
Interpolating Between Sampling and Variational Inference with Infinite Stochastic Mixtures (Supplementary material)

Richard D. Lange¹

Ari S. Benjamin¹

Ralf M. Haefner^{*2}

Xaq Pitkow^{*3}

¹Dept. of Neurobiology, University of Pennsylvania, Philadelphia, Pennsylvania, USA

²Dept. of Brain and Cognitive Sciences, University of Rochester, Rochester, New York, USA

³Baylor College of Medicine, Rice University Houston, Texas, USA

^{*}equal contribution

A PROOFS AND DERIVATIONS

Throughout, we assume that θ forms a minimal statistical manifold [Amari, 2016], so that the degrees of freedom of q match the dimensionality of θ , and whenever $q(\mathbf{x}; \theta^{(i)}) = q(\mathbf{x}; \theta^{(j)})$ for all \mathbf{x} , it must be that $\theta^{(i)} = \theta^{(j)}$.

Recall that in the main text, we defined the following objective:

$$\mathcal{L}(\psi, \lambda) \equiv \mathcal{I}[\mathbf{x}; \theta] - \lambda \mathbb{E}_{\psi(\theta)} [\text{KL}(q(\mathbf{x}; \theta) \| p^*(\mathbf{x}))], \quad ((5) \text{ restated})$$

where $\lambda \in [1, \infty)$ is a hyper-parameter, and ψ is a probability density on θ . We also introduced an **approximate objective** in which $\mathcal{I}[\mathbf{x}; \theta]$ is replaced with

$$\mathcal{I}_{\mathcal{F}}[\mathbf{x}; \theta] \equiv \mathcal{H}[\theta] - \frac{1}{2} \mathbb{E}_{\psi(\theta)} [\log |2\pi e \mathcal{F}(\theta)^{-1}|]. \quad ((8) \text{ restated})$$

This approximate objective is

$$\mathcal{L}_{\mathcal{F}}(\psi, \lambda) = \mathcal{H}[\theta] + \mathbb{E}_{\psi(\theta)} \left[\frac{1}{2} \log |\mathcal{F}(\theta)| - \lambda \text{KL}(q(\mathbf{x}; \theta) \| p^*(\mathbf{x})) \right], \quad ((9) \text{ restated})$$

and it is maximized for a given λ by

$$\psi(\theta) = \frac{1}{Z(\lambda)} \exp \left(\frac{1}{2} \log |\mathcal{F}(\theta)| - \lambda \text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x})) \right) \quad ((10) \text{ restated})$$

where $Z(\lambda) = \int_{\theta} \exp \left(\frac{1}{2} \log |\mathcal{F}(\theta)| - \lambda \text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x})) \right) d\theta$.

A.1 CHARACTERIZING THE PARETO FRONT

Let us begin with a set of results regarding the shape of the Pareto front that connects VI to Sampling in Figure 2.

Lemma 1 $\mathcal{L}(\psi, \lambda)$ is concave in ψ , i.e. $\mathcal{L}(\omega\psi_1 + (1 - \omega)\psi_2, \lambda) \geq \omega\mathcal{L}(\psi_1, \lambda) + (1 - \omega)\mathcal{L}(\psi_2, \lambda)$ for $0 \leq \omega \leq 1$. Further, $\mathcal{L}_{\mathcal{F}}(\psi, \lambda)$ is strictly concave in ψ .

Proof: The proof for \mathcal{L} follows from the fact that $\mathbb{E}_{\psi(\theta)} [\text{KL}(q(\mathbf{x}; \theta) \| p^*(\mathbf{x}))]$ is linear in ψ , and $\mathcal{I}[\mathbf{x}; \theta]$ is known to be concave in the marginal distribution of either variable [Braverman and Bhowmick, 2011]. The proof for $\mathcal{L}_{\mathcal{F}}$ is similar: the $\mathbb{E}_{\psi(\theta)} [\frac{1}{2} \log |\mathcal{F}(\theta)|]$ term is linear in ψ , and $\mathcal{H}[\theta]$ is strictly concave in ψ . This can be seen, for instance, by taking the

second variational derivative of $\mathcal{H}[\theta]$ with respect to ψ :

$$\begin{aligned}\nabla_{\psi}^2 \mathcal{H}[\theta] |_{\theta_i \theta_j} &= \nabla_{\psi} \left(\nabla_{\psi} \mathcal{H}[\theta] |_{\theta_i} \right) |_{\theta_j} \\ &= \nabla_{\psi} \left(-\nabla_{\psi} \int_{\theta} \psi(\theta) \log \psi(\theta) d\theta |_{\theta_i} \right) |_{\theta_j} \\ &= \nabla_{\psi} (-1 - \log \psi(\theta_i)) |_{\theta_j} \\ &= \begin{cases} -\frac{1}{\psi(\theta_i)} & \text{if } \theta_i = \theta_j \\ 0 & \text{otherwise.} \end{cases}\end{aligned}$$

Since $\psi(\theta) \geq 0$ everywhere, this implies that the curvature of $\mathcal{H}[\theta]$ is strictly negative at all values of θ . \blacksquare

Lemma 2 Let $\mathcal{I}^*(\lambda)$ and $\mathbb{E}[\text{KL}]^*(\lambda)$ denote the values of Mutual Information and Expected KL achieved by optima of \mathcal{L} for a given λ . Then, λ defines the slope of the Pareto front:

$$\lambda = \frac{d\mathcal{I}^*/d\lambda}{d\mathbb{E}[\text{KL}]^*/d\lambda}.$$

Or, in the case of $\mathcal{L}_{\mathcal{F}}$, λ similarly defines the slope of

$$\lambda = \frac{d\mathcal{I}_{\mathcal{F}}^*/d\lambda}{d\mathbb{E}[\text{KL}]^*/d\lambda},$$

with $\mathcal{I}_{\mathcal{F}}$ in place of \mathcal{I} .

Proof: This follows from viewing \mathcal{L} as the Lagrangian of a constrained optimization problem, with λ as a Lagrange multiplier. The same argument applies to both \mathcal{L} and \mathcal{I} as to $\mathcal{L}_{\mathcal{F}}$ and $\mathcal{I}_{\mathcal{F}}$, so we will just give the proof for one. Consider the constrained optimization problem of maximizing \mathcal{I} (or $\mathcal{I}_{\mathcal{F}}$) subject to the constraint that $\mathbb{E}[\text{KL}(q||p)] = C$. The Lagrangian for this problem is identical to (5), but with C added:

$$\mathcal{L}(\psi, \lambda) \equiv \mathcal{I}[\mathbf{x}; \theta] - \lambda (\mathbb{E}_{\psi(\theta)} [\text{KL}(q(\mathbf{x}; \theta)||p^*(\mathbf{x}))] - C)$$

Optimizing with respect to ψ , this is a concave maximization problem with a linear constraint. A well-known property of such problems is that, at the solution, the Lagrange multiplier (λ) is equal to the change in the objective (\mathcal{I}^*) per change in the constraint (C), or $\lambda = \frac{d\mathcal{I}^*}{dC}$. Since C is the constrained value of $\mathbb{E}[\text{KL}(q||p)]$, we also immediately have $\frac{d\mathbb{E}[\text{KL}]^*}{dC} = 1$. This implies that

$$\lambda = \frac{d\mathcal{I}_{\mathcal{F}}^*/dC}{d\mathbb{E}[\text{KL}]^*/dC}.$$

So far, we have treated λ as a function of C , but for all values of λ that correspond to a unique C , we can invert this relationship and treat C as a function of λ . Then, assuming $\frac{dC}{d\lambda} \neq 0$ for all $1 \leq \lambda < \infty$ that we are interested in, we have

$$\lambda = \frac{d\mathcal{I}_{\mathcal{F}}^*/dC \times dC/d\lambda}{d\mathbb{E}[\text{KL}]^*/dC \times dC/d\lambda} = \frac{d\mathcal{I}_{\mathcal{F}}^*/d\lambda}{d\mathbb{E}[\text{KL}]^*/d\lambda}.$$

Again using the fact that $C = \mathbb{E}[\text{KL}]^*$ by construction, the assumption that $\frac{dC}{d\lambda} \neq 0$ is equivalent to $\frac{d\mathbb{E}[\text{KL}]^*}{d\lambda} \neq 0$. In other words, as long as changing λ has some effect on $\mathbb{E}[\text{KL}]^*$, the combined effect on \mathcal{I}^* and $\mathbb{E}[\text{KL}]^*$ will be such that $\lambda = \frac{d\mathcal{I}^*}{d\mathbb{E}[\text{KL}]^*}$. \blacksquare

A.2 SAMPLING-LIKE BEHAVIOR OF OUR METHOD

Recall our definition of sampling:

Definition 1 (Sampling) A stochastic mixture, defined by the component family $q(\mathbf{x}; \theta)$ and mixing distribution $\psi(\theta)$, is considered to be “sampling” if it is **unbiased** and it consists of **non-overlapping components**.

An **unbiased** mixture is one where $m(\mathbf{x}) = p(\mathbf{x})$.

A mixture consists of T **non-overlapping components** if $\sum_{t=1}^T q(\mathbf{x}; \theta^{(t)}) \approx \max_t q(\mathbf{x}; \theta^{(t)})$ with high probability.

We will assume throughout this section that q is a location-scale family, and in particular Gaussian for Lemma 4, but it may be fruitful for future work to consider other families of mixture components.

Lemma 3 *Sampling is an optimum of the original objective, \mathcal{L} , when $\lambda = 1$.*

Proof: When $\lambda = 1$, \mathcal{L} simplifies back to $\text{KL}(m||p)$. Any **unbiased** mixture is a minimum of $\text{KL}(m||p)$. ■

Note, however, that this does not imply sampling is the unique optimum. In general there may be other unbiased mixing distributions $\psi(\theta)$ such that $m(\mathbf{x}) = p(\mathbf{x})$. For instance, if q is Gaussian and $p(\mathbf{x})$ is itself a finite mixture of Gaussians, then $\psi(\theta)$ could concentrate on exactly those modes in p . In any case where there two such unbiased ψ s, there are in fact infinitely many unbiased, since any mixture of them, $\alpha\psi_1(\theta) + (1 - \alpha)\psi_2(\theta)$, will also be unbiased. Among all unbiased mixtures, sampling may in some sense be the worst choice – we conjecture that it has the highest variance of all unbiased mixtures.

Lemma 4 *When q is Gaussian and $\lambda = 1$, the optimal ψ that maximizes the approximate objective $\mathcal{L}_{\mathcal{F}}$ is both **unbiased** and has **non-overlapping components**.*

In other words, Lemma 4 states that the solution to the approximate objective $\mathcal{L}_{\mathcal{F}}$ “looks like” sampling when $\lambda = 1$, in the sense of Definition 1.

Proof: Without loss of generality, let us assume that θ is already parameterized in terms of its location and scale, $[\boldsymbol{\mu}, \boldsymbol{\sigma}]$, where $\boldsymbol{\mu}$ determines the mean of q and $\boldsymbol{\sigma}$ determines its covariance. Then, the Fisher Information Matrix is a block-diagonal matrix:¹

$$\mathcal{F}(\theta) = \begin{bmatrix} \mathcal{F}(\boldsymbol{\mu}) & 0 \\ 0 & \mathcal{F}(\boldsymbol{\sigma}) \end{bmatrix}$$

where

$$\begin{aligned} \mathcal{F}(\boldsymbol{\mu}) &= \Lambda \\ \mathcal{F}(\boldsymbol{\sigma})_{ij} &= \frac{1}{2} \text{Tr} \left(\Lambda \frac{\partial \Sigma}{\partial \sigma_i} \Lambda \frac{\partial \Sigma}{\partial \sigma_j} \right). \end{aligned}$$

Λ and Σ are the precision matrix and covariance matrix of q , respectively. Both Λ and Σ are functions of the parameters $\boldsymbol{\sigma}$ but not of $\boldsymbol{\mu}$. To simplify further, consider a coordinate system where the covariance of q is diagonal, and that σ_i is the log standard deviation of the i th dimension of \mathbf{x} :

$$\Sigma(\boldsymbol{\sigma})_{ij} = \begin{cases} e^{2\sigma_i} & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

We emphasize that this simplification is for notational convenience only, and other parameterizations of $\Sigma(\boldsymbol{\sigma})$ are permissible (e.g. additionally parameterizing the orientation of Σ with a rotation matrix). With this assumption, $\mathcal{F}(\boldsymbol{\sigma})$ becomes the identity matrix, and the log determinant of $\mathcal{F}(\theta)$ becomes simply

$$\log |\mathcal{F}(\theta)| = \log |\Lambda|.$$

So, for Gaussian q , the expression for ψ becomes

$$\log \psi(\theta) = \log \psi(\boldsymbol{\mu}, \boldsymbol{\sigma}) = \frac{1}{2} \log |\Lambda(\boldsymbol{\sigma})| - \lambda \text{KL}(q(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma}) || p(\mathbf{x})).$$

Next, we will split $\text{KL}(q||p)$ into separate entropy and cross-entropy terms:

$$\begin{aligned} \text{KL}(q||p) &= \mathbb{E}_{q(\mathbf{x};\theta)} [\log q(\mathbf{x}; \theta)] - \mathbb{E}_{q(\mathbf{x};\theta)} [\log p(\mathbf{x})] \\ &= -\mathcal{H}[q] + \mathcal{CE}[q||p]. \end{aligned}$$

¹https://en.wikipedia.org/wiki/Fisher_information#Multivariate_normal_distribution

And note that when q is Gaussian, its entropy is given by

$$\mathcal{H}[q] = \frac{1}{2} \log |2\pi e \Sigma| = \frac{1}{2} \log |\Sigma| + \text{constants}.$$

Taking $\lambda = 1$ and using the fact that $\log |\Sigma| = -\log |\Sigma^{-1}| = -\log |\Lambda|$ and combining the above three equations, the $\mathcal{H}[q]$ and $\log |\mathcal{F}(\boldsymbol{\mu})|$ terms cancel in ψ and we are left – up to additive constants – with

$$\log \psi(\theta) = -\mathcal{CE}[q||p] = \mathbb{E}_{q(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma})} [\log p(\mathbf{x})]. \quad (\text{A.1})$$

To summarize, equation (A.1) says that, using Gaussian components and letting $\lambda \rightarrow 1$, our method, derived from the $\mathcal{I}_{\mathcal{F}}$ approximation to \mathcal{I} , selects components simply according to the *cross entropy* between $q(\mathbf{x}; \theta)$ and $p(\mathbf{x})$.

Note that (A.1) is not a proper distribution over θ . To see this, consider any sufficiently narrow component such that q behaves like a Dirac delta, or $\mathbb{E}_{q(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma})} [\log p(\mathbf{x})] \approx \log p(\boldsymbol{\mu})$. Wherever this holds for some $\boldsymbol{\sigma}$, it will additionally hold for all *narrower* components at the same $\boldsymbol{\mu}$.² Therefore, below a particular scale where q behaves like a Dirac delta, (A.1) places uniform mass on the infinitely many q s that are at least as narrow. This effect is visible in the top-right panel of Figure 2. Also note that ψ is only improper for $\lambda = 1$; for all other $\lambda > 1$, a $(\lambda - 1)\mathcal{H}[q]$ term remains, and ψ cannot place arbitrarily much mass on arbitrarily narrow components.

Despite its impropriety, we are free to draw samples of θ from this improper ψ when $\lambda = 1$ [Besag et al., 1995, Hobert and Casella, 1996]. We will then find that with probability approaching 1 we only ever see components that “look like” Dirac-deltas. This phenomenon is seen empirically in all of our experiments where we set $\lambda = 1$ and run HMC dynamics drawing $\theta \sim \psi(\theta)$ (in practice, we set a lower bound on $\log \sigma$ for numerical stability). Since components will become arbitrarily narrow with high probability, we have that $q(\mathbf{x}; \theta^{(j)}) \ll q(\mathbf{x}; \theta^{(i)})$ in the region where $q(\mathbf{x}; \theta^{(i)})$ has appreciable mass. This means that the mixture will consist of **non-overlapping components** when $\lambda = 1$.

The fact that each component shrinks towards a Dirac delta with high probability then implies that the mixture will be unbiased. To see this, consider decomposing $\psi(\theta)$ into $\psi(\boldsymbol{\sigma})\psi(\boldsymbol{\mu}|\boldsymbol{\sigma})$. The previous paragraph establishes that the marginal distribution $\psi(\boldsymbol{\sigma})$ will allocate effectively all samples to parts of θ -space where components behave like Dirac deltas. This implies

$$\begin{aligned} \log \psi(\boldsymbol{\mu}|\boldsymbol{\sigma} = \text{narrow}) &= \mathbb{E}_{q(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\sigma})} [\log p(\mathbf{x})] \\ &= \log p(\boldsymbol{\mu}). \end{aligned}$$

In other words, when components are narrow, the distribution of means $\boldsymbol{\mu}$ according to ψ will match the true distribution p . Hence, $m(\mathbf{x})$ will be a mixture of Dirac-delta-like components, each of which is chosen in proportion to the true probability of its mean, $p(\boldsymbol{\mu})$. This means that $m(\mathbf{x})$ will be **unbiased** when $\lambda = 1$. ■

Theorem 1 (Improve on sampling) *If a mixture is sampling as in Definition 1, then $\frac{d}{d\lambda} \text{KL bias} = 0$ and $\frac{d}{d\lambda} \text{KL variance} < 0$. Thus, $\frac{d}{d\lambda} \text{KL error} < 0$.*

Proof: Our approach will be to calculate the variational derivatives of KL bias and KL variance with respect to ψ , then take the inner product (directional derivative) with the change in ψ per change in λ .

First, we need the sensitivity of ψ to changes in λ . Recall that the closed-form solution for ψ we get from solving $\mathcal{L}_{\mathcal{F}}$ is

$$\log \psi(\theta) = \frac{1}{2} \log |\mathcal{F}(\theta)| - \lambda \text{KL}(q(\mathbf{x}; \theta)||p(\mathbf{x})) - \log Z(\lambda).$$

The sensitivity of $\log \psi$ to λ is

$$\begin{aligned} \frac{d}{d\lambda} \log \psi(\theta) &= -\text{KL}(q||p) + \frac{1}{Z} \int_{\theta'} e^{\frac{1}{2} \log |\mathcal{F}(\theta)| - \lambda \text{KL}(q||p)} \text{KL}(q||p) d\theta' \\ &= \mathbb{E}_{\psi} [\text{KL}(q||p)] - \text{KL}(q||p). \end{aligned}$$

²There is an implicit assumption here that $\log p(\mathbf{x})$ is almost everywhere smooth, so that there is some small enough scale at which $p(\mathbf{x})$ appears locally linear under q .

Converting from $\log \psi$ to ψ , we get

$$\frac{d}{d\lambda} \psi(\theta) = \psi(\theta) (\mathbb{E}_{\psi}[\text{KL}(q||p)] - \text{KL}(q||p)) \quad (\text{A.2})$$

Recall that we defined KL bias = $\text{KL}(m||p)$ and KL variance = $\mathbb{E}[\text{KL}(m_T||m)]$. The variational derivative of KL bias with respect to ψ , evaluated at θ is³

$$\begin{aligned} \nabla_{\psi} \text{KL}(m||p) &= \nabla_{\psi} \int_{\mathbf{x}} (\mathbb{E}_{\psi}[q(\mathbf{x}; \theta)]) \log \frac{(\mathbb{E}_{\psi}[q(\mathbf{x}; \theta)])}{p(\mathbf{x})} d\mathbf{x} \\ &= \int_{\mathbf{x}} \left(m(\mathbf{x}) \frac{p(\mathbf{x})}{m(\mathbf{x})} \frac{q(\mathbf{x}; \theta)}{p(\mathbf{x})} + q(\mathbf{x}; \theta) \log \frac{m(\mathbf{x})}{p(\mathbf{x})} \right) d\mathbf{x} \\ &= 1 + \mathbb{E}_{q(\mathbf{x}; \theta)} \left[\log \frac{m(\mathbf{x})}{p(\mathbf{x})} \right]. \end{aligned} \quad (\text{A.3})$$

To get the sensitivity of KL bias to λ at the point $\lambda = 1$, we will take the inner-product of (A.2) with (A.3). This is

$$\begin{aligned} \frac{d}{d\lambda} \text{KL bias} &= \left\langle \frac{d\text{KL bias}}{d\psi}, \frac{d\psi}{d\lambda} \right\rangle \\ &= \int_{\theta} \left(1 + \mathbb{E}_{q(\mathbf{x}; \theta)} \left[\log \frac{m(\mathbf{x})}{p(\mathbf{x})} \right] \right) \psi(\theta) (\mathbb{E}_{\psi}[\text{KL}(q||p)] - \text{KL}(q||p)) d\theta \\ &= \int_{\theta} (1 + 0) \psi(\theta) (\mathbb{E}_{\psi}[\text{KL}(q||p)] - \text{KL}(q||p)) d\theta \quad (\text{unbiased}) \\ &= \mathbb{E}_{\psi}[\text{KL}(q||p)] - \mathbb{E}_{\psi}[\text{KL}(q||p)] \\ &= 0. \end{aligned}$$

So, we can conclude that in the sampling limit, small changes in λ have no effect on KL bias. Geometrically, this tells us the Pareto front is tangent to the $y=x$ line in that limit, as illustrated in Figure 2.

Next we will consider the variational derivative of KL variance with respect to ψ , where

$$\begin{aligned} \text{KL variance} &\equiv \mathbb{E}_{1..T}[\text{KL}(m_T||m)] \\ &= \mathbb{E}_{1..T} \left[\int_{\mathbf{x}} \left(\frac{1}{T} \sum_{i=1}^T q(\mathbf{x}; \theta^{(i)}) \right) \log \frac{\left(\frac{1}{T} \sum_{j=1}^T q(\mathbf{x}; \theta^{(j)}) \right)}{m(\mathbf{x})} d\mathbf{x} \right] \end{aligned}$$

using the shorthand $\mathbb{E}_{1..T}[\dots]$ to denote an expectation over independent draws of $\{\theta^{(t)}\} \sim \psi(\theta)$, for each of $t = \{1..T\}$. Applying the assumption of **non-overlapping components**, the sum inside the log is dominated by its maximum. We can therefore approximate KL variance as

$$\text{KL variance} \approx \mathbb{E}_{1..T} \left[\int_{\mathbf{x}} \frac{1}{T} \sum_{i=1}^T q(\mathbf{x}; \theta^{(i)}) \log \frac{\frac{1}{T} q(\mathbf{x}; \theta^{(i)})}{m(\mathbf{x})} d\mathbf{x} \right],$$

since the maximum of the sum over j inside the log will be the i th component from outside the log. This step is most applicable for small to moderate T , since when T grows sufficiently large, even narrow components will overlap each other with appreciable probability. Let us continue assuming that T is sufficiently small and that components are sufficiently non-overlapping. By symmetry, we can remove the sum over i and simplify the outer expectation to a single θ :

$$\dots = \mathbb{E}_i \left[\int_{\mathbf{x}} q(\mathbf{x}; \theta^{(i)}) \log \frac{\frac{1}{T} q(\mathbf{x}; \theta^{(i)})}{m(\mathbf{x})} d\mathbf{x} \right],$$

³This section is best viewed in color. Our notation uses red θ to indicate the value where the variational derivative is evaluated, which is distinct from from black θ s, which are integrated out.

which simplifies to

$$\text{KL variance} \approx \mathcal{I}[\mathbf{x}; \theta] - \log T.$$

The variational derivative of $\mathcal{I}[\mathbf{x}; \theta]$ with respect to ψ is

$$\begin{aligned} \nabla_{\psi} \mathcal{I}[\mathbf{x}; \theta] \Big|_{\theta} &\approx \nabla_{\psi} \int_{\theta} \psi(\theta) \int_{\mathbf{x}} q(\mathbf{x}; \theta) \log \frac{q(\mathbf{x}; \theta)}{m(\mathbf{x})} d\mathbf{x} d\theta \Big|_{\theta} \\ &= - \int_{\theta} \psi(\theta) \int_{\mathbf{x}} q(\mathbf{x}; \theta) \frac{m(\mathbf{x})}{q(\mathbf{x}; \theta)} \frac{q(\mathbf{x}; \theta)}{m(\mathbf{x})^2} q(\mathbf{x}; \theta) d\mathbf{x} d\theta + \int_{\mathbf{x}} q(\mathbf{x}; \theta) \log \frac{q(\mathbf{x}; \theta)}{m(\mathbf{x})} d\mathbf{x} \\ &= -1 + \text{KL}(q(\mathbf{x}; \theta) \| m(\mathbf{x})). \end{aligned}$$

Taking the inner product with $\frac{d}{d\lambda} \psi$, and applying the assumptions from the definition of sampling,

$$\begin{aligned} \frac{d}{d\lambda} \text{KL variance} &= \left\langle \frac{d\text{KL variance}}{d\psi}, \frac{d\psi}{d\lambda} \right\rangle \\ &\approx \left\langle \frac{d\mathcal{I}[\mathbf{x}; \theta]}{d\psi}, \frac{d\psi}{d\lambda} \right\rangle && \text{(non-overlapping)} \\ &= \int_{\theta} (-1 + \text{KL}(q \| m)) \psi(\theta) (\mathbb{E}_{\psi}[\text{KL}(q \| p)] - \text{KL}(q \| p)) d\theta \\ &= \int_{\theta} (-1 + \text{KL}(q \| p)) \psi(\theta) (\mathbb{E}_{\psi}[\text{KL}(q \| p)] - \text{KL}(q \| p)) d\theta && \text{(unbiased)} \\ &= -\mathbb{E}_{\psi(\theta)} [(\text{KL}(q \| p) - \mathbb{E}_{\psi}[\text{KL}(q \| p)]) \text{KL}(q \| p)] \\ &= -\text{var}(\text{KL}(q \| p)). \end{aligned}$$

In other words, this says that the change in the (upper bound on) KL variance is *non-positive*, and its magnitude is given by the variance of the values taken by $\text{KL}(q \| p)$ across all θ .

To summarize, we have shown that, in the sampling limit, where $\lambda = 1$, we have $\frac{d}{d\lambda} \text{KL bias} = 0$ and $\frac{d}{d\lambda} \text{KL variance} \leq 0$, which proves the theorem. \blacksquare

A.3 VI-LIKE BEHAVIOR OF OUR METHOD

Definition 2 (VI limit) We model the large λ limit of our method using a Laplace approximation around the optimal $\theta^* = \arg \min_{\theta} \text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x}))$:

$$\begin{aligned} \psi(\theta) &\approx \mathcal{N}(\theta; \theta^*, \Sigma^*) \\ \text{where } \Sigma^{*-1} &= \lambda \nabla_{\theta}^2 \text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x})) \Big|_{\theta^*}. \end{aligned} \tag{A.4}$$

In other words, we approximate ψ by a normal distribution whose mean is θ^* and whose precision is set by the curvature of $\text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x}))$ and scales with λ . We will assume, for the purposes of proofs related to the VI limit, that there is a single optimal θ^* . As long as $\nabla^2 \text{KL}(q \| p)$ is positive definite, which is guaranteed by the assumption that θ^* is unique, the accuracy of this Laplace approximation can be made arbitrarily good by considering larger and larger λ .

Theorem 2 (Improve on VI) Assume that $q(\mathbf{x}; \theta^*)$ is poorly matched to $p(\mathbf{x})$, in the sense that $\text{Tr}((\nabla_{\theta}^2 \text{KL}(q \| p))^{-1} \mathcal{F}) > |\theta|$, and that λ is sufficiently large to use a Laplace approximation to ψ around θ^* . Then, there exists some finite $T_0 > 1$ such that for all $T \geq T_0$, $\frac{d}{d\lambda} \text{KL error} > 0$.

Proof: As λ grows, the Laplace approximation in (A.4) becomes increasingly accurate, and increasingly narrow. Thus, for sufficiently large λ , we can accurately approximate expectations under ψ using a second order Taylor approximation to the integrand. The general rule for multivariate Gaussians is

$$\mathbb{E}_{\mathcal{N}(\mathbf{y}; \mu, \Sigma)} [f(\mathbf{y})] \approx f(\mu) + \frac{1}{2} \text{Tr}(\Sigma \nabla_{\mathbf{y}}^2 f) \Big|_{\mu}$$

Recall that we defined KL error as $\mathbb{E}_{1..T}[\text{KL}(m_T(\mathbf{x})||p(\mathbf{x}))]$. Approximating each $\psi(\theta^{(t)})$ as a multivariate Gaussian, their product is also a multivariate Gaussian whose collective covariance is block-diagonal⁴ containing T copies of Σ^* from (A.4), and whose collective mean is θ^* for each component. At this mean value where all T components' parameters are equal to θ^* , $m_T(\mathbf{x})$ becomes $q(\mathbf{x}; \theta^*)$. Hence, applying the Taylor series approximation to KL error, the $f(\mu)$ term is just $\text{KL}(q(\mathbf{x}; \theta^*)||p(\mathbf{x}))$. The second term is

$$\frac{1}{2} \text{Tr} \left(\begin{bmatrix} \Sigma^* & & & 0 \\ & \Sigma^* & & \\ & & \ddots & \\ 0 & & & \Sigma^* \end{bmatrix} \nabla_{\theta_1, \dots, \theta^{(t)}}^2 \text{KL}(m_T||p) \right).$$

First, note that the zeros in the off-block-diagonal terms on the left mean that we can ignore interactions between θ s across different mixture components in the Hessian term on the right. Second, there is T -fold symmetry between all components. So, this simplifies to

$$\frac{T}{2} \text{Tr} (\Sigma^* \nabla_{\theta_1}^2 \text{KL}(m_T||p)) = \frac{T}{2\lambda} \text{Tr} ((\nabla_{\theta}^2 \text{KL}(q||p))^{-1} \nabla_{\theta_1}^2 \text{KL}(m_T||p)).$$

Next, since this Hessian is being evaluated around the point θ^* , all of $\theta_2, \dots, \theta^{(t)}$ are equal to θ^* , and we can write the mixture as a function only of the component parameters we are varying in the Hessian. Call this mixture, with $T - 1$ components set to the variational solution, " m_T^* ," defined as

$$m_T^*(\mathbf{x}; \theta) = \frac{T-1}{T} q(\mathbf{x}; \theta^*) + \frac{1}{T} q(\mathbf{x}; \theta).$$

Next we will calculate the Hessian of $\text{KL}(m_T^*(\mathbf{x}; \theta)||p(\mathbf{x}))$. Note that the derivatives are with respect to θ , not θ^* . First, the Hessian of $\text{KL}(q||p)$ is

$$\begin{aligned} \frac{\partial^2}{\partial \theta_j \partial \theta_i} \text{KL}(q(\mathbf{x}; \theta)||p(\mathbf{x})) &= \frac{\partial^2}{\partial \theta_j \partial \theta_i} \int_{\mathbf{x}} q(\mathbf{x}; \theta) \log \frac{q(\mathbf{x}; \theta)}{p(\mathbf{x})} d\mathbf{x} \\ &= \frac{\partial}{\partial \theta_j} \int_{\mathbf{x}} \left[\left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) \left(1 + \log \frac{q(\mathbf{x}; \theta)}{p(\mathbf{x})} \right) \right] d\mathbf{x} \\ &= \int_{\mathbf{x}} \left[\left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) \left(\frac{\frac{\partial}{\partial \theta_j} q(\mathbf{x}; \theta)}{q(\mathbf{x}; \theta^*)} \right) + \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} q(\mathbf{x}; \theta) \right) \left(1 + \log \frac{q(\mathbf{x}; \theta^*)}{p(\mathbf{x})} \right) \right] d\mathbf{x} \\ (*) &= \int_{\mathbf{x}} \frac{\left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) \left(\frac{\partial}{\partial \theta_j} q(\mathbf{x}; \theta) \right)}{q(\mathbf{x}; \theta^*)} d\mathbf{x} + \int_{\mathbf{x}} \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} q(\mathbf{x}; \theta) \right) \log \frac{q(\mathbf{x}; \theta^*)}{p(\mathbf{x})} d\mathbf{x} \\ &= \mathcal{F}(\theta^*) + M(\theta^*). \end{aligned} \tag{A.5}$$

In line (*) we used the fact that $\int_{\mathbf{x}} \nabla_{\theta}^2 q(\mathbf{x}; \theta) d\mathbf{x} = \nabla_{\theta}^2 \int_{\mathbf{x}} q(\mathbf{x}; \theta) d\mathbf{x} = \nabla_{\theta}^2 1 = 0$. In the last line, we recognized the first term as the Fisher Information Matrix $\mathcal{F}(\theta^*)$, and we have defined $M(\theta) = \int_{\mathbf{x}} \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} q(\mathbf{x}; \theta) \right) \log \frac{q(\mathbf{x}; \theta)}{p(\mathbf{x})} d\mathbf{x}$ as a placeholder.

⁴This assumes the T components are statistically independent draws from $\psi(\theta)$. The approach outlined here could be generalized to include correlation between θ s in the off-block-diagonals to model variance of an autocorrelated chain of θ values.

Following a similar derivation, the Hessian of $\text{KL}(m_T^*(\mathbf{x}; \theta) \| p(\mathbf{x}))$ is

$$\begin{aligned}
\frac{\partial^2}{\partial \theta_j \partial \theta_i} \text{KL}(m_T^*(\mathbf{x}; \theta) \| p(\mathbf{x})) &= \frac{\partial^2}{\partial \theta_j \partial \theta_i} \int_{\mathbf{x}} \left(\frac{T-1}{T} q(\mathbf{x}; \theta^*) + \frac{1}{T} q(\mathbf{x}; \theta) \right) \log \frac{\left(\frac{T-1}{T} q(\mathbf{x}; \theta^*) + \frac{1}{T} q(\mathbf{x}; \theta) \right)}{p(\mathbf{x})} d\mathbf{x} \\
&= \frac{\partial}{\partial \theta_j} \int_{\mathbf{x}} \left[\frac{1}{T} \left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) + \frac{1}{T} \left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) \log \frac{\left(\frac{T-1}{T} q(\mathbf{x}; \theta^*) + \frac{1}{T} q(\mathbf{x}; \theta) \right)}{p(\mathbf{x})} \right] d\mathbf{x} \\
&= \frac{1}{T} \frac{\partial}{\partial \theta_j} \int_{\mathbf{x}} \left[\left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) \left(1 + \log \frac{\left(\frac{T-1}{T} q(\mathbf{x}; \theta^*) + \frac{1}{T} q(\mathbf{x}; \theta) \right)}{p(\mathbf{x})} \right) \right] d\mathbf{x} \\
&= \frac{1}{T} \int_{\mathbf{x}} \left[\left(\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \right) \left(\frac{\frac{1}{T} \frac{\partial}{\partial \theta_j} q(\mathbf{x}; \theta)}{m_T^*(\mathbf{x}; \theta)} \right) + \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} q(\mathbf{x}; \theta) \right) \left(1 + \log \frac{m_T^*(\mathbf{x}; \theta)}{p(\mathbf{x})} \right) \right] d\mathbf{x} \\
(**) &= \frac{1}{T^2} \int_{\mathbf{x}} \left(\frac{\frac{\partial}{\partial \theta_i} q(\mathbf{x}; \theta) \frac{\partial}{\partial \theta_j} q(\mathbf{x}; \theta)}{q(\mathbf{x}; \theta^*)} \right) d\mathbf{x} + \frac{1}{T} \int_{\mathbf{x}} \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} q(\mathbf{x}; \theta) \right) \log \frac{q(\mathbf{x}; \theta^*)}{p(\mathbf{x})} d\mathbf{x} \\
&= \frac{1}{T^2} \mathcal{F}(\theta^*) + \frac{1}{T} M(\theta^*) \\
&= \frac{1}{T} (\mathcal{F}(\theta^*) + M(\theta^*)) + \left(\frac{1}{T^2} - \frac{1}{T} \right) \mathcal{F}(\theta^*) \\
&= \frac{1}{T} \frac{\partial^2}{\partial \theta_j \partial \theta_i} \text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x})) + \mathcal{F}(\theta) \left(\frac{1-T}{T^2} \right) \tag{A.6}
\end{aligned}$$

Here, in (**), we additionally used the fact that $m_T^*(\mathbf{x}; \theta^*) = q(\mathbf{x}; \theta^*)$. We then wrote the final line in terms of the Hessian of $\text{KL}(q \| p)$ in (A.5).

To summarize, near the variational limit we have that the KL error is approximately

$$\text{KL}(q(\mathbf{x}; \theta^*) \| p(\mathbf{x})) + \frac{T}{2\lambda} \text{Tr} \left(\underbrace{(\nabla_{\theta}^2 \text{KL}(q \| p))^{-1}}_{(A.5)} \underbrace{(\nabla_{\theta}^2 \text{KL}(m_T^* \| p))}_{(A.6)} \right),$$

and we found that (A.6) could be written in terms of (A.5). To reduce clutter temporarily, let $\mathbf{H} = \nabla_{\theta}^2 \text{KL}(q \| p)$. Combining terms, we have

$$\begin{aligned}
\text{KL error} &\approx \text{KL}(q(\mathbf{x}; \theta^*) \| p(\mathbf{x})) + \frac{T}{2\lambda} \text{Tr} \left(\mathbf{H}^{-1} \left(\frac{1}{T} \mathbf{H} + \frac{1-T}{T^2} \mathcal{F}(\theta^*) \right) \right) \\
&= \text{KL}(q(\mathbf{x}; \theta^*) \| p(\mathbf{x})) + \frac{1}{2\lambda} \text{Tr} \left(\mathbf{I} + \frac{1-T}{T} \mathbf{H}^{-1} \mathcal{F} \right) \\
&= \text{KL}(q(\mathbf{x}; \theta^*) \| p(\mathbf{x})) + \frac{d}{2\lambda} - \frac{1}{2\lambda} \text{Tr} \left(\frac{T-1}{T} \mathbf{H}^{-1} \mathcal{F} \right)
\end{aligned}$$

where \mathbf{I} is the identity matrix and $d = \text{Tr}(\mathbf{I})$ is the dimensionality of θ . Consider the case where $T = 1$: the KL error simplifies to $\text{KL}(q(\mathbf{x}; \theta^*) \| p(\mathbf{x})) + \frac{d}{2\lambda}$. Therefore when $T = 1$, KL error is only reduced by further increasing λ . This is an intuitive result: we cannot reduce bias compared to VI when using a single component, and any stochasticity only adds variance.

Now consider the case where $T \geq 2$. We are interested in cases where KL error *increases* with λ near the VI limit, as this would imply that using a finite λ would improve on VI. This is equivalent to asking when the following inequality holds:

$$\frac{\text{Tr}(\mathbf{H}^{-1} \mathcal{F})}{d} > \frac{T}{T-1}.$$

Recall that \mathbf{H} was defined as the Hessian of $\text{KL}(q(\mathbf{x}; \theta) \| p(\mathbf{x}))$, so this is

$$\frac{\text{Tr} \left((\nabla_{\theta}^2 \text{KL}(q \| p))^{-1} \mathcal{F} \right)}{d} > \frac{T}{T-1}.$$

The Fisher Information Matrix can also be derived from a local quadratic approximation to $\text{KL}(q||q)$; this means that in the case where the VI solution is exact, or $q(\mathbf{x}; \theta^*) = p(\mathbf{x})$, the trace term becomes $\text{Tr}(\mathcal{F}^{-1}\mathcal{F})$, and the inequality is $1 > \frac{T}{T-1}$. This inequality is not satisfied by any positive integer T , and so this expression captures the intuitive condition that VI cannot be improved upon by reducing λ – for any finite T – if the VI solution is already exact.

Conversely, the ratio

$$R \equiv \frac{\text{Tr}((\nabla_{\theta}^2 \text{KL}(q||p))^{-1}\mathcal{F})}{d}$$

may be thought of as quantifying the extent to which $p(\mathbf{x})$ is over-dispersed relative to the VI solution $q(\mathbf{x}; \theta^*)$. If the curvature of $\text{KL}(q||p)$ is low, then many “nearby” q s would fit p almost as well, and this will be reflected in this ratio being larger than 1. Then, the minimum T for which reducing λ improves KL error relative to VI can be found by solving the above inequality, giving which gives

$$T_0 = \left\lceil \frac{R}{R-1} \right\rceil, \tag{A.7}$$

so that for all $T > T_0$, we have the desired property that $\frac{d}{d\lambda} \text{KL error} > 0$, which implies that using some finite $\lambda < \infty$ will reduce error relative to VI. ■

B NUMERICAL DETAILS

All code to generate the figures in this paper is available at <https://github.com/wrongu/sampling-variational-demos>.

We implemented (10) in Stan [Carpenter et al., 2017]. For q , we used the family of multivariate Gaussians with diagonal covariance, parameterized as $\theta = [\mu_1, \dots, \mu_n, \log \sigma_1, \dots, \log \sigma_n]$ where n is the number of unconstrained parameters (i.e the dimensionality of \mathbf{x}). In this parameterization, $\frac{1}{2} \log \mathcal{F}(\theta)$ is simply $-\sum_{i=1}^n \log \sigma_i$. We sampled θ from $\psi(\theta)$ using Stan’s default implementation of the No U-Turn Sampler (NUTS) with automatic step-size adaptation [Hoffman and Gelman, 2014], and we set the mass equal to λ times the identity matrix. NUTS requires both $\text{KL}(q||p)$ and its gradient, which we computed using Monte Carlo samples from q and the reparameterization trick. The reparameterized samples were frozen for each trajectory of NUTS and resampled between trajectories.

We used two toy distributions in the main paper:

- The “banana” distribution over \mathbb{R}^2 , defined as

$$\log p(x, y) = -(y - (x/2)^2)^2 - (x/2)^2.$$

- The “Laplace mixture” distribution over \mathbb{R}^1 , defined as

$$p(x) \propto 0.4e^{\frac{|x+1.5|}{0.75}} + 0.6e^{\frac{|x-1.5|}{0.75}}.$$

We also tested our method on three reference problems taken from posteriordb [Magnusson et al., 2021], a database of reference problems for testing and validating inference methods. These were `arK`, `eight schools centered`, and `garch11`. These problems are 7-, 10-, and 4-dimensional problems, respectively. Finally, we synthesized ground-truth data from a hierarchical regression problem with 30 regressors (a total of 32 parameters) to test how our algorithm would scale to an even higher dimensional problem. Results for all of these additional experiments are shown in Figure B.1.

In our experiments, all functions integrated are sums of sinusoids,

$$f(\mathbf{x}) = \sum_{\omega=1}^N a \sin(\omega \mathbf{t}^T \mathbf{x} + \phi_{\omega})$$

where \mathbf{t} is a random unit vector. This is a convenient target distribution as the integral of a sinusoid under a Gaussian is known analytically:

$$\int_{\mathbf{x}} \sin(\omega \mathbf{t}^T \mathbf{x} + \phi_{\omega}) \mathcal{N}(\mu, \Sigma) = \sin(\omega \mathbf{t}^T \mu + \phi_{\omega}) \exp\left(-\frac{\omega^2}{2} \mathbf{t}^T \Sigma \mathbf{t}\right)$$

The capability for exact integration of $\int_{\mathbf{x}} m_T(\mathbf{x})f(\mathbf{x})d\mathbf{x}$ ensures that no additional variance is introduced in plots; all variance is due to the selection of the components q . In general this integral can be computed with MC methods or, in low enough dimensions, Gaussian quadrature.

In our experiments (Figures 3 and 4) we used $N = 100$ sinusoidal components in $f(\mathbf{x})$, and calculated bias using $T = 5,000$ components thinned from 4 MCMC chains of length 50,000. To calculate variance, we subsampled $T = 10$ components from these chains, and computed variance over these random instantiations of $m_{10}(\mathbf{x})$. The NUTS samples over \mathbf{x} treated as ground truth derive from 4 chains of length 1,000,000. Because variance scales like $1/T$, we estimated variance for other values of T simply by scaling the variance from the $T = 10$ case. We chose $T = 30$ (scaling variance from $T = 10$ by a factor of $1/3$) for the plots using the banana distribution in the main paper because this gave roughly equal magnitude to the variance of NUTS and the bias of ADVI, which makes the effects of the trade-off between them most pronounced. In Figure B.1, we chose T separately for each problem using the same strategy.

C DETAILED COMPARISON TO PRIOR WORK

Author (Year)	Component parameters	Component family	Auxiliary Optimization	Time; Space Complexity	Recovers Sampling
Jaakkola and Jordan [1998] "Mixture of Mean Field"	optimized	mean field	minor	$\mathcal{O}(T^2); \mathcal{O}(T)$	$T \rightarrow \infty$
Salimans et al. [2015] "Bridging the Gap"	sampled (implicit)	flexible	NN optimization (k steps)	$\mathcal{O}(k+T); \mathcal{O}(T)$	unknown
Gershman et al. [2012] "Nonparametric VI"	optimized	Gaussian	none	$\mathcal{O}(T^2); \mathcal{O}(T)$	$T \rightarrow \infty$
Zobay [2014] "VI with Gaussian Mixtures"	optimized	Gaussian	none	$\mathcal{O}(T^2); \mathcal{O}(T)$	$T \rightarrow \infty$
Guo et al. [2016], Miller et al. [2017] "Boosting VI"	optimized	flexible	none	$\mathcal{O}(T^2); \mathcal{O}(T)$	no
Nalisnick and Smyth [2017] "Stein Mixtures"	optimized	flexible	none	$\mathcal{O}(T^2); \mathcal{O}(T)^*$	$T \rightarrow \infty$
Yin and Zhou [2018] "SIVI"	sampled (implicit)	Gaussian*	NN optimization (k steps)	$\mathcal{O}(kT^2); \mathcal{O}(T)^*$	unknown
Acerbi [2018, 2020] "VBMC"	optimized	Gaussian	fit GP (k steps)	$\mathcal{O}(kT^2); \mathcal{O}(T)$	no
Ours	sampled (closed-form)	Gaussian*	none	$\mathcal{O}(T); \mathcal{O}(T)^*$	$\lambda \rightarrow 1$

Table C.1: Comparison of our proposed algorithm to a number of existing methods that, in some way or another, use a mixture of “simple” component distributions for approximate inference. **Component parameters:** how are the parameters of the individual mixture components chosen? In terms of minimizing bias and variance, it is best to jointly optimize the location of all T components together. The approach of Salimans et al. [2015] and Yin and Zhou [2018] is similar to ours in that mixture components are stochastically sampled, though ours is the only stochastic method for which we have the mixing distribution $\psi(\theta)$ explicitly in closed form. **Component family:** what is the allowable form of $q(\mathbf{x}; \theta)$? Methods marked “Gaussian*” (including ours) are in principle applicable to a wider class of distributions, but so far only demonstrated empirically using Gaussian mixtures. **Auxiliary optimization:** Are there additional parameters of the inference process itself that need to be fit or optimized at inference-time? The advantage of our method is that, since we derived $\psi(\theta)$ in closed form, we can begin sampling from it immediately without further optimizations. **Time; Space Complexity:** Note that all methods for which component parameters are “optimized” incur a $\mathcal{O}(T^2)$ cost in time complexity, since the optimal location of each component depends on the location of other components. Methods that require auxiliary optimization (such as training a neural network or NN) incur additional runtime costs. Methods marked $\mathcal{O}(T)^*$ space-complexity can in principle be “streamed,” in which case they use constant $\mathcal{O}(1)$ space. **Recovers sampling:** some methods “look like” sampling (unbiased and with narrow components) only in the $T \rightarrow \infty$ limit, which is infeasible given $\mathcal{O}(T^2)$ time complexity. Ours is the only method we know of for which sampling-like behavior can be recovered, independent of T , by setting $\lambda = 1$.

References

- Luigi Acerbi. Variational Bayesian Monte Carlo. *Advances in Neural Information Processing Systems*, 2018.
- Luigi Acerbi. Variational Bayesian Monte Carlo with Noisy Likelihoods. *arXiv*, 2020. ISSN 10495258.
- S Amari. *Information Geometry and Its Applications*. Applied Mathematical Sciences. Springer Japan, 2016. ISBN 9784431559788. URL <https://books.google.com/books?id=UkSFCwAAQBAJ>.
- Julian Besag, Peter Green, David Higdon, and Kerrie Mengersen. Bayesian computation and stochastic systems. *Statistical Science*, 10(1):3–44, 1995. ISSN 08834237. doi: 10.1214/ss/1177010123.
- Mark Braverman and Abhishek Bhowmick. Convexity/concavity of mutual information, September 2011. URL <https://www.cs.princeton.edu/courses/archive/fall11/cos597D/L04.pdf>.
- Bob Carpenter, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus A. Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of Statistical Software*, 76(1), 2017. ISSN 15487660. doi: 10.18637/jss.v076.i01.
- Samuel J. Gershman, Matthew D. Hoffman, and David M. Blei. Nonparametric Variational Inference. *Proceedings of the 29th International Conference on Machine Learning*, pages 235–242, 2012. ISSN 0899-7667. doi: 10.1162/089976699300016331. URL <https://icml.cc/Conferences/2012/papers/360.pdf>.
- Fangjian Guo, Xiangyu Wang, Kai Fan, Tamara Broderick, and David B. Dunson. Boosting Variational Inference. *arXiv*, 2016. URL <http://arxiv.org/abs/1611.05559>.
- James P. Hobert and George Casella. The Effect of Improper Priors on Gibbs Sampling in Hierarchical Linear Mixed Models. *Journal of the American Statistical Association*, 91(436):1461–1473, 1996. ISSN 1537274X. doi: 10.1080/01621459.1996.10476714.
- Matthew D Hoffman and Andrew Gelman. The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15:1351–1381, 2014.
- Tommi S. Jaakkola and Michael I. Jordan. Improving the Mean Field Approximation via the Use of Mixture Distributions. In Michael I. Jordan, editor, *Learning in Graphical Models*. Kluwer Academic Publishers, 1998.
- M. Magnusson, Paul-Christian Bürkner, and Aki Vehtari. posteriordb: A database of Bayesian posterior inference, 2021. URL <https://github.com/stan-dev/posteriordb>.
- Andrew C. Miller, Nicholas J. Foti, and Ryan P. Adams. Variational Boosting: Iteratively Refining Posterior Approximations. *arXiv*, 2017. URL <http://arxiv.org/abs/1611.06585>.
- Eric Nalisnick and Padhraic Smyth. Variational Inference with Stein Mixtures. *NIPS2017 (Workshop)*, 2017. ISSN 00368075. doi: 10.1126/science.1070850. URL [https://www.ics.uci.edu/~\\$sim\\$nalisni/AABI_paper30-Stein_Mixtures.pdf](https://www.ics.uci.edu/~simnalisni/AABI_paper30-Stein_Mixtures.pdf).
- Tim Salimans, Diederik P. Kingma, and Max Welling. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. *Proceedings of the 32nd International Conference on Machine Learning*, pages 1218–1226, 2015. URL <http://arxiv.org/abs/1410.6460>.
- Mingzhang Yin and Mingyuan Zhou. Semi-Implicit Variational Inference. *International Conference on Machine Learning*, 35, 2018.
- O. Zoby. Variational Bayesian inference with Gaussian-mixture approximations. *Electronic Journal of Statistics*, 8(1): 355–389, 2014. ISSN 19357524. doi: 10.1214/14-EJS887.

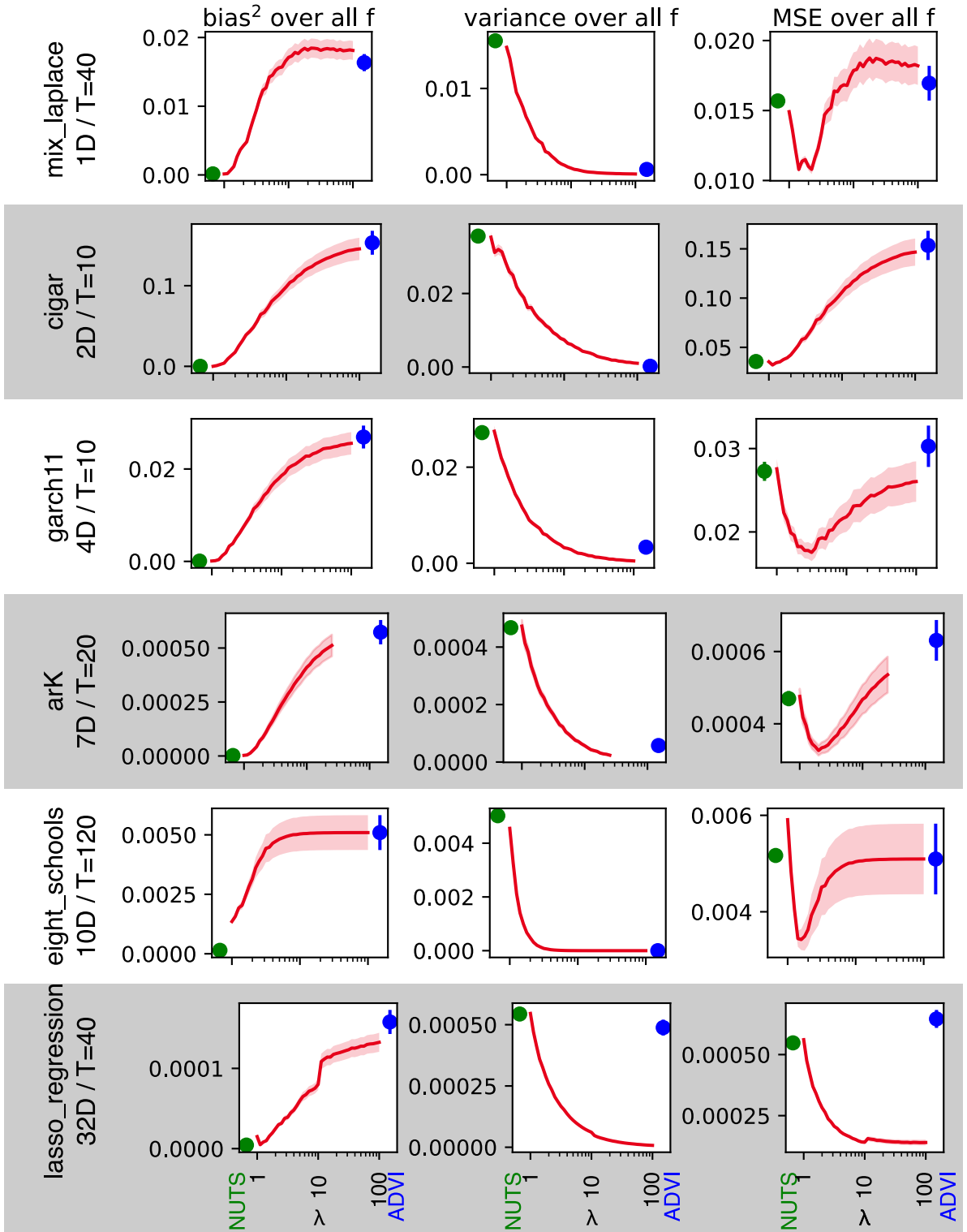


Figure B.1: Bias/variance decomposition across many random f_s (all with $\alpha = -1.5$) for additional test problems, plotted in the same format as Figure 3d-f. Each row corresponds to a different inference problem. We chose T separately for each problem such that the variance of NUTS and bias of ADVI were of similar orders of magnitude. **mix laplace** is the 1D mixture of two Laplace distributions shown in Figure 2 of the main paper. **cigar** is a 2D Gaussian with a correlation of 0.99, oriented along the $y = x$ axis. **garch11**, **eight_schools**, and **arK** refer to problems taken from posteriondb [Magnusson et al., 2021]. **lasso regression** is a regression problem with 30 regressors and a double-exponential (Laplace) prior on each regression weight. We report the variance of ADVI using all default settings, which fails to converge reliably for this problem, and hence has significant variance. Optimizing the ADVI settings for this problem could potentially result in much lower variance.