}\mathbf{H}_1 = \mathbf{0}$ , we have  $\mathbf{C} = \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{G}$ . The proof is completed by substituting  $(\mathbf{0}_{k \times k_1}, \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M})\mathbf{H}^{-1}$  in (20) with  $\boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{G}$ .  $\square$

*Alternative Proof of Theorem 1.* To solve the problem in (3), one may solve an equivalent problem in (6) by defining an invertible transformation matrix  $\mathbf{H}$  that satisfies  $\mathbf{G}\mathbf{H}_1 = \mathbf{0}_{k_2 \times k_1}$ . Let us denote  $\boldsymbol{\Lambda} = \mathbf{H}^T\boldsymbol{\Sigma}^{-1}\mathbf{H}$  as a precision matrix that can be partitioned as in (7). To simply the problem in (6), we choose the transformation matrix  $\mathbf{H}$  to make  $\mathbf{z}_1$  and  $\mathbf{z}_2$  be independent to each other. Since  $\mathbf{z}$  follows a MVN distribution,  $\mathbf{z}_1$  and  $\mathbf{z}_2$  are independent to each other if and only if

$$\boldsymbol{\Lambda}_{12} = \mathbf{H}_1^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_2 = \mathbf{0}.$$

Since  $\mathbf{H}_1^T\mathbf{G}^T = \mathbf{0}$  by definition, to make  $\boldsymbol{\Lambda}_{12} = \mathbf{0}$ , if and only if we have  $\mathbf{H}_2$  as

$$\mathbf{H}_2 = \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M},$$

where  $\mathbf{M} \in \mathbb{R}^{k_2 \times k_2}$  is an arbitrary full rank matrix. Accordingly, we have

$$\mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^{-1})^T = \begin{bmatrix} (\mathbf{H}_1^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_1)^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{H}_2^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_2)^{-1} \end{bmatrix}.$$

Thus with  $\mathbf{H}$  satisfying  $\mathbf{G}\mathbf{H}_1 = \mathbf{0}_{k_2 \times k_1}$  and  $\mathbf{H}_2 = \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M}$ , one can transform the original problem in (3) to that in (6), where  $\mathbf{z}_1$  and  $\mathbf{z}_2$  are independent and the restrictions  $\mathbf{G}\mathbf{x} = \mathbf{r}$  and  $\mathbf{z}_2 = (\mathbf{G}\mathbf{H}_2)^{-1}\mathbf{r}$  imply each other. Following the naive approach shown in Algorithm 1, one can generate  $\mathbf{x}$  from (3) as follows

- Find  $\mathbf{H} = (\mathbf{H}_1, \mathbf{H}_2)$  with  $\mathbf{H}_2 = \boldsymbol{\Sigma}\mathbf{G}^T\mathbf{M}$  and with  $\mathbf{H}_1$  satisfying  $\mathbf{G}\mathbf{H}_1 = \mathbf{0}_{k_2 \times k_1}$ ;
- Sample  $\mathbf{z}_1 \sim \mathcal{N}[(\mathbf{I}_{k_1}, \mathbf{0}_{k_1 \times k_2})\mathbf{H}^{-1}\boldsymbol{\mu}, (\mathbf{H}_1^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_1)^{-1}]$ ;
- Return  $\mathbf{x} = \mathbf{H} \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{M}^{-1}(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} \end{bmatrix}$ .

However, this naive approach contains intermediate variables that could be computationally expensive to compute. Below we present a method to bypass these intermediate variables. Since the last step could be reexpressed as

$$\begin{aligned} \mathbf{x} &= \mathbf{H}_1\mathbf{z}_1 + \mathbf{H}_2\mathbf{M}^{-1}(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} \\ &= \mathbf{H}_1\mathbf{z}_1 + \mathbf{H}_2\mathbf{z}_2 + \mathbf{H}_2[\mathbf{M}^{-1}(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} - \mathbf{z}_2] \\ &= \mathbf{H}\mathbf{z} + \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} - \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{M}\mathbf{z}_2), \end{aligned}$$

where  $\mathbf{z} = (\mathbf{z}_1^T, \mathbf{z}_2^T)^T$  and  $\mathbf{z}_2 \in \mathbb{R}^{k_2}$  is a vector whose values can be chosen arbitrarily. In addition, since  $\mathbf{M}$  is an arbitrary full-rank matrix, we can let

$$\mathbf{z}_2 \sim \mathcal{N}[(\mathbf{0}_{k_2 \times k_1}, \mathbf{I}_{k_2})\mathbf{H}^{-1}\boldsymbol{\mu}, (\mathbf{H}_2^T\boldsymbol{\Sigma}^{-1}\mathbf{H}_2)^{-1}],$$

which means  $\mathbf{z} \sim \mathcal{N}[\mathbf{H}^{-1}\boldsymbol{\mu}, \mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^{-1})^T]$ , and choose  $\mathbf{M}$  to make

$$\mathbf{G}\mathbf{x} = \mathbf{G}\mathbf{H}\mathbf{z} + \mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}\mathbf{r} - \mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{M}\mathbf{z}_2) = \mathbf{r},$$

which means  $\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{M}\mathbf{z}_2) = \mathbf{G}\mathbf{H}\mathbf{z}$ . Thus we have

$$\mathbf{x} = \mathbf{H}\mathbf{z} + \boldsymbol{\Sigma}\mathbf{G}^T(\mathbf{G}\boldsymbol{\Sigma}\mathbf{G}^T)^{-1}(\mathbf{r} - \mathbf{G}\mathbf{H}\mathbf{z}).$$

In addition, since if  $\mathbf{z} \sim \mathcal{N}[\mathbf{H}^{-1}\boldsymbol{\mu}, \mathbf{H}^{-1}\boldsymbol{\Sigma}(\mathbf{H}^{-1})^T]$ , then  $\mathbf{y} = \mathbf{H}\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Therefore, without the need to compute any intermediate variables, one may use Algorithm 2 to generate  $\mathbf{x}$  from the hyperplane truncated MVN distribution.  $\square$

*Proof of Theorem 2.* Using the matrix inversion lemma on (8), we have

$$\begin{aligned} p(\mathbf{x}_1) &\propto \exp\left[-\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu}_1)^T (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1)\right], \\ &\propto \exp\left\{-\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \left[\boldsymbol{\Sigma}_{11}^{-1} + \boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} (\boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})^{-1} \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\right] (\mathbf{x}_1 - \boldsymbol{\mu}_1)\right\}. \end{aligned} \quad (21)$$

Using (11), we have  $\mathbf{G}\boldsymbol{\mu} = \mathbf{G}_1\boldsymbol{\mu}_1 + \mathbf{G}_2\boldsymbol{\mu}_2 = \mathbf{r}$ . Since  $\mathbf{G}\mathbf{x} = \mathbf{r}$ , we further have  $\mathbf{G}(\mathbf{x} - \boldsymbol{\mu}) = 0$  and hence  $\mathbf{G}_1(\mathbf{x}_1 - \boldsymbol{\mu}_1) = -\mathbf{G}_2(\mathbf{x}_2 - \boldsymbol{\mu}_2)$ . Therefore, given the construction of  $\boldsymbol{\mu}_2$  as in (11), we can replace the equality constraint  $\mathbf{G}\mathbf{x} = \mathbf{r}$  on  $\mathbf{x}$  by requiring  $(\mathbf{x}_2 - \boldsymbol{\mu}_2) = -\mathbf{G}_2^{-1}\mathbf{G}_1(\mathbf{x}_1 - \boldsymbol{\mu}_1)$ . Using this equivalent constraint together with (3), we have

$$\begin{aligned} &p(\mathbf{x} \mid \boldsymbol{\mu}, \tilde{\boldsymbol{\Sigma}}, \mathbf{G}, \mathbf{r}) \\ &\propto \exp\left[-\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_{11}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1) - \frac{1}{2}(\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \tilde{\boldsymbol{\Sigma}}_{22}^{-1}(\mathbf{x}_2 - \tilde{\boldsymbol{\mu}}_2)\right] \delta(\mathbf{G}\mathbf{x} = \mathbf{r}) \\ &\propto \exp\left\{-\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \left[\boldsymbol{\Sigma}_{11}^{-1} + \mathbf{G}_1^T(\mathbf{G}_2^{-1})^T \tilde{\boldsymbol{\Sigma}}_{22}^{-1} \mathbf{G}_2^{-1} \mathbf{G}_1\right] (\mathbf{x}_1 - \boldsymbol{\mu}_1)\right\} \\ &\quad \times \delta[\mathbf{x}_2 = \boldsymbol{\mu}_2 - \mathbf{G}_2^{-1}\mathbf{G}_1(\mathbf{x}_1 - \boldsymbol{\mu}_1)] \\ &= \mathcal{N}\left\{\mathbf{x}_1; \boldsymbol{\mu}_1, \left[\boldsymbol{\Sigma}_{11}^{-1} + \mathbf{G}_1^T(\mathbf{G}_2^{-1})^T \tilde{\boldsymbol{\Sigma}}_{22}^{-1} \mathbf{G}_2^{-1} \mathbf{G}_1\right]^{-1}\right\} \\ &\quad \times \delta[\mathbf{x}_2 = \boldsymbol{\mu}_2 - \mathbf{G}_2^{-1}\mathbf{G}_1(\mathbf{x}_1 - \boldsymbol{\mu}_1)] \end{aligned} \quad (22)$$

It is clear that the marginal distribution of  $\mathbf{x}_1$  in (22) matches the conditional distribution of  $\mathbf{x}_1$  in (21) if we further construct  $\tilde{\boldsymbol{\Sigma}}_{22}$  using (12).  $\square$



*Proof of Corollary 4.* Applying Theorem 1 to Corollary 3, we can generate  $\mathbf{x}$  with

- Sample  $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \tilde{\boldsymbol{\Sigma}})$ ;
- Return  $\mathbf{x} = \mathbf{y} + \tilde{\boldsymbol{\Sigma}}\mathbf{G}^T[\mathbf{G}\tilde{\boldsymbol{\Sigma}}\mathbf{G}^T]^{-1}(\mathbf{r} - \mathbf{G}\mathbf{y})$ .

Since  $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \mathbf{0} \end{bmatrix}$ ,  $\tilde{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \end{bmatrix}$ ,  $\mathbf{G} = (\mathbf{G}_1, \mathbf{G}_2) = (\boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}, \mathbf{I}_{k_2})$ , and  $\mathbf{r} = \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1$ , we have  $\tilde{\boldsymbol{\Sigma}}\mathbf{G}^T = \begin{bmatrix} \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \end{bmatrix}$  and  $[\mathbf{G}\tilde{\boldsymbol{\Sigma}}\mathbf{G}^T]^{-1} = \boldsymbol{\Sigma}_{22}^{-1}$ . Since  $\tilde{\boldsymbol{\Sigma}}$  is block diagonal, we can independently sample  $\mathbf{y}_1$  and  $\mathbf{y}_2$  as  $\mathbf{y}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$  and  $\mathbf{y}_2 \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})$ , respectively, with which we can further sample  $\mathbf{x}$  as

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12} \end{bmatrix} \boldsymbol{\Sigma}_{22}^{-1}(\boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\mathbf{y}_1 - \mathbf{y}_2).$$

Thus we can let  $\mathbf{x}_1 = \boldsymbol{\mu}_1 + \mathbf{y}'_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\mathbf{y}'_1 + \mathbf{y}_2)$ , where  $\mathbf{y}'_1 = \mathbf{y}_1 - \boldsymbol{\mu}_1 \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{11})$ .  $\square$

*Proof of Corollary 5.* Using the matrix inversion lemma, we have

$$\boldsymbol{\Sigma}_\beta = \mathbf{A}^{-1} - \mathbf{A}^{-1}\boldsymbol{\Phi}^T(\boldsymbol{\Omega}^{-1} + \boldsymbol{\Phi}\mathbf{A}^{-1}\boldsymbol{\Phi}^T)^{-1}\boldsymbol{\Phi}\mathbf{A}^{-1}. \quad (23)$$

The proof is completed by using Corollary 4 with  $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_\beta$ ,  $\boldsymbol{\Sigma}_{11} = \mathbf{A}^{-1}$ ,  $\boldsymbol{\Sigma}_{12} = \mathbf{A}^{-1}\boldsymbol{\Phi}^T$ , and  $\boldsymbol{\Sigma}_{22} = \boldsymbol{\Omega}^{-1} + \boldsymbol{\Phi}\mathbf{A}^{-1}\boldsymbol{\Phi}^T$ .  $\square$

## Computational Complexity

We present the computational complexities of all proposed algorithms in the following tables, where we highlight with bold the lowest complexity in each row.

Table 1: Computational complexity of Algorithm 1.

Calculation	Computational complexity	
	Non-diagonal $\boldsymbol{\Sigma}$	Diagonal $\boldsymbol{\Sigma}$
$\mathbf{H}$	$\mathcal{O}(k_2k^2)$	$\mathcal{O}(k_2k^2)$
$\mathbf{z}_2$	$\mathcal{O}(k_2^2k)$	$\mathcal{O}(k_2^2k)$
$\boldsymbol{\Sigma}^{-1}$	$\mathcal{O}(k^3)$	$\mathcal{O}(\mathbf{k})$
$\boldsymbol{\Lambda}_{11}$	$\mathcal{O}(k_1k^2)$	$\mathcal{O}(\mathbf{k}_1^2\mathbf{k})$
$\boldsymbol{\Lambda}_{12}$	$\mathcal{O}(k_1k_2k)$	$\mathcal{O}(k_1k_2k)$
$\boldsymbol{\mu}_{\mathbf{z}_1}$	$\mathcal{O}(\max(k^2, k_1^3, k_1^2k_2))$	$\mathcal{O}(\max(k^2, k_1^3, k_1^2k_2))$
$\mathbf{z}_1$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1^3)$
$\mathbf{x}$	$\mathcal{O}(\max(k_1k, k_2k))$	$\mathcal{O}(\max(k_1k, k_2k))$
Summary	$\mathcal{O}(k^3)$	$\mathcal{O}(\max(k_2k^2, \mathbf{k}_1^2\mathbf{k}))$

Table 2: Computational complexity of Algorithm 2.

Calculation	Computational complexity	
	Non-diagonal $\Sigma$	Diagonal $\Sigma$
$\mathbf{y}$	$\mathcal{O}(k^3)$	$\mathcal{O}(k)$
$\mathbf{G}\Sigma\mathbf{G}^T$	$\mathcal{O}(k_2k^2)$	$\mathcal{O}(k_2^2k)$
$\alpha$	$\mathcal{O}(\max(k_2k, k_2^3))$	$\mathcal{O}(\max(k_2k, k_2^3))$
$\mathbf{x}$	$\mathcal{O}(k_2k)$	$\mathcal{O}(k_2k)$
Summary	$\mathcal{O}(k^3)$	$\mathcal{O}(k_2^2k)$

Table 3: Computational complexity of naive simulation in Algorithm 3.

Calculation	Computational complexity			
	Non-diagonal $\Sigma_{11}$	Diagonal $\Sigma_{11}$	Non-diagonal $\Sigma_{11}$	Diagonal $\Sigma_{11}$
	Non-diagonal $\Sigma_{22}$	Non-diagonal $\Sigma_{22}$	Diagonal $\Sigma_{22}$	Diagonal $\Sigma_{22}$
$\Sigma_{22}^{-1}$	$\mathcal{O}(k_2^3)$	$\mathcal{O}(k_2^3)$	$\mathcal{O}(k_2)$	$\mathcal{O}(k_2)$
$\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$	$\mathcal{O}(\max(k_1^2k_2, k_1k_2^2))$	$\mathcal{O}(\max(k_1^2k_2, k_1k_2^2))$	$\mathcal{O}(k_1^2k_2)$	$\mathcal{O}(k_1^2k_2)$
$\mathbf{x}_1$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1^3)$
Summary	$\mathcal{O}(\max(k_1^3, k_2^3))$	$\mathcal{O}(\max(k_1^3, k_2^3))$	$\mathcal{O}(\max(k_1^3, k_1^2k_2))$	$\mathcal{O}(\max(k_1^3, k_1^2k_2))$

Table 4: Computational complexity of Algorithm 3.

Calculation	Computational complexity			
	Non-diagonal $\Sigma_{11}$	Diagonal $\Sigma_{11}$	Non-diagonal $\Sigma_{11}$	Diagonal $\Sigma_{11}$
	Non-diagonal $\Sigma_{22}$	Non-diagonal $\Sigma_{22}$	Diagonal $\Sigma_{22}$	Diagonal $\Sigma_{22}$
$\mathbf{y}_1$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1)$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1)$
$\Sigma_{11}^{-1}$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1)$	$\mathcal{O}(k_1^3)$	$\mathcal{O}(k_1)$
$\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$	$\mathcal{O}(\max(k_1k_2^2, k_1^2k_2))$	$\mathcal{O}(k_1k_2^2)$	$\mathcal{O}(\max(k_1k_2^2, k_1^2k_2))$	$\mathcal{O}(k_1k_2^2)$
$\mathbf{y}_2$	$\mathcal{O}(k_2^3)$	$\mathcal{O}(k_2^3)$	$\mathcal{O}(k_2^3)$	$\mathcal{O}(k_2^3)$
$\alpha$	$\mathcal{O}(\max(k_1k_2, k_2^3))$	$\mathcal{O}(\max(k_1k_2, k_2^3))$	$\mathcal{O}(k_1k_2)$	$\mathcal{O}(k_1k_2)$
$\mathbf{x}_1$	$\mathcal{O}(k_1k_2)$	$\mathcal{O}(k_1k_2)$	$\mathcal{O}(k_1k_2)$	$\mathcal{O}(k_1k_2)$
Summary	$\mathcal{O}(\max(k_1^3, k_2^3))$	$\mathcal{O}(\max(k_1k_2^2, k_2^3))$	$\mathcal{O}(\max(k_1^3, k_2^3))$	$\mathcal{O}(\max(k_1k_2^2, k_2^3))$

Table 5: Computational complexity of naive simulation in Algorithm 4.

Calculation	Computational complexity			
	Non-diagonal $\mathbf{A}$	Diagonal $\mathbf{A}$	Non-diagonal $\mathbf{A}$	Diagonal $\mathbf{A}$
	Non-diagonal $\Omega$	Non-diagonal $\Omega$	Diagonal $\Omega$	Diagonal $\Omega$
$\Phi^T\Omega\Phi$	$\mathcal{O}(\max(n^2p, np^2))$	$\mathcal{O}(\max(n^2p, np^2))$	$\mathcal{O}(np^2)$	$\mathcal{O}(np^2)$
$(\mathbf{A} + \Phi^T\Omega\Phi)^{-1}$	$\mathcal{O}(p^3)$	$\mathcal{O}(p^3)$	$\mathcal{O}(p^3)$	$\mathcal{O}(p^3)$
$\beta$	$\mathcal{O}(p^3)$	$\mathcal{O}(p^3)$	$\mathcal{O}(p^3)$	$\mathcal{O}(p^3)$
Summary	$\mathcal{O}(\max(n^2p, p^3))$	$\mathcal{O}(\max(n^2p, p^3))$	$\mathcal{O}(\max(np^2, p^3))$	$\mathcal{O}(\max(np^2, p^3))$

Table 6: Computational complexity of Algorithm 4.

Calculation	Computational complexity			
	Non-diagonal $\mathbf{A}$	Diagonal $\mathbf{A}$	Non-diagonal $\mathbf{A}$	Diagonal $\mathbf{A}$
	Non-diagonal $\mathbf{\Omega}$	Non-diagonal $\mathbf{\Omega}$	Diagonal $\mathbf{\Omega}$	Diagonal $\mathbf{\Omega}$
$\mathbf{A}^{-1}$	$\mathcal{O}(p^3)$	$\mathcal{O}(p)$	$\mathcal{O}(p^3)$	$\mathcal{O}(p)$
$\mathbf{y}_1$	$\mathcal{O}(p^3)$	$\mathcal{O}(p)$	$\mathcal{O}(p^3)$	$\mathcal{O}(p)$
$\mathbf{\Omega}^{-1}$	$\mathcal{O}(n^3)$	$\mathcal{O}(n^3)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
$\mathbf{y}_2$	$\mathcal{O}(n^3)$	$\mathcal{O}(n^3)$	$\mathcal{O}(n)$	$\mathcal{O}(n)$
$\mathbf{\Omega}^{-1} + \mathbf{\Phi}\mathbf{A}^{-1}\mathbf{\Phi}^T$	$\mathcal{O}(\max(np^2, n^2p))$	$\mathcal{O}(n^2p)$	$\mathcal{O}(\max(np^2, n^2p))$	$\mathcal{O}(n^2p)$
$\alpha$	$\mathcal{O}(\max(np, n^3))$	$\mathcal{O}(\max(np, n^3))$	$\mathcal{O}(\max(np, n^3))$	$\mathcal{O}(\max(np, n^3))$
$\beta$	$\mathcal{O}(np)$	$\mathcal{O}(np)$	$\mathcal{O}(np)$	$\mathcal{O}(np)$
Summary	$\mathcal{O}(\max(n^3, p^3))$	$\mathcal{O}(\max(n^2p, n^3))$	$\mathcal{O}(\max(n^3, p^3))$	$\mathcal{O}(\max(n^2p, n^3))$

### Brief derivation of SG-MCMC for a simplex-constrained vector

Based on a comprehensive framework for constructing SG-MCMC algorithms in Ma et al. (2015), we have the mini-batch update rule for a global variable  $\mathbf{z}$  as

$$\begin{aligned} \mathbf{z}_{t+1} = & \mathbf{z}_t + \varepsilon_t \left\{ -[\mathbf{D}(\mathbf{z}_t) + \mathbf{Q}(\mathbf{z}_t)] \nabla \tilde{H}(\mathbf{z}_t) + \Gamma(\mathbf{z}_t) \right\} \\ & + \mathcal{N} \left( \mathbf{0}, \varepsilon_t \left[ 2\mathbf{D}(\mathbf{z}_t) - \varepsilon_t \hat{\mathbf{B}}_t \right] \right), \end{aligned} \quad (24)$$

where  $\varepsilon_t$  are annealed step sizes,  $\mathbf{D}(\mathbf{z})$  is a positive semidefinite diffusion matrix,  $\mathbf{Q}(\mathbf{z})$  is a skew-symmetric curl matrix,  $\hat{\mathbf{B}}_t$  is an estimate of the stochastic gradient noise variance satisfying a positive definite constraint as  $2\mathbf{D}(\mathbf{z}_t) - \varepsilon_t \hat{\mathbf{B}}_t \succ \mathbf{0}$ , and  $\Gamma_i(\mathbf{z})$ , the  $i^{\text{th}}$  element of the compensation vector  $\Gamma(\mathbf{z})$ , is defined as  $\Gamma_i(\mathbf{z}) = \sum_j \frac{\partial}{\partial z_j} [\mathbf{D}_{ij}(\mathbf{z}) + \mathbf{Q}_{ij}(\mathbf{z})]$ . The mini-batch estimation of energy function is defined as  $\tilde{H}(\mathbf{z}) = -\ln p(\mathbf{z}) - \rho \sum_{\mathbf{x} \in \tilde{X}} \ln p(\mathbf{x} | \mathbf{z})$ , with  $\tilde{X}$  the mini-batch and  $\rho$  the ratio of the dataset size to the mini-batch size.

For simplicity, we adopt the same specifications that lead to the stochastic gradient Riemannian Langevin dynamics (SGRLD) inference algorithm for simplex-constrained model parameters (Patterson and Teh 2013; Ma et al. 2015), namely  $\mathbf{D}(\mathbf{z}) = \mathbf{G}(\mathbf{z})^{-1}$ ,  $\mathbf{Q}(\mathbf{z}) = \mathbf{0}$ , and  $\hat{\mathbf{B}}_t = \mathbf{0}$ , where  $\mathbf{G}(\mathbf{z})$  denotes the Fisher information matrix (FIM) (Giro-lami and Calderhead 2011).

With the multinomial likelihood  $\mathbf{n}_j \sim \text{Mult}(n_j, \phi)$ , and the reduced-mean parameterization  $\varphi \in \mathbb{R}_+^{V-1}$ , where  $j \in \{1, \dots, N\}$  with  $N$  the dataset size, it is straight to derive the FIM as

$$\begin{aligned} \mathbf{G}(\varphi) &= -\mathbb{E} \left[ \frac{\partial^2}{\partial \varphi^2} \ln \left( \prod_j \text{Mult}[\mathbf{n}_j; n_j, \phi] \right) \right] \\ &= M \left[ \text{diag}(1/\varphi) + \mathbf{1}\mathbf{1}^T / (1 - \varphi) \right], \end{aligned} \quad (25)$$

where  $M := \mathbb{E} \left[ \sum_{j=1}^N n_j \right]$  is approximated along the updating as  $M = (1 - \varepsilon_t) M + \varepsilon_t \rho E[n..]$ . Further with the Dirichlet prior  $\phi \sim \text{Dir}(\eta \mathbf{1}_V)$ , we have the conditional

posterior of  $\phi$  as  $(\phi|-) \sim \text{Dir}(\sum_j n_{1j} + \eta, \dots, \sum_j n_{Vj} + \eta)$ . Taking the gradient with respect to the reduced-mean parameterization  $\varphi \in \mathbb{R}_+^{V-1}$  on the mini-batch estimation of the negative energy function, we have

$$\nabla_{\varphi} \left[ -\tilde{H}(\varphi) \right] = \frac{\rho \bar{\mathbf{n}}_{\cdot\cdot} + \eta - 1}{\varphi} - \frac{\rho n_{V\cdot} + \eta - 1}{1 - \varphi}. \quad (26)$$

Substituting both (25) and (26) into (24), we have (16) as

$$\varphi_{t+1} = \left[ \varphi_t + \frac{\varepsilon_t}{M} [(\rho \bar{\mathbf{n}}_{\cdot\cdot} + \eta) - (\rho n_{\cdot\cdot} + \eta V) \varphi_t] + \mathcal{N} \left( \mathbf{0}, \frac{2\varepsilon_t}{M} [\text{diag}(\varphi_t) - \varphi_t \varphi_t^T] \right) \right]_{\Delta},$$

where  $[\cdot]_{\Delta}$ , denoting the constraint that  $\varphi \in \mathbb{R}_+^{V-1}$  and  $\varphi_{\cdot} \leq 1$ , ensures  $\varphi$  to be valid.

Next we prove that equation (16) can be equivalently represented as (17), namely

$$\phi_{t+1} = \left[ \phi_t + \frac{\varepsilon_t}{M} [(\rho \mathbf{n}_{\cdot\cdot} + \eta) - (\rho n_{\cdot\cdot} + \eta V) \phi_t] + \mathcal{N} \left( \mathbf{0}, \frac{2\varepsilon_t}{M} \text{diag}(\phi_t) \right) \right]_{\angle},$$

where  $[\cdot]_{\angle}$  represents the constraint that  $\phi \in \mathbb{R}_+^V$  and  $\mathbf{1}^T \phi = 1$ . By substituting  $\phi = (\varphi^T, 1 - \mathbf{1}^T \varphi)^T$  into (17), one can easily verify that the MVN simulation in (17) is identical to that in (16). By further pointing out the fact that  $[\cdot]_{\Delta}$  is the same as  $[\cdot]_{\angle}$  under the reduced-mean parameterization, we conclude the proof.