

# Extrinsic Priors for Bayesian Modeling with Parameter Constraints

Leo Duan, Akihiko Nishimura, David Dunson

**Abstract:** Prior information often takes for the form of parameter constraints. Bayesian methods include such information through prior distributions having constrained support. By using posterior sampling algorithms, one can quantify uncertainty without relying on asymptotic approximations. However, outside of narrow settings, parameter constraints make it difficult to develop efficient posterior sampling algorithms. We propose a general solution, which relaxes the constraint through the use of an *extrinsic prior*, which is concentrated close to the constrained space. General off the shelf posterior sampling algorithms, such as Hamiltonian Monte Carlo (HMC), can then be used directly. We illustrate this approach through multiple examples involving equality and inequality constraints. While existing methods tend to rely on conjugate families, our proposed approach frees us up to define new classes of hierarchical models for constrained problems. We illustrate this through application to a variety of simulated and real datasets.

KEY WORDS: Constraint relaxation; Euclidean Embedding; Monotone Dirichlet; Soft Constraint; Stiefel Manifold; Projected Markov chain

## 1 Introduction

It is extremely common to have prior information available on parameter constraints in statistical models. For example, one may have prior knowledge that a vector of parameters lies on the probability simplex or satisfies a particular set of inequality constraints. Other common examples include shape constraints on functions, positive semidefiniteness of matrices and orthogonality. There is a very rich literature on optimization subject to parameter constraints. One common approach is to rely on Lagrange and Karush-Kuhn-Tucker multipliers (Boyd and Vandenberghe, 2004). However, simply producing a point estimate is often insufficient, as uncertainty quantification (UQ) is a key component of most statistical analyses. Usual large sample asymptotic theory, for example showing asymptotic normality of statistical estimators, tends to break down in constrained inference problems. Instead, limiting distributions may have a complex form that needs to be rederived for each new type of constraint, and may be intractable. An appealing alternative is to rely on Bayesian methods for UQ, including the constraint through a prior distribution having restricted support, and then applying Markov chain Monte Carlo (MCMC) to avoid the need for large sample approximations.

Conceptually MCMC can be applied in a broad class of constrained parameter problems without complications Gelfand et al. (1992). **Aki: I think sampling from a constrained space, say a manifold, can be challenging even without the computational efficiency issues?** However, in practice, a primary difficulty is designing a Markov transition kernel that leads to an MCMC algorithm with sufficient computational efficiency to be practically useful. Common default transition kernels correspond to Gibbs sampling, random walk Metropolis-Hastings, and (more recently) Hamiltonian Monte Carlo (HMC). Gibbs sampling relies on alternately sampling from the full conditional posterior distributions for the different parameters, ideally in blocks to improve mixing. Gibbs requires the conditional distributions to be available in a form that is tractable to sample from directly, limiting consideration to specialized models. In constrained problems, block updating is typically either not possible or very inefficient (e.g. relying on rejection sampling with a high rejection probability), and one-at-a-time updating can lead to extremely slow mixing. Random walk algorithms provide an alternative, but each step of the random walk must maintain the parameter constraint. A common approach is to apply a normal random walk and simply reject proposals that violate the constraint, but this can have very high rejection rates even if using an adaptive approach that learns the covariance based on the history of the chain. An alternative is to rely on HMC. In simple settings in which a reparameterization can be applied to remove the constraint, HMC can be applied easily. Otherwise, HMC will generate proposals that violate the constraint, and hence face problems with high rejection rates in heavily constrained problems.

Due to the above hurdles, most of the focus in the literature has been on customized solutions developed for specific constraints. One popular strategy is to carefully pick a prior and likelihood such that posterior sampling is tractable. For example, for modeling of data on manifolds, it is typical to restrict attention to specific models, such as the Bingham-von Mises-Fisher distribution for Stiefel manifolds (Khatri and Mardia, 1977; Hoff, 2009). For data on the probability simplex, one instead relies on the Dirichlet distribution. An alternative is to reparameterize the model to eliminate or simplify the constraint. For example, when faced with a monotonicity constraint, one may reparameterize in terms of differences as the resulting positivity constraint leads to much easier sampling (REFs). In the literature on modeling of data on manifolds, there are two strategies: (i) *intrinsic* methods that define a statistical model directly on the manifold, and (ii) *extrinsic* methods that indirectly induce a model on the manifold through embedding the manifold in a Euclidean space, defining a model in the Euclidean space, and then projecting back onto the manifold. Essentially all of the current strategies for Bayesian modeling with constraints take an intrinsic-style approach. However, by strictly maintaining the constraint at all stages of the modeling and computation process, one limits the possibilities in terms of defining general methods to deal with parameter constraints.

These drawbacks motivate the development of *extrinsic* approaches that define an unconstrained model and/or computational algorithm, and then somehow adjust for the constraint. A related idea is Gelfand

et al. (1992), who suggested running Gibbs sampling ignoring the constraint but only accepting the draws satisfying the constraint. Unfortunately, such an approach is highly inefficient, as motivated above. An alternative is to run MCMC ignoring the constraint, and then project draws from the unconstrained posterior to the appropriately constrained space. Such an approach was proposed for generalized linear models with order constraints by Dunson and Neelon (2003), extended to functional data with monotone or unimodal constraints Gunn and Dunson (2005), and recently modified to nonparametric regression with monotonicity Lin and Dunson (2014) or manifold Lin et al. (2016) constraints.

An alternative idea is to *relax* a sharp parameter constraint by defining a prior that has unrestricted support but places small probability outside of the constrained region. Neal (2011) suggested such an approach to apply HMC in settings involving a simple truncation constraint, while Pakman and Paninski (2014) applied a related idea to improve sampling from truncated multivariate normal distributions.

The goal of this article is to dramatically generalize these specific approaches to develop a broad class of *extrinsic priors* for parameter constrained problems. These priors are defined to place small probability outside of the constrained region, while permitting use of efficient and general use MCMC algorithms; in particular, HMC. When the constraints need to be upheld strictly, the approximation can be corrected with a simple projection, followed by a Metropolis-Hastings step with high acceptance probability. Unlike intrinsic methods, such as Riemannian and geodesic HMC (Girolami and Calderhead, 2011; Byrne and Girolami, 2013), our approach is relatively efficient and simple to implement in general settings using automatic algorithms. **Aki: I am not sure the "relatively efficient" claim is warranted without some empirical comparisons. Perhaps better just to say "simpler and more generally applicable"?** The generality frees up a much broader spectrum of Bayesian models, as one no longer needs to focus on very specific computationally tractable models. Theoretic studies are conducted and original models are shown in simulations and data applications.

## 2 Extrinsic Bayes Methodology

### 2.1 Conditioned Intrinsic Distribution

Let  $\theta \in \mathcal{D}$  denote the parameters in likelihood function  $L(\theta; y)$ , with  $y$  the data. The support  $\mathcal{D}$  is a constrained space. The usual Bayesian approach assigns a prior density  $\pi_{0,\mathcal{D}}(\theta)$  for  $\theta$  having support  $\mathcal{D}$ . Traditional strategy is to reparameterize  $\theta$  with  $f(\theta)$  to always satisfy  $\theta \in \mathcal{D}$ , and find an *intrinsic* distribution on  $\mathcal{D}$ . Diaconis et al. (2013) provide a detailed review on these methods. Although some of them are successful, in general, reparameterization does not always exist and achieving desired property (such as uniformity on  $\mathcal{D}$ ) can be difficult due to variable transformation.

We now present a different intrinsic strategy that does not involve transformation. Assuming that  $\mathcal{D} \subset \mathcal{R}$ , with  $\mathcal{R}$  denoting a ‘less constrained’ space, we start with an unconstrained prior  $\pi_{0,\mathcal{R}}(\theta)$ , conditioning it on a

random variable  $\omega = \mathbb{1}_{\theta \in \mathcal{D}}$  to obtain conditional density  $\pi_{0,\mathcal{D}}(\theta) = \pi_{0,\mathcal{R}}(\theta \mid \omega = 1)$ . We call this ‘conditioned intrinsic distribution’. In order to utilize the various existing distribution families and geometric measure theory, we focus on  $\mathcal{R}$  being Euclidean space  $\mathbb{R}^p$  or its truncated subspace.

Let  $(\Omega, \mathcal{B}, P)$  be a probability space and sub- $\sigma$ -field  $\mathcal{A} \subseteq \mathcal{B}$ , the function  $P(\cdot \mid \mathcal{A})$  is a *regular conditional probability* (r.c.p.) (Kolmogorov, 1950) on  $\mathcal{B}$  given  $\mathcal{A}$ , if:

1. For each  $\omega \in \Omega$ ,  $P(\cdot \mid \mathcal{A})(\omega)$  is a probability measure on  $\mathcal{B}$ .
2. For each  $B \in \mathcal{B}$ , the mapping  $P(B \mid \mathcal{A})(\cdot)$  is  $\mathcal{A}$ -measurable.
3. For each  $A \in \mathcal{A}$ ,  $B \in \mathcal{B}$ ,  $P(A \cap B) = \int_A P(B \mid \mathcal{A})(\omega) P(d\omega)$ .

In our case, recall  $\omega = 1$  if  $\theta \in \mathcal{D}$  and equal 0 otherwise. Let  $Z(\mathcal{D}) = \int_{\mathcal{D}} \pi_{0,\mathcal{R}}(\theta) d\theta \in [0, 1]$ , we first assign:

$$P(\omega) = wZ(\mathcal{D}) + (1-w)(1-Z(\mathcal{D})), \quad P(B \mid \mathcal{A})(\omega) = P(B \cap \mathcal{D} \mid \omega = 1)\omega + P(B \cap (\mathcal{R} \setminus \mathcal{D}) \mid \omega = 0)(1-\omega), \quad (1)$$

where  $P(B \cap (\mathcal{R} \setminus \mathcal{D}) \mid \omega = 0)$  is a probability measure that we are not interested in. When  $Z(\mathcal{D}) > 0$ , one can simply define:

$$\int_B \pi_{0,\mathcal{D}}(\theta) d\theta = P(B \cap \mathcal{D} \mid \omega = 1) = \int_B \frac{\pi_{0,\mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}}}{Z(\mathcal{D})} d\theta. \quad (2)$$

When  $Z(\mathcal{D}) = 0$ , in (1), any proper  $P(B \cap \mathcal{D} \mid \omega = 1)$  would satisfy the r.c.p criteria. However, one can still utilize  $\pi_{0,\mathcal{R}}(\theta)$  by considering measure at a lower dimension  $m < n$ . As an intuitive example, a surface can have Lebesgue volume  $\lambda^3(\mathcal{D}) = 0$ , but positive area on a 2-dimensional manifold.

Formally, for a set  $\mathcal{D}$ , if it is the Lipschitz image of a bounded subset in  $\mathbb{R}^m$  (known as ‘ $m$ -rectifiable’), it has Minkowski content:

$$\mathcal{M}^m(\mathcal{D}) = \lim_{r \rightarrow 0^+} \frac{\lambda^n(D_r)}{\alpha_{n-m} r^{n-m}} \quad (3)$$

where  $\alpha_t = \frac{\Gamma(\frac{1}{2})^t}{\Gamma(1+\frac{t}{2})}$ , with  $D_r = \{x \in \mathbb{R}^n : \text{dist}(x, y) < r, y \in \mathcal{D}\}$  known as the open  $r$ -parallel set,  $\text{dist}(\cdot, \cdot)$  is a distance;  $\lambda^n(D_r) = \int_{D_r} \pi_{0,\mathcal{R}}(\theta) d\theta$  is its Lebesgue measure in  $\mathbb{R}^n$ .

When  $m$  is chosen such that  $0 < \mathcal{M}^m(\mathcal{D}) < \infty$ , then  $\mathcal{D}$  is called  $m$ -dimensional Minkowski measurable, and  $\mathcal{M}^m(\mathcal{D})$  is proportional to the Hausdorff measure  $\mathcal{H}^m(\mathcal{D})$  for  $m$ -rectifiable  $\mathcal{D}$ . Fortunately, we do not need to find  $m$  in this case. Letting  $P(B \cap \mathcal{D} \mid \omega = 1) \propto \mathcal{M}^m(B \cap \mathcal{D})$  and  $P(\mathcal{D} \mid \omega = 1) = 1$  yields:

$$\int_B \pi_{0,\mathcal{D}}(\theta) d\theta = P(B \cap \mathcal{D} \mid \omega = 1) = \frac{\mathcal{M}^m(B \cap \mathcal{D})}{\mathcal{M}^m(\mathcal{D})} = \lim_{r \rightarrow 0^+} \frac{\lambda^n(B \cap D_r)}{\lambda^n(D_r)} = \lim_{r \rightarrow 0^+} \frac{\int_B \pi_{0,\mathcal{R}}(\theta) \mathbb{1}_{\theta \in D_r} d\theta}{\int_{D_r} \pi_{0,\mathcal{R}}(\theta)} \quad (4)$$

where the terms involving  $m$  are canceled, due to  $\mathcal{M}^m(B \cap \mathcal{D}) = \lim_{r \rightarrow 0^+} \frac{\lambda^n(B \cap \mathcal{D}_r)}{\alpha_{n-m} r^{n-m}}$  (Winter, 2016). Clearly, this result is extensible to unbounded set  $\mathcal{D}$  provided  $0 < \mathcal{M}^m(\mathcal{D}) < \infty$  for certain  $m$ , although the equivalence of Minkowski content and Hausdorff measure may not hold.

Letting the likelihood function be  $L(y; \theta)$ , the posterior of  $\theta$  is  $\pi_{\mathcal{D}}(\theta | y) = \frac{L(y; \theta) \pi_{0, \mathcal{D}}(\theta)}{\int_{\mathcal{D}} L(y; \theta) \pi_{0, \mathcal{D}}(\theta) d\theta}$ . If  $\int_{\mathcal{D}} \pi_{0, \mathcal{R}}(\theta) d\theta > 0$  and  $\int_{\mathcal{R}} L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}} d\theta > 0$ ,

$$\pi_{\mathcal{D}}(\theta | y) = \frac{L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}}}{\int_{\mathcal{R}} L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}} d\theta}; \quad (5)$$

otherwise,

$$\int_B \pi_{\mathcal{D}}(\theta | y) d\theta = \lim_{r \rightarrow 0^+} \frac{\int_B L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}_r} d\theta}{\int_{\mathcal{R}} L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}_r} d\theta}. \quad (6)$$

Therefore, both prior and posterior can be obtained intrinsically via this conditioning approach, although the posterior sampling is generally difficult and (4)(6) may be intractable.

We now use an example to illustrate this strategy approach when  $Z(\mathcal{D}) = 0$ .

**Example 1A: Two Gaussians with Sum Constraint (Conditioned Intrinsic)**

Consider a bivariate Gaussian random vector  $[\theta_1, \theta_2]' \sim \text{No}(0, I)$  in constrained space  $\mathcal{D} = \{(\theta_1, \theta_2) : \theta_1 + \theta_2 = 1\}$ . Denoting  $\phi(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$ ,  $Z(\mathcal{D}) = \int_{\mathcal{D}} \phi(\theta_1) \phi(\theta_2) d\theta = 0$ . Using  $\mathcal{D}_r = \{(\theta_1, \theta_2) : \theta_1 + \theta_2 \in (1-r, 1+r)\}$  and  $B = \{(\theta_1, \theta_2) : \theta_1 < x, \theta_2 \in \mathbb{R}\}$ ,

$$\begin{aligned} \int_B \pi_{0, \mathcal{D}}(\theta) d\theta &= \lim_{r \rightarrow 0^+} \frac{\int_{(1-r)}^{(1+r)} \int_{-\infty}^x \phi(\theta_1) \phi(z - \theta_1) d\theta_1 dz}{\int_{(1-r)/\sqrt{2}}^{(1+r)/\sqrt{2}} \phi(z) dz} \\ &= \lim_{r \rightarrow 0^+} \frac{\frac{\partial}{\partial r} \int_{(1-r)}^{(1+r)} \int_{-\infty}^x \phi(\theta_1) \phi(z - \theta_1) d\theta_1 dz}{\frac{\partial}{\partial r} \int_{(1-r)/\sqrt{2}}^{(1+r)/\sqrt{2}} \phi(z) dz} \\ &= \lim_{r \rightarrow 0^+} \frac{\int_{-\infty}^x \phi(\theta_1) (\phi(1+r-\theta_1) d\theta_1 + \phi(1-r-\theta_1)) d\theta_1}{\phi(\frac{1+r}{\sqrt{2}})/\sqrt{2} + \phi(\frac{1-r}{\sqrt{2}})/\sqrt{2}} \\ &= \int_{-\infty}^x \frac{\sqrt{2}}{\sqrt{2\pi}} \exp(-\frac{(\theta_1 - \frac{1}{2})^2}{2/2}) d\theta_1, \end{aligned}$$

which corresponds to  $\theta_1 | (\theta_1 + \theta_2 = 1) \sim \text{No}(1/2, 1/2)$ ,  $\theta_2 | \theta_1 \sim \delta_{1-\theta_1}$  with  $\delta_x$  denoting a point mass at  $x$ .

Marginally, this is a degenerate bivariate Gaussian distribution:

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \sim \text{No} \left( \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}, \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix} \right).$$

## 2.2 Extrinsic Bayes

The conditioned intrinsic distribution relies on the sharp indicator function  $\mathbb{1}_{\theta \in \mathcal{D}}$  and limit of open r-parallel set  $\mathcal{D}_r$ , leading to computational difficulty in prior and posterior. We now propose an extrinsic prior that builds on (2) and (4), approximating  $\mathbb{1}_{\theta \in \mathcal{D}}$  with a *smooth* alternative  $\mathcal{K}(\theta; \mathcal{D})$  with less constrained support:

$$\tilde{\pi}_{0, \mathcal{D}}(\theta) = \frac{\pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D})}{\int_{\mathcal{R}} \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}) d\theta} \quad (7)$$

where  $\mathcal{K}(\theta; \mathcal{D})$  satisfies  $\mathcal{K}(\theta; \mathcal{D}) = 1$  if  $\theta \in \mathcal{D}$  and  $\int_{\mathcal{R}} \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}) d\theta > 0$ . When  $Z(\mathcal{D}) = 0$ ,  $\mathcal{K}(\theta; \mathcal{D}) = \lim_{r \rightarrow 0^+} \mathcal{K}(\theta; \mathcal{D}_r)$  with  $\mathcal{K}(\theta; \mathcal{D}_r)$  approximating  $\mathbb{1}_{\theta \in \mathcal{D}_r}$ , then  $\int_B \tilde{\pi}_{0, \mathcal{D}}(\theta) d\theta = \lim_{r \rightarrow 0^+} \frac{\int_B \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta}{\int_{\mathcal{R}} \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta} = \int_B \frac{\pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D})}{\int_{\mathcal{R}} \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}) d\theta} d\theta$ . Therefore, the limit forms in (4) and (6) are no longer needed. The posterior takes similar form

$$\tilde{\pi}_{\mathcal{D}}(\theta | y) = \frac{L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D})}{\int_{\mathcal{R}} L(y; \theta) \pi_{0, \mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}) d\theta}. \quad (8)$$

We now define  $\mathcal{K}(\theta; \mathcal{D})$ . Assuming there are  $m$  constraints with each defining a constrained subspace  $\mathcal{D}_k$ ,  $\mathcal{D} = \bigcap_{k=1}^m \mathcal{D}_k$ , we have

$$\mathcal{K}(\theta; \mathcal{D}) = \prod_{k=1}^m K_k(v_k(\theta)), \quad (9)$$

where  $v_k : \mathcal{R} \rightarrow [0, \infty)$  is a measurable function and quantifies the distance to space  $\mathcal{D}_k$ . For example, one can use  $v(\theta) = |f(\theta)|$  as a distance to equality-constrained space  $\{\theta : f(\theta) = 0\}$ ;  $v(\theta) = |f(\theta)|_+$ , where  $(x)_+ =$

$\begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$  as a distance to inequality-constrained space  $\{\theta : f(\theta) \leq 0\}$ . The function  $K_k : [0, \infty) \rightarrow [0, 1]$

penalizes  $v_k(\theta) > 0$  and decreases in  $v_k(\theta)$ . For example,  $K_k(v) = \exp(-v)$  or  $K_k(v) = \frac{1}{1+v^2}$ . For illustration

simplicity, in this paper, we focus on exponential smoothing function  $K_k(v(\theta)) = \exp(-\frac{v(\theta)}{\lambda_k})$  with  $\lambda_k > 0$  as

a tuning parameter. **Aki:** We need to assume the growth rate on  $v(\theta)$  to ensure that an extrinsic prior is integrable and proper. There can be a problem if  $v(\theta) = \log \|\theta\|$  for example. In fact, the geometric ergodicity of HMC would require  $-\log K(v(\theta)) = O(\|\theta\|^\alpha)$  for  $1 \leq \alpha \leq 2$ . **Leo:** can you provide the justification or reference for this result?

This framework is applicable to a variety of complicated case. For example,  $\theta$  can have only some of parameters constrained; parameters can be in multiple constraints simultaneously; constraints can be dependent. Regardless, it is generally simple to adapt (9) via appropriate mapping through  $v(\theta)$ .

We now provide some additional regularity conditions on  $v_k$  and  $K_k$ : for  $k = 1, \dots, m$ ,  $v_k(\theta) = 0$  only if  $\theta \in \mathcal{D}_k$  and  $K_k(v) = 1$  only if  $v = 0$ . Therefore,  $\theta \in \mathcal{D} \Leftrightarrow \text{all } v_k(\theta) = 0 \Leftrightarrow \mathcal{K}(\theta; \mathcal{D}) = 1$ , this ensures  $\mathcal{K}(\theta; \mathcal{D})$  is the same as the indicator function for  $\theta \in \mathcal{D}$ .

To illustrate, consider a truncated normal prior  $\text{No}_{(-\infty, 5)}(0, 5^2)$ . The unnormalized density of intrinsic prior is  $\exp(-\theta^2/2 \cdot 5^2) \mathbb{1}_{\theta \in (-\infty, 5)}$ , which can be approximated by extrinsic prior  $\exp(-\theta^2/2 \cdot 5^2) \exp(-v(\theta))$ . We plot two distances  $v(\theta) = (\theta - 5)_+$  and  $v(\theta) = (\theta - 5)_+^2$  in Figure 1. Inside  $\mathcal{D} = (-\infty, 5)$ , both intrinsic and extrinsic distribution are the same, except for a different normalizing constant; outside  $\mathcal{D}$ , intrinsic one drops directly to 0 at the boundary, whereas the extrinsic one decreases smoothly.

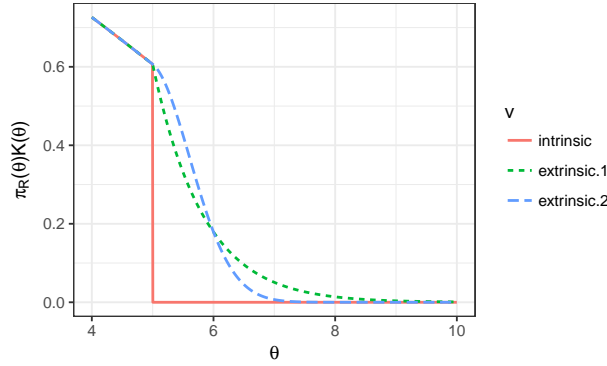


Figure 1: Unnormalized densities for truncated normal  $\text{No}_{(-\infty, 5)}(0, 5^2)$  under exact intrinsic prior and approximating extrinsic prior. Inside  $(-\infty, 5)$ , the priors are the same up to a constant difference. The intrinsic prior abruptly drops to 0 on the boundary, while the approximating ones drop smoothly. Intrinsic prior based on first-order  $v(\theta)$  drops faster than the one based on second order when  $v(\theta) \in (0, 1)$ .

We now apply extrinsic Bayes in the previous example of two Gaussians under sum constraint.

#### Example 1B: Two Gaussians with Sum Constraint (Extrinsic Approach)

We now use extrinsic prior  $\tilde{\pi}_{0, \mathcal{D}}(\theta) \propto \exp(-\frac{\theta_1^2 + \theta_2^2}{2}) \exp(-\frac{v(\theta)}{\lambda})$ . Choosing  $v(\theta) = (\theta_1 + \theta_2 - 1)^2$  allows us to obtain closed-form for the extrinsic prior  $\theta_1 \sim \text{No}(\frac{2}{\lambda+4}, \frac{\lambda+2}{\lambda+4})$ ,  $\theta_2 \mid \theta_1 \sim \text{No}(\frac{2}{\lambda+2}(1 - \theta_1), \frac{\lambda}{\lambda+2})$ . Marginally,

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \sim \text{No} \left( \begin{bmatrix} \frac{2}{\lambda+4} \\ \frac{2}{\lambda+4} \end{bmatrix}, \begin{bmatrix} \frac{\lambda+2}{\lambda+4} & -\frac{2}{\lambda+4} \\ -\frac{2}{\lambda+4} & \frac{\lambda+2}{\lambda+4} \end{bmatrix} \right).$$

As  $\lambda \rightarrow 0$ , the extrinsic prior becomes the same as the degenerate bivariate Gaussian in intrinsic approach.

Obviously, the approximation based on extrinsic approach is much simpler to derive compared to the exact conditioned intrinsic distribution.

## 2.3 Approximation Error

We now study the properties of the extrinsic prior and posterior. One important task is to quantify the difference between extrinsic and intrinsic ones. Due to similar reasoning for prior and posterior, we use some

general notation. Let  $\pi_{\mathcal{R}}(\theta)$  be an unnormalized density in  $\mathcal{R}$ , which is  $\pi_{0,\mathcal{R}}(\theta)$  when studying prior and  $L(y; \theta)\pi_{0,\mathcal{R}}(\theta)$  when studying posterior;  $\Pi(\cdot)$  and  $\tilde{\Pi}(\cdot)$  to represent the measures under intrinsic and extrinsic methods associated with either prior or posterior distribution.

We first explore the case in (2) and (5), where  $\int_{\mathcal{D}} \pi_{\mathcal{R}}(\theta) d\theta > 0$ .

**Remark 1.** Let  $M_1 = \int_{\mathcal{R}} \pi_{0,\mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}} d\theta$  and  $M_2 = \int_{\mathcal{R}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}) d\theta$ , when  $M_1 > 0$ , the total variation distance between the measures of extrinsic and intrinsic distributions is

$$\|\Pi(\cdot), \tilde{\Pi}(\cdot)\|_{TV} = 1 - \frac{M_1}{M_2} \leq \frac{\int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}) d\theta}{M_1}.$$

proof: Via definition of total variation distance and  $\mathcal{K}(\theta; \mathcal{D}) = 1$  when  $\theta \in \mathcal{D}$ .

Using exponential smoothing function  $\mathcal{K}(\theta; \mathcal{D}) = \prod_{k=1}^m \exp(-v_k(\theta)/\lambda_k)$ , as all  $\lambda_k \rightarrow 0$ , the total variation  $\|\Pi(\cdot), \tilde{\Pi}(\cdot)\|_{TV} \rightarrow 0$ , if  $\int_{\mathcal{R}} \pi_{\mathcal{R}}(\theta) d\theta < \infty$ . This is a direct result of dominated convergence theorem and  $\mathcal{K}(\theta; \mathcal{D}) \leq 1$ .

With some mild conditions, we further quantify the non-asymptotic rate. Letting  $\lambda = \sup_k \lambda_k$ ,  $M = \int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) d\theta < \infty$ . As a linear combination of measurable functions,  $v(\theta) = \lambda \sum_{k=1}^m \frac{v_k(\theta)}{\lambda_k}$  is measurable. Let  $f(v)$  be the density of  $v(\theta)$ . **Aki:** The density of  $v(\theta)$  crucially depends on the level set of  $v(\theta)$  and is a pretty complicated (i.e. intractable) object in general.

**Remark 2.** If there exists a  $t < \infty$  such that  $f(v) < \infty$ , for any  $t > 0$ ,

$$\int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) \prod_{k=1}^m \exp(-v_k(\theta)/\lambda_k) d\theta \leq M \exp(-\frac{t}{\lambda}) + M \sup_{t^* \in (0, t)} f(t^*) \lambda$$

**Aki:** How is this bound based on  $v(\theta) < t$  useful? We need  $v(\theta) \rightarrow \infty$  as  $\|\theta\| \rightarrow \infty$  for the extrinsic prior to be integrable. **Leo:** I rephrased the statement



proof:

$$\begin{aligned}
\frac{1}{M} \int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) \prod_{k=1}^m \exp(-v_k(\theta)/\lambda_k) d\theta &= \mathbb{E} \exp(-\frac{v}{\lambda}) \\
&= \mathbb{E} \mathbb{1}_{(0,t)} \exp(-\frac{v}{\lambda}) + \mathbb{E} \mathbb{1}_{(t,\infty)} \exp(-\frac{v}{\lambda}) \\
&\leq \int_0^t f(v) \exp(-\frac{v}{\lambda}) dv + \exp(-\frac{t}{\lambda}) \\
&\leq \sup_{t^* \in (0,t)} f(t^*) \int_0^t \exp(-\frac{v}{\lambda}) dv + \exp(-\frac{t}{\lambda}) \\
&= \sup_{t^* \in (0,t)} f(t^*) \lambda (1 - \exp(-t/\lambda)) + \exp(-\frac{t}{\lambda}) \\
&\leq \sup_{t^* \in (0,t)} f(t^*) \lambda + \exp(-\frac{t}{\lambda})
\end{aligned}$$

(10)

Rearranging terms yields the result. ■

For  $\lambda$  close to 0, the extrinsic measure approaches intrinsic one in total variation distance in  $O(\lambda)$ . This rate is quantified under very general assumption. We expect it can be sharpened under special cases.

We now examine the second case in (4) and (6) where  $\int_{\mathcal{D}} \pi_{\mathcal{R}}(\theta) d\theta = 0$ .

**Remark 3.** If  $\lim_{r \rightarrow 0^+} \frac{\int_{\mathcal{A}} \pi_{\mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}_r^+} d\theta}{\int_{\mathcal{D}_r^+} \pi_{\mathcal{R}}(\theta) d\theta}$  and  $\lim_{r \rightarrow 0^+} \frac{\int_{\mathcal{A}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r^+) d\theta}{\int_{\mathcal{R}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r^+) d\theta}$  converge uniformly in  $\mathcal{A}$ , the total variation distance between the measures of extrinsic and intrinsic distribution

$$\begin{aligned}
\|\Pi(\cdot), \tilde{\Pi}(\cdot)\|_{TV} &= \sup_{\mathcal{A}} \|\Pi(\mathcal{A}) - \tilde{\Pi}(\mathcal{A})\| \\
&= \sup_{\mathcal{A}} \left\| \lim_{r \rightarrow 0^+} \frac{\int_{\mathcal{A}} \pi_{\mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}_r^+} d\theta}{\int_{\mathcal{D}_r^+} \pi_{\mathcal{R}}(\theta) d\theta} - \lim_{r \rightarrow 0^+} \frac{\int_{\mathcal{A}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r^+) d\theta}{\int_{\mathcal{R}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r^+) d\theta} \right\| \\
&= \lim_{r \rightarrow 0^+} \sup_{\mathcal{A}} \left\| \frac{\int_{\mathcal{A}} \pi_{\mathcal{R}}(\theta) \mathbb{1}_{\theta \in \mathcal{D}_r^+} d\theta}{\int_{\mathcal{D}_r^+} \pi_{\mathcal{R}}(\theta) d\theta} - \frac{\int_{\mathcal{A}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r^+) d\theta}{\int_{\mathcal{R}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r^+) d\theta} \right\| \\
&= \lim_{r \rightarrow 0^+} \frac{\int_{\mathcal{R} \setminus \mathcal{D}_r} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta}{\int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta}.
\end{aligned}$$

(11)

Let  $g(r, \lambda) = \frac{\int_{\mathcal{R} \setminus \mathcal{D}_r} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta}{\int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta}$ , using exponential smoothing function and  $\lambda = \sup_k \lambda_k$ ,  $\lim_{\lambda \rightarrow 0^+} g(r, \lambda) =$

0 pointwise in  $r$ . And we have  $\|\Pi(\cdot), \tilde{\Pi}(\cdot)\|_{TV} \leq \lim_{r \rightarrow 0^+} \frac{M_r \exp(-\frac{t}{\lambda}) + M_r \sup_{t^* \in (0, t)} f(t^*)\lambda}{\int_{\mathcal{R} \setminus \mathcal{D}} \pi_{\mathcal{R}}(\theta) \mathcal{K}(\theta; \mathcal{D}_r) d\theta}$  assuming  $M_r = \int_{\mathcal{R} \setminus \mathcal{D}_r} \pi_{\mathcal{R}}(\theta) d\theta < \infty$ . For every  $r$ , there is a small  $\lambda$  to make the total variation distance 0. However, with  $\lambda = \lambda_0$  fixed and non-negligible,  $\lim_{r \rightarrow 0^+} g(r, \lambda) = 1$ . This is an inherent weakness of total variation distance in comparing continuous and degenerate distributions. Therefore, to quantify distance with non-zero  $\lambda$ , we use Wasserstein distance instead:

$$W_p(\Pi, \tilde{\Pi}) = \left( \inf_{\gamma \in \Gamma(\Pi, \tilde{\Pi})} \int d(x, y)^p d\gamma(x, y) \right)^{1/p}$$

where  $d(x, y)$  is the metric on  $\mathcal{R}$  and  $\Gamma(\Pi, \tilde{\Pi})$  the set of couplings of  $\Pi$  and  $\tilde{\Pi}$ .

We continue in Example 1, denoting the bivariate Gaussian distribution obtained extrinsically as  $\text{No}(\mu_\lambda, \Sigma_\lambda)$ , the intrinsic one can be viewed as  $\text{No}(\mu_0, \Sigma_0)$ . The Wasserstein distance between the intrinsic and extrinsic

distributions in Example 1 (Dowson and Landau, 1982) is  $W_2(\Pi, \tilde{\Pi}) = \left( \|\mu_\lambda - \mu_0\|_2^2 + \text{tr}(\Sigma_\lambda + \Sigma_0 - 2(\Sigma_\lambda^{1/2} \Sigma_0 \Sigma_\lambda^{1/2})^{1/2}) \right)^{1/2}$ .

Figure 2 plots  $W_2$  under different values of  $\lambda$ , with  $W_2 \approx 0$  with  $\lambda = 10^{-3}$ .

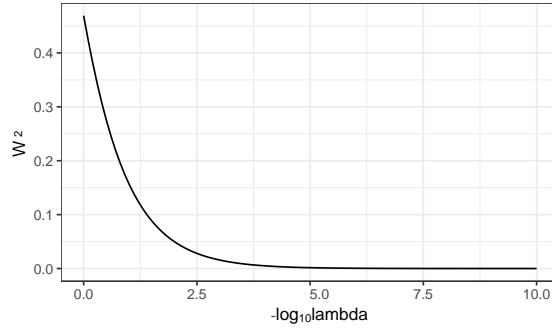


Figure 2: Wasserstein distance of 2nd order  $W_2$  between the intrinsic and extrinsic prior for two normal random variables under sum constraint.  $W_2$  declines rapidly as  $\lambda$  (shown as  $-\log_{10}(\lambda)$ ) decreases.

## 2.4 Approximation Error (Aki's proposal)

**Aki:** The convergence in total variation is too strong, so I propose something like below.

**Leo:** Aki, I like your use of Feder's co-area formula, can we get some non-asymptotic result when  $\lambda > m$ ? The point of total variation being less ideal is that although it DOES work for  $\lambda \rightarrow 0$  (Remark 1 and 3) with almost no assumption, but doesn't work in non-asymptotic case with  $Z(\mathcal{D}) = 0$ . So probably with some assumption, we can get some result of

We establish that the inference based on an intrinsic prior is well approximated by the one based on an extrinsic prior for sufficiently small  $\lambda$ . In particular, for a posterior summary of interest  $g(\theta)$ , we have  $\mathbb{E}_{\tilde{\pi}_{\mathcal{D}, \lambda}}[g(\theta)] \rightarrow \mathbb{E}_{\pi_{\mathcal{D}}}[g(\theta)]$  as  $\lambda \rightarrow 0$  provided  $g(\theta)$  is continuous and integrable under the extrinsic posterior

for some  $\lambda < \infty$ . In case the constraint takes the form  $\mathcal{D} = \{f(\theta) = 0\}$  and the extrinsic posterior is given by  $\tilde{\pi}_{\mathcal{D},\lambda}(\theta) = \pi_{\mathcal{R}}(\theta) \exp(-\lambda^{-1}|f(\theta)|) / Z_{\lambda}$ , the proof is as follows. **Aki: From the proof, it appears that we also need to choose the constraint so that  $\|\nabla f\| \equiv 1$  on  $\mathcal{D}$ . This makes sense from geometric intuition too, I think.**

*Proof.* By the co-area formula of Federer (Diaconis et al., 2013), we have

$$\lambda^{-1} Z_{\lambda} = \lambda^{-1} \int_{\mathbb{R}^d} \pi_{\mathcal{R}}(\theta) \exp(-\lambda^{-1}|f(\theta)|) d\theta = \int_{-\infty}^{\infty} \left[ \int_{f^{-1}(\alpha)} \frac{g(\theta) \pi_{\mathcal{R}}(\theta)}{\|\nabla f(\theta)\|} \mathcal{H}^{d-1}(d\theta) \right] \lambda^{-1} \exp(-\lambda^{-1}|\alpha|) d\alpha \quad (12)$$

where  $\mathcal{H}^k(d\theta)$  denotes a  $k$ -dimensional Hausdorff measure. Since the measure  $\exp(-\lambda^{-1}|\alpha|) d\alpha$  concentrates around  $\alpha = 0$  as  $\lambda \rightarrow 0$ , we obtain

$$\lim_{\lambda \rightarrow 0} \lambda^{-1} Z_{\lambda} = \int_{f^{-1}(0)} \frac{\pi_{\mathcal{R}}(\theta)}{\|\nabla f(\theta)\|} \mathcal{H}^{d-1}(d\theta) = \int_{f^{-1}(0)} \pi_{\mathcal{D}}(d\theta) \quad (13)$$

by the assumption  $\|\nabla f(\theta)\| = 1$  and the fact that  $\pi_{\mathcal{D}}(d\theta)$  coincides with the conditional density of  $\pi_{\mathcal{R}}(\theta)$  on  $\mathcal{D}$ . Also by the co-area formula,

$$\int_{\mathbb{R}^d} g(\theta) \pi_{\mathcal{R}}(\theta) \exp(-\lambda^{-1}|f(\theta)|) / Z_{\lambda} d\theta = \int_{-\infty}^{\infty} \left[ \int_{f^{-1}(\alpha)} \frac{g(\theta) \pi_{\mathcal{R}}(\theta)}{\|\nabla f(\theta)\|} \mathcal{H}^{d-1}(d\theta) \right] \exp(-\lambda^{-1}|\alpha|) / Z_{\lambda} d\alpha \quad (14)$$

Again, since the measure  $\exp(-\lambda^{-1}|\alpha|) d\alpha$  concentrates around  $\alpha = 0$  as  $\lambda \rightarrow 0$ , we obtain

$$\lim_{\lambda \rightarrow 0} \int_{\mathbb{R}^d} g(\theta) \pi_{\mathcal{R}}(\theta) \exp(-\lambda^{-1}|f(\theta)|) / Z_{\lambda} d\theta = \int_{f^{-1}(0)} \frac{g(\theta) \pi_{\mathcal{R}}(\theta)}{\|\nabla f(\theta)\|} d\mathcal{H}^{d-1}(\theta) = \frac{\int_{f^{-1}(0)} g(\theta) \pi_{\mathcal{D}}(d\theta)}{\int_{f^{-1}(0)} \pi_{\mathcal{D}}(d\theta)} \quad \square \quad (15)$$

### 3 Posterior Computation

One particular appeal of the extrinsic approach is its advantage in posterior computation. As it is supported on a less restrictive space  $\mathcal{R}$ , one can exploit conventional sampling tools such as slice sampling, adaptive Metropolis-Hastings and Hamiltonian Monte Carlo (HMC). In this section, we focus on HMC for its easiness to use and good performance in sampling with high-dimensional parameters.

#### 3.1 Hamiltonian Monte Carlo for Extrinsic Posterior Sampling

We provide a brief overview of HMC for continuous  $\theta$  under extrinsic prior. Discrete extension is possible via recent work of Nishimura et al. (2017).

In order to sample from  $\theta \in \mathcal{R} \subset \mathbb{R}^d$ , HMC introduces an auxillary momentum variable  $p \in \mathbb{R}^d$  (commonly generated from  $\text{No}(0, \Sigma)$  with  $\Sigma$  pre-specified), and sample from the joint target density  $\pi(\theta, p) \propto \exp(-H(\theta, p))$  where, in the case of an extrinsic posterior (8),

$$H(\theta, p) = U(\theta) + K(p),$$

$$\text{where } U(\theta) = -\log \{L(\theta; y)\pi_{0, \mathcal{R}}(\theta)\mathcal{K}(\theta; \mathcal{D})\},$$

$$K(p) = \frac{p'\Sigma^{-1}p}{2}.$$
(16)

From the current state  $(\theta^{(0)}, p^{(0)})$ , HMC generates a proposal for Metropolis-Hastings algorithm by simulating Hamiltonian dynamics, which is defined by a differential equation:

$$\frac{\partial \theta^{(t)}}{\partial t} = \frac{\partial H(\theta, p)}{\partial p} = \Sigma^{-1}p,$$

$$\frac{\partial p^{(t)}}{\partial t} = -\frac{\partial H(\theta, p)}{\partial \theta} = -\frac{\partial U(\theta)}{\partial \theta}.$$
(17)

The exact solution to (17) is typically intractable but a valid Metropolis proposal can be generated by numerically approximating (17) with a reversible and volume-preserving integrator (Neal, 2011). The standard choice is the *leapfrog* integrator which approximates the evolution  $(\theta^{(t)}, p^{(t)}) \rightarrow (\theta^{(t+\eta)}, p^{(t+\eta)})$  through the following update equations:

$$p \leftarrow p - \frac{\eta}{2} \frac{\partial U}{\partial \theta}, \quad \theta \leftarrow \theta + \eta \Sigma^{-1}p, \quad p \leftarrow p - \frac{\eta}{2} \frac{\partial U}{\partial \theta}$$
(18)

**Aki:**  $\eta$  is a really unusual notation an HMC integrator. We should plan to replace it with a more standard notation e.g.  $\epsilon$  (most standard though already used elsewhere) or  $\Delta t$ . Taking  $L$  leapfrog steps from the current state  $(\theta^{(0)}, p^{(0)})$  generates a proposal  $(\theta^*, p^*) \approx (\theta^{(L\eta)}, p^{(L\eta)})$ , which is accepted with the probability

$$1 \wedge \exp \left( -H(\theta^*, p^*) + H(\theta^{(0)}, p^{(0)}) \right)$$

### 3.2 Support Expansion and Computing Efficiency

**Aki:** This section requires a major revision; many statements are not quite right. I would need some more time to work on it.

As a Markov chain Monte Carlo, one would hope to use HMC to build a chain that rapidly converges. The convergence rate is associated with the maximal correlation (Liu, 2008),  $\gamma(\theta, \theta^*) = \sup_{g \in L^2(\Pi)} \text{corr}(g(\theta), g(\theta^*))$ , where  $L^2(\Pi) = \{g(\theta) : \text{var}(g(\theta)) < \infty\}$ . Smaller  $\gamma(\theta, \theta^*)$  corresponds to less correlation and faster conver-

gence. Given  $\eta$ , one can adaptively choosing  $L$  so that  $\gamma(\theta, \theta^*)$  falls under certain desired rate (Hoffman and Gelman, 2014).

Since each integrator step is often the bottleneck, one would use  $L$  as small as possible; this minimum step number depends on the step size  $\eta$ . Given a desired convergence rate  $\gamma^*$ , an optimal  $\eta$  would be:

$$\eta^* = \arg \inf_{\eta: \eta \leq \eta_{max}} \inf_L \{L : \gamma(\theta^{(0)}, \theta^{(\eta L)}) \leq \gamma^*\},$$

where  $\eta_{max}$  is the stability bound to prevent the integrator from diverging in approximation, often  $\eta^* \approx \eta_{max}$ .

The stability bound  $\eta_{max}$  is roughly determined by the width of support in the most constrained direction (Neal, 2011). Formally, letting  $\mathcal{S}_\eta$  be the outer contour region of the effective support  $\mathcal{S}_\delta = \{x : \pi(x) \in (0, \delta)\}$  with  $\delta$  arbitrarily small and  $\pi(x)$  denoting the density, the minimum support width is:

$$w_\pi = \inf_{x_1 \in \mathcal{S}_\delta} \sup_{x_2 \in \mathcal{S}_\delta} \{d(x_1, x_2) : \pi(\alpha x_1 + (1 - \alpha)x_2) > 0 \text{ for } \alpha \in [0, 1]\}$$

where  $d(\cdot, \cdot)$  is the appropriate metric in space  $\mathcal{R}$  and is Euclidean distance in continuous HMC.

To provide an intuition for why  $\eta_{max}$  depends on  $w_\pi$ , we focus on one-step update  $L = 1$ . Each update in leap-frog algorithm corresponds to  $\theta^{(\eta)} = \theta^{(0)} + \eta p^{(0)} - \eta^2/2 \frac{\partial U}{\partial \theta} = \theta^{(0)} + \eta p^{(0)} + O(\eta^2)$ . If  $\eta \gg w_\pi$ , a random move in  $\eta p^{(0)}$  can end outside the support with  $U(\theta^{(\eta)}) = \infty$ . This would violate the condition that discrete integrator approximately preserves  $U(\theta^{(\eta)}) + M(p^{(\eta)}) = U(\theta^{(0)}) + M(p^{(0)}) < \infty$ , resulting in divergence.

Since HMC with extrinsic prior operates in  $\mathcal{R}$ , whether constrained space  $\mathcal{D}$  has a small  $w_\pi$  or not should determine how to specify  $\mathcal{K}(\theta; \mathcal{D})$  in extrinsic prior. When  $w_\pi$  is away from 0, one can let  $\mathcal{K}(\theta; \mathcal{D})$  rapidly fall to 0 when  $\theta \notin \mathcal{D}$ . In this case, one can simply use very large  $\lambda$  in exponential smoothing function  $K(v(\theta)) = \exp(-v(\theta)/\lambda)$ , without affecting computing efficiency. On the other hand, if  $w_\pi \approx 0$  for  $\mathcal{D}$ , one needs to let  $\mathcal{K}(\theta; \mathcal{D})$  induce some support expansion, so that the computation can be efficiently carried out. This involves some trade-off between computing time and approximation accuracy. Empirically, we found that  $\lambda = 10^{-3}$  provide a good balance point. These will be illustrated via two examples in the next section.

## 4 Simulated Examples and Application

We now use examples to illustrate the properties of extrinsic priors and their utility in common scenarios.

### 4.1 Simulations

#### Example 2: Linear Regression Under Inequality Constraint

When the support on constrained space  $\mathcal{D}$  has  $w_\pi$  away from 0, one can use extrinsic prior with almost no support expansion. This applies a large class of inequality constraints that has wide support. For example,

consider a linearly constrained regression model:

$$y_i \sim \text{No}(x_i\theta, \sigma^2) \text{ for } i = 1, \dots, n, \quad \text{with } A\theta \leq c$$

where parameter  $\theta$  is a  $p$ -dimensional vector; the constraint parameters  $A$ , a  $d \times p$  matrix, and  $c$ , a  $d$ -dimensional vector, are both given. The inequalities form one or multiple polyhedrons in  $\mathbb{R}^p$ , with  $\mathcal{D}$  as the interior or exterior space.

We consider simple bivariate case  $\theta \in (0, 1)^2$  subject to  $\theta_1 + \theta_2 \leq 1$ , making  $\mathcal{D}$  a triangle. To simulate data, we use  $\sigma^2 = 0.1^2$ ,  $x_i \sim \text{No}([0, 0]', I)$  for  $i = 1, \dots, n$ . We then generate two datasets using different values of  $\theta$  and  $n$ . In the first experiment, we use  $\theta = [0.3, 0.3]'$  with  $n = 10$ , so that the posterior has wide spread and centered in the interior of  $\mathcal{D}$ ; in the second experiment, we use  $\theta = [0.7, 0.3]'$  with  $n = 10^4$  so that the posterior is concentrated on the boundary. In both cases, we assign weakly informative prior for  $\theta \sim \text{No}([0.5, 0.5]', I10^2)$  and inverse-Gamma prior  $\sigma^2 \sim \text{IG}(2, 1)$ .

We use  $\mathcal{K}(\theta) = \exp(-\frac{v(\theta)}{\lambda})$  with  $v(\theta) = |\theta_1 + \theta_2 - 1|_+$  in the extrinsic prior. We choose  $\lambda = 10^{-8}$  leading to almost no support expansion. We collect 10,000 posterior samples efficiently via HMC. The Markov chain mixes rapidly, generating 10,000 effective sample size in both experiments. Figure 3 plots the posterior sample and its contour. There is no posterior that fall outside  $\mathcal{D}$ , thanks to small  $\lambda$ .

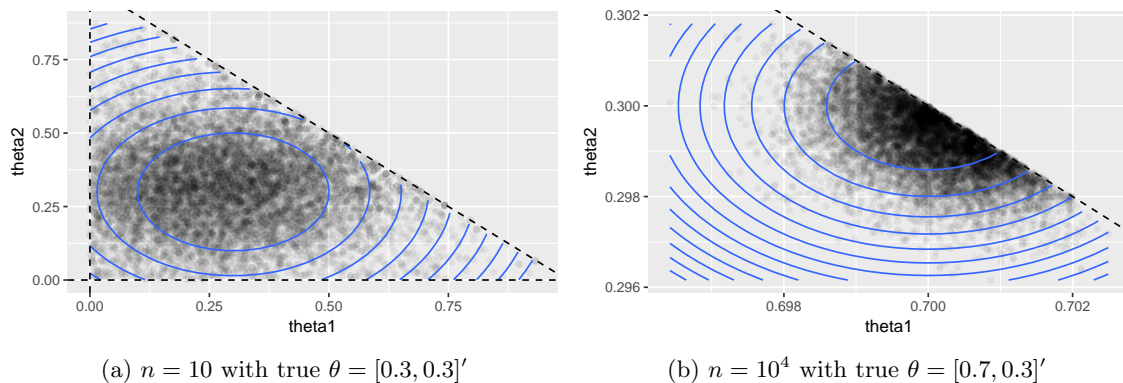


Figure 3: Extrinsic posterior distribution of the normal mean  $\theta$ , with approximation to constraint  $\theta_1 + \theta_2 \leq 1$ . Posterior is either loosely distributed near the center (panel (a)) or concentrated on the boundary (panel (b)) of the region. The extrinsic posterior has no samples outside of the region due to almost no relaxation.

### Example 3: Unit Circle

If the constrained support in  $\mathcal{D}$  has a minimum support width  $w_\pi$  close to 0, it limits the stability bound and computing efficiency in continuous HMC. In this example, we consider  $\mathcal{D}$  is a two-dimensional unit circle  $\{(\theta_1, \theta_2) : \theta_1^2 + \theta_2^2 = 1\}$  (alternatively,  $V(2, 1)$ , a  $(2, 1)$ -Stiefel manifold). In this space,  $w_\pi = 0$  hence support expansion is necessary.

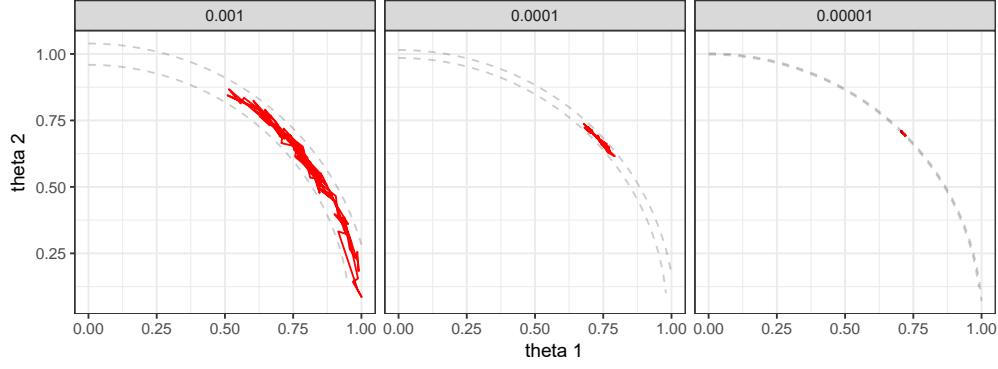
Let data  $y_i \in \mathbb{R}^2$  for  $i = 1, \dots, n$  be noisy realization from one point on unit circle:

$$y_i \sim \text{No}(\theta, I_2 \sigma^2), \text{ with } \theta' \theta = 1,$$

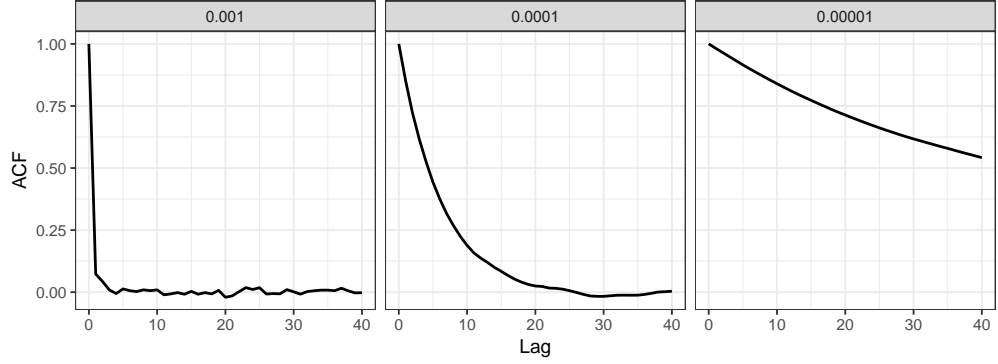
where  $\theta \in \mathcal{D}$  is assigned a von Mises–Fisher prior  $\pi_{0,\mathcal{D}}(\theta) \propto \exp(F' \theta_i)$ .

To generate data, we use  $\theta = (\sqrt{3}/2, 1/2)$ ,  $\sigma^2 = 0.5^2$  and small  $n = 5$ , in order to induce widely spread-out posterior  $\theta$  on the manifold. We then use  $F = (1, 1)$  to induce a weakly informative prior for  $\theta$  and an inverse-Gamma prior  $\text{IG}(2, 1)$  for  $\sigma^2$ . To assign extrinsic prior, we use  $v(\theta) = |\theta' \theta - 1|$  as the distance to circle and extrinsic prior  $\tilde{\pi}_{0,\mathcal{D}}(\theta) = \exp(F' \theta_i) \exp(-\frac{|\theta' \theta - 1|}{\lambda})$ .

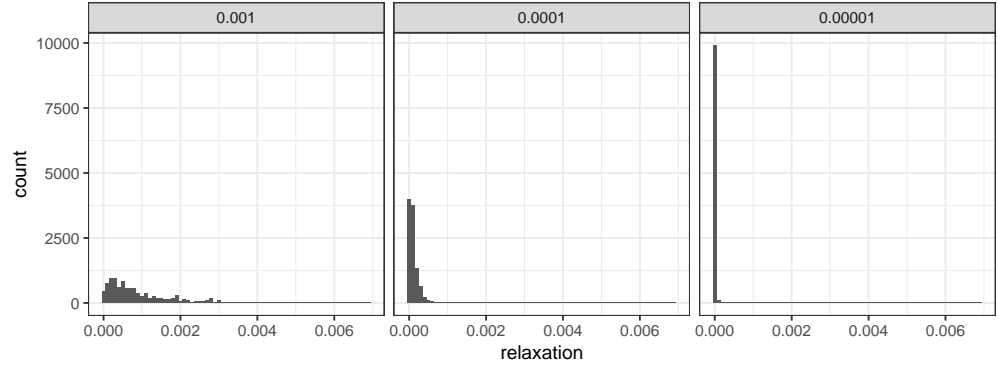
We test  $\lambda = 0.001, 0.0001$  and  $0.00001$  in three experiments. We run each algorithm for 20,000 steps, with first 10,000 discarded. To compare the computing efficiency, we restrict the maximum leap-frog steps  $L$  to be 100 and visualize how much space each algorithm can explore within one HMC iteration. Figure 4(a) plots the path of  $L = 100$  leap-frog steps. Larger  $\lambda$  leads to wider support expansion and larger stability bound  $\eta_{max}$ . This makes  $\theta^{(\eta^L)}$  much less correlated with  $\theta^{(0)}$  at each iteration, measured by autocorrelation (ACF) based on the posterior sample of  $\theta_1$  (Figure 4(b)). Even though we use somewhat small  $\lambda$ , the posterior distance  $v(\theta)$  is small for all three settings (Figure 4(c)), with smaller  $\lambda$  associated with smaller  $v(\theta)$  but lower computing efficiency.



(a) Path of 100 integrator steps in one HMC iteration



(b) Autocorrelation of  $\theta_1$



(c) Posterior distribution of  $|\theta'\theta - 1|$

Figure 4: HMC Sampling on a unit circle, using extrinsic prior with  $\mathcal{K}(\theta) = \exp(-\frac{|\theta'\theta - 1|}{\lambda})$ , with  $\lambda = 0.001$ ,  $0.0001$  and  $0.00001$ . Panel (a) shows the larger relaxation in the narrowest direction of support (orthogonal vector to the circle) can result in more efficient space exploration within 100 leap-frog steps; panel (b) shows the autocorrelation of the posterior sample; panel (c) shows the posterior distribution of the distance to the constraint.

## 4.2 Applications

In statistical modeling, it is common to encounter parameters or latent variables that are non-identifiable. Imposing constraints can often solve or reduce such issue, although they tend to make the computation difficult. We now illustrate utility of extrinsic prior for two applications.

### Example 3: Ordered Dirichlet Prior



We first consider ordered simplex in finite mixture model. A  $(J-1)$ -simplex is a vector  $w = \{w_1, \dots, w_J\}$  with  $1 > w_1 \geq \dots \geq w_J > 0$  and  $\sum_{j=1}^J w_j = 1$ .

For probability simplex, standard practice assigns Dirichlet prior  $Dir(\alpha)$ , with  $\pi_{0,\mathcal{D}}(w) = \prod_{j=1}^J w_j^{\alpha-1} \mathbb{1}_{\sum_{j=1}^J w_j=1}$ . However, this does not accomodate ordering; therefore, the index  $j$  is exchangeable and permutating  $j$ 's does not change the density. This commonly leads to label-switching problem in mixture model estimation (reviewed in Jasra et al. (2005)).

Imposing order constraint yields an ordered Dirichlet prior:

$$\pi_{0,\mathcal{D}}(w_1, \dots, w_J) \propto \prod_{j=1}^J w_j^{\alpha-1} \cdot \mathbb{1}_{\sum_{j=1}^J w_j=1} \cdot \prod_{j=1}^{J-1} \mathbb{1}_{w_j \geq w_{j+1}}. \quad (19)$$

where  $w_j \in (0, 1)$  for  $j = 1, \dots, J$ . The ordered Dirichlet prior can be approximated by extrinsic prior:

$$\tilde{\pi}_{0,\mathcal{D}}(w_1, \dots, w_J) \propto \prod_{j=1}^J w_j^{\alpha-1} \cdot \exp\left(-\frac{\sum_{j=1}^J (w_{j+1} - w_j)_+}{\lambda_1}\right) \exp\left(-\frac{|\sum_{j=1}^J w_j - 1|}{\lambda_2}\right)$$

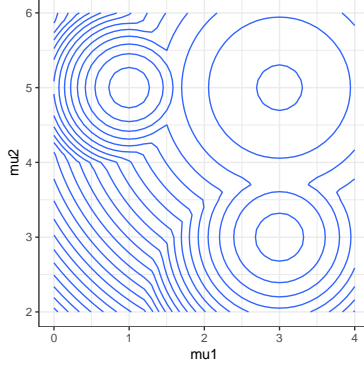
We now adopt this simplex distribution in a normal mixture model with mixture means and common variance, for data  $y_i \in \mathbb{R}^d$  indexed by  $i = 1, \dots, n$ :

$$y_i \overset{indep}{\sim} \text{No}(\mu_i, \Sigma), \quad \mu_i \overset{iid}{\sim} G, \quad G(\cdot) = \sum_{j=1}^J w_j \delta_{\mu_j}(\cdot),$$

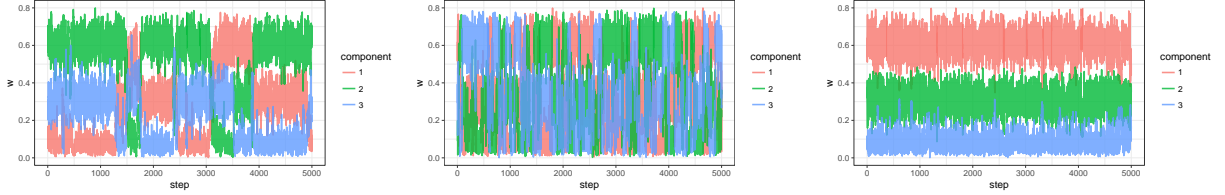
where  $\delta_b(a) = 1$  if  $a = b$  and 0 otherwise.

We generate  $n = 100$  samples from 3 components with true  $\{w_1, w_2, w_3\} = \{0.6, 0.3, 0.1\}$  and two-dimensional means  $\{\mu_1, \mu_2, \mu_3\} = \{[1, 5], [3, 3], [3, 5]\}$  with identity covariance  $\Sigma = I_2$ . We assign weakly informative priors  $\text{No}(0, 10I_2)$  for each  $\mu_j$  and inverse Gamma prior for the diagonal element in  $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2)$  with  $\sigma_1^2, \sigma_2^2 \sim IG(2, 1)$ . We use  $\lambda_1 = 10^{-6}$  to induce almost no relaxation on the ordering and  $\lambda_2 = 10^{-3}$  to allow efficient mixing in embedding a simplex in  $(0, 1)^J$ . To illustrate the benefit of ordered Dirichlet, we also test Gibbs sampling and extrinsic prior method on canonical Dirichlet prior without order constraint.

Figure 5(a) shows the contour of true posterior density of  $\mu_j$ 's. The small component sample size leads to large overlap among the posterior of  $\mu_j$ 's, generating in significant label-switching in both Gibbs and HMC under canonical Dirichlet prior. Figure 5(b,c,d) show the traceplot of  $w$ . Ordered Dirichlet has clearly better convergence due to ordering.



(a) Posterior density of the component means.



(b) Gibbs sampling under canonical Dirichlet (c) HMC sampling under canonical Dirichlet, with extrinsic prior (d) HMC sampling under ordered Dirichlet, with extrinsic prior

Figure 5: Contour of the posterior density of component means and traceplot of the posterior sample for the component weights  $w$ , in a 3-component normal mixture model. Panel (a) shows that there is significant overlap among component means  $\mu_j$ 's, creating label-switching issues in both Gibbs sampling (b) and HMC using canonical prior (c). The ordered Dirichlet prior significantly reducing label-switching (d).

#### Example 4: Orthonormal Tensor Factorization of Multiple Undirected Networks

We now consider a real data application in brain network analysis. The brain connectivity structures are obtained in the data set KKI-42 (Landman et al. 2011), which consists of 21 healthy subjects without any history of neurological disease. Each subject has two brain network observations from scan-rescan, yielding a total of  $n = 42$ . Each observation is a  $V \times V$  symmetric network  $A_i$ , recorded as adjacency matrix  $A_i$  for  $i = 1, \dots, n$ . For the  $i$ th matrix  $A_i$ ,  $A_{i,k,l} \in \{0, 1\}$  is element on the  $k$ th row and  $l$ th column of  $A_i$ , with  $A_{i,k,l} = 1$  indicating there is an connection between  $k$ th and  $l$ th region,  $A_{i,k,l} = 0$  if there is no connection. The regions are constructed via the Desikan et al. (2006) atlas, for a total of  $V = 68$  nodes.

The ambient dimension of observation is  $V(V - 1)/2 = 2,278$ , which is significantly larger than sample size  $n = 40$ . They potentially contain observational error in recording connectivity, and the diagonal in each  $A_i$  is missing due to the lack of interpretable self-connectivity. These facts motivate a probabilistic low-rank model approach. We consider a symmetric tensor decomposition model:

$$A_{i,k,l} \sim \text{Bern}\left(\frac{1}{1 + \exp(-\psi_{i,k,l} - Z_{k,l})}\right)$$

$$\psi_{i,k,l} = \sum_{r_1=1}^{d_1} \sum_{r_2=1}^{d_2} D_{r_1,r_2} W_{i,r_2} U_{k,r_1} U_{l,r_1}$$

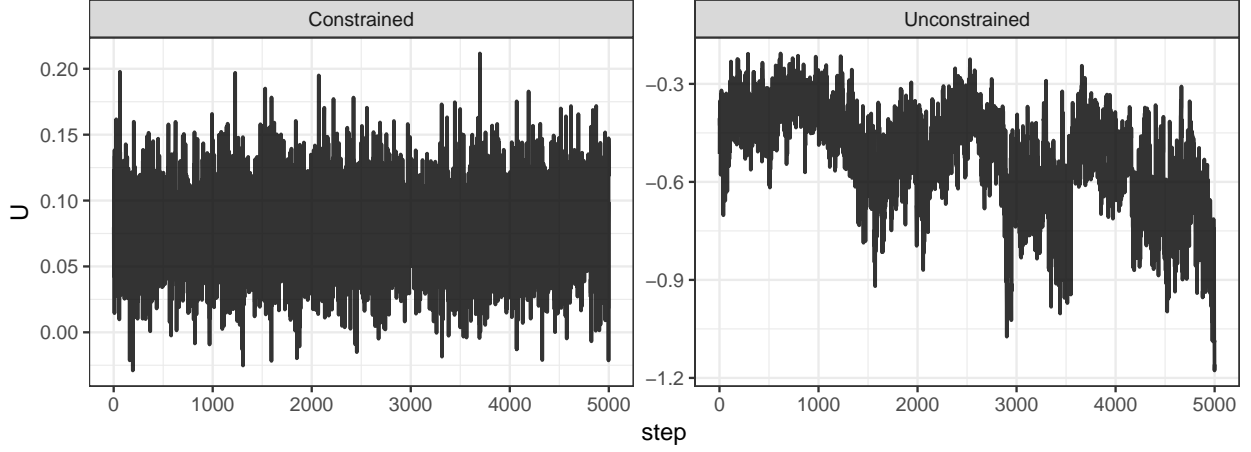
for  $k > l$ ,  $k = 2, \dots, V$ ,  $i = 1, \dots, n$ ;  $U$  is  $V \times d_1$  matrix,  $W$  is  $n \times d_2$  matrix;  $D$  is a  $d_1 \times d_2$  array. The  $V \times V$  matrix  $Z$  is almost unstructural except symmetric  $Z_{k,l} = Z_{l,k}$ , which is commonly used to induce low-rank in the decomposition (Durante et al., 2016).

This model is a special Tucker decomposition with a sparse core tensor, whose diagonal plane is equal to  $D$  and 0 for other elements. The Tucker decomposition is more flexible than another routinely used decomposition, namely parallel factor analysis (PARAFAC). The PARAFAC assumes all ranks are equal and the core tensor  $D$  only has non-zero value when all its sub-indices are equal. In this case, PARAFAC would assume  $d_1 = d_2$ . The additional flexibility in the Tucker is appealing, as one would utilize the varying rank over different sub-direction (mode) of the tensor. On the other hand, a completely unconstrained Tucker decomposition is not identifiable in the matrices and core tensor, due scaling. For example, one can multiply a  $d_1 \times d_1$  non-zero diagonal matrix  $R$ , to  $U$  and obtain  $U^* = UR$  obtain  $D = R^{-2}D_{r_1,\dots}$ . This leaves the likelihood unchanged, creating identifiability issue.

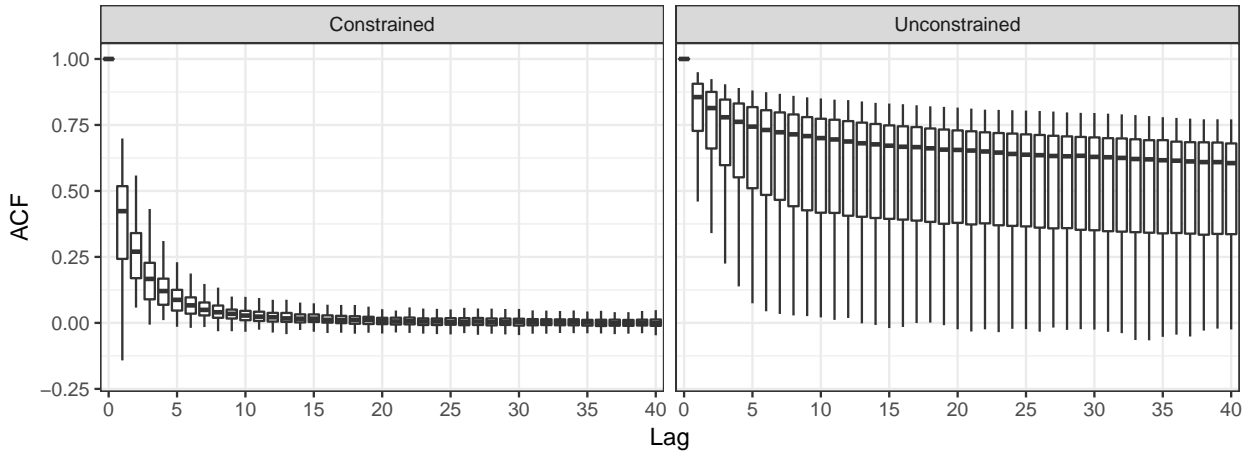
Therefore, we consider applying some constraint on the Tucker decomposition, while still maintaining its varying rank property over different modes. Motivated by high-order singular value decomposition, we impose orthonormality constraints  $U'U = I_{d_1}$  and  $W'W = I_{d_2}$ .

We assign normal prior for  $U_{k,r_2} \sim \text{No}(0, \phi_1)$ ,  $W_{i,r_1} \sim \text{No}(0, \phi_2)$ ,  $Z_{k,l} \sim \text{No}(0, \phi_3)$ ,  $D_{r_1,r_2} \sim \text{No}(0, \phi_{4,r_1,r_2})$  for all  $i, k, l, r_1, r_2$ , and inverse-Gamma prior  $\phi_1, \phi_2, \phi_3 \stackrel{\text{indep}}{\sim} \text{IG}(2, 1)$ ,  $\phi_{4,r_1,r_2} = \tau_{r_1}\tau_{r_2}$ , with  $\tau_{r_1}, \tau_{r_2} \stackrel{\text{indep}}{\sim} \text{IG}(2, 1)$  for all  $r_1, r_2$ .

To allow estimation for model with orthonormality constraint, we use extrinsic prior with  $\mathcal{K}(\theta) = \exp(-\frac{(U'U - I_{d_1})^2 + (W'W - I_{d_2})^2}{\lambda})$  and set  $\lambda = 10^{-3}$ . To compare, we also test with the same model configuration without the orthonormality constraint. We run both models for 10,000 steps and discard the first 5,000 steps. Figure 6 plots the traceplot and autocorrelation for matrix  $U$ . Unconstrained model has severe convergence issue due to the non-identifiability, while constrained model converges and show low autocorrelation for all the parameters.



(a) Traceplot of  $U_{1,1}$ .



(b) ACF of all elements in  $U$

Figure 6: Orthonormality constraint in the tensor decomposition model allows convergence and rapid mixing on the factor matrix (left column); whereas unconstrained model does not converge due to free scaling.

## 5 Discussion

## References

- Boyd, S. and L. Vandenberghe (2004). *Convex optimization*. Cambridge university press.
- Byrne, S. and M. Girolami (2013). Geodesic monte carlo on embedded manifolds. *Scandinavian Journal of Statistics* 40(4), 825–845.
- Diaconis, P., S. Holmes, M. Shahshahani, et al. (2013). Sampling from a manifold. In *Advances in Modern Statistical Theory and Applications: A Festschrift in honor of Morris L. Eaton*, pp. 102–125. Institute of Mathematical Statistics.
- Dowson, D. and B. Landau (1982). The fr chet distance between multivariate normal distributions. *Journal of multivariate analysis* 12(3), 450–455.

- Dunson, D. B. and B. Neelon (2003). Bayesian inference on order-constrained parameters in generalized linear models. *Biometrics* 59(2), 286–295.
- Durante, D., D. B. Dunson, and J. T. Vogelstein (2016). Nonparametric bayes modeling of populations of networks. *Journal of the American Statistical Association* (In press).
- Gelfand, A. E., A. F. Smith, and T.-M. Lee (1992). Bayesian analysis of constrained parameter and truncated data problems using gibbs sampling. *Journal of the American Statistical Association* 87(418), 523–532.
- Girolami, M. and B. Calderhead (2011). Riemann manifold langevin and hamiltonian monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73(2), 123–214.
- Gunn, L. H. and D. B. Dunson (2005). A transformation approach for incorporating monotone or unimodal constraints. *Biostatistics* 6(3), 434–449.
- Hoff, P. D. (2009). Simulation of the matrix bingham–von mises–fisher distribution, with applications to multivariate and relational data. *Journal of Computational and Graphical Statistics* 18(2), 438–456.
- Hoffman, M. D. and A. Gelman (2014). The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *Journal of Machine Learning Research* 15(1), 1593–1623.
- Jasra, A., C. C. Holmes, and D. A. Stephens (2005). Markov chain monte carlo methods and the label switching problem in bayesian mixture modeling. *Statistical Science*, 50–67.
- Khatri, C. and K. Mardia (1977). The von mises-fisher matrix distribution in orientation statistics. *Journal of the Royal Statistical Society. Series B (Methodological)*, 95–106.
- Kolmogorov, A. N. (1950). Foundations of the theory of probability.
- Lin, L. and D. B. Dunson (2014). Bayesian monotone regression using gaussian process projection. *Biometrika*.
- Lin, L., B. St Thomas, H. Zhu, and D. B. Dunson (2016). Extrinsic local regression on manifold-valued data. *Journal of the American Statistical Association* (just-accepted).
- Liu, J. S. (2008). *Monte Carlo strategies in scientific computing*. Springer Science & Business Media.
- Neal, R. M. (2011). Mcmc using hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo* 2, 113–162.
- Nishimura, A., D. Dunson, and J. Lu (2017). Discontinuous hamiltonian monte carlo for sampling discrete parameters. *arXiv preprint arXiv:1705.08510*.

- Pakman, A. and L. Paninski (2014). Exact hamiltonian monte carlo for truncated multivariate gaussians. *Journal of Computational and Graphical Statistics* 23(2), 518–542.
- Winter, S. (2016). Localization results for minkowski contents. *arXiv preprint arXiv:1610.03117*.