

Pathwise Derivatives Beyond the Reparameterization Trick

Martin Jankowiak ^{*1} Fritz Obermeyer ^{*1}

Abstract

We observe that gradients computed via the reparameterization trick are in direct correspondence with solutions of the transport equation in the formalism of optimal transport. We use this perspective to compute (approximate) pathwise gradients for probability distributions not directly amenable to the reparameterization trick: Gamma, Beta, and Dirichlet. We further observe that when the reparameterization trick is applied to the Cholesky-factorized multivariate Normal distribution, the resulting gradients are suboptimal in the sense of optimal transport. We derive the optimal gradients and show that they have reduced variance in a Gaussian Process regression task. We demonstrate with a variety of synthetic experiments and stochastic variational inference tasks that our pathwise gradients are competitive with other methods.

1. Introduction

Maximizing objective functions via gradient methods is ubiquitous in machine learning. When the objective function \mathcal{L} is defined as an expectation of a (differentiable) test function $f_\theta(z)$ w.r.t. a probability distribution $q_\theta(z)$,

$$\mathcal{L} = \mathbb{E}_{q_\theta(z)} [f_\theta(z)] \quad (1)$$

computing exact gradients w.r.t. the parameters θ is often unfeasible so that optimization methods must instead make due with stochastic gradient estimates. If the gradient estimator is unbiased, then stochastic gradient descent with an appropriately chosen sequence of step sizes can be shown to have nice convergence properties (Robbins & Monroe, 1951). If, however, the gradient estimator exhibits large variance, stochastic optimization algorithms may be impractically slow. Thus it is of general interest to develop gradient estimators with reduced variance.

^{*}Equal contribution ¹Uber AI Labs, San Francisco, USA. Correspondence to: <jankowiak@uber.com>, <fritzo@uber.com>.

We revisit the class of gradient estimators popularized in (Kingma & Welling, 2013; Rezende et al., 2014; Titsias & Lázaro-Gredilla, 2014), which go under the name of the pathwise derivative or the reparameterization trick. While this class of gradient estimators is not applicable to all choices of probability distribution $q_\theta(z)$, empirically it has been shown to yield suitably low variance in many cases of practical interest and thus has seen wide use. We show that the pathwise derivative in the literature is in fact a particular instance of a continuous family of gradient estimators. Drawing a connection to tangent fields in the field of optimal transport,¹ we show that one can define a unique pathwise gradient that is optimal in the sense of optimal transport. For the purposes of this paper, we will refer to these optimal gradients as OMT (optimal mass transport) gradients.

The resulting geometric picture is particularly intriguing in the case of multivariate distributions, where each choice of gradient estimator specifies a velocity field on the sample space. To make this picture more concrete, in Figure 1 we show the velocity fields that correspond to two different gradient estimators for the off-diagonal element of the Cholesky factor parameterizing a bivariate Normal distribution. We note that the velocity field that corresponds to the reparameterization trick has a large rotational component that makes it suboptimal in the sense of optimal transport. In Sec. 7 we show that this suboptimality can result in reduced performance when fitting a Gaussian Process to data.

The rest of this paper is organized as follows. In Sec. 2 we provide a brief overview of stochastic gradient variational inference. In Sec. 3 we show how to compute pathwise gradients for univariate distributions. In Sec. 4 we expand our discussion of pathwise gradients to the case of multivariate distributions, introduce the connection to the transport equation, and provide an analytic formula for the OMT gradient in the case of the multivariate Normal. In Sec. 5 we discuss how we can compute high precision approximate pathwise gradients for the Gamma, Beta, and Dirichlet distributions. In Sec. 6 we place our work in the context of related research. In Sec. 7 we demonstrate the performance of our gradient estimators with a variety of synthetic experiments and experiments on real world datasets. Finally, in Sec. 8 we conclude with a discussion of directions for future work.

¹See (Villani, 2003; Ambrosio et al., 2008) for a review.

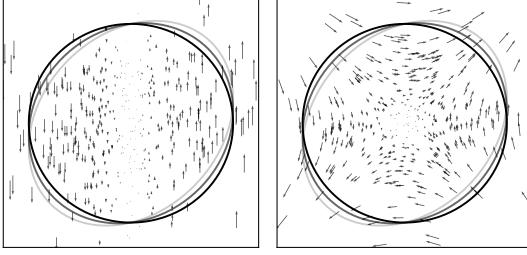


Figure 1. Velocity fields for a bivariate Normal distribution parameterized by a Cholesky factor $\mathbf{L} = \mathbb{1}_2$. The gradient is w.r.t. the off-diagonal element L_{21} . On the left we depict the velocity field corresponding to the reparameterization trick and on the right we depict the velocity field that is optimal in the sense of optimal transport. The solid black circle denotes the 1σ covariance ellipse, with the gray ellipses denoting displaced covariance ellipses that result from small increases in L_{21} . Note that the ellipses evolve the same way under both velocity fields, but *individual* particles flow differently to effect the same global displacement of mass.

2. Stochastic Gradient Variational Inference

One area where stochastic gradient estimators play a particularly central role is stochastic variational inference (Hoffman et al., 2013). This is especially the case for black-box methods (Wingate & Weber, 2013; Ranganath et al., 2014), where conjugacy and other simplifying structural assumptions are unavailable, with the consequence that Monte Carlo estimators become necessary. For concreteness, we will refer to this class of methods as Stochastic Gradient Variational Inference (SGVI). In this section we give a brief overview of this line of research, as it serves as the motivating use case for our work. Furthermore, in Sec. 7 SGVI will serve as the main testbed for our proposed methods.

Let $p(\mathbf{x}, \mathbf{z})$ define a joint probability distribution over observed data \mathbf{x} and latent random variables \mathbf{z} . One of the main tasks in Bayesian inference is to compute the posterior distribution $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})}$. For many models of interest, this is an intractably hard problem and so approximate methods become necessary. Variational inference recasts Bayesian inference as an optimization problem. Specifically we define a variational family of distributions $q_\theta(\mathbf{z})$ parameterized by θ and seek to find a value of θ that minimizes the KL divergence between $q_\theta(\mathbf{z})$ and the (unknown) posterior $p(\mathbf{z}|\mathbf{x})$. This is equivalent to maximizing the ELBO (Jordan et al., 1999), defined as

$$\text{ELBO} = \mathbb{E}_{q_\theta(\mathbf{z})} [\log p(\mathbf{x}, \mathbf{z}) - \log q_\theta(\mathbf{z})] \quad (2)$$

For general choices of $p(\mathbf{x}, \mathbf{z})$ and $q_\theta(\mathbf{z})$, this expectation—much less its gradients—cannot be computed analytically. In these circumstances a natural approach is to build a Monte

Carlo estimator of the ELBO and its gradient w.r.t. θ . The properties of the chosen gradient estimator—especially its bias and variance—play a critical role in determining the viability of the resulting stochastic optimization. Next, we review two commonly used gradient estimators; we leave a brief discussion of more elaborate variants to Sec. 6.

2.1. Score Function Estimator

The score function estimator, also referred to as the log-derivative trick or REINFORCE (Glynn, 1990; Williams, 1992), provides a simple and broadly applicable recipe for estimating ELBO gradients (Paisley et al., 2012). The score function estimator expresses the gradient as an expectation with respect to $q_\theta(\mathbf{z})$, with the simplest variant given by

$$\nabla_\theta \text{ELBO} = \mathbb{E}_{q_\theta(\mathbf{z})} [\nabla_\theta \log r + \log r \nabla_\theta \log q_\theta(\mathbf{z})] \quad (3)$$

where $\log r = \log p(\mathbf{x}, \mathbf{z}) - \log q_\theta(\mathbf{z})$. Monte Carlo estimates of Eqn. 3 can be formed by drawing samples from $q_\theta(\mathbf{z})$ and computing the term in the square brackets. Although the score function estimator is very general (e.g. it applies to discrete random variables) it typically suffers from high variance, although this can be mitigated with the use of variance reduction techniques such as Rao-Blackwellization (Casella & Robert, 1996) and control variates (Ross, 2006).

2.2. Pathwise Gradient Estimator

The pathwise gradient estimator, a.k.a. the reparameterization trick (RT), is not as broadly applicable as the score function estimator, but it generally exhibits lower variance (Price, 1958; Salimans et al., 2013; Kingma & Welling, 2013; Glasserman, 2013; Rezende et al., 2014; Titsias & Lázaro-Gredilla, 2014). It is applicable to continuous random variables whose probability density $q_\theta(\mathbf{z})$ can be reparameterized such that we can rewrite expectations

$$\mathbb{E}_{q_\theta(\mathbf{z})} [f_\theta(\mathbf{z})] \longrightarrow \mathbb{E}_{q_0(\epsilon)} [f_\theta(\mathcal{T}(\epsilon; \theta))] \quad (4)$$

where $q_0(\mathbf{z})$ is a fixed distribution with no dependence on θ and $\mathcal{T}(\epsilon; \theta)$ is a differentiable θ -dependent transformation. Since the expectation w.r.t. $q_0(\epsilon)$ has no θ dependence, gradients w.r.t. θ can be computed by pushing ∇_θ through the expectation. This reparameterization can be done for a number of distributions, including for example the Normal distribution. Unfortunately the reparameterization trick is non-trivial to apply to a number of commonly used distributions, e.g. the Gamma and Beta distributions, since the required shape transformations $\mathcal{T}(\epsilon; \theta)$ inevitably involve special functions.

3. Univariate Pathwise Gradients

Consider an objective function given as the expectation of a test function $f_\theta(z)$ with respect to a distribution $q_\theta(z)$,

where z is a continuous one-dimensional random variable:

$$\mathcal{L} = \mathbb{E}_{q_{\theta}(z)} [f_{\theta}(z)] \quad (5)$$

Here $q_{\theta}(z)$ and $f_{\theta}(z)$ are parameterized by θ , and we would like to compute (stochastic) gradients of \mathcal{L} w.r.t. θ , where θ is a scalar component of θ :

$$\nabla_{\theta} \mathcal{L} = \nabla_{\theta} \mathbb{E}_{q_{\theta}(z)} [f_{\theta}(z)] \quad (6)$$

Crucially we would like to avoid the log-derivative trick, which yields a gradient estimator that tends to have high variance. Doing so will be easy if we can rewrite the expectation in terms of a fixed distribution that does not depend on θ . A natural choice is to use the standard uniform distribution \mathcal{U} ,

$$\mathcal{L} = \mathbb{E}_{\mathcal{U}(u)} [f_{\theta}(F_{\theta}^{-1}(u))] \quad (7)$$

where the transformation $F_{\theta}^{-1} : u \rightarrow z$ is the inverse CDF of $q_{\theta}(z)$. As desired, all dependence on θ is now inside the expectation. Unfortunately, for many continuous univariate distributions of interest (e.g. the Gamma and Beta distributions) the transformation F_{θ}^{-1} (as well as its derivative w.r.t. θ) does not admit a simple analytic expression.

Fortunately, by making use of implicit differentiation we can compute the gradient in Eqn. 6 without explicitly introducing F_{θ}^{-1} . To complete the derivation define u by

$$u \equiv F_{\theta}(z) = \int_{-\infty}^z q_{\theta}(z') dz' \quad (8)$$

and differentiate both sides of Eqn. 8 w.r.t. θ and make use of the fact that $u \sim \mathcal{U}$ does not depend on θ to obtain

$$0 = \frac{dz}{d\theta} q_{\theta}(z) + \int_{-\infty}^z \frac{\partial}{\partial \theta} q_{\theta}(z') dz' \quad (9)$$

This then yields our master formula for the univariate case

$$\frac{dz}{d\theta} = -\frac{\frac{\partial F_{\theta}}{\partial \theta}(z)}{q_{\theta}(z)} \quad (10)$$

where the corresponding gradient estimator is given by

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{q_{\theta}(z)} \left[\frac{df_{\theta}(z)}{dz} \frac{dz}{d\theta} + \frac{\partial f_{\theta}(z)}{\partial \theta} \right] \quad (11)$$

While this derivation is elementary, it helps to clarify things: the key ingredient needed to compute pathwise gradients in Eqn. 6 is the ability to compute (or approximate) the derivative of the CDF, i.e. $\frac{\partial}{\partial \theta} F_{\theta}(z)$. In the supplementary materials we verify that Eqn. 11 results in correct gradients.

It is worth emphasizing how this approach differs from a closely related alternative. Suppose we construct a (differentiable) approximation of the *inverse* CDF, $\hat{F}_{\theta}^{-1}(u) \approx F_{\theta}^{-1}(u)$. For example, we might train a neural network

$\text{nn}(u, \theta) \approx F_{\theta}^{-1}(u)$. We can then push samples $u \sim \mathcal{U}$ through $\text{nn}(u, \theta)$ and obtain approximate samples from $q_{\theta}(z)$ as well as approximate derivatives $\frac{dz}{d\theta}$ via the chain rule; in this case, there will be a mismatch between the probability $q_{\theta}(z)$ assigned to samples z and the actual distribution over z . By contrast, if we use the construction of Eqn. 10, our samples z will still be exact² and the fidelity of our approximation of (the derivatives of) $F_{\theta}(z)$ will only affect the accuracy of our approximation for $\frac{dz}{d\theta}$.

4. Multivariate Pathwise Gradients

In the previous section we focused on continuous univariate distributions. Pathwise gradients can also be constructed for continuous multivariate distributions, although the analysis is in general expected to be much more complicated than in the univariate case—directly analogous to the difference between ordinary and partial differential equations. Before constructing estimators for particular distributions, we introduce the connection to the transport equation.

4.1. The Transport Equation

Consider a multivariate distribution $q_{\theta}(z)$ in D dimensions and consider differentiating $\mathbb{E}_{q_{\theta}(z)} [f(z)]$ with respect to the parameter θ .³ As we vary θ we move $q_{\theta}(z)$ along a curve in the space of distributions over the sample space. Alternatively, we can think of each distribution as a cloud of particles; as we vary θ from θ to $\theta + \Delta\theta$ each particle undergoes an infinitesimal displacement dz . Any set of displacements that ensures that the displaced particles are distributed according to the displaced distribution $q_{\theta+\Delta\theta}(z)$ is allowed. This intuitive picture can be formalized with the transport a.k.a. continuity equation:⁴

$$\frac{\partial}{\partial \theta} q_{\theta} + \nabla_z \cdot (q_{\theta} \mathbf{v}^{\theta}) = 0 \quad (12)$$

Here the *velocity field* \mathbf{v}^{θ} is a vector field defined on the sample space that displaces samples (i.e. particles) z as we vary θ infinitesimally. Note that there is a velocity field \mathbf{v}^{θ} for each component θ of θ . This equation is readily interpreted in the language of fluid dynamics. In order for the total probability to be conserved, the term $\frac{\partial}{\partial \theta} q_{\theta}(z)$ —which is the rate of change of the number of particles in the infinitesimal volume element at z —has to be counterbalanced by the in/out-flow of particles—as given by the divergence term.

²Or rather their exactness will be determined by the quality of our sampler for $q_{\theta}(z)$, which is fully decoupled from how we compute derivatives $\frac{dz}{d\theta}$.

³Here without loss of generality we assume that $f(z)$ has no dependence on θ , since computing $\mathbb{E}_{q_{\theta}(z)} [\nabla_{\theta} f_{\theta}(z)]$ presents no difficulty; the difficulty stems from the dependence on θ in $q_{\theta}(z)$.

⁴We refer the reader to (Villani, 2003) and (Ambrosio et al., 2008) for details.

4.2. Gradient Estimator

Given a solution to Eqn. 12, we can form the gradient estimator

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{q_{\theta}(z)} [\mathbf{v}^{\theta} \cdot \nabla_{\mathbf{z}} f] \quad (13)$$

which generalizes Eqn. 11 to the multivariate case. That this is an unbiased gradient estimator follows directly from the divergence theorem (see the supplementary materials).

4.3. Tangent Fields

In general Eqn. 12 admits an infinite dimensional space of solutions. In the context of our derivation of Eqn. 10, we might loosely say that different solutions of Eqn. 12 correspond to different ways of specifying quantiles of $q_{\theta}(z)$. To determine a *unique*⁵ solution—the tangent field from the theory of optimal transport—we require that

$$\frac{\partial v_i^{\text{OMT}}}{\partial z_j} = \frac{\partial v_j^{\text{OMT}}}{\partial z_i} \quad \forall i, j \quad (14)$$

In this case it can be shown that \mathbf{v}^{OMT} minimizes the total kinetic energy, which is given by⁶

$$K(\mathbf{v}) = \frac{1}{2} \int d\mathbf{z} q_{\theta}(z) \|\mathbf{v}\|^2 \quad (15)$$

4.4. Gradient variance

The $\|\mathbf{v}\|^2$ term that appears in Eqn. 15 might lead one to hope that \mathbf{v}^{OMT} provides gradients that minimize gradient variance. Unfortunately, the situation is more complicated. Denoting the (mean) gradient by $\mathbf{g} = \mathbb{E}_{q_{\theta}(z)}[\mathbf{v} \cdot \nabla_{\mathbf{z}} f(z)]$ the total gradient variance is given by

$$\mathbb{E}_{q_{\theta}(z)} [\|\mathbf{v} \cdot \nabla_{\mathbf{z}} f\|^2] - \|\mathbf{g}\|^2 \quad (16)$$

Since \mathbf{g} is the same for all unbiased gradient estimators, the gradient estimator that minimizes the total variance is the one that minimizes the first term in Eqn. 16. For test functions $f(z)$ that approximately satisfy $\nabla_{\mathbf{z}} f \propto 1$ over the bulk of the support of $q_{\theta}(z)$, the first term in Eqn. 16 term is approximately proportional to the kinetic energy. In this case the OMT gradient estimator will be (nearly) optimal. Note that the kinetic energy weighs contributions from different components of \mathbf{v} equally, whereas \mathbf{g} scales different components of \mathbf{v} with $\nabla_{\mathbf{z}} f$. Thus we can think of the OMT gradient estimator as a good choice for generic choices of $f(z)$ that are relatively flat and isotropic (or, alternatively, for choices of $f(z)$ where we have little *a priori* knowledge about the detailed structure of $\nabla_{\mathbf{z}} f$). So for any particular choice of a generic $f(z)$ there will be some gradient

⁵We refer the reader to Ch. 8 of (Ambrosio et al., 2008) for details.

⁶Note that the univariate solution, Eqn. 10, is automatically the OMT solution.

estimator that has lower variance than the OMT gradient estimator. Still, for *many* choices of $f(z)$ we expect the OMT gradient estimator to have lower variance than the RT gradient estimator, since the latter has no particular optimality guarantees (at least not in any coordinate system that we expect to be well adapted to $f(z)$).

4.5. The Multivariate Normal

In the case of a (zero mean) multivariate Normal distribution parameterized by a Cholesky factor \mathbf{L} via $\mathbf{z} = \mathbf{L}\tilde{\mathbf{z}}$, where $\tilde{\mathbf{z}}$ is white noise, the reparameterization trick yields the following velocity field for L_{ab} :⁷

$$v_i^{\text{RT}} = \frac{\partial z_i}{\partial L_{ab}} = \delta_{ia}(L^{-1}\mathbf{z})_b \quad (17)$$

Note that Eqn. 17 is just a particular instance of the solution to the transport equation that is implicitly provided by the reparameterization trick, namely

$$\mathbf{v}^{\theta} = \left. \frac{\partial \mathcal{T}(\epsilon; \theta)}{\partial \theta} \right|_{\epsilon=\mathcal{T}^{-1}(\mathbf{z}; \theta)} \quad (18)$$

In the supplementary materials we verify that Eqn. 17 satisfies the transport equation Eqn. 12. However, it is evidently *not* optimal in the sense of optimal transport, since $\frac{\partial v_i^{\text{RT}}}{\partial z_j} = \delta_{ia}L_{bj}^{-1}$ is not symmetric in i and j . In fact the tangent field takes the form

$$v_i^{\text{OMT}} = \frac{1}{2} (\delta_{ia}(L^{-1}\mathbf{z})_b + z_a L_{bi}^{-1}) + (S^{ab}\mathbf{z})_i \quad (19)$$

where S^{ab} is a symmetric matrix whose precise form we give in the supplementary materials. We note that computing gradients with Eqn. 19 is $\mathcal{O}(D^3)$, since it involves a singular value decomposition of the covariance matrix. In Sec. 7 we show that the resulting gradient estimator can lead to reduced variance.

5. Numerical Recipes

In this section we show how Eqn. 10 can be used to obtain pathwise gradients in practice. In many cases of interest we will need to derive approximations to $\frac{\partial}{\partial \theta} F(z)$ that balance the need for high accuracy (thus yielding gradient estimates with negligible bias) with the need for computational efficiency. In particular we will derive accurate approximations to Eqn. 10 for the Gamma, Beta, and Dirichlet distributions. These approximations will involve three basic components:

1. Elementary Taylor expansions
2. The Lugannani-Rice saddlepoint expansion (Lugannani & Rice, 1980; Butler, 2007)

⁷Note that the reparameterization trick already yields the OMT gradient for the location parameter μ .

3. Rational polynomial approximations in regions of (z, θ) that are analytically intractable

5.1. Gamma

The CDF of the Gamma distribution involves the (lower) incomplete gamma function $\gamma(\cdot)$: $F_{\alpha, \beta}(z) = \frac{\gamma(\alpha, \beta z)}{\Gamma(\alpha)}$. Unfortunately $\gamma(\cdot)$ does not admit simple analytic expressions for derivatives w.r.t. its first argument, and so we must resort to numerical approximations. Since $z \sim \text{Gamma}(\alpha, \beta = 1) \Leftrightarrow z/\beta \sim \text{Gamma}(\alpha, \beta)$ it is sufficient to consider $\frac{dz}{d\alpha}$ for the standard Gamma distribution with $\beta = 1$.

5.1.1. $z \ll 1$

To give a flavor for the kinds of approximations we use, consider how we can approximate $\frac{\partial}{\partial \alpha} \gamma(\alpha, z)$ in the limit $z \ll 1$. We simply do a Taylor series in powers of z :

$$\begin{aligned} \frac{\partial}{\partial \alpha} \gamma(\alpha, z) &= \frac{\partial}{\partial \alpha} \int_0^z (z')^\alpha (1/z' - 1 + \frac{1}{2}z' + \dots) dz' \\ &= \frac{\partial}{\partial \alpha} z^\alpha \left(\frac{1}{\alpha} - \frac{z}{\alpha+1} + \frac{\frac{1}{2}z^2}{\alpha+2} + \dots \right) \end{aligned}$$

In practice we use 6 terms in this expansion, which is accurate for $z < 0.8$. Details for the remaining approximations can be found in the supplementary materials.

5.2. Beta

The CDF of the Beta distribution, F_{Beta} , is the (regularized) incomplete beta function; just like in the case of the Gamma distribution, its derivatives do not admit simple analytic expressions. We describe the numerical approximations we used in the supplementary materials.

5.3. Dirichlet

Let $z \sim \text{Dir}(\alpha)$ be Dirichlet distributed with n components. Noting that the z_i are constrained to lie within the unit $(n-1)$ -simplex, we proceed by representing z in terms of $n-1$ mutually independent Beta variates (Wilks, 1962):

$$\begin{aligned} \tilde{z}_i &\sim \text{Beta}(\alpha_i, \sum_{j=i+1}^n \alpha_j) \quad \text{for } i = 1, \dots, n-1 \\ z_1 &= \tilde{z}_1 \quad z_n = \prod_{j=1}^{n-1} (1 - \tilde{z}_j) \\ z_i &= \tilde{z}_i \prod_{j=1}^{i-1} (1 - \tilde{z}_j) \quad \text{for } i = 2, \dots, n-1 \end{aligned}$$

Without loss of generality, we will compute $\frac{d}{d\alpha_1} z_i$ for $i = 1, \dots, n$. Crucially, the only dependence on α_1 in Eqn. 20 is through \tilde{z}_1 . We find:

$$\frac{d\mathbf{z}}{d\alpha_1} = -\frac{\frac{\partial F_{\text{Beta}}}{\partial \alpha_1}(z_1 | \alpha_1, \alpha_{\text{tot}} - \alpha_1)}{\text{Beta}(z_1 | \alpha_1, \alpha_{\text{tot}} - \alpha_1)} \times \left(1, \frac{-z_2}{1-z_1}, \dots, \frac{-z_n}{1-z_1} \right) \quad (20)$$

Note that Eqn. 20 implies that $\frac{d}{d\alpha} \sum_i z_i = 0$, as it must because of the simplex constraint. Since we have already

developed an approximation for $\frac{\partial F_{\text{Beta}}}{\partial \theta}$, Eqn. 20 provides a complete recipe for pathwise Dirichlet gradients. Note that although we have used a stick-breaking construction to derive Eqn. 20, this in no way dictates the sampling scheme we use when generating $z \sim \text{Dir}(\alpha)$. In the supplementary materials we verify that Eqn. 20 satisfies the transport equation.

5.4. Implementation

It is worth emphasizing that pathwise gradient estimators of the form in Eqn. 13 have the advantage of being ‘plug-and-play.’ We simply plug an approximate or exact velocity field into our favorite automatic differentiation engine⁸ so that samples z and $f_\theta(z)$ are differentiable w.r.t. θ . There is no need to construct a surrogate objective function to form the gradient estimator.

6. Related Work

A number of lines of research bears upon our work. There is a large body of work on constructing gradient estimators with reduced variance, much of which can be understood in terms of control variates (Ross, 2006): for example, (Mnih & Gregor, 2014) construct neural baselines for score-function gradients; (Schulman et al., 2015) discuss gradient estimators for stochastic computation graphs and their Rao-Blackwellization; and (Tucker et al., 2017; Grathwohl et al., 2017) construct adaptive control variates for discrete random variables. Another example of this line of work is reference (Miller et al., 2017), where the authors construct control variates that are applicable when $q_\theta(z)$ is a diagonal Normal distribution. While our OMT gradient for the multivariate Normal distribution, Eqn. 19, can also be understood in the language of control variates,⁹ (Miller et al., 2017) relies on Taylor expansions of the test function $f_\theta(z)$.¹⁰

In (Graves, 2016), the author derives formula Eqn. 10 and uses it to construct gradient estimators for mixture distributions. Unfortunately, the resulting gradient estimator is expensive, relying on a recursive computation that scales with the dimension of the sample space.

Another line of work constructs partially reparameterized gradient estimators for cases where the reparameterization trick is difficult to apply. The generalized reparameterization gradient (G-REP) (Ruiz et al., 2016) uses standardization

⁸Our approximations for pathwise gradients for the Gamma, Beta, and Dirichlet distributions are available in the 0.4 release of PyTorch (Paszke et al., 2017).

⁹See Sec. 8 and the supplementary materials for a brief discussion.

¹⁰In addition, note that in their approach variance reduction for gradients w.r.t. the scale parameter σ necessitates a multi-sample estimator (at least for high-dimensional models where computing the diagonal of the Hessian is prohibitively expensive).

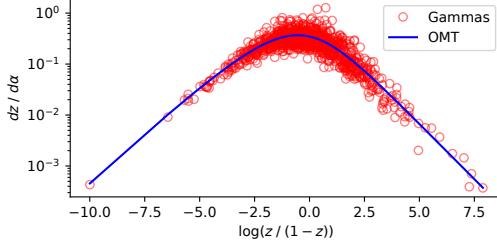


Figure 2. Derivatives $\frac{dz}{d\alpha}$ for samples $z \sim \text{Beta}(1, 1)$. We compare the OMT gradient to the gradient that is obtained when samples $z \sim \text{Beta}(\alpha, \beta)$ are represented as the ratio of two Gamma variates (each with its own pathwise derivative). The OMT derivative has a deterministic value for each sample z , whereas the Gamma representation induces a higher variance stochastic derivative due to the presence of an auxiliary random variable.

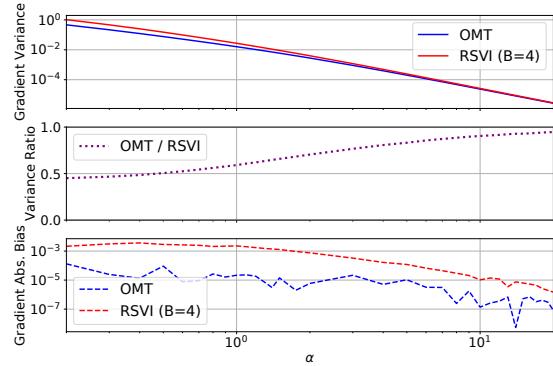


Figure 3. We compare the OMT gradient to the RSVI gradient with $B = 4$ for the test function $f(z) = z^3$ and $q_\theta(z) = \text{Beta}(z|\alpha, \alpha)$. In the bottom panel we depict finite-sample bias for 25 million samples (this also includes effects from finite numerical precision).

via sufficient statistics to obtain a transformation $\mathcal{T}(\epsilon; \theta)$ that minimizes the dependence of $q(\epsilon)$ on θ . This results in a partially reparameterized gradient estimator that also includes a score function-like term.¹¹ In RSVI (Naesseth et al., 2017) the authors consider gradient estimators in the case that $q_\theta(z)$ can be sampled from efficiently via rejection sampling. This results in a gradient estimator with the same generic structure as G-REP, although in the case of RSVI the score function-like term can often be dropped in practice at the cost of small bias (with the benefit of reduced variance). Besides the fact that this gradient estimator is not fully pathwise, one key difference with our approach is that for many distributions of interest (e.g. the Beta and Dirichlet distributions), rejection sampling introduces auxiliary random variables, which results in additional stochasticity and thus higher variance (cf. Figure 2). In contrast our pathwise gradients for the Beta and Dirichlet distributions are *deterministic* for a given z and θ . Finally, (Knowles, 2015) uses (somewhat imprecise) approximations to the inverse CDF to derive gradient estimators for Gamma random variables.

As the final version of this manuscript was being prepared, we became aware of (Figurnov et al., 2018), which has some overlap with this work. In particular, (Figurnov et al., 2018) derives Eqn. 10 and an interesting generalization to the multivariate case. This allows the authors to construct pathwise derivatives for the Gamma, Beta, and Dirichlet distributions. For the latter two distributions, however, the derivatives include additional stochasticity that our pathwise derivatives avoid. Also, the authors do not draw the connection to the transport equation and optimal transport or consider the multivariate Normal distribution in any detail.

¹¹That is a term in the gradient estimator that is proportional to the test function $f_\theta(z)$.

7. Experiments

All experiments in this section use single-sample gradient estimators.

7.1. Synthetic Experiments

In this section we validate our pathwise gradients for the Beta, Dirichlet, and multivariate Normal distributions. Where appropriate we compare to the RT gradient, the score function gradient, or RSVI.

7.1.1. BETA DISTRIBUTION

In Fig. 3 we compare the performance of our OMT gradient for Beta random variables to the RSVI gradient estimator. We use a test function $f(z) = z^3$ for which we can compute the gradient exactly. We see that the OMT gradient performs favorably over the entire range of parameter α that defines the distribution $\text{Beta}(\alpha, \alpha)$ used to compute \mathcal{L} . For smaller α , where \mathcal{L} exhibits larger curvature, the variance of the estimator is noticeably reduced. Notice that one reason for the reduced variance of the OMT estimator as compared to the RSVI estimator is the presence of an auxiliary random variable in the latter case (cf. Figure 2).

7.1.2. DIRICHLET DISTRIBUTION

In Fig. 4 we compare the variance of our pathwise gradient for the Dirichlet distribution to the RSVI gradient estimator. We compute stochastic gradients of the ELBO for a Multinomial-Dirichlet model initialized at the exact posterior (where the exact gradient is zero). The Dirichlet distribution has 1995 components, and the single data point is a bag of words from a natural language document. We see that the pathwise gradient performs favorably over the entire

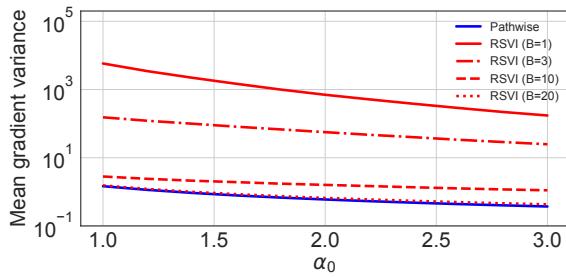


Figure 4. Gradient variance for the ELBO of a conjugate Multinomial-Dirichlet model. We compare the pathwise gradient to RSVI for different boosts B . See Sec. 7.1.2 for details.

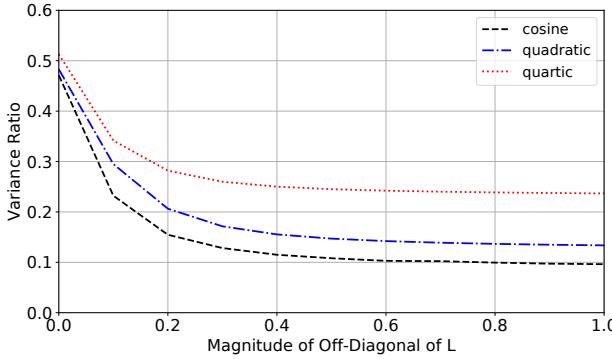


Figure 5. We compare the OMT gradient estimator for the multivariate Normal distribution to the RT estimator for three test functions. The horizontal axis controls the magnitude of the off-diagonal elements of the Cholesky factor L . The vertical axis depicts the ratio of the mean variance of the OMT estimator to that of the RT estimator for the off-diagonal elements of L .

range of the model hyperparameter α_0 considered. Note that as we crank up the shape augmentation setting B , the RSVI variance approaches that of the pathwise gradient.¹²

7.1.3. MULTIVARIATE NORMAL

In Fig. 5 we use synthetic test functions to illustrate the amount of variance reduction that can be achieved with the OMT gradient estimator for the multivariate Normal distribution. The dimension is $D = 50$; the results are qualitatively similar for different dimensions.

¹²As discussed in Sec. 6, the variance of the RSVI gradient estimator can also be reduced by dropping the score function-like term (at the cost of some bias).

7.2. Real World Datasets

In this section we investigate the performance of our gradient estimators for the Gamma, Beta, and multivariate Normal distributions in two variational inference tasks on real world datasets. Note that we include an additional experiment for the multivariate Normal distribution in the supplementary materials, see Sec. 9.11. All the experiments in this section were implemented in the Pyro¹³ probabilistic programming language.

7.2.1. SPARSE GAMMA DEF

The Sparse Gamma DEF (Ranganath et al., 2015) is a probabilistic model with multiple layers of local latent random variables $z_{nk}^{(\ell)}$ and global random weights $w_{kk'}^{(\ell)}$ that mimics the architecture of a deep neural network. Here each n corresponds to an observed data point x_n , ℓ indexes the layer, and k and k' run over the latent components. We consider Poisson-distributed observations x_{nd} for each dimension d . Concretely, the model is specified as¹⁴

$$z_{nk}^{(\ell)} \sim \text{Gamma} \left(\alpha_z, \frac{\alpha_z}{\sum_{k'} z_{nk'}^{(\ell+1)} w_{k'k}^{(\ell)}} \right) \quad \ell = 1, 2, \dots, L-1$$

$$x_{nd} \sim \text{Poisson} \left(\sum_{k'} z_{nk'}^{(1)} w_{k'd}^{(0)} \right) \quad z_{nk}^L \sim \text{Gamma} (\alpha_z, \alpha_z)$$

We set $\alpha_z = 0.1$ and use $L = 3$ layers with 100, 40, and 15 latent factors per data point (for $\ell = 1, 2, 3$, respectively). We consider two model variants that differ in the prior placed on the weights. In the first variant we place Gamma priors over the weights with $\alpha = 0.3$ and $\beta = 0.1$. In the second variant we place β' priors over the weights with the same means and variances as in the first variant.¹⁵ The dataset we consider is the Olivetti faces dataset,¹⁶ which consists of 64×64 grayscale images of human faces. In Fig. 6 we depict how the training set ELBO increases during the course of optimization. We find that on this task the performance of the OMT gradient estimator is nearly identical to RSVI.¹⁷ Figure 6 suggests that gradient variance is not the limiting factor for this particular task and dataset.

¹³<http://pyro.ai>

¹⁴Note that this experiment closely follows the setup in (Ruiz et al., 2016) and (Naesseth et al., 2017).

¹⁵If $z \sim \text{Beta}(\alpha, \beta)$ then $\frac{z}{1-z} \sim \beta'(\alpha, \beta)$. Thus like the Gamma distribution the Beta prime distribution has support on the positive real line.

¹⁶<http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>

¹⁷Note that we do not compare to any alternative estimators such as G-REP, since (Naesseth et al., 2017) shows that RSVI has superior performance on this task.

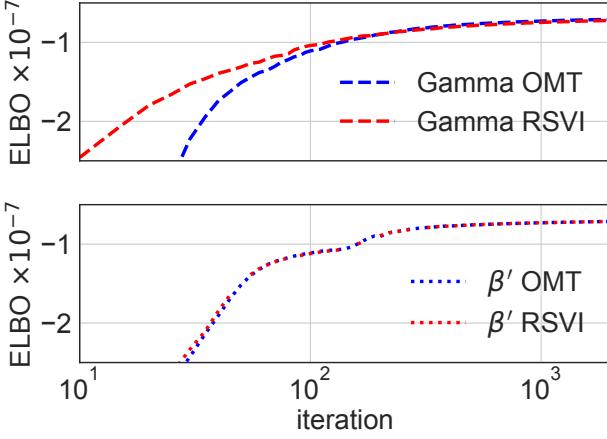


Figure 6. ELBO during training for two variants of the Sparse Gamma DEF, one with and one without Beta random variables. We compare the OMT gradient to RSVI. At each iteration we depict a multi-sample estimate of the ELBO with $N = 100$ samples.

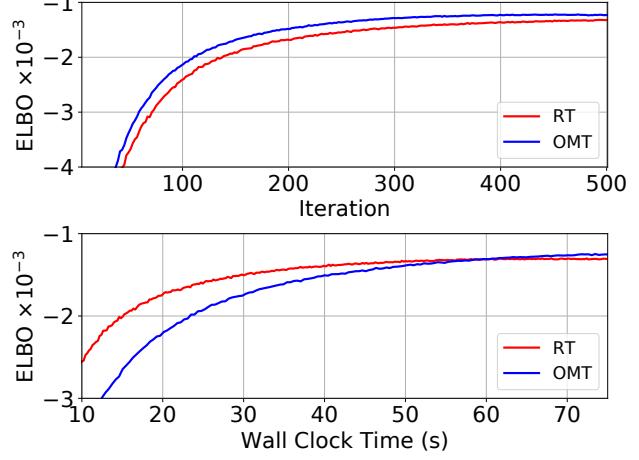


Figure 7. ELBO during training for the Gaussian Process regression task in Sec. 7.2.2. At each iteration we depict a multi-sample estimate of the ELBO with $N = 100$ samples. We compare the OMT gradient estimator to the RT estimator.

7.2.2. GAUSSIAN PROCESS REGRESSION

In this section we investigate the performance of our OMT gradient for the multivariate Normal distribution, Eqn. 19, in the context of a Gaussian Process regression task. We model the Mauna Loa CO₂ data from (Keeling & Whorf, 2004) considered in (Rasmussen, 2004). We use a structured kernel that accommodates a long term linear trend as well as a periodic component. We fit the GP using a single-sample Monte Carlo ELBO gradient estimator and all $D = 468$ data points. The variational family is a multivariate Normal distribution with a Cholesky parameterization for the covariance matrix. Progress on the ELBO during the course of training is depicted in Fig. 7. We can see that the OMT gradient estimator has superior sample efficiency due to its lower variance. By iteration 270 the OMT gradient estimator has attained the same ELBO that the RT estimator attains at iteration 500. Since each iteration of the OMT estimator is ~ 1.9 x slower than the corresponding RT iteration, the superior sample efficiency of the OMT estimator is largely canceled when judged by wall clock time. Nevertheless, the lower variance of the OMT estimator results in a higher ELBO than that obtained by the RT estimator.

8. Discussion and Future Work

We have seen that optimal transport offers a fruitful perspective on pathwise gradients. On the one hand it has helped us formulate pathwise gradients in situations where this was assumed to be impractical. On the other hand it has focused our attention on a particular notion of optimality, which led us to develop a new gradient estimator for the multivariate Normal distribution. A better understanding of this notion

of optimality and, more broadly, a better understanding of when pathwise gradients are preferable over score function gradients (or vice versa) would be useful in guiding the practical application of these methods.

Since each solution of the transport equation Eqn. 12 yields an unbiased gradient estimator, the difference between any two such estimators can be thought of as a control variate. In the case of the multivariate Normal distribution, where computing the OMT gradient has a cost $\mathcal{O}(D^3)$, an attractive alternative to using v^{OMT} is to adaptively choose v during the course of optimization in direct analogy to adaptive control variate techniques. In future work we will explore this approach in detail, which promises lower variance than the OMT estimator at reduced computational cost.

The geometric picture from optimal transport—and thus the potential for non-trivial derivative applications—is especially rich for multivariate distributions. Here we have explored the multivariate Normal and Dirichlet distributions in some detail, but this just scratches the surface of multivariate distributions. It would be of general interest to develop pathwise gradients for a broader class of multivariate distributions, including for example mixture distributions. Rich distributions with low variance gradient estimators are of special interest in the context of SGVI, where the need to approximate complex posteriors demands rich families of distributions that lend themselves to stochastic optimization. In future work we intend to explore this connection further.

Acknowledgements

We thank Peter Dayan and Zoubin Ghahramani for feedback on a draft manuscript and other colleagues at Uber AI Labs—especially Noah Goodman and Theofanis Karaletsos—for stimulating conversations during the course of this work. We also thank Christian Naesseth for clarifying details of the experimental setup for the deep exponential family experiment in (Naesseth et al., 2017).

References

- Ambrosio, Luigi, Gigli, Nicola, and Savaré, Giuseppe. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- Butler, Ronald W. *Saddlepoint approximations with applications*, volume 22. Cambridge University Press, 2007.
- Casella, George and Robert, Christian P. Rao-blackwellisation of sampling schemes. *Biometrika*, 83(1):81–94, 1996.
- Duchi, John, Hazan, Elad, and Singer, Yoram. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Efron, Bradley and Morris, Carl. Data analysis using stein’s estimator and its generalizations. *Journal of the American Statistical Association*, 70(350):311–319, 1975.
- Figurnov, Michael, Mohamed, Shakir, and Mnih, Andriy. Implicit reparameterization gradients. *arXiv preprint arXiv:1805.08498*, 2018.
- Glasserman, Paul. *Monte Carlo methods in financial engineering*, volume 53. Springer Science & Business Media, 2013.
- Glynn, Peter W. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33(10):75–84, 1990.
- Grathwohl, Will, Choi, Dami, Wu, Yuhuai, Roeder, Geoff, and Duvenaud, David. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. *arXiv preprint arXiv:1711.00123*, 2017.
- Graves, Alex. Stochastic backpropagation through mixture density distributions. *arXiv preprint arXiv:1607.05690*, 2016.
- Hoffman, Matthew D, Blei, David M, Wang, Chong, and Paisley, John. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- Jordan, Michael I, Ghahramani, Zoubin, Jaakkola, Tommi S, and Saul, Lawrence K. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- Keeling, Charles David and Whorf, Timothy P. Atmospheric co₂ concentrations derived from flask air samples at sites in the sio network. *Trends: a compendium of data on Global Change*, 2004.
- Kingma, Diederik P and Ba, Jimmy. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kingma, Diederik P and Welling, Max. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Knowles, David A. Stochastic gradient variational bayes for gamma approximating distributions. *arXiv preprint arXiv:1509.01631*, 2015.
- Kucukelbir, Alp, Tran, Dustin, Ranganath, Rajesh, Gelman, Andrew, and Blei, David M. Automatic differentiation variational inference. *arXiv preprint arXiv:1603.00788*, 2016.
- Lugannani, Robert and Rice, Stephen. Saddle point approximation for the distribution of the sum of independent random variables. *Advances in applied probability*, 12(2):475–490, 1980.
- Miller, Andrew, Foti, Nick, D’Amour, Alexander, and Adams, Ryan P. Reducing reparameterization gradient variance. In *Advances in Neural Information Processing Systems*, pp. 3711–3721, 2017.
- Mnih, Andriy and Gregor, Karol. Neural variational inference and learning in belief networks. *arXiv preprint arXiv:1402.0030*, 2014.
- Naesseth, Christian, Ruiz, Francisco, Linderman, Scott, and Blei, David. Reparameterization gradients through acceptance-rejection sampling algorithms. In *Artificial Intelligence and Statistics*, pp. 489–498, 2017.
- Paisley, John, Blei, David, and Jordan, Michael. Variational bayesian inference with stochastic search. *arXiv preprint arXiv:1206.6430*, 2012.
- Paszke, Adam, Gross, Sam, Chintala, Soumith, Chanan, Gregory, Yang, Edward, DeVito, Zachary, Lin, Zeming, Desmaison, Alban, Antiga, Luca, and Lerer, Adam. Automatic differentiation in pytorch. 2017.
- Price, Robert. A useful theorem for nonlinear devices having gaussian inputs. *IRE Transactions on Information Theory*, 4(2):69–72, 1958.

- Ranganath, Rajesh, Gerrish, Sean, and Blei, David. Black box variational inference. In *Artificial Intelligence and Statistics*, pp. 814–822, 2014.
- Ranganath, Rajesh, Tang, Linpeng, Charlin, Laurent, and Blei, David. Deep exponential families. In *Artificial Intelligence and Statistics*, pp. 762–771, 2015.
- Rasmussen, Carl Edward. Gaussian processes in machine learning. In *Advanced lectures on machine learning*, pp. 63–71. Springer, 2004.
- Rezende, Danilo Jimenez, Mohamed, Shakir, and Wierstra, Daan. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- Robbins, Herbert and Monro, Sutton. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Ross, Sheldon M. *Simulation*. Academic Press, San Diego, 2006.
- Ruiz, Francisco R, AUEB, Michalis Titsias RC, and Blei, David. The generalized reparameterization gradient. In *Advances in Neural Information Processing Systems*, pp. 460–468, 2016.
- Salimans, Tim, Knowles, David A, et al. Fixed-form variational posterior approximation through stochastic linear regression. *Bayesian Analysis*, 8(4):837–882, 2013.
- Schulman, John, Heess, Nicolas, Weber, Theophane, and Abbeel, Pieter. Gradient estimation using stochastic computation graphs. In *Advances in Neural Information Processing Systems*, pp. 3528–3536, 2015.
- Stan Manual. Stan modeling language users guide and reference manual, version 2.17.0. <http://mc-stan.org/users/documentation/case-studies/pool-binary-trials.html>, 2017.
- Tieleman, T. and Hinton, G. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
- Titsias, Michalis and Lázaro-Gredilla, Miguel. Doubly stochastic variational bayes for non-conjugate inference. In *International Conference on Machine Learning*, pp. 1971–1979, 2014.
- Tucker, George, Mnih, Andriy, Maddison, Chris J, Lawson, John, and Sohl-Dickstein, Jascha. Rebar: Low-variance, unbiased gradient estimates for discrete latent variable models. In *Advances in Neural Information Processing Systems*, pp. 2624–2633, 2017.
- Villani, Cédric. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.
- Wilks, S.S. *Mathematical Statistics*. John Wiley and Sons Inc., 1962.
- Williams, Ronald J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Wingate, David and Weber, Theophane. Automated variational inference in probabilistic programming. *arXiv preprint arXiv:1301.1299*, 2013.

9. Supplementary Materials

9.1. The Univariate Case

For completeness we show explicitly that the formula

$$\frac{dz}{d\theta} = -\frac{\partial F_\theta(z)}{\partial \theta} \quad (21)$$

yields the correct gradient. Without loss of generality we assume that $f(z)$ has no explicit dependence on θ . Substituting Eqn. 21 for $\frac{dz}{d\theta}$ we have

$$\begin{aligned} \mathbb{E}_{q_\theta(z)} \left[\frac{\partial f}{\partial z} \frac{\partial z}{\partial \theta} \right] &= - \int_{-\infty}^{\infty} \frac{q_\theta(z)}{q_\theta(z)} \frac{\partial f}{\partial z} \int_{-\infty}^z \frac{\partial q_\theta(z')}{\partial \theta} dz' dz \\ &= - \int_{-\infty}^{\infty} \frac{\partial q_\theta(z')}{\partial \theta} \int_{z'}^{\infty} \frac{\partial f}{\partial z} dz dz' \\ &= - \int_{-\infty}^{\infty} \frac{\partial q_\theta(z')}{\partial \theta} (-f(z')) dz' \\ &= \frac{d}{d\theta} E_{q_\theta(z)}[f(z)] \end{aligned} \quad (22)$$

In the second line we changed the order of integration and in the third we appealed to the fundamental theorem of calculus, assuming that $f(z)$ is sufficiently regular that we can drop the boundary term at infinity.

Note that Eqn. 21 is the unique solution $v = \frac{dz}{d\theta}$ to the one-dimensional version of the transport equation that satisfies the boundary condition $\lim_{z \rightarrow \infty} q_\theta v = 0$:

$$\frac{\partial q_\theta}{\partial \theta} + \frac{\partial}{\partial z} (q_\theta v) = 0 \quad (23)$$

9.1.1. EXAMPLE: TRUNCATED UNIT NORMAL

We consider an illustrative case where Eqn. 21 can be computed in closed form. For simplicity we consider the unit Normal distribution truncated¹⁸ to the interval $[0, \kappa]$ with κ

¹⁸As one would expect, Eqn. 21 yields the standard reparameterized gradient in the case of an non-truncated Normal distribution.

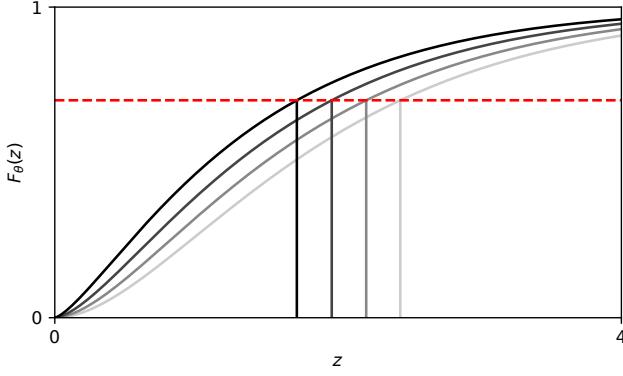


Figure 8. We illustrate how the pathwise derivative is obtained from the CDF in the univariate case. The black curves depict the CDF of the Gamma distribution with $\beta = 1$ and α varying between 1.4 and 2.0. The red line corresponds to a fixed quantile u . As we vary α the point z where the CDF intersects the red line varies. The rate of this variation is precisely the derivative $\frac{dz}{d\alpha}$.

as the only free parameter. A simple computation yields

$$\frac{dz}{d\kappa} = e^{\frac{1}{2}(z^2 - \kappa^2)} \frac{\text{erf}(\frac{z}{\sqrt{2}})}{\text{erf}(\frac{\kappa}{\sqrt{2}})} \quad (24)$$

First, notice that for $z = \kappa$ we have $\frac{dz}{d\kappa} = 1$, which is what we would expect, since $u = 1$ is mapped to the rightmost edge of the interval at $z = \kappa$, i.e. $F_\kappa^{-1}(1) = \kappa$. Similarly we have $\frac{dz}{d\kappa} = 0$ for $z = 0$. For $z \in (0, \kappa)$ the derivative $\frac{dz}{d\kappa}$ interpolates smoothly between 0 and 1. This makes sense, since for a fixed value of u as we get further into the tails of the distribution, nudging κ to the right has a correspondingly larger effect on $z = F_\kappa^{-1}(u)$, while it has a correspondingly smaller effect for u in the bulk of the distribution.

9.1.2. EXAMPLE: UNIVARIATE MIXTURE DISTRIBUTIONS

Consider a mixture of univariate distributions:

$$q_{\theta}(z) = \sum_{k=1}^K \pi_k q_{\theta_k}(z) \quad (25)$$

If we have analytic control over the individual CDFs (or know how to approximate them and their derivatives w.r.t. the parameters) then we can immediately appeal to Eqn. 21. Concretely for derivatives w.r.t. the parameters of each component distribution we have:

$$\frac{\partial z}{\partial \theta_i} = -\frac{\pi_i \frac{\partial F_{\theta_i}}{\partial \theta_i}(z)}{q_{\theta}(z)} \quad (26)$$

Also note that the truncated unit normal is amenable to the reparameterization trick provided that one can compute the inverse error function erf^{-1} .

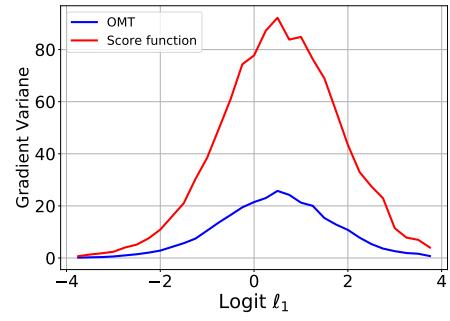


Figure 9. We compare the OMT gradient to the score function gradient for the test function $f(z) = z^4$ where $q_{\theta}(z)$ is a mixture with two components. Depicted is the variance of the gradient w.r.t. the logit ℓ_1 that governs the mixture probability of the first component. The logit of the second component is fixed to be zero.

from which we can get, for example

$$\frac{\partial z}{\partial \mu_i} = \frac{\pi_i q_{\mu_i, \sigma_i}(z)}{q_{\theta}(z)} \quad (27)$$

for a mixture of univariate Normal distributions.

In Fig. 9 we demonstrate that the OMT gradient for a mixture of univariate Normal distributions can have much lower variance than the corresponding score function gradient. Here the mixture has two components with $\mu = (0, 1)$ and $\sigma = (1, 1)$. Note that using the reparameterization trick in this setting would be impractical.

9.2. The Multivariate Case

Suppose we are given a velocity field that satisfies the transport equation:

$$\frac{\partial}{\partial \theta} q_{\theta} + \nabla_z \cdot (q_{\theta} v^{\theta}) = 0 \quad (28)$$

Then, as discussed in the main text, we can form the gradient estimator

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{q_{\theta}(z)} [\mathbf{v}^{\theta} \cdot \nabla_z f] \quad (29)$$

That this gradient estimator is unbiased follows directly from the transport equation and divergence theorem:

$$\begin{aligned} \nabla_{\theta} \mathcal{L} &= \int dz \frac{\partial q_{\theta}(z)}{\partial \theta} f(z) = - \int dz \nabla_z \cdot (q_{\theta} v^{\theta}) f(z) = \\ &= \int dz q_{\theta}(z) \nabla_z f \cdot v^{\theta} = \mathbb{E}_{q_{\theta}(z)} [\nabla_z f \cdot v^{\theta}] \end{aligned} \quad (30)$$

where we appeal to the identity

$$\begin{aligned} \int_V f \nabla_z \cdot (q_{\theta} v^{\theta}) dV &= - \int_V \nabla_z f \cdot (q_{\theta} v^{\theta}) dV + \\ &\quad \oint_S (q_{\theta} f v^{\theta}) \cdot \hat{\mathbf{n}} dS \end{aligned} \quad (31)$$

and assume that $q_{\theta} f \mathbf{v}^{\theta}$ is sufficiently well-behaved that we can drop the surface integral. This is just the multivariate generalization of the derivation in the previous section.

9.3. Multivariate Normal

9.3.1. WHITENED COORDINATES

First we take a look at gradient estimators in whitened coordinates $\tilde{\mathbf{z}} = L^{-1}\mathbf{z}$. The reparameterization trick ansatz for the velocity field can be obtained by transforming the solution in Eqn. 46 (which is also given in the main text) to the new coordinates:

$$\tilde{v}_i \equiv \frac{\partial \tilde{z}_i}{\partial L_{ab}} = L_{ia}^{-1} \tilde{z}_b \quad (32)$$

Note that the transport equation for the multivariate distribution can be written in the form

$$\frac{\partial}{\partial L_{ab}} \log q + \nabla \cdot \tilde{\mathbf{v}} + \tilde{\mathbf{v}} \cdot \nabla \log q = 0 \quad (33)$$

The homogenous equation (i.e. the transport equation without the source term $\frac{\partial \log q}{\partial L_{ab}}$) is then given by

$$\nabla \cdot \tilde{\mathbf{v}} = \tilde{\mathbf{v}} \cdot \tilde{\mathbf{z}} \quad (34)$$

In these coordinates it is evident that infinitesimal rotations, i.e. vector fields of the form

$$\tilde{w}_i = (A\tilde{\mathbf{z}})_i \quad \text{with} \quad A_{ij} = -A_{ji} \quad (35)$$

satisfy¹⁹ the homogenous equation, since

$$\nabla \cdot \tilde{\mathbf{w}} = \text{Tr } A = 0 = \sum_{ij} \tilde{z}_i A_{ij} \tilde{z}_j = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{z}} \quad (36)$$

Finally, if we make the specific choice

$$A_{ij} = \frac{1}{2} (\delta_{ib} L_{ja}^{-1} - \delta_{jb} L_{ia}^{-1}) \quad (37)$$

we find that $\tilde{v}_i + \tilde{w}_i$ (which automatically satisfies the transport equation) and which is given by

$$\tilde{v}_i + \tilde{w}_i \equiv \left(\frac{\partial \tilde{z}_i}{\partial L_{ab}} \right)^{\text{OMT}} = \frac{1}{2} \left(L_{ia}^{-1} \tilde{z}_b + \delta_{ib} \sum_k L_{ka}^{-1} \tilde{z}_k \right)$$

satisfies the symmetry condition

$$\frac{\partial}{\partial \tilde{z}_j} \left(\frac{\partial \tilde{z}_i}{\partial L_{ab}} \right)^{\text{OMT}} = \frac{\partial}{\partial \tilde{z}_i} \left(\frac{\partial \tilde{z}_j}{\partial L_{ab}} \right)^{\text{OMT}} \quad (38)$$

since

$$\frac{\partial}{\partial \tilde{z}_j} \left(\frac{\partial \tilde{z}_i}{\partial L_{ab}} \right)^{\text{OMT}} = \frac{1}{2} (L_{ia}^{-1} \delta_{jb} + L_{ja}^{-1} \delta_{ib}) \quad (39)$$

¹⁹These are in fact not the only solutions; in addition there are non-linear solutions.

which is symmetric in i and j . This implies that the velocity field can be specified as the gradient of a scalar field (this is generally true for the OMT solution), i.e.

$$\left(\frac{\partial \tilde{z}_i}{\partial L_{ab}} \right)^{\text{OMT}} = \frac{\partial}{\partial \tilde{z}_i} \tilde{T}^{ab}(\tilde{\mathbf{z}}) \quad (40)$$

for some $\tilde{T}^{ab}(\tilde{\mathbf{z}})$, which is evidently given by²⁰

$$\tilde{T}^{ab}(\tilde{\mathbf{z}}) = \frac{1}{2} (L^{-\text{T}} \tilde{\mathbf{z}})_a \tilde{z}_b \quad (41)$$

Note, however, that this is not the OMT solution we care about: it minimizes a *different* kinetic energy functional to the one we care about (namely it minimizes the kinetic energy functional in whitened coordinates and not in natural coordinates).

We now explicitly show that solutions of the transport equation that are modified by the addition of an infinitesimal rotation (as in Eqn. 38) still yield valid gradient estimators. Consider a test statistic $f(\tilde{\mathbf{z}})$ that is a monomial in $\tilde{\mathbf{z}}$:

$$f(\tilde{\mathbf{z}}) = \kappa \prod_{i=1}^n \tilde{z}_i^{n_i} \quad (42)$$

It is enough to show that the following expectation vanishes:²¹

$$\mathbb{E}_{q_{\theta}(\tilde{\mathbf{z}})} \left[\sum_{ij} \frac{\partial f}{\partial \tilde{z}_i} A_{ij} \tilde{z}_j \right] \quad (43)$$

where A_{ij} is an antisymmetric matrix. The sum in Eqn. 43 splits up into a sum of paired terms of the form

$$\mathbb{E}_{q_{\theta}(\tilde{\mathbf{z}})} \left[A_{ij} \left(\frac{\partial f}{\partial \tilde{z}_i} \tilde{z}_j - \frac{\partial f}{\partial \tilde{z}_j} \tilde{z}_i \right) \right] \quad (44)$$

We can easily show that each of these paired terms has zero expectation. First note that the expectation is zero if either of i or j is even (since $\mathbb{E}_{q_{\theta}(\tilde{\mathbf{z}})} [\tilde{z}_l^{2k-1}] = 0$). If both i and j are odd we get (using $\mathbb{E}_{q_{\theta}(\tilde{\mathbf{z}})} [\tilde{z}_l^{2k}] = (2k-1)!!$, where $!!$ is the double factorial)

$$\kappa A_{ij} [n_i(n_i-2)!! n_j !! - n_j(n_j-2)!! n_i !!] = 0 \quad (45)$$

Thus, solutions of the transport equation that are modified by the addition of an infinitesimal rotation still yield the same gradient $\nabla_{L_{ab}} \mathbb{E}_{q_{\theta}(\tilde{\mathbf{z}})} [f(\tilde{\mathbf{z}})]$ in expectation.

9.4. Natural Coordinates

We first show that the velocity field \mathbf{v}^{RT} that follows from the reparameterization trick satisfies the transport equation in the (given) coordinates \mathbf{z} , where we have

$$v_i^{\text{RT}} \equiv \frac{\partial z_i}{\partial L_{ab}} = \delta_{ia} (L^{-1} \mathbf{z})_b \quad (46)$$

²⁰Up to an unspecified additive constant.

²¹Note that we can thus think of this term as a control variate.

We have that

$$\begin{aligned}\frac{\partial \log q}{\partial L_{ab}} &= \frac{\partial}{\partial L_{ab}} \left(-\log \det L - \frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z} \right) \\ &= -L_{ba}^{-1} + (\Sigma^{-1} \mathbf{z})_a (L^{-1} \mathbf{z})_b\end{aligned}\quad (47)$$

and

$$\nabla \cdot \mathbf{v}^{\text{RT}} = L_{ba}^{-1} \quad (48)$$

and

$$\mathbf{v}^{\text{RT}} \cdot \nabla \log q = -\mathbf{v}^{\text{RT}} \cdot (\Sigma^{-1} \mathbf{z}) = -(\Sigma^{-1} \mathbf{z})_a (L^{-1} \mathbf{z})_b$$

Thus, the terms cancel term by term and the transport equation is satisfied.

What about the OMT gradient in the natural (given) coordinates \mathbf{z} ? To proceed we represent \mathbf{v} as a linear vector field with symmetric and antisymmetric parts. Imposing the OMT condition determines the antisymmetric part. Imposing the transport equation determines the symmetric part. We find that

$$v_i^{\text{OMT}} = \frac{1}{2} (\delta_{ia} (L^{-1} \mathbf{z})_b + z_a L_{bi}^{-1}) + (S^{ab} \mathbf{z})_i \quad (49)$$

where S^{ab} is the unique symmetric matrix that satisfies the equation

$$\Sigma^{-1} S^{ab} + S^{ab} \Sigma^{-1} = \Xi^{ab} \text{ with } \Xi^{ab} \equiv \xi^{ab} + (\xi^{ab})^T$$

where we define

$$\xi_{ij}^{ab} = \frac{1}{2} (L_{bi}^{-1} \Sigma_{aj}^{-1} - \delta_{ai} (L^{-1} \Sigma^{-1})_{bj}) \quad (50)$$

To explicitly solve Eqn. 50 for S^{ab} we use SVD to write

$$\Sigma^{-1} = U D U^T \quad \text{and} \quad \tilde{\Xi}^{ab} = U^T \Xi^{ab} U \quad (51)$$

where D and U are diagonal and orthogonal matrices, respectively. Then we have that

$$S^{ab} = U \left(\tilde{\Xi}^{ab} \div (D \otimes \mathbb{1} + \mathbb{1} \otimes D) \right) U^T \quad (52)$$

where \div represents elementwise division and \otimes is the outer product. Note that a naive implementation of a gradient estimator based on Eqn. 49 would explicitly construct ξ_{ij}^{ab} , which has size quartic in the dimension. A more efficient implementation will instead make use of ξ_{ij}^{ab} 's structure as a sum of products and never explicitly constructs ξ_{ij}^{ab} .²²

9.4.1. BIVARIATE NORMAL DISTRIBUTION

In Fig. 10 we compare the performance of our OMT gradient for a bivariate Normal distribution to the reparameterization trick gradient estimator. We use a test function $f_\theta(\mathbf{z})$ for which we can compute the gradient exactly. We see that the OMT gradient estimator performs favorably over the entire range of parameters considered.

²²Our implementation can be found here:

https://github.com/uber/pyro/blob/0.2.1/pyro/distributions/omt_mvnp.py

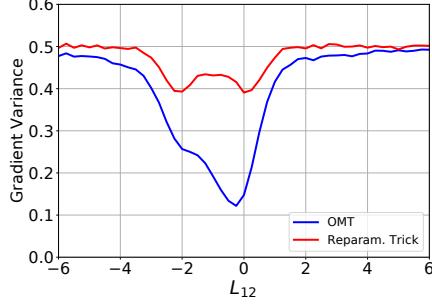


Figure 10. We compare the OMT gradient to the gradient from the reparameterization trick for a bivariate Normal distribution and the test function $f_\theta(\mathbf{z}) = \cos \omega \cdot \mathbf{z}$ with $\omega = (1, 1)$. The Cholesky factor L has diagonal elements $(1, 1)$ and off-diagonal element L_{21} . The gradient is with respect to L_{21} . The variance for the OMT gradient is everywhere lower than for the reparameterization trick gradient.

9.5. Gradient Variance for Linear Test Functions

We use the following example to give more intuition for when we expect OMT gradients for the multivariate Normal distribution to be lower variance than RT gradients. Let $q_\theta(\mathbf{z})$ be the unit normal distribution in D dimensions. Consider the test function

$$f(\mathbf{z}) = \sum_{i=1}^D \kappa_i z_i \quad \mathcal{L} = \mathbb{E}_{q_\theta(\mathbf{z})} [f(\mathbf{z})] \quad (53)$$

and the derivative w.r.t. the off-diagonal elements of the Cholesky factor L . A simple computation yields the total variance of the RT estimator:

$$\sum_{a>b} \text{Var} \left(\frac{\partial \mathcal{L}}{\partial L_{ab}} \right) = \sum_{a>b} \kappa_a^2 \quad (54)$$

Similarly for the OMT estimator we find

$$\sum_{a>b} \text{Var} \left(\frac{\partial \mathcal{L}}{\partial L_{ab}} \right) = \frac{1}{4} \sum_{a>b} (\kappa_a^2 + \kappa_b^2) \quad (55)$$

So if we draw the parameters κ_i from a generic prior we expect the variance of the OMT estimator to be about half of that of the RT estimator. Concretely, if $\kappa_i \sim \mathcal{N}(0, 1)$ then the variance of the OMT estimator will be exactly half that of the RT estimator in expectation. While this computation is for a very specific case—a linear test function and a unit normal $q_\theta(\mathbf{z})$ —we find that this magnitude of variance reduction is typical.

9.6. The Lugannani-Rice Approximation

Saddlepoint approximation methods take advantage of cumulant generating functions (CGFs) to construct (often very

accurate) approximations to probability density functions in situations where full analytic control is intractable.²³ These methods are also directly applicable to CDFs, where a particularly useful approximation—often used by statisticians to estimate various tail probabilities—has been developed by Lugannani and Rice (Lugannani & Rice, 1980). This approximation—after additional differentiation w.r.t. the parameters of the distribution $q_\theta(z)$ —forms the basis of our approximate formulas for pathwise gradients for the Gamma, Beta and Dirichlet distributions in regions of (z, θ) where the (marginal) density is approximately gaussian. As we will see these approximations attain high accuracy.

For completeness we briefly describe the Lugannani-Rice approximation. It is given by:

$$F(z) \approx \begin{cases} \Phi(\hat{w}) + \phi(\hat{w})(1/\hat{w} - 1/\hat{u}) & \text{if } z \neq \mu \\ \frac{1}{2} + \frac{K'''(0)}{6\sqrt{2\pi K''(0)^3}} & \text{if } z = \mu \end{cases} \quad (56)$$

where

$$\hat{w} = \text{sgn}(\hat{s})\sqrt{2\{\hat{s}z - K(\hat{s})\}} \quad \hat{u} = \hat{s}\sqrt{K''(\hat{s})} \quad (57)$$

and where \hat{w} and \hat{u} are functions of z and the saddlepoint \hat{s} , with the saddlepoint defined implicitly by the equation $K'(\hat{s}) = z$. Here $K(s) = \log \mathbb{E}_{q_\theta(z)}[\exp(sz)]$ is the CGF of $q_\theta(z)$, μ is the mean of $q_\theta(z)$, and $\Phi(\cdot)$ and $\phi(\cdot)$ are the CDFs and probability densities of the unit normal distribution. Note that Eqn. 56 appears to have a singularity at $z = \mu$; it can be shown, however, that Eqn. 56 is in fact smooth at $z = \mu$. Nevertheless, in our numerical recipes we will need to take care to avoid numerical instabilities near $z = \mu$ that result from finite numerical precision.

9.7. Gamma Distribution

Our numerical recipe for $\frac{dz}{d\alpha}$ for the standard Gamma distribution with $\beta = 1$ divides (z, α) space into three regions. If $z < 0.8$ we use the Taylor series expansion given in the main text. If $\alpha > 8$ we use the following set of expressions derived from the Lugannani-Rice approximation. Away from the singularity, for $z \gtrless \alpha \pm \delta \cdot \alpha$, we use:

$$\frac{dz}{d\alpha} = \frac{\sqrt{\frac{2}{\alpha}} \frac{\alpha+z}{(\alpha-z)^2} + \log \frac{z}{\alpha} \left(\frac{\sqrt{8\alpha}}{z-\alpha} \pm (z-\alpha-\alpha \log \frac{z}{\alpha})^{-\frac{3}{2}} \right)}{\sqrt{8\alpha}/(z\mathcal{S}_\alpha)} \quad (58)$$

where

$$\mathcal{S}_\alpha \equiv 1 + \frac{1}{12\alpha} + \frac{1}{288\alpha^2}$$

Near the singularity, i.e. for $|z - \alpha| \leq \delta \cdot \alpha$, we use:

$$\frac{dz}{d\alpha} = \frac{1440\alpha^3 + 6\alpha z(53 - 120z) - 65z^2 + \alpha^2(107 + 3600z)}{1244160\alpha^5/(1 + 24\alpha + 288\alpha^2)} \quad (59)$$

²³We refer the reader to (Butler, 2007) for an overview.

Note that Eqn. 59 is derived from Eqn. 58 by a Taylor expansion in powers of $(z - \alpha)$. We set $\delta = 0.1$, which is chosen to balance use of Eqn. 58 (which is more accurate) and Eqn. 59 (which is more numerically stable for $z \approx \alpha$). Finally, in the remaining region ($z > 0.8$ and $\alpha < 8$) we use a bivariate rational polynomial approximation $f(z, \alpha) = \exp\left(\frac{p(z, \alpha)}{q(z, \alpha)}\right)$ where p, q are polynomials in the coordinates $\log(z/\alpha)$ and $\log(\alpha)$, with terms up to order 2 in $\log(z/\alpha)$ and order 3 in $\log(\alpha)$. We fit the rational approximation using least squares on 15696 random (z, α) pairs with α sampled log uniformly between 0.00001 and 10, and z sampled conditioned on α . Our complete approximation for $\frac{dz}{d\alpha}$ is unit tested to have relative accuracy of 0.0005 on a wide range of inputs.

9.8. Beta Distribution

The CDF of the Beta distribution is given by

$$F_{\alpha, \beta}(z) = \frac{B(z; \alpha, \beta)}{B(\alpha, \beta)} \quad (60)$$

where $B(z; \alpha, \beta)$ and $B(\alpha, \beta)$ are the incomplete beta function and beta function, respectively. Our numerical recipe for computing $\frac{dz}{d\alpha}$ and $\frac{dz}{d\beta}$ for the Beta distribution divides (z, α, β) space into three sets of regions. First suppose that $z \ll 1$. Then just like for the Gamma distribution, we can compute a Taylor series of $B(z; \alpha, \beta)$ in powers of z

$$B(z; \alpha, \beta) = z^\alpha \left(\frac{1}{\alpha} + \frac{1-\beta}{1+\alpha}z + \frac{1-\frac{3\beta}{2} + \frac{\beta^2}{2}}{2+\alpha}z^2 + \dots \right) \quad (61)$$

that can readily be differentiated w.r.t. either α or β . Combined with the derivatives of the beta function,

$$\begin{aligned} \frac{d}{d\alpha} B(\alpha, \beta) &= B(\alpha, \beta) (\psi(\alpha) - \psi(\alpha + \beta)) \\ \frac{d}{d\beta} B(\alpha, \beta) &= B(\alpha, \beta) (\psi(\beta) - \psi(\alpha + \beta)) \end{aligned} \quad (62)$$

this gives a complete recipe for approximating $\frac{dz}{d\alpha}$ and $\frac{dz}{d\beta}$ for small z .²⁴ By appealing to the symmetry of the Beta distribution

$$\text{Beta}(z|\alpha, \beta) = \text{Beta}(1-z|\beta, \alpha) \quad (63)$$

we immediately gain approximations to $\frac{dz}{d\alpha}$ and $\frac{dz}{d\beta}$ for $1-z \ll 1$. It remains to specify when these various approximations are applicable. Let us define $\xi = z(1-z)(\alpha + \beta)$. Empirically we find that these approximations are accurate for $\frac{dz}{d\alpha}$ if

1. $z \leq 0.5$ and $\xi < 2.5$; or

²⁴Here $\psi(\cdot)$ is the digamma function, which is available in most advanced tensor libraries.

2. $z \geq 0.5$ and $\xi < 0.75$

with the conditions flipped for $\frac{dz}{d\beta}$. Depending on the precise region, we use 8 to 10 terms in the Taylor series.

Next we describe the set of approximations we derived from the Lugannani-Rice approximation and that we find to be accurate for $\alpha > 6$ and $\beta > 6$. By Eqn. 63 it is sufficient to describe our approximation for $\frac{dz}{d\alpha}$. First define $\sigma = \frac{\sqrt{\alpha\beta}}{(\alpha+\beta)\sqrt{\alpha+\beta+1}}$, the standard deviation of the Beta distribution. Then away from the singularity, for $z \gtrless \frac{\alpha}{\alpha+\beta} \pm \epsilon \cdot \sigma$, we use:

$$\frac{dz}{d\alpha} = \frac{z(1-z) \left(\mathcal{A} + \log \frac{\alpha}{z(\alpha+\beta)} \mathcal{B}_\pm \right)}{\sqrt{\frac{2\alpha\beta}{\alpha+\beta}} \frac{S_{\alpha\beta}}{S_\alpha S_\beta}} \quad (64)$$

with

$$\mathcal{A} = \frac{\beta(2\alpha^2(1-z) + \alpha\beta(1-z) + \beta^2z)}{\sqrt{2\alpha\beta}(\alpha+\beta)^{3/2}(\alpha(1-z) - \beta z)^2}$$

and

$$\mathcal{B}_\pm = \frac{\sqrt{\frac{2\alpha\beta}{\alpha+\beta}}}{\alpha(1-z) - \beta z} \pm \frac{1}{2} \left(\alpha \log \frac{\alpha}{(\alpha+\beta)(1-z)} + \beta \log \frac{\beta}{(\alpha+\beta)z} \right)^{-3/2}$$

Near the singularity, i.e. for $|z - \frac{\alpha}{\alpha+\beta}| \leq \epsilon \cdot \sigma$, we use:

$$\frac{dz}{d\alpha} = \frac{(12\alpha+1)(12\beta+1)(\mathcal{H} + \mathcal{I} + \mathcal{J} + \mathcal{K})}{12960\alpha^3\beta^2(\alpha+\beta)^2(12\alpha+12\beta+1)} \quad (65)$$

with

$$\begin{aligned} \mathcal{H} &= 8\alpha^4(135\beta - 11)(1-z) \\ \mathcal{I} &= \alpha^3\beta(453 - 455z + 1620\beta(1-z)) \\ \mathcal{J} &= 3\alpha^2\beta^2(180\beta - 90z + 59) \\ \mathcal{K} &= \alpha\beta^3(20z(27\beta + 16) + 43) + 47\beta^4z \end{aligned}$$

We set $\epsilon = 0.1$, which is chosen to balance numerical accuracy and numerical stability (just as in the case of the Gamma distribution).

Finally, in the remaining region we use a rational multivariate polynomial approximation

$$f(z, \alpha, \beta) = \frac{p(z, \alpha, \beta)}{q(z, \alpha, \beta)} \frac{z(1-z)}{\beta} (\psi(\alpha + \beta) - \psi(\alpha))$$

where p, q are polynomials in the three coordinates $\log(z)$, $\log(\alpha/z)$, and $\log((\alpha + \beta)z/\alpha)$ with terms up to order 2, 2, and 3 in the respective coordinates. The rational approximation was minimax fit to 2842 points in the remaining region for $0.01 < \alpha, \beta < 1000$. Test points were randomly sampled using log uniform sampling of α, β and stratified sampling of z conditioned on α, β . Minimax fitting achieved about half the maximum error of simple least squares fitting. Our complete approximation for $\frac{dz}{d\alpha}$ and $\frac{dz}{d\beta}$ is unit tested to have relative accuracy of 0.001 on a wide range of inputs.

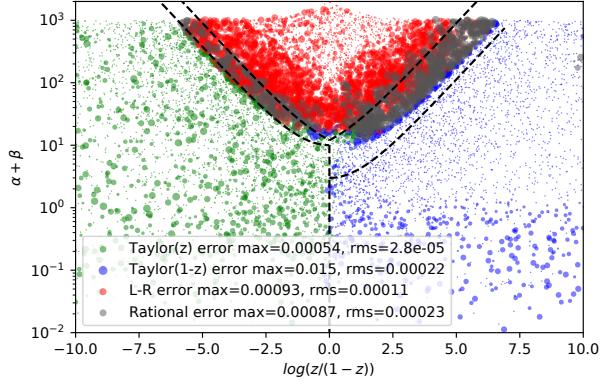


Figure 11. Relative error of our four approximations for $\frac{dz}{d\alpha}$ for the Beta distribution in their respective regions. Note that the region boundaries are in the three-dimensional z, α, β space, so the upper boundaries are only cross-sections.

9.9. Dirichlet Distribution

For completeness we record the general version of the formula for the pathwise gradient (given implicitly in the main text):

$$\frac{dz_i}{d\alpha_j} = -\frac{\frac{\partial F_{\text{Beta}}}{\partial \alpha_j}(z_j | \alpha_j, \alpha_{\text{tot}} - \alpha_j)}{\text{Beta}(z_j | \alpha_j, \alpha_{\text{tot}} - \alpha_j)} \times \left(\frac{\delta_{ij} - z_i}{1 - z_j} \right) \quad (66)$$

We want to confirm that Eqn. 66 satisfies the transport equation for each choice of $j = 1, \dots, n$:

$$\frac{\partial}{\partial \alpha_j} \log q + \nabla \cdot \mathbf{v} + \mathbf{v} \cdot \nabla \log q = 0 \quad (67)$$

Treating z_j as a function of $\mathbf{z}_{-j} = (z_1, \dots, z_{j-1}, z_{j+1}, \dots, z_n)$ everywhere and introducing obvious shorthand for $F_{\text{Beta}}(\cdot)$ and $\text{Beta}(\cdot)$ we have:

$$\begin{aligned} \nabla \cdot \mathbf{v} &= \sum_{i \neq j} \frac{\partial}{\partial z_i} \left(\frac{\frac{\partial F_{\text{Beta}}}{\partial \alpha_j}(z_j | \alpha_j, \alpha_{\text{tot}} - \alpha_j)}{\text{Beta}(z_j | \alpha_j, \alpha_{\text{tot}} - \alpha_j)} \frac{z_i}{\sum_{k \neq j} z_k} \right) \\ &= \frac{\frac{\partial F}{\partial \alpha_j}}{B} \frac{n-2}{1-z_j} - \frac{\partial \log B}{\partial \alpha_j} + \frac{\partial F}{\partial \alpha_j} \frac{(\log B)'}{B} \end{aligned}$$

where $(\log B)'$ is differentiated w.r.t. the argument of $B(z_j)$. We further have that

$$\mathbf{v} \cdot \nabla \log q = \frac{\frac{\partial F}{\partial \alpha_j}}{B} \left(\sum_{i \neq j} \frac{\alpha_i - 1}{1 - z_j} - \frac{\alpha_j - 1}{z_j} \right)$$

and

$$\frac{\partial}{\partial \alpha_j} \log q = \psi(\alpha_j) - \psi(\alpha_{\text{tot}}) + \log z_j$$

Since we have

$$\frac{\partial \log B}{\partial \alpha_j} = \psi(\alpha_j) - \psi(\alpha_{\text{tot}}) + \log z_j$$

and

$$(\log B)' = \frac{\alpha_j - 1}{z_j} - \frac{\alpha_{\text{tot}} - \alpha_j - 1}{1 - z_j}$$

it becomes clear by comparing the individual terms that everything cancels identically and so Eqn. 67 is in fact satisfied by the velocity field in Eqn. 66.

Finally, we note that Eqn. 66 is *not* the OMT solution in the coordinates z_{-j} . It *is* the OMT solution in some coordinate system, but it is not readily apparent which coordinate system that might be.

9.10. Student's t-Distribution

As another example of how to compute pathwise gradients consider Student's t-distribution. Although we have not done so ourselves, it should be straightforward to compute an accurate approximation to Eqn. 21. In the absence of such an approximation, however, we can still get a pathwise gradient for the Student's t-distribution by composing the Normal and Gamma distributions:

$$\begin{aligned} \tau &\sim \text{Gamma}(\nu/2, 1) & x|\tau &\sim \mathcal{N}(0, \tau^{-\frac{1}{2}}) \\ \Rightarrow z &\equiv \sqrt{\frac{\nu}{2}}x \sim \text{Student}(\nu) \end{aligned} \quad (68)$$

Since sampling z like this introduces an auxiliary random degree of freedom, pathwise gradients $\frac{dz}{d\nu}$ computed using Eqn. 68 will exhibit a larger variance than a direct computation of Eqn. 21 would yield.²⁵ The point is that *no additional work* is needed to obtain this particular form of the pathwise gradient: just use pathwise gradients for the Gamma and Normal distributions and the sampling procedure in Eqn. 68.

9.11. Baseball Experiment

To gain more insight into when we expect the OMT gradient estimator for the multivariate Normal distribution to outperform the RT gradient estimator, we conduct an additional experiment. We consider a model for repeated binary trial data (baseball players at bat) using the data in (Efron & Morris, 1975) and the modeling setup in (Stan Manual, 2017) with partial pooling. There are 18 baseball players and the data consists of 45 hits/misses for each player. The model has two global latent variables and 18 local latent variables so that the posterior is 20-dimensional. Specifically, the two global latent random variables are ϕ and κ , with priors $\text{Uniform}(0, 1)$ and $\text{Pareto}(1, 1.5) \propto \kappa^{-5/2}$, respectively. The local latent random variables are given by θ_i for

²⁵Note, however, that this additional variance will decrease as ν increases.

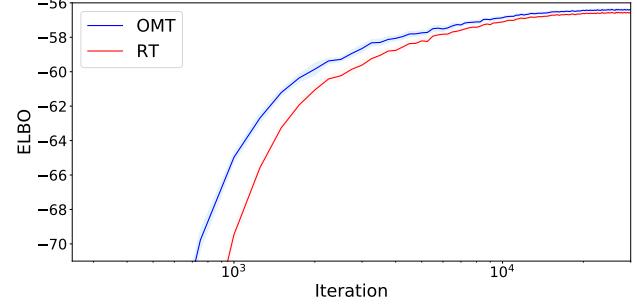


Figure 12. ELBO training curves for the experiment in Sec. 9.11 for the case where the Cholesky factor is initialized far from the identity. Depicted is the mean ELBO for 10 runs with $1 - \sigma$ uncertainty bands around the mean. The OMT gradient estimator learns more quickly than the RT estimator and attains a higher ELBO.

$i = 0, \dots, 17$, with $p(\theta_i) = \text{Beta}(\theta_i | \alpha = \phi\kappa, \beta = (1 - \phi)\kappa)$. The data likelihood factorizes into 45 Bernoulli observations with mean chance of success θ_i for each player i . The variational approximation is formed in the unconstrained space $\{\text{logit}(\phi), \log(\kappa - 1), \text{logit}(\theta_i)\}$ and consists of a multivariate Normal distribution with a full-rank Cholesky factor L . We use the Adam optimizer for training with a learning rate of 5×10^{-3} (Kingma & Ba, 2014).

For this particular model mean field SGVI performs reasonably well, since correlations between the latent random variables are not particularly strong. If we initialize L near the identity, we find that the OMT and RT gradient estimators perform nearly identically, with the difference that the former has an increased computational cost of about 25% per iteration. If, however, we initialize L far from the identity—so that the optimizer has to traverse a considerable distance in L space where the covariance matrix exhibits strong correlations—we find that the OMT estimator makes progress more quickly than the RT estimator and converges to a higher ELBO, see Fig. 12. Generalizing from this, we expect the OMT gradient estimator for the multivariate Normal distribution to exhibit better sample efficiency than the RT estimator in problems where the covariance matrix exhibits strong correlations. This is indeed the case for the GP experiment in the main text, where the learned kernel induces strong temporal correlations.

9.12. Experimental Details

As noted in the main text, we use single-sample gradient estimators in all experiments. Unless noted otherwise, we always include the score function term for RSVI.

9.12.1. MULTIVARIATE NORMAL SYNTHETIC TEST FUNCTION EXPERIMENT

We describe the setup for the experiment corresponding to Fig. 5 in the main text. The dimension is fixed to $D = 50$ and the mean of q_θ is fixed to the zero vector. The Cholesky factor \mathbf{L} that enters into q_θ is constructed as follows. The diagonal of \mathbf{L} consists of all ones. To construct the off-diagonal terms we proceed as follows. We populate the entries below the diagonal of a matrix $\Delta\mathbf{L}$ by drawing each entry from the uniform distribution on the unit interval. Then we define $\mathbf{L} = \mathbb{1}_D + r\Delta\mathbf{L}$. Here r controls the magnitude of off-diagonal terms of \mathbf{L} and appears on the horizontal axis of Fig. 5 in the main text. The three test functions are constructed as follows. First we construct a strictly lower diagonal matrix \mathbf{Q}' by drawing each entry from a bernoulli distribution with probability 0.5. We then define $\mathbf{Q} = \mathbf{Q}' + \mathbf{Q}'^T$. The cosine test function is then given by

$$f(\mathbf{z}) = \cos \left(\sum_{i,j} Q_{ij} z_i / D \right) \quad (69)$$

The quadratic test function is given by

$$f(\mathbf{z}) = \mathbf{z}^T \mathbf{Q} \mathbf{z} \quad (70)$$

The quartic test function is given by

$$f(\mathbf{z}) = (\mathbf{z}^T \mathbf{Q} \mathbf{z})^2 \quad (71)$$

In all cases the gradients can be computed analytically, which makes it easier to reliably estimate the variance of the gradient estimators.

9.12.2. SPARSE GAMMA DEF

Following (Naesseth et al., 2017), we use analytic expressions for each entropy term (as opposed to using the sampling estimate). We use the adaptive step sequence ρ^n proposed by (Kucukelbir et al., 2016) and also used in (Naesseth et al., 2017), which combines RMSPROP (Tieleman & Hinton, 2012) and Adagrad (Duchi et al., 2011):

$$\begin{aligned} \rho^n &= \eta \cdot n^{-1/2+\delta} \cdot \left(1 + \sqrt{s^n} \right)^{-1}. \\ s^n &= t (\hat{g}^n)^2 + (1-t)s^{n-1} \end{aligned} \quad (72)$$

Here $n = 1, 2, \dots$ is the iteration number and the operations in Eqn. 72 are to be understood element-wise. In our case the gradient \hat{g}^n is always a single-sample estimate. We fix $\delta = 10^{-16}$ and $t = 0.1$. In contrast to (Kucukelbir et al., 2016) but in line with (Naesseth et al., 2017) we initialize s_0 at zero. To choose η we did a grid search for each gradient estimator and each of the two model variants. Specifically, for each η we did 100 training iterations for three trials with different random seeds and then chose the η that yielded the

highest mean ELBO after 100 iterations. This procedure led to the selection of $\eta = 4.5$ for the first model variant and $\eta = 30$ for the second model variant (note that within each model variant the gradient estimators preferred the same value of η). For the first model variant we included the score function-like term in the RSVI gradient estimator, while we did not include it for the second model variant, as we found that this hurt performance. In both cases we used the shape augmentation setting $B = 4$, which was also used for the results reported in (Naesseth et al., 2017). After fixing η we trained the model for 2000 iterations, initializing with another random number seed. The figure in the main text shows the training curves for that single run. We confirmed that other random number seeds give similar results. A reference implementation can be found here:

https://github.com/uber/pyro/blob/0.2.1/examples/sparse_gamma_def.py

9.12.3. GAUSSIAN PROCESS REGRESSION

We used the Adam optimizer (Kingma & Ba, 2014) to optimize the ELBO with single-sample gradient estimates. We chose the Adam hyperparameters by doing a grid search over the learning rate and β_1 . For each combination (lr, β_1) we did 20 training iterations for three trials with different random seeds and then chose the combination that yielded the highest mean ELBO after 20 iterations. This procedure led to the selection of a learning rate of 0.030 and $\beta_1 = 0.50$ for both gradient estimators (OMT and reparameterization trick). We then trained the model for 500 iterations, initializing with another random number seed. The figure in the main text shows the training curves for that single run. We confirmed that other random number seeds give similar results.